# Information Technology —
# Portable Operating System Interface (POSIX®)

# Rationale

Sponsor

**Portable Applications Standards Committee**
of the
**IEEE Computer Society**

and

**The Open Group**

◆IEEE   THE *Open* GROUP

[This page intentionally left blank]

**Abstract**

This standard is simultaneously ISO/IEC 9945:2002, IEEE Std 1003.1-2001, and forms the core of the Single UNIX Specification, Version 3.

The IEEE Std 1003.1, 2003 Edition includes IEEE Std 1003.1-2001/Cor 1-2002 incorporated into IEEE Std 1003.1-2001 (base document). The Corrigendum addresses problems discovered since the approval of IEEE Std 1003.1-2001. These changes are mainly due to resolving integration issues raised by the merger of the base documents that were incorporated into IEEE Std 1003.1-2001, which is the single common revision to IEEE Std 1003.1™-1996, IEEE Std 1003.2™-1992, ISO/IEC 9945-1:1996, ISO/IEC 9945-2:1993, and the Base Specifications of The Open Group Single UNIX® Specification, Version 2.

This standard defines a standard operating system interface and environment, including a command interpreter (or ''shell''), and common utility programs to support applications portability at the source code level. This standard is intended to be used by both applications developers and system implementors and comprises four major components (each in an associated volume):

- General terms, concepts, and interfaces common to all volumes of this standard, including utility conventions and C-language header definitions, are included in the Base Definitions volume.

- Definitions for system service functions and subroutines, language-specific system services for the C programming language, function issues, including portability, error handling, and error recovery, are included in the System Interfaces volume.

- Definitions for a standard source code-level interface to command interpretation services (a ''shell'') and common utility programs for application programs are included in the Shell and Utilities volume.

- Extended rationale that did not fit well into the rest of the document structure, which contains historical information concerning the contents of this standard and why features were included or discarded by the standard developers, is included in the Rationale (Informative) volume.

The following areas are outside the scope of this standard:

- Graphics interfaces

- Database management system interfaces

- Record I/O considerations

- Object or binary code portability

- System configuration and resource availability

This standard describes the external characteristics and facilities that are of importance to applications developers, rather than the internal construction techniques employed to achieve these capabilities. Special emphasis is placed on those functions and facilities that are needed in a wide variety of commercial applications.

**Keywords**

application program interface (API), argument, asynchronous, basic regular expression (BRE), batch job, batch system, built-in utility, byte, child, command language interpreter, CPU, extended regular expression (ERE), FIFO, file access control mechanism, input/output (I/O), job control, network, portable operating system interface (POSIX®), parent, shell, stream, string, synchronous, system, thread, X/Open System Interface (XSI)

Rationale (Informative)

**Permissions**

**Feedback**

This standard has been prepared by the Austin Group. Feedback relating to the material contained in this standard may be submitted using the Austin Group web site at *http://www.opengroup.org/austin/defectform.html*.

**IEEE**

IEEE Standards documents are developed within the IEEE Societies and the Standards Coordinating Committees of the IEEE Standards Association (IEEE-SA) Standards Board. The IEEE develops its standards through a consensus development process, approved by the American National Standards Institute, which brings together volunteers representing varied viewpoints and interests to achieve the final product. Volunteers are not necessarily members of the Institute and serve without compensation. While the IEEE administers the process and establishes rules to promote fairness in the consensus development process, the IEEE does not independently evaluate, test, or verify the accuracy of any of the information contained in its standards.

Use of an IEEE Standard is wholly voluntary. The IEEE disclaims liability for any personal injury, property, or other damage, of any nature whatsoever, whether special, indirect, consequential, or compensatory, directly or indirectly resulting from the publication, use of, or reliance upon this, or any other IEEE Standard document.

**The IEEE does not warrant or represent the accuracy or content of the material contained herein, and expressly disclaims any express or implied warranty, including any implied warranty of merchantability or fitness for a specific purpose, or that the use of the material contained herein is free from patent infringement. IEEE Standards documents are supplied ''AS IS''.**

The existence of an IEEE Standard does not imply that there are no other ways to produce, test, measure, purchase, market, or provide other goods and services related to the scope of the IEEE Standard. Furthermore, the viewpoint expressed at the time a standard is approved and issued is subject to change brought about through developments in the state of the art and comments received from users of the standard. Every IEEE Standard is subjected to review at least every five years for revision or reaffirmation. When a document is more than five years old and has not been reaffirmed, it is reasonable to conclude that its contents, although still of some value, do not wholly reflect the present state of the art. Users are cautioned to check to determine that they have the latest edition of any IEEE Standard.

In publishing and making this document available, the IEEE is not suggesting or rendering professional or other services for, or on behalf of, any person or entity. Nor is the IEEE undertaking to perform any duty owed by any other person or entity to another. Any person utilizing this, and any other IEEE Standards document, should rely upon the advice of a competent professional in determining the exercise of reasonable care in any given circumstances.

Interpretations: Occasionally questions may arise regarding the meaning of portions of standards as they relate to specific applications. When the need for interpretations is brought to the attention of the IEEE, the Institute will initiate action to prepare appropriate responses. Since IEEE Standards represent a consensus of concerned interests, it is important to ensure that any interpretation has also received the concurrence of a balance of interests. For this reason, IEEE and the members of its societies and Standards Coordinating Committees are not able to provide an instant response to interpretation requests except in those cases where the matter has previously received formal consideration.

Comments for revision of IEEE Standards are welcome from any interested party, regardless of membership affiliation with the IEEE.[1] Suggestions for changes in documents should be in the form of a proposed change of text, together with appropriate supporting comments. Comments on standards and requests for interpretations should be addressed to:

Secretary, IEEE-SA Standards Board, 445 Hoes Lane, P.O. Box 1331, Piscataway, NJ 08855-1331, U.S.A.

---

Attention is called to the possibility that implementation of this standard may require use of subject matter covered by patent rights. By publication of this standard, no position is taken with respect to the existence or validity of any patent rights in connection therewith. The IEEE shall not be responsible for identifying patents for which a license may be required by an IEEE Standard or for conducting inquiries into the legal validity or scope of those patents that are brought to its attention.

A patent holder has filed a statement of assurance that it will grant licenses under these rights without compensation or under reasonable rates and non-discriminatory, reasonable terms and conditions to all applicants desiring to obtain such licenses. The IEEE makes no representation as to the reasonableness of rates and/or terms and conditions of the license agreements offered by patent holders. Further information may be obtained from the IEEE Standards Department.

---

Authorization to photocopy portions of any individual standard for internal or personal use is granted in the U.S. by the Institute of Electrical and Electronics Engineers, Inc., provided that the appropriate fee is paid to the Copyright Clearance Center.[2] Permission to photocopy portions of any individual standard for educational classroom use can also be obtained through the Copyright Clearance Center. To arrange for payment of the licensing fee, please contact:

Copyright Clearance Center, Customer Service, 222 Rosewood Drive, Danvers, MA 01923, U.S.A., Tel.: +1 978 750 8400

Amendments, corrigenda, and interpretations for this standard, or information about the IEEE standards development process, may be found at *http://standards.ieee.org*.

Full catalog and ordering information on all IEEE publications is available from the IEEE Online Catalog & Store at *http://shop.ieee.org/store*.

---

1. For this standard, please send comments via the Austin Group as requested on page iii.
2. Please refer to the special provisions for this standard on page iii concerning permissions from both copyright holders and arrangements to cover photocopying and reproduction across the world, as well as by commercial organizations wishing to license the material for use in product documentation.

**The Open Group**

The Open Group, a vendor and technology-neutral consortium, is committed to delivering greater business efficiency by bringing together buyers and suppliers of information technology to lower the time, cost, and risks associated with integrating new technology across the enterprise.

The Open Group's mission is to offer all organizations concerned with open information infrastructures a forum to share knowledge, integrate open initiatives, and certify approved products and processes in a manner in which they continue to trust our impartiality.

In the global eCommerce world of today, no single economic entity can achieve independence while still ensuring interoperability. The assurance that products will interoperate with each other across differing systems and platforms is essential to the success of eCommerce and business workflow. The Open Group, with its proven testing and certification program, is the international guarantor of interoperability in the new century.

The Open Group provides opportunities to exchange information and shape the future of IT. The Open Group's members include some of the largest and most influential organizations in the world. The flexible structure of The Open Groups membership allows for almost any organization, no matter what their size, to join and have a voice in shaping the future of the IT world.

More information is available on The Open Group web site at *http://www.opengroup.org*.

The Open Group has over 15 years' experience in developing and operating certification programs and has extensive experience developing and facilitating industry adoption of test suites used to validate conformance to an open standard or specification. The Open Group portfolio of test suites includes the *Westwood* family of tests for this standard and the associated certification program for Version 3 of the Single UNIX Specification, as well tests for CDE, CORBA, Motif, Linux, LDAP, POSIX.1, POSIX.2, POSIX Realtime, Sockets, UNIX, XPG4, XNFS, XTI, and X11. The Open Group test tools are essential for proper development and maintenance of standards-based products, ensuring conformance of products to industry-standard APIs, applications portability, and interoperability. In-depth testing identifies defects at the earliest possible point in the development cycle, saving costs in development and quality assurance.

More information is available at *http://www.opengroup.org/testing*.

The Open Group publishes a wide range of technical documentation, the main part of which is focused on development of Technical and Product Standards and Guides, but which also includes white papers, technical studies, branding and testing documentation, and business titles. Full details and a catalog are available at *http://www.opengroup.org/pubs*.

As with all *live* documents, Technical Standards and Specifications require revision to align with new developments and associated international standards. To distinguish between revised specifications which are fully backwards compatible and those which are not:

- A new *Version* indicates there is no change to the definitive information contained in the previous publication of that title, but additions/extensions are included. As such, it *replaces* the previous publication.

- A new *Issue* indicates there is substantive change to the definitive information contained in the previous publication of that title, and there may also be additions/extensions. As such, both previous and new documents are maintained as current publications.

Readers should note that Corrigenda may apply to any publication. Corrigenda information is published at *http://www.opengroup.org/corrigenda*.

Full catalog and ordering information on all Open Group publications is available at *http://www.opengroup.org/pubs*.

# *Contents*

# Contents

*Contents*

*Contents*

### List of Figures

### List of Tables

# *Foreword*

**Structure of the Standard**

This standard was originally developed by the Austin Group, a joint working group of members of the IEEE, members of The Open Group, and members of ISO/IEC Joint Technical Committee 1, as one of the four volumes of IEEE Std 1003.1-2001. The standard was approved by ISO and IEC and published in four parts, correlating to the original volumes.

A mapping of the parts to the volumes is shown below:

| ISO/IEC 9945 Part | IEEE Std 1003.1 Volume | Description |
|---|---|---|
| 9945-1 | Base Definitions | Includes general terms, concepts, and interfaces common to all parts of ISO/IEC 9945, including utility conventions and C-language header definitions. |
| 9945-2 | System Interfaces | Includes definitions for system service functions and subroutines, language-specific system services for the C programming language, function issues, including portability, error handling, and error recovery. |
| 9945-3 | Shell and Utilities | Includes definitions for a standard source code-level interface to command interpretation services (a ''shell'') and common utility programs for application programs. |
| 9945-4 | Rationale | Includes extended rationale that did not fit well into the rest of the document structure, containing historical information concerning the contents of ISO/IEC 9945 and why features were included or discarded by the standard developers. |

All four parts comprise the entire standard, and are intended to be used together to accommodate significant internal referencing among them. POSIX-conforming systems are required to support all four parts.

# *Introduction*

**Note:**   This introduction is not part of IEEE Std 1003.1-2001, Standard for Information Technology — Portable Operating System Interface (POSIX).

This standard has been jointly developed by the IEEE and The Open Group. It is simultaneously an IEEE Standard, an ISO/IEC Standard, and an Open Group Technical Standard.

**The Austin Group**

This standard was developed, and is maintained, by a joint working group of members of the IEEE Portable Applications Standards Committee, members of The Open Group, and members of ISO/IEC Joint Technical Committee 1. This joint working group is known as the Austin Group.[3] The Austin Group arose out of discussions amongst the parties which started in early 1998, leading to an initial meeting and formation of the group in September 1998. The purpose of the Austin Group has been to revise, combine, and update the following standards: ISO/IEC 9945-1, ISO/IEC 9945-2, IEEE Std 1003.1, IEEE Std 1003.2, and the Base Specifications of The Open Group Single UNIX Specification.

After two initial meetings, an agreement was signed in July 1999 between The Open Group and the Institute of Electrical and Electronics Engineers (IEEE), Inc., to formalize the project with the first draft of the revised specifications being made available at the same time. Under this agreement, The Open Group and IEEE agreed to share joint copyright of the resulting work. The Open Group has provided the chair and secretariat for the Austin Group.

The base document for the revision was The Open Group's Base volumes of its Single UNIX Specification, Version 2. These were selected since they were a superset of the existing POSIX.1 and POSIX.2 specifications and had some organizational aspects that would benefit the audience for the new revision.

The approach to specification development has been one of ''write once, adopt everywhere'', with the deliverables being a set of specifications that carry the IEEE POSIX designation, The Open Group's Technical Standard designation, and an ISO/IEC designation. This set of specifications forms the core of the Single UNIX Specification, Version 3.

This unique development has combined both the industry-led efforts and the formal standardization activities into a single initiative, and included a wide spectrum of participants. The Austin Group continues as the maintenance body for this document.

Anyone wishing to participate in the Austin Group should contact the chair with their request. There are no fees for participation or membership. You may participate as an observer or as a contributor. You do not have to attend face-to-face meetings to participate; electronic participation is most welcome. For more information on the Austin Group and how to participate, see *http://www.opengroup.org/austin*.

_____

3.  The Austin Group is named after the location of the inaugural meeting held at the IBM facility in Austin, Texas in September 1998.

**Background**

The developers of this standard represent a cross section of hardware manufacturers, vendors of operating systems and other software development tools, software designers, consultants, academics, authors, applications programmers, and others.

Conceptually, this standard describes a set of fundamental services needed for the efficient construction of application programs. Access to these services has been provided by defining an interface, using the C programming language, a command interpreter, and common utility programs that establish standard semantics and syntax. Since this interface enables application writers to write portable applications—it was developed with that goal in mind—it has been designated POSIX,[4] an acronym for Portable Operating System Interface.

Although originated to refer to the original IEEE Std 1003.1-1988, the name POSIX more correctly refers to a *family* of related standards: IEEE Std 1003.*n* and the parts of ISO/IEC 9945. In earlier editions of the IEEE standard, the term POSIX was used as a synonym for IEEE Std 1003.1-1988. A preferred term, POSIX.1, emerged. This maintained the advantages of readability of the symbol ''POSIX'' without being ambiguous with the POSIX family of standards.

**Audience**

The intended audience for this standard is all persons concerned with an industry-wide standard operating system based on the UNIX system. This includes at least four groups of people:

1. Persons buying hardware and software systems

2. Persons managing companies that are deciding on future corporate computing directions

3. Persons implementing operating systems, and especially

4. Persons developing applications where portability is an objective

**Purpose**

Several principles guided the development of this standard:

- Application-Oriented

  The basic goal was to promote portability of application programs across UNIX system environments by developing a clear, consistent, and unambiguous standard for the interface specification of a portable operating system based on the UNIX system documentation. This standard codifies the common, existing definition of the UNIX system.

- Interface, Not Implementation

  This standard defines an interface, not an implementation. No distinction is made between library functions and system calls; both are referred to as functions. No details of the implementation of any function are given (although historical practice is sometimes indicated in the RATIONALE section). Symbolic names are given for constants (such as signals and error numbers) rather than numbers.

---

4. The name POSIX was suggested by Richard Stallman. It is expected to be pronounced *pahz-icks*, as in *positive*, not *poh-six*, or other variations. The pronunciation has been published in an attempt to promulgate a standardized way of referring to a standard operating system interface.

- Source, Not Object, Portability

  This standard has been written so that a program written and translated for execution on one conforming implementation may also be translated for execution on another conforming implementation. This standard does not guarantee that executable (object or binary) code will execute under a different conforming implementation than that for which it was translated, even if the underlying hardware is identical.

- The C Language

  The system interfaces and header definitions are written in terms of the standard C language as specified in the ISO C standard.

- No Superuser, No System Administration

  There was no intention to specify all aspects of an operating system. System administration facilities and functions are excluded from this standard, and functions usable only by the superuser have not been included. Still, an implementation of the standard interface may also implement features not in this standard. This standard is also not concerned with hardware constraints or system maintenance.

- Minimal Interface, Minimally Defined

  In keeping with the historical design principles of the UNIX system, the mandatory core facilities of this standard have been kept as minimal as possible. Additional capabilities have been added as optional extensions.

- Broadly Implementable

  The developers of this standard endeavored to make all specified functions implementable across a wide range of existing and potential systems, including:

  1. All of the current major systems that are ultimately derived from the original UNIX system code (Version 7 or later)

  2. Compatible systems that are not derived from the original UNIX system code

  3. Emulations hosted on entirely different operating systems

  4. Networked systems

  5. Distributed systems

  6. Systems running on a broad range of hardware

  No direct references to this goal appear in this standard, but some results of it are mentioned in the Rationale (Informative) volume.

- Minimal Changes to Historical Implementations

  When the original version of IEEE Std 1003.1 was published, there were no known historical implementations that did not have to change. However, there was a broad consensus on a set of functions, types, definitions, and concepts that formed an interface that was common to most historical implementations.

  The adoption of the 1988 and 1990 IEEE system interface standards, the 1992 IEEE shell and utilities standard, the various Open Group (formerly X/Open) specifications, and the subsequent revisions and addenda to all of them have consolidated this consensus, and this revision reflects the significantly increased level of consensus arrived at since the original versions. The earlier standards and their modifications specified a number of areas where consensus had not been reached before, and these are now reflected in this revision. The authors of the original versions tried, as much as possible, to follow the principles below

when creating new specifications:

1. By standardizing an interface like one in an historical implementation; for example, directories

2. By specifying an interface that is readily implementable in terms of, and backwards-compatible with, historical implementations, such as the extended *tar* format defined in the *pax* utility

3. By specifying an interface that, when added to an historical implementation, will not conflict with it; for example, the *sigaction*( ) function

This revision tries to minimize the number of changes required to implementations which conform to the earlier versions of the approved standards to bring them into conformance with the current standard. Specifically, the scope of this work excluded doing any ''new'' work, but rather collecting into a single document what had been spread across a number of documents, and presenting it in what had been proven in practice to be a more effective way. Some changes to prior conforming implementations were unavoidable, primarily as a consequence of resolving conflicts found in prior revisions, or which became apparent when bringing the various pieces together.

However, since it references the 1999 version of the ISO C standard, and no longer supports ''Common Usage C'', there are a number of unavoidable changes. Applications portability is similarly affected.

This standard is specifically not a codification of a particular vendor's product.

It should be noted that implementations will have different kinds of extensions. Some will reflect ''historical usage'' and will be preserved for execution of pre-existing applications. These functions should be considered ''obsolescent'' and the standard functions used for new applications. Some extensions will represent functions beyond the scope of this standard. These need to be used with careful management to be able to adapt to future extensions of this standard and/or port to implementations that provide these services in a different manner.

- Minimal Changes to Existing Application Code

A goal of this standard was to minimize additional work for the developers of applications. However, because every known historical implementation will have to change at least slightly to conform, some applications will have to change.

**This Standard**

This standard defines the Portable Operating System Interface (POSIX) requirements and consists of the following volumes:

- Base Definitions

- Shell and Utilities

- System Interfaces

- Rationale (Informative) (this volume)

**This Volume**

This volume is being published to assist in the process of review. It contains historical information concerning the contents of this standard and why features were included or discarded by the standard developers. It also contains notes of interest to application programmers on recommended programming practices, emphasizing the consequences of some aspects of this standard that may not be immediately apparent.

This volume is organized in parallel to the normative volumes of this standard, with a separate part for each of the three normative volumes.

Within this volume, the following terms are used:

**base standard**
   The portions of this standard that are not optional, equivalent to the definitions of *classic* POSIX.1 and POSIX.2.

**POSIX.0**
   Although this term is not used in the normative text of this standard, it is used in this volume to refer to IEEE Std 1003.0-1995.

**POSIX.1b**
   Although this term is not used in the normative text of this standard, it is used in this volume to refer to the elements of the POSIX Realtime Extension amendment. (This was earlier referred to as POSIX.4 during the standard development process.)

**POSIX.1c**
   Although this term is not used in the normative text of this standard, it is used in this volume to refer to the POSIX Threads Extension amendment. (This was earlier referred to as POSIX.4a during the standard development process.)

**standard developers**
   The individuals and companies in the development organizations responsible for this standard: the IEEE P1003.1 working groups, The Open Group Base working group, advised by the hundreds of individual technical experts who balloted the draft standards within the Austin Group, and the member bodies and technical experts of ISO/IEC JTC 1/SC22/WG15.

**XSI extension**
   The portions of this standard addressing the extension added for support of the Single UNIX Specification.

# *Participants*

IEEE Std 1003.1-2001 was prepared by the Austin Group, sponsored by the Portable Applications Standards Committee of the IEEE Computer Society, The Open Group, and ISO/SC22 WG15.

**The Austin Group**

At the time of approval, the membership of the Austin Group was as follows:

**Andrew Josey**, Chair
**Donald W. Cragun**, Organizational Representative, IEEE PASC
**Nicholas Stoughton**, Organizational Representative, ISO/SC22 WG15
**Mark Brown**, Organizational Representative, The Open Group
**Cathy Hughes**, Technical Editor

**Austin Group Technical Reviewers**

| | | |
|---|---|---|
| Peter Anvin | Michael Gonzalez | Sandra O'Donnell |
| Bouazza Bachar | Joseph M. Gwinn | Frank Prindle |
| Theodore P. Baker | Jon Hitchcock | Curtis Royster Jr. |
| Walter Briscoe | Yvette Ho Sang | Glen Seeds |
| Mark Brown | Cathy Hughes | Keld Jorn Simonsen |
| Dave Butenhof | Lowell G. Johnson | Raja Srinivasan |
| Geoff Clare | Andrew Josey | Nicholas Stoughton |
| Donald W. Cragun | Michael Kavanaugh | Donn S. Terry |
| Lee Damico | David Korn | Fred Tydeman |
| Ulrich Drepper | Marc Aurele La France | Peter Van Der Veen |
| Paul Eggert | Jim Meyering | James Youngman |
| Joanna Farley | Gary Miller | Jim Zepeda |
| Clive D.W. Feather | Finnbarr P. Murphy | Jason Zions |
| Andrew Gollan | Joseph S. Myers | |

**Austin Group Working Group Members**

| | | |
|---|---|---|
| Harold C. Adams | Michael Gonzalez | Sandra O'Donnell |
| Peter Anvin | Karen D. Gordon | Frank Prindle |
| Pierre-Jean Arcos | Joseph M. Gwinn | Francois Riche |
| Jay Ashford | Steven A. Haaser | John D. Riley |
| Bouazza Bachar | Charles E. Hammons | Andrew K. Roach |
| Theodore P. Baker | Chris J. Harding | Helmut Roth |
| Robert Barned | Barry Hedquist | Jaideep Roy |
| Joel Berman | Vincent E. Henley | Curtis Royster Jr. |
| David J. Blackwood | Karl Heubaum | Stephen C. Schwarm |
| Shirley Bockstahler-Brandt | Jon Hitchcock | Glen Seeds |
| James Bottomley | Yvette Ho Sang | Richard Seibel |
| Walter Briscoe | Niklas Holsti | David L. Shroads Jr. |
| Andries Brouwer | Thomas Hosmer | W. Olin Sibert |
| Mark Brown | Cathy Hughes | Keld Jorn Simonsen |
| Eric W. Burger | Jim D. Isaak | Curtis Smith |
| Alan Burns | Lowell G. Johnson | Raja Srinivasan |
| Andries Brouwer | Michael B. Jones | Nicholas Stoughton |
| Dave Butenhof | Andrew Josey | Marc J. Teller |
| Keith Chow | Michael J. Karels | Donn S. Terry |
| Geoff Clare | Michael Kavanaugh | Fred Tydeman |
| Donald W. Cragun | David Korn | Mark-Rene Uchida |
| Lee Damico | Steven Kramer | Scott A. Valcourt |
| Juan Antonio De La Puente | Thomas M. Kurihara | Peter Van Der Veen |
| Ming De Zhou | Marc Aurele La France | Michael W. Vannier |
| Steven J. Dovich | C. Douglass Locke | Eric Vought |
| Richard P. Draves | Nick Maclaren | Frederick N. Webb |
| Ulrich Drepper | Roger J. Martin | Paul A.T. Wolfgang |
| Paul Eggert | Craig H. Meyer | Garrett A. Wollman |
| Philip H. Enslow | Jim Meyering | James Youngman |
| Joanna Farley | Gary Miller | Oren Yuen |
| Clive D.W. Feather | Finnbarr P. Murphy | Janusz Zalewski |
| Pete Forman | Joseph S. Myers | Jim Zepeda |
| Mark Funkenhauser | John Napier | Jason Zions |
| Lois Goldthwaite | Peter E. Obermayer | |
| Andrew Gollan | James T. Oblinger | |

**The Open Group**

When The Open Group approved the Base Specifications, Issue 6 on 12 September 2001, the membership of The Open Group Base Working Group was as follows:

**Andrew Josey**, Chair
**Finnbarr P. Murphy**, Vice-Chair
**Mark Brown**, Austin Group Liaison
**Cathy Hughes**, Technical Editor

**Base Working Group Members**

| | | |
|---|---|---|
| Bouazza Bachar | Joanna Farley | Frank Prindle |
| Mark Brown | Andrew Gollan | Andrew K. Roach |
| Dave Butenhof | Karen D. Gordon | Curtis Royster Jr. |
| Donald W. Cragun | Gary Miller | Nicholas Stoughton |
| Larry Dwyer | Finnbarr P. Murphy | Kenjiro Tsuji |

**IEEE**

When the IEEE Standards Board approved IEEE Std 1003.1-2001 on 6 December 2001, the membership of the committees was as follows:

**Portable Applications Standards Committee (PASC)**

**Lowell G. Johnson**, Chair
**Joseph M. Gwinn**, Vice-Chair
**Jay Ashford**, Functional Chair
**Andrew Josey**, Functional Chair
**Curtis Royster Jr.**, Functional Chair
**Nicholas Stoughton**, Secretary

**Balloting Committee**

The following members of the balloting committee voted on IEEE Std 1003.1-2001. Balloters may have voted for approval, disapproval, or abstention:

| | | |
|---|---|---|
| Harold C. Adams | Steven A. Haaser | Frank Prindle |
| Pierre-Jean Arcos | Charles E. Hammons | Francois Riche |
| Jay Ashford | Chris J. Harding | John D. Riley |
| Theodore P. Baker | Barry Hedquist | Andrew K. Roach |
| Robert Barned | Vincent E. Henley | Helmut Roth |
| David J. Blackwood | Karl Heubaum | Jaideep Roy |
| Shirley Bockstahler-Brandt | Niklas Holsti | Curtis Royster Jr. |
| James Bottomley | Thomas Hosmer | Stephen C. Schwarm |
| Mark Brown | Jim D. Isaak | Richard Seibel |
| Eric W. Burger | Lowell G. Johnson | David L. Shroads Jr. |
| Alan Burns | Michael B. Jones | W. Olin Sibert |
| Dave Butenhof | Andrew Josey | Keld Jorn Simonsen |
| Keith Chow | Michael J. Karels | Nicholas Stoughton |
| Donald W. Cragun | Steven Kramer | Donn S. Terry |
| Juan Antonio De La Puente | Thomas M. Kurihara | Mark-Rene Uchida |
| Ming De Zhou | C. Douglass Locke | Scott A. Valcourt |
| Steven J. Dovich | Roger J. Martin | Michael W. Vannier |
| Richard P. Draves | Craig H. Meyer | Frederick N. Webb |
| Philip H. Enslow | Finnbarr P. Murphy | Paul A.T. Wolfgang |
| Michael Gonzalez | John Napier | Oren Yuen |
| Karen D. Gordon | Peter E. Obermayer | Janusz Zalewski |
| Joseph M. Gwinn | James T. Oblinger | |

The following organizational representative voted on this standard:

**Andrew Josey**, X/Open Company Ltd.

**IEEE-SA Standards Board**

When the IEEE-SA Standards Board approved IEEE Std 1003.1-2001 on 6 December 2001, it had the following membership:

**Donald N. Heirman**, Chair
**James T. Carlo**, Vice-Chair
**Judith Gorman**, Secretary

| | | |
|---|---|---|
| Satish K. Aggarwal | James H. Gurney | James W. Moore |
| Mark D. Bowman | Richard J. Holleman | Robert F. Munzner |
| Gary R. Engmann | Lowell G. Johnson | Ronald C. Petersen |
| Harold E. Epstein | Robert J. Kennelly | Gerald H. Peterson |
| H. Landis Floyd | Joseph L. Koepfinger* | John B. Posey |
| Jay Forster* | Peter H. Lips | Gary S. Robinson |
| Howard M. Frazier | L. Bruce McClung | Akio Tojo |
| Ruben D. Garzon | Daleep C. Mohla | Donald W. Zipse |

Also included are the following non-voting IEEE-SA Standards Board liaisons:

**Alan Cookson**, NIST Representative
**Donald R. Volzka**, TAB Representative
**Yvette Ho Sang**, **Don Messina**, **Savoula Amanatidis**, IEEE Project Editors

_____

\*   Member Emeritus

IEEE Std 1003.1-2001/Cor 1-2002 was prepared by the Austin Group, sponsored by the Portable Applications Standards Committee of the IEEE Computer Society, The Open Group, and ISO/IEC JTC 1/SC22/WG15.

**The Austin Group**

At the time of approval, the membership of the Austin Group was as follows:

**Andrew Josey**, Chair
**Donald W. Cragun**, Organizational Representative, IEEE PASC
**Nicholas Stoughton**, Organizational Representative, ISO/IEC JTC 1/SC22/WG15
**Mark Brown**, Organizational Representative, The Open Group
**Cathy Fox**, Technical Editor

**Austin Group Technical Reviewers**

| | | |
|---|---|---|
| Theodore P. Baker | Mark Funkenhauser | Frank Prindle |
| Julian Blake | Lois Goldthwaite | Kenneth Raeburn |
| Andries Brouwer | Andrew Gollan | Tim Robbins |
| Mark Brown | Michael Gonzalez | Glen Seeds |
| Dave Butenhof | Bruno Haible | Matthew Seitz |
| Geoff Clare | Ben Harris | Keld Jorn Simonsen |
| Donald W. Cragun | Jon Hitchcock | Nicholas Stoughton |
| Ken Dawson | Andreas Jaeger | Alexander Terekhov |
| Ulrich Drepper | Andrew Josey | Donn S. Terry |
| Larry Dwyer | Jonathan Lennox | Mike Wilson |
| Paul Eggert | Nick Maclaren | Garrett A. Wollman |
| Joanna Farley | Jack McCann | Mark Ziegast |
| Clive D.W. Feather | Wilhelm Mueller | |
| Cathy Fox | Joseph S. Myers | |

**Austin Group Working Group Members**

| | | |
|---|---|---|
| Harold C. Adams | Clive D.W. Feather | Wilhelm Mueller |
| Alejandro Alonso | Yaacov Fenster | Finnbarr P. Murphy |
| Jay Ashford | Cathy Fox | Joseph S. Myers |
| Theodore P. Baker | Mark Funkenhauser | Alexey Neyman |
| David J. Blackwood | Lois Goldthwaite | Charles Ngethe |
| Julian Blake | Andrew Gollan | Peter Petrov |
| Mitchell Bonnett | Michael Gonzalez | Frank Prindle |
| Andries Brouwer | Karen D. Gordon | Vikram Punj |
| Mark Brown | Scott Gudgel | Kenneth Raeburn |
| Eric W. Burger | Joseph M. Gwinn | Francois Riche |
| Alan Burns | Steven A. Haaser | Tim Robbins |
| Dave Butenhof | Bruno Haible | Curtis Royster Jr. |
| Keith Chow | Charles E. Hammons | Diane Schleicher |
| Geoff Clare | Bryan Harold | Gil Shultz |
| Luis Cordova | Ben Harris | Stephen C. Schwarm |
| Donald W. Cragun | Barry Hedquist | Glen Seeds |
| Dragan Cvetkovic | Karl Heubaum | Matthew Seitz |
| Lee Damico | Jon Hitchcock | Keld Jorn Simonsen |
| Ken Dawson | Andreas Jaeger | Doug Stevenson |
| Jeroen Dekkers | Andrew Josey | Nicholas Stoughton |
| Juan Antonio De La Puente | Kenneth Lang | Alexander Terekhov |
| Steven J. Dovich | Pi-Cheng Law | Donn S. Terry |
| Ulrich Drepper | Jonathan Lennox | Mike Wilson |
| Dr. Sourav Dutta | Nick Maclaren | Garrett A. Wollman |
| Larry Dwyer | Roger J. Martin | Oren Yuen |
| Paul Eggert | Jack McCann | Mark Ziegast |
| Joanna Farley | George Miao | |

**The Open Group**

When The Open Group approved the Base Specifications, Issue 6, Technical Corrigendum 1 on 7 February 2003, the membership of The Open Group Base Working Group was as follows:

**Andrew Josey**, Chair
**Finnbarr P. Murphy**, Vice-Chair
**Mark Brown**, Austin Group Liaison
**Cathy Fox**, Technical Editor

**Base Working Group Members**

| | | |
|---|---|---|
| Mark Brown | Joanna Farley | Curtis Royster Jr. |
| Dave Butenhof | Andrew Gollan | Nicholas Stoughton |
| Donald W. Cragun | Finnbarr P. Murphy | Kenjiro Tsuji |
| Larry Dwyer | Frank Prindle | |
| Ulrich Drepper | Andrew K. Roach | |

**IEEE**

When the IEEE Standards Board approved IEEE Std 1003.1-2001/Cor 1-2002 on 11 December 2002, the membership of the committees was as follows:

**Portable Applications Standards Committee (PASC)**

**Lowell G. Johnson**, Chair
**Joseph M. Gwinn**, Vice-Chair
**Jay Ashford**, Functional Chair
**Andrew Josey**, Functional Chair
**Curtis Royster Jr.**, Functional Chair
**Nicholas Stoughton**, Secretary

**Balloting Committee**

The following members of the balloting committee voted on IEEE Std 1003.1-2001/Cor 1-2002. Balloters may have voted for approval, disapproval, or abstention:

| | | |
|---|---|---|
| Alejandro Alonso | Michael Gonzalez | Charles Ngethe |
| Jay Ashford | Scott Gudgel | Peter Petrov |
| David J. Blackwood | Charles E. Hammons | Frank Prindle |
| Julian Blake | Bryan Harold | Vikram Punj |
| Mitchell Bonnett | Barry Hedquist | Francois Riche |
| Mark Brown | Karl Heubaum | Curtis Royster Jr. |
| Dave Butenhof | Lowell G. Johnson | Diane Schleicher |
| Keith Chow | Andrew Josey | Stephen C. Schwarm |
| Luis Cordova | Kenneth Lang | Gil Shultz |
| Donald W. Cragun | Pi-Cheng Law | Nicholas Stoughton |
| Steven J. Dovich | George Miao | Donn S. Terry |
| Dr. Sourav Dutta | Roger J. Martin | Oren Yuen |
| Yaacov Fenster | Finnbarr P. Murphy | Juan A. de la Puente |

**IEEE-SA Standards Board**

When the IEEE-SA Standards Board approved IEEE Std 1003.1-2001/Cor 1-2002 on 11 December 2002, the membership was as follows:

**James T. Carlo**, Chair
**James H. Gurney**, Vice-Chair
**Judith Gorman**, Secretary

| | | |
|---|---|---|
| Sid Bennett | Arnold M. Greenspan | Daleep C. Mohla |
| H. Stephen Berger | Raymond Hapeman | William J. Moylan |
| Clyde R. Camp | Donald M. Heirman | Malcolm V. Thaden |
| Richard DeBlasio | Richard H. Hulett | Geoffrey O. Thompson |
| Harold E. Epstein | Lowell G. Johnson | Howard L. Wolfman |
| Julian Forster* | Joseph L. Koepfinger* | Don Wright |
| Howard M. Frazier | Peter H. Lips | |
| Toshio Fukuda | Nader Mehravari | |

Also included are the following non-voting IEEE-SA Standards Board liaisons:

**Alan Cookson**, NIST Representative
**Satish K. Aggarwal**, NRC Representative
**Savoula Amanatidis**, IEEE Standards Managing Editor

_____

\*   Member Emeritus

# *Trademarks*

The following information is given for the convenience of users of this standard and does not constitute endorsement of these products by The Open Group or the IEEE. There may be other products mentioned in the text that might be covered by trademark protection and readers are advised to verify them independently.

1003.1™ is a trademark of the Institute of Electrical and Electronic Engineers, Inc.

AIX® is a registered trademark of IBM Corporation.

AT&T® is a registered trademark of AT&T in the U.S.A. and other countries.

BSD™ is a trademark of the University of California, Berkeley, U.S.A.

Hewlett-Packard®, HP®, and HP-UX® are registered trademarks of Hewlett-Packard Company.

IBM® is a registered trademark of International Business Machines Corporation.

The Open Group and Boundaryless Information Flow are trademarks and UNIX is a registered trademark of The Open Group in the United States and other countries. All other trademarks are the property of their respective owners.

POSIX® is a registered trademark of the Institute of Electrical and Electronic Engineers, Inc.

Sun® and Sun Microsystems® are registered trademarks of Sun Microsystems, Inc.

/usr/group® is a registered trademark of UniForum, the International Network of UNIX System Users.

# *Acknowledgements*

The contributions of the following organizations to the development of IEEE Std 1003.1-2001 are gratefully acknowledged:

- AT&T for permission to reproduce portions of its copyrighted System V Interface Definition (SVID) and material from the UNIX System V Release 2.0 documentation.
- The SC22 WG14 Committees.

This standard was prepared by the Austin Group, a joint working group of the IEEE, The Open Group, and ISO SC22 WG15.

# *Referenced Documents*

**Normative References**

Normative references for this standard are defined in the Base Definitions volume.

**Informative References**

The following documents are referenced in this standard:

1984 /usr/group Standard
> /usr/group Standards Committee, Santa Clara, CA, UniForum 1984.

Almasi and Gottlieb
> George S. Almasi and Allan Gottlieb, *Highly Parallel Computing*, The Benjamin/Cummings Publishing Company, Inc., 1989, ISBN: 0-8053-0177-1.

ANSI C
> American National Standard for Information Systems: Standard X3.159-1989, Programming Language C.

ANSI X3.226-1994
> American National Standard for Information Systems: Standard X3.226-1994, Programming Language Common LISP.

Brawer
> Steven Brawer, *Introduction to Parallel Programming*, Academic Press, 1989, ISBN: 0-12-128470-0.

DeRemer and Pennello Article
> DeRemer, Frank and Pennello, Thomas J., *Efficient Computation of LALR(1) Look-Ahead Sets*, SigPlan Notices, Volume 15, No. 8, August 1979.

Draft ANSI X3J11.1
> IEEE Floating Point draft report of ANSI X3J11.1 (NCEG).

FIPS 151-1
> Federal Information Procurement Standard (FIPS) 151-1. Portable Operating System Interface (POSIX)—Part 1: System Application Program Interface (API) [C Language].

FIPS 151-2
> Federal Information Procurement Standards (FIPS) 151-2, Portable Operating System Interface (POSIX)— Part 1: System Application Program Interface (API) [C Language].

HP-UX Manual
> Hewlett-Packard HP-UX Release 9.0 Reference Manual, Third Edition, August 1992.

IEC 60559: 1989
> IEC 60559: 1989, Binary Floating-Point Arithmetic for Microprocessor Systems (previously designated IEC 559: 1989).

IEEE Std 754-1985
> IEEE Std 754-1985, IEEE Standard for Binary Floating-Point Arithmetic.

IEEE Std 854-1987
> IEEE Std 854-1987, IEEE Standard for Radix-Independent Floating-Point Arithmetic.

IEEE Std 1003.9-1992
   IEEE Std 1003.9-1992, IEEE Standard for Information Technology — POSIX FORTRAN 77 Language Interfaces — Part 1: Binding for System Application Program Interface API.

IETF RFC 791
   Internet Protocol, Version 4 (IPv4), September 1981.

IETF RFC 819
   The Domain Naming Convention for Internet User Applications, Z. Su, J. Postel, August 1982.

IETF RFC 822
   Standard for the Format of ARPA Internet Text Messages, D.H. Crocker, August 1982.

IETF RFC 919
   Broadcasting Internet Datagrams, J. Mogul, October 1984.

IETF RFC 920
   Domain Requirements, J. Postel, J. Reynolds, October 1984.

IETF RFC 921
   Domain Name System Implementation Schedule, J. Postel, October 1984.

IETF RFC 922
   Broadcasting Internet Datagrams in the Presence of Subnets, J. Mogul, October 1984.

IETF RFC 1034
   Domain Names — Concepts and Facilities, P. Mockapetris, November 1987.

IETF RFC 1035
   Domain Names — Implementation and Specification, P. Mockapetris, November 1987.

IETF RFC 1123
   Requirements for Internet Hosts — Application and Support, R. Braden, October 1989.

IETF RFC 1886
   DNS Extensions to Support Internet Protocol, Version 6 (IPv6), C. Huitema, S. Thomson, December 1995.

IETF RFC 2045
   Multipurpose Internet Mail Extensions (MIME), Part 1: Format of Internet Message Bodies, N. Freed, N. Borenstein, November 1996.

IETF RFC 2181
   Clarifications to the DNS Specification, R. Elz, R. Bush, July 1997.

IETF RFC 2373
   Internet Protocol, Version 6 (IPv6) Addressing Architecture, S. Deering, R. Hinden, July 1998.

IETF RFC 2460
   Internet Protocol, Version 6 (IPv6), S. Deering, R. Hinden, December 1998.

Internationalisation Guide
   Guide, July 1993, Internationalisation Guide, Version 2 (ISBN: 1-859120-02-4, G304), published by The Open Group.

ISO C (1990)
   ISO/IEC 9899:1990, Programming Languages — C, including Amendment 1:1995 (E), C Integrity (Multibyte Support Extensions (MSE) for ISO C).

ISO 2375: 1985
> ISO 2375: 1985, Data Processing — Procedure for Registration of Escape Sequences.

ISO 8652: 1987
> ISO 8652: 1987, Programming Languages — Ada (technically identical to ANSI standard 1815A-1983).

ISO/IEC 1539: 1990
> ISO/IEC 1539: 1990, Information Technology — Programming Languages — Fortran (technically identical to the ANSI X3.9-1978 standard [FORTRAN 77]).

ISO/IEC 4873: 1991
> ISO/IEC 4873: 1991, Information Technology — ISO 8-bit Code for Information Interchange — Structure and Rules for Implementation.

ISO/IEC 6429: 1992
> ISO/IEC 6429: 1992, Information Technology — Control Functions for Coded Character Sets.

ISO/IEC 6937: 1994
> ISO/IEC 6937: 1994, Information Technology — Coded Character Set for Text Communication — Latin Alphabet.

ISO/IEC 8802-3: 1996
> ISO/IEC 8802-3: 1996, Information Technology — Telecommunications and Information Exchange Between Systems — Local and Metropolitan Area Networks — Specific Requirements — Part 3: Carrier Sense Multiple Access with Collision Detection (CSMA/CD) Access Method and Physical Layer Specifications.

ISO/IEC 8859
> ISO/IEC 8859, Information Technology — 8-Bit Single-Byte Coded Graphic Character Sets:
>
> Part 1: Latin Alphabet No. 1
> Part 2: Latin Alphabet No. 2
> Part 3: Latin Alphabet No. 3
> Part 4: Latin Alphabet No. 4
> Part 5: Latin/Cyrillic Alphabet
> Part 6: Latin/Arabic Alphabet
> Part 7: Latin/Greek Alphabet
> Part 8: Latin/Hebrew Alphabet
> Part 9: Latin Alphabet No. 5
> Part 10: Latin Alphabet No. 6
> Part 13: Latin Alphabet No. 7
> Part 14: Latin Alphabet No. 8
> Part 15: Latin Alphabet No. 9

ISO POSIX-1: 1996
> ISO/IEC 9945-1: 1996, Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language] (identical to ANSI/IEEE Std 1003.1-1996). Incorporating ANSI/IEEE Stds 1003.1-1990, 1003.1b-1993, 1003.1c-1995, and 1003.1i-1995.

ISO POSIX-2: 1993
> ISO/IEC 9945-2: 1993, Information Technology — Portable Operating System Interface (POSIX) — Part 2: Shell and Utilities (identical to ANSI/IEEE Std 1003.2-1992, as amended by ANSI/IEEE Std 1003.2a-1992).

Issue 1
:   X/Open Portability Guide, July 1985 (ISBN: 0-444-87839-4).

Issue 2
:   X/Open Portability Guide, January 1987:

    - Volume 1: XVS Commands and Utilities (ISBN: 0-444-70174-5)

    - Volume 2: XVS System Calls and Libraries (ISBN: 0-444-70175-3)

Issue 3
:   X/Open Specification, 1988, 1989, February 1992:

    - Commands and Utilities, Issue 3 (ISBN: 1-872630-36-7, C211); this specification was formerly X/Open Portability Guide, Issue 3, Volume 1, January 1989, XSI Commands and Utilities (ISBN: 0-13-685835-X, XO/XPG/89/002)

    - System Interfaces and Headers, Issue 3 (ISBN: 1-872630-37-5, C212); this specification was formerly X/Open Portability Guide, Issue 3, Volume 2, January 1989, XSI System Interface and Headers (ISBN: 0-13-685843-0, XO/XPG/89/003)

    - Curses Interface, Issue 3, contained in Supplementary Definitions, Issue 3 (ISBN: 1-872630-38-3, C213), Chapters 9 to 14 inclusive; this specification was formerly X/Open Portability Guide, Issue 3, Volume 3, January 1989, XSI Supplementary Definitions (ISBN: 0-13-685850-3, XO/XPG/89/004)

    - Headers Interface, Issue 3, contained in Supplementary Definitions, Issue 3 (ISBN: 1-872630-38-3, C213), Chapter 19, Cpio and Tar Headers; this specification was formerly X/Open Portability Guide Issue 3, Volume 3, January 1989, XSI Supplementary Definitions (ISBN: 0-13-685850-3, XO/XPG/89/004)

Issue 4
:   CAE Specification, July 1992, published by The Open Group:

    - System Interface Definitions (XBD), Issue 4 (ISBN: 1-872630-46-4, C204)

    - Commands and Utilities (XCU), Issue 4 (ISBN: 1-872630-48-0, C203)

    - System Interfaces and Headers (XSH), Issue 4 (ISBN: 1-872630-47-2, C202)

Issue 4, Version 2
:   CAE Specification, August 1994, published by The Open Group:

    - System Interface Definitions (XBD), Issue 4, Version 2 (ISBN: 1-85912-036-9, C434)

    - Commands and Utilities (XCU), Issue 4, Version 2 (ISBN: 1-85912-034-2, C436)

    - System Interfaces and Headers (XSH), Issue 4, Version 2 (ISBN: 1-85912-037-7, C435)

Issue 5
:   Technical Standard, February 1997, published by The Open Group:

    - System Interface Definitions (XBD), Issue 5 (ISBN: 1-85912-186-1, C605)

    - Commands and Utilities (XCU), Issue 5 (ISBN: 1-85912-191-8, C604)

    - System Interfaces and Headers (XSH), Issue 5 (ISBN: 1-85912-181-0, C606)

Knuth Article
:   Knuth, Donald E., *On the Translation of Languages from Left to Right*, Information and Control, Volume 8, No. 6, October 1965.

KornShell

> Bolsky, Morris I. and Korn, David G., *The New KornShell Command and Programming Language*, March 1995, Prentice Hall.

MSE Working Draft

> Working draft of ISO/IEC 9899: 1990/Add3: Draft, Addendum 3 — Multibyte Support Extensions (MSE) as documented in the ISO Working Paper SC22/WG14/N205 dated 31 March 1992.

POSIX.0: 1995

> IEEE Std 1003.0-1995, IEEE Guide to the POSIX Open System Environment (OSE) (identical to ISO/IEC TR 14252).

POSIX.1: 1988

> IEEE Std 1003.1-1988, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language].

POSIX.1: 1990

> IEEE Std 1003.1-1990, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language].

POSIX.1a

> P1003.1a, Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — (C Language) Amendment.

POSIX.1d: 1999

> IEEE Std 1003.1d-1999, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 4: Additional Realtime Extensions [C Language].

POSIX.1g: 2000

> IEEE Std 1003.1g-2000, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 6: Protocol-Independent Interfaces (PII).

POSIX.1j: 2000

> IEEE Std 1003.1j-2000, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 5: Advanced Realtime Extensions [C Language].

POSIX.1q: 2000

> IEEE Std 1003.1q-2000, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 7: Tracing [C Language].

POSIX.2b

> P1003.2b, Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 2: Shell and Utilities — Amendment.

POSIX.2d:-1994

> IEEE Std 1003.2d-1994, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 2: Shell and Utilities — Amendment 1: Batch Environment.

POSIX.13:-1998

IEEE Std 1003.13:1998, IEEE Standard for Information Technology — Standardized Application Environment Profile (AEP) — POSIX Realtime Application Support.

Sarwate Article

Sarwate, Dilip V., *Computation of Cyclic Redundancy Checks via Table Lookup*, Communications of the ACM, Volume 30, No. 8, August 1988.

Sprunt, Sha, and Lehoczky

Sprunt, B., Sha, L., and Lehoczky, J.P., *Aperiodic Task Scheduling for Hard Real-Time Systems*, The Journal of Real-Time Systems, Volume 1, 1989, Pages 27-60.

SVID, Issue 1

American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue 1; Morristown, NJ, UNIX Press, 1985.

SVID, Issue 2

American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue 2; Morristown, NJ, UNIX Press, 1986.

SVID, Issue 3

American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue 3; Morristown, NJ, UNIX Press, 1989.

The AWK Programming Language

Aho, Alfred V., Kernighan, Brian W., and Weinberger, Peter J., *The AWK Programming Language*, Reading, MA, Addison-Wesley 1988.

UNIX Programmer's Manual

American Telephone and Telegraph Company, *UNIX Time-Sharing System: UNIX Programmer's Manual*, 7th Edition, Murray Hill, NJ, Bell Telephone Laboratories, January 1979.

XNS, Issue 4

CAE Specification, August 1994, Networking Services, Issue 4 (ISBN: 1-85912-049-0, C438), published by The Open Group.

XNS, Issue 5

CAE Specification, February 1997, Networking Services, Issue 5 (ISBN: 1-85912-165-9, C523), published by The Open Group.

XNS, Issue 5.2

Technical Standard, January 2000, Networking Services (XNS), Issue 5.2 (ISBN: 1-85912-241-8, C808), published by The Open Group.

X/Open Curses, Issue 4, Version 2

CAE Specification, May 1996, X/Open Curses, Issue 4, Version 2 (ISBN: 1-85912-171-3, C610), published by The Open Group.

Yacc

*Yacc: Yet Another Compiler Compiler*, Stephen C. Johnson, 1978.

**Source Documents**

Parts of the following documents were used to create the base documents for this standard:

AIX 3.2 Manual
> AIX Version 3.2 For RISC System/6000, Technical Reference: Base Operating System and Extensions, 1990, 1992 (Part No. SC23-2382-00).

OSF/1
> OSF/1 Programmer's Reference, Release 1.2 (ISBN: 0-13-020579-6).

OSF AES
> Application Environment Specification (AES) Operating System Programming Interfaces Volume, Revision A (ISBN: 0-13-043522-8).

System V Release 2.0

> — UNIX System V Release 2.0 Programmer's Reference Manual (April 1984 - Issue 2).

> — UNIX System V Release 2.0 Programming Guide (April 1984 - Issue 2).

System V Release 4.2
> Operating System API Reference, UNIX SVR4.2 (1992) (ISBN: 0-13-017658-3).

1 # *Rationale (Informative)*

2 **Part A:**

3 **Base Definitions**

4 *The Open Group*
5 *The Institute of Electrical and Electronics Engineers, Inc.*

# Rationale for Base Definitions

## A.1 Introduction

### A.1.1 Scope

IEEE Std 1003.1-2001 is one of a family of standards known as POSIX. The family of standards extends to many topics; IEEE Std 1003.1-2001 is known as POSIX.1 and consists of both operating system interfaces and shell and utilities. IEEE Std 1003.1-2001 is technically identical to The Open Group Base Specifications, Issue 6, which comprise the core volumes of the Single UNIX Specification, Version 3.

**Scope of IEEE Std 1003.1-2001**

The (paraphrased) goals of this development were to produce a single common revision to the overlapping POSIX.1 and POSIX.2 standards, and the Single UNIX Specification, Version 2. As such, the scope of the revision includes the scopes of the original documents merged.

Since the revision includes merging the Base volumes of the Single UNIX Specification, many features that were previously not ''adopted'' into earlier revisions of POSIX.1 and POSIX.2 are now included in IEEE Std 1003.1-2001. In most cases, these additions are part of the XSI extension; in other cases the standard developers decided that now was the time to migrate these to the base standard.

The Single UNIX Specification programming environment provides a broad-based functional set of interfaces to support the porting of existing UNIX applications and the development of new applications. The environment also supports a rich set of tools for application development.

The majority of the obsolescent material from the existing POSIX.1 and POSIX.2 standards, and material marked LEGACY from The Open Group's Base specifications, has been removed in this revision. New members of the Legacy Option Group have been added, reflecting the advance in understanding of what is required.

The following IEEE standards have been added to the base documents in this revision:

- IEEE Std 1003.1d-1999
- IEEE Std 1003.1j-2000
- IEEE Std 1003.1q-2000
- IEEE P1003.1a draft standard
- IEEE Std 1003.2d-1994
- IEEE P1003.2b draft standard
- Selected parts of IEEE Std 1003.1g-2000

Only selected parts of IEEE Std 1003.1g-2000 were included. This was because there is much duplication between the XNS, Issue 5.2 specification (another base document) and the material from IEEE Std 1003.1g-2000, the former document being aligned with the latest networking specifications for IPv6. Only the following sections of IEEE Std 1003.1g-2000 were considered for inclusion:

44    • General terms related to sockets (Section 2.2.2)

45    • Socket concepts (Sections 5.1 through 5.3 inclusive)

46    • The *pselect*( ) function (Sections 6.2.2.1 and 6.2.3)

47    • The **<sys/select.h>** header (Section 6.2)

48    The following were requirements on IEEE Std 1003.1-2001:

49    • Backward-compatibility

50    It was agreed that there should be no breakage of functionality in the existing base
51    documents. This requirement was tempered by changes introduced due to interpretations
52    and corrigenda on the base documents, and any changes introduced in the
53    ISO/IEC 9899:1999 standard (C Language).

54    • Architecture and n-bit neutral

55    The common standard should not make any implicit assumptions about the system
56    architecture or size of data types; for example, previously some 32-bit implicit assumptions
57    had crept into the standards.

58    • Extensibility

59    It should be possible to extend the common standard without breaking backwards-
60    compatibility. For example, the name space should be reserved and structured to avoid
61    duplication of names between the standard and extensions to it.

62    **POSIX.1 and the ISO C Standard**

63    Previous revisions of POSIX.1 built upon the ISO C standard by reference only. This revision
64    takes a different approach.

65    The standard developers believed it essential for a programmer to have a single complete
66    reference place, but recognized that deference to the formal standard had to be addressed for the
67    duplicate interface definitions between the ISO C standard and the Single UNIX Specification.

68    It was agreed that where an interface has a version in the ISO C standard, the DESCRIPTION
69    section should describe the relationship to the ISO C standard and markings should be added as
70    appropriate to show where the ISO C standard has been extended in the text.

71    A block of text was added to the start of each affected reference page stating whether the page is
72    aligned with the ISO C standard or extended. Each page was parsed for additions beyond the
73    ISO C standard (that is, including both POSIX and UNIX extensions), and these extensions are
74    marked as CX extensions (for C Extensions).

75    **FIPS Requirements**

76    The Federal Information Processing Standards (FIPS) are a series of U.S. government
77    procurement standards managed and maintained on behalf of the U.S. Department of
78    Commerce by the National Institute of Standards and Technology (NIST).

79    The following restrictions have been made in this version of IEEE Std 1003.1 in order to align
80    with FIPS 151-2 requirements:

81    • The implementation supports _POSIX_CHOWN_RESTRICTED.

82    • The limit {NGROUPS_MAX} is now greater than or equal to 8.

83    • The implementation supports the setting of the group ID of a file (when it is created) to that
84    of the parent directory.

85       • The implementation supports _POSIX_SAVED_IDS.

86       • The implementation supports _POSIX_VDISABLE.

87       • The implementation supports _POSIX_JOB_CONTROL.

88       • The implementation supports _POSIX_NO_TRUNC.

89       • The *read*( ) function returns the number of bytes read when interrupted by a signal and does
90          not return −1.

91       • The *write*( ) function returns the number of bytes written when interrupted by a signal and
92          does not return −1.

93       • In the environment for the login shell, the environment variables *LOGNAME* and *HOME* are
94          defined and have the properties described in IEEE Std 1003.1-2001.

95       • The value of {CHILD_MAX} is now greater than or equal to 25.

96       • The value of {OPEN_MAX} is now greater than or equal to 20.

97       • The implementation supports the functionality associated with the symbols CS7, CS8,
98          CSTOPB, PARODD, and PARENB defined in <**termios.h**>.

## A.1.2   Conformance

100     See Section A.2 (on page 9).

## A.1.3   Normative References

102     There is no additional rationale provided for this section.

## A.1.4   Terminology

104     The meanings specified in IEEE Std 1003.1-2001 for the words *shall*, *should*, and *may* are
105     mandated by ISO/IEC directives.

106     In the Rationale (Informative) volume of IEEE Std 1003.1-2001, the words *shall*, *should*, and *may*
107     are sometimes used to illustrate similar usages in IEEE Std 1003.1-2001. However, the rationale
108     itself does not specify anything regarding implementations or applications.

**conformance document**
110        As a practical matter, the conformance document is effectively part of the system
111        documentation. Conformance documents are distinguished by IEEE Std 1003.1-2001 so that
112        they can be referred to distinctly.

**implementation-defined**
114        This definition is analogous to that of the ISO C standard and, together with ''undefined''
115        and ''unspecified'', provides a range of specification of freedom allowed to the interface
116        implementor.

**may**
118        The use of *may* has been limited as much as possible, due both to confusion stemming from
119        its ordinary English meaning and to objections regarding the desirability of having as few
120        options as possible and those as clearly specified as possible.

121        The usage of *can* and *may* were selected to contrast optional application behavior (can)
122        against optional implementation behavior (may).

**shall**

Declarative sentences are sometimes used in IEEE Std 1003.1-2001 as if they included the word *shall*, and facilities thus specified are no less required. For example, the two statements:

1. The *foo*( ) function shall return zero.

2. The *foo*( ) function returns zero.

are meant to be exactly equivalent.

**should**

In IEEE Std 1003.1-2001, the word *should* does not usually apply to the implementation, but rather to the application. Thus, the important words regarding implementations are *shall*, which indicates requirements, and *may*, which indicates options.

**obsolescent**

The term ''obsolescent'' means ''do not use this feature in new applications''. The obsolescence concept is not an ideal solution, but was used as a method of increasing consensus: many more objections would be heard from the user community if some of these historical features were suddenly withdrawn without the grace period obsolescence implies. The phrase ''may be considered for withdrawal in future revisions'' implies that the result of that consideration might in fact keep those features indefinitely if the predominance of applications do not migrate away from them quickly.

**legacy**

The term ''legacy'' was added for compatibility with the Single UNIX Specification. It means ''this feature is historic and optional; do not use this feature in new applications. There are alternative interfaces that are more suitable.''. It is used exclusively for XSI extensions, and includes facilities that were mandatory in previous versions of the base document but are optional in this revision. This is a way to ''sunset'' the usage of certain functions. Application writers should not rely on the existence of these facilities in new applications, but should follow the migration path detailed in the APPLICATION USAGE sections of the relevant pages.

The terms ''legacy'' and ''obsolescent'' are different: a feature marked LEGACY is not recommended for new work and need not be present on an implementation (if the XSI Legacy Option Group is not supported). A feature noted as obsolescent is supported by all implementations, but may be removed in a future revision; new applications should not use these features.

**system documentation**

The system documentation should normally describe the whole of the implementation, including any extensions provided by the implementation. Such documents normally contain information at least as detailed as the specifications in IEEE Std 1003.1-2001. Few requirements are made on the system documentation, but the term is needed to avoid a dangling pointer where the conformance document is permitted to point to the system documentation.

**undefined**

See *implementation-defined*.

**unspecified**

See *implementation-defined*.

The definitions for ''unspecified'' and ''undefined'' appear nearly identical at first examination, but are not. The term ''unspecified'' means that a conforming application may deal with the unspecified behavior, and it should not care what the outcome is. The term

170
171
172
173
''undefined'' says that a conforming application should not do it because no definition is provided for what it does (and implicitly it would care what the outcome was if it tried it). It is important to remember, however, that if the syntax permits the statement at all, it must have some outcome in a real implementation.

174
175
176
177
Thus, the terms ''undefined'' and ''unspecified'' apply to the way the application should think about the feature. In terms of the implementation, it is always ''defined''—there is always some result, even if it is an error. The implementation is free to choose the behavior it prefers.

178
179
This also implies that an implementation, or another standard, could specify or define the result in a useful fashion. The terms apply to IEEE Std 1003.1-2001 specifically.

180
181
182
183
184
185
186
187
188
189
190
191
The term ''implementation-defined'' implies requirements for documentation that are not required for ''undefined'' (or ''unspecified''). Where there is no need for a conforming program to know the definition, the term ''undefined'' is used, even though ''implementation-defined'' could also have been used in this context. There could be a fourth term, specifying ''this standard does not say what this does; it is acceptable to define it in an implementation, but it does not need to be documented'', and undefined would then be used very rarely for the few things for which any definition is not useful. In particular, implementation-defined is used where it is believed that certain classes of application will need to know such details to determine whether the application can be successfully ported to the implementation. Such applications are not always strictly portable, but nevertheless are common and useful; often the requirements met by the application cannot be met without dealing with the issues implied by ''implementation-defined''.

192
193
194
195
196
197
198
In many places IEEE Std 1003.1-2001 is silent about the behavior of some possible construct. For example, a variable may be defined for a specified range of values and behaviors are described for those values; nothing is said about what happens if the variable has any other value. That kind of silence can imply an error in the standard, but it may also imply that the standard was intentionally silent and that any behavior is permitted. There is a natural tendency to infer that if the standard is silent, a behavior is prohibited. That is not the intent. Silence is intended to be equivalent to the term ''unspecified''.

199
200
The term ''application'' is not defined in IEEE Std 1003.1-2001; it is assumed to be a part of general computer science terminology.

201
202
Three terms used within IEEE Std 1003.1-2001 overlap in meaning: ''macro'', ''symbolic name'', and ''symbolic constant''.

203
**macro**

204
205
206
This usually describes a C preprocessor symbol, the result of the **#define** operator, with or without an argument. It may also be used to describe similar mechanisms in editors and text processors.

207
**symbolic name**

208
209
210
This can also refer to a C preprocessor symbol (without arguments), but is also used to refer to the names for characters in character sets. In addition, it is sometimes used to refer to host names and even filenames.

211
**symbolic constant**

212
This also refers to a C preprocessor symbol (also without arguments).

213
214
In most cases, the difference in semantic content is negligible to nonexistent. Readers should not attempt to read any meaning into the various usages of these terms.

215 **A.1.5    Portability**

216   To aid the identification of options within IEEE Std 1003.1-2001, a notation consisting of margin
217   codes and shading is used. This is based on the notation used in previous revisions of The Open
218   Group's Base specifications.

219   The benefit of this approach is a reduction in the number of *if* statements within the running
220   text, that makes the text easier to read, and also an identification to the programmer that they
221   need to ensure that their target platforms support the underlying options. For example, if
222   functionality is marked with THR in the margin, it will be available on all systems supporting
223   the Threads option, but may not be available on some others.

224 *A.1.5.1    Codes*

225   This section includes codes for options defined in the Base Definitions volume of
226   IEEE Std 1003.1-2001, Section 2.1.6, Options, and the following additional codes for other
227   purposes:

228   CX        This margin code is used to denote extensions beyond the ISO C standard. For
229             interfaces that are duplicated between IEEE Std 1003.1-2001 and the ISO C standard, a
230             CX introduction block describes the nature of the duplication, with any extensions
231             appropriately CX marked and shaded.

232             Where an interface is added to an ISO C standard header, within the header the
233             interface has an appropriate margin marker and shading (for example, CX, XSI, TSF,
234             and so on) and the same marking appears on the reference page in the SYNOPSIS
235             section. This enables a programmer to easily identify that the interface is extending an
236             ISO C standard header.

237   MX        This margin code is used to denote IEC 60559:1989 standard floating-point extensions.

238   OB        This margin code is used to denote obsolescent behavior and thus flag a possible future
239             applications portability warning.

240   OH        The Single UNIX Specification has historically tried to reduce the number of headers an
241             application has had to include when using a particular interface. Sometimes this was
242             fewer than the base standard, and hence a notation is used to flag which headers are
243             optional if you are using a system supporting the XSI extension.

244   XSI       This code is used to denote interfaces and facilities within interfaces only required on
245             systems supporting the XSI extension. This is introduced to support the Single UNIX
246             Specification.

247   XSR       This code is used to denote interfaces and facilities within interfaces only required on
248             systems supporting STREAMS. This is introduced to support the Single UNIX
249             Specification, although it is defined in a way so that it can stand alone from the XSI
250             notation.

251 *A.1.5.2    Margin Code Notation*

252   Since some features may depend on one or more options, or require more than one option, a
253   notation is used. Where a feature requires support of a single option, a single margin code will
254   occur in the margin. If it depends on two options and both are required, then the codes will
255   appear with a <space> separator. If either of two options are required, then a logical OR is
256   denoted using the ' | ' symbol. If more than two codes are used, a special notation is used.

## A.2 Conformance

The terms ''profile'' and ''profiling'' are used throughout this section.

A profile of a standard or standards is a codified set of option selections, such that by being conformant to a profile, particular classes of users are specifically supported.

These conformance definitions are descended from those in the ISO POSIX-1: 1996 standard, but with changes for the following:

- The addition of profiling options, allowing larger profiles of options such as the XSI extension used by the Single UNIX Specification. In effect, it has profiled itself (that is, created a self-profile).

- The addition of a definition of subprofiling considerations, to allow smaller profiles of options.

- The addition of a hierarchy of super-options for XSI; these were formerly known as ''Feature Groups'' in the System Interfaces and Headers, Issue 5 specification.

- Options from the ISO POSIX-2: 1993 standard are also now included, as IEEE Std 1003.1-2001 merges the functionality from it.

### A.2.1 Implementation Conformance

These definitions allow application developers to know what to depend on in an implementation.

There is no definition of a ''strictly conforming implementation''; that would be an implementation that provides *only* those facilities specified by POSIX.1 with no extensions whatsoever. This is because no actual operating system implementation can exist without system administration and initialization facilities that are beyond the scope of POSIX.1.

#### A.2.1.1 Requirements

The word ''support'' is used in certain instances, rather than ''provide'', in order to allow an implementation that has no resident software development facilities, but that supports the execution of a *Strictly Conforming POSIX.1 Application*, to be a *conforming implementation*.

#### A.2.1.2 Documentation

The conformance documentation is required to use the same numbering scheme as POSIX.1 for purposes of cross-referencing. All options that an implementation chooses are reflected in **<limits.h>** and **<unistd.h>**.

Note that the use of ''may'' in terms of where conformance documents record where implementations may vary, implies that it is not required to describe those features identified as undefined or unspecified.

Other aspects of systems must be evaluated by purchasers for suitability. Many systems incorporate buffering facilities, maintaining updated data in volatile storage and transferring such updates to non-volatile storage asynchronously. Various exception conditions, such as a power failure or a system crash, can cause this data to be lost. The data may be associated with a file that is still open, with one that has been closed, with a directory, or with any other internal system data structures associated with permanent storage. This data can be lost, in whole or part, so that only careful inspection of file contents could determine that an update did not occur.

298   Also, interrelated file activities, where multiple files and/or directories are updated, or where
299   space is allocated or released in the file system structures, can leave inconsistencies in the
300   relationship between data in the various files and directories, or in the file system itself. Such
301   inconsistencies can break applications that expect updates to occur in a specific sequence, so that
302   updates in one place correspond with related updates in another place.

303   For example, if a user creates a file, places information in the file, and then records this action in
304   another file, a system or power failure at this point followed by restart may result in a state in
305   which the record of the action is permanently recorded, but the file created (or some of its
306   information) has been lost. The consequences of this to the user may be undesirable. For a user
307   on such a system, the only safe action may be to require the system administrator to have a
308   policy that requires, after any system or power failure, that the entire file system must be
309   restored from the most recent backup copy (causing all intervening work to be lost).

310   The characteristics of each implementation will vary in this respect and may or may not meet the
311   requirements of a given application or user. Enforcement of such requirements is beyond the
312   scope of POSIX.1. It is up to the purchaser to determine what facilities are provided in an
313   implementation that affect the exposure to possible data or sequence loss, and also what
314   underlying implementation techniques and/or facilities are provided that reduce or limit such
315   loss or its consequences.

316   *A.2.1.3   POSIX Conformance*

317   This really means conformance to the base standard; however, since this revision includes the
318   core material of the Single UNIX Specification, the standard developers decided that it was
319   appropriate to segment the conformance requirements into two, the former for the base
320   standard, and the latter for the Single UNIX Specification.

321   Within POSIX.1 there are some symbolic constants that, if defined, indicate that a certain option
322   is enabled. Other symbolic constants exist in POSIX.1 for other reasons.

323   As part of the revision some alignment has occurred of the options with the FIPS 151-2 profile on
324   the POSIX.1-1990 standard. The following options from the POSIX.1-1990 standard are now
325   mandatory:

326   • _POSIX_JOB_CONTROL

327   • _POSIX_SAVED_IDS

328   • _POSIX_VDISABLE

329   A POSIX-conformant system may support the XSI extensions of the Single UNIX Specification.
330   This was intentional since the standard developers intend them to be upwards-compatible, so
331   that a system conforming to the Single UNIX Specification can also conform to the base standard
332   at the same time.

333   *A.2.1.4   XSI Conformance*

334   This section is added since the revision merges in the base volumes of the Single UNIX
335   Specification.

336   XSI conformance can be thought of as a profile, selecting certain options from
337   IEEE Std 1003.1-2001.

338     *A.2.1.5   Option Groups*

339     The concept of ''Option Groups'' is introduced to IEEE Std 1003.1-2001 to allow collections of |
340     related functions or options to be grouped together. This has been used as follows: the ''XSI |
341     Option Groups'' have been created to allow super-options, collections of underlying options and |
342     related functions, to be collectively supported by XSI-conforming systems. These reflect the |
343     ''Feature Groups'' from the System Interfaces and Headers, Issue 5 specification. |

344     The standard developers considered the matter of subprofiling and decided it was better to
345     include an enabling mechanism rather than detailed normative requirements. A set of
346     subprofiling options was developed and included later in this volume of IEEE Std 1003.1-2001 as
347     an informative illustration.

348     **Subprofiling Considerations**

349     The goal of not simultaneously fixing maximums and minimums was to allow implementations
350     of the base standard or standards to support multiple profiles without conflict.

351     The following summarizes the rules for the limit types:

| Limit<br>Type | Fixed<br>Value | Minimum Acceptable<br>Value | Maximum Acceptable<br>Value |
|---|---|---|---|
| Standard<br>Profile | $Xs$<br>$Xp == Xs$<br>(No change) | $Ys$<br>$Yp >= Ys$<br>(May increase the limit) | $Zs$<br>$Zp <= Zs$<br>(May decrease the limit) |

357     The intent is that ranges specified by limits in profiles be entirely contained within the
358     corresponding ranges of the base standard or standards being profiled, and that the unlimited
359     end of a range in a base standard must remain unlimited in any profile of that standard.

360     Thus, the fixed _POSIX_* limits are constants and must not be changed by a profile. The variable
361     counterparts (typically without the leading _POSIX_) can be changed but still remain
362     semantically the same; that is, they still allow implementation values to vary as long as they
363     meet the requirements for that value (be it a minimum or maximum).

364     Where a profile does not provide a feature upon which a limit is based, the limit is not relevant.
365     Applications written to that profile should be written to operate independently of the value of
366     the limit.

367     An example which has previously allowed implementations to support both the base standard
368     and two other profiles in a compatible manner follows:

```
369     Base standard (POSIX.1-1996): _POSIX_CHILD_MAX 6
370     Base standard: CHILD_MAX   minimum maximum _POSIX_CHILD_MAX
371         FIPS profile/SUSv2  CHILD_MAX   25 (minimum maximum)
```

372     Another example:

```
373     Base standard (POSIX.1-1996): _POSIX_NGROUPS_MAX 0
374     Base standard: NGROUPS_MAX   minimum maximum _POSIX_NGROUP_MAX
375         FIPS profile/SUSv2  NGROUPS_MAX   8
```

376     A profile may lower a minimum maximum below the equivalent _POSIX value:

```
377     Base standard: _POSIX_foo_MAX   Z
378     Base standard: foo_MAX   _POSIX_foo_MAX
379         profile standard : foo_MAX   X  (X can be less than, equal to,
380                                           or greater than _POSIX_foo_MAX)
```

381 In this case an implementation conforming to the profile may not conform to the base standard,
382 but an implementation to the base standard will conform to the profile.

383 *A.2.1.6 Options*

384 The final subsections within *Implementation Conformance* list the core options within
385 IEEE Std 1003.1-2001. This includes both options for the System Interfaces volume of
386 IEEE Std 1003.1-2001 and the Shell and Utilities volume of IEEE Std 1003.1-2001.

387 **A.2.2    Application Conformance**

388 These definitions guide users or adaptors of applications in determining on which
389 implementations an application will run and how much adaptation would be required to make
390 it run on others. These definitions are modeled after related ones in the ISO C standard.

391 POSIX.1 occasionally uses the expressions ''portable application'' or ''conforming application''.
392 As they are used, these are synonyms for any of these terms. The differences between the classes
393 of application conformance relate to the requirements for other standards, the options supported
394 (such as the XSI extension) or, in the case of the Conforming POSIX.1 Application Using
395 Extensions, to implementation extensions. When one of the less explicit expressions is used, it
396 should be apparent from the context of the discussion which of the more explicit names is
397 appropriate

398 *A.2.2.1 Strictly Conforming POSIX Application*

399 This definition is analogous to that of an ISO C standard ''conforming program''.

400 The major difference between a Strictly Conforming POSIX Application and an ISO C standard
401 strictly conforming program is that the latter is not allowed to use features of POSIX that are not
402 in the ISO C standard.

403 *A.2.2.2 Conforming POSIX Application*

404 Examples of <National Bodies> include ANSI, BSI, and AFNOR.

405 *A.2.2.3 Conforming POSIX Application Using Extensions*

406 Due to possible requirements for configuration or implementation characteristics in excess of the
407 specifications in **<limits.h>** or related to the hardware (such as array size or file space), not every
408 Conforming POSIX Application Using Extensions will run on every conforming
409 implementation.

410 *A.2.2.4 Strictly Conforming XSI Application*

411 This is intended to be upwards-compatible with the definition of a Strictly Conforming POSIX
412 Application, with the addition of the facilities and functionality included in the XSI extension.

413 *A.2.2.5 Conforming XSI Application Using Extensions*

414 Such applications may use extensions beyond the facilities defined by IEEE Std 1003.1-2001
415 including the XSI extension, but need to document the additional requirements.

### A.2.3 Language-Dependent Services for the C Programming Language

POSIX.1 is, for historical reasons, both a specification of an operating system interface, shell and utilities, and a C binding for that specification. Efforts had been previously undertaken to generate a language-independent specification; however, that had failed, and the fact that the ISO C standard is the *de facto* primary language on POSIX and the UNIX system makes this a necessary and workable situation.

### A.2.4 Other Language-Related Specifications

There is no additional rationale provided for this section.

## A.3    Definitions

The definitions in this section are stated so that they can be used as exact substitutes for the terms in text. They should not contain requirements or cross-references to sections within IEEE Std 1003.1-2001; that is accomplished by using an informative note. In addition, the term should not be included in its own definition. Where requirements or descriptions need to be addressed but cannot be included in the definitions, due to not meeting the above criteria, these occur in the General Concepts chapter.

In this revision, the definitions have been reworked extensively to meet style requirements and to include terms from the base documents (see the Scope).

Many of these definitions are necessarily circular, and some of the terms (such as ''process'') are variants of basic computing science terms that are inherently hard to define. Where some definitions are more conceptual and contain requirements, these appear in the General Concepts chapter. Those listed in this section appear in an alphabetical glossary format of terms.

Some definitions must allow extension to cover terms or facilities that are not explicitly mentioned in IEEE Std 1003.1-2001. For example, the definition of ''Extended Security Controls'' permits implementations beyond those defined in IEEE Std 1003.1-2001.

Some terms in the following list of notes do not appear in IEEE Std 1003.1-2001; these are marked suffixed with an asterisk (*). Many of them have been specifically excluded from IEEE Std 1003.1-2001 because they concern system administration, implementation, or other issues that are not specific to the programming interface. Those are marked with a reason, such as ''implementation-defined''.

**Appropriate Privileges**

One of the fundamental security problems with many historical UNIX systems has been that the privilege mechanism is monolithic—a user has either no privileges or *all* privileges. Thus, a successful ''trojan horse'' attack on a privileged process defeats all security provisions. Therefore, POSIX.1 allows more granular privilege mechanisms to be defined. For many historical implementations of the UNIX system, the presence of the term ''appropriate privileges'' in POSIX.1 may be understood as a synonym for ''superuser'' (UID 0). However, other systems have emerged where this is not the case and each discrete controllable action has *appropriate privileges* associated with it. Because this mechanism is implementation-defined, it must be described in the conformance document. Although that description affects several parts of POSIX.1 where the term ''appropriate privilege'' is used, because the term ''implementation-defined'' only appears here, the description of the entire mechanism and its effects on these other sections belongs in this equivalent section of the conformance document. This is especially convenient for implementations with a single mechanism that applies in all areas, since it only needs to be described once.

**Byte**

461 The restriction that a byte is now exactly eight bits was a conscious decision by the standard
462 developers. It came about due to a combination of factors, primarily the use of the type **int8_t**
463 within the networking functions and the alignment with the ISO/IEC 9899:1999 standard, where
464 the **intN_t** types are now defined.

465 According to the ISO/IEC 9899:1999 standard:

466 • The **[u]intN_t** types must be two's complement with no padding bits and no illegal values.

467 • All types (apart from bit fields, which are not relevant here) must occupy an integral number
468   of bytes.

469 • If a type with width *W* occupies *B* bytes with *C* bits per byte (*C* is the value of {CHAR_BIT}),
470   then it has *P* padding bits where $P+W=B*C$.

471 • Therefore, for **int8_t** $P=0$, $W=8$. Since $B{\geq}1$, $C{\geq}8$, the only solution is $B=1$, $C=8$.

472 The standard developers also felt that this was not an undue restriction for the current state-of-
473 the-art for this version of IEEE Std 1003.1, but recognize that if industry trends continue, a wider
474 character type may be required in the future.

**Character**

476 The term ''character'' is used to mean a sequence of one or more bytes representing a single
477 graphic symbol. The deviation in the exact text of the ISO C standard definition for ''byte'' meets
478 the intent of the rationale of the ISO C standard also clears up the ambiguity raised by the term
479 ''basic execution character set''. The octet-minimum requirement is a reflection of the
480 {CHAR_BIT} value.

**Clock Tick**

482 The ISO C standard defines a similar interval for use by the *clock*() function. There is no
483 requirement that these intervals be the same. In historical implementations these intervals are
484 different.

**Command**

486 The terms ''command'' and ''utility'' are related but have distinct meanings. Command is
487 defined as ''a directive to a shell to perform a specific task''. The directive can be in the form of a
488 single utility name (for example, *ls*), or the directive can take the form of a compound command
489 (for example, `"ls | grep name | pr"`). A utility is a program that can be called by name
490 from a shell. Issuing only the name of the utility to a shell is the equivalent of a one-word
491 command. A utility may be invoked as a separate program that executes in a different process
492 than the command language interpreter, or it may be implemented as a part of the command
493 language interpreter. For example, the *echo* command (the directive to perform a specific task)
494 may be implemented such that the *echo* utility (the logic that performs the task of echoing) is in a
495 separate program; therefore, it is executed in a process that is different from the command
496 language interpreter. Conversely, the logic that performs the *echo* utility could be built into the
497 command language interpreter; therefore, it could execute in the same process as the command
498 language interpreter.

499 The terms ''tool'' and ''application'' can be thought of as being synonymous with ''utility'' from
500 the perspective of the operating system kernel. Tools, applications, and utilities historically have
501 run, typically, in processes above the kernel level. Tools and utilities historically have been a part
502 of the operating system non-kernel code and have performed system-related functions, such as
503 listing directory contents, checking file systems, repairing file systems, or extracting system

504  status information. Applications have not generally been a part of the operating system, and
505  they perform non-system-related functions, such as word processing, architectural design,
506  mechanical design, workstation publishing, or financial analysis. Utilities have most frequently
507  been provided by the operating system distributor, applications by third-party software
508  distributors, or by the users themselves. Nevertheless, IEEE Std 1003.1-2001 does not
509  differentiate between tools, utilities, and applications when it comes to receiving services from
510  the system, a shell, or the standard utilities. (For example, the *xargs* utility invokes another
511  utility; it would be of fairly limited usefulness if the users could not run their own applications
512  in place of the standard utilities.) Utilities are not applications in the sense that they are not
513  themselves subject to the restrictions of IEEE Std 1003.1-2001 or any other standard—there is no
514  requirement for *grep*, *stty*, or any of the utilities defined here to be any of the classes of
515  conforming applications.

**Column Positions**

517  In most 1-byte character sets, such as ASCII, the concept of column positions is identical to
518  character positions and to bytes. Therefore, it has been historically acceptable for some
519  implementations to describe line folding or tab stops or table column alignment in terms of bytes
520  or character positions. Other character sets pose complications, as they can have internal
521  representations longer than one octet and they can have display characters that have different
522  widths on the terminal screen or printer.

523  In IEEE Std 1003.1-2001 the term ''column positions'' has been defined to mean character—not
524  byte—positions in input files (such as ''column position 7 of the FORTRAN input''). Output files
525  describe the column position in terms of the display width of the narrowest printable character
526  in the character set, adjusted to fit the characteristics of the output device. It is very possible that
527  *n* column positions will not be able to hold *n* characters in some character sets, unless all of those
528  characters are of the narrowest width. It is assumed that the implementation is aware of the
529  width of the various characters, deriving this information from the value of *LC_CTYPE*, and thus
530  can determine how many column positions to allot for each character in those utilities where it is
531  important.

532  The term ''column position'' was used instead of the more natural ''column'' because the latter is
533  frequently used in the different contexts of columns of figures, columns of table values, and so
534  on. Wherever confusion might result, these latter types of columns are referred to as ''text
535  columns''.

**Controlling Terminal**

537  The question of which of possibly several special files referring to the terminal is meant is not
538  addressed in POSIX.1. The filename **/dev/tty** is a synonym for the controlling terminal associated
539  with a process.

**Device Number***

541  The concept is handled in *stat*( ) as *ID of device*.

542       **Direct I/O**

543       Historically, direct I/O refers to the system bypassing intermediate buffering, but may be
544       extended to cover implementation-defined optimizations.

545       **Directory**

546       The format of the directory file is implementation-defined and differs radically between
547       System V and 4.3 BSD. However, routines (derived from 4.3 BSD) for accessing directories and
548       certain constraints on the format of the information returned by those routines are described in
549       the **<dirent.h>** header.

550       **Directory Entry**

551       Throughout IEEE Std 1003.1-2001, the term ''link'' is used (about the *link*( ) function, for
552       example) in describing the objects that point to files from directories.

553       **Display**

554       The Shell and Utilities volume of IEEE Std 1003.1-2001 assigns precise requirements for the
555       terms ''display'' and ''write''. Some historical systems have chosen to implement certain utilities
556       without using the traditional file descriptor model. For example, the *vi* editor might employ
557       direct screen memory updates on a personal computer, rather than a *write*( ) system call. An
558       instance of user prompting might appear in a dialog box, rather than with standard error. When
559       the Shell and Utilities volume of IEEE Std 1003.1-2001 uses the term ''display'', the method of
560       outputting to the terminal is unspecified; many historical implementations use *termcap* or
561       *terminfo*, but this is not a requirement. The term ''write'' is used when the Shell and Utilities
562       volume of IEEE Std 1003.1-2001 mandates that a file descriptor be used and that the output can
563       be redirected. However, it is assumed that when the writing is directly to the terminal (it has not
564       been redirected elsewhere), there is no practical way for a user or test suite to determine whether
565       a file descriptor is being used. Therefore, the use of a file descriptor is mandated only for the
566       redirection case and the implementation is free to use any method when the output is not
567       redirected. The verb *write* is used almost exclusively, with the very few exceptions of those
568       utilities where output redirection need not be supported: *tabs*, *talk*, *tput*, and *vi*.

569       **Dot**

570       The symbolic name *dot* is carefully used in POSIX.1 to distinguish the working directory
571       filename from a period or a decimal point.

572       **Dot-Dot**

573       Historical implementations permit the use of these filenames without their special meanings.
574       Such use precludes any meaningful use of these filenames by a Conforming POSIX.1
575       Application. Therefore, such use is considered an extension, the use of which makes an
576       implementation non-conforming; see also Section A.4.11 (on page 37).

577        **Epoch**

578        Historically, the origin of UNIX system time was referred to as ''00:00:00 GMT, January 1, 1970''.
579        Greenwich Mean Time is actually not a term acknowledged by the international standards
580        community; therefore, this term, ''Epoch'', is used to abbreviate the reference to the actual
581        standard, Coordinated Universal Time.

582        **FIFO Special File**

583        See **Pipe** (on page 24).

584        **File**

585        It is permissible for an implementation-defined file type to be non-readable or non-writable.

586        **File Classes**

587        These classes correspond to the historical sets of permission bits. The classes are general to
588        allow implementations flexibility in expanding the access mechanism for more stringent security
589        environments. Note that a process is in one and only one class, so there is no ambiguity.

590        **Filename**

591        At the present time, the primary responsibility for truncating filenames containing multi-byte
592        characters must reside with the application. Some industry groups involved in
593        internationalization believe that in the future the responsibility must reside with the kernel. For
594        the moment, a clearer understanding of the implications of making the kernel responsible for
595        truncation of multi-byte filenames is needed.

596        Character-level truncation was not adopted because there is no support in POSIX.1 that advises
597        how the kernel distinguishes between single and multi-byte characters. Until that time, it must
598        be incumbent upon application writers to determine where multi-byte characters must be
599        truncated.

600        **File System**

601        Historically, the meaning of this term has been overloaded with two meanings: that of the
602        complete file hierarchy, and that of a mountable subset of that hierarchy; that is, a mounted file
603        system. POSIX.1 uses the term ''file system'' in the second sense, except that it is limited to the
604        scope of a process (and a process' root directory). This usage also clarifies the domain in which a
605        file serial number is unique.

606        **Graphic Character**

607        This definition is made available for those definitions (in particular, *TZ*) which must exclude
608        control characters.

609        **Group Database**

610        See **User Database** (on page 32).

611          **Group File**\*

612          Implementation-defined; see **User Database** (on page 32).


613          **Historical Implementations**\*

614          This refers to previously existing implementations of programming interfaces and operating
615          systems that are related to the interface specified by POSIX.1.


616          **Hosted Implementation**\*

617          This refers to a POSIX.1 implementation that is accomplished through interfaces from the
618          POSIX.1 services to some alternate form of operating system kernel services. Note that the line
619          between a hosted implementation and a native implementation is blurred, since most
620          implementations will provide some services directly from the kernel and others through some
621          indirect path. (For example, *fopen*( ) might use *open*( ); or *mkfifo*( ) might use *mknod*( ).) There is
622          no necessary relationship between the type of implementation and its correctness, performance,
623          and/or reliability.


624          **Implementation**\*

625          This term is generally used instead of its synonym, ''system'', to emphasize the consequences of
626          decisions to be made by system implementors. Perhaps if no options or extensions to POSIX.1
627          were allowed, this usage would not have occurred.

628          The term ''specific implementation'' is sometimes used as a synonym for ''implementation''.
629          This should not be interpreted too narrowly; both terms can represent a relatively broad group
630          of systems. For example, a hardware vendor could market a very wide selection of systems that
631          all used the same instruction set, with some systems desktop models and others large multi-user
632          minicomputers. This wide range would probably share a common POSIX.1 operating system,
633          allowing an application compiled for one to be used on any of the others; this is a [*specific*]
634          *implementation*. However, such a wide range of machines probably has some differences
635          between the models. Some may have different clock rates, different file systems, different
636          resource limits, different network connections, and so on, depending on their sizes or intended
637          usages. Even on two identical machines, the system administrators may configure them
638          differently. Each of these different systems is known by the term ''a specific instance of a specific
639          implementation''. This term is only used in the portions of POSIX.1 dealing with runtime
640          queries: *sysconf*( ) and *pathconf*( ).


641          **Incomplete Pathname**\*

642          Absolute pathname has been adequately defined.


643          **Job Control**

644          In order to understand the job control facilities in POSIX.1 it is useful to understand how they
645          are used by a job control-cognizant shell to create the user interface effect of job control.

646          While the job control facilities supplied by POSIX.1 can, in theory, support different types of
647          interactive job control interfaces supplied by different types of shells, there was historically one
648          particular interface that was most common when the standard was originally developed
649          (provided by BSD C Shell).                                                                      |

650          This discussion describes that interface as a means of illustrating how the POSIX.1 job control  |
651          facilities can be used.

652  Job control allows users to selectively stop (suspend) the execution of processes and continue
653  (resume) their execution at a later point. The user typically employs this facility via the
654  interactive interface jointly supplied by the terminal I/O driver and a command interpreter
655  (shell).

656  The user can launch jobs (command pipelines) in either the foreground or background. When
657  launched in the foreground, the shell waits for the job to complete before prompting for
658  additional commands. When launched in the background, the shell does not wait, but
659  immediately prompts for new commands.

660  If the user launches a job in the foreground and subsequently regrets this, the user can type the
661  suspend character (typically set to <control>-Z), which causes the foreground job to stop and the
662  shell to begin prompting for new commands. The stopped job can be continued by the user (via
663  special shell commands) either as a foreground job or as a background job. Background jobs can
664  also be moved into the foreground via shell commands.

665  If a background job attempts to access the login terminal (controlling terminal), it is stopped by
666  the terminal driver and the shell is notified, which, in turn, notifies the user. (Terminal access
667  includes *read*( ) and certain terminal control functions, and conditionally includes *write*( ).) The
668  user can continue the stopped job in the foreground, thus allowing the terminal access to
669  succeed in an orderly fashion. After the terminal access succeeds, the user can optionally move
670  the job into the background via the suspend character and shell commands.

671  *Implementing Job Control Shells*

672  The interactive interface described previously can be accomplished using the POSIX.1 job
673  control facilities in the following way.

674  The key feature necessary to provide job control is a way to group processes into jobs. This
675  grouping is necessary in order to direct signals to a single job and also to identify which job is in
676  the foreground. (There is at most one job that is in the foreground on any controlling terminal at
677  a time.)

678  The concept of process groups is used to provide this grouping. The shell places each job in a
679  separate process group via the *setpgid*( ) function. To do this, the *setpgid*( ) function is invoked by
680  the shell for each process in the job. It is actually useful to invoke *setpgid*( ) twice for each
681  process: once in the child process, after calling *fork*( ) to create the process, but before calling one
682  of the *exec* family of functions to begin execution of the program, and once in the parent shell
683  process, after calling *fork*( ) to create the child. The redundant invocation avoids a race condition
684  by ensuring that the child process is placed into the new process group before either the parent
685  or the child relies on this being the case. The process group ID for the job is selected by the shell
686  to be equal to the process ID of one of the processes in the job. Some shells choose to make one
687  process in the job be the parent of the other processes in the job (if any). Other shells (for
688  example, the C Shell) choose to make themselves the parent of all processes in the pipeline (job).
689  In order to support this latter case, the *setpgid*( ) function accepts a process group ID parameter
690  since the correct process group ID cannot be inherited from the shell. The shell itself is
691  considered to be a job and is the sole process in its own process group.

692  The shell also controls which job is currently in the foreground. A foreground and background
693  job differ in two ways: the shell waits for a foreground command to complete (or stop) before
694  continuing to read new commands, and the terminal I/O driver inhibits terminal access by
695  background jobs (causing the processes to stop). Thus, the shell must work cooperatively with
696  the terminal I/O driver and have a common understanding of which job is currently in the
697  foreground. It is the user who decides which command should be currently in the foreground,
698  and the user informs the shell via shell commands. The shell, in turn, informs the terminal I/O
699  driver via the *tcsetpgrp*( ) function. This indicates to the terminal I/O driver the process group ID

700  of the foreground process group (job). When the current foreground job either stops or
701  terminates, the shell places itself in the foreground via *tcsetpgrp*( ) before prompting for
702  additional commands. Note that when a job is created the new process group begins as a
703  background process group. It requires an explicit act of the shell via *tcsetpgrp*( ) to move a
704  process group (job) into the foreground.

705  When a process in a job stops or terminates, its parent (for example, the shell) receives
706  synchronous notification by calling the *waitpid*( ) function with the WUNTRACED flag set.
707  Asynchronous notification is also provided when the parent establishes a signal handler for
708  SIGCHLD and does not specify the SA_NOCLDSTOP flag. Usually all processes in a job stop as
709  a unit since the terminal I/O driver always sends job control stop signals to all processes in the
710  process group.

711  To continue a stopped job, the shell sends the SIGCONT signal to the process group of the job. In
712  addition, if the job is being continued in the foreground, the shell invokes *tcsetpgrp*( ) to place the
713  job in the foreground before sending SIGCONT. Otherwise, the shell leaves itself in the
714  foreground and reads additional commands.

715  There is additional flexibility in the POSIX.1 job control facilities that allows deviations from the
716  typical interface. Clearing the TOSTOP terminal flag allows background jobs to perform *write*( )
717  functions without stopping. The same effect can be achieved on a per-process basis by having a
718  process set the signal action for SIGTTOU to SIG_IGN.

719  Note that the terms ''job'' and ''process group'' can be used interchangeably. A login session that
720  is not using the job control facilities can be thought of as a large collection of processes that are
721  all in the same job (process group). Such a login session may have a partial distinction between
722  foreground and background processes; that is, the shell may choose to wait for some processes
723  before continuing to read new commands and may not wait for other processes. However, the
724  terminal I/O driver will consider all these processes to be in the foreground since they are all
725  members of the same process group.

726  In addition to the basic job control operations already mentioned, a job control-cognizant shell
727  needs to perform the following actions.

728  When a foreground (not background) job stops, the shell must sample and remember the current
729  terminal settings so that it can restore them later when it continues the stopped job in the
730  foreground (via the *tcgetattr*( ) and *tcsetattr*( ) functions).

731  Because a shell itself can be spawned from a shell, it must take special action to ensure that
732  subshells interact well with their parent shells.

733  A subshell can be spawned to perform an interactive function (prompting the terminal for
734  commands) or a non-interactive function (reading commands from a file). When operating non-
735  interactively, the job control shell will refrain from performing the job control-specific actions
736  described above. It will behave as a shell that does not support job control. For example, all jobs
737  will be left in the same process group as the shell, which itself remains in the process group
738  established for it by its parent. This allows the shell and its children to be treated as a single job
739  by a parent shell, and they can be affected as a unit by terminal keyboard signals.

740  An interactive subshell can be spawned from another job control-cognizant shell in either the
741  foreground or background. (For example, from the C Shell, the user can execute the command,
742  csh *&*.) Before the subshell activates job control by calling *setpgid*( ) to place itself in its own
743  process group and *tcsetpgrp*( ) to place its new process group in the foreground, it needs to
744  ensure that it has already been placed in the foreground by its parent. (Otherwise, there could
745  be multiple job control shells that simultaneously attempt to control mediation of the terminal.)
746  To determine this, the shell retrieves its own process group via *getpgrp*( ) and the process group
747  of the current foreground job via *tcgetpgrp*( ). If these are not equal, the shell sends SIGTTIN to

748 its own process group, causing itself to stop. When continued later by its parent, the shell
749 repeats the process group check. When the process groups finally match, the shell is in the
750 foreground and it can proceed to take control. After this point, the shell ignores all the job
751 control stop signals so that it does not inadvertently stop itself.

752 *Implementing Job Control Applications*

753 Most applications do not need to be aware of job control signals and operations; the intuitively
754 correct behavior happens by default. However, sometimes an application can inadvertently
755 interfere with normal job control processing, or an application may choose to overtly effect job
756 control in cooperation with normal shell procedures.

757 An application can inadvertently subvert job control processing by ''blindly'' altering the
758 handling of signals. A common application error is to learn how many signals the system
759 supports and to ignore or catch them all. Such an application makes the assumption that it does
760 not know what this signal is, but knows the right handling action for it. The system may
761 initialize the handling of job control stop signals so that they are being ignored. This allows
762 shells that do not support job control to inherit and propagate these settings and hence to be
763 immune to stop signals. A job control shell will set the handling to the default action and
764 propagate this, allowing processes to stop. In doing so, the job control shell is taking
765 responsibility for restarting the stopped applications. If an application wishes to catch the stop
766 signals itself, it should first determine their inherited handling states. If a stop signal is being
767 ignored, the application should continue to ignore it. This is directly analogous to the
768 recommended handling of SIGINT described in the referenced UNIX Programmer's Manual.

769 If an application is reading the terminal and has disabled the interpretation of special characters
770 (by clearing the ISIG flag), the terminal I/O driver will not send SIGTSTP when the suspend
771 character is typed. Such an application can simulate the effect of the suspend character by
772 recognizing it and sending SIGTSTP to its process group as the terminal driver would have
773 done. Note that the signal is sent to the process group, not just to the application itself; this
774 ensures that other processes in the job also stop. (Note also that other processes in the job could
775 be children, siblings, or even ancestors.) Applications should not assume that the suspend
776 character is <control>-Z (or any particular value); they should retrieve the current setting at
777 startup.

778 *Implementing Job Control Systems*

779 The intent in adding 4.2 BSD-style job control functionality was to adopt the necessary 4.2 BSD
780 programmatic interface with only minimal changes to resolve syntactic or semantic conflicts
781 with System V or to close recognized security holes. The goal was to maximize the ease of
782 providing both conforming implementations and Conforming POSIX.1 Applications.

783 It is only useful for a process to be affected by job control signals if it is the descendant of a job
784 control shell. Otherwise, there will be nothing that continues the stopped process.

785 POSIX.1 does not specify how controlling terminal access is affected by a user logging out (that
786 is, by a controlling process terminating). 4.2 BSD uses the *vhangup*( ) function to prevent any
787 access to the controlling terminal through file descriptors opened prior to logout. System V does
788 not prevent controlling terminal access through file descriptors opened prior to logout (except
789 for the case of the special file, **/dev/tty**). Some implementations choose to make processes
790 immune from job control after logout (that is, such processes are always treated as if in the
791 foreground); other implementations continue to enforce foreground/background checks after
792 logout. Therefore, a Conforming POSIX.1 Application should not attempt to access the
793 controlling terminal after logout since such access is unreliable. If an implementation chooses to
794 deny access to a controlling terminal after its controlling process exits, POSIX.1 requires a certain
795 type of behavior (see **Controlling Terminal** (on page 15)).

796      **Kernel**\*

797      See **System Call**\* (on page 30).

800      **Library Routine**\*

799      See **System Call**\* (on page 30).

800      **Logical Device**\*

801      Implementation-defined.

802      **Map**

803      The definition of map is included to clarify the usage of mapped pages in the description of the
804      behavior of process memory locking.

805      **Memory-Resident**

806      The term ''memory-resident'' is historically understood to mean that the so-called resident
807      pages are actually present in the physical memory of the computer system and are immune from
808      swapping, paging, copy-on-write faults, and so on. This is the actual intent of
809      IEEE Std 1003.1-2001 in the process memory locking section for implementations where this is
810      logical. But for some implementations—primarily mainframes—actually locking pages into
811      primary storage is not advantageous to other system objectives, such as maximizing throughput.
812      For such implementations, memory locking is a ''hint'' to the implementation that the
813      application wishes to avoid situations that would cause long latencies in accessing memory.
814      Furthermore, there are other implementation-defined issues with minimizing memory access
815      latencies that ''memory residency'' does not address—such as MMU reload faults. The definition
816      attempts to accommodate various implementations while allowing conforming applications to
817      specify to the implementation that they want or need the best memory access times that the
818      implementation can provide.

819      **Memory Object**\*

820      The term ''memory object'' usually implies shared memory. If the object is the same as a
821      filename in the file system name space of the implementation, it is expected that the data written
822      into the memory object be preserved on disk. A memory object may also apply to a physical
823      device on an implementation. In this case, writes to the memory object are sent to the controller
824      for the device and reads result in control registers being returned.

825      **Mount Point**\*

826      The directory on which a ''mounted file system'' is mounted. This term, like *mount*() and
827      *umount*(), was not included because it was implementation-defined.

828      **Mounted File System**\*

829      See **File System** (on page 17).

830      **Name**

831      There are no explicit limits in IEEE Std 1003.1-2001 on the sizes of names, words (see the
832      definition of word in the Base Definitions volume of IEEE Std 1003.1-2001), lines, or other
833      objects. However, other implicit limits do apply: shell script lines produced by many of the
834      standard utilities cannot exceed {LINE_MAX} and the sum of exported variables comes under
835      the {ARG_MAX} limit. Historical shells dynamically allocate memory for names and words and
836      parse incoming lines a character at a time. Lines cannot have an arbitrary {LINE_MAX} limit
837      because of historical practice, such as makefiles, where *make* removes the <newline>s associated
838      with the commands for a target and presents the shell with one very long line. The text on
839      INPUT FILES in the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 1.11, Utility
840      Description Defaults does allow a shell to run out of memory, but it cannot have arbitrary
841      programming limits.

842      **Native Implementation**\*

843      This refers to an implementation of POSIX.1 that interfaces directly to an operating system
844      kernel; see also *hosted implementation*. A similar concept is a native UNIX system, which would      |
845      be a kernel derived from one of the original UNIX system products.

846      **Nice Value**

847      This definition is not intended to suggest that all processes in a system have priorities that are
848      comparable. Scheduling policy extensions, such as adding realtime priorities, make the notion of
849      a single underlying priority for all scheduling policies problematic. Some implementations may
850      implement the features related to *nice* to affect all processes on the system, others to affect just
851      the general time-sharing activities implied by IEEE Std 1003.1-2001, and others may have no
852      effect at all. Because of the use of ''implementation-defined'' in *nice* and *renice*, a wide range of
853      implementation strategies is possible.

854      **Open File Description**

855      An ''open file description'', as it is currently named, describes how a file is being accessed. What
856      is currently called a ''file descriptor'' is actually just an identifier or ''handle''; it does not actually
857      describe anything.

858      The following alternate names were discussed:

859      • For ''open file description'':
860        ''open instance'', ''file access description'', ''open file information'', and ''file access
861        information''.

862      • For ''file descriptor'':
863        ''file handle'', ''file number'' (cf., *fileno*()). Some historical implementations use the term ''file
864        table entry''.

865      **Orphaned Process Group**

866      Historical implementations have a concept of an orphaned process, which is a process whose
867      parent process has exited. When job control is in use, it is necessary to prevent processes from
868      being stopped in response to interactions with the terminal after they no longer are controlled by
869      a job control-cognizant program. Because signals generated by the terminal are sent to a process
870      group and not to individual processes, and because a signal may be provoked by a process that
871      is not orphaned, but sent to another process that is orphaned, it is necessary to define an
872      orphaned process group. The definition assumes that a process group will be manipulated as a
873      group and that the job control-cognizant process controlling the group is outside of the group

874    and is the parent of at least one process in the group (so that state changes may be reported via
875    *waitpid*( )).  Therefore, a group is considered to be controlled as long as at least one process in the
876    group has a parent that is outside of the process group, but within the session.

877    This definition of orphaned process groups ensures that a session leader's process group is
878    always considered to be orphaned, and thus it is prevented from stopping in response to
879    terminal signals.

880    **Page**

881    The term ''page'' is defined to support the description of the behavior of memory mapping for
882    shared memory and memory mapped files, and the description of the behavior of process
883    memory locking. It is not intended to imply that shared memory/file mapping and memory
884    locking are applicable only to ''paged'' architectures. For the purposes of IEEE Std 1003.1-2001,
885    whatever the granularity on which an architecture supports mapping or locking, this is
886    considered to be a ''page'' . If an architecture cannot support the memory mapping or locking
887    functions specified by IEEE Std 1003.1-2001 on any granularity, then these options will not be
888    implemented on the architecture.

889    **Passwd File**\*

890    Implementation-defined; see **User Database** (on page 32).

891    **Parent Directory**

892    There may be more than one directory entry pointing to a given directory in some
893    implementations. The wording here identifies that exactly one of those is the parent directory. In
894    pathname resolution, dot-dot is identified as the way that the unique directory is identified.
895    (That is, the parent directory is the one to which dot-dot points.) In the case of a remote file
896    system, if the same file system is mounted several times, it would appear as if they were distinct
897    file systems (with interesting synchronization properties).

898    **Pipe**

899    It proved convenient to define a pipe as a special case of a FIFO, even though historically the
900    latter was not introduced until System III and does not exist at all in 4.3 BSD.

901    **Portable Filename Character Set**

902    The encoding of this character set is not specified—specifically, ASCII is not required. But the
903    implementation must provide a unique character code for each of the printable graphics
904    specified by POSIX.1; see also Section A.4.6 (on page 34).

905    Situations where characters beyond the portable filename character set (or historically ASCII or
906    the ISO/IEC 646:1991 standard) would be used (in a context where the portable filename
907    character set or the ISO/IEC 646:1991 standard is required by POSIX.1) are expected to be
908    common.  Although such a situation renders the use technically non-compliant, mutual
909    agreement among the users of an extended character set will make such use portable between
910    those users. Such a mutual agreement could be formalized as an optional extension to POSIX.1.
911    (Making it required would eliminate too many possible systems, as even those systems using the
912    ISO/IEC 646:1991 standard as a base character set extend their character sets for Western
913    Europe and the rest of the world in different ways.)

914    Nothing in POSIX.1 is intended to preclude the use of extended characters where interchange is
915    not required or where mutual agreement is obtained. It has been suggested that in several places
916    ''should'' be used instead of ''shall''. Because (in the worst case) use of any character beyond the

917   portable filename character set would render the program or data not portable to all possible
918   systems, no extensions are permitted in this context.

**Regular File**

919

920   POSIX.1 does not intend to preclude the addition of structuring data (for example, record
921   lengths) in the file, as long as such data is not visible to an application that uses the features
922   described in POSIX.1.

**Root Directory**

923

924   This definition permits the operation of *chroot*( ), even though that function is not in POSIX.1; see
925   also Section A.4.5 (on page 33).

**Root File System***

926

927   Implementation-defined.

**Root of a File System***

928

929   Implementation-defined; see **Mount Point*** (on page 22).

**Signal**

930

931   The definition implies a double meaning for the term. Although a signal is an event, common
932   usage implies that a signal is an identifier of the class of event.

**Superuser***

933

934   This concept, with great historical significance to UNIX system users, has been replaced with the
935   notion of appropriate privileges.

**Supplementary Group ID**

936

937   The POSIX.1-1990 standard is inconsistent in its treatment of supplementary groups. The
938   definition of supplementary group ID explicitly permits the effective group ID to be included in
939   the set, but wording in the description of the *setuid*( ) and *setgid*( ) functions states: ''Any
940   supplementary group IDs of the calling process remain unchanged by these function calls''. In
941   the case of *setgid*( ) this contradicts that definition. In addition, some felt that the unspecified
942   behavior in the definition of supplementary group IDs adds unnecessary portability problems.
943   The standard developers considered several solutions to this problem:

944   1.  Reword the description of *setgid*( ) to permit it to change the supplementary group IDs to
945       reflect the new effective group ID. A problem with this is that it adds more ''may''s to the
946       wording and does not address the portability problems of this optional behavior.

947   2.  Mandate the inclusion of the effective group ID in the supplementary set (giving
948       {NGROUPS_MAX} a minimum value of 1). This is the behavior of 4.4 BSD. In that system,
949       the effective group ID is the first element of the array of supplementary group IDs (there is
950       no separate copy stored, and changes to the effective group ID are made only in the
951       supplementary group set). By convention, the initial value of the effective group ID is
952       duplicated elsewhere in the array so that the initial value is not lost when executing a set-
953       group-ID program.

954   3.  Change the definition of supplementary group ID to exclude the effective group ID and
955       specify that the effective group ID does not change the set of supplementary group IDs.
956       This is the behavior of 4.2 BSD, 4.3 BSD, and System V Release 4.

957      4.   Change the definition of supplementary group ID to exclude the effective group ID, and
958           require that *getgroups*( ) return the union of the effective group ID and the supplementary
959           group IDs.

960      5.   Change the definition of {NGROUPS_MAX} to be one more than the number of
961           supplementary group IDs, so it continues to be the number of values returned by
962           *getgroups*( ) and existing applications continue to work. This alternative is effectively the
963           same as the second (and might actually have the same implementation).

964   The standard developers decided to permit either 2 or 3. The effective group ID is orthogonal to
965   the set of supplementary group IDs, and it is implementation-defined whether *getgroups*( )
966   returns this. If the effective group ID is returned with the set of supplementary group IDs, then
967   all changes to the effective group ID affect the supplementary group set returned by *getgroups*( ).
968   It is permissible to eliminate duplicates from the list returned by *getgroups*( ). However, if a
969   group ID is contained in the set of supplementary group IDs, setting the group ID to that value
970   and then to a different value should not remove that value from the supplementary group IDs.

971   The definition of supplementary group IDs has been changed to not include the effective group
972   ID. This simplifies permanent rationale and makes the relevant functions easier to understand.
973   The *getgroups*( ) function has been modified so that it can, on an implementation-defined basis,
974   return the effective group ID. By making this change, functions that modify the effective group
975   ID do not need to discuss adding to the supplementary group list; the only view into the
976   supplementary group list that the application writer has is through the *getgroups*( ) function.

977   **Symbolic Link**

978   Many implementations associate no attributes, including ownership with symbolic links.
979   Security experts encouraged consideration for defining these attributes as optional.
980   Consideration was given to changing *utime*( ) to allow modification of the times for a symbolic
981   link, or as an alternative adding an *lutime*( ) interface. Modifications to *chown*( ) were also
982   considered: allow changing symbolic link ownership or alternatively adding *lchown*( ). As a
983   result of alignment with the Single UNIX Specification, the *lchown*( ) function is included as part
984   of the XSI extension and XSI-conformant systems may support an owner and a group associated
985   with a symbolic link. There was no consensus to define further attributes for symbolic links, and
986   for systems not supporting the XSI extension only the file type bits in the *st_mode* member and
987   the *st_size* member of the **stat** structure are required to be applicable to symbolic links.

988   Historical implementations were followed when determining which interfaces should apply to
989   symbolic links. Interfaces that historically followed symbolic links include *chmod*( ), *link*( ), and
990   *utime*( ). Interfaces that historically do not follow symbolic links include *chown*( ), *lstat*( ),
991   *readlink*( ), *rename*( ), *remove*( ), *rmdir*( ), and *unlink*( ). IEEE Std 1003.1-2001 deviates from
992   historical practice only in the case of *chown*( ). Because there is no requirement for systems not
993   supporting the XSI extension that there is an association of ownership with symbolic links, there
994   was no interface in the base standard to change ownership. In addition, other implementations
995   of symbolic links have modified *chown*( ) to follow symbolic links.

996   In the case of symbolic links, IEEE Std 1003.1-2001 states that a trailing slash is considered to be
997   the final component of a pathname rather than the pathname component that preceded it. This is
998   the behavior of historical implementations. For example, for **/a/b** and **/a/b/**, if **/a/b** is a symbolic
999   link to a directory, then **/a/b** refers to the symbolic link, and **/a/b/** is the same as **/a/b/.**, which is the
1000  directory to which the symbolic link points.

1001  For multi-level security purposes, it is possible to have the link read mode govern permission for
1002  the *readlink*( ) function. It is also possible that the read permissions of the directory containing
1003  the link be used for this purpose. Implementations may choose to use either of these methods;
1004  however, this is not current practice and neither method is specified.

1005 Several reasons were advanced for requiring that when a symbolic link is used as the source
1006 argument to the *link*( ) function, the resulting link will apply to the file named by the contents of
1007 the symbolic link rather than to the symbolic link itself. This is the case in historical
1008 implementations. This action was preferred, as it supported the traditional idea of persistence
1009 with respect to the target of a hard link. This decision is appropriate in light of a previous
1010 decision not to require association of attributes with symbolic links, thereby allowing
1011 implementations which do not use inodes. Opposition centered on the lack of symmetry on the
1012 part of the *link*( ) and *unlink*( ) function pair with respect to symbolic links.

1013 Because a symbolic link and its referenced object coexist in the file system name space, confusion
1014 can arise in distinguishing between the link itself and the referenced object. Historically, utilities
1015 and system calls have adopted their own link following conventions in a somewhat *ad hoc*
1016 fashion. Rules for a uniform approach are outlined here, although historical practice has been
1017 adhered to as much as was possible. To promote consistent system use, user-written utilities are
1018 encouraged to follow these same rules.

1019 Symbolic links are handled either by operating on the link itself, or by operating on the object
1020 referenced by the link. In the latter case, an application or system call is said to ''follow'' the link.
1021 Symbolic links may reference other symbolic links, in which case links are dereferenced until an
1022 object that is not a symbolic link is found, a symbolic link that references a file that does not exist
1023 is found, or a loop is detected. (Current implementations do not detect loops, but have a limit on
1024 the number of symbolic links that they will dereference before declaring it an error.)

1025 There are four domains for which default symbolic link policy is established in a system. In
1026 almost all cases, there are utility options that override this default behavior. The four domains
1027 are as follows:

1028     1.   Symbolic links specified to system calls that take filename arguments

1029     2.   Symbolic links specified as command line filename arguments to utilities that are not
1030         performing a traversal of a file hierarchy

1031     3.   Symbolic links referencing files not of type directory, specified to utilities that are
1032         performing a traversal of a file hierarchy

1033     4.   Symbolic links referencing files of type directory, specified to utilities that are performing a
1034         traversal of a file hierarchy

1035 *First Domain*

1036 The first domain is considered in earlier rationale.

1037 *Second Domain*

1038 The reason this category is restricted to utilities that are not traversing the file hierarchy is that
1039 some standard utilities take an option that specifies a hierarchical traversal, but by default
1040 operate on the arguments themselves. Generally, users specifying the option for a file hierarchy
1041 traversal wish to operate on a single, physical hierarchy, and therefore symbolic links, which
1042 may reference files outside of the hierarchy, are ignored. For example, *chown owner file* is a
1043 different operation from the same command with the −**R** option specified. In this example, the
1044 behavior of the command *chown owner file* is described here, while the behavior of the command
1045 *chown* −**R** *owner file* is described in the third and fourth domains.

1046 The general rule is that the utilities in this category follow symbolic links named as arguments.

1047 Exceptions in the second domain are:

1048     • The *mv* and *rm* utilities do not follow symbolic links named as arguments, but respectively
1049       attempt to rename or delete them.

1050    • The *ls* utility is also an exception to this rule. For compatibility with historical systems, when
1051    the −**R** option is not specified, the *ls* utility follows symbolic links named as arguments if the
1052    −**L** option is specified or if the −**F**, −**d**, or −**l** options are not specified. (If the −**L** option is
1053    specified, *ls* always follows symbolic links; it is the only utility where the −**L** option affects its
1054    behavior even though a tree walk is not being performed.)

1055    All other standard utilities, when not traversing a file hierarchy, always follow symbolic links
1056    named as arguments.

1057    Historical practice is that the −**h** option is specified if standard utilities are to act upon symbolic
1058    links instead of upon their targets. Examples of commands that have historically had a −**h** option
1059    for this purpose are the *chgrp*, *chown*, *file*, and *test* utilities.

1060    *Third Domain*

1061    The third domain is symbolic links, referencing files not of type directory, specified to utilities
1062    that are performing a traversal of a file hierarchy. (This includes symbolic links specified as
1063    command line filename arguments or encountered during the traversal.)

1064    The intention of the Shell and Utilities volume of IEEE Std 1003.1-2001 is that the operation that
1065    the utility is performing is applied to the symbolic link itself, if that operation is applicable to
1066    symbolic links. The reason that the operation is not required is that symbolic links in some
1067    implementations do not have such attributes as a file owner, and therefore the *chown* operation
1068    would be meaningless. If symbolic links on the system have an owner, it is the intention that the
1069    utility *chown* cause the owner of the symbolic link to change. If symbolic links do not have an
1070    owner, the symbolic link should be ignored. Specifically, by default, no change should be made
1071    to the file referenced by the symbolic link.

1072    *Fourth Domain*

1073    The fourth domain is symbolic links referencing files of type directory, specified to utilities that
1074    are performing a traversal of a file hierarchy. (This includes symbolic links specified as
1075    command line filename arguments or encountered during the traversal.)

1076    Most standard utilities do not, by default, indirect into the file hierarchy referenced by the
1077    symbolic link. (The Shell and Utilities volume of IEEE Std 1003.1-2001 uses the informal term
1078    ''physical walk'' to describe this case. The case where the utility does indirect through the
1079    symbolic link is termed a ''logical walk''.)

1080    There are three reasons for the default to be a physical walk:

1081    1.    With very few exceptions, a physical walk has been the historical default on UNIX systems
1082          supporting symbolic links. Because some utilities (that is, *rm*) must default to a physical
1083          walk, regardless, changing historical practice in this regard would be confusing to users
1084          and needlessly incompatible.

1085    2.    For systems where symbolic links have the historical file attributes (that is, *owner*, *group*,
1086          *mode*), defaulting to a logical traversal would require the addition of a new option to the
1087          commands to modify the attributes of the link itself. This is painful and more complex
1088          than the alternatives.

1089    3.    There is a security issue with defaulting to a logical walk. Historically, the command
1090          *chown* −**R** *user file* has been safe for the superuser because *setuid* and *setgid* bits were lost
1091          when the ownership of the file was changed. If the walk were logical, changing ownership
1092          would no longer be safe because a user might have inserted a symbolic link pointing to any
1093          file in the tree. Again, this would necessitate the addition of an option to the commands
1094          doing hierarchy traversal to not indirect through the symbolic links, and historical scripts
1095          doing recursive walks would instantly become security problems. While this is mostly an

1096      issue for system administrators, it is preferable to not have different defaults for different
1097      classes of users.

1098    However, the standard developers agreed to leave it unspecified to achieve consensus.

1099    As consistently as possible, users may cause standard utilities performing a file hierarchy
1100    traversal to follow any symbolic links named on the command line, regardless of the type of file
1101    they reference, by specifying the –**H** (for half logical) option. This option is intended to make the
1102    command line name space look like the logical name space.

1103    As consistently as possible, users may cause standard utilities performing a file hierarchy
1104    traversal to follow any symbolic links named on the command line as well as any symbolic links
1105    encountered during the traversal, regardless of the type of file they reference, by specifying the
1106    –**L** (for logical) option. This option is intended to make the entire name space look like the
1107    logical name space.

1108    For consistency, implementors are encouraged to use the –**P** (for ''physical'') flag to specify the
1109    physical walk in utilities that do logical walks by default for whatever reason. The only standard
1110    utilities that require the –**P** option are *cd* and *pwd*; see the note below.

1111    When one or more of the –**H**, –**L**, and –**P** flags can be specified, the last one specified determines
1112    the behavior of the utility. This permits users to alias commands so that the default behavior is a
1113    logical walk and then override that behavior on the command line.

1114    *Exceptions in the Third and Fourth Domains*

1115    The *ls* and *rm* utilities are exceptions to these rules. The *rm* utility never follows symbolic links
1116    and does not support the –**H**, –**L**, or –**P** options. Some historical versions of *ls* always followed
1117    symbolic links given on the command line whether the –**L** option was specified or not. Historical
1118    versions of *ls* did not support the –**H** option. In IEEE Std 1003.1-2001, unless one of the –**H** or –**L**
1119    options is specified, the *ls* utility only follows symbolic links to directories that are given as
1120    operands. The *ls* utility does not support the –**P** option.

1121    The Shell and Utilities volume of IEEE Std 1003.1-2001 requires that the standard utilities *ls*, *find*,
1122    and *pax* detect infinite loops when doing logical walks; that is, a directory, or more commonly a
1123    symbolic link, that refers to an ancestor in the current file hierarchy. If the file system itself is
1124    corrupted, causing the infinite loop, it may be impossible to recover. Because *find* and *ls* are often
1125    used in system administration and security applications, they should attempt to recover and
1126    continue as best as they can. The *pax* utility should terminate because the archive it was creating
1127    is by definition corrupted. Other, less vital, utilities should probably simply terminate as well.
1128    Implementations are strongly encouraged to detect infinite loops in all utilities.

1129    Historical practice is shown in Table A-1 (on page 30). The heading **SVID3** stands for the Third
1130    Edition of the System V Interface Definition.

1131    Historically, several shells have had built-in versions of the *pwd* utility. In some of these shells,
1132    *pwd* reported the physical path, and in others, the logical path. Implementations of the shell
1133    corresponding to IEEE Std 1003.1-2001 must report the logical path by default. Earlier versions of
1134    IEEE Std 1003.1-2001 did not require the *pwd* utility to be a built-in utility. Now that *pwd* is
1135    required to set an environment variable in the current shell execution environment, it must be a
1136    built-in utility.

1137    The *cd* command is required, by default, to treat the filename dot-dot logically. Implementors are
1138    required to support the –**P** flag in *cd* so that users can have their current environment handled
1139    physically. In 4.3 BSD, *chgrp* during tree traversal changed the group of the symbolic link, not
1140    the target. Symbolic links in 4.4 BSD do not have *owner*, *group*, *mode*, or other standard UNIX
1141    system file attributes.

**Table A-1** Historical Practice for Symbolic Links

| Utility | SVID3 | 4.3 BSD | 4.4 BSD | POSIX | Comments |
|---------|-------|---------|---------|-------|----------|
| *cd* | | | | –**L** | Treat "`..`" logically. |
| *cd* | | | | –**P** | Treat "`..`" physically. |
| *chgrp* | | | –**H** | –**H** | Follow command line symlinks. |
| *chgrp* | | | –**h** | –**L** | Follow symlinks. |
| *chgrp* | –**h** | | | –**h** | Affect the symlink. |
| *chmod* | | | | | Affect the symlink. |
| *chmod* | | | –**H** | | Follow command line symlinks. |
| *chmod* | | | –**h** | | Follow symlinks. |
| *chown* | | | –**H** | –**H** | Follow command line symlinks. |
| *chown* | | | –**h** | –**L** | Follow symlinks. |
| *chown* | –**h** | | | –**h** | Affect the symlink. |
| *cp* | | | –**H** | –**H** | Follow command line symlinks. |
| *cp* | | | –**h** | –**L** | Follow symlinks. |
| *cpio* | –**L** | | –**L** | | Follow symlinks. |
| *du* | | | –**H** | –**H** | Follow command line symlinks. |
| *du* | | | –**h** | –**L** | Follow symlinks. |
| *file* | –**h** | | | –**h** | Affect the symlink. |
| *find* | | | –**H** | –**H** | Follow command line symlinks. |
| *find* | | | –**h** | –**L** | Follow symlinks. |
| *find* | –**follow** | | –**follow** | | Follow symlinks. |
| *ln* | –**s** | –**s** | –**s** | –**s** | Create a symbolic link. |
| *ls* | –**L** | –**L** | –**L** | –**L** | Follow symlinks. |
| *ls* | | | | –**H** | Follow command line symlinks. |
| *mv* | | | | | Operates on the symlink. |
| *pax* | | | –**H** | –**H** | Follow command line symlinks. |
| *pax* | | | –**h** | –**L** | Follow symlinks. |
| *pwd* | | | | –**L** | Printed path may contain symlinks. |
| *pwd* | | | | –**P** | Printed path will not contain symlinks. |
| *rm* | | | | | Operates on the symlink. |
| *tar* | | | –**H** | | Follow command line symlinks. |
| *tar* | | –**h** | –**h** | | Follow symlinks. |
| *test* | –**h** | | –**h** | –**h** | Affect the symlink. |

**Synchronously-Generated Signal**

Those signals that may be generated synchronously include SIGABRT, SIGBUS, SIGILL, SIGFPE, SIGPIPE, and SIGSEGV.

Any signal sent via the *raise*( ) function or a *kill*( ) function targeting the current process is also considered synchronous.

**System Call**\*

The distinction between a ''system call'' and a ''library routine'' is an implementation detail that may differ between implementations and has thus been excluded from POSIX.1.

See ''Interface, Not Implementation'' in **Introduction** (on page xiii).

1185    **System Reboot**

1186    A ''system reboot'' is an event initiated by an unspecified circumstance that causes all processes
1187    (other than special system processes) to be terminated in an implementation-defined manner,
1188    after which any changes to the state and contents of files created or written to by a Conforming
1189    POSIX.1 Application prior to the event are implementation-defined.

1190    **Synchronized I/O Data (and File) Integrity Completion**

1191    These terms specify that for synchronized read operations, pending writes must be successfully
1192    completed before the read operation can complete. This is motivated by two circumstances.
1193    Firstly, when synchronizing processes can access the same file, but not share common buffers
1194    (such as for a remote file system), this requirement permits the reading process to guarantee that
1195    it can read data written remotely. Secondly, having data written synchronously is insufficient to
1196    guarantee the order with respect to a subsequent write by a reading process, and thus this extra
1197    read semantic is necessary.

1198    **Text File**

1199    The term ''text file'' does not prevent the inclusion of control or other non-printable characters
1200    (other than NUL). Therefore, standard utilities that list text files as inputs or outputs are either
1201    able to process the special characters or they explicitly describe their limitations within their
1202    individual descriptions. The definition of ''text file'' has caused controversy. The only difference
1203    between text and binary files is that text files have lines of less than {LINE_MAX} bytes, with no
1204    NUL characters, each terminated by a <newline>. The definition allows a file with a single
1205    <newline>, but not a totally empty file, to be called a text file. If a file ends with an incomplete
1206    line it is not strictly a text file by this definition. The <newline> referred to in
1207    IEEE Std 1003.1-2001 is not some generic line separator, but a single character; files created on
1208    systems where they use multiple characters for ends of lines are not portable to all conforming
1209    systems without some translation process unspecified by IEEE Std 1003.1-2001.

1210    **Thread**

1211    IEEE Std 1003.1-2001 defines a thread to be a flow of control within a process. Each thread has a
1212    minimal amount of private state; most of the state associated with a process is shared among all
1213    of the threads in the process. While most multi-thread extensions to POSIX have taken this
1214    approach, others have made different decisions.

1215    **Note:**      The choice to put threads within a process does not constrain implementations to implement
1216              threads in that manner. However, all functions have to behave as though threads share the
1217              indicated state information with the process from which they were created.

1218    Threads need to share resources in order to cooperate. Memory has to be widely shared between
1219    threads in order for the threads to cooperate at a fine level of granularity. Threads keep data
1220    structures and the locks protecting those data structures in shared memory. For a data structure
1221    to be usefully shared between threads, such structures should not refer to any data that can only
1222    be interpreted meaningfully by a single thread. Thus, any system resources that might be
1223    referred to in data structures need to be shared between all threads. File descriptors, pathnames,
1224    and pointers to stack variables are all things that programmers want to share between their
1225    threads. Thus, the file descriptor table, the root directory, the current working directory, and the
1226    address space have to be shared.

1227    Library implementations are possible as long as the effective behavior is as if system services
1228    invoked by one thread do not suspend other threads. This may be difficult for some library
1229    implementations on systems that do not provide asynchronous facilities.

1230        See Section B.2.9 (on page 150) for additional rationale.

1231        **Thread ID**

1232        See Section B.2.9.2 (on page 166) for additional rationale.

1233        **Thread-Safe Function**

1234        All functions required by IEEE Std 1003.1-2001 need to be thread-safe; see Section A.4.16 (on
1235        page 40) and Section B.2.9.1 (on page 163) for additional rationale.

1236        **User Database**

1237        There are no references in IEEE Std 1003.1-2001 to a ''passwd file'' or a ''group file'', and there is
1238        no requirement that the *group* or *passwd* databases be kept in files containing editable text. Many
1239        large timesharing systems use *passwd* databases that are hashed for speed. Certain security
1240        classifications prohibit certain information in the *passwd* database from being publicly readable.

1241        The term ''encoded'' is used instead of ''encrypted'' in order to avoid the implementation
1242        connotations (such as reversibility or use of a particular algorithm) of the latter term.

1243        The *getgrent*( ), *setgrent*( ), *endgrent*( ), *getpwent*( ), *setpwent*( ), and *endpwent*( ) functions are not
1244        included as part of the base standard because they provide a linear database search capability
1245        that is not generally useful (the *getpwuid*( ), *getpwnam*( ), *getgrgid*( ), and *getgrnam*( ) functions are
1246        provided for keyed lookup) and because in certain distributed systems, especially those with
1247        different authentication domains, it may not be possible or desirable to provide an application
1248        with the ability to browse the system databases indiscriminately. They are provided on XSI-
1249        conformant systems due to their historical usage by many existing applications.

1250        A change from historical implementations is that the structures used by these functions have
1251        fields of the types **gid_t** and **uid_t**, which are required to be defined in the **<sys/types.h>** header.
1252        IEEE Std 1003.1-2001 requires implementations to ensure that these types are defined by
1253        inclusion of **<grp.h>** and **<pwd.h>**, respectively, without imposing any name space pollution or
1254        errors from redefinition of types.

1255        IEEE Std 1003.1-2001 is silent about the content of the strings containing user or group names.
1256        These could be digit strings. IEEE Std 1003.1-2001 is also silent as to whether such digit strings
1257        bear any relationship to the corresponding (numeric) user or group ID.

1258        *Database Access*

1259        The thread-safe versions of the user and group database access functions return values in user-
1260        supplied buffers instead of possibly using static data areas that may be overwritten by each call.

1261        **Virtual Processor***

1262        The term ''virtual processor'' was chosen as a neutral term describing all kernel-level
1263        schedulable entities, such as processes, Mach tasks, or lightweight processes. Implementing
1264        threads using multiple processes as virtual processors, or implementing multiplexed threads
1265        above a virtual processor layer, should be possible, provided some mechanism has also been
1266        implemented for sharing state between processes or virtual processors. Many systems may also
1267        wish to provide implementations of threads on systems providing ''shared processes'' or
1268        ''variable-weight processes''. It was felt that exposing such implementation details would
1269        severely limit the type of systems upon which the threads interface could be supported and
1270        prevent certain types of valid implementations. It was also determined that a virtual processor
1271        interface was out of the scope of the Rationale (Informative) volume of IEEE Std 1003.1-2001.

1272 **XSI**

1273 This is introduced to allow IEEE Std 1003.1-2001 to be adopted as an IEEE standard and an Open
1274 Group Technical Standard, serving both the POSIX and the Single UNIX Specification in a core
1275 set of volumes.

1276 The term ''XSI'' has been used for 10 years in connection with the XPG series and the first and
1277 second versions of the base volumes of the Single UNIX Specification. The XSI margin code was
1278 introduced to denote the extended or more restrictive semantics beyond POSIX that are
1279 applicable to UNIX systems.

1280 ## A.4    General Concepts

1281 ### A.4.1    Concurrent Execution

1282 There is no additional rationale provided for this section.

1283 ### A.4.2    Directory Protection

1284 There is no additional rationale provided for this section.

1285 ### A.4.3    Extended Security Controls

1286 Allowing an implementation to define extended security controls enables the use of
1287 IEEE Std 1003.1-2001 in environments that require different or more rigorous security than that
1288 provided in POSIX.1. Extensions are allowed in two areas: privilege and file access permissions.
1289 The semantics of these areas have been defined to permit extensions with reasonable, but not
1290 exact, compatibility with all existing practices. For example, the elimination of the superuser
1291 definition precludes identifying a process as privileged or not by virtue of its effective user ID.

1292 ### A.4.4    File Access Permissions

1293 A process should not try to anticipate the result of an attempt to access data by *a priori* use of
1294 these rules. Rather, it should make the attempt to access data and examine the return value (and
1295 possibly *errno* as well), or use *access*(). An implementation may include other security
1296 mechanisms in addition to those specified in POSIX.1, and an access attempt may fail because of
1297 those additional mechanisms, even though it would succeed according to the rules given in this
1298 section. (For example, the user's security level might be lower than that of the object of the access
1299 attempt.)  The supplementary group IDs provide another reason for a process to not attempt to
1300 anticipate the result of an access attempt.

1301 ### A.4.5    File Hierarchy

1302 Though the file hierarchy is commonly regarded to be a tree, POSIX.1 does not define it as such
1303 for three reasons:

1304 1.  Links may join branches.

1305 2.  In some network implementations, there may be no single absolute root directory; see
1306     *pathname resolution*.

1307 3.  With symbolic links, the file system need not be a tree or even a directed acyclic graph.

1308 **A.4.6    Filenames**

1309    Historically, certain filenames have been reserved. This list includes **core**, **/etc/passwd**, and so
1310    on. Conforming applications should avoid these.

1311    Most historical implementations prohibit case folding in filenames; that is, treating uppercase
1312    and lowercase alphabetic characters as identical. However, some consider case folding desirable:

1313        • For user convenience

1314        • For ease-of-implementation of the POSIX.1 interface as a hosted system on some popular
1315          operating systems

1316    Variants, such as maintaining case distinctions in filenames, but ignoring them in comparisons,
1317    have been suggested. Methods of allowing escaped characters of the case opposite the default
1318    have been proposed.

1319    Many reasons have been expressed for not allowing case folding, including:

1320        • No solid evidence has been produced as to whether case-sensitivity or case-insensitivity is
1321          more convenient for users.

1322        • Making case-insensitivity a POSIX.1 implementation option would be worse than either
1323          having it or not having it, because:

1324            — More confusion would be caused among users.

1325            — Application developers would have to account for both cases in their code.

1326            — POSIX.1 implementors would still have other problems with native file systems, such as
1327              short or otherwise constrained filenames or pathnames, and the lack of hierarchical
1328              directory structure.

1329        • Case folding is not easily defined in many European languages, both because many of them
1330          use characters outside the US ASCII alphabetic set, and because:

1331            — In Spanish, the digraph `"ll"` is considered to be a single letter, the capitalized form of
1332              which may be either `"Ll"` or `"LL"`, depending on context.

1333            — In French, the capitalized form of a letter with an accent may or may not retain the accent,
1334              depending on the country in which it is written.

1335            — In German, the sharp ess may be represented as a single character resembling a Greek
1336              beta (β) in lowercase, but as the digraph `"SS"` in uppercase.

1337            — In Greek, there are several lowercase forms of some letters; the one to use depends on its
1338              position in the word. Arabic has similar rules.

1339        • Many East Asian languages, including Japanese, Chinese, and Korean, do not distinguish
1340          case and are sometimes encoded in character sets that use more than one byte per character.

1341        • Multiple character codes may be used on the same machine simultaneously. There are
1342          several ISO character sets for European alphabets. In Japan, several Japanese character codes
1343          are commonly used together, sometimes even in filenames; this is evidently also the case in
1344          China. To handle case insensitivity, the kernel would have to at least be able to distinguish
1345          for which character sets the concept made sense.

1346        • The file system implementation historically deals only with bytes, not with characters, except
1347          for slash and the null byte.

1348        • The purpose of POSIX.1 is to standardize the common, existing definition, not to change it.
1349          Mandating case-insensitivity would make all historical implementations non-standard.

1350 • Not only the interface, but also application programs would need to change, counter to the
1351 purpose of having minimal changes to existing application code.

1352 • At least one of the original developers of the UNIX system has expressed objection in the
1353 strongest terms to either requiring case-insensitivity or making it an option, mostly on the
1354 basis that POSIX.1 should not hinder portability of application programs across related
1355 implementations in order to allow compatibility with unrelated operating systems.

1356 Two proposals were entertained regarding case folding in filenames:

1357 1. Remove all wording that previously permitted case folding.

1358 Rationale     Case folding is inconsistent with portable filename character set definition
1359 and filename definition (all characters except slash and null). No known
1360 implementations allowing all characters except slash and null also do case
1361 folding.

1362 2. Change ''though this practice is not recommended:'' to ''although this practice is strongly
1363 discouraged.''

1364 Rationale     If case folding must be included in POSIX.1, the wording should be stronger
1365 to discourage the practice.

1366 The consensus selected the first proposal. Otherwise, a conforming application would have to
1367 assume that case folding would occur when it was not wanted, but that it would not occur when
1368 it was wanted.

## A.4.7 File Times Update

1370 This section reflects the actions of historical implementations. The times are not updated
1371 immediately, but are only marked for update by the functions. An implementation may update
1372 these times immediately.

1373 The accuracy of the time update values is intentionally left unspecified so that systems can
1374 control the bandwidth of a possible covert channel.

1375 The wording was carefully chosen to make it clear that there is no requirement that the
1376 conformance document contain information that might incidentally affect file update times. Any
1377 function that performs pathname resolution might update several *st_atime* fields. Functions such
1378 as *getpwnam*( ) and *getgrnam*( ) might update the *st_atime* field of some specific file or files. It is
1379 intended that these are not required to be documented in the conformance document, but they
1380 should appear in the system documentation.

## A.4.8 Host and Network Byte Order

1382 There is no additional rationale provided for this section.

## A.4.9 Measurement of Execution Time

1384 The methods used to measure the execution time of processes and threads, and the precision of
1385 these measurements, may vary considerably depending on the software architecture of the
1386 implementation, and on the underlying hardware. Implementations can also make tradeoffs
1387 between the scheduling overhead and the precision of the execution time measurements.
1388 IEEE Std 1003.1-2001 does not impose any requirement on the accuracy of the execution time; it
1389 instead specifies that the measurement mechanism and its precision are implementation-
1390 defined.

1391 **A.4.10   Memory Synchronization**

1392   In older multi-processors, access to memory by the processors was strictly multiplexed. This
1393   meant that a processor executing program code interrogates or modifies memory in the order
1394   specified by the code and that all the memory operation of all the processors in the system
1395   appear to happen in some global order, though the operation histories of different processors are
1396   interleaved arbitrarily. The memory operations of such machines are said to be sequentially
1397   consistent. In this environment, threads can synchronize using ordinary memory operations. For
1398   example, a producer thread and a consumer thread can synchronize access to a circular data
1399   buffer as follows:

```
1400        int rdptr = 0;
1401        int wrptr = 0;
1402        data_t buf[BUFSIZE];

1403        Thread 1:
1404            while (work_to_do) {
1405                int next;

1406                buf[wrptr] = produce();
1407                next = (wrptr + 1) % BUFSIZE;
1408                while (rdptr == next)
1409                    ;
1410                wrptr = next;
1411        }

1412        Thread 2:
1413            while (work_to_do) {
1414                while (rdptr == wrptr)
1415                    ;
1416                consume(buf[rdptr]);
1417                rdptr = (rdptr + 1) % BUFSIZE;
1418            }
```

1419   In modern multi-processors, these conditions are relaxed to achieve greater performance. If one
1420   processor stores values in location A and then location B, then other processors loading data
1421   from location B and then location A may see the new value of B but the old value of A. The
1422   memory operations of such machines are said to be weakly ordered. On these machines, the
1423   circular buffer technique shown in the example will fail because the consumer may see the new
1424   value of *wrptr* but the old value of the data in the buffer. In such machines, synchronization can
1425   only be achieved through the use of special instructions that enforce an order on memory
1426   operations. Most high-level language compilers only generate ordinary memory operations to
1427   take advantage of the increased performance. They usually cannot determine when memory
1428   operation order is important and generate the special ordering instructions. Instead, they rely on
1429   the programmer to use synchronization primitives correctly to ensure that modifications to a
1430   location in memory are ordered with respect to modifications and/or access to the same location
1431   in other threads. Access to read-only data need not be synchronized. The resulting program is
1432   said to be data race-free.

1433   Synchronization is still important even when accessing a single primitive variable (for example,
1434   an integer). On machines where the integer may not be aligned to the bus data width or be larger
1435   than the data width, a single memory load may require multiple memory cycles.  This means
1436   that it may be possible for some parts of the integer to have an old value while other parts have a
1437   newer value. On some processor architectures this cannot happen, but portable programs cannot
1438   rely on this.

1439      In summary, a portable multi-threaded program, or a multi-process program that shares
1440      writable memory between processes, has to use the synchronization primitives to synchronize
1441      data access. It cannot rely on modifications to memory being observed by other threads in the
1442      order written in the application or even on modification of a single variable being seen
1443      atomically.

1444      Conforming applications may only use the functions listed to synchronize threads of control
1445      with respect to memory access. There are many other candidates for functions that might also be
1446      used. Examples are: signal sending and reception, or pipe writing and reading. In general, any
1447      function that allows one thread of control to wait for an action caused by another thread of
1448      control is a candidate. IEEE Std 1003.1-2001 does not require these additional functions to
1449      synchronize memory access since this would imply the following:

1450      • All these functions would have to be recognized by advanced compilation systems so that
1451        memory operations and calls to these functions are not reordered by optimization.

1452      • All these functions would potentially have to have memory synchronization instructions
1453        added, depending on the particular machine.

1454      • The additional functions complicate the model of how memory is synchronized and make
1455        automatic data race detection techniques impractical.

1456      Formal definitions of the memory model were rejected as unreadable by the vast majority of
1457      programmers. In addition, most of the formal work in the literature has concentrated on the
1458      memory as provided by the hardware as opposed to the application programmer through the
1459      compiler and runtime system. It was believed that a simple statement intuitive to most
1460      programmers would be most effective. IEEE Std 1003.1-2001 defines functions that can be used
1461      to synchronize access to memory, but it leaves open exactly how one relates those functions to
1462      the semantics of each function as specified elsewhere in IEEE Std 1003.1-2001.
1463      IEEE Std 1003.1-2001 also does not make a formal specification of the partial ordering in time
1464      that the functions can impose, as that is implied in the description of the semantics of each
1465      function. It simply states that the programmer has to ensure that modifications do not occur
1466      ''simultaneously'' with other access to a memory location.

1467      IEEE Std 1003.1-2001/Cor 1-2002, item XBD/TC1/D6/4 is applied, adding a new paragraph    |
1468      beneath the table of functions: ''The *pthread_once*() function shall synchronize memory for the  |
1469      first call in each thread for a given **pthread_once_t** object.''.                             |

## 1470      A.4.11    Pathname Resolution

1471      It is necessary to differentiate between the definition of pathname and the concept of pathname
1472      resolution with respect to the handling of trailing slashes. By specifying the behavior here, it is
1473      not possible to provide an implementation that is conforming but extends all interfaces that
1474      handle pathnames to also handle strings that are not legal pathnames (because they have trailing
1475      slashes).

1476      Pathnames that end with one or more trailing slash characters must refer to directory paths.
1477      Previous versions of IEEE Std 1003.1-2001 were not specific about the distinction between
1478      trailing slashes on files and directories, and both were permitted.

1479      Two types of implementation have been prevalent; those that ignored trailing slash characters
1480      on all pathnames regardless, and those that permitted them only on existing directories.

1481      IEEE Std 1003.1-2001 requires that a pathname with a trailing slash character be treated as if it
1482      had a trailing "/." everywhere.

1483      Note that this change does not break any conforming applications; since there were two
1484      different types of implementation, no application could have portably depended on either

1485   behavior. This change does however require some implementations to be altered to remain
1486   compliant. Substantial discussion over a three-year period has shown that the benefits to
1487   application developers outweighs the disadvantages for some vendors.

1488   On a historical note, some early applications automatically appended a ′/′ to every path.
1489   Rather than fix the applications, the system implementation was modified to accept this
1490   behavior by ignoring any trailing slash.

1491   Each directory has exactly one parent directory which is represented by the name **dot-dot** in the
1492   first directory. No other directory, regardless of linkages established by symbolic links, is
1493   considered the parent directory by IEEE Std 1003.1-2001.

1494   There are two general categories of interfaces involving pathname resolution: those that follow
1495   the symbolic link, and those that do not. There are several exceptions to this rule; for example,
1496   *open*(*path*,O_CREAT | O_EXCL) will fail when *path* names a symbolic link. However, in all other
1497   situations, the *open*( ) function will follow the link.

1498   What the filename **dot-dot** refers to relative to the root directory is implementation-defined. In
1499   Version 7 it refers to the root directory itself; this is the behavior mentioned in
1500   IEEE Std 1003.1-2001. In some networked systems the construction **/. . /hostname/** is used to refer
1501   to the root directory of another host, and POSIX.1 permits this behavior.

1502   Other networked systems use the construct **//hostname** for the same purpose; that is, a double
1503   initial slash is used. There is a potential problem with existing applications that create full
1504   pathnames by taking a trunk and a relative pathname and making them into a single string
1505   separated by ′/′, because they can accidentally create networked pathnames when the trunk is
1506   ′/′. This practice is not prohibited because such applications can be made to conform by
1507   simply changing to use "//" as a separator instead of ′/′:

1508   • If the trunk is ′/′, the full pathname will begin with "///" (the initial ′/′ and the
1509       separator "//"). This is the same as ′/′, which is what is desired. (This is the general case
1510       of making a relative pathname into an absolute one by prefixing with "///" instead of ′/′.)

1511   • If the trunk is "/A", the result is "/A//. . ."; since non-leading sequences of two or more
1512       slashes are treated as a single slash, this is equivalent to the desired "/A/. . .".

1513   • If the trunk is "//A", the implementation-defined semantics will apply. (The multiple slash
1514       rule would apply.)

1515   Application developers should avoid generating pathnames that start with "//".
1516   Implementations are strongly encouraged to avoid using this special interpretation since a
1517   number of applications currently do not follow this practice and may inadvertently generate
1518   "//. . .".

1519   The term ''root directory'' is only defined in POSIX.1 relative to the process. In some
1520   implementations, there may be no absolute root directory. The initialization of the root directory
1521   of a process is implementation-defined.

### 1522 **A.4.12  Process ID Reuse**

1523    There is no additional rationale provided for this section.

### 1524 **A.4.13  Scheduling Policy**

1525    There is no additional rationale provided for this section.

### 1526 **A.4.14  Seconds Since the Epoch**

1527    Coordinated Universal Time (UTC) includes leap seconds. However, in POSIX time (seconds
1528    since the Epoch), leap seconds are ignored (not applied) to provide an easy and compatible
1529    method of computing time differences. Broken-down POSIX time is therefore not necessarily
1530    UTC, despite its appearance.

1531    As of September 2000, 24 leap seconds had been added to UTC since the Epoch, 1 January, 1970.
1532    Historically, one leap second is added every 15 months on average, so this offset can be expected
1533    to grow steadily with time.

1534    Most systems' notion of ''time'' is that of a continuously increasing value, so this value should
1535    increase even during leap seconds.  However, not only do most systems not keep track of leap
1536    seconds, but most systems are probably not synchronized to any standard time reference.
1537    Therefore, it is inappropriate to require that a time represented as seconds since the Epoch
1538    precisely represent the number of seconds between the referenced time and the Epoch.

1539    It is sufficient to require that applications be allowed to treat this time as if it represented the
1540    number of seconds between the referenced time and the Epoch. It is the responsibility of the
1541    vendor of the system, and the administrator of the system, to ensure that this value represents
1542    the number of seconds between the referenced time and the Epoch as closely as necessary for the
1543    application being run on that system.

1544    It is important that the interpretation of time names and seconds since the Epoch values be
1545    consistent across conforming systems; that is, it is important that all conforming systems
1546    interpret ''536 457 599 seconds since the Epoch'' as 59 seconds, 59 minutes, 23 hours 31 December
1547    1986, regardless of the accuracy of the system's idea of the current time.  The expression is given
1548    to ensure a consistent interpretation, not to attempt to specify the calendar. The relationship
1549    between *tm_yday* and the day of week, day of month, and month is in accordance with the
1550    Gregorian calendar, and so is not specified in POSIX.1.

1551    Consistent interpretation of seconds since the Epoch can be critical to certain types of distributed
1552    applications that rely on such timestamps to synchronize events. The accrual of leap seconds in
1553    a time standard is not predictable. The number of leap seconds since the Epoch will likely
1554    increase. POSIX.1 is more concerned about the synchronization of time between applications of
1555    astronomically short duration.

1556    Note that *tm_yday* is zero-based, not one-based, so the day number in the example above is 364.
1557    Note also that the division is an integer division (discarding remainder) as in the C language.

1558    Note also that the meaning of *gmtime*( ), *localtime*( ), and *mktime*( ) is specified in terms of this
1559    expression. However, the ISO C standard computes *tm_yday* from *tm_mday*, *tm_mon*, and
1560    *tm_year* in *mktime*( ).  Because it is stated as a (bidirectional) relationship, not a function, and
1561    because the conversion between month-day-year and day-of-year dates is presumed well known
1562    and is also a relationship, this is not a problem.

1563    Implementations that implement **time_t** as a signed 32-bit integer will overflow in 2 038. The
1564    data size for **time_t** is as per the ISO C standard definition, which is implementation-defined.

1565        See also **Epoch** (on page 17).

1566        The topic of whether seconds since the Epoch should account for leap seconds has been debated
1567        on a number of occasions, and each time consensus was reached (with acknowledged dissent
1568        each time) that the majority of users are best served by treating all days identically. (That is, the
1569        majority of applications were judged to assume a single length—as measured in seconds since
1570        the Epoch—for all days. Thus, leap seconds are not applied to seconds since the Epoch.) Those
1571        applications which do care about leap seconds can determine how to handle them in whatever
1572        way those applications feel is best. This was particularly emphasized because there was
1573        disagreement about what the best way of handling leap seconds might be. It is a practical
1574        impossibility to mandate that a conforming implementation must have a fixed relationship to
1575        any particular official clock (consider isolated systems, or systems performing ''reruns'' by
1576        setting the clock to some arbitrary time).

1577        Note that as a practical consequence of this, the length of a second as measured by some external
1578        standard is not specified. This unspecified second is nominally equal to an International System
1579        (SI) second in duration. Applications must be matched to a system that provides the particular
1580        handling of external time in the way required by the application.

## 1581   A.4.15   Semaphore

1582        There is no additional rationale provided for this section.

## 1583   A.4.16   Thread-Safety

1584        Where the interface of a function required by IEEE Std 1003.1-2001 precludes thread-safety, an
1585        alternate thread-safe form is provided. The names of these thread-safe forms are the same as the
1586        non-thread-safe forms with the addition of the suffix ''_r''. The suffix ''_r'' is historical, where
1587        the `'r'` stood for ''reentrant''.

1588        In some cases, thread-safety is provided by restricting the arguments to an existing function.

1589        See also Section B.2.9.1 (on page 163).

## 1590   A.4.17   Tracing

1591        Refer to Section B.2.11 (on page 179).

## 1592   A.4.18   Treatment of Error Conditions for Mathematical Functions

1593        There is no additional rationale provided for this section.

## 1594   A.4.19   Treatment of NaN Arguments for Mathematical Functions

1595        There is no additional rationale provided for this section.

## 1596   A.4.20   Utility

1597        There is no additional rationale provided for this section.

1598 **A.4.21  Variable Assignment**

1599     There is no additional rationale provided for this section.

1600 # A.5    **File Format Notation**

1601     The notation for spaces allows some flexibility for application output. Note that an empty
1602     character position in *format* represents one or more <blank>s on the output (not *white space*,
1603     which can include <newline>s). Therefore, another utility that reads that output as its input
1604     must be prepared to parse the data using *scanf*(), *awk*, and so on. The '∆' character is used when
1605     exactly one <space> is output.

1606     The treatment of integers and spaces is different from the *printf*() function in that they can be
1607     surrounded with <blank>s. This was done so that, given a format such as:

1608         `"%d\n",<foo>`

1609     the implementation could use a *printf*() call such as:

1610         `printf("%6d\n", foo);`

1611     and still conform. This notation is thus somewhat like *scanf*() in addition to *printf*().

1612     The *printf*() function was chosen as a model because most of the standard developers were
1613     familiar with it. One difference from the C function *printf*() is that the `l` and `h` conversion
1614     specifier characters are not used. As expressed by the Shell and Utilities volume of
1615     IEEE Std 1003.1-2001, there is no differentiation between decimal values for type **int**, type **long**,
1616     or type **short**. The conversion specifications `%d` or `%i` should be interpreted as an arbitrary
1617     length sequence of digits. Also, no distinction is made between single precision and double
1618     precision numbers (**float** or **double** in C). These are simply referred to as floating-point numbers.

1619     Many of the output descriptions in the Shell and Utilities volume of IEEE Std 1003.1-2001 use the
1620     term ''line'', such as:

1621         `"%s", <input line>`

1622     Since the definition of *line* includes the trailing <newline> already, there is no need to include a
1623     '\n' in the format; a double <newline> would otherwise result.

1624 # A.6    **Character Set**

1625 **A.6.1    Portable Character Set**

1626     The portable character set is listed in full so there is no dependency on the ISO/IEC 646:1991
1627     standard (or historically ASCII) encoded character set, although the set is identical to the
1628     characters defined in the International Reference version of the ISO/IEC 646:1991 standard.

1629     IEEE Std 1003.1-2001 poses no requirement that multiple character sets or codesets be supported,
1630     leaving this as a marketing differentiation for implementors. Although multiple charmap files
1631     are supported, it is the responsibility of the implementation to provide the file(s); if only one is
1632     provided, only that one will be accessible using the *localedef* **−f** option.

1633     The statement about invariance in codesets for the portable character set is worded to avoid
1634     precluding implementations where multiple incompatible codesets are available (for instance,
1635     ASCII and EBCDIC). The standard utilities cannot be expected to produce predictable results if
1636     they access portable characters that vary on the same implementation.

1637  Not all character sets need include the portable character set, but each locale must include it. For
1638  example, a Japanese-based locale might be supported by a mixture of character sets: JIS X 0201
1639  Roman (a Japanese version of the ISO/IEC 646:1991 standard), JIS X 0208, and JIS X 0201
1640  Katakana. Not all of these character sets include the portable characters, but at least one does
1641  (JIS X 0201 Roman).

## A.6.2  Character Encoding

1643  Encoding mechanisms based on single shifts, such as the EUC encoding used in some Asian and
1644  other countries, can be supported via the current charmap mechanism. With single-shift
1645  encoding, each character is preceded by a shift code (SS2 or SS3). A complete EUC code,
1646  consisting of the portable character set (G0) and up to three additional character sets (G1, G2,
1647  G3), can be described using the current charmap mechanism; the encoding for each character in
1648  additional character sets G2 and G3 must then include their single-shift code. Other mechanisms
1649  to support locales based on encoding mechanisms such as locking shift are not addressed by this
1650  volume of IEEE Std 1003.1-2001.

## A.6.3  C Language Wide-Character Codes

1652  There is no additional rationale provided for this section.

## A.6.4  Character Set Description File

1654  IEEE PASC Interpretation 1003.2 #196 is applied, removing three lines of text dealing with
1655  ranges of symbolic names using position constant values which had been erroneously included
1656  in the final IEEE P1003.2b draft standard.

### A.6.4.1  *State-Dependent Character Encodings*

1658  A requirement was considered that would force utilities to eliminate any redundant locking
1659  shifts, but this was left as a quality of implementation issue.

1660  This change satisfies the following requirement from the ISO POSIX-2:1993 standard, Annex
1661  H.1:

1662      *The support of state-dependent (shift encoding) character sets should be addressed fully. See*
1663      *descriptions of these in the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.2, Character*
1664      *Encoding. If such character encodings are supported, it is expected that this will impact the Base*
1665      *Definitions volume of IEEE Std 1003.1-2001, Section 6.2, Character Encoding, the Base Definitions*
1666      *volume of IEEE Std 1003.1-2001, Chapter 7, Locale, the Base Definitions volume of*
1667      *IEEE Std 1003.1-2001, Chapter 9, Regular Expressions , and the comm, cut, diff, grep, head, join,*
1668      *paste, and tail utilities.*

1669  The character set description file provides:

1670  • The capability to describe character set attributes (such as collation order or character
1671     classes) independent of character set encoding, and using only the characters in the portable
1672     character set. This makes it possible to create generic *localedef* source files for all codesets that
1673     share the portable character set (such as the ISO 8859 family or IBM Extended ASCII).

1674  • Standardized symbolic names for all characters in the portable character set, making it
1675     possible to refer to any such character regardless of encoding.

1676  Implementations are free to choose their own symbolic names, as long as the names identified
1677  by the Base Definitions volume of IEEE Std 1003.1-2001 are also defined; this provides support
1678  for already existing ''character names''.

1679 The names selected for the members of the portable character set follow the
1680 ISO/IEC 8859-1:1998 standard and the ISO/IEC 10646-1:2000 standard. However, several
1681 commonly used UNIX system names occur as synonyms in the list:

1682 • The historical UNIX system names are used for control characters.

1683 • The word ''slash'' is given in addition to ''solidus''.

1684 • The word ''backslash'' is given in addition to ''reverse-solidus''.

1685 • The word ''hyphen'' is given in addition to ''hyphen-minus''.

1686 • The word ''period'' is given in addition to ''full-stop''.

1687 • For digits, the word ''digit'' is eliminated.

1688 • For letters, the words ''Latin Capital Letter'' and ''Latin Small Letter'' are eliminated.

1689 • The words ''left brace'' and ''right brace'' are given in addition to ''left-curly-bracket'' and
1690 ''right-curly-bracket''.

1691 • The names of the digits are preferred over the numbers to avoid possible confusion between
1692 '0' and 'O', and between '1' and 'l' (one and the letter ell).

1693 The names for the control characters in the Base Definitions volume of IEEE Std 1003.1-2001,
1694 Chapter 6, Character Set were taken from the ISO/IEC 4873:1991 standard.

1695 The charmap file was introduced to resolve problems with the portability of, especially, *localedef*
1696 sources. IEEE Std 1003.1-2001 assumes that the portable character set is constant across all
1697 locales, but does not prohibit implementations from supporting two incompatible codings, such
1698 as both ASCII and EBCDIC. Such dual-support implementations should have all charmaps and
1699 *localedef* sources encoded using one portable character set, in effect cross-compiling for the other
1700 environment. Naturally, charmaps (and *localedef* sources) are only portable without
1701 transformation between systems using the same encodings for the portable character set. They
1702 can, however, be transformed between two sets using only a subset of the actual characters (the
1703 portable character set). However, the particular coded character set used for an application or an
1704 implementation does not necessarily imply different characteristics or collation; on the contrary,
1705 these attributes should in many cases be identical, regardless of codeset. The charmap provides
1706 the capability to define a common locale definition for multiple codesets (the same *localedef*
1707 source can be used for codesets with different extended characters; the ability in the charmap to
1708 define empty names allows for characters missing in certain codesets).

1709 The **<escape_char>** declaration was added at the request of the international community to ease
1710 the creation of portable charmap files on terminals not implementing the default backslash
1711 escape. The **<comment_char>** declaration was added at the request of the international
1712 community to eliminate the potential confusion between the number sign and the pound sign.

1713 The octal number notation with no leading zero required was selected to match those of *awk* and
1714 *tr* and is consistent with that used by *localedef*. To avoid confusion between an octal constant
1715 and the back-references used in *localedef* source, the octal, hexadecimal, and decimal constants
1716 must contain at least two digits. As single-digit constants are relatively rare, this should not
1717 impose any significant hardship. Provision is made for more digits to account for systems in
1718 which the byte size is larger than 8 bits. For example, a Unicode (ISO/IEC 10646-1:2000
1719 standard) system that has defined 16-bit bytes may require six octal, four hexadecimal, and five
1720 decimal digits.

1721 The decimal notation is supported because some newer international standards define character
1722 values in decimal, rather than in the old column/row notation.

1723 The charmap identifies the coded character sets supported by an implementation. At least one
1724 charmap must be provided, but no implementation is required to provide more than one.
1725 Likewise, implementations can allow users to generate new charmaps (for instance, for a new
1726 version of the ISO 8859 family of coded character sets), but does not have to do so. If users are
1727 allowed to create new charmaps, the system documentation describes the rules that apply (for
1728 instance, ''only coded character sets that are supersets of the ISO/IEC 646:1991 standard IRV, no
1729 multi-byte characters'').

1730 This addition of the **WIDTH** specification satisfies the following requirement from the
1731 ISO POSIX-2:1993 standard, Annex H.1:

1732 (9) *The definition of column position relies on the implementation's knowledge of the integral width*
1733 *of the characters. The charmap or LC_CTYPE locale definitions should be enhanced to allow*
1734 *application specification of these widths.*

1735 The character ''width'' information was first considered for inclusion under *LC_CTYPE* but was
1736 moved because it is more closely associated with the information in the charmap than
1737 information in the locale source (cultural conventions information). Concerns were raised that
1738 formalizing this type of information is moving the locale source definition from the codeset-
1739 independent entity that it was designed to be to a repository of codeset-specific information. A
1740 similar issue occurred with the **<code_set_name>**, **<mb_cur_max>**, and **<mb_cur_min>**
1741 information, which was resolved to reside in the charmap definition.

1742 The width definition was added to the IEEE P1003.2b draft standard with the intent that the
1743 *wcswidth*() and/or *wcwidth*() functions (currently specified in the System Interfaces volume of
1744 IEEE Std 1003.1-2001) be the mechanism to retrieve the character width information.

## 1745 A.7    Locale

### 1746 A.7.1    General

1747 The description of locales is based on work performed in the UniForum Technical Committee,
1748 Subcommittee on Internationalization. Wherever appropriate, keywords are taken from the
1749 ISO C standard or the X/Open Portability Guide.

1750 The value used to specify a locale with environment variables is the name specified as the *name*
1751 operand to the *localedef* utility when the locale was created. This provides a verifiable method to
1752 create and invoke a locale.

1753 The ''object'' definitions need not be portable, as long as ''source'' definitions are. Strictly
1754 speaking, source definitions are portable only between implementations using the same
1755 character set(s). Such source definitions, if they use symbolic names only, easily can be ported
1756 between systems using different codesets, as long as the characters in the portable character set
1757 (see the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.1, Portable Character Set)
1758 have common values between the codesets; this is frequently the case in historical
1759 implementations. Of source, this requires that the symbolic names used for characters outside
1760 the portable character set be identical between character sets. The definition of symbolic names
1761 for characters is outside the scope of IEEE Std 1003.1-2001, but is certainly within the scope of
1762 other standards organizations.

1763 Applications can select the desired locale by invoking the *setlocale*() function (or equivalent)
1764 with the appropriate value. If the function is invoked with an empty string, the value of the
1765 corresponding environment variable is used. If the environment variable is not set or is set to the
1766 empty string, the implementation sets the appropriate environment as defined in the Base

1767    Definitions volume of IEEE Std 1003.1-2001, Chapter 8, Environment Variables.

## 1768  A.7.2  POSIX Locale

1769    The POSIX locale is equal to the C locale. To avoid being classified as a C-language function, the
1770    name has been changed to the POSIX locale; the environment variable value can be either
1771    `"POSIX"` or, for historical reasons, `"C"`.

1772    The POSIX definitions mirror the historical UNIX system behavior.

1773    The use of symbolic names for characters in the tables does not imply that the POSIX locale must
1774    be described using symbolic character names, but merely that it may be advantageous to do so.

## 1775  A.7.3  Locale Definition

1776    The decision to separate the file format from the *localedef* utility description was only partially
1777    editorial. Implementations may provide other interfaces than *localedef*. Requirements on ''the
1778    utility'', mostly concerning error messages, are described in this way because they are meant to
1779    affect the other interfaces implementations may provide as well as *localedef*.

1780    The text about POSIX2_LOCALEDEF does not mean that internationalization is optional; only
1781    that the functionality of the *localedef* utility is. REs, for instance, must still be able to recognize,
1782    for example, character class expressions such as `"[[:alpha:]]"`. A possible analogy is with
1783    an applications development environment; while all conforming implementations must be
1784    capable of executing applications, not all need to have the development environment installed.
1785    The assumption is that the capability to modify the behavior of utilities (and applications) via
1786    locale settings must be supported. If the *localedef* utility is not present, then the only choice is to
1787    select an existing (presumably implementation-documented) locale. An implementation could,
1788    for example, choose to support only the POSIX locale, which would in effect limit the amount of
1789    changes from historical implementations quite drastically. The *localedef* utility is still required,
1790    but would always terminate with an exit code indicating that no locale could be created.
1791    Supported locales must be documented using the syntax defined in this chapter. (This ensures
1792    that users can accurately determine what capabilities are provided. If the implementation
1793    decides to provide additional capabilities to the ones in this chapter, that is already provided
1794    for.)

1795    If the option is present (that is, locales can be created), then the *localedef* utility must be capable
1796    of creating locales based on the syntax and rules defined in this chapter. This does not mean that
1797    the implementation cannot also provide alternate means for creating locales.

1798    The octal, decimal, and hexadecimal notations are the same employed by the charmap facility
1799    (see the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.4, Character Set Description
1800    File). To avoid confusion between an octal constant and a back-reference, the octal, hexadecimal,
1801    and decimal constants must contain at least two digits. As single-digit constants are relatively
1802    rare, this should not impose any significant hardship. Provision is made for more digits to
1803    account for systems in which the byte size is larger than 8 bits. For example, a Unicode (see the
1804    ISO/IEC 10646-1:2000 standard) system that has defined 16-bit bytes may require six octal, four
1805    hexadecimal, and five decimal digits. As with the charmap file, multi-byte characters are
1806    described in the locale definition file using ''big-endian'' notation for reasons of portability.
1807    There is no requirement that the internal representation in the computer memory be in this same
1808    order.

1809    One of the guidelines used for the development of this volume of IEEE Std 1003.1-2001 is that
1810    characters outside the invariant part of the ISO/IEC 646:1991 standard should not be used in
1811    portable specifications. The backslash character is not in the invariant part; the number sign is,
1812    but with multiple representations: as a number sign, and as a pound sign. As far as general

1813   usage of these symbols, they are covered by the ''grandfather clause'', but for newly defined
1814   interfaces, the WG15 POSIX working group has requested that POSIX provide alternate
1815   representations. Consequently, while the default escape character remains the backslash and the
1816   default comment character is the number sign, implementations are required to recognize
1817   alternative representations, identified in the applicable source file via the **<escape_char>** and
1818   **<comment_char>** keywords.

## A.7.3.1   LC_CTYPE

1820   The *LC_CTYPE* category is primarily used to define the encoding-independent aspects of a
1821   character set, such as character classification. In addition, certain encoding-dependent
1822   characteristics are also defined for an application via the *LC_CTYPE* category.
1823   IEEE Std 1003.1-2001 does not mandate that the encoding used in the locale is the same as the
1824   one used by the application because an implementation may decide that it is advantageous to
1825   define locales in a system-wide encoding rather than having multiple, logically identical locales
1826   in different encodings, and to convert from the application encoding to the system-wide
1827   encoding on usage. Other implementations could require encoding-dependent locales.

1828   In either case, the *LC_CTYPE* attributes that are directly dependent on the encoding, such as
1829   **<mb_cur_max>** and the display width of characters, are not user-specifiable in a locale source
1830   and are consequently not defined as keywords.

1831   Implementations may define additional keywords or extend the *LC_CTYPE* mechanism to allow
1832   application-defined keywords.

1833   The text ''The ellipsis specification shall only be valid within a single encoded character set'' is
1834   present because it is possible to have a locale supported by multiple character encodings, as
1835   explained in the rationale for the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.1,
1836   Portable Character Set. An example given there is of a possible Japanese-based locale supported
1837   by a mixture of the character sets JIS X 0201 Roman, JIS X 0208, and JIS X 0201 Katakana.
1838   Attempting to express a range of characters across these sets is not logical and the
1839   implementation is free to reject such attempts.

1840   As the *LC_CTYPE* character classes are based on the ISO C standard character class definition,
1841   the category does not support multi-character elements. For instance, the German character
1842   <sharp-s> is traditionally classified as a lowercase letter. There is no corresponding uppercase
1843   letter; in proper capitalization of German text, the <sharp-s> will be replaced by `"SS"`; that is, by
1844   two characters. This kind of conversion is outside the scope of the **toupper** and **tolower**
1845   keywords.

1846   Where IEEE Std 1003.1-2001 specifies that only certain characters can be specified, as for the
1847   keywords **digit** and **xdigit**, the specified characters must be from the portable character set, as
1848   shown. As an example, only the Arabic digits 0 through 9 are acceptable as digits.

1849   The character classes **digit**, **xdigit**, **lower**, **upper**, and **space** have a set of automatically included
1850   characters. These only need to be specified if the character values (that is, encoding) differs from
1851   the implementation default values. It is not possible to define a locale without these
1852   automatically included characters unless some implementation extension is used to prevent
1853   their inclusion. Such a definition would not be a proper superset of the C locale, and thus, it
1854   might not be possible for the standard utilities to be implemented as programs conforming to
1855   the ISO C standard.

1856   The definition of character class **digit** requires that only ten characters—the ones defining
1857   digits—can be specified; alternate digits (for example, Hindi or Kanji) cannot be specified here.
1858   However, the encoding may vary if an implementation supports more than one encoding.

1859 The definition of character class **xdigit** requires that the characters included in character class
1860 **digit** are included here also and allows for different symbols for the hexadecimal digits 10
1861 through 15.

1862 The inclusion of the **charclass** keyword satisfies the following requirement from the
1863 ISO POSIX-2: 1993 standard, Annex H.1:

1864 (3) *The LC_CTYPE (2.5.2.1) locale definition should be enhanced to allow user-specified additional*
1865 *character classes, similar in concept to the ISO C standard Multibyte Support Extension (MSE)*
1866 *iswctype( ) function.*

1867 This keyword was previously included in The Open Group specifications and is now mandated
1868 in the Shell and Utilities volume of IEEE Std 1003.1-2001.

1869 The symbolic constant {CHARCLASS_NAME_MAX} was also adopted from The Open Group
1870 specifications. Applications portability is enhanced by the use of symbolic constants.

1871 *A.7.3.2 LC_COLLATE*

1872 The rules governing collation depend to some extent on the use. At least five different levels of
1873 increasingly complex collation rules can be distinguished:

1874 1. *Byte/machine code order*: This is the historical collation order in the UNIX system and many
1875 proprietary operating systems. Collation is here performed character by character, without
1876 any regard to context. The primary virtue is that it usually is quite fast and also
1877 completely deterministic; it works well when the native machine collation sequence
1878 matches the user expectations.

1879 2. *Character order*: On this level, collation is also performed character by character, without
1880 regard to context. The order between characters is, however, not determined by the code
1881 values, but on the expectations by the user of the ''correct'' order between characters. In
1882 addition, such a (simple) collation order can specify that certain characters collate equally
1883 (for example, uppercase and lowercase letters).

1884 3. *String ordering*: On this level, entire strings are compared based on relatively
1885 straightforward rules. Several ''passes'' may be required to determine the order between
1886 two strings. Characters may be ignored in some passes, but not in others; the strings may
1887 be compared in different directions; and simple string substitutions may be performed
1888 before strings are compared. This level is best described as ''dictionary'' ordering; it is
1889 based on the spelling, not the pronunciation, or meaning, of the words.

1890 4. *Text search ordering*: This is a further refinement of the previous level, best described as
1891 ''telephone book ordering''; some common homonyms (words spelled differently but with
1892 the same pronunciation) are collated together; numbers are collated as if they were spelled
1893 out, and so on.

1894 5. *Semantic-level ordering*: Words and strings are collated based on their meaning; entire words
1895 (such as ''the'') are eliminated; the ordering is not deterministic. This usually requires
1896 special software and is highly dependent on the intended use.

1897 While the historical collation order formally is at level 1, for the English language it corresponds
1898 roughly to elements at level 2. The user expects to see the output from the *ls* utility sorted very
1899 much as it would be in a dictionary. While telephone book ordering would be an optimal goal
1900 for standard collation, this was ruled out as the order would be language-dependent.
1901 Furthermore, a requirement was that the order must be determined solely from the text string
1902 and the collation rules; no external information (for example, ''pronunciation dictionaries'')
1903 could be required.

As a result, the goal for the collation support is at level 3. This also matches the requirements for the Canadian collation order, as well as other, known collation requirements for alphabetic scripts. It specifically rules out collation based on pronunciation rules or based on semantic analysis of the text.

The syntax for the *LC_COLLATE* category source meets the requirements for level 3 and has been verified to produce the correct result with examples based on French, Canadian, and Danish collation order. Because it supports multi-character collating elements, it is also capable of supporting collation in codesets where a character is expressed using non-spacing characters followed by the base character (such as the ISO/IEC 6937: 1994 standard).

The directives that can be specified in an operand to the **order_start** keyword are based on the requirements specified in several proposed standards and in customary use. The following is a rephrasing of rules defined for ''lexical ordering in English and French'' by the Canadian Standards Association (the text in square brackets is rephrased):

- Once special characters [punctuation] have been removed from original strings, the ordering is determined by scanning forwards (left to right) [disregarding case and diacriticals].

- In case of equivalence, special characters are once again removed from original strings and the ordering is determined by scanning backwards (starting from the rightmost character of the string and back), character by character [disregarding case but considering diacriticals].

- In case of repeated equivalence, special characters are removed again from original strings and the ordering is determined by scanning forwards, character by character [considering both case and diacriticals].

- If there is still an ordering equivalence after the first three rules have been applied, then only special characters and the position they occupy in the string are considered to determine ordering. The string that has a special character in the lowest position comes first. If two strings have a special character in the same position, the character [with the lowest collation value] comes first. In case of equality, the other special characters are considered until there is a difference or until all special characters have been exhausted.

It is estimated that this part of IEEE Std 1003.1-2001 covers the requirements for all European languages, and no particular problems are anticipated with Slavic or Middle East character sets.

The Far East (particularly Japanese/Chinese) collations are often based on contextual information and pronunciation rules (the same ideogram can have different meanings and different pronunciations). Such collation, in general, falls outside the desired goal of IEEE Std 1003.1-2001. There are, however, several other collation rules (stroke/radical or ''most common pronunciation'') that can be supported with the mechanism described here.

The character order is defined by the order in which characters and elements are specified between the **order_start** and **order_end** keywords. Weights assigned to the characters and elements define the collation sequence; in the absence of weights, the character order is also the collation sequence.

The **position** keyword provides the capability to consider, in a compare, the relative position of characters not subject to **IGNORE**. As an example, consider the two strings `"o-ring"` and `"or-ing"`. Assuming the hyphen is subject to **IGNORE** on the first pass, the two strings compare equal, and the position of the hyphen is immaterial. On second pass, all characters except the hyphen are subject to **IGNORE**, and in the normal case the two strings would again compare equal. By taking position into account, the first collates before the second.

1948 *A.7.3.3  LC_MONETARY*

1949   The currency symbol does not appear in *LC_MONETARY* because it is not defined in the C locale
1950   of the ISO C standard.

1951   The ISO C standard limits the size of decimal points and thousands delimiters to single-byte
1952   values. In locales based on multi-byte coded character sets, this cannot be enforced;
1953   IEEE Std 1003.1-2001 does not prohibit such characters, but makes the behavior unspecified (in
1954   the text ''In contexts where other standards ...'').

1955   The grouping specification is based on, but not identical to, the ISO C standard. The −1 indicates
1956   that no further grouping is performed; the equivalent of {CHAR_MAX} in the ISO C standard.

1957   The text ''the value is not available in the locale'' is taken from the ISO C standard and is used
1958   instead of the ''unspecified'' text in early proposals. There is no implication that omitting these
1959   keywords or assigning them values of " " or −1 produces unspecified results; such omissions or
1960   assignments eliminate the effects described for the keyword or produce zero-length strings, as
1961   appropriate.

1962   The locale definition is an extension of the ISO C standard *localeconv*() specification. In
1963   particular, rules on how **currency_symbol** is treated are extended to also cover **int_curr_symbol**,
1964   and **p_set_by_space** and **n_sep_by_space** have been augmented with the value 2, which places
1965   a <space> between the sign and the symbol (if they are adjacent; otherwise, it should be treated
1966   as a 0). The following table shows the result of various combinations:

1967
1968

| | | **p_sep_by_space** | | |
| | | **2** | **1** | **0** |
|---|---|---|---|---|
| **p_cs_precedes** = 1 | **p_sign_posn** = 0 | `($1.25)` | `($ 1.25)` | `($1.25)` |
| | **p_sign_posn** = 1 | `+ $1.25` | `+$ 1.25` | `+$1.25` |
| | **p_sign_posn** = 2 | `$1.25 +` | `$ 1.25+` | `$1.25+` |
| | **p_sign_posn** = 3 | `+ $1.25` | `+$ 1.25` | `+$1.25` |
| | **p_sign_posn** = 4 | `$ +1.25` | `$+ 1.25` | `$+1.25` |
| **p_cs_precedes** = 0 | **p_sign_posn** = 0 | `(1.25 $)` | `(1.25 $)` | `(1.25$)` |
| | **p_sign_posn** = 1 | `+1.25 $` | `+1.25 $` | `+1.25$` |
| | **p_sign_posn** = 2 | `1.25$ +` | `1.25 $+` | `1.25$+` |
| | **p_sign_posn** = 3 | `1.25+ $` | `1.25 +$` | `1.25+$` |
| | **p_sign_posn** = 4 | `1.25$ +` | `1.25 $+` | `1.25$+` |

1969–1978

1979   The following is an example of the interpretation of the **mon_grouping** keyword. Assuming that
1980   the value to be formatted is 123 456 789 and the **mon_thousands_sep** is ' ' ' ', then the following
1981   table shows the result. The third column shows the equivalent string in the ISO C standard that
1982   would be used by the *localeconv*() function to accommodate this grouping.

1983

| **mon_grouping** | **Formatted Value** | **ISO C String** |
|---|---|---|
| 3;−1 | 123456'789 | `"\3\177"` |
| 3 | 123'456'789 | `"\3"` |
| 3;2;−1 | 1234'56'789 | `"\3\2\177"` |
| 3;2 | 12'34'56'789 | `"\3\2"` |
| −1 | 123456789 | `"\177"` |

1984–1988

1989   In these examples, the octal value of {CHAR_MAX} is 177.

1990   IEEE Std 1003.1-2001/Cor 1-2002, item XBD/TC1/D6/6 adds a correction that permits the Euro    |
1991   currency symbol and addresses extensibility. The correction is stated using the term ''should''    |
1992   intentionally, in order to make this a recommendation rather than a restriction on    |
1993   implementations. This allows for flexibility in implementations on how they handle future    |

1994              currency symbol additions.                                                                      |

1995              IEEE Std 1003.1-2001/Cor 1-2002, tem XBD/TC1/D6/5 is applied, adding the **int_[np]_\*** values    |
1996              to the POSIX locale definition of *LC_MONETARY*.                                                  |

### A.7.3.4   *LC_NUMERIC*

1997

1998              See the rationale for *LC_MONETARY* for a description of the behavior of grouping.

### A.7.3.5   *LC_TIME*

1999

2000              Although certain of the conversion specifications in the POSIX locale (such as the name of the
2001              month) are shown with initial capital letters, this need not be the case in other locales. Programs
2002              using these conversion specifications may need to adjust the capitalization if the output is going
2003              to be used at the beginning of a sentence.

2004              The *LC_TIME* descriptions of **abday**, **day**, **mon**, and **abmon** imply a Gregorian style calendar (7-
2005              day weeks, 12-month years, leap years, and so on). Formatting time strings for other types of
2006              calendars is outside the scope of IEEE Std 1003.1-2001.

2007              While the ISO 8601: 2000 standard numbers the weekdays starting with Monday, historical
2008              practice is to use the Sunday as the first day. Rather than change the order and introduce
2009              potential confusion, the days must be specified beginning with Sunday; previous references to
2010              ''first day'' have been removed. Note also that the Shell and Utilities volume of
2011              IEEE Std 1003.1-2001 *date* utility supports numbering compliant with the ISO 8601: 2000
2012              standard.

2013              As specified under *date* in the Shell and Utilities volume of IEEE Std 1003.1-2001 and *strftime*( ) in
2014              the System Interfaces volume of IEEE Std 1003.1-2001, the conversion specifications
2015              corresponding to the optional keywords consist of a modifier followed by a traditional
2016              conversion specification (for instance, %Ex). If the optional keywords are not supported by the
2017              implementation or are unspecified for the current locale, these modified conversion
2018              specifications are treated as the traditional conversion specifications. For example, assume the
2019              following keywords:

2020              ```
                  alt_digits     "0th";"1st";"2nd";"3rd";"4th";"5th";\
2021                             "6th";"7th";"8th";"9th";"10th"
2022                  d_fmt          "The %Od day of %B in %Y"
                  ```

2023              On July 4th 1776, the %x conversion specifications would result in "The 4th day of July
2024              in 1776", while on July 14th 1789 it would result in "The 14 day of July in 1789". It
2025              can be noted that the above example is for illustrative purposes only; the %O modifier is
2026              primarily intended to provide for Kanji or Hindi digits in *date* formats.

2027              The following is an example for Japan that supports the current plus last three Emperors and
2028              reverts to Western style numbering for years prior to the Meiji era. The example also allows for
2029              the custom of using a special name for the first year of an era instead of using 1. (The examples
2030              substitute romaji where kanji should be used.)

```
2031          era_d_fmt "%EY%mgatsu%dnichi (%a)"

2032          era    "+:2:1990/01/01:+*:Heisei:%EC%Eynen";\
2033                 "+:1:1989/01/08:1989/12/31:Heisei:%ECgannen";\
2034                 "+:2:1927/01/01:1989/01/07:Shouwa:%EC%Eynen";\
2035                 "+:1:1926/12/25:1926/12/31:Shouwa:%ECgannen";\
2036                 "+:2:1913/01/01:1926/12/24:Taishou:%EC%Eynen";\
2037                 "+:1:1912/07/30:1912/12/31:Taishou:%ECgannen";\
2038                 "+:2:1869/01/01:1912/07/29:Meiji:%EC%Eynen";\
2039                 "+:1:1868/09/08:1868/12/31:Meiji:%ECgannen";\
2040                 "-:1868:1868/09/07:-*::%Ey"
```

2041   Assuming that the current date is September 21, 1991, a request to *date* or *strftime*( ) would yield
2042   the following results:

```
2043          %Ec - Heisei3nen9gatsu21nichi (Sat) 14:39:26
2044          %EC - Heisei
2045          %Ex - Heisei3nen9gatsu21nichi (Sat)
2046          %Ey - 3
2047          %EY - Heisei3nen
```

2048   Example era definitions for the Republic of China:

```
2049          era    "+:2:1913/01/01:+*:ChungHwaMingGuo:%EC%EyNen";\
2050                 "+:1:1912/1/1:1912/12/31:ChungHwaMingGuo:%ECYuenNen";\
2051                 "+:1:1911/12/31:-*:MingChien:%EC%EyNen"
```

2052   Example definitions for the Christian Era:

```
2053          era    "+:1:0001/01/01:+*:AD:%EC %Ey";\
2054                 "+:1:-0001/12/31:-*:BC:%Ey %EC"
```

2055   *A.7.3.6   LC_MESSAGES*

2056   The **yesstr** and **nostr** locale keywords and the YESSTR and NOSTR *langinfo* items were formerly
2057   used to match user affirmative and negative responses. In IEEE Std 1003.1-2001, the **yesexpr**,
2058   **noexpr**, YESEXPR, and NOEXPR extended regular expressions have replaced them.
2059   Applications should use the general locale-based messaging facilities to issue prompting
2060   messages which include sample desired responses.

2061   **A.7.4     Locale Definition Grammar**

2062   There is no additional rationale provided for this section.

2063   *A.7.4.1   Locale Lexical Conventions*

2064   There is no additional rationale provided for this section.

2065   *A.7.4.2   Locale Grammar*

2066   There is no additional rationale provided for this section.

2067 **A.7.5    Locale Definition Example**

2068    The following is an example of a locale definition file that could be used as input to the *localedef*
2069    utility. It assumes that the utility is executed with the **–f** option, naming a charmap file with (at
2070    least) the following content:

```
2071        CHARMAP
2072        <space>       \x20
2073        <dollar>      \x24
2074        <A>           \101
2075        <a>           \141
2076        <A-acute>     \346
2077        <a-acute>     \365
2078        <A-grave>     \300
2079        <a-grave>     \366
2080        <b>           \142
2081        <C>           \103
2082        <c>           \143
2083        <c-cedilla>   \347
2084        <d>           \x64
2085        <H>           \110
2086        <h>           \150
2087        <eszet>       \xb7
2088        <s>           \x73
2089        <z>           \x7a
2090        END CHARMAP
```

2091    It should not be taken as complete or to represent any actual locale, but only to illustrate the
2092    syntax.

```
2093        #
2094        LC_CTYPE
2095        lower   <a>;<b>;<c>;<c-cedilla>;<d>;...;<z>
2096        upper   A;B;C;Ç;...;Z
2097        space   \x20;\x09;\x0a;\x0b;\x0c;\x0d
2098        blank   \040;\011
2099        toupper (<a>,<A>);(b,B);(c,C);(ç,Ç);(d,D);(z,Z)
2100        END LC_CTYPE
2101        #
2102        LC_COLLATE
2103        #
2104        # The following example of collation is based on
2105        # Canadian standard Z243.4.1-1998, "Canadian Alphanumeric
2106        # Ordering Standard for Character Sets of CSA Z234.4 Standard".
2107        # (Other parts of this example locale definition file do not
2108        # purport to relate to Canada, or to any other real culture.)
2109        # The proposed standard defines a 4-weight collation, such that
2110        # in the first pass, characters are compared without regard to
2111        # case or accents; in the second pass, backwards-compare without
2112        # regard to case; in the third pass, forwards-compare without
2113        # regard to diacriticals. In the 3 first passes, non-alphabetic
2114        # characters are ignored; in the fourth pass, only special
2115        # characters are considered, such that "The string that has a
2116        # special character in the lowest position comes first. If two
```

```
2117        # strings have a special character in the same position, the
2118        # collation value of the special character determines ordering.
2119        #
2120        # Only a subset of the character set is used here; mostly to
2121        # illustrate the set-up.
2122        #
2123        collating-symbol <NULL>
2124        collating-symbol <LOW_VALUE>
2125        collating-symbol <LOWER-CASE>
2126        collating-symbol <SUBSCRIPT-LOWER>
2127        collating-symbol <SUPERSCRIPT-LOWER>
2128        collating-symbol <UPPER-CASE>
2129        collating-symbol <NO-ACCENT>
2130        collating-symbol <PECULIAR>
2131        collating-symbol <LIGATURE>
2132        collating-symbol <ACUTE>
2133        collating-symbol <GRAVE>
2134        # Further collating-symbols follow.
2135        #
2136        # Properly, the standard does not include any multi-character
2137        # collating elements; the one below is added for completeness.
2138        #
2139        collating_element <ch> from "<c><h>"
2140        collating_element <CH> from "<C><H>"
2141        collating_element <Ch> from "<C><h>"
2142        #
2143        order_start forward;backward;forward;forward,position
2144        #
2145        # Collating symbols are specified first in the sequence to allocate
2146        # basic collation values to them, lower than that of any character.
2147        <NULL>
2148        <LOW_VALUE>
2149        <LOWER-CASE>
2150        <SUBSCRIPT-LOWER>
2151        <SUPERSCRIPT-LOWER>
2152        <UPPER-CASE>
2153        <NO-ACCENT>
2154        <PECULIAR>
2155        <LIGATURE>
2156        <ACUTE>
2157        <GRAVE>
2158        <RING-ABOVE>
2159        <DIAERESIS>
2160        <TILDE>
2161        # Further collating symbols are given a basic collating value here.
2162        #
2163        # Here follow special characters.
2164        <space>          IGNORE;IGNORE;IGNORE;<space>
2165        # Other special characters follow here.
2166        #
2167        # Here follow the regular characters.
2168        <a>         <a>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
```

```
2169        <A>          <a>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2170        <a-acute>  <a>;<ACUTE>;<LOWER-CASE>;IGNORE
2171        <A-acute>  <a>;<ACUTE>;<UPPER-CASE>;IGNORE
2172        <a-grave>  <a>;<GRAVE>;<LOWER-CASE>;IGNORE
2173        <A-grave>  <a>;<GRAVE>;<UPPER-CASE>;IGNORE
2174        <ae>         "<a><e>";"<LIGATURE><LIGATURE>";\
2175                     "<LOWER-CASE><LOWER-CASE>";IGNORE
2176        <AE>         "<a><e>";"<LIGATURE><LIGATURE>";\
2177                     "<UPPER-CASE><UPPER-CASE>";IGNORE
2178        <b>           <b>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2179        <B>           <b>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2180        <c>           <c>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2181        <C>           <c>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2182        <ch>          <ch>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2183        <Ch>          <ch>;<NO-ACCENT>;<PECULIAR>;IGNORE
2184        <CH>          <ch>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2185        #
2186        # As an example, the strings "Bach" and "bach" could be encoded (for
2187        # compare purposes) as:
2188        # "Bach"  <b>;<a>;<ch>;<LOW_VALUE>;<NO_ACCENT>;<NO_ACCENT>;\
2189        #         <NO_ACCENT>;<LOW_VALUE>;<UPPER-CASE>;<LOWER-CASE>;\
2190        #         <LOWER-CASE>;<NULL>
2191        # "bach"  <b>;<a>;<ch>;<LOW_VALUE>;<NO_ACCENT>;<NO_ACCENT>;\
2192        #         <NO_ACCENT>;<LOW_VALUE>;<LOWER-CASE>;<LOWER-CASE>;\
2193        #         <LOWER-CASE>;<NULL>
2194        #
2195        # The two strings are equal in pass 1 and 2, but differ in pass 3.
2196        #
2197        # Further characters follow.
2198        #
2199        UNDEFINED    IGNORE;IGNORE;IGNORE;IGNORE
2200        #
2201        order_end
2202        #
2203        END LC_COLLATE
2204        #
2205        LC_MONETARY
2206        int_curr_symbol    "USD "
2207        currency_symbol    "$"
2208        mon_decimal_point  "."
2209        mon_grouping       3;0
2210        positive_sign      ""
2211        negative_sign      "-"
2212        p_cs_precedes      1
2213        n_sign_posn        0
2214        END LC_MONETARY
2215        #
2216        LC_NUMERIC
2217        copy "US_en.ASCII"
2218        END LC_NUMERIC
2219        #
2220        LC_TIME
```

```
2221        abday   "Sun";"Mon";"Tue";"Wed";"Thu";"Fri";"Sat"
2222        #
2223        day     "Sunday";"Monday";"Tuesday";"Wednesday";\
2224                "Thursday";"Friday";"Saturday"
2225        #
2226        abmon   "Jan";"Feb";"Mar";"Apr";"May";"Jun";\
2227                 "Jul";"Aug";"Sep";"Oct";"Nov";"Dec"
2228        #
2229        mon     "January";"February";"March";"April";\
2230                "May";"June";"July";"August";"September";\
2231                "October";"November";"December"
2232        #
2233        d_t_fmt "%a %b %d %T %Z %Y\n"
2234        END LC_TIME
2235        #
2236        LC_MESSAGES
2237        yesexpr "^([yY][[:alpha:]]*)|(OK)"
2238        #
2239        noexpr  "^[nN][[:alpha:]]*"
2240        END LC_MESSAGES
```

# 2241   A.8   Environment Variables

## 2242   A.8.1   Environment Variable Definition

2243  The variable *environ* is not intended to be declared in any header, but rather to be declared by the
2244  user for accessing the array of strings that is the environment. This is the traditional usage of the
2245  symbol. Putting it into a header could break some programs that use the symbol for their own
2246  purposes.

2247  The decision to restrict conforming systems to the use of digits, uppercase letters, and
2248  underscores for environment variable names allows applications to use lowercase letters in their
2249  environment variable names without conflicting with any conforming system.

2250  In addition to the obvious conflict with the shell syntax for positional parameter substitution,
2251  some historical applications (including some shells) exclude names with leading digits from the
2252  environment.

## 2253   A.8.2   Internationalization Variables

2254  The text about locale implies that any utilities written in standard C and conforming to
2255  IEEE Std 1003.1-2001 must issue the following call:

2256  ```
     setlocale(LC_ALL, "")
```

2257  If this were omitted, the ISO C standard specifies that the C locale would be used.

2258  If any of the environment variables are invalid, it makes sense to default to an implementation-
2259  defined, consistent locale environment. It is more confusing for a user to have partial settings
2260  occur in case of a mistake. All utilities would then behave in one language/cultural
2261  environment. Furthermore, it provides a way of forcing the whole environment to be the
2262  implementation-defined default. Disastrous results could occur if a pipeline of utilities partially
2263  uses the environment variables in different ways. In this case, it would be appropriate for
2264  utilities that use *LANG* and related variables to exit with an error if any of the variables are

2265   invalid. For example, users typing individual commands at a terminal might want *date* to work if
2266   *LC_MONETARY* is invalid as long as *LC_TIME* is valid. Since these are conflicting reasonable
2267   alternatives, IEEE Std 1003.1-2001 leaves the results unspecified if the locale environment
2268   variables would not produce a complete locale matching the specification of the user.

2269   The locale settings of individual categories cannot be truly independent and still guarantee
2270   correct results. For example, when collating two strings, characters must first be extracted from
2271   each string (governed by *LC_CTYPE*) before being mapped to collating elements (governed by
2272   *LC_COLLATE*) for comparison. That is, if *LC_CTYPE* is causing parsing according to the rules of
2273   a large, multi-byte code set (potentially returning 20 000 or more distinct character codeset
2274   values), but *LC_COLLATE* is set to handle only an 8-bit codeset with 256 distinct characters,
2275   meaningful results are obviously impossible.

2276   The *LC_MESSAGES* variable affects the language of messages generated by the standard
2277   utilities.

2278   The description of the environment variable names starting with the characters ''LC_''
2279   acknowledges the fact that the interfaces presented may be extended as new international
2280   functionality is required. In the ISO C standard, names preceded by ''LC_'' are reserved in the
2281   name space for future categories.

2282   To avoid name clashes, new categories and environment variables are divided into two
2283   classifications: ''implementation-independent'' and ''implementation-defined''.

2284   Implementation-independent names will have the following format:

2285       LC_*NAME*

2286   where *NAME* is the name of the new category and environment variable. Capital letters must be
2287   used for implementation-independent names.

2288   Implementation-defined names must be in lowercase letters, as below:

2289       LC_*name*

## 2290   A.8.3   Other Environment Variables

2291   **COLUMNS, LINES**

2292   The default values for the number of column positions, *COLUMNS*, and screen height, *LINES*,
2293   are unspecified because historical implementations use different methods to determine values
2294   corresponding to the size of the screen in which the utility is run. This size is typically known to
2295   the implementation through the value of *TERM*, or by more elaborate methods such as
2296   extensions to the *stty* utility or knowledge of how the user is dynamically resizing windows on a
2297   bit-mapped display terminal. Users should not need to set these variables in the environment
2298   unless there is a specific reason to override the default behavior of the implementation, such as
2299   to display data in an area arbitrarily smaller than the terminal or window. Values for these
2300   variables that are not decimal integers greater than zero are implicitly undefined values; it is
2301   unnecessary to enumerate all of the possible values outside of the acceptable set.

2302        **LOGNAME**

2303        In most implementations, the value of such a variable is easily forged, so security-critical
2304        applications should rely on other means of determining user identity. *LOGNAME* is required to
2305        be constructed from the portable filename character set for reasons of interchange. No diagnostic
2306        condition is specified for violating this rule, and no requirement for enforcement exists. The
2307        intent of the requirement is that if extended characters are used, the ''guarantee'' of portability
2308        implied by a standard is void.

2309        **PATH**

2310        Many historical implementations of the Bourne shell do not interpret a trailing colon to represent
2311        the current working directory and are thus non-conforming. The C Shell and the KornShell
2312        conform to IEEE Std 1003.1-2001 on this point. The usual name of dot may also be used to refer
2313        to the current working directory.

2314        Many implementations historically have used a default value of **/bin** and **/usr/bin** for the *PATH*
2315        variable. IEEE Std 1003.1-2001 does not mandate this default path be identical to that retrieved
2316        from *getconf* _CS_PATH because it is likely that the standardized utilities may be provided in
2317        another directory separate from the directories used by some historical applications.

2318        **SHELL**

2319        The *SHELL* variable names the preferred shell of the user; it is a guide to applications. There is
2320        no direct requirement that that shell conform to IEEE Std 1003.1-2001; that decision should rest
2321        with the user. It is the intention of the standard developers that alternative shells be permitted, if
2322        the user chooses to develop or acquire one. An operating system that builds its shell into the
2323        ''kernel'' in such a manner that alternative shells would be impossible does not conform to the
2324        spirit of IEEE Std 1003.1-2001.

2325        **TZ**

2326        The quoted form of the timezone variable allows timezone names of the form UTC+1 (or any
2327        name that contains the character plus (' + '), the character minus (' − '), or digits), which may be
2328        appropriate for countries that do not have an official timezone name. It would be coded as
2329        <UTC+1>+1<UTC+2>, which would cause *std* to have a value of UTC+1 and *dst* a value of
2330        UTC+2, each with a length of 5 characters. This does not appear to conflict with any existing
2331        usage. The characters ' < ' and ' > ' were chosen for quoting because they are easier to parse
2332        visually than a quoting character that does not provide some sense of bracketing (and in a string
2333        like this, such bracketing is helpful). They were also chosen because they do not need special
2334        treatment when assigning to the *TZ* variable. Users are often confused by embedding quotes in a
2335        string. Because ' < ' and ' > ' are meaningful to the shell, the whole string would have to be
2336        quoted, but that is easily explained. (Parentheses would have presented the same problems.)
2337        Although the ' > ' symbol could have been permitted in the string by either escaping it or
2338        doubling it, it seemed of little value to require that. This could be provided as an extension if
2339        there was a need. Timezone names of this new form lead to a requirement that the value of
2340        {_POSIX_TZNAME_MAX} change from 3 to 6.

2341        Since the *TZ* environment variable is usually inherited by all applications started by a user after
2342        the value of the *TZ* environment variable is changed and since many applications run using the
2343        C or POSIX locale, using characters that are not in the portable character set in the *std* and *dst*
2344        fields could cause unexpected results.

2345        The format of the *TZ* environment variable is changed in Issue 6 to allow for the quoted form, as
2346        defined in previous versions of the ISO POSIX-1 standard.

2347    IEEE Std 1003.1-2001/Cor 1-2002, item XBD/TC1/D6/7 is applied, adding the *ctime_r*() and  |
2348    *localtime_r*() functions to the list of functions that use the *TZ* environment variable.            |

## 2349  A.9    Regular Expressions

2350    Rather than repeating the description of REs for each utility supporting REs, the standard
2351    developers preferred a common, comprehensive description of regular expressions in one place.
2352    The most common behavior is described here, and exceptions or extensions to this are
2353    documented for the respective utilities, as appropriate.

2354    The BRE corresponds to the *ed* or historical *grep* type, and the ERE corresponds to the historical
2355    *egrep* type (now *grep* −**E**).

2356    The text is based on the *ed* description and substantially modified, primarily to aid developers
2357    and others in the understanding of the capabilities and limitations of REs. Much of this was
2358    influenced by internationalization requirements.

2359    It should be noted that the definitions in this section do not cover the *tr* utility; the *tr* syntax does
2360    not employ REs.

2361    The specification of REs is particularly important to internationalization because pattern
2362    matching operations are very basic operations in business and other operations. The syntax and
2363    rules of REs are intended to be as intuitive as possible to make them easy to understand and use.
2364    The historical rules and behavior do not provide that capability to non-English language users,
2365    and do not provide the necessary support for commonly used characters and language
2366    constructs. It was necessary to provide extensions to the historical RE syntax and rules to
2367    accommodate other languages.

2368    As they are limited to bracket expressions, the rationale for these modifications is in the Base
2369    Definitions volume of IEEE Std 1003.1-2001, Section 9.3.5, RE Bracket Expression.

### 2370  A.9.1    Regular Expression Definitions

2371    It is possible to determine what strings correspond to subexpressions by recursively applying
2372    the leftmost longest rule to each subexpression, but only with the proviso that the overall match
2373    is leftmost longest. For example, matching `"\(ac*\)c*d[ac]*\1"` against *acdacaaa* matches
2374    *acdacaaa* (with \1=*a*); simply matching the longest match for `"\(ac*\)"` would yield \1=*ac*, but
2375    the overall match would be smaller (*acdac*). Conceptually, the implementation must examine
2376    every possible match and among those that yield the leftmost longest total matches, pick the one
2377    that does the longest match for the leftmost subexpression, and so on. Note that this means that
2378    matching by subexpressions is context-dependent: a subexpression within a larger RE may
2379    match a different string from the one it would match as an independent RE, and two instances of
2380    the same subexpression within the same larger RE may match different lengths even in similar
2381    sequences of characters. For example, in the ERE `"(a.*b)(a.*b)"`, the two identical
2382    subexpressions would match four and six characters, respectively, of *accbacccb*.

2383    The definition of single character has been expanded to include also collating elements
2384    consisting of two or more characters; this expansion is applicable only when a bracket
2385    expression is included in the BRE or ERE. An example of such a collating element may be the
2386    Dutch *ij*, which collates as a ′y′. In some encodings, a ligature ''i with j'' exists as a character
2387    and would represent a single-character collating element. In another encoding, no such ligature
2388    exists, and the two-character sequence *ij* is defined as a multi-character collating element.
2389    Outside brackets, the *ij* is treated as a two-character RE and matches the same characters in a
2390    string. Historically, a bracket expression only matched a single character. The ISO POSIX-2: 1993
2391    standard required bracket expressions like `"[^[:lower:]]"` to match multi-character collating

2392    elements such as `"ij"`. However, this requirement led to behavior that many users did not
2393    expect and that could not feasibly be mimicked in user code, and it was rarely if ever
2394    implemented correctly. The current standard leaves it unspecified whether a bracket expression
2395    matches a multi-character collating element, allowing both historical and ISO POSIX-2:1993
2396    standard implementations to conform.

2397    Also, in the current standard, it is unspecified whether character class expressions like
2398    `"[:lower:]"` can include multi-character collating elements like `"ij"`; hence
2399    `"[[:lower:]]"` can match `"ij"`, and `"[^[:lower:]]"` can fail to match `"ij"`. Common
2400    practice is for a character class expression to match a collating element if it matches the collating
2401    element's first character.

## 2402    A.9.2    Regular Expression General Requirements

2403    The definition of which sequence is matched when several are possible is based on the leftmost-
2404    longest rule historically used by deterministic recognizers. This rule is easier to define and
2405    describe, and arguably more useful, than the first-match rule historically used by non-
2406    deterministic recognizers. It is thought that dependencies on the choice of rule are rare; carefully
2407    contrived examples are needed to demonstrate the difference.

2408    A formal expression of the leftmost-longest rule is:

2409        The search is performed as if all possible suffixes of the string were tested for a prefix
2410        matching the pattern; the longest suffix containing a matching prefix is chosen, and the
2411        longest possible matching prefix of the chosen suffix is identified as the matching sequence.

2412    Historically, most RE implementations only match lines, not strings. However, that is more an
2413    effect of the usage than of an inherent feature of REs themselves. Consequently,
2414    IEEE Std 1003.1-2001 does not regard <newline>s as special; they are ordinary characters, and
2415    both a period and a non-matching list can match them. Those utilities (like *grep*) that do not
2416    allow <newline>s to match are responsible for eliminating any <newline> from strings before
2417    matching against the RE. The *regcomp*() function, however, can provide support for such
2418    processing without violating the rules of this section.

2419    Some implementations of *egrep* have had very limited flexibility in handling complex EREs.
2420    IEEE Std 1003.1-2001 does not attempt to define the complexity of a BRE or ERE, but does place a
2421    lower limit on it—any RE must be handled, as long as it can be expressed in 256 bytes or less. (Of
2422    course, this does not place an upper limit on the implementation.) There are historical programs
2423    using a non-deterministic-recognizer implementation that should have no difficulty with this
2424    limit. It is possible that a good approach would be to attempt to use the faster, but more limited,
2425    deterministic recognizer for simple expressions and to fall back on the non-deterministic
2426    recognizer for those expressions requiring it. Non-deterministic implementations must be
2427    careful to observe the rules on which match is chosen; the longest match, not the first match,
2428    starting at a given character is used.

2429    The term ''invalid'' highlights a difference between this section and some others:
2430    IEEE Std 1003.1-2001 frequently avoids mandating of errors for syntax violations because they
2431    can be used by implementors to trigger extensions. However, the authors of the
2432    internationalization features of REs wanted to mandate errors for certain conditions to identify
2433    usage problems or non-portable constructs. These are identified within this rationale as
2434    appropriate. The remaining syntax violations have been left implicitly or explicitly undefined.
2435    For example, the BRE construct `"\{1,2,3\}"` does not comply with the grammar. A
2436    conforming application cannot rely on it producing an error nor matching the literal characters
2437    `"\{1,2,3\}"`.

2438        The term ''undefined'' was used in favor of ''unspecified'' because many of the situations are
2439        considered errors on some implementations, and the standard developers considered that
2440        consistency throughout the section was preferable to mixing undefined and unspecified.

2441  **A.9.3     Basic Regular Expressions**

2442        There is no additional rationale provided for this section.

2443  *A.9.3.1   BREs Matching a Single Character or Collating Element*

2444        There is no additional rationale provided for this section.

2445  *A.9.3.2   BRE Ordinary Characters*

2446        There is no additional rationale provided for this section.

2447  *A.9.3.3   BRE Special Characters*

2448        There is no additional rationale provided for this section.

2449  *A.9.3.4   Periods in BREs*

2450        There is no additional rationale provided for this section.

2451  *A.9.3.5   RE Bracket Expression*

2452        Range expressions are, historically, an integral part of REs. However, the requirements of
2453        ''natural language behavior'' and portability do conflict. In the POSIX locale, ranges must be
2454        treated according to the collating sequence and include such characters that fall within the range
2455        based on that collating sequence, regardless of character values. In other locales, ranges have
2456        unspecified behavior.

2457        Some historical implementations allow range expressions where the ending range point of one
2458        range is also the starting point of the next (for instance, `"[a−m−o]"`). This behavior should not
2459        be permitted, but to avoid breaking historical implementations, it is now *undefined* whether it is a
2460        valid expression and how it should be interpreted.

2461        Current practice in *awk* and *lex* is to accept escape sequences in bracket expressions as per the
2462        Base Definitions volume of IEEE Std 1003.1-2001, Table 5-1, Escape Sequences and Associated
2463        Actions, while the normal ERE behavior is to regard such a sequence as consisting of two
2464        characters. Allowing the *awk/lex* behavior in EREs would change the normal behavior in an
2465        unacceptable way; it is expected that *awk* and *lex* will decode escape sequences in EREs before
2466        passing them to *regcomp*() or comparable routines. Each utility describes the escape sequences it
2467        accepts as an exception to the rules in this section; the list is not the same, for historical reasons.

2468        As noted previously, the new syntax and rules have been added to accommodate other
2469        languages than English. The remainder of this section describes the rationale for these
2470        modifications.

2471        In the POSIX locale, a regular expression that starts with a range expression matches a set of
2472        strings that are contiguously sorted, but this is not necessarily true in other locales. For example,
2473        a French locale might have the following behavior:

```
2474    $ ls
2475    alpha   Alpha   estimé   ESTIMÉ   été   eurêka
2476    $ ls [a-e]*
2477    alpha   Alpha   estimé   eurêka
```

        

2478     Such disagreements between matching and contiguous sorting are unavoidable because POSIX
2479     sorting cannot be implemented in terms of a deterministic finite-state automaton (DFA), but
2480     range expressions by design are implementable in terms of DFAs.

2481     Historical implementations used native character order to interpret range expressions. The
2482     ISO POSIX-2: 1993 standard instead required collating element order (CEO): the order that
2483     collating elements were specified between the **order_start** and **order_end** keywords in the
2484     *LC_COLLATE* category of the current locale. CEO had some advantages in portability over the
2485     native character order, but it also had some disadvantages:

2486     • CEO could not feasibly be mimicked in user code, leading to inconsistencies between POSIX
2487       matchers and matchers in popular user programs like Emacs, *ksh*, and Perl.

2488     • CEO caused range expressions to match accented and capitalized letters contrary to many
2489       users' expectations. For example, "`[a-e]`" typically matched both 'E' and 'á' but neither
2490       'A' nor 'é'.

2491     • CEO was not consistent across implementations. In practice, CEO was often less portable
2492       than native character order. For example, it was common for the CEOs of two
2493       implementation-supplied locales to disagree, even if both locales were named "`da_DK`".

2494     Because of these problems, some implementations of regular expressions continued to use
2495     native character order. Others used the collation sequence, which is more consistent with sorting
2496     than either CEO or native order, but which departs further from the traditional POSIX semantics
2497     because it generally requires "`[a-e]`" to match either 'A' or 'E' but not both. As a result of
2498     this kind of implementation variation, programmers who wanted to write portable regular
2499     expressions could not rely on the ISO POSIX-2: 1993 standard guarantees in practice.

2500     While revising the standard, lengthy consideration was given to proposals to attack this problem
2501     by adding an API for querying the CEO to allow user-mode matchers, but none of these
2502     proposals had implementation experience and none achieved consensus. Leaving the standard
2503     alone was also considered, but rejected due to the problems described above.

2504     The current standard leaves unspecified the behavior of a range expression outside the POSIX
2505     locale. This makes it clearer that conforming applications should avoid range expressions
2506     outside the POSIX locale, and it allows implementations and compatible user-mode matchers to
2507     interpret range expressions using native order, CEO, collation sequence, or other, more
2508     advanced techniques. The concerns which led to this change were raised in IEEE PASC
2509     interpretation 1003.2 #43 and others, and related to ambiguities in the specification of how
2510     multi-character collating elements should be handled in range expressions. These ambiguities
2511     had led to multiple interpretations of the specification, in conflicting ways, which led to varying
2512     implementations. As noted above, efforts were made to resolve the differences, but no solution
2513     has been found that would be specific enough to allow for portable software while not
2514     invalidating existing implementations.

2515     The standard developers recognize that collating elements are important, such elements being
2516     common in several European languages; for example, 'ch' or 'll' in traditional Spanish; 'aa'
2517     in several Scandinavian languages. Existing internationalized implementations have processed,
2518     and continue to process, these elements in range expressions. Efforts are expected to continue in
2519     the future to find a way to define the behavior of these elements precisely and portably.

2520     The ISO POSIX-2: 1993 standard required "`[b-a]`" to be an invalid expression in the POSIX
2521     locale, but this requirement has been relaxed in this version of the standard so that "`[b-a]`" can
2522     instead be treated as a valid expression that does not match any string.

2523   *A.9.3.6   BREs Matching Multiple Characters*

2524   The limit of nine back-references to subexpressions in the RE is based on the use of a single-digit
2525   identifier; increasing this to multiple digits would break historical applications. This does not
2526   imply that only nine subexpressions are allowed in REs. The following is a valid BRE with ten
2527   subexpressions:

2528   `\(\(\(\(ab\)*c\)*d\)\(ef\)*\(gh\)\{2\}\(ij\)*\(kl\)*\(mn\)*\(op\)*\(qr\)*`

2529   The standard developers regarded the common historical behavior, which supported `"\n*"`, but
2530   not `"\n\{min,max\}"`, `"\(...\)*"`, or `"\(...\)\{min,max\}"`, as a non-intentional
2531   result of a specific implementation, and they supported both duplication and interval
2532   expressions following subexpressions and back-references.

2533   The changes to the processing of the back-reference expression remove an unspecified or
2534   ambiguous behavior in the Shell and Utilities volume of IEEE Std 1003.1-2001, aligning it with
2535   the requirements specified for the *regcomp*() expression, and is the result of PASC Interpretation
2536   1003.2-92 #43 submitted for the ISO POSIX-2: 1993 standard.

2537   *A.9.3.7   BRE Precedence*

2538   There is no additional rationale provided for this section.

2539   *A.9.3.8   BRE Expression Anchoring*

2540   Often, the dollar sign is viewed as matching the ending <newline> in text files. This is not
2541   strictly true; the <newline> is typically eliminated from the strings to be matched, and the dollar
2542   sign matches the terminating null character.

2543   The ability of '^', '$', and '*' to be non-special in certain circumstances may be confusing to
2544   some programmers, but this situation was changed only in a minor way from historical practice
2545   to avoid breaking many historical scripts. Some consideration was given to making the use of
2546   the anchoring characters undefined if not escaped and not at the beginning or end of strings.
2547   This would cause a number of historical BREs, such as `"2^10"`, `"$HOME"`, and `"$1.35"`, that
2548   relied on the characters being treated literally, to become invalid.

2549   However, one relatively uncommon case was changed to allow an extension used on some
2550   implementations. Historically, the BREs `"^foo"` and `"\(^foo\)"` did not match the same
2551   string, despite the general rule that subexpressions and entire BREs match the same strings. To
2552   increase consensus, IEEE Std 1003.1-2001 has allowed an extension on some implementations to
2553   treat these two cases in the same way by declaring that anchoring *may* occur at the beginning or
2554   end of a subexpression. Therefore, portable BREs that require a literal circumflex at the
2555   beginning or a dollar sign at the end of a subexpression must escape them. Note that a BRE such
2556   as `"a\(^bc\)"` will either match `"a^bc"` or nothing on different systems under the rules.

2557   ERE anchoring has been different from BRE anchoring in all historical systems. An unescaped
2558   anchor character has never matched its literal counterpart outside a bracket expression. Some
2559   implementations treated `"foo$bar"` as a valid expression that never matched anything; others
2560   treated it as invalid. IEEE Std 1003.1-2001 mandates the former, valid unmatched behavior.

2561   Some implementations have extended the BRE syntax to add alternation. For example, the
2562   subexpression `"\(foo$\|bar\)"` would match either `"foo"` at the end of the string or `"bar"`
2563   anywhere. The extension is triggered by the use of the undefined `"\|"` sequence. Because the
2564   BRE is undefined for portable scripts, the extending system is free to make other assumptions,
2565   such that the '$' represents the end-of-line anchor in the middle of a subexpression. If it were
2566   not for the extension, the '$' would match a literal dollar sign under the rules.

2567 **A.9.4** **Extended Regular Expressions**

2568 As with BREs, the standard developers decided to make the interpretation of escaped ordinary
2569 characters undefined.

2570 The right parenthesis is not listed as an ERE special character because it is only special in the
2571 context of a preceding left parenthesis. If found without a preceding left parenthesis, the right
2572 parenthesis has no special meaning.

2573 The interval expression, "{m,n}", has been added to EREs. Historically, the interval expression
2574 has only been supported in some ERE implementations. The standard developers estimated that
2575 the addition of interval expressions to EREs would not decrease consensus and would also make
2576 BREs more of a subset of EREs than in many historical implementations.

2577 It was suggested that, in addition to interval expressions, back-references ('\n') should also be
2578 added to EREs. This was rejected by the standard developers as likely to decrease consensus.

2579 In historical implementations, multiple duplication symbols are usually interpreted from left to
2580 right and treated as additive. As an example, "a+*b" matches zero or more instances of 'a'
2581 followed by a 'b'. In IEEE Std 1003.1-2001, multiple duplication symbols are undefined; that is,
2582 they cannot be relied upon for conforming applications. One reason for this is to provide some
2583 scope for future enhancements.

2584 The precedence of operations differs between EREs and those in *lex*; in *lex*, for historical reasons,
2585 interval expressions have a lower precedence than concatenation.

2586 *A.9.4.1* *EREs Matching a Single Character or Collating Element*

2587 There is no additional rationale provided for this section.

2588 *A.9.4.2* *ERE Ordinary Characters*

2589 There is no additional rationale provided for this section.

2590 *A.9.4.3* *ERE Special Characters*

2591 There is no additional rationale provided for this section.

2592 *A.9.4.4* *Periods in EREs*

2593 There is no additional rationale provided for this section.

2594 *A.9.4.5* *ERE Bracket Expression*

2595 There is no additional rationale provided for this section.

2596 *A.9.4.6* *EREs Matching Multiple Characters*

2597 There is no additional rationale provided for this section.

2598 *A.9.4.7* *ERE Alternation*

2599 There is no additional rationale provided for this section.

2600  *A.9.4.8   ERE Precedence*

2601      There is no additional rationale provided for this section.

2602  *A.9.4.9   ERE Expression Anchoring*

2603      There is no additional rationale provided for this section.

2604  **A.9.5        Regular Expression Grammar**

2605      The grammars are intended to represent the range of acceptable syntaxes available to
2606      conforming applications. There are instances in the text where undefined constructs are
2607      described; as explained previously, these allow implementation extensions. There is no intended
2608      requirement that an implementation extension must somehow fit into the grammars shown
2609      here.

2610      The BRE grammar does not permit L_ANCHOR or R_ANCHOR inside "\(" and "\)" (which
2611      implies that '^' and '$' are ordinary characters). This reflects the semantic limits on the
2612      application, as noted in the Base Definitions volume of IEEE Std 1003.1-2001, Section 9.3.8, BRE
2613      Expression Anchoring. Implementations are permitted to extend the language to interpret '^'
2614      and '$' as anchors in these locations, and as such, conforming applications cannot use
2615      unescaped '^' and '$' in positions inside "\(" and "\)" that might be interpreted as anchors.

2616      The ERE grammar does not permit several constructs that the Base Definitions volume of
2617      IEEE Std 1003.1-2001, Section 9.4.2, ERE Ordinary Characters and the Base Definitions volume of
2618      IEEE Std 1003.1-2001, Section 9.4.3, ERE Special Characters specify as having undefined results:

2619      • ORD_CHAR preceded by '\'

2620      • *ERE_dupl_symbol*(s) appearing first in an ERE, or immediately following '|', '^', or '('

2621      • '{' not part of a valid *ERE_dupl_symbol*

2622      • '|' appearing first or last in an ERE, or immediately following '|' or '(', or immediately
2623         preceding ')'

2624      Implementations are permitted to extend the language to allow these. Conforming applications
2625      cannot use such constructs.

2626  *A.9.5.1   BRE/ERE Grammar Lexical Conventions*

2627      There is no additional rationale provided for this section.

2628  *A.9.5.2   RE and Bracket Expression Grammar*

2629      The removal of the *Back_open_paren Back_close_paren* option from the *nondupl_RE* specification is
2630      the result of PASC Interpretation 1003.2-92 #43 submitted for the ISO POSIX-2:1993 standard.
2631      Although the grammar required support for null subexpressions, this section does not describe
2632      the meaning of, and historical practice did not support, this construct.

2633  *A.9.5.3   ERE Grammar*

2634      There is no additional rationale provided for this section.

## A.10    Directory Structure and Devices

2635

### A.10.1    Directory Structure and Files

2636

2637  A description of the historical **/usr/tmp** was omitted, removing any concept of differences in
2638  emphasis between the / and **/usr** directories. The descriptions of **/bin**, **/usr/bin**, **/lib**, and **/usr/lib**
2639  were omitted because they are not useful for applications. In an early draft, a distinction was
2640  made between system and application directory usage, but this was not found to be useful.

2641  The directories / and **/dev** are included because the notion of a hierarchical directory structure is
2642  key to other information presented elsewhere in IEEE Std 1003.1-2001. In early drafts, it was
2643  argued that special devices and temporary files could conceivably be handled without a
2644  directory structure on some implementations. For example, the system could treat the characters
2645  `"/tmp"` as a special token that would store files using some non-POSIX file system structure.
2646  This notion was rejected by the standard developers, who required that all the files in this
2647  section be implemented via POSIX file systems.

2648  The **/tmp** directory is retained in IEEE Std 1003.1-2001 to accommodate historical applications
2649  that assume its availability. Implementations are encouraged to provide suitable directory
2650  names in the environment variable *TMPDIR* and applications are encouraged to use the contents
2651  of *TMPDIR* for creating temporary files.

2652  The standard files **/dev/null** and **/dev/tty** are required to be both readable and writable to allow
2653  applications to have the intended historical access to these files.

2654  The standard file **/dev/console** has been added for alignment with the Single UNIX Specification.

### A.10.2    Output Devices and Terminal Types

2655

2656  There is no additional rationale provided for this section.

## A.11    General Terminal Interface

2657

2658  If the implementation does not support this interface on any device types, it should behave as if
2659  it were being used on a device that is not a terminal device (in most cases *errno* will be set to
2660  [ENOTTY] on return from functions defined by this interface). This is based on the fact that
2661  many applications are written to run both interactively and in some non-interactive mode, and
2662  they adapt themselves at runtime. Requiring that they all be modified to test an environment
2663  variable to determine whether they should try to adapt is unnecessary. On a system that
2664  provides no general terminal interface, providing all the entry points as stubs that return
2665  [ENOTTY] (or an equivalent, as appropriate) has the same effect and requires no changes to the
2666  application.

2667  Although the needs of both interface implementors and application developers were addressed
2668  throughout IEEE Std 1003.1-2001, this section pays more attention to the needs of the latter. This
2669  is because, while many aspects of the programming interface can be hidden from the user by the
2670  application developer, the terminal interface is usually a large part of the user interface.
2671  Although to some extent the application developer can build missing features or work around
2672  inappropriate ones, the difficulties of doing that are greater in the terminal interface than
2673  elsewhere. For example, efficiency prohibits the average program from interpreting every
2674  character passing through it in order to simulate character erase, line kill, and so on. These
2675  functions should usually be done by the operating system, possibly at the interrupt level.

2676  The *tc*\*() functions were introduced as a way of avoiding the problems inherent in the
2677  traditional *ioctl*() function and in variants of it that were proposed. For example, *tcsetattr*() is

2678  specified in place of the use of the TCSETA *ioctl*( ) command function. This allows specification
2679  of all the arguments in a manner consistent with the ISO C standard unlike the varying third
2680  argument of *ioctl*( ), which is sometimes a pointer (to any of many different types) and
2681  sometimes an **int**.

2682  The advantages of this new method include:

2683  • It allows strict type checking.

2684  • The direction of transfer of control data is explicit.

2685  • Portable capabilities are clearly identified.

2686  • The need for a general interface routine is avoided.

2687  • Size of the argument is well-defined (there is only one type).

2688  The disadvantages include:

2689  • No historical implementation used the new method.

2690  • There are many small routines instead of one general-purpose one.

2691  • The historical parallel with *fcntl*( ) is broken.

2692  The issue of modem control was excluded from IEEE Std 1003.1-2001 on the grounds that:

2693  • It was concerned with setting and control of hardware timers.

2694  • The appropriate timers and settings vary widely internationally.

2695  • Feedback from European computer manufacturers indicated that this facility was not
2696    consistent with European needs and that specification of such a facility was not a
2697    requirement for portability.

## 2698 A.11.1  Interface Characteristics

### 2699 A.11.1.1  Opening a Terminal Device File

2700  There is no additional rationale provided for this section.

### 2701 A.11.1.2  Process Groups

2702  There is a potential race when the members of the foreground process group on a terminal leave
2703  that process group, either by exit or by changing process groups. After the last process exits the
2704  process group, but before the foreground process group ID of the terminal is changed (usually
2705  by a job control shell), it would be possible for a new process to be created with its process ID
2706  equal to the terminal's foreground process group ID. That process might then become the
2707  process group leader and accidentally be placed into the foreground on a terminal that was not
2708  necessarily its controlling terminal. As a result of this problem, the controlling terminal is
2709  defined to not have a foreground process group during this time.

2710  The cases where a controlling terminal has no foreground process group occur when all
2711  processes in the foreground process group either terminate and are waited for or join other
2712  process groups via *setpgid*( ) or *setsid*( ). If the process group leader terminates, this is the first
2713  case described; if it leaves the process group via *setpgid*( ), this is the second case described (a
2714  process group leader cannot successfully call *setsid*( )). When one of those cases causes a
2715  controlling terminal to have no foreground process group, it has two visible effects on
2716  applications. The first is the value returned by *tcgetpgrp*( ). The second (which occurs only in the
2717  case where the process group leader terminates) is the sending of signals in response to special
2718  input characters. The intent of IEEE Std 1003.1-2001 is that no process group be wrongly

2719    identified as the foreground process group by *tcgetpgrp*() or unintentionally receive signals
2720    because of placement into the foreground.

2721    In 4.3 BSD, the old process group ID continues to be used to identify the foreground process
2722    group and is returned by the function equivalent to *tcgetpgrp*(). In that implementation it is
2723    possible for a newly created process to be assigned the same value as a process ID and then form
2724    a new process group with the same value as a process group ID. The result is that the new
2725    process group would receive signals from this terminal for no apparent reason, and
2726    IEEE Std 1003.1-2001 precludes this by forbidding a process group from entering the foreground
2727    in this way. It would be more direct to place part of the requirement made by the last sentence
2728    under *fork*(), but there is no convenient way for that section to refer to the value that *tcgetpgrp*()
2729    returns, since in this case there is no process group and thus no process group ID.

2730    One possibility for a conforming implementation is to behave similarly to 4.3 BSD, but to
2731    prevent this reuse of the ID, probably in the implementation of *fork*(), as long as it is in use by
2732    the terminal.

2733    Another possibility is to recognize when the last process stops using the terminal's foreground
2734    process group ID, which is when the process group lifetime ends, and to change the terminal's
2735    foreground process group ID to a reserved value that is never used as a process ID or process
2736    group ID. (See the definition of *process group lifetime* in the definitions section.) The process ID
2737    can then be reserved until the terminal has another foreground process group.

2738    The 4.3 BSD implementation permits the leader (and only member) of the foreground process
2739    group to leave the process group by calling the equivalent of *setpgid*() and to later return,
2740    expecting to return to the foreground. There are no known application needs for this behavior,
2741    and IEEE Std 1003.1-2001 neither requires nor forbids it (except that it is forbidden for session
2742    leaders) by leaving it unspecified.

2743    *A.11.1.3  The Controlling Terminal*

2744    IEEE Std 1003.1-2001 does not specify a mechanism by which to allocate a controlling terminal.
2745    This is normally done by a system utility (such as *getty*) and is considered an administrative
2746    feature outside the scope of IEEE Std 1003.1-2001.

2747    Historical implementations allocate controlling terminals on certain *open*() calls. Since *open*() is
2748    part of POSIX.1, its behavior had to be dealt with. The traditional behavior is not required
2749    because it is not very straightforward or flexible for either implementations or applications.
2750    However, because of its prevalence, it was not practical to disallow this behavior either. Thus, a
2751    mechanism was standardized to ensure portable, predictable behavior in *open*().

2752    Some historical implementations deallocate a controlling terminal on the last system-wide close.
2753    This behavior in neither required nor prohibited. Even on implementations that do provide this
2754    behavior, applications generally cannot depend on it due to its system-wide nature.

2755    *A.11.1.4  Terminal Access Control*

2756    The access controls described in this section apply only to a process that is accessing its
2757    controlling terminal. A process accessing a terminal that is not its controlling terminal is
2758    effectively treated the same as a member of the foreground process group. While this may seem
2759    unintuitive, note that these controls are for the purpose of job control, not security, and job
2760    control relates only to a process' controlling terminal. Normal file access permissions handle
2761    security.

2762    If the process calling *read*() or *write*() is in a background process group that is orphaned, it is not
2763    desirable to stop the process group, as it is no longer under the control of a job control shell that
2764    could put it into the foreground again. Accordingly, calls to *read*() or *write*() functions by such

2765  processes receive an immediate error return. This is different from 4.2 BSD, which kills orphaned
2766  processes that receive terminal stop signals.

2767  The foreground/background/orphaned process group check performed by the terminal driver
2768  must be repeatedly performed until the calling process moves into the foreground or until the
2769  process group of the calling process becomes orphaned. That is, when the terminal driver
2770  determines that the calling process is in the background and should receive a job control signal,
2771  it sends the appropriate signal (SIGTTIN or SIGTTOU) to every process in the process group of
2772  the calling process and then it allows the calling process to immediately receive the signal. The
2773  latter is typically performed by blocking the process so that the signal is immediately noticed.
2774  Note, however, that after the process finishes receiving the signal and control is returned to the
2775  driver, the terminal driver must re-execute the foreground/background/orphaned process
2776  group check. The process may still be in the background, either because it was continued in the
2777  background by a job control shell, or because it caught the signal and did nothing.

2778  The terminal driver repeatedly performs the foreground/background/orphaned process group
2779  checks whenever a process is about to access the terminal. In the case of *write*() or the control
2780  *tc*\*() functions, the check is performed at the entry of the function. In the case of *read*(), the check
2781  is performed not only at the entry of the function, but also after blocking the process to wait for
2782  input characters (if necessary). That is, once the driver has determined that the process calling
2783  the *read*() function is in the foreground, it attempts to retrieve characters from the input queue. If
2784  the queue is empty, it blocks the process waiting for characters. When characters are available
2785  and control is returned to the driver, the terminal driver must return to the repeated
2786  foreground/background/orphaned process group check again. The process may have moved
2787  from the foreground to the background while it was blocked waiting for input characters.

2788  *A.11.1.5  Input Processing and Reading Data*

2789  There is no additional rationale provided for this section.

2790  *A.11.1.6  Canonical Mode Input Processing*

2791  The term ''character'' is intended here. ERASE should erase the last character, not the last byte.
2792  In the case of multi-byte characters, these two may be different.

2793  4.3 BSD has a WERASE character that erases the last ''word'' typed (but not any preceding
2794  <blank>s or <tab>s). A word is defined as a sequence of non-<blank>s, with <tab>s counted as
2795  <blank>s. Like ERASE, WERASE does not erase beyond the beginning of the line. This
2796  WERASE feature has not been specified in POSIX.1 because it is difficult to define in the
2797  international environment. It is only useful for languages where words are delimited by
2798  <blank>s. In some ideographic languages, such as Japanese and Chinese, words are not
2799  delimited at all. The WERASE character should presumably go back to the beginning of a
2800  sentence in those cases; practically, this means it would not be used much for those languages.

2801  It should be noted that there is a possible inherent deadlock if the application and
2802  implementation conflict on the value of {MAX_CANON}. With ICANON set (if IXOFF is
2803  enabled) and more than {MAX_CANON} characters transmitted without a <linefeed>,
2804  transmission will be stopped, the <linefeed> (or <carriage-return> when ICRLF is set) will never
2805  arrive, and the *read*() will never be satisfied.

2806  An application should not set IXOFF if it is using canonical mode unless it knows that (even in
2807  the face of a transmission error) the conditions described previously cannot be met or unless it is
2808  prepared to deal with the possible deadlock in some other way, such as timeouts.

2809  It should also be noted that this can be made to happen in non-canonical mode if the trigger
2810  value for sending IXOFF is less than VMIN and VTIME is zero.

2811  *A.11.1.7  Non-Canonical Mode Input Processing*

2812       Some points to note about MIN and TIME:

2813       1.  The interactions of MIN and TIME are not symmetric. For example, when MIN>0 and
2814           TIME=0, TIME has no effect. However, in the opposite case where MIN=0 and TIME>0,
2815           both MIN and TIME play a role in that MIN is satisfied with the receipt of a single
2816           character.

2817       2.  Also note that in case A (MIN>0, TIME>0), TIME represents an inter-character timer, while
2818           in case C (MIN=0, TIME>0), TIME represents a read timer.

2819       These two points highlight the dual purpose of the MIN/TIME feature. Cases A and B, where
2820       MIN>0, exist to handle burst-mode activity (for example, file transfer programs) where a
2821       program would like to process at least MIN characters at a time. In case A, the inter-character
2822       timer is activated by a user as a safety measure; in case B, it is turned off.

2823       Cases C and D exist to handle single-character timed transfers. These cases are readily adaptable
2824       to screen-based applications that need to know if a character is present in the input queue before
2825       refreshing the screen. In case C, the read is timed; in case D, it is not.

2826       Another important note is that MIN is always just a minimum. It does not denote a record
2827       length. That is, if a program does a read of 20 bytes, MIN is 10, and 25 characters are present, 20
2828       characters are returned to the user. In the special case of MIN=0, this still applies: if more than
2829       one character is available, they all will be returned immediately.

2830  *A.11.1.8  Writing Data and Output Processing*

2831       There is no additional rationale provided for this section.

2832  *A.11.1.9  Special Characters*

2833       There is no additional rationale provided for this section.

2834  *A.11.1.10 Modem Disconnect*

2835       There is no additional rationale provided for this section.

2836  *A.11.1.11 Closing a Terminal Device File*

2837       IEEE Std 1003.1-2001 does not specify that a *close*() on a terminal device file include the
2838       equivalent of a call to *tcflow*(*fd*,TCOON).

2839       An implementation that discards output at the time *close*() is called after reporting the return
2840       value to the *write*() call that data was written does not conform with IEEE Std 1003.1-2001. An
2841       application has functions such as *tcdrain*(), *tcflush*(), and *tcflow*() available to obtain the detailed
2842       behavior it requires with respect to flushing of output.

2843       At the time of the last close on a terminal device, an application relinquishes any ability to exert
2844       flow control via *tcflow*().

2845 **A.11.2   Parameters that Can be Set**

2846 *A.11.2.1 The termios Structure*

2847       This structure is part of an interface that, in general, retains the historic grouping of flags.
2848       Although a more optimal structure for implementations may be possible, the degree of change
2849       to applications would be significantly larger.

2850 *A.11.2.2 Input Modes*

2851       Some historical implementations treated a long break as multiple events, as many as one per
2852       character time. The wording in POSIX.1 explicitly prohibits this.

2853       Although the ISTRIP flag is normally superfluous with today's terminal hardware and software,
2854       it is historically supported. Therefore, applications may be using ISTRIP, and there is no
2855       technical problem with supporting this flag. Also, applications may wish to receive only 7-bit
2856       input bytes and may not be connected directly to the hardware terminal device (for example,
2857       when a connection traverses a network).

2858       Also, there is no requirement in general that the terminal device ensures that high-order bits
2859       beyond the specified character size are cleared. ISTRIP provides this function for 7-bit
2860       characters, which are common.

2861       In dealing with multi-byte characters, the consequences of a parity error in such a character, or in
2862       an escape sequence affecting the current character set, are beyond the scope of POSIX.1 and are
2863       best dealt with by the application processing the multi-byte characters.

2864 *A.11.2.3 Output Modes*

2865       POSIX.1 does not describe postprocessing of output to a terminal or detailed control of that from
2866       a conforming application. (That is, translation of <newline> to <carriage-return> followed by
2867       <linefeed> or <tab> processing.) There is nothing that a conforming application should do to its
2868       output for a terminal because that would require knowledge of the operation of the terminal. It
2869       is the responsibility of the operating system to provide postprocessing appropriate to the output
2870       device, whether it is a terminal or some other type of device.

2871       Extensions to POSIX.1 to control the type of postprocessing already exist and are expected to
2872       continue into the future. The control of these features is primarily to adjust the interface between
2873       the system and the terminal device so the output appears on the display correctly. This should
2874       be set up before use by any application.

2875       In general, both the input and output modes should not be set absolutely, but rather modified
2876       from the inherited state.

2877 *A.11.2.4 Control Modes*

2878       This section could be misread that the symbol ''CSIZE'' is a title in the **termios** *c_cflag* field.
2879       Although it does serve that function, it is also a required symbol, as a literal reading of POSIX.1
2880       (and the caveats about typography) would indicate.

2881 *A.11.2.5 Local Modes*

2882       Non-canonical mode is provided to allow fast bursts of input to be read efficiently while still
2883       allowing single-character input.

2884       The ECHONL function historically has been in many implementations. Since there seems to be
2885       no technical problem with supporting ECHONL, it is included in POSIX.1 to increase consensus.

2886  The alternate behavior possible when ECHOK or ECHOE are specified with ICANON is
2887  permitted as a compromise depending on what the actual terminal hardware can do. Erasing
2888  characters and lines is preferred, but is not always possible.

2889  *A.11.2.6  Special Control Characters*

2890  Permitting VMIN and VTIME to overlap with VEOF and VEOL was a compromise for historical
2891  implementations. Only when backwards-compatibility of object code is a serious concern to an
2892  implementor should an implementation continue this practice. Correct applications that work
2893  with the overlap (at the source level) should also work if it is not present, but not the reverse.


2894 # A.12  Utility Conventions

2895 ## A.12.1  Utility Argument Syntax

2896  The standard developers considered that recent trends toward diluting the SYNOPSIS sections
2897  of historical reference pages to the equivalent of:

2898  `command [`*`options`*`][`*`operands`*`]`

2899  were a disservice to the reader. Therefore, considerable effort was placed into rigorous
2900  definitions of all the command line arguments and their interrelationships. The relationships
2901  depicted in the synopses are normative parts of IEEE Std 1003.1-2001; this information is
2902  sometimes repeated in textual form, but that is only for clarity within context.

2903  The use of ''undefined'' for conflicting argument usage and for repeated usage of the same
2904  option is meant to prevent conforming applications from using conflicting arguments or
2905  repeated options unless specifically allowed (as is the case with *ls*, which allows simultaneous,
2906  repeated use of the −**C**, −**l**, and −**1** options). Many historical implementations will tolerate this
2907  usage, choosing either the first or the last applicable argument. This tolerance can continue, but
2908  conforming applications cannot rely upon it. (Other implementations may choose to print usage
2909  messages instead.)

2910  The use of ''undefined'' for conflicting argument usage also allows an implementation to make
2911  reasonable extensions to utilities where the implementor considers mutually-exclusive options
2912  according to IEEE Std 1003.1-2001 to have a sensible meaning and result.

2913  IEEE Std 1003.1-2001 does not define the result of a command when an option-argument or
2914  operand is not followed by ellipses and the application specifies more than one of that option-
2915  argument or operand. This allows an implementation to define valid (although non-standard)
2916  behavior for the utility when more than one such option or operand is specified.

2917  The following table summarizes the requirements for option-arguments:

| | SYNOPSIS Shows: | | |
|---|---|---|---|
| | −a *arg* | −b*arg* | −c [*arg*] |
| Conforming application uses: | −a *arg* | −b*arg* | −c*arg* or −c |
| System supports: | −a *arg* and −a*arg* | −b *arg* and −b*arg* | −c*arg* and −c |
| Non-conforming applications may use: | −a*arg* | −b *arg* | N/A |

2925  Allowing <blank>s after an option (that is, placing an option and its option-argument into
2926  separate argument strings) when IEEE Std 1003.1-2001 does not require it encourages portability
2927  of users, while still preserving backwards-compatibility of scripts. Inserting <blank>s between

2928    the option and the option-argument is preferred; however, historical usage has not been
2929    consistent in this area; therefore, <blank>s are required to be handled by all implementations,
2930    but implementations are also allowed to handle the historical syntax. Another justification for
2931    selecting the multiple-argument method was that the single-argument case is inherently
2932    ambiguous when the option-argument can legitimately be a null string.

2933    IEEE Std 1003.1-2001 explicitly states that digits are permitted as operands and option-
2934    arguments. The lower and upper bounds for the values of the numbers used for operands and
2935    option-arguments were derived from the ISO C standard values for {LONG_MIN} and
2936    {LONG_MAX}. The requirement on the standard utilities is that numbers in the specified range
2937    do not cause a syntax error, although the specification of a number need not be semantically
2938    correct for a particular operand or option-argument of a utility. For example, the specification of:

2939        dd obs=3000000000

2940    would yield undefined behavior for the application and could be a syntax error because the
2941    number 3 000 000 000 is outside of the range −2 147 483 647 to +2 147 483 647. On the other hand:

2942        dd obs=2000000000

2943    may cause some error, such as ''blocksize too large'', rather than a syntax error.

## 2944    A.12.2    Utility Syntax Guidelines

2945    This section is based on the rules listed in the SVID. It was included for two reasons:

2946    1.  The individual utility descriptions in the Shell and Utilities volume of
2947        IEEE Std 1003.1-2001, Chapter 4, Utilities needed a set of common (although not universal)
2948        actions on which they could anchor their descriptions of option and operand syntax. Most
2949        of the standard utilities actually do use these guidelines, and many of their historical
2950        implementations use the *getopt*( ) function for their parsing. Therefore, it was simpler to
2951        cite the rules and merely identify exceptions.

2952    2.  Writers of conforming applications need suggested guidelines if the POSIX community is
2953        to avoid the chaos of historical UNIX system command syntax.

2954    It is recommended that all *future* utilities and applications use these guidelines to enhance ''user
2955    portability''. The fact that some historical utilities could not be changed (to avoid breaking
2956    historical applications) should not deter this future goal.

2957    The voluntary nature of the guidelines is highlighted by repeated uses of the word *should*
2958    throughout. This usage should not be misinterpreted to imply that utilities that claim
2959    conformance in their OPTIONS sections do not always conform.

2960    Guidelines 1 and 2 are offered as guidance for locales using Latin alphabets. No
2961    recommendations are made by IEEE Std 1003.1-2001 concerning utility naming in other locales.

2962    In the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.9.1, Simple Commands, it is
2963    further stated that a command used in the Shell Command Language cannot be named with a
2964    trailing colon.

2965    Guideline 3 was changed to allow alphanumeric characters (letters and digits) from the character
2966    set to allow compatibility with historical usage. Historical practice allows the use of digits
2967    wherever practical, and there are no portability issues that would prohibit the use of digits. In
2968    fact, from an internationalization viewpoint, digits (being non-language-dependent) are
2969    preferable over letters (a −**2** is intuitively self-explanatory to any user, while in the −**f** *filename* the
2970    letter 'f' is a mnemonic aid only to speakers of Latin-based languages where ''filename''
2971    happens to translate to a word that begins with 'f'. Since guideline 3 still retains the word
2972    ''single'', multi-digit options are not allowed. Instances of historical utilities that used them have

been marked obsolescent, with the numbers being changed from option names to option-arguments.

It was difficult to achieve a satisfactory solution to the problem of name space in option characters. When the standard developers desired to extend the historical *cc* utility to accept ISO C standard programs, they found that all of the portable alphabet was already in use by various vendors. Thus, they had to devise a new name, *c89* (now superseded by *c99*), rather than something like *cc* –**X**. There were suggestions that implementors be restricted to providing extensions through various means (such as using a plus sign as the option delimiter or using option characters outside the alphanumeric set) that would reserve all of the remaining alphanumeric characters for future POSIX standards. These approaches were resisted because they lacked the historical style of UNIX systems. Furthermore, if a vendor-provided option should become commonly used in the industry, it would be a candidate for standardization. It would be desirable to standardize such a feature using historical practice for the syntax (the semantics can be standardized with any syntax). This would not be possible if the syntax was one reserved for the vendor. However, since the standardization process may lead to minor changes in the semantics, it may prove to be better for a vendor to use a syntax that will not be affected by standardization.

Guideline 8 includes the concept of comma-separated lists in a single argument. It is up to the utility to parse such a list itself because *getopt*() just returns the single string. This situation was retained so that certain historical utilities would not violate the guidelines. Applications preparing for international use should be aware of an occasional problem with comma-separated lists: in some locales, the comma is used as the radix character. Thus, if an application is preparing operands for a utility that expects a comma-separated list, it should avoid generating non-integer values through one of the means that is influenced by setting the *LC_NUMERIC* variable (such as *awk*, *bc*, *printf*, or *printf*()).

Applications calling any utility with a first operand starting with ′−′ should usually specify −−, as indicated by Guideline 10, to mark the end of the options. This is true even if the SYNOPSIS in the Shell and Utilities volume of IEEE Std 1003.1-2001 does not specify any options; implementations may provide options as extensions to the Shell and Utilities volume of IEEE Std 1003.1-2001. The standard utilities that do not support Guideline 10 indicate that fact in the OPTIONS section of the utility description.

Guideline 11 was modified to clarify that the order of different options should not matter relative to one another. However, the order of repeated options that also have option-arguments may be significant; therefore, such options are required to be interpreted in the order that they are specified. The *make* utility is an instance of a historical utility that uses repeated options in which the order is significant. Multiple files are specified by giving multiple instances of the −**f** option; for example:

```
make −f common_header −f specific_rules target
```

Guideline 13 does not imply that all of the standard utilities automatically accept the operand ′−′ to mean standard input or output, nor does it specify the actions of the utility upon encountering multiple ′−′ operands. It simply says that, by default, ′−′ operands are not used for other purposes in the file reading or writing (but not when using *stat*(), *unlink*(), *touch*, and so on) utilities. All information concerning actual treatment of the ′−′ operand is found in the individual utility sections.

An area of concern was that as implementations mature, implementation-defined utilities and implementation-defined utility options will result. The idea was expressed that there needed to be a standard way, say an environment variable or some such mechanism, to identify implementation-defined utilities separately from standard utilities that may have the same name. It was decided that there already exist several ways of dealing with this situation and that

3022 it is outside of the scope to attempt to standardize in the area of non-standard items. A method
3023 that exists on some historical implementations is the use of the so-called **/local/bin** or
3024 **/usr/local/bin** directory to separate local or additional copies or versions of utilities. Another
3025 method that is also used is to isolate utilities into completely separate domains. Still another
3026 method to ensure that the desired utility is being used is to request the utility by its full
3027 pathname. There are many approaches to this situation; the examples given above serve to
3028 illustrate that there is more than one.

3029 ## A.13   Headers

3030 ### A.13.1   Format of Entries

3031 Each header reference page has a common layout of sections describing the interface. This layout
3032 is similar to the manual page or ''man'' page format shipped with most UNIX systems, and each
3033 header has sections describing the SYNOPSIS and DESCRIPTION. These are the two sections
3034 that relate to conformance.

3035 Additional sections are informative, and add considerable information for the application
3036 developer. APPLICATION USAGE sections provide additional caveats, issues, and
3037 recommendations to the developer. RATIONALE sections give additional information on the
3038 decisions made in defining the interface.

3039 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
3040 the future, and often cautions the developer to architect the code to account for a change in this
3041 area. Note that a future directions statement should not be taken as a commitment to adopt a
3042 feature or interface in the future.

3043 The CHANGE HISTORY section describes when the interface was introduced, and how it has
3044 changed.

3045 Option labels and margin markings in the page can be useful in guiding the application
3046 developer.

# Rationale (Informative)

3048 **Part B:**

3049 **System Interfaces**

3050 *The Open Group*
3051 *The Institute of Electrical and Electronics Engineers, Inc.*

# Rationale for System Interfaces

3052

## B.1 Introduction

### B.1.1 Scope

Refer to Section A.1.1 (on page 3).

### B.1.2 Conformance

Refer to Section A.2 (on page 9).

### B.1.3 Normative References

There is no additional rationale provided for this section.

### B.1.4 Change History

The change history is provided as an informative section, to track changes from previous issues of IEEE Std 1003.1-2001.

The following sections describe changes made to the System Interfaces volume of IEEE Std 1003.1-2001 since Issue 5 of the base document. The CHANGE HISTORY section for each entry details the technical changes that have been made to that entry from Issue 5. Changes between earlier issues of the base document and Issue 5 are not included.

The change history between Issue 5 and Issue 6 also lists the changes since the ISO POSIX-1: 1996 standard.

**Changes from Issue 5 to Issue 6 (IEEE Std 1003.1-2001)**

The following list summarizes the major changes that were made in the System Interfaces volume of IEEE Std 1003.1-2001 from Issue 5 to Issue 6:

- This volume of IEEE Std 1003.1-2001 is extensively revised so that it can be both an IEEE POSIX Standard and an Open Group Technical Standard.
- The POSIX System Interfaces requirements incorporate support of FIPS 151-2.
- The POSIX System Interfaces requirements are updated to align with some features of the Single UNIX Specification.
- A RATIONALE section is added to each reference page.
- Networking interfaces from the XNS, Issue 5.2 specification are incorporated.
- IEEE Std 1003.1d-1999 is incorporated.
- IEEE Std 1003.1j-2000 is incorporated.
- IEEE Std 1003.1q-2000 is incorporated.
- IEEE P1003.1a draft standard is incorporated.

3083        • Existing functionality is aligned with the ISO/IEC 9899: 1999 standard.

3084        • New functionality from the ISO/IEC 9899: 1999 standard is incorporated.

3085        • IEEE PASC Interpretations are applied.

3086        • The Open Group corrigenda and resolutions are applied.

3087    **New Features in Issue 6**

3088    The functions first introduced in Issue 6 (over the Issue 5 Base document) are listed in the table
3089    below:

3090

| New Functions in Issue 6 | | |
|---|---|---|
| *acosf*( ) | *catanhl*( ) | *cprojf*( ) |
| *acoshf*( ) | *catanl*( ) | *cprojl*( ) |
| *acoshl*( ) | *cbrtf*( ) | *creal*( ) |
| *acosl*( ) | *cbrtl*( ) | *crealf*( ) |
| *asinf*( ) | *ccos*( ) | *creall*( ) |
| *asinhf*( ) | *ccosf*( ) | *csin*( ) |
| *asinhl*( ) | *ccosh*( ) | *csinf*( ) |
| *asinl*( ) | *ccoshf*( ) | *csinh*( ) |
| *atan2f*( ) | *ccoshl*( ) | *csinhf*( ) |
| *atan2l*( ) | *ccosl*( ) | *csinhl*( ) |
| *atanf*( ) | *ceilf*( ) | *csinl*( ) |
| *atanhf*( ) | *ceill*( ) | *csqrt*( ) |
| *atanhl*( ) | *cexp*( ) | *csqrtf*( ) |
| *atanl*( ) | *cexpf*( ) | *csqrtl*( ) |
| *atoll*( ) | *cexpl*( ) | *ctan*( ) |
| *cabs*( ) | *cimag*( ) | *ctanf*( ) |
| *cabsf*( ) | *cimagf*( ) | *ctanh*( ) |
| *cabsl*( ) | *cimagl*( ) | *ctanhf*( ) |
| *cacos*( ) | *clock_getcpuclockid*( ) | *ctanhl*( ) |
| *cacosf*( ) | *clock_nanosleep*( ) | *ctanl*( ) |
| *cacosh*( ) | *clog*( ) | *erfcf*( ) |
| *cacoshf*( ) | *clogf*( ) | *erfcl*( ) |
| *cacoshl*( ) | *clogl*( ) | *erff*( ) |
| *cacosl*( ) | *conj*( ) | *erfl*( ) |
| *carg*( ) | *conjf*( ) | *exp2*( ) |
| *cargf*( ) | *conjl*( ) | *exp2f*( ) |
| *cargl*( ) | *copysign*( ) | *exp2l*( ) |
| *casin*( ) | *copysignf*( ) | *expf*( ) |
| *casinf*( ) | *copysignl*( ) | *expl*( ) |
| *casinh*( ) | *cosf*( ) | *expm1f*( ) |
| *casinhf*( ) | *coshf*( ) | *expm1l*( ) |
| *casinhl*( ) | *coshl*( ) | *fabsf*( ) |
| *casinl*( ) | *cosl*( ) | *fabsl*( ) |
| *catan*( ) | *cpow*( ) | *fdim*( ) |
| *catanf*( ) | *cpowf*( ) | *fdimf*( ) |
| *catanh*( ) | *cpowl*( ) | *fdiml*( ) |
| *catanhf*( ) | *cproj*( ) | *feclearexcept*( ) |

| | New Functions in Issue 6 | |
|---|---|---|
| *fegetenv*( ) | *ldexpl*( ) | *posix_fallocate*( ) |
| *fegetexceptflag*( ) | *lgammaf*( ) | *posix_madvise*( ) |
| *fegetround*( ) | *lgammal*( ) | *posix_mem_offset*( ) |
| *feholdexcept*( ) | *llabs*( ) | *posix_memalign*( ) |
| *feraiseexcept*( ) | *lldiv*( ) | *posix_openpt*( ) |
| *fesetenv*( ) | *llrint*( ) | *posix_spawn*( ) |
| *fesetexceptflag*( ) | *llrintf*( ) | *posix_spawn_file_actions_addclose*( ) |
| *fesetround*( ) | *llrintl*( ) | *posix_spawn_file_actions_adddup2*( ) |
| *fetestexcept*( ) | *llround*( ) | *posix_spawn_file_actions_addopen*( ) |
| *feupdateenv*( ) | *llroundf*( ) | *posix_spawn_file_actions_destroy*( ) |
| *floorf*( ) | *llroundl*( ) | *posix_spawn_file_actions_init*( ) |
| *floorl*( ) | *log10f*( ) | *posix_spawnattr_destroy*( ) |
| *fma*( ) | *log10l*( ) | *posix_spawnattr_getflags*( ) |
| *fmaf*( ) | *log1pf*( ) | *posix_spawnattr_getpgroup*( ) |
| *fmal*( ) | *log1pl*( ) | *posix_spawnattr_getschedparam*( ) |
| *fmax*( ) | *log2*( ) | *posix_spawnattr_getschedpolicy*( ) |
| *fmaxf*( ) | *log2f*( ) | *posix_spawnattr_getsigdefault*( ) |
| *fmaxl*( ) | *log2l*( ) | *posix_spawnattr_getsigmask*( ) |
| *fmin*( ) | *logbf*( ) | *posix_spawnattr_init*( ) |
| *fminf*( ) | *logbl*( ) | *posix_spawnattr_setflags*( ) |
| *fminl*( ) | *logf*( ) | *posix_spawnattr_setpgroup*( ) |
| *fmodf*( ) | *logl*( ) | *posix_spawnattr_setschedparam*( ) |
| *fmodl*( ) | *lrint*( ) | *posix_spawnattr_setschedpolicy*( ) |
| *fpclassify*( ) | *lrintf*( ) | *posix_spawnattr_setsigdefault*( ) |
| *frexpf*( ) | *lrintl*( ) | *posix_spawnattr_setsigmask*( ) |
| *frexpl*( ) | *lround*( ) | *posix_spawnp*( ) |
| *hypotf*( ) | *lroundf*( ) | *posix_trace_attr_destroy*( ) |
| *hypotl*( ) | *lroundl*( ) | *posix_trace_attr_getclockres*( ) |
| *ilogbf*( ) | *modff*( ) | *posix_trace_attr_getcreatetime*( ) |
| *ilogbl*( ) | *modfl*( ) | *posix_trace_attr_getgenversion*( ) |
| *imaxabs*( ) | *mq_timedreceive*( ) | *posix_trace_attr_getinherited*( ) |
| *imaxdiv*( ) | *mq_timedsend*( ) | *posix_trace_attr_getlogfullpolicy*( ) |
| *isblank*( ) | *nan*( ) | *posix_trace_attr_getlogsize*( ) |
| *isfinite*( ) | *nanf*( ) | *posix_trace_attr_getmaxdatasize*( ) |
| *isgreater*( ) | *nanl*( ) | *posix_trace_attr_getmaxsystemeventsize*( ) |
| *isgreaterequal*( ) | *nearbyint*( ) | *posix_trace_attr_getmaxusereventsize*( ) |
| *isinf*( ) | *nearbyintf*( ) | *posix_trace_attr_getname*( ) |
| *isless*( ) | *nearbyintl*( ) | *posix_trace_attr_getstreamfullpolicy*( ) |
| *islessequal*( ) | *nextafterf*( ) | *posix_trace_attr_getstreamsize*( ) |
| *islessgreater*( ) | *nextafterl*( ) | *posix_trace_attr_init*( ) |
| *isnormal*( ) | *nexttoward*( ) | *posix_trace_attr_setinherited*( ) |
| *isunordered*( ) | *nexttowardf*( ) | *posix_trace_attr_setlogfullpolicy*( ) |
| *iswblank*( ) | *nexttowardl*( ) | *posix_trace_attr_setlogsize*( ) |
| *ldexpf*( ) | *posix_fadvise*( ) | *posix_trace_create*( ) |

3175

| New Functions in Issue 6 | | |
|---|---|---|
| *posix_trace_attr_setmaxdatasize*( ) | *pthread_barrier_destroy*( ) | *signbit*( ) |
| *posix_trace_attr_setname*( ) | *pthread_barrier_init*( ) | *sinf*( ) |
| *posix_trace_attr_setstreamfullpolicy*( ) | *pthread_barrier_wait*( ) | *sinhf*( ) |
| *posix_trace_attr_setstreamsize*( ) | *pthread_barrierattr_destroy*( ) | *sinhl*( ) |
| *posix_trace_clear*( ) | *pthread_barrierattr_getpshared*( ) | *sinl*( ) |
| *posix_trace_close*( ) | *pthread_barrierattr_init*( ) | *sockatmark*( ) |
| *posix_trace_create_withlog*( ) | *pthread_barrierattr_setpshared*( ) | *sqrtf*( ) |
| *posix_trace_event*( ) | *pthread_condattr_getclock*( ) | *sqrtl*( ) |
| *posix_trace_eventid_equal*( ) | *pthread_condattr_setclock*( ) | *strerror_r*( ) |
| *posix_trace_eventid_get_name*( ) | *pthread_getcpuclockid*( ) | *strtoimax*( ) |
| *posix_trace_eventid_open*( ) | *pthread_mutex_timedlock*( ) | *strtoll*( ) |
| *posix_trace_eventset_add*( ) | *pthread_rwlock_timedrdlock*( ) | *strtoull*( ) |
| *posix_trace_eventset_del*( ) | *pthread_rwlock_timedwrlock*( ) | *strtoumax*( ) |
| *posix_trace_eventset_empty*( ) | *pthread_setschedprio*( ) | *tanf*( ) |
| *posix_trace_eventset_fill*( ) | *pthread_spin_destroy*( ) | *tanhf*( ) |
| *posix_trace_eventset_ismember*( ) | *pthread_spin_init*( ) | *tanhl*( ) |
| *posix_trace_eventtypelist_getnext_id*( ) | *pthread_spin_lock*( ) | *tanl*( ) |
| *posix_trace_eventtypelist_rewind*( ) | *pthread_spin_trylock*( ) | *tgamma*( ) |
| *posix_trace_flush*( ) | *pthread_spin_unlock*( ) | *tgammaf*( ) |
| *posix_trace_get_attr*( ) | *remainderf*( ) | *tgammal*( ) |
| *posix_trace_get_filter*( ) | *remainderl*( ) | *trunc*( ) |
| *posix_trace_get_status*( ) | *remquo*( ) | *truncf*( ) |
| *posix_trace_getnext_event*( ) | *remquof*( ) | *truncl*( ) |
| *posix_trace_open*( ) | *remquol*( ) | *unsetenv*( ) |
| *posix_trace_rewind*( ) | *rintf*( ) | *vfprintf*( ) |
| *posix_trace_set_filter*( ) | *rintl*( ) | *vfscanf*( ) |
| *posix_trace_shutdown*( ) | *round*( ) | *vfwscanf*( ) |
| *posix_trace_start*( ) | *roundf*( ) | *vprintf*( ) |
| *posix_trace_stop*( ) | *roundl*( ) | *vscanf*( ) |
| *posix_trace_timedgetnext_event*( ) | *scalbln*( ) | *vsnprintf*( ) |
| *posix_trace_trid_eventid_open*( ) | *scalblnf*( ) | *vsprintf*( ) |
| *posix_trace_trygetnext_event*( ) | *scalblnl*( ) | *vsscanf*( ) |
| *posix_typed_mem_get_info*( ) | *scalbn*( ) | *vswscanf*( ) |
| *posix_typed_mem_open*( ) | *scalbnf*( ) | *vwscanf*( ) |
| *powf*( ) | *scalbnl*( ) | *wcstoimax*( ) |
| *powl*( ) | *sem_timedwait*( ) | *wcstoll*( ) |
| *pselect*( ) | *setegid*( ) | *wcstoull*( ) |
| *pthread_attr_getstack*( ) | *setenv*( ) | *wcstoumax*( ) |
| *pthread_attr_setstack*( ) | *seteuid*( ) | |

3216    The following new headers are introduced in Issue 6:

3217

| New Headers in Issue 6 | | |
|---|---|---|
| **<complex.h>** | **<spawn.h>** | **<tgmath.h>** |
| **<fenv.h>** | **<stdbool.h>** | **<trace.h>** |
| **<net/if.h>** | **<stdint.h>** | |

3222   The following table lists the functions and symbols from the XSI extension. These are new since
3223   the ISO POSIX-1: 1996 standard.
3224

| New XSI Functions and Symbols in Issue 6 | | | |
|---|---|---|---|
| *_longjmp*( ) | *getcontext*( ) | *msgget*( ) | *setstate*( ) |
| *_setjmp*( ) | *getdate*( ) | *msgrcv*( ) | *setutxent*( ) |
| *_tolower*( ) | *getgrent*( ) | *msgsnd*( ) | *shmat*( ) |
| *_toupper*( ) | *gethostid*( ) | *nftw*( ) | *shmctl*( ) |
| *a64l*( ) | *getitimer*( ) | *nice*( ) | *shmdt*( ) |
| *basename*( ) | *getpgid*( ) | *nl_langinfo*( ) | *shmget*( ) |
| *bcmp*( ) | *getpmsg*( ) | *nrand48*( ) | *sigaltstack*( ) |
| *bcopy*( ) | *getpriority*( ) | *openlog*( ) | *sighold*( ) |
| *bzero*( ) | *getpwent*( ) | *poll*( ) | *sigignore*( ) |
| *catclose*( ) | *getrlimit*( ) | *posix_openpt*( ) | *siginterrupt*( ) |
| *catgets*( ) | *getrusage*( ) | *pread*( ) | *sigpause*( ) |
| *catopen*( ) | *getsid*( ) | *pthread_attr_getguardsize*( ) | *sigrelse*( ) |
| *closelog*( ) | *getsubopt*( ) | *pthread_attr_setguardsize*( ) | *sigset*( ) |
| *crypt*( ) | *gettimeofday*( ) | *pthread_attr_setstack*( ) | *srand48*( ) |
| *daylight* | *getutxent*( ) | *pthread_getconcurrency*( ) | *srandom*( ) |
| *dbm_clearerr*( ) | *getutxid*( ) | *pthread_mutexattr_gettype*( ) | *statvfs*( ) |
| *dbm_close*( ) | *getutxline*( ) | *pthread_mutexattr_settype*( ) | *strcasecmp*( ) |
| *dbm_delete*( ) | *getwd*( ) | *pthread_rwlockattr_init*( ) | *strdup*( ) |
| *dbm_error*( ) | *grantpt*( ) | *pthread_rwlockattr_setpshared*( ) | *strfmon*( ) |
| *dbm_fetch*( ) | *hcreate*( ) | *pthread_setconcurrency*( ) | *strncasecmp*( ) |
| *dbm_firstkey*( ) | *hdestroy*( ) | *ptsname*( ) | *strptime*( ) |
| *dbm_nextkey*( ) | *hsearch*( ) | *putenv*( ) | *swab*( ) |
| *dbm_open*( ) | *iconv*( ) | *pututxline*( ) | *swapcontext*( ) |
| *dbm_store*( ) | *iconv_close*( ) | *pwrite*( ) | *sync*( ) |
| *dirname*( ) | *iconv_open*( ) | *random*( ) | *syslog*( ) |
| *dlclose*( ) | *index*( ) | *readv*( ) | *tcgetsid*( ) |
| *dlerror*( ) | *initstate*( ) | *realpath*( ) | *tdelete*( ) |
| *dlopen*( ) | *insque*( ) | *remque*( ) | *telldir*( ) |
| *dlsym*( ) | *isascii*( ) | *rindex*( ) | *tempnam*( ) |
| *drand48*( ) | *jrand48*( ) | *seed48*( ) | *tfind*( ) |
| *ecvt*( ) | *killpg*( ) | *seekdir*( ) | *timezone* |
| *encrypt*( ) | *l64a*( ) | *semctl*( ) | *toascii*( ) |
| *endgrent*( ) | *lchown*( ) | *semget*( ) | *truncate*( ) |
| *endpwent*( ) | *lcong48*( ) | *semop*( ) | *tsearch*( ) |
| *endutxent*( ) | *lfind*( ) | *setcontext*( ) | *twalk*( ) |
| *erand48*( ) | *lockf*( ) | *setgrent*( ) | *ulimit*( ) |
| *fchdir*( ) | *lrand48*( ) | *setitimer*( ) | *unlockpt*( ) |
| *fcvt*( ) | *lsearch*( ) | *setkey*( ) | *utimes*( ) |
| *ffs*( ) | *makecontext*( ) | *setlogmask*( ) | *waitid*( ) |
| *fmtmsg*( ) | *memccpy*( ) | *setpgrp*( ) | *wcswcs*( ) |
| *fstatvfs*( ) | *mknod*( ) | *setpriority*( ) | *wcswidth*( ) |
| *ftime*( ) | *mkstemp*( ) | *setpwent*( ) | *wcwidth*( ) |
| *ftok*( ) | *mktemp*( ) | *setregid*( ) | *writev*( ) |
| *ftw*( ) | *mrand48*( ) | *setreuid*( ) | |
| *gcvt*( ) | *msgctl*( ) | *setrlimit*( ) | |

3271         The following table lists the headers from the XSI extension. These are new since the
3272         ISO POSIX-1:1996 standard.

3273
3274

| New XSI Headers in Issue 6 | | |
|---|---|---|
| **<cpio.h>** | **<poll.h>** | **<sys/statvfs.h>** |
| **<dlfcn.h>** | **<search.h>** | **<sys/time.h>** |
| **<fmtmsg.h>** | **<strings.h>** | **<sys/timeb.h>** |
| **<ftw.h>** | **<stropts.h>** | **<sys/uio.h>** |
| **<iconv.h>** | **<sys/ipc.h>** | **<syslog.h>** |
| **<langinfo.h>** | **<sys/mman.h>** | **<ucontext.h>** |
| **<libgen.h>** | **<sys/msg.h>** | **<ulimit.h>** |
| **<monetary.h>** | **<sys/resource.h>** | **<utmpx.h>** |
| **<ndbm.h>** | **<sys/sem.h>** | |
| **<nl_types.h>** | **<sys/shm.h>** | |

(line numbers 3275–3284 correspond to the table rows)

3285 **B.1.5 Terminology**

3286         Refer to Section A.1.4 (on page 5).

3287 **B.1.6 Definitions**

3288         Refer to Section A.3 (on page 13).

3289 **B.1.7 Relationship to Other Formal Standards**

3290         There is no additional rationale provided for this section.

3291 **B.1.8 Portability**

3292         Refer to Section A.1.5 (on page 8).

3293 *B.1.8.1 Codes*

3294         Refer to Section A.1.5.1 (on page 8).

3295 **B.1.9 Format of Entries**

3296         Each system interface reference page has a common layout of sections describing the interface.
3297         This layout is similar to the manual page or ''man'' page format shipped with most UNIX
3298         systems, and each header has sections describing the SYNOPSIS, DESCRIPTION, RETURN
3299         VALUE, and ERRORS. These are the four sections that relate to conformance.

3300         Additional sections are informative, and add considerable information for the application
3301         developer. EXAMPLES sections provide example usage. APPLICATION USAGE sections
3302         provide additional caveats, issues, and recommendations to the developer. RATIONALE
3303         sections give additional information on the decisions made in defining the interface.

3304         FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
3305         the future, and often cautions the developer to architect the code to account for a change in this
3306         area. Note that a future directions statement should not be taken as a commitment to adopt a
3307         feature or interface in the future.

3308         The CHANGE HISTORY section describes when the interface was introduced, and how it has
3309         changed.

3310    Option labels and margin markings in the page can be useful in guiding the application
3311    developer.

## B.2    General Information

### B.2.1    Use and Implementation of Functions

3314    The information concerning the use of functions was adapted from a description in the ISO C
3315    standard. Here is an example of how an application program can protect itself from functions
3316    that may or may not be macros, rather than true functions:

3317    The *atoi*() function may be used in any of several ways:

3318        • By use of its associated header (possibly generating a macro expansion):

```
3319        #include <stdlib.h>
3320        /* ... */
3321        i = atoi(str);
```

3322        • By use of its associated header (assuredly generating a true function call):

```
3323        #include <stdlib.h>
3324        #undef atoi
3325        /* ... */
3326        i = atoi(str);
```

3327        or:

```
3328        #include <stdlib.h>
3329        /* ... */
3330        i = (atoi) (str);
```

3331        • By explicit declaration:

```
3332        extern int atoi (const char *);
3333        /* ... */
3334        i = atoi(str);
```

3335        • By implicit declaration:

```
3336        /* ... */
3337        i = atoi(str);
```

3338        (Assuming no function prototype is in scope. This is not allowed by the ISO C standard for
3339        functions with variable arguments; furthermore, parameter type conversion ''widening'' is
3340        subject to different rules in this case.)

3341    Note that the ISO C standard reserves names starting with '_' for the compiler. Therefore, the
3342    compiler could, for example, implement an intrinsic, built-in function *_asm_builtin_atoi*(), which
3343    it recognized and expanded into inline assembly code. Then, in **<stdlib.h>**, there could be the
3344    following:

```
3345        #define atoi(X) _asm_builtin_atoi(X)
```

3346    The user's ''normal'' call to *atoi*() would then be expanded inline, but the implementor would
3347    also be required to provide a callable function named *atoi*() for use when the application
3348    requires it; for example, if its address is to be stored in a function pointer variable.

3349 **B.2.2**     **The Compilation Environment**

3350 *B.2.2.1*    *POSIX.1 Symbols*

3351 This and the following section address the issue of ''name space pollution''. The ISO C standard
3352 requires that the name space beyond what it reserves not be altered except by explicit action of
3353 the application writer. This section defines the actions to add the POSIX.1 symbols for those
3354 headers where both the ISO C standard and POSIX.1 need to define symbols, and also where the
3355 XSI extension extends the base standard.

3356 When headers are used to provide symbols, there is a potential for introducing symbols that the
3357 application writer cannot predict. Ideally, each header should only contain one set of symbols,
3358 but this is not practical for historical reasons. Thus, the concept of feature test macros is
3359 included. Two feature test macros are explicitly defined by IEEE Std 1003.1-2001; it is expected
3360 that future revisions may add to this.

3361 **Note:**      Feature test macros allow an application to announce to the implementation its desire to have
3362                   certain symbols and prototypes exposed. They should not be confused with the version test
3363                   macros and constants for options in **<unistd.h>** which are the implementation's way of
3364                   announcing functionality to the application.

3365 It is further intended that these feature test macros apply only to the headers specified by
3366 IEEE Std 1003.1-2001. Implementations are expressly permitted to make visible symbols not
3367 specified by IEEE Std 1003.1-2001, within both POSIX.1 and other headers, under the control of
3368 feature test macros that are not defined by IEEE Std 1003.1-2001.

3369 **The _POSIX_C_SOURCE Feature Test Macro**

3370 Since _POSIX_SOURCE specified by the POSIX.1-1990 standard did not have a value associated
3371 with it, the _POSIX_C_SOURCE macro replaces it, allowing an application to inform the system
3372 of the revision of the standard to which it conforms. This symbol will allow implementations to
3373 support various revisions of IEEE Std 1003.1-2001 simultaneously. For instance, when either
3374 _POSIX_SOURCE is defined or _POSIX_C_SOURCE is defined as 1, the system should make
3375 visible the same name space as permitted and required by the POSIX.1-1990 standard. When
3376 _POSIX_C_SOURCE is defined, the state of _POSIX_SOURCE is completely irrelevant.

3377 It is expected that C bindings to future POSIX standards will define new values for
3378 _POSIX_C_SOURCE, with each new value reserving the name space for that new standard, plus
3379 all earlier POSIX standards.

3380 **The _XOPEN_SOURCE Feature Test Macro**

3381 The feature test macro _XOPEN_SOURCE is provided as the announcement mechanism for the
3382 application that it requires functionality from the Single UNIX Specification. _XOPEN_SOURCE
3383 must be defined to the value 600 before the inclusion of any header to enable the functionality in
3384 the Single UNIX Specification. Its definition subsumes the use of _POSIX_SOURCE and
3385 _POSIX_C_SOURCE.

3386 An extract of code from a conforming application, that appears before any **#include** statements,
3387 is given below:

3388 ```
#define _XOPEN_SOURCE 600 /* Single UNIX Specification, Version 3 */
```

3389 ```
#include ...
```

3390 Note that the definition of _XOPEN_SOURCE with the value 600 makes the definition of
3391 _POSIX_C_SOURCE redundant and it can safely be omitted.

*B.2.2.2    The Name Space*

3393 The reservation of identifiers is paraphrased from the ISO C standard. The text is included
3394 because it needs to be part of IEEE Std 1003.1-2001, regardless of possible changes in future
3395 versions of the ISO C standard.

3396 These identifiers may be used by implementations, particularly for feature test macros.
3397 Implementations should not use feature test macro names that might be reasonably used by a
3398 standard.

3399 Including headers more than once is a reasonably common practice, and it should be carried
3400 forward from the ISO C standard. More significantly, having definitions in more than one
3401 header is explicitly permitted. Where the potential declaration is ''benign'' (the same definition
3402 twice) the declaration can be repeated, if that is permitted by the compiler. (This is usually true
3403 of macros, for example.) In those situations where a repetition is not benign (for example,
3404 **typedefs**), conditional compilation must be used. The situation actually occurs both within the
3405 ISO C standard and within POSIX.1: **time_t** should be in **<sys/types.h>**, and the ISO C standard
3406 mandates that it be in **<time.h>**.

3407 The area of name space pollution *versus* additions to structures is difficult because of the macro
3408 structure of C. The following discussion summarizes all the various problems with and
3409 objections to the issue.

3410 Note the phrase ''user-defined macro''. Users are not permitted to define macro names (or any
3411 other name) beginning with `"_[A-Z_]"`. Thus, the conflict cannot occur for symbols reserved
3412 to the vendor's name space, and the permission to add fields automatically applies, without
3413 qualification, to those symbols.

3414    1.   Data structures (and unions) need to be defined in headers by implementations to meet
3415         certain requirements of POSIX.1 and the ISO C standard.

3416    2.   The structures defined by POSIX.1 are typically minimal, and any practical
3417         implementation would wish to add fields to these structures either to hold additional
3418         related information or for backwards-compatibility (or both). Future standards (and *de
3419         facto* standards) would also wish to add to these structures. Issues of field alignment make
3420         it impractical (at least in the general case) to simply omit fields when they are not defined
3421         by the particular standard involved.

3422         The **dirent** structure is an example of such a minimal structure (although one could argue
3423         about whether the other fields need visible names). The *st_rdev* field of most
3424         implementations' **stat** structure is a common example where extension is needed and
3425         where a conflict could occur.

3426    3.   Fields in structures are in an independent name space, so the addition of such fields
3427         presents no problem to the C language itself in that such names cannot interact with
3428         identically named user symbols because access is qualified by the specific structure name.

3429    4.   There is an exception to this: macro processing is done at a lexical level. Thus, symbols
3430         added to a structure might be recognized as user-provided macro names at the location
3431         where the structure is declared. This only can occur if the user-provided name is declared
3432         as a macro before the header declaring the structure is included. The user's use of the name
3433         after the declaration cannot interfere with the structure because the symbol is hidden and
3434         only accessible through access to the structure. Presumably, the user would not declare
3435         such a macro if there was an intention to use that field name.

3436    5.   Macros from the same or a related header might use the additional fields in the structure,
3437         and those field names might also collide with user macros. Although this is a less frequent
3438         occurrence, since macros are expanded at the point of use, no constraint on the order of use

3439          of names can apply.

3440    6.    An ''obvious'' solution of using names in the reserved name space and then redefining
3441          them as macros when they should be visible does not work because this has the effect of
3442          exporting the symbol into the general name space. For example, given a (hypothetical)
3443          system-provided header <**h.h**>, and two parts of a C program in **a.c** and **b.c**, in header
3444          <**h.h**>:

3445          ```
              struct foo {
3446              int __i;
3447          }
              ```
3448          ```
              #ifdef _FEATURE_TEST
3449          #define i __i;
3450          #endif
              ```

3451          In file **a.c**:

3452          ```
              #include h.h
3453          extern int i;
3454          ...
              ```

3455          In file **b.c**:

3456          ```
              extern int i;
3457          ...
              ```

3458          The symbol that the user thinks of as *i* in both files has an external name of _ _*i* in **a.c**; the
3459          same symbol *i* in **b.c** has an external name *i* (ignoring any hidden manipulations the
3460          compiler might perform on the names). This would cause a mysterious name resolution
3461          problem when **a.o** and **b.o** are linked.

3462          Simply avoiding definition then causes alignment problems in the structure.

3463          A structure of the form:

3464          ```
              struct foo {
3465              union {
3466                  int __i;
3467          #ifdef _FEATURE_TEST
3468                  int i;
3469          #endif
3470              } __ii;
3471          }
              ```

3472          does not work because the name of the logical field *i* is _ _*ii.i*, and introduction of a macro
3473          to restore the logical name immediately reintroduces the problem discussed previously
3474          (although its manifestation might be more immediate because a syntax error would result
3475          if a recursive macro did not cause it to fail first).

7.    A more workable solution would be to declare the structure:

```
struct foo {
#ifdef _FEATURE_TEST
    int i;
#else
    int __i;
#endif
}
```

However, if a macro (particularly one required by a standard) is to be defined that uses this field, two must be defined: one that uses *i*, the other that uses *__i*. If more than one additional field is used in a macro and they are conditional on distinct combinations of features, the complexity goes up as $2^n$.

All this leaves a difficult situation: vendors must provide very complex headers to deal with what is conceptually simple and safe—adding a field to a structure. It is the possibility of user-provided macros with the same name that makes this difficult.

Several alternatives were proposed that involved constraining the user's access to part of the name space available to the user (as specified by the ISO C standard). In some cases, this was only until all the headers had been included. There were two proposals discussed that failed to achieve consensus:

1.    Limiting it for the whole program.

2.    Restricting the use of identifiers containing only uppercase letters until after all system headers had been included. It was also pointed out that because macros might wish to access fields of a structure (and macro expansion occurs totally at point of use) restricting names in this way would not protect the macro expansion, and thus the solution was inadequate.

It was finally decided that reservation of symbols would occur, but as constrained.

The current wording also allows the addition of fields to a structure, but requires that user macros of the same name not interfere. This allows vendors to do one of the following:

- Not create the situation (do not extend the structures with user-accessible names or use the solution in (7) above)

- Extend their compilers to allow some way of adding names to structures and macros safely

There are at least two ways that the compiler might be extended: add new preprocessor directives that turn off and on macro expansion for certain symbols (without changing the value of the macro) and a function or lexical operation that suppresses expansion of a word. The latter seems more flexible, particularly because it addresses the problem in macros as well as in declarations.

The following seems to be a possible implementation extension to the C language that will do this: any token that during macro expansion is found to be preceded by three '#' symbols shall not be further expanded in exactly the same way as described for macros that expand to their own name as in Section 3.8.3.4 of the ISO C standard. A vendor may also wish to implement this as an operation that is lexically a function, which might be implemented as:

```
#define __safe_name(x) ###x
```

Using a function notation would insulate vendors from changes in standards until such a functionality is standardized (if ever). Standardization of such a function would be valuable because it would then permit third parties to take advantage of it portably in software they may

3521     supply.

3522     The symbols that are ''explicitly permitted, but not required by IEEE Std 1003.1-2001'' include
3523     those classified below. (That is, the symbols classified below might, but are not required to, be
3524     present when _POSIX_C_SOURCE is defined to have the value 200112L.)

3525       • Symbols in **<limits.h>** and **<unistd.h>** that are defined to indicate support for options or
3526         limits that are constant at compile-time

3527       • Symbols in the name space reserved for the implementation by the ISO C standard

3528       • Symbols in a name space reserved for a particular type of extension (for example, type names
3529         ending with **_t** in **<sys/types.h>**)

3530       • Additional members of structures or unions whose names do not reduce the name space
3531         reserved for applications

3532     Since both implementations and future revisions of IEEE Std 1003.1 and other POSIX standards
3533     may use symbols in the reserved spaces described in these tables, there is a potential for name
3534     space clashes. To avoid future name space clashes when adding symbols, implementations
3535     should not use the posix_, POSIX_, or _POSIX_ prefixes.

3536     IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/2 is applied, deleting the entries POSIX_,  |
3537     _POSIX_, and posix_ from the column of allowed name space prefixes for use by an  |
3538     implementation in the first table. The presence of these prefixes was contradicting later text  |
3539     which states that: ''The prefixes posix_, POSIX_, and _POSIX are reserved for use by Shell and  |
3540     Utilities volume of IEEE Std 1003.1-2001, Chapter 2, Shell Command Language and other POSIX  |
3541     standards. Implementations may add symbols to the headers shown in the following table,  |
3542     provided the identifiers . . . do not use the reserved prefixes posix_, POSIX_, or _POSIX.''.  |

3543     IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/3 is applied, correcting the reserved  |
3544     macro prefix from: ''PRI[a-z], SCN[a-z]'' to: ''PRI[Xa-z], SCN[Xa-z]'' in the second table. The  |
3545     change was needed since the ISO C standard allows implementations to define macros of the  |
3546     form PRI or SCN followed by any lowercase letter or ′X′ in **<inttypes.h>**. (The  |
3547     ISO/IEC 9899: 1999 standard, Subclause 7.26.4.)  |

3548     IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/4 is applied, adding a new section listing  |
3549     reserved names for the **<stdint.h>** header. This change is for alignment with the ISO C standard.  |

## 3550   B.2.3     Error Numbers

3551     It was the consensus of the standard developers that to allow the conformance document to
3552     state that an error occurs and under what conditions, but to disallow a statement that it never
3553     occurs, does not make sense. It could be implied by the current wording that this is allowed, but
3554     to reduce the possibility of future interpretation requests, it is better to make an explicit
3555     statement.

3556     The ISO C standard requires that *errno* be an assignable lvalue. Originally, the definition in
3557     POSIX.1 was stricter than that in the ISO C standard, **extern int** *errno*, in order to support
3558     historical usage. In a multi-threaded environment, implementing *errno* as a global variable
3559     results in non-deterministic results when accessed. It is required, however, that *errno* work as a
3560     per-thread error reporting mechanism. In order to do this, a separate *errno* value has to be
3561     maintained for each thread. The following section discusses the various alternative solutions
3562     that were considered.

3563     In order to avoid this problem altogether for new functions, these functions avoid using *errno*
3564     and, instead, return the error number directly as the function return value; a return value of zero
3565     indicates that no error was detected.

3566 For any function that can return errors, the function return value is not used for any purpose
3567 other than for reporting errors. Even when the output of the function is scalar, it is passed
3568 through a function argument. While it might have been possible to allow some scalar outputs to
3569 be coded as negative function return values and mixed in with positive error status returns, this
3570 was rejected—using the return value for a mixed purpose was judged to be of limited use and
3571 error prone.

3572 Checking the value of *errno* alone is not sufficient to determine the existence or type of an error,
3573 since it is not required that a successful function call clear *errno*. The variable *errno* should only
3574 be examined when the return value of a function indicates that the value of *errno* is meaningful.
3575 In that case, the function is required to set the variable to something other than zero.

3576 The variable *errno* is never set to zero by any function call; to do so would contradict the ISO C
3577 standard.

3578 POSIX.1 requires (in the ERRORS sections of function descriptions) certain error values to be set
3579 in certain conditions because many existing applications depend on them. Some error numbers,
3580 such as [EFAULT], are entirely implementation-defined and are noted as such in their
3581 description in the ERRORS section. This section otherwise allows wide latitude to the
3582 implementation in handling error reporting.

3583 Some of the ERRORS sections in IEEE Std 1003.1-2001 have two subsections. The first:

3584     ''The function shall fail if:''

3585 could be called the ''mandatory'' section.

3586 The second:

3587     ''The function may fail if:''

3588 could be informally known as the ''optional'' section.

3589 Attempting to infer the quality of an implementation based on whether it detects optional error
3590 conditions is not useful.

3591 Following each one-word symbolic name for an error, there is a description of the error. The
3592 rationale for some of the symbolic names follows:

3593 [ECANCELED]   This spelling was chosen as being more common.

3594 [EFAULT]       Most historical implementations do not catch an error and set *errno* when an
3595                  invalid address is given to the functions *wait*( ), *time*( ), or *times*( ). Some
3596                  implementations cannot reliably detect an invalid address. And most systems
3597                  that detect invalid addresses will do so only for a system call, not for a library
3598                  routine.

3599 [EFTYPE]       This error code was proposed in earlier proposals as ''Inappropriate operation
3600                  for file type'', meaning that the operation requested is not appropriate for the
3601                  file specified in the function call. This code was proposed, although the same
3602                  idea was covered by [ENOTTY], because the connotations of the name would
3603                  be misleading. It was pointed out that the *fcntl*( ) function uses the error code
3604                  [EINVAL] for this notion, and hence all instances of [EFTYPE] were changed
3605                  to this code.

3606 [EINTR]        POSIX.1 prohibits conforming implementations from restarting interrupted
3607                  system calls of conforming applications unless the SA_RESTART flag is in
3608                  effect for the signal. However, it does not require that [EINTR] be returned
3609                  when another legitimate value may be substituted; for example, a partial
3610                  transfer count when *read*( ) or *write*( ) are interrupted. This is only given when

       

| 3611 | | the signal-catching function returns normally as opposed to returns by |
| 3612 | | mechanisms like *longjmp*( ) or *siglongjmp*( ). |
| 3613 | [ELOOP] | In specifying conditions under which implementations would generate this |
| 3614 | | error, the following goals were considered: |

- 3615 • To ensure that actual loops are detected, including loops that result from
- 3616 symbolic links across distributed file systems.

- 3617 • To ensure that during pathname resolution an application can rely on the
- 3618 ability to follow at least {SYMLOOP_MAX} symbolic links in the absence
- 3619 of a loop.

- 3620 • To allow implementations to provide the capability of traversing more
- 3621 than {SYMLOOP_MAX} symbolic links in the absence of a loop.

- 3622 • To allow implementations to detect loops and generate the error prior to
- 3623 encountering {SYMLOOP_MAX} symbolic links.

3624 [ENAMETOOLONG]
3625 When a symbolic link is encountered during pathname resolution, the
3626 contents of that symbolic link are used to create a new pathname. The
3627 standard developers intended to allow, but not require, that implementations
3628 enforce the restriction of {PATH_MAX} on the result of this pathname
3629 substitution.

3630 [ENOMEM]  The term ''main memory'' is not used in POSIX.1 because it is
3631 implementation-defined.

3632 [ENOTSUP]  This error code is to be used when an implementation chooses to implement
3633 the required functionality of IEEE Std 1003.1-2001 but does not support
3634 optional facilities defined by IEEE Std 1003.1-2001. The return of [ENOSYS] is
3635 to be taken to indicate that the function of the interface is not supported at all;
3636 the function will always fail with this error code.

3637 [ENOTTY]  The symbolic name for this error is derived from a time when device control
3638 was done by *ioctl*( ) and that operation was only permitted on a terminal
3639 interface. The term ''TTY'' is derived from ''teletypewriter'', the devices to
3640 which this error originally applied.

3641 [EOVERFLOW] Most of the uses of this error code are related to large file support. Typically,
3642 these cases occur on systems which support multiple programming
3643 environments with different sizes for **off_t**, but they may also occur in
3644 connection with remote file systems.

3645 In addition, when different programming environments have different widths
3646 for types such as **int** and **uid_t**, several functions may encounter a condition
3647 where a value in a particular environment is too wide to be represented. In
3648 that case, this error should be raised. For example, suppose the currently
3649 running process has 64-bit **int**, and file descriptor 9 223 372 036 854 775 807 is
3650 open and does not have the close-on-*exec* flag set. If the process then uses
3651 *execl*( ) to *exec* a file compiled in a programming environment with 32-bit **int**,
3652 the call to *execl*( ) can fail with *errno* set to [EOVERFLOW]. A similar failure
3653 can occur with *execl*( ) if any of the user IDs or any of the group IDs to be
3654 assigned to the new process image are out of range for the executed file's
3655 programming environment.

| | | |
|---|---|---|
| 3656 | | Note, however, that this condition cannot occur for functions that are |
| 3657 | | explicitly described as always being successful, such as *getpid*( ). |
| 3658 | [EPIPE] | This condition normally generates the signal SIGPIPE; the error is returned if |
| 3659 | | the signal does not terminate the process. |
| 3660 | [EROFS] | In historical implementations, attempting to *unlink*( ) or *rmdir*( ) a mount point |
| 3661 | | would generate an [EBUSY] error. An implementation could be envisioned |
| 3662 | | where such an operation could be performed without error. In this case, if |
| 3663 | | *either* the directory entry or the actual data structures reside on a read-only file |
| 3664 | | system, [EROFS] is the appropriate error to generate. (For example, changing |
| 3665 | | the link count of a file on a read-only file system could not be done, as is |
| 3666 | | required by *unlink*( ), and thus an error should be reported.) |

3667 Three error numbers, [EDOM], [EILSEQ], and [ERANGE], were added to this section primarily
3668 for consistency with the ISO C standard.

**Alternative Solutions for Per-Thread errno**

3670 The usual implementation of *errno* as a single global variable does not work in a multi-threaded
3671 environment. In such an environment, a thread may make a POSIX.1 call and get a −1 error
3672 return, but before that thread can check the value of *errno*, another thread might have made a
3673 second POSIX.1 call that also set *errno*. This behavior is unacceptable in robust programs. There
3674 were a number of alternatives that were considered for handling the *errno* problem:

3675 • Implement *errno* as a per-thread integer variable.

3676 • Implement *errno* as a service that can access the per-thread error number.

3677 • Change all POSIX.1 calls to accept an extra status argument and avoid setting *errno*.

3678 • Change all POSIX.1 calls to raise a language exception.

3679 The first option offers the highest level of compatibility with existing practice but requires
3680 special support in the linker, compiler, and/or virtual memory system to support the new
3681 concept of thread private variables. When compared with current practice, the third and fourth
3682 options are much cleaner, more efficient, and encourage a more robust programming style, but
3683 they require new versions of all of the POSIX.1 functions that might detect an error. The second
3684 option offers compatibility with existing code that uses the **<errno.h>** header to define the
3685 symbol *errno*. In this option, *errno* may be a macro defined:

```
3686    #define errno    (*__errno())
3687    extern int       *__errno();
```

3688 This option may be implemented as a per-thread variable whereby an *errno* field is allocated in
3689 the user space object representing a thread, and whereby the function *__errno*( ) makes a system
3690 call to determine the location of its user space object and returns the address of the *errno* field of
3691 that object. Another implementation, one that avoids calling the kernel, involves allocating
3692 stacks in chunks. The stack allocator keeps a side table indexed by chunk number containing a
3693 pointer to the thread object that uses that chunk. The *__errno*( ) function then looks at the stack
3694 pointer, determines the chunk number, and uses that as an index into the chunk table to find its
3695 thread object and thus its private value of *errno*. On most architectures, this can be done in four
3696 to five instructions. Some compilers may wish to implement *__errno*( ) inline to improve
3697 performance.

3698 **Disallowing Return of the [EINTR] Error Code**

3699 Many blocking interfaces defined by IEEE Std 1003.1-2001 may return [EINTR] if interrupted
3700 during their execution by a signal handler. Blocking interfaces introduced under the Threads
3701 option do not have this property. Instead, they require that the interface appear to be atomic
3702 with respect to interruption. In particular, clients of blocking interfaces need not handle any
3703 possible [EINTR] return as a special case since it will never occur. If it is necessary to restart
3704 operations or complete incomplete operations following the execution of a signal handler, this is
3705 handled by the implementation, rather than by the application.

3706 Requiring applications to handle [EINTR] errors on blocking interfaces has been shown to be a
3707 frequent source of often unreproducible bugs, and it adds no compelling value to the available
3708 functionality. Thus, blocking interfaces introduced for use by multi-threaded programs do not
3709 use this paradigm. In particular, in none of the functions *flockfile*(), *pthread_cond_timedwait*(),
3710 *pthread_cond_wait*(), *pthread_join*(), *pthread_mutex_lock*(), and *sigwait*() did providing [EINTR]
3711 returns add value, or even particularly make sense. Thus, these functions do not provide for an
3712 [EINTR] return, even when interrupted by a signal handler. The same arguments can be applied
3713 to *sem_wait*(), *sem_trywait*(), *sigwaitinfo*(), and *sigtimedwait*(), but implementations are
3714 permitted to return [EINTR] error codes for these functions for compatibility with earlier
3715 versions of IEEE Std 1003.1. Applications cannot rely on calls to these functions returning
3716 [EINTR] error codes when signals are delivered to the calling thread, but they should allow for
3717 the possibility.

3718 *B.2.3.1    Additional Error Numbers*

3719 The ISO C standard defines the name space for implementations to add additional error
3720 numbers.

3721 **B.2.4    Signal Concepts**

3722 Historical implementations of signals, using the *signal*() function, have shortcomings that make
3723 them unreliable for many application uses. Because of this, a new signal mechanism, based very
3724 closely on the one of 4.2 BSD and 4.3 BSD, was added to POSIX.1.

3725 **Signal Names**

3726 The restriction on the actual type used for **sigset_t** is intended to guarantee that these objects can
3727 always be assigned, have their address taken, and be passed as parameters by value. It is not
3728 intended that this type be a structure including pointers to other data structures, as that could
3729 impact the portability of applications performing such operations. A reasonable implementation
3730 could be a structure containing an array of some integer type.

3731 The signals described in IEEE Std 1003.1-2001 must have unique values so that they may be
3732 named as parameters of **case** statements in the body of a C-language **switch** clause. However,
3733 implementation-defined signals may have values that overlap with each other or with signals
3734 specified in IEEE Std 1003.1-2001. An example of this is SIGABRT, which traditionally overlaps
3735 some other signal, such as SIGIOT.

3736 SIGKILL, SIGTERM, SIGUSR1, and SIGUSR2 are ordinarily generated only through the explicit
3737 use of the *kill*() function, although some implementations generate SIGKILL under
3738 extraordinary circumstances. SIGTERM is traditionally the default signal sent by the *kill*
3739 command.

3740 The signals SIGBUS, SIGEMT, SIGIOT, SIGTRAP, and SIGSYS were omitted from POSIX.1
3741 because their behavior is implementation-defined and could not be adequately categorized.
3742 Conforming implementations may deliver these signals, but must document the circumstances

3743   under which they are delivered and note any restrictions concerning their delivery. The signals
3744   SIGFPE, SIGILL, and SIGSEGV are similar in that they also generally result only from
3745   programming errors. They were included in POSIX.1 because they do indicate three relatively
3746   well-categorized conditions. They are all defined by the ISO C standard and thus would have to
3747   be defined by any system with an ISO C standard binding, even if not explicitly included in
3748   POSIX.1.

3749   There is very little that a Conforming POSIX.1 Application can do by catching, ignoring, or
3750   masking any of the signals SIGILL, SIGTRAP, SIGIOT, SIGEMT, SIGBUS, SIGSEGV, SIGSYS, or
3751   SIGFPE. They will generally be generated by the system only in cases of programming errors.
3752   While it may be desirable for some robust code (for example, a library routine) to be able to
3753   detect and recover from programming errors in other code, these signals are not nearly sufficient
3754   for that purpose. One portable use that does exist for these signals is that a command interpreter
3755   can recognize them as the cause of a process' termination (with *wait*()) and print an appropriate
3756   message. The mnemonic tags for these signals are derived from their PDP-11 origin.

3757   The signals SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU, and SIGCONT are provided for job control
3758   and are unchanged from 4.2 BSD. The signal SIGCHLD is also typically used by job control
3759   shells to detect children that have terminated or, as in 4.2 BSD, stopped.

3760   Some implementations, including System V, have a signal named SIGCLD, which is similar to
3761   SIGCHLD in 4.2 BSD. POSIX.1 permits implementations to have a single signal with both
3762   names. POSIX.1 carefully specifies ways in which conforming applications can avoid the
3763   semantic differences between the two different implementations. The name SIGCHLD was
3764   chosen for POSIX.1 because most current application usages of it can remain unchanged in
3765   conforming applications. SIGCLD in System V has more cases of semantics that POSIX.1 does
3766   not specify, and thus applications using it are more likely to require changes in addition to the
3767   name change.

3768   The signals SIGUSR1 and SIGUSR2 are commonly used by applications for notification of
3769   exceptional behavior and are described as ''reserved as application-defined'' so that such use is
3770   not prohibited. Implementations should not generate SIGUSR1 or SIGUSR2, except when
3771   explicitly requested by *kill*(). It is recommended that libraries not use these two signals, as such
3772   use in libraries could interfere with their use by applications calling the libraries. If such use is
3773   unavoidable, it should be documented. It is prudent for non-portable libraries to use non-
3774   standard signals to avoid conflicts with use of standard signals by portable libraries.

3775   There is no portable way for an application to catch or ignore non-standard signals. Some
3776   implementations define the range of signal numbers, so applications can install signal-catching
3777   functions for all of them. Unfortunately, implementation-defined signals often cause problems
3778   when caught or ignored by applications that do not understand the reason for the signal. While
3779   the desire exists for an application to be more robust by handling all possible signals (even those
3780   only generated by *kill*()), no existing mechanism was found to be sufficiently portable to include
3781   in POSIX.1. The value of such a mechanism, if included, would be diminished given that
3782   SIGKILL would still not be catchable.

3783   A number of new signal numbers are reserved for applications because the two user signals
3784   defined by POSIX.1 are insufficient for many realtime applications. A range of signal numbers is
3785   specified, rather than an enumeration of additional reserved signal names, because different
3786   applications and application profiles will require a different number of application signals. It is
3787   not desirable to burden all application domains and therefore all implementations with the
3788   maximum number of signals required by all possible applications. Note that in this context,
3789   signal numbers are essentially different signal priorities.

3790   The relatively small number of required additional signals, {_POSIX_RTSIG_MAX}, was chosen
3791   so as not to require an unreasonably large signal mask/set. While this number of signals defined

3792      in POSIX.1 will fit in a single 32-bit word signal mask, it is recognized that most existing
3793      implementations define many more signals than are specified in POSIX.1 and, in fact, many
3794      implementations have already exceeded 32 signals (including the ''null signal''). Support of
3795      {_POSIX_RTSIG_MAX} additional signals may push some implementation over the single 32-bit
3796      word line, but is unlikely to push any implementations that are already over that line beyond the
3797      64-signal line.

3798 *B.2.4.1*     *Signal Generation and Delivery*

3799      The terms defined in this section are not used consistently in documentation of historical
3800      systems. Each signal can be considered to have a lifetime beginning with generation and ending
3801      with delivery or acceptance. The POSIX.1 definition of ''delivery'' does not exclude ignored
3802      signals; this is considered a more consistent definition. This revised text in several parts of
3803      IEEE Std 1003.1-2001 clarifies the distinct semantics of asynchronous signal delivery and
3804      synchronous signal acceptance. The previous wording attempted to categorize both under the
3805      term ''delivery'', which led to conflicts over whether the effects of asynchronous signal delivery
3806      applied to synchronous signal acceptance.

3807      Signals generated for a process are delivered to only one thread. Thus, if more than one thread is
3808      eligible to receive a signal, one has to be chosen. The choice of threads is left entirely up to the
3809      implementation both to allow the widest possible range of conforming implementations and to
3810      give implementations the freedom to deliver the signal to the ''easiest possible'' thread should
3811      there be differences in ease of delivery between different threads.

3812      Note that should multiple delivery among cooperating threads be required by an application,
3813      this can be trivially constructed out of the provided single-delivery semantics. The construction
3814      of a *sigwait_multiple*() function that accomplishes this goal is presented with the rationale for
3815      *sigwaitinfo*().

3816      Implementations should deliver unblocked signals as soon after they are generated as possible.
3817      However, it is difficult for POSIX.1 to make specific requirements about this, beyond those in
3818      *kill*() and *sigprocmask*(). Even on systems with prompt delivery, scheduling of higher priority
3819      processes is always likely to cause delays.

3820      In general, the interval between the generation and delivery of unblocked signals cannot be
3821      detected by an application. Thus, references to pending signals generally apply to blocked,
3822      pending signals. An implementation registers a signal as pending on the process when no thread
3823      has the signal unblocked and there are no threads blocked in a *sigwait*() function for that signal.
3824      Thereafter, the implementation delivers the signal to the first thread that unblocks the signal or
3825      calls a *sigwait*() function on a signal set containing this signal rather than choosing the recipient
3826      thread at the time the signal is sent.

3827      In the 4.3 BSD system, signals that are blocked and set to SIG_IGN are discarded immediately
3828      upon generation. For a signal that is ignored as its default action, if the action is SIG_DFL and
3829      the signal is blocked, a generated signal remains pending. In the 4.1 BSD system and in System V
3830      Release 3 (two other implementations that support a somewhat similar signal mechanism), all
3831      ignored blocked signals remain pending if generated. Because it is not normally useful for an
3832      application to simultaneously ignore and block the same signal, it was unnecessary for POSIX.1
3833      to specify behavior that would invalidate any of the historical implementations.

3834      There is one case in some historical implementations where an unblocked, pending signal does
3835      not remain pending until it is delivered. In the System V implementation of *signal*(), pending
3836      signals are discarded when the action is set to SIG_DFL or a signal-catching routine (as well as to
3837      SIG_IGN). Except in the case of setting SIGCHLD to SIG_DFL, implementations that do this do
3838      not conform completely to POSIX.1. Some earlier proposals for POSIX.1 explicitly stated this,
3839      but these statements were redundant due to the requirement that functions defined by POSIX.1

3840    not change attributes of processes defined by POSIX.1 except as explicitly stated.

3841    POSIX.1 specifically states that the order in which multiple, simultaneously pending signals are
3842    delivered is unspecified. This order has not been explicitly specified in historical
3843    implementations, but has remained quite consistent and been known to those familiar with the
3844    implementations. Thus, there have been cases where applications (usually system utilities) have
3845    been written with explicit or implicit dependencies on this order. Implementors and others
3846    porting existing applications may need to be aware of such dependencies.

3847    When there are multiple pending signals that are not blocked, implementations should arrange
3848    for the delivery of all signals at once, if possible. Some implementations stack calls to all pending
3849    signal-catching routines, making it appear that each signal-catcher was interrupted by the next
3850    signal. In this case, the implementation should ensure that this stacking of signals does not
3851    violate the semantics of the signal masks established by *sigaction*( ). Other implementations
3852    process at most one signal when the operating system is entered, with remaining signals saved
3853    for later delivery. Although this practice is widespread, this behavior is neither standardized
3854    nor endorsed. In either case, implementations should attempt to deliver signals associated with
3855    the current state of the process (for example, SIGFPE) before other signals, if possible.

3856    In 4.2 BSD and 4.3 BSD, it is not permissible to ignore or explicitly block SIGCONT, because if
3857    blocking or ignoring this signal prevented it from continuing a stopped process, such a process
3858    could never be continued (only killed by SIGKILL). However, 4.2 BSD and 4.3 BSD do block
3859    SIGCONT during execution of its signal-catching function when it is caught, creating exactly
3860    this problem. A proposal was considered to disallow catching SIGCONT in addition to ignoring
3861    and blocking it, but this limitation led to objections. The consensus was to require that
3862    SIGCONT always continue a stopped process when generated. This removed the need to
3863    disallow ignoring or explicit blocking of the signal; note that SIG_IGN and SIG_DFL are
3864    equivalent for SIGCONT.

3865    *B.2.4.2    Realtime Signal Generation and Delivery*

3866    The Realtime Signals Extension option to POSIX.1 signal generation and delivery behavior is
3867    required for the following reasons:

3868    • The **sigevent** structure is used by other POSIX.1 functions that result in asynchronous event
3869      notifications to specify the notification mechanism to use and other information needed by
3870      the notification mechanism. IEEE Std 1003.1-2001 defines only three symbolic values for the
3871      notification mechanism:                                                                              |

3872      — SIGEV_NONE is used to indicate that no notification is required when the event occurs.    |
3873        This is useful for applications that use asynchronous I/O with polling for completion.     |

3874      — SIGEV_SIGNAL indicates that a signal is generated when the event occurs.                    |

3875      — SIGEV_THREAD provides for ''callback functions'' for asynchronous notifications done   |
3876        by a function call within the context of a new thread. This provides a multi-threaded   |
3877        process with a more natural means of notification than signals.                          |

3878    The primary difficulty with previous notification approaches has been to specify the    |
3879    environment of the notification routine.

3880      — One approach is to limit the notification routine to call only functions permitted in a
3881        signal handler. While the list of permissible functions is clearly stated, this is overly
3882        restrictive.

3883      — A second approach is to define a new list of functions or classes of functions that are
3884        explicitly permitted or not permitted. This would give a programmer more lists to deal
3885        with, which would be awkward.

— The third approach is to define completely the environment for execution of the notification function. A clear definition of an execution environment for notification is provided by executing the notification function in the environment of a newly created thread.

Implementations may support additional notification mechanisms by defining new values for *sigev_notify*.

For a notification type of SIGEV_SIGNAL, the other members of the **sigevent** structure defined by IEEE Std 1003.1-2001 specify the realtime signal—that is, the signal number and application-defined value that differentiates between occurrences of signals with the same number—that will be generated when the event occurs. The structure is defined in **<signal.h>**, even though the structure is not directly used by any of the signal functions, because it is part of the signals interface used by the POSIX.1b ''client functions''. When the client functions include **<signal.h>** to define the signal names, the **sigevent** structure will also be defined.

An application-defined value passed to the signal handler is used to differentiate between different ''events'' instead of requiring that the application use different signal numbers for several reasons:

— Realtime applications potentially handle a very large number of different events. Requiring that implementations support a correspondingly large number of distinct signal numbers will adversely impact the performance of signal delivery because the signal masks to be manipulated on entry and exit to the handlers will become large.

— Event notifications are prioritized by signal number (the rationale for this is explained in the following paragraphs) and the use of different signal numbers to differentiate between the different event notifications overloads the signal number more than has already been done. It also requires that the application writer make arbitrary assignments of priority to events that are logically of equal priority.

A union is defined for the application-defined value so that either an integer constant or a pointer can be portably passed to the signal-catching function. On some architectures a pointer cannot be cast to an **int** and *vice versa*.

Use of a structure here with an explicit notification type discriminant rather than explicit parameters to realtime functions, or embedded in other realtime structures, provides for future extensions to IEEE Std 1003.1-2001. Additional, perhaps more efficient, notification mechanisms can be supported for existing realtime function interfaces, such as timers and asynchronous I/O, by extending the **sigevent** structure appropriately. The existing realtime function interfaces will not have to be modified to use any such new notification mechanism. The revised text concerning the SIGEV_SIGNAL value makes consistent the semantics of the members of the **sigevent** structure, particularly in the definitions of *lio_listio*( ) and *aio_fsync*( ). For uniformity, other revisions cause this specification to be referred to rather than inaccurately duplicated in the descriptions of functions and structures using the **sigevent** structure. The revised wording does not relax the requirement that the signal number be in the range SIGRTMIN to SIGRTMAX to guarantee queuing and passing of the application value, since that requirement is still implied by the signal names.

• IEEE Std 1003.1-2001 is intentionally vague on whether ''non-realtime'' signal-generating mechanisms can result in a **siginfo_t** being supplied to the handler on delivery. In one existing implementation, a **siginfo_t** is posted on signal generation, even though the implementation does not support queuing of multiple occurrences of a signal. It is not the intent of IEEE Std 1003.1-2001 to preclude this, independent of the mandate to define signals that do support queuing. Any interpretation that appears to preclude this is a mistake in the

3934     reading or writing of the standard.

3935 • Signals handled by realtime signal handlers might be generated by functions or conditions
3936     that do not allow the specification of an application-defined value and do not queue.
3937     IEEE Std 1003.1-2001 specifies the *si_code* member of the **siginfo_t** structure used in existing
3938     practice and defines additional codes so that applications can detect whether an application-
3939     defined value is present or not. The code SI_USER for *kill*( )-generated signals is adopted
3940     from existing practice.

3941 • The *sigaction*( ) *sa_flags* value SA_SIGINFO tells the implementation that the signal-catching
3942     function expects two additional arguments. When the flag is not set, a single argument, the
3943     signal number, is passed as specified by IEEE Std 1003.1-2001. Although IEEE Std 1003.1-2001
3944     does not explicitly allow the *info* argument to the handler function to be NULL, this is
3945     existing practice. This provides for compatibility with programs whose signal-catching
3946     functions are not prepared to accept the additional arguments. IEEE Std 1003.1-2001 is
3947     explicitly unspecified as to whether signals actually queue when SA_SIGINFO is not set for a
3948     signal, as there appear to be no benefits to applications in specifying one behavior or another.
3949     One existing implementation queues a **siginfo_t** on each signal generation, unless the signal
3950     is already pending, in which case the implementation discards the new **siginfo_t**; that is, the
3951     queue length is never greater than one. This implementation only examines SA_SIGINFO on
3952     signal delivery, discarding the queued **siginfo_t** if its delivery was not requested.

3953     IEEE Std 1003.1-2001 specifies several new values for the *si_code* member of the **siginfo_t**
3954     structure. In existing practice, a *si_code* value of less than or equal to zero indicates that the
3955     signal was generated by a process via the *kill*( ) function. In existing practice, values of *si_code*
3956     that provide additional information for implementation-generated signals, such as SIGFPE or
3957     SIGSEGV, are all positive. Thus, if implementations define the new constants specified in
3958     IEEE Std 1003.1-2001 to be negative numbers, programs written to use existing practice will
3959     not break. IEEE Std 1003.1-2001 chose not to attempt to specify existing practice values of
3960     *si_code* other than SI_USER both because it was deemed beyond the scope of
3961     IEEE Std 1003.1-2001 and because many of the values in existing practice appear to be
3962     platform and implementation-defined. But, IEEE Std 1003.1-2001 does specify that if an
3963     implementation—for example, one that does not have existing practice in this area—chooses
3964     to define additional values for *si_code*, these values have to be different from the values of the
3965     symbols specified by IEEE Std 1003.1-2001. This will allow conforming applications to
3966     differentiate between signals generated by one of the POSIX.1b asynchronous events and
3967     those generated by other implementation events in a manner compatible with existing
3968     practice.

3969     The unique values of *si_code* for the POSIX.1b asynchronous events have implications for
3970     implementations of, for example, asynchronous I/O or message passing in user space library
3971     code. Such an implementation will be required to provide a hidden interface to the signal
3972     generation mechanism that allows the library to specify the standard values of *si_code*.

3973     Existing practice also defines additional members of **siginfo_t**, such as the process ID and
3974     user ID of the sending process for *kill*( )-generated signals. These members were deemed not
3975     necessary to meet the requirements of realtime applications and are not specified by
3976     IEEE Std 1003.1-2001. Neither are they precluded.

3977     The third argument to the signal-catching function, *context*, is left undefined by
3978     IEEE Std 1003.1-2001, but is specified in the interface because it matches existing practice for
3979     the SA_SIGINFO flag. It was considered undesirable to require a separate implementation
3980     for SA_SIGINFO for POSIX conformance on implementations that already support the two
3981     additional parameters.

3982      • The requirement to deliver lower numbered signals in the range SIGRTMIN to SIGRTMAX
3983        first, when multiple unblocked signals are pending, results from several considerations:

3984        — A method is required to prioritize event notifications. The signal number was chosen
3985           instead of, for instance, associating a separate priority with each request, because an
3986           implementation has to check pending signals at various points and select one for delivery
3987           when more than one is pending. Specifying a selection order is the minimal additional
3988           semantic that will achieve prioritized delivery. If a separate priority were to be associated
3989           with queued signals, it would be necessary for an implementation to search all non-
3990           empty, non-blocked signal queues and select from among them the pending signal with
3991           the highest priority. This would significantly increase the cost of and decrease the
3992           determinism of signal delivery.

3993        — Given the specified selection of the lowest numeric unblocked pending signal,
3994           preemptive priority signal delivery can be achieved using signal numbers and signal
3995           masks by ensuring that the *sa_mask* for each signal number blocks all signals with a
3996           higher numeric value.

3997           For realtime applications that want to use only the newly defined realtime signal numbers
3998           without interference from the standard signals, this can be achieved by blocking all of the
3999           standard signals in the thread signal mask and in the *sa_mask* installed by the signal   |
4000           action for the realtime signal handlers.

4001        IEEE Std 1003.1-2001 explicitly leaves unspecified the ordering of signals outside of the range
4002        of realtime signals and the ordering of signals within this range with respect to those outside
4003        the range. It was believed that this would unduly constrain implementations or standards in
4004        the future definition of new signals.

4005  *B.2.4.3*   *Signal Actions*

4006        Early proposals mentioned SIGCONT as a second exception to the rule that signals are not
4007        delivered to stopped processes until continued. Because IEEE Std 1003.1-2001 now specifies that
4008        SIGCONT causes the stopped process to continue when it is generated, delivery of SIGCONT is
4009        not prevented because a process is stopped, even without an explicit exception to this rule.

4010        Ignoring a signal by setting the action to SIG_IGN (or SIG_DFL for signals whose default action
4011        is to ignore) is not the same as installing a signal-catching function that simply returns. Invoking
4012        such a function will interrupt certain system functions that block processes (for example, *wait*(),
4013        *sigsuspend*(), *pause*(), *read*(), *write*()) while ignoring a signal has no such effect on the process.

4014        Historical implementations discard pending signals when the action is set to SIG_IGN.
4015        However, they do not always do the same when the action is set to SIG_DFL and the default
4016        action is to ignore the signal. IEEE Std 1003.1-2001 requires this for the sake of consistency and
4017        also for completeness, since the only signal this applies to is SIGCHLD, and IEEE Std 1003.1-2001
4018        disallows setting its action to SIG_IGN.

4019        Some implementations (System V, for example) assign different semantics for SIGCLD
4020        depending on whether the action is set to SIG_IGN or SIG_DFL. Since POSIX.1 requires that the
4021        default action for SIGCHLD be to ignore the signal, applications should always set the action to
4022        SIG_DFL in order to avoid SIGCHLD.

4023        Whether or not an implementation allows SIG_IGN as a SIGCHLD disposition to be inherited
4024        across a call to one of the *exec* family of functions or *posix_spawn*() is explicitly left as
4025        unspecified. This change was made as a result of IEEE PASC Interpretation 1003.1 #132, and
4026        permits the implementation to decide between the following alternatives:

- Unconditionally leave SIGCHLD set to SIG_IGN, in which case the implementation would not allow applications that assume inheritance of SIG_DFL to conform to IEEE Std 1003.1-2001 without change. The implementation would, however, retain an ability to control applications that create child processes but never call on the *wait* family of functions, potentially filling up the process table.

- Unconditionally reset SIGCHLD to SIG_DFL, in which case the implementation would allow applications that assume inheritance of SIG_DFL to conform. The implementation would, however, lose an ability to control applications that spawn child processes but never reap them.

- Provide some mechanism, not specified in IEEE Std 1003.1-2001, to control inherited SIGCHLD dispositions.

Some implementations (System V, for example) will deliver a SIGCLD signal immediately when a process establishes a signal-catching function for SIGCLD when that process has a child that has already terminated. Other implementations, such as 4.3 BSD, do not generate a new SIGCHLD signal in this way. In general, a process should not attempt to alter the signal action for the SIGCHLD signal while it has any outstanding children. However, it is not always possible for a process to avoid this; for example, shells sometimes start up processes in pipelines with other processes from the pipeline as children. Processes that cannot ensure that they have no children when altering the signal action for SIGCHLD thus need to be prepared for, but not depend on, generation of an immediate SIGCHLD signal.

The default action of the stop signals (SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU) is to stop a process that is executing. If a stop signal is delivered to a process that is already stopped, it has no effect. In fact, if a stop signal is generated for a stopped process whose signal mask blocks the signal, the signal will never be delivered to the process since the process must receive a SIGCONT, which discards all pending stop signals, in order to continue executing.

The SIGCONT signal continues a stopped process even if SIGCONT is blocked (or ignored). However, if a signal-catching routine has been established for SIGCONT, it will not be entered until SIGCONT is unblocked.

If a process in an orphaned process group stops, it is no longer under the control of a job control shell and hence would not normally ever be continued. Because of this, orphaned processes that receive terminal-related stop signals (SIGTSTP, SIGTTIN, SIGTTOU, but not SIGSTOP) must not be allowed to stop. The goal is to prevent stopped processes from languishing forever. (As SIGSTOP is sent only via *kill*(), it is assumed that the process or user sending a SIGSTOP can send a SIGCONT when desired.) Instead, the system must discard the stop signal. As an extension, it may also deliver another signal in its place. 4.3 BSD sends a SIGKILL, which is overly effective because SIGKILL is not catchable. Another possible choice is SIGHUP. 4.3 BSD also does this for orphaned processes (processes whose parent has terminated) rather than for members of orphaned process groups; this is less desirable because job control shells manage process groups. POSIX.1 also prevents SIGTTIN and SIGTTOU signals from being generated for processes in orphaned process groups as a direct result of activity on a terminal, preventing infinite loops when *read*() and *write*() calls generate signals that are discarded; see Section A.11.1.4 (on page 67). A similar restriction on the generation of SIGTSTP was considered, but that would be unnecessary and more difficult to implement due to its asynchronous nature.

Although POSIX.1 requires that signal-catching functions be called with only one argument, there is nothing to prevent conforming implementations from extending POSIX.1 to pass additional arguments, as long as Strictly Conforming POSIX.1 Applications continue to compile and execute correctly. Most historical implementations do, in fact, pass additional, signal-specific arguments to certain signal-catching routines.

4075    There was a proposal to change the declared type of the signal handler to:

4076        `void func (int sig, ...);`

4077    The usage of ellipses ("`...`") is ISO C standard syntax to indicate a variable number of
4078    arguments. Its use was intended to allow the implementation to pass additional information to
4079    the signal handler in a standard manner.

4080    Unfortunately, this construct would require all signal handlers to be defined with this syntax
4081    because the ISO C standard allows implementations to use a different parameter passing
4082    mechanism for variable parameter lists than for non-variable parameter lists. Thus, all existing
4083    signal handlers in all existing applications would have to be changed to use the variable syntax
4084    in order to be standard and portable. This is in conflict with the goal of Minimal Changes to
4085    Existing Application Code.

4086    When terminating a process from a signal-catching function, processes should be aware of any
4087    interpretation that their parent may make of the status returned by *wait*() or *waitpid*(). In
4088    particular, a signal-catching function should not call *exit*(0) or *_exit*(0) unless it wants to indicate
4089    successful termination. A non-zero argument to *exit*() or *_exit*() can be used to indicate
4090    unsuccessful termination. Alternatively, the process can use *kill*() to send itself a fatal signal
4091    (first ensuring that the signal is set to the default action and not blocked). See also the
4092    RATIONALE section of the *_exit*() function.

4093    The behavior of *unsafe* functions, as defined by this section, is undefined when they are invoked
4094    from signal-catching functions in certain circumstances. The behavior of reentrant functions, as
4095    defined by this section, is as specified by POSIX.1, regardless of invocation from a signal-
4096    catching function. This is the only intended meaning of the statement that reentrant functions
4097    may be used in signal-catching functions without restriction. Applications must still consider all
4098    effects of such functions on such things as data structures, files, and process state.  In particular,
4099    application writers need to consider the restrictions on interactions when interrupting *sleep*()
4100    (see *sleep*()) and interactions among multiple handles for a file description. The fact that any
4101    specific function is listed as reentrant does not necessarily mean that invocation of that function
4102    from a signal-catching function is recommended.

4103    In order to prevent errors arising from interrupting non-reentrant function calls, applications
4104    should protect calls to these functions either by blocking the appropriate signals or through the
4105    use of some programmatic semaphore. POSIX.1 does not address the more general problem of
4106    synchronizing access to shared data structures. Note in particular that even the ''safe'' functions
4107    may modify the global variable *errno*; the signal-catching function may want to save and restore
4108    its value.  The same principles apply to the reentrancy of application routines and asynchronous
4109    data access.

4110    Note that *longjmp*() and *siglongjmp*() are not in the list of reentrant functions. This is because the
4111    code executing after *longjmp*() or *siglongjmp*() can call any unsafe functions with the same
4112    danger as calling those unsafe functions directly from the signal handler. Applications that use
4113    *longjmp*() or *siglongjmp*() out of signal handlers require rigorous protection in order to be
4114    portable. Many of the other functions that are excluded from the list are traditionally
4115    implemented using either the C language *malloc*() or *free*() functions or the ISO C standard I/O
4116    library, both of which traditionally use data structures in a non-reentrant manner. Because any
4117    combination of different functions using a common data structure can cause reentrancy
4118    problems, POSIX.1 does not define the behavior when any unsafe function is called in a signal
4119    handler that interrupts any unsafe function.

4120    The only realtime extension to signal actions is the addition of the additional parameters to the
4121    signal-catching function. This extension has been explained and motivated in the previous
4122    section. In making this extension, though, developers of POSIX.1b ran into issues relating to

4123   function prototypes. In response to input from the POSIX.1 standard developers, members were
4124   added to the **sigaction** structure to specify function prototypes for the newer signal-catching
4125   function specified by POSIX.1b. These members follow changes that are being made to POSIX.1.
4126   Note that IEEE Std 1003.1-2001 explicitly states that these fields may overlap so that a union can
4127   be defined. This enabled existing implementations of POSIX.1 to maintain binary-compatibility
4128   when these extensions were added.

4129   The **siginfo_t** structure was adopted for passing the application-defined value to match existing
4130   practice, but the existing practice has no provision for an application-defined value, so this was
4131   added. Note that POSIX normally reserves the ''_t'' type designation for opaque types. The
4132   **siginfo_t** structure breaks with this convention to follow existing practice and thus promote
4133   portability. Standardization of the existing practice for the other members of this structure may
4134   be addressed in the future.

4135   Although it is not explicitly visible to applications, there are additional semantics for signal
4136   actions implied by queued signals and their interaction with other POSIX.1b realtime functions.
4137   Specifically:

4138        • It is not necessary to queue signals whose action is SIG_IGN.

4139        • For implementations that support POSIX.1b timers, some interaction with the timer functions
4140          at signal delivery is implied to manage the timer overrun count.

4141   IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/5 is applied, reordering the RTS shaded   |
4142   text under the third and fourth paragraphs of the SIG_DFL description. This corrects an earlier   |
4143   editorial error in this section.                                                            |

4144   IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/6 is applied, adding the *abort*() function   |
4145   to the list of async-cancel-safe functions.                                                  |

4146   *B.2.4.4*   *Signal Effects on Other Functions*

4147   The most common behavior of an interrupted function after a signal-catching function returns is
4148   for the interrupted function to give an [EINTR] error unless the SA_RESTART flag is in effect for
4149   the signal. However, there are a number of specific exceptions, including *sleep*() and certain
4150   situations with *read*() and *write*().

4151   The historical implementations of many functions defined by IEEE Std 1003.1-2001 are not
4152   interruptible, but delay delivery of signals generated during their execution until after they
4153   complete. This is never a problem for functions that are guaranteed to complete in a short
4154   (imperceptible to a human) period of time. It is normally those functions that can suspend a
4155   process indefinitely or for long periods of time (for example, *wait*(), *pause*(), *sigsuspend*(), *sleep*(),
4156   or *read*()/*write*() on a slow device like a terminal) that are interruptible. This permits
4157   applications to respond to interactive signals or to set timeouts on calls to most such functions
4158   with *alarm*(). Therefore, implementations should generally make such functions (including ones
4159   defined as extensions) interruptible.

4160   Functions not mentioned explicitly as interruptible may be so on some implementations,
4161   possibly as an extension where the function gives an [EINTR] error. There are several functions
4162   (for example, *getpid*(), *getuid*()) that are specified as never returning an error, which can thus
4163   never be extended in this way.

4164   If a signal-catching function returns while the SA_RESTART flag is in effect, an interrupted
4165   function is restarted at the point it was interrupted. Conforming applications cannot make
4166   assumptions about the internal behavior of interrupted functions, even if the functions are
4167   async-signal-safe. For example, suppose the *read*() function is interrupted with SA_RESTART in
4168   effect, the signal-catching function closes the file descriptor being read from and returns, and the

4169   *read*( ) function is then restarted; in this case the application cannot assume that the *read*( )
4170   function will give an [EBADF] error, since *read*( ) might have checked the file descriptor for
4171   validity before being interrupted.

## B.2.5   Standard I/O Streams

4173   *B.2.5.1   Interaction of File Descriptors and Standard I/O Streams*

4174   There is no additional rationale provided for this section.

4175   *B.2.5.2   Stream Orientation and Encoding Rules*

4176   There is no additional rationale provided for this section.

## B.2.6   STREAMS

4178   STREAMS are introduced into IEEE Std 1003.1-2001 as part of the alignment with the Single
4179   UNIX Specification, but marked as an option in recognition that not all systems may wish to
4180   implement the facility. The option within IEEE Std 1003.1-2001 is denoted by the XSR margin
4181   marker. The standard developers made this option independent of the XSI option.

4182   STREAMS are a method of implementing network services and other character-based
4183   input/output mechanisms, with the STREAM being a full-duplex connection between a process
4184   and a device. STREAMS provides direct access to protocol modules, and optional protocol
4185   modules can be interposed between the process-end of the STREAM and the device-driver at the
4186   device-end of the STREAM. Pipes can be implemented using the STREAMS mechanism, so they
4187   can provide process-to-process as well as process-to-device communications.

4188   This section introduces STREAMS I/O, the message types used to control them, an overview of
4189   the priority mechanism, and the interfaces used to access them.

4190   *B.2.6.1   Accessing STREAMS*

4191   There is no additional rationale provided for this section.

## B.2.7   XSI Interprocess Communication

4193   There are two forms of IPC supported as options in IEEE Std 1003.1-2001. The traditional
4194   System V IPC routines derived from the SVID—that is, the *msg*\*( ), *sem*\*( ), and *shm*\*( )
4195   interfaces—are mandatory on XSI-conformant systems. Thus, all XSI-conformant systems
4196   provide the same mechanisms for manipulating messages, shared memory, and semaphores.

4197   In addition, the POSIX Realtime Extension provides an alternate set of routines for those systems
4198   supporting the appropriate options.

4199   The application writer is presented with a choice: the System V interfaces or the POSIX
4200   interfaces (loosely derived from the Berkeley interfaces). The XSI profile prefers the System V
4201   interfaces, but the POSIX interfaces may be more suitable for realtime or other performance-
4202   sensitive applications.

4203 *B.2.7.1*   *IPC General Information*

4204    General information that is shared by all three mechanisms is described in this section. The
4205    common permissions mechanism is briefly introduced, describing the mode bits, and how they
4206    are used to determine whether or not a process has access to read or write/alter the appropriate
4207    instance of one of the IPC mechanisms. All other relevant information is contained in the
4208    reference pages themselves.

4209    The semaphore type of IPC allows processes to communicate through the exchange of
4210    semaphore values. A semaphore is a positive integer. Since many applications require the use of
4211    more than one semaphore, XSI-conformant systems have the ability to create sets or arrays of
4212    semaphores.

4213    Calls to support semaphores include:

4214        *semctl*( ), *semget*( ), *semop*( )

4215    Semaphore sets are created by using the *semget*( ) function.

4216    The message type of IPC allows processes to communicate through the exchange of data stored
4217    in buffers. This data is transmitted between processes in discrete portions known as messages.

4218    Calls to support message queues include:

4219        *msgctl*( ), *msgget*( ), *msgrcv*( ), *msgsnd*( )

4220    The shared memory type of IPC allows two or more processes to share memory and
4221    consequently the data contained therein. This is done by allowing processes to set up access to a
4222    common memory address space.  This sharing of memory provides a fast means of exchange of
4223    data between processes.

4224    Calls to support shared memory include:

4225        *shmctl*( ), *shmdt*( ), *shmget*( )

4226    The *ftok*( ) interface is also provided.

4227 **B.2.8**   **Realtime**

4228    **Advisory Information**

4229    POSIX.1b contains an Informative Annex with proposed interfaces for ''realtime files''. These
4230    interfaces could determine groups of the exact parameters required to do ''direct I/O'' or
4231    ''extents''. These interfaces were objected to by a significant portion of the balloting group as too
4232    complex. A conforming application had little chance of correctly navigating the large parameter
4233    space to match its desires to the system. In addition, they only applied to a new type of file
4234    (realtime files) and they told the implementation exactly what to do as opposed to advising the
4235    implementation on application behavior and letting it optimize for the system the (portable)
4236    application was running on. For example, it was not clear how a system that had a disk array
4237    should set its parameters.

4238    There seemed to be several overall goals:

4239    • Optimizing sequential access

4240    • Optimizing caching behavior

4241    • Optimizing I/O data transfer

4242    • Preallocation

4243   The advisory interfaces, *posix_fadvise*() and *posix_madvise*(), satisfy the first two goals. The
4244   POSIX_FADV_SEQUENTIAL     and     POSIX_MADV_SEQUENTIAL     advice     tells     the
4245   implementation to expect serial access. Typically the system will prefetch the next several serial
4246   accesses in order to overlap I/O. It may also free previously accessed serial data if memory is
4247   tight. If the application is not doing serial access it can use POSIX_FADV_WILLNEED and
4248   POSIX_MADV_WILLNEED to accomplish I/O overlap, as required. When the application
4249   advises POSIX_FADV_RANDOM or POSIX_MADV_RANDOM behavior, the implementation
4250   usually tries to fetch a minimum amount of data with each request and it does not expect much
4251   locality. POSIX_FADV_DONTNEED and POSIX_MADV_DONTNEED allow the system to free
4252   up caching resources as the data will not be required in the near future.

4253   POSIX_FADV_NOREUSE tells the system that caching the specified data is not optimal. For file
4254   I/O, the transfer should go directly to the user buffer instead of being cached internally by the
4255   implementation. To portably perform direct disk I/O on all systems, the application must
4256   perform its I/O transfers according to the following rules:

   1.   The user buffer should be aligned according to the {POSIX_REC_XFER_ALIGN} *pathconf*()
4257        variable.
4258

   2.   The number of bytes transferred in an I/O operation should be a multiple of the
4259        {POSIX_ALLOC_SIZE_MIN} *pathconf*() variable.
4260

   3.   The offset into the file at the start of an I/O operation should be a multiple of the
4261        {POSIX_ALLOC_SIZE_MIN} *pathconf*() variable.
4262

   4.   The application should ensure that all threads which open a given file specify
4263        POSIX_FADV_NOREUSE to be sure that there is no unexpected interaction between
4264        threads using buffered I/O and threads using direct I/O to the same file.
4265

4266   In some cases, a user buffer must be properly aligned in order to be transferred directly to/from
4267   the device. The {POSIX_REC_XFER_ALIGN} *pathconf*() variable tells the application the proper
4268   alignment.

4269   The preallocation goal is met by the space control function, *posix_fallocate*().  The application can
4270   use *posix_fallocate*() to guarantee no [ENOSPC] errors and to improve performance by prepaying
4271   any overhead required for block allocation.

4272   Implementations may use information conveyed by a previous *posix_fadvise*() call to influence
4273   the manner in which allocation is performed. For example, if an application did the following
4274   calls:

```
4275       fd = open("file");
4276       posix_fadvise(fd, offset, len, POSIX_FADV_SEQUENTIAL);
4277       posix_fallocate(fd, len, size);
```

4278   an implementation might allocate the file contiguously on disk.

4279   Finally,        the        *pathconf*()        variables        {POSIX_REC_MIN_XFER_SIZE},
4280   {POSIX_REC_MAX_XFER_SIZE}, and {POSIX_REC_INCR_XFER_SIZE} tell the application a
4281   range of transfer sizes that are recommended for best I/O performance.

4282   Where bounded response time is required, the vendor can supply the appropriate settings of the
4283   advisories to achieve a guaranteed performance level.

4284   The interfaces meet the goals while allowing applications using regular files to take advantage of
4285   performance optimizations. The interfaces tell the implementation expected application
4286   behavior which the implementation can use to optimize performance on a particular system
4287   with a particular dynamic load.

4288  The *posix_memalign*( ) function was added to allow for the allocation of specifically aligned
4289  buffers; for example, for {POSIX_REC_XFER_ALIGN}.

4290  The working group also considered the alternative of adding a function which would return an
4291  aligned pointer to memory within a user-supplied buffer. This was not considered to be the best
4292  method, because it potentially wastes large amounts of memory when buffers need to be aligned
4293  on large alignment boundaries.

4294  **Message Passing**

4295  This section provides the rationale for the definition of the message passing interface in
4296  IEEE Std 1003.1-2001. This is presented in terms of the objectives, models, and requirements
4297  imposed upon this interface.

4298  • Objectives

4299  Many applications, including both realtime and database applications, require a means of
4300  passing arbitrary amounts of data between cooperating processes comprising the overall
4301  application on one or more processors. Many conventional interfaces for interprocess
4302  communication are insufficient for realtime applications in that efficient and deterministic
4303  data passing methods cannot be implemented. This has prompted the definition of message
4304  passing interfaces providing these facilities:

4305  — Open a message queue.

4306  — Send a message to a message queue.

4307  — Receive a message from a queue, either synchronously or asynchronously.

4308  — Alter message queue attributes for flow and resource control.

4309  It is assumed that an application may consist of multiple cooperating processes and that
4310  these processes may wish to communicate and coordinate their activities. The message
4311  passing facility described in IEEE Std 1003.1-2001 allows processes to communicate through
4312  system-wide queues. These message queues are accessed through names that may be
4313  pathnames. A message queue can be opened for use by multiple sending and/or multiple
4314  receiving processes.

4315  • Background on Embedded Applications

4316  Interprocess communication utilizing message passing is a key facility for the construction of
4317  deterministic, high-performance realtime applications. The facility is present in all realtime
4318  systems and is the framework upon which the application is constructed. The performance of
4319  the facility is usually a direct indication of the performance of the resulting application.

4320  Realtime applications, especially for embedded systems, are typically designed around the
4321  performance constraints imposed by the message passing mechanisms. Applications for
4322  embedded systems are typically very tightly constrained. Application writers expect to
4323  design and control the entire system. In order to minimize system costs, the writer will
4324  attempt to use all resources to their utmost and minimize the requirement to add additional
4325  memory or processors.

4326  The embedded applications usually share address spaces and only a simple message passing
4327  mechanism is required. The application can readily access common data incurring only
4328  mutual-exclusion overheads. The models desired are the simplest possible with the
4329  application building higher-level facilities only when needed.

4330          • Requirements

4331          The following requirements determined the features of the message passing facilities defined
4332          in IEEE Std 1003.1-2001:

4333          — Naming of Message Queues

4334          The mechanism for gaining access to a message queue is a pathname evaluated in a
4335          context that is allowed to be a file system name space, or it can be independent of any file
4336          system. This is a specific attempt to allow implementations based on either method in
4337          order to address both embedded systems and to also allow implementation in larger
4338          systems.

4339          The interface of *mq_open*( ) is defined to allow but not require the access control and name
4340          conflicts resulting from utilizing a file system for name resolution.  All required behavior
4341          is specified for the access control case. Yet a conforming implementation, such as an
4342          embedded system kernel, may define that there are no distinctions between users and
4343          may define that all processes have all access privileges.

4344          — Embedded System Naming

4345          Embedded systems need to be able to utilize independent name spaces for accessing the
4346          various system objects. They typically do not have a file system, precluding its utilization
4347          as a common name resolution mechanism. The modularity of an embedded system limits
4348          the connections between separate mechanisms that can be allowed.

4349          Embedded systems typically do not have any access protection. Since the system does not
4350          support the mixing of applications from different areas, and usually does not even have
4351          the concept of an authorization entity, access control is not useful.

4352          — Large System Naming

4353          On systems with more functionality, the name resolution must support the ability to use
4354          the file system as the name resolution mechanism/object storage medium and to have
4355          control over access to the objects. Utilizing the pathname space can result in further errors
4356          when the names conflict with other objects.

4357          — Fixed Size of Messages

4358          The interfaces impose a fixed upper bound on the size of messages that can be sent to a
4359          specific message queue. The size is set on an individual queue basis and cannot be
4360          changed dynamically.

4361          The purpose of the fixed size is to increase the ability of the system to optimize the
4362          implementation of *mq_send*( ) and *mq_receive*( ).  With fixed sizes of messages and fixed
4363          numbers of messages, specific message blocks can be pre-allocated. This eliminates a
4364          significant amount of checking for errors and boundary conditions. Additionally, an
4365          implementation can optimize data copying to maximize performance.  Finally, with a
4366          restricted range of message sizes, an implementation is better able to provide
4367          deterministic operations.

4368          — Prioritization of Messages

4369          Message prioritization allows the application to determine the order in which messages
4370          are received. Prioritization of messages is a key facility that is provided by most realtime
4371          kernels and is heavily utilized by the applications. The major purpose of having priorities
4372          in message queues is to avoid priority inversions in the message system, where a high-
4373          priority message is delayed behind one or more lower-priority messages. This allows the
4374          applications to be designed so that they do not need to be interrupted in order to change

4375 the flow of control when exceptional conditions occur. The prioritization does add
4376 additional overhead to the message operations in those cases it is actually used but a
4377 clever implementation can optimize for the FIFO case to make that more efficient.

4378 — Asynchronous Notification

4379 The interface supports the ability to have a task asynchronously notified of the
4380 availability of a message on the queue. The purpose of this facility is to allow the task to
4381 perform other functions and yet still be notified that a message has become available on
4382 the queue.

4383 To understand the requirement for this function, it is useful to understand two models of
4384 application design: a single task performing multiple functions and multiple tasks
4385 performing a single function.  Each of these models has advantages.

4386 Asynchronous notification is required to build the model of a single task performing
4387 multiple operations. This model typically results from either the expectation that
4388 interruption is less expensive than utilizing a separate task or from the growth of the
4389 application to include additional functions.

4390 **Semaphores**

4391 Semaphores are a high-performance process synchronization mechanism.  Semaphores are
4392 named by null-terminated strings of characters.

4393 A semaphore is created using the *sem_init*( ) function or the *sem_open*( ) function with the
4394 O_CREAT flag set in *oflag*.

4395 To use a semaphore, a process has to first initialize the semaphore or inherit an open descriptor
4396 for the semaphore via *fork*( ).

4397 A semaphore preserves its state when the last reference is closed. For example, if a semaphore
4398 has a value of 13 when the last reference is closed, it will have a value of 13 when it is next
4399 opened.

4400 When a semaphore is created, an initial state for the semaphore has to be provided. This value is
4401 a non-negative integer. Negative values are not possible since they indicate the presence of
4402 blocked processes. The persistence of any of these objects across a system crash or a system
4403 reboot is undefined. Conforming applications must not depend on any sort of persistence across
4404 a system reboot or a system crash.

4405 • Models and Requirements

4406 A realtime system requires synchronization and communication between the processes
4407 comprising the overall application. An efficient and reliable synchronization mechanism has
4408 to be provided in a realtime system that will allow more than one schedulable process
4409 mutually-exclusive access to the same resource. This synchronization mechanism has to
4410 allow for the optimal implementation of synchronization or systems implementors will
4411 define other, more cost-effective methods.

4412 At issue are the methods whereby multiple processes (tasks) can be designed and
4413 implemented to work together in order to perform a single function. This requires
4414 interprocess communication and synchronization. A semaphore mechanism is the lowest
4415 level of synchronization that can be provided by an operating system.

4416 A semaphore is defined as an object that has an integral value and a set of blocked processes
4417 associated with it. If the value is positive or zero, then the set of blocked processes is empty;
4418 otherwise, the size of the set is equal to the absolute value of the semaphore value.  The value
4419 of the semaphore can be incremented or decremented by any process with access to the

4420    semaphore and must be done as an indivisible operation. When a semaphore value is less
4421    than or equal to zero, any process that attempts to lock it again will block or be informed that
4422    it is not possible to perform the operation.

4423    A semaphore may be used to guard access to any resource accessible by more than one
4424    schedulable task in the system. It is a global entity and not associated with any particular
4425    process. As such, a method of obtaining access to the semaphore has to be provided by the
4426    operating system. A process that wants access to a critical resource (section) has to wait on
4427    the semaphore that guards that resource. When the semaphore is locked on behalf of a
4428    process, it knows that it can utilize the resource without interference by any other
4429    cooperating process in the system. When the process finishes its operation on the resource,
4430    leaving it in a well-defined state, it posts the semaphore, indicating that some other process
4431    may now obtain the resource associated with that semaphore.

4432    In this section, mutexes and condition variables are specified as the synchronization
4433    mechanisms between threads.

4434    These primitives are typically used for synchronizing threads that share memory in a single
4435    process. However, this section provides an option allowing the use of these synchronization
4436    interfaces and objects between processes that share memory, regardless of the method for
4437    sharing memory.

4438    Much experience with semaphores shows that there are two distinct uses of synchronization:
4439    locking, which is typically of short duration; and waiting, which is typically of long or
4440    unbounded duration. These distinct usages map directly onto mutexes and condition
4441    variables, respectively.

4442    Semaphores are provided in IEEE Std 1003.1-2001 primarily to provide a means of
4443    synchronization for processes; these processes may or may not share memory. Mutexes and
4444    condition variables are specified as synchronization mechanisms between threads; these
4445    threads always share (some) memory. Both are synchronization paradigms that have been in
4446    widespread use for a number of years. Each set of primitives is particularly well matched to
4447    certain problems.

4448    With respect to binary semaphores, experience has shown that condition variables and
4449    mutexes are easier to use for many synchronization problems than binary semaphores. The
4450    primary reason for this is the explicit appearance of a Boolean predicate that specifies when
4451    the condition wait is satisfied. This Boolean predicate terminates a loop, including the call to
4452    *pthread_cond_wait*(). As a result, extra wakeups are benign since the predicate governs
4453    whether the thread will actually proceed past the condition wait. With stateful primitives,
4454    such as binary semaphores, the wakeup in itself typically means that the wait is satisfied. The
4455    burden of ensuring correctness for such waits is thus placed on *all* signalers of the semaphore
4456    rather than on an *explicitly coded* Boolean predicate located at the condition wait. Experience
4457    has shown that the latter creates a major improvement in safety and ease-of-use.

4458    Counting semaphores are well matched to dealing with producer/consumer problems,
4459    including those that might exist between threads of different processes, or between a signal
4460    handler and a thread. In the former case, there may be little or no memory shared by the
4461    processes; in the latter case, one is not communicating between co-equal threads, but
4462    between a thread and an interrupt-like entity. It is for these reasons that IEEE Std 1003.1-2001
4463    allows semaphores to be used by threads.

4464    Mutexes and condition variables have been effectively used with and without priority
4465    inheritance, priority ceiling, and other attributes to synchronize threads that share memory.
4466    The efficiency of their implementation is comparable to or better than that of other
4467    synchronization primitives that are sometimes harder to use (for example, binary

4468  semaphores). Furthermore, there is at least one known implementation of Ada tasking that
4469  uses these primitives. Mutexes and condition variables together constitute an appropriate,
4470  sufficient, and complete set of inter-thread synchronization primitives.

4471  Efficient multi-threaded applications require high-performance synchronization primitives.
4472  Considerations of efficiency and generality require a small set of primitives upon which more
4473  sophisticated synchronization functions can be built.

4474  • Standardization Issues

4475  It is possible to implement very high-performance semaphores using test-and-set
4476  instructions on shared memory locations. The library routines that implement such a high-
4477  performance interface have to properly ensure that a *sem_wait*() or *sem_trywait*() operation
4478  that cannot be performed will issue a blocking semaphore system call or properly report the
4479  condition to the application. The same interface to the application program would be
4480  provided by a high-performance implementation.

4481  *B.2.8.1  Realtime Signals*

4482  **Realtime Signals Extension**

4483  This portion of the rationale presents models, requirements, and standardization issues relevant
4484  to the Realtime Signals Extension. This extension provides the capability required to support
4485  reliable, deterministic, asynchronous notification of events. While a new mechanism,
4486  unencumbered by the historical usage and semantics of POSIX.1 signals, might allow for a more
4487  efficient implementation, the application requirements for event notification can be met with a
4488  small number of extensions to signals. Therefore, a minimal set of extensions to signals to
4489  support the application requirements is specified.

4490  The realtime signal extensions specified in this section are used by other realtime functions
4491  requiring asynchronous notification:

4492  • Models

4493  The model supported is one of multiple cooperating processes, each of which handles
4494  multiple asynchronous external events. Events represent occurrences that are generated as
4495  the result of some activity in the system. Examples of occurrences that can constitute an
4496  event include:

4497  — Completion of an asynchronous I/O request

4498  — Expiration of a POSIX.1b timer

4499  — Arrival of an interprocess message

4500  — Generation of a user-defined event

4501  Processing of these events may occur synchronously via polling for event notifications or
4502  asynchronously via a software interrupt mechanism. Existing practice for this model is well
4503  established for traditional proprietary realtime operating systems, realtime executives, and
4504  realtime extended POSIX-like systems.

4505  A contrasting model is that of ''cooperating sequential processes'' where each process
4506  handles a single priority of events via polling. Each process blocks while waiting for events,
4507  and each process depends on the preemptive, priority-based process scheduling mechanism
4508  to arbitrate between events of different priority that need to be processed concurrently.
4509  Existing practice for this model is also well established for small realtime executives that
4510  typically execute in an unprotected physical address space, but it is just emerging in the

4511 context of a fuller function operating system with multiple virtual address spaces.

4512 It could be argued that the cooperating sequential process model, and the facilities supported
4513 by the POSIX Threads Extension obviate a software interrupt model. But, even with the
4514 cooperating sequential process model, the need has been recognized for a software interrupt
4515 model to handle exceptional conditions and process aborting, so the mechanism must be
4516 supported in any case. Furthermore, it is not the purview of IEEE Std 1003.1-2001 to attempt
4517 to convince realtime practitioners that their current application models based on software
4518 interrupts are ''broken'' and should be replaced by the cooperating sequential process model.
4519 Rather, it is the charter of IEEE Std 1003.1-2001 to provide standard extensions to
4520 mechanisms that support existing realtime practice.

4521 • Requirements

4522 This section discusses the following realtime application requirements for asynchronous
4523 event notification:

4524 — Reliable delivery of asynchronous event notification

4525 The events notification mechanism guarantees delivery of an event notification.
4526 Asynchronous operations (such as asynchronous I/O and timers) that complete
4527 significantly after they are invoked have to guarantee that delivery of the event
4528 notification can occur at the time of completion.

4529 — Prioritized handling of asynchronous event notifications

4530 The events notification mechanism supports the assigning of a user function as an event
4531 notification handler. Furthermore, the mechanism supports the preemption of an event
4532 handler function by a higher priority event notification and supports the selection of the
4533 highest priority pending event notification when multiple notifications (of different
4534 priority) are pending simultaneously.

4535 The model here is based on hardware interrupts. Asynchronous event handling allows
4536 the application to ensure that time-critical events are immediately processed when
4537 delivered, without the indeterminism of being at a random location within a polling loop.
4538 Use of handler priority allows the specification of how handlers are interrupted by other
4539 higher priority handlers.

4540 — Differentiation between multiple occurrences of event notifications of the same type

4541 The events notification mechanism passes an application-defined value to the event
4542 handler function. This value can be used for a variety of purposes, such as enabling the
4543 application to identify which of several possible events of the same type (for example,
4544 timer expirations) has occurred.

4545 — Polled reception of asynchronous event notifications

4546 The events notification mechanism supports blocking and non-blocking polls for
4547 asynchronous event notification.

4548 The polled mode of operation is often preferred over the interrupt mode by those
4549 practitioners accustomed to this model. Providing support for this model facilitates the
4550 porting of applications based on this model to POSIX.1b conforming systems.

4551 — Deterministic response to asynchronous event notifications

4552 The events notification mechanism does not preclude implementations that provide
4553 deterministic event dispatch latency and minimizes the number of system calls needed to
4554 use the event facilities during realtime processing.

4555 • Rationale for Extension

4556 POSIX.1 signals have many of the characteristics necessary to support the asynchronous
4557 handling of event notifications, and the Realtime Signals Extension addresses the following
4558 deficiencies in the POSIX.1 signal mechanism:

4559 — Signals do not support reliable delivery of event notification. Subsequent occurrences of
4560 a pending signal are not guaranteed to be delivered.

4561 — Signals do not support prioritized delivery of event notifications. The order of signal
4562 delivery when multiple unblocked signals are pending is undefined.

4563 — Signals do not support the differentiation between multiple signals of the same type.

4564 *B.2.8.2    Asynchronous I/O*

4565 Many applications need to interact with the I/O subsystem in an asynchronous manner. The
4566 asynchronous I/O mechanism provides the ability to overlap application processing and I/O
4567 operations initiated by the application. The asynchronous I/O mechanism allows a single
4568 process to perform I/O simultaneously to a single file multiple times or to multiple files
4569 multiple times.

4570 **Overview**

4571 Asynchronous I/O operations proceed in logical parallel with the processing done by the
4572 application after the asynchronous I/O has been initiated. Other than this difference,
4573 asynchronous I/O behaves similarly to normal I/O using *read*(), *write*(), *lseek*(), and *fsync*().
4574 The effect of issuing an asynchronous I/O request is as if a separate thread of execution were to
4575 perform atomically the implied *lseek*() operation, if any, and then the requested I/O operation
4576 (either *read*(), *write*(), or *fsync*()). There is no seek implied with a call to *aio_fsync*(). Concurrent
4577 asynchronous operations and synchronous operations applied to the same file update the file as
4578 if the I/O operations had proceeded serially.

4579 When asynchronous I/O completes, a signal can be delivered to the application to indicate the
4580 completion of the I/O. This signal can be used to indicate that buffers and control blocks used
4581 for asynchronous I/O can be reused. Signal delivery is not required for an asynchronous
4582 operation and may be turned off on a per-operation basis by the application. Signals may also be
4583 synchronously polled using *aio_suspend*(), *sigtimedwait*(), or *sigwaitinfo*().

4584 Normal I/O has a return value and an error status associated with it. Asynchronous I/O returns
4585 a value and an error status when the operation is first submitted, but that only relates to whether
4586 the operation was successfully queued up for servicing. The I/O operation itself also has a
4587 return status and an error value. To allow the application to retrieve the return status and the
4588 error value, functions are provided that, given the address of an asynchronous I/O control
4589 block, yield the return and error status associated with the operation. Until an asynchronous I/O
4590 operation is done, its error status is [EINPROGRESS]. Thus, an application can poll for
4591 completion of an asynchronous I/O operation by waiting for the error status to become equal to
4592 a value other than [EINPROGRESS]. The return status of an asynchronous I/O operation is
4593 undefined so long as the error status is equal to [EINPROGRESS].

4594 Storage for asynchronous operation return and error status may be limited. Submission of
4595 asynchronous I/O operations may fail if this storage is exceeded. When an application retrieves
4596 the return status of a given asynchronous operation, therefore, any system-maintained storage
4597 used for this status and the error status may be reclaimed for use by other asynchronous
4598 operations.

4599 Asynchronous I/O can be performed on file descriptors that have been enabled for POSIX.1b
4600 synchronized I/O. In this case, the I/O operation still occurs asynchronously, as defined herein;
4601 however, the asynchronous operation I/O in this case is not completed until the I/O has reached
4602 either the state of synchronized I/O data integrity completion or synchronized I/O file integrity
4603 completion, depending on the sort of synchronized I/O that is enabled on the file descriptor.

4604 **Models**

4605 Three models illustrate the use of asynchronous I/O: a journalization model, a data acquisition
4606 model, and a model of the use of asynchronous I/O in supercomputing applications.

4607 • Journalization Model

4608 Many realtime applications perform low-priority journalizing functions. Journalizing
4609 requires that logging records be queued for output without blocking the initiating process.

4610 • Data Acquisition Model

4611 A data acquisition process may also serve as a model. The process has two or more channels
4612 delivering intermittent data that must be read within a certain time. The process issues one
4613 asynchronous read on each channel. When one of the channels needs data collection, the
4614 process reads the data and posts it through an asynchronous write to secondary memory for
4615 future processing.

4616 • Supercomputing Model

4617 The supercomputing community has used asynchronous I/O much like that specified in
4618 POSIX.1 for many years. This community requires the ability to perform multiple I/O
4619 operations to multiple devices with a minimal number of entries to ''the system''; each entry
4620 to ''the system'' provokes a major delay in operations when compared to the normal progress
4621 made by the application. This existing practice motivated the use of combined *lseek*( ) and
4622 *read*( ) or *write*( ) calls, as well as the *lio_listio*( ) call. Another common practice is to disable
4623 signal notification for I/O completion, and simply poll for I/O completion at some interval
4624 by which the I/O should be completed. Likewise, interfaces like *aio_cancel*( ) have been in
4625 successful commercial use for many years. Note also that an underlying implementation of
4626 asynchronous I/O will require the ability, at least internally, to cancel outstanding
4627 asynchronous I/O, at least when the process exits. (Consider an asynchronous read from a
4628 terminal, when the process intends to exit immediately.)

4629 **Requirements**

4630 Asynchronous input and output for realtime implementations have these requirements:

4631 • The ability to queue multiple asynchronous read and write operations to a single open
4632 instance. Both sequential and random access should be supported.

4633 • The ability to queue asynchronous read and write operations to multiple open instances.

4634 • The ability to obtain completion status information by polling and/or asynchronous event
4635 notification.

4636 • Asynchronous event notification on asynchronous I/O completion is optional.

4637 • It has to be possible for the application to associate the event with the *aiocbp* for the operation
4638 that generated the event.

4639 • The ability to cancel queued requests.

4640 • The ability to wait upon asynchronous I/O completion in conjunction with other types of
4641 events.

4642       • The ability to accept an *aio_read*( ) and an *aio_cancel*( ) for a device that accepts a *read*( ), and
4643         the ability to accept an *aio_write*( ) and an *aio_cancel*( ) for a device that accepts a *write*( ). This
4644         does not imply that the operation is asynchronous.

4645      **Standardization Issues**

4646      The following issues are addressed by the standardization of asynchronous I/O:

4647      • Rationale for New Interface

4648        Non-blocking I/O does not satisfy the needs of either realtime or high-performance
4649        computing models; these models require that a process overlap program execution and I/O
4650        processing. Realtime applications will often make use of direct I/O to or from the address
4651        space of the process, or require synchronized (unbuffered) I/O; they also require the ability
4652        to overlap this I/O with other computation. In addition, asynchronous I/O allows an
4653        application to keep a device busy at all times, possibly achieving greater throughput.
4654        Supercomputing and database architectures will often have specialized hardware that can
4655        provide true asynchrony underlying the logical asynchrony provided by this interface. In
4656        addition, asynchronous I/O should be supported by all types of files and devices in the same
4657        manner.

4658      • Effect of Buffering

4659        If asynchronous I/O is performed on a file that is buffered prior to being actually written to
4660        the device, it is possible that asynchronous I/O will offer no performance advantage over
4661        normal I/O; the cycles *stolen* to perform the asynchronous I/O will be taken away from the
4662        running process and the I/O will occur at interrupt time. This potential lack of gain in
4663        performance in no way obviates the need for asynchronous I/O by realtime applications,
4664        which very often will use specialized hardware support, multiple processors, and/or
4665        unbuffered, synchronized I/O.

4666  *B.2.8.3*   *Memory Management*

4667      All memory management and shared memory definitions are located in the **<sys/mman.h>**
4668      header. This is for alignment with historical practice.

4669      IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/7 is applied, correcting the shading and   |
4670      margin markers in the introduction to Section 2.8.3.1.                               |

4671      **Memory Locking Functions**

4672      This portion of the rationale presents models, requirements, and standardization issues relevant
4673      to process memory locking.

4674      • Models

4675        Realtime systems that conform to IEEE Std 1003.1-2001 are expected (and desired) to be
4676        supported on systems with demand-paged virtual memory management, non-paged
4677        swapping memory management, and physical memory systems with no memory
4678        management hardware. The general case, however, is the demand-paged, virtual memory
4679        system with each POSIX process running in a virtual address space. Note that this includes
4680        architectures where each process resides in its own virtual address space and architectures
4681        where the address space of each process is only a portion of a larger global virtual address
4682        space.

4683        The concept of memory locking is introduced to eliminate the indeterminacy introduced by
4684        paging and swapping, and to support an upper bound on the time required to access the
4685        memory mapped into the address space of a process. Ideally, this upper bound will be the

        

4686 same as the time required for the processor to access ''main memory'', including any address
4687 translation and cache miss overheads. But some implementations—primarily on
4688 mainframes—will not actually force locked pages to be loaded and held resident in main
4689 memory. Rather, they will handle locked pages so that accesses to these pages will meet the
4690 performance metrics for locked process memory in the implementation. Also, although it is
4691 not, for example, the intention that this interface, as specified, be used to lock process
4692 memory into ''cache'', it is conceivable that an implementation could support a large static
4693 RAM memory and define this as ''main memory'' and use a large[r] dynamic RAM as
4694 ''backing store''. These interfaces could then be interpreted as supporting the locking of
4695 process memory into the static RAM. Support for multiple levels of backing store would
4696 require extensions to these interfaces.

4697 Implementations may also use memory locking to guarantee a fixed translation between
4698 virtual and physical addresses where such is beneficial to improving determinancy for
4699 direct-to/from-process input/output. IEEE Std 1003.1-2001 does not guarantee to the
4700 application that the virtual-to-physical address translations, if such exist, are fixed, because
4701 such behavior would not be implementable on all architectures on which implementations of
4702 IEEE Std 1003.1-2001 are expected. But IEEE Std 1003.1-2001 does mandate that an
4703 implementation define, for the benefit of potential users, whether or not locking guarantees
4704 fixed translations.

4705 Memory locking is defined with respect to the address space of a process. Only the pages
4706 mapped into the address space of a process may be locked by the process, and when the
4707 pages are no longer mapped into the address space—for whatever reason—the locks
4708 established with respect to that address space are removed. Shared memory areas warrant
4709 special mention, as they may be mapped into more than one address space or mapped more
4710 than once into the address space of a process; locks may be established on pages within these
4711 areas with respect to several of these mappings. In such a case, the lock state of the
4712 underlying physical pages is the logical OR of the lock state with respect to each of the
4713 mappings. Only when all such locks have been removed are the shared pages considered
4714 unlocked.

4715 In recognition of the page granularity of Memory Management Units (MMU), and in order to
4716 support locking of ranges of address space, memory locking is defined in terms of ''page''
4717 granularity. That is, for the interfaces that support an address and size specification for the
4718 region to be locked, the address must be on a page boundary, and all pages mapped by the
4719 specified range are locked, if valid. This means that the length is implicitly rounded up to a
4720 multiple of the page size. The page size is implementation-defined and is available to
4721 applications as a compile-time symbolic constant or at runtime via *sysconf*( ).

4722 A ''real memory'' POSIX.1b implementation that has no MMU could elect not to support
4723 these interfaces, returning [ENOSYS]. But an application could easily interpret this as
4724 meaning that the implementation would unconditionally page or swap the application when
4725 such is not the case. It is the intention of IEEE Std 1003.1-2001 that such a system could define
4726 these interfaces as ''NO-OPs'', returning success without actually performing any function
4727 except for mandated argument checking.

4728 • Requirements

4729 For realtime applications, memory locking is generally considered to be required as part of
4730 application initialization. This locking is performed after an application has been loaded (that
4731 is, *exec*'d) and the program remains locked for its entire lifetime. But to support applications
4732 that undergo major mode changes where, in one mode, locking is required, but in another it
4733 is not, the specified interfaces allow repeated locking and unlocking of memory within the
4734 lifetime of a process.

When a realtime application locks its address space, it should not be necessary for the application to then ''touch'' all of the pages in the address space to guarantee that they are resident or else suffer potential paging delays the first time the page is referenced. Thus, IEEE Std 1003.1-2001 requires that the pages locked by the specified interfaces be resident when the locking functions return successfully.

Many architectures support system-managed stacks that grow automatically when the current extent of the stack is exceeded. A realtime application has a requirement to be able to ''preallocate'' sufficient stack space and lock it down so that it will not suffer page faults to grow the stack during critical realtime operation. There was no consensus on a portable way to specify how much stack space is needed, so IEEE Std 1003.1-2001 supports no specific interface for preallocating stack space. But an application can portably lock down a specific amount of stack space by specifying MCL_FUTURE in a call to *mlockall*( ) and then calling a dummy function that declares an automatic array of the desired size.

Memory locking for realtime applications is also generally considered to be an ''all or nothing'' proposition. That is, the entire process, or none, is locked down. But, for applications that have well-defined sections that need to be locked and others that do not, IEEE Std 1003.1-2001 supports an optional set of interfaces to lock or unlock a range of process addresses. Reasons for locking down a specific range include:

— An asynchronous event handler function that must respond to external events in a deterministic manner such that page faults cannot be tolerated

— An input/output ''buffer'' area that is the target for direct-to-process I/O, and the overhead of implicit locking and unlocking for each I/O call cannot be tolerated

Finally, locking is generally viewed as an ''application-wide'' function. That is, the application is globally aware of which regions are locked and which are not over time. This is in contrast to a function that is used temporarily within a ''third party'' library routine whose function is unknown to the application, and therefore must have no ''side effects''. The specified interfaces, therefore, do not support ''lock stacking'' or ''lock nesting'' within a process. But, for pages that are shared between processes or mapped more than once into a process address space, ''lock stacking'' is essentially mandated by the requirement that unlocking of pages that are mapped by more that one process or more than once by the same process does not affect locks established on the other mappings.

There was some support for ''lock stacking'' so that locking could be transparently used in functions or opaque modules. But the consensus was not to burden all implementations with lock stacking (and reference counting), and an implementation option was proposed. There were strong objections to the option because applications would have to support both options in order to remain portable. The consensus was to eliminate lock stacking altogether, primarily through overwhelming support for the System V ''m[un]lock[all]'' interface on which IEEE Std 1003.1-2001 is now based.

Locks are not inherited across *fork*( )s because some implementations implement *fork*( ) by creating new address spaces for the child. In such an implementation, requiring locks to be inherited would lead to new situations in which a fork would fail due to the inability of the system to lock sufficient memory to lock both the parent and the child. The consensus was that there was no benefit to such inheritance. Note that this does not mean that locks are removed when, for instance, a thread is created in the same address space.

Similarly, locks are not inherited across *exec* because some implementations implement *exec* by unmapping all of the pages in the address space (which, by definition, removes the locks on these pages), and maps in pages of the *exec*'d image. In such an implementation, requiring locks to be inherited would lead to new situations in which *exec* would fail. Reporting this

4783    failure would be very cumbersome to detect in time to report to the calling process, and no
4784    appropriate mechanism exists for informing the *exec*'d process of its status.

4785    It was determined that, if the newly loaded application required locking, it was the
4786    responsibility of that application to establish the locks. This is also in keeping with the
4787    general view that it is the responsibility of the application to be aware of all locks that are
4788    established.

4789    There was one request to allow (not mandate) locks to be inherited across *fork*(), and a
4790    request for a flag, MCL_INHERIT, that would specify inheritance of memory locks across
4791    *exec*s. Given the difficulties raised by this and the general lack of support for the feature in
4792    IEEE Std 1003.1-2001, it was not added. IEEE Std 1003.1-2001 does not preclude an
4793    implementation from providing this feature for administrative purposes, such as a ''run''
4794    command that will lock down and execute a specified application. Additionally, the rationale
4795    for the objection equated *fork*() with creating a thread in the address space.
4796    IEEE Std 1003.1-2001 does not mandate releasing locks when creating additional threads in
4797    an existing process.

4798    • Standardization Issues

4799    One goal of IEEE Std 1003.1-2001 is to define a set of primitives that provide the necessary
4800    functionality for realtime applications, with consideration for the needs of other application
4801    domains where such were identified, which is based to the extent possible on existing
4802    industry practice.

4803    The Memory Locking option is required by many realtime applications to tune performance.
4804    Such a facility is accomplished by placing constraints on the virtual memory system to limit
4805    paging of time of the process or of critical sections of the process. This facility should not be
4806    used by most non-realtime applications.

4807    Optional features provided in IEEE Std 1003.1-2001 allow applications to lock selected
4808    address ranges with the caveat that the process is responsible for being aware of the page
4809    granularity of locking and the unnested nature of the locks.

4810    **Mapped Files Functions**

4811    The Memory Mapped Files option provides a mechanism that allows a process to access files by
4812    directly incorporating file data into its address space. Once a file is ''mapped'' into a process
4813    address space, the data can be manipulated by instructions as memory. The use of mapped files
4814    can significantly reduce I/O data movement since file data does not have to be copied into
4815    process data buffers as in *read*() and *write*(). If more than one process maps a file, its contents
4816    are shared among them. This provides a low overhead mechanism by which processes can
4817    synchronize and communicate.

4818    • Historical Perspective

4819    Realtime applications have historically been implemented using a collection of cooperating
4820    processes or tasks. In early systems, these processes ran on bare hardware (that is, without an
4821    operating system) with no memory relocation or protection. The application paradigms that
4822    arose from this environment involve the sharing of data between the processes.

4823    When realtime systems were implemented on top of vendor-supplied operating systems, the
4824    paradigm or performance benefits of direct access to data by multiple processes was still
4825    deemed necessary. As a result, operating systems that claim to support realtime applications
4826    must support the shared memory paradigm.

4827    Additionally, a number of realtime systems provide the ability to map specific sections of the
4828    physical address space into the address space of a process. This ability is required if an

4829   application is to obtain direct access to memory locations that have specific properties (for
4830   example, refresh buffers or display devices, dual ported memory locations, DMA target
4831   locations). The use of this ability is common enough to warrant some degree of
4832   standardization of its interface. This ability overlaps the general paradigm of shared
4833   memory in that, in both instances, common global objects are made addressable by
4834   individual processes or tasks.

4835   Finally, a number of systems also provide the ability to map process addresses to files. This
4836   provides both a general means of sharing persistent objects, and using files in a manner that
4837   optimizes memory and swapping space usage.

4838   Simple shared memory is clearly a special case of the more general file mapping capability.
4839   In addition, there is relatively widespread agreement and implementation of the file
4840   mapping interface. In these systems, many different types of objects can be mapped (for
4841   example, files, memory, devices, and so on) using the same mapping interfaces. This
4842   approach both minimizes interface proliferation and maximizes the generality of programs
4843   using the mapping interfaces.

4844   • Memory Mapped Files Usage

4845   A memory object can be concurrently mapped into the address space of one or more
4846   processes. The *mmap*() and *munmap*() functions allow a process to manipulate their address
4847   space by mapping portions of memory objects into it and removing them from it. When
4848   multiple processes map the same memory object, they can share access to the underlying
4849   data. Implementations may restrict the size and alignment of mappings to be on *page*-size
4850   boundaries. The page size, in bytes, is the value of the system-configurable variable
4851   {PAGESIZE}, typically accessed by calling *sysconf*() with a *name* argument of
4852   _SC_PAGESIZE. If an implementation has no restrictions on size or alignment, it may
4853   specify a 1-byte page size.

4854   To map memory, a process first opens a memory object. The *ftruncate*() function can be used
4855   to contract or extend the size of the memory object even when the object is currently
4856   mapped. If the memory object is extended, the contents of the extended areas are zeros.

4857   After opening a memory object, the application maps the object into its address space using
4858   the *mmap*() function call. Once a mapping has been established, it remains mapped until
4859   unmapped with *munmap*(), even if the memory object is closed. The *mprotect*() function can
4860   be used to change the memory protections initially established by *mmap*().

4861   A *close*() of the file descriptor, while invalidating the file descriptor itself, does not unmap
4862   any mappings established for the memory object. The address space, including all mapped
4863   regions, is inherited on *fork*(). The entire address space is unmapped on process termination
4864   or by successful calls to any of the *exec* family of functions.

4865   The *msync*() function is used to force mapped file data to permanent storage.

4866   • Effects on Other Functions

4867   When the Memory Mapped Files option is supported, the operation of the *open*(), *creat*(), and
4868   *unlink*() functions are a natural result of using the file system name space to map the global
4869   names for memory objects.

4870   The *ftruncate*() function can be used to set the length of a sharable memory object.

4871   The meaning of *stat*() fields other than the size and protection information is undefined on
4872   implementations where memory objects are not implemented using regular files. When
4873   regular files are used, the times reflect when the implementation updated the file image of
4874   the data, not when a process updated the data in memory.

4875    The operations of *fdopen*( ), *write*( ), *read*( ), and *lseek*( ) were made unspecified for objects
4876    opened with *shm_open*( ), so that implementations that did not implement memory objects as
4877    regular files would not have to support the operation of these functions on shared memory
4878    objects.

4879    The behavior of memory objects with respect to *close*( ), *dup*( ), *dup2*( ), *open*( ), *close*( ), *fork*( ),
4880    *_exit*( ), and the *exec* family of functions is the same as the behavior of the existing practice of
4881    the *mmap*( ) function.

4882    A memory object can still be referenced after a close. That is, any mappings made to the file
4883    are still in effect, and reads and writes that are made to those mappings are still valid and are
4884    shared with other processes that have the same mapping. Likewise, the memory object can
4885    still be used if any references remain after its name(s) have been deleted. Any references that
4886    remain after a close must not appear to the application as file descriptors.

4887    This is existing practice for *mmap*( ) and *close*( ). In addition, there are already mappings
4888    present (text, data, stack) that do not have open file descriptors. The text mapping in
4889    particular is considered a reference to the file containing the text. The desire was to treat all
4890    mappings by the process uniformly. Also, many modern implementations use *mmap*( ) to
4891    implement shared libraries, and it would not be desirable to keep file descriptors for each of
4892    the many libraries an application can use. It was felt there were many other existing
4893    programs that used this behavior to free a file descriptor, and thus IEEE Std 1003.1-2001
4894    could not forbid it and still claim to be using existing practice.

4895    For implementations that implement memory objects using memory only, memory objects
4896    will retain the memory allocated to the file after the last close and will use that same memory
4897    on the next open. Note that closing the memory object is not the same as deleting the name,
4898    since the memory object is still defined in the memory object name space.

4899    The locks of *fcntl*( ) do not block any read or write operation, including read or write access to
4900    shared memory or mapped files. In addition, implementations that only support shared
4901    memory objects should not be required to implement record locks. The reference to *fcntl*( ) is
4902    added to make this point explicitly. The other *fcntl*( ) commands are useful with shared
4903    memory objects.

4904    The size of pages that mapping hardware may be able to support may be a configurable
4905    value, or it may change based on hardware implementations. The addition of the
4906    _SC_PAGESIZE parameter to the *sysconf*( ) function is provided for determining the mapping
4907    page size at runtime.

4908    **Shared Memory Functions**

4909    Implementations may support the Shared Memory Objects option without supporting a general
4910    Memory Mapped Files option. Shared memory objects are named regions of storage that may be
4911    independent of the file system and can be mapped into the address space of one or more
4912    processes to allow them to share the associated memory.

4913    • Requirements

4914    Shared memory is used to share data among several processes, each potentially running at
4915    different priority levels, responding to different inputs, or performing separate tasks. Shared
4916    memory is not just simply providing common access to data, it is providing the fastest
4917    possible communication between the processes. With one memory write operation, a process
4918    can pass information to as many processes as have the memory region mapped.

4919    As a result, shared memory provides a mechanism that can be used for all other interprocess
4920    communication facilities. It may also be used by an application for implementing more

4921        sophisticated mechanisms than semaphores and message queues.

4922        The need for a shared memory interface is obvious for virtual memory systems, where the
4923        operating system is directly preventing processes from accessing each other's data. However,
4924        in unprotected systems, such as those found in some embedded controllers, a shared
4925        memory interface is needed to provide a portable mechanism to allocate a region of memory
4926        to be shared and then to communicate the address of that region to other processes.

4927        This, then, provides the minimum functionality that a shared memory interface must have in
4928        order to support realtime applications: to allocate and name an object to be mapped into
4929        memory for potential sharing (*open*() or *shm_open*()), and to make the memory object
4930        available within the address space of a process (*mmap*()). To complete the interface, a
4931        mechanism to release the claim of a process on a shared memory object (*munmap*()) is also
4932        needed, as well as a mechanism for deleting the name of a sharable object that was
4933        previously created (*unlink*() or *shm_unlink*()).

4934        After a mapping has been established, an implementation should not have to provide
4935        services to maintain that mapping. All memory writes into that area will appear immediately
4936        in the memory mapping of that region by any other processes.

4937        Thus, requirements include:

4938        — Support creation of sharable memory objects and the mapping of these objects into the
4939             address space of a process.

4940        — Sharable memory objects should be accessed by global names accessible from all
4941             processes.

4942        — Support the mapping of specific sections of physical address space (such as a memory
4943             mapped device) into the address space of a process. This should not be done by the
4944             process specifying the actual address, but again by an implementation-defined global
4945             name (such as a special device name) dedicated to this purpose.

4946        — Support the mapping of discrete portions of these memory objects.

4947        — Support for minimum hardware configurations that contain no physical media on which
4948             to store shared memory contents permanently.

4949        — The ability to preallocate the entire shared memory region so that minimum hardware
4950             configurations without virtual memory support can guarantee contiguous space.

4951        — The maximizing of performance by not requiring functionality that would require
4952             implementation interaction above creating the shared memory area and returning the
4953             mapping.

4954        Note that the above requirements do not preclude:

4955        — The sharable memory object from being implemented using actual files on an actual file
4956             system.

4957        — The global name that is accessible from all processes being restricted to a file system area
4958             that is dedicated to handling shared memory.

4959        — An implementation not providing implementation-defined global names for the purpose
4960             of physical address mapping.

4961    • Shared Memory Objects Usage

4962        If the Shared Memory Objects option is supported, a shared memory object may be created,
4963        or opened if it already exists, with the *shm_open*() function. If the shared memory object is
4964        created, it has a length of zero. The *ftruncate*() function can be used to set the size of the

4965      shared memory object after creation. The *shm_unlink*( ) function removes the name for a
4966      shared memory object created by *shm_open*( ).

4967      • Shared Memory Overview

4968      The shared memory facility defined by IEEE Std 1003.1-2001 usually results in memory
4969      locations being added to the address space of the process. The implementation returns the
4970      address of the new space to the application by means of a pointer. This works well in
4971      languages like C. However, in languages without pointer types it will not work. In the
4972      bindings for such a language, either a special COMMON section will need to be defined
4973      (which is unlikely), or the binding will have to allow existing structures to be mapped. The
4974      implementation will likely have to place restrictions on the size and alignment of such
4975      structures or will have to map a suitable region of the address space of the process into the
4976      memory object, and thus into other processes. These are issues for that particular language
4977      binding. For IEEE Std 1003.1-2001, however, the practice will not be forbidden, merely
4978      undefined.

4979      Two potentially different name spaces are used for naming objects that may be mapped into
4980      process address spaces. When the Memory Mapped Files option is supported, files may be
4981      accessed via *open*( ). When the Shared Memory Objects option is supported, sharable
4982      memory objects that might not be files may be accessed via the *shm_open*( ) function. These
4983      options are not mutually-exclusive.

4984      Some implementations supporting the Shared Memory Objects option may choose to
4985      implement the shared memory object name space as part of the file system name space.
4986      There are several reasons for this:

4987      — It allows applications to prevent name conflicts by use of the directory structure.

4988      — It uses an existing mechanism for accessing global objects and prevents the creation of a
4989          new mechanism for naming global objects.

4990      In such implementations, memory objects can be implemented using regular files, if that is
4991      what the implementation chooses. The *shm_open*( ) function can be implemented as an *open*( )
4992      call in a fixed directory followed by a call to *fcntl*( ) to set FD_CLOEXEC. The *shm_unlink*( )
4993      function can be implemented as an *unlink*( ) call.

4994      On the other hand, it is also expected that small embedded systems that support the Shared
4995      Memory Objects option may wish to implement shared memory without having any file
4996      systems present. In this case, the implementations may choose to use a simple string valued
4997      name space for shared memory regions. The *shm_open*( ) function permits either type of
4998      implementation.

4999      Some implementations have hardware that supports protection of mapped data from certain
5000      classes of access and some do not. Systems that supply this functionality can support the
5001      Memory Protection option.

5002      Some implementations restrict size, alignment, and protections to be on *page*-size
5003      boundaries. If an implementation has no restrictions on size or alignment, it may specify a 1-
5004      byte page size. Applications on implementations that do support larger pages must be
5005      cognizant of the page size since this is the alignment and protection boundary.

5006      Simple embedded implementations may have a 1-byte page size and only support the Shared
5007      Memory Objects option. This provides simple shared memory between processes without
5008      requiring mapping hardware.

5009      IEEE Std 1003.1-2001 specifically allows a memory object to remain referenced after a close
5010      because that is existing practice for the *mmap*( ) function.

5011    **Typed Memory Functions**

5012    Implementations may support the Typed Memory Objects option without supporting either the
5013    Shared Memory option or the Memory Mapped Files option. Typed memory objects are pools of
5014    specialized storage, different from the main memory resource normally used by a processor to
5015    hold code and data, that can be mapped into the address space of one or more processes.

5016        • Model

5017        Realtime systems conforming to one of the POSIX.13 realtime profiles are expected (and
5018        desired) to be supported on systems with more than one type or pool of memory (for
5019        example, SRAM, DRAM, ROM, EPROM, EEPROM), where each type or pool of memory may
5020        be accessible by one or more processors via one or more busses (ports). Memory mapped
5021        files, shared memory objects, and the language-specific storage allocation operators (*malloc*( )
5022        for the ISO C standard, *new* for ISO Ada) fail to provide application program interfaces
5023        versatile enough to allow applications to control their utilization of such diverse memory
5024        resources. The typed memory interfaces *posix_typed_mem_open*( ), *posix_mem_offset*( ),
5025        *posix_typed_mem_get_info*( ), *mmap*( ), and *munmap*( ) defined herein support the model of
5026        typed memory described below.

5027        For purposes of this model, a system comprises several processors (for example, $P_1$ and $P_2$),
5028        several physical memory pools (for example, $M_1$, $M_2$, $M_{2a}$, $M_{2b}$, $M_3$, $M_4$, and $M_5$), and several
5029        busses or ''ports'' (for example, $B_1$, $B_2$, $B_3$, and $B_4$) interconnecting the various processors and
5030        memory pools in some system-specific way. Notice that some memory pools may be
5031        contained in others (for example, $M_{2a}$ and $M_{2b}$ are contained in $M_2$).

5032        Figure B-1 shows an example of such a model. In a system like this, an application should be
5033        able to perform the following operations:

5034



            * All addresses in pool $M_2$ (comprising pools $M_{2a}$ and $M_{2b}$) accessible via port $B_1$.
              Addresses in pool $M_{2b}$ are also accessible via port $B_2$.
              Addresses in pool $M_{2a}$ are *not* accessible via port $B_2$.

5035                        **Figure B-1**  Example of a System with Typed Memory

5036        — Typed Memory Allocation

5037        An application should be able to allocate memory dynamically from the desired pool
5038        using the desired bus, and map it into a process' address space. For example, processor $P_1$
5039        can allocate some portion of memory pool $M_1$ through port $B_1$, treating all unmapped
5040        subareas of $M_1$ as a heap-storage resource from which memory may be allocated. This
5041        portion of memory is mapped into the process' address space, and subsequently
5042        deallocated when unmapped from all processes.

5043  — Using the Same Storage Region from Different Busses

5044  An application process with a mapped region of storage that is accessed from one bus
5045  should be able to map that same storage area at another address (subject to page size
5046  restrictions detailed in *mmap*( )), to allow it to be accessed from another bus. For example,
5047  processor $P_1$ may wish to access the same region of memory pool $M_{2b}$ both through ports
5048  $B_1$ and $B_2$.

5049  — Sharing Typed Memory Regions

5050  Several application processes running on the same or different processors may wish to
5051  share a particular region of a typed memory pool. Each process or processor may wish to
5052  access this region through different busses. For example, processor $P_1$ may want to share
5053  a region of memory pool $M_4$ with processor $P_2$, and they may be required to use busses $B_2$
5054  and $B_3$, respectively, to minimize bus contention. A problem arises here when a process
5055  allocates and maps a portion of fragmented memory and then wants to share this region
5056  of memory with another process, either in the same processor or different processors. The
5057  solution adopted is to allow the first process to find out the memory map (offsets and
5058  lengths) of all the different fragments of memory that were mapped into its address
5059  space, by repeatedly calling *posix_mem_offset*( ). Then, this process can pass the offsets
5060  and lengths obtained to the second process, which can then map the same memory
5061  fragments into its address space.

5062  — Contiguous Allocation

5063  The problem of finding the memory map of the different fragments of the memory pool
5064  that were mapped into logically contiguous addresses of a given process can be solved by
5065  requesting contiguous allocation. For example, a process in $P_1$ can allocate 10 Kbytes of
5066  physically contiguous memory from $M_3$-$B_1$, and obtain the offset (within pool $M_3$) of this
5067  block of memory. Then, it can pass this offset (and the length) to a process in $P_2$ using
5068  some interprocess communication mechanism. The second process can map the same
5069  block of memory by using the offset transferred and specifying $M_3$-$B_2$.

5070  — Unallocated Mapping

5071  Any subarea of a memory pool that is mapped to a process, either as the result of an
5072  allocation request or an explicit mapping, is normally unavailable for allocation. Special
5073  processes such as debuggers, however, may need to map large areas of a typed memory
5074  pool, yet leave those areas available for allocation.

5075  Typed memory allocation and mapping has to coexist with storage allocation operators like
5076  *malloc*( ), but systems are free to choose how to implement this coexistence. For example, it
5077  may be system configuration-dependent if all available system memory is made part of one
5078  of the typed memory pools or if some part will be restricted to conventional allocation
5079  operators. Equally system configuration-dependent may be the availability of operators like
5080  *malloc*( ) to allocate storage from certain typed memory pools. It is not excluded to configure
5081  a system such that a given named pool, $P_1$, is in turn split into non-overlapping named
5082  subpools. For example, $M_1$-$B_1$, $M_2$-$B_1$, and $M_3$-$B_1$ could also be accessed as one common pool
5083  $M_{123}$-$B_1$. A call to *malloc*( ) on $P_1$ could work on such a larger pool while full optimization of
5084  memory usage by $P_1$ would require typed memory allocation at the subpool level.

5085  • Existing Practice

5086  OS-9 provides for the naming (numbering) and prioritization of memory types by a system
5087  administrator. It then provides APIs to request memory allocation of typed (colored)
5088  memory by number, and to generate a bus address from a mapped memory address
5089  (translate). When requesting colored memory, the user can specify type 0 to signify allocation

5090          from the first available type in priority order.

5091          HP-RT presents interfaces to map different kinds of storage regions that are visible through a
5092          VME bus, although it does not provide allocation operations. It also provides functions to
5093          perform address translation between VME addresses and virtual addresses. It represents a
5094          VME-bus unique solution to the general problem.

5095          The PSOS approach is similar (that is, based on a pre-established mapping of bus address
5096          ranges to specific memories) with a concept of segments and regions (regions dynamically
5097          allocated from a heap which is a special segment). Therefore, PSOS does not fully address the
5098          general allocation problem either. PSOS does not have a ''process''-based model, but more of
5099          a ''thread''-only-based model of multi-tasking. So mapping to a process address space is not
5100          an issue.

5101          QNX uses the System V approach of opening specially named devices (shared memory
5102          segments) and using *mmap*( ) to then gain access from the process. They do not address
5103          allocation directly, but once typed shared memory can be mapped, an ''allocation manager''
5104          process could be written to handle requests for allocation.

5105          The System V approach also included allocation, implemented by opening yet other special
5106          ''devices'' which allocate, rather than appearing as a whole memory object.

5107          The Orkid realtime kernel interface definition has operations to manage memory ''regions''
5108          and ''pools'', which are areas of memory that may reflect the differing physical nature of the
5109          memory. Operations to allocate memory from these regions and pools are also provided.

5110      • Requirements

5111          Existing practice in SVID-derived UNIX systems relies on functionality similar to *mmap*( )
5112          and its related interfaces to achieve mapping and allocation of typed memory. However, the
5113          issue of sharing typed memory (allocated or mapped) and the complication of multiple ports
5114          are not addressed in any consistent way by existing UNIX system practice. Part of this
5115          functionality is existing practice in specialized realtime operating systems. In order to
5116          solidify the capabilities implied by the model above, the following requirements are imposed
5117          on the interface:

5118          — Identification of Typed Memory Pools and Ports

5119              All processes (running in all processors) in the system are able to identify a particular
5120              (system configured) typed memory pool accessed through a particular (system
5121              configured) port by a name. That name is a member of a name space common to all these
5122              processes, but need not be the same name space as that containing ordinary filenames.
5123              The association between memory pools/ports and corresponding names is typically
5124              established when the system is configured. The ''open'' operation for typed memory
5125              objects should be distinct from the *open*( ) function, for consistency with other similar
5126              services, but implementable on top of *open*( ). This implies that the handle for a typed
5127              memory object will be a file descriptor.

5128          — Allocation and Mapping of Typed Memory

5129              Once a typed memory object has been identified by a process, it is possible to both map
5130              user-selected subareas of that object into process address space and to map system-
5131              selected (that is, dynamically allocated) subareas of that object, with user-specified
5132              length, into process address space. It is also possible to determine the maximum length of
5133              memory allocation that may be requested from a given typed memory object.

5134 — Sharing Typed Memory

5135 Two or more processes are able to share portions of typed memory, either user-selected or
5136 dynamically allocated. This requirement applies also to dynamically allocated regions of
5137 memory that are composed of several non-contiguous pieces.

5138 — Contiguous Allocation

5139 For dynamic allocation, it is the user's option whether the system is required to allocate a
5140 contiguous subarea within the typed memory object, or whether it is permitted to allocate
5141 discontiguous fragments which appear contiguous in the process mapping. Contiguous
5142 allocation simplifies the process of sharing allocated typed memory, while discontiguous
5143 allocation allows for potentially better recovery of deallocated typed memory.

5144 — Accessing Typed Memory Through Different Ports

5145 Once a subarea of a typed memory object has been mapped, it is possible to determine the
5146 location and length corresponding to a user-selected portion of that object within the
5147 memory pool. This location and length can then be used to remap that portion of memory
5148 for access from another port. If the referenced portion of typed memory was allocated
5149 discontiguously, the length thus determined may be shorter than anticipated, and the
5150 user code must adapt to the value returned.

5151 — Deallocation

5152 When a previously mapped subarea of typed memory is no longer mapped by any
5153 process in the system—as a result of a call or calls to *munmap*()—that subarea becomes
5154 potentially reusable for dynamic allocation; actual reuse of the subarea is a function of the
5155 dynamic typed memory allocation policy.

5156 — Unallocated Mapping

5157 It must be possible to map user-selected subareas of a typed memory object without
5158 marking that subarea as unavailable for allocation. This option is not the default behavior,
5159 and requires appropriate privilege.

5160 • Scenario

5161 The following scenario will serve to clarify the use of the typed memory interfaces.

5162 Process A running on $P_1$ (see Figure B-1 (on page 122)) wants to allocate some memory from
5163 memory pool $M_2$, and it wants to share this portion of memory with process B running on $P_2$.
5164 Since $P_2$ only has access to the lower part of $M_2$, both processes will use the memory pool
5165 named $M_{2b}$ which is the part of $M_2$ that is accessible both from $P_1$ and $P_2$. The operations that
5166 both processes need to perform are shown below:

5167 — Allocating Typed Memory

5168 Process A calls *posix_typed_mem_open*() with the name **/typed.m2b-b1** and a *tflag* of
5169 POSIX_TYPED_MEM_ALLOCATE to get a file descriptor usable for allocating from pool
5170 $M_{2b}$ accessed through port $B_1$. It then calls *mmap*() with this file descriptor requesting a
5171 length of 4 096 bytes. The system allocates two discontiguous blocks of sizes 1 024 and
5172 3 072 bytes within $M_{2b}$. The *mmap*() function returns a pointer to a 4 096-byte array in
5173 process A's logical address space, mapping the allocated blocks contiguously. Process A
5174 can then utilize the array, and store data in it.

5175 — Determining the Location of the Allocated Blocks

5176 Process A can determine the lengths and offsets (relative to $M_{2b}$) of the two blocks
5177 allocated, by using the following procedure: First, process A calls *posix_mem_offset*() with

5178   the address of the first element of the array and length 4 096. Upon return, the offset and
5179   length (1 024 bytes) of the first block are returned. A second call to *posix_mem_offset*( ) is
5180   then made using the address of the first element of the array plus 1 024 (the length of the
5181   first block), and a new length of 4 096–1 024. If there were more fragments allocated, this
5182   procedure could have been continued within a loop until the offsets and lengths of all the
5183   blocks were obtained. Notice that this relatively complex procedure can be avoided if
5184   contiguous allocation is requested (by opening the typed memory object with the *tflag*
5185   POSIX_TYPED_MEM_ALLOCATE_CONTIG).

5186   — Sharing Data Across Processes

5187   Process A passes the two offset values and lengths obtained from the *posix_mem_offset*( )
5188   calls to process B running on P$_2$, via some form of interprocess communication. Process B
5189   can gain access to process A's data by calling *posix_typed_mem_open*( ) with the name
5190   /**typed.m2b-b2** and a *tflag* of zero, then using two *mmap*( ) calls on the resulting file
5191   descriptor to map the two subareas of that typed memory object to its own address space.

5192   • Rationale for no *mem_alloc*( ) and *mem_free*( )

5193   The standard developers had originally proposed a pair of new flags to *mmap*( ) which, when
5194   applied to a typed memory object descriptor, would cause *mmap*( ) to allocate dynamically
5195   from an unallocated and unmapped area of the typed memory object. Deallocation was
5196   similarly accomplished through the use of *munmap*( ). This was rejected by the ballot group
5197   because it excessively complicated the (already rather complex) *mmap*( ) interface and
5198   introduced semantics useful only for typed memory, to a function which must also map
5199   shared memory and files. They felt that a memory allocator should be built on top of *mmap*( )
5200   instead of being incorporated within the same interface, much as the ISO C standard libraries
5201   build *malloc*( ) on top of the virtual memory mapping functions *brk*( ) and *sbrk*( ). This would
5202   eliminate the complicated semantics involved with unmapping only part of an allocated
5203   block of typed memory.

5204   To attempt to achieve ballot group consensus, typed memory allocation and deallocation was
5205   first migrated from *mmap*( ) and *munmap*( ) to a pair of complementary functions modeled on
5206   the ISO C standard *malloc*( ) and *free*( ). The *mem_alloc*( ) function specified explicitly the
5207   typed memory object (typed memory pool ∕ access port) from which allocation takes place,
5208   unlike *malloc*( ) where the memory pool and port are unspecified. The *mem_free*( ) function
5209   handled deallocation. These new semantics still met all of the requirements detailed above
5210   without modifying the behavior of *mmap*( ) except to allow it to map specified areas of typed
5211   memory objects. An implementation would have been free to implement *mem_alloc*( ) and
5212   *mem_free*( ) over *mmap*( ), through *mmap*( ), or independently but cooperating with *mmap*( ).

5213   The ballot group was queried to see if this was an acceptable alternative, and while there was
5214   some agreement that it achieved the goal of removing the complicated semantics of
5215   allocation from the *mmap*( ) interface, several balloters realized that it just created two
5216   additional functions that behaved, in great part, like *mmap*( ). These balloters proposed an
5217   alternative which has been implemented here in place of a separate *mem_alloc*( ) and
5218   *mem_free*( ). This alternative is based on four specific suggestions:

5219   1. The *posix_typed_mem_open*( ) function should provide a flag which specifies ''allocate
5220      on *mmap*( )'' (otherwise, *mmap*( ) just maps the underlying object). This allows things
5221      roughly similar to /**dev**∕**zero** *versus* /**dev**∕**swap**. Two such flags have been implemented,
5222      one of which forces contiguous allocation.

5223   2. The *posix_mem_offset*( ) function is acceptable because it can be applied usefully to
5224      mapped objects in general. It should return the file descriptor of the underlying object.

5225        3.   The *mem_get_info*( ) function in an earlier draft should be renamed
5226             *posix_typed_mem_get_info*( ) because it is not generally applicable to memory objects. It
5227             should probably return the file descriptor's allocation attribute. The renaming of the
5228             function has been implemented, but having it return a piece of information which is
5229             readily known by an application without this function has been rejected. Its whole
5230             purpose is to query the typed memory object for attributes that are not user-specified,
5231             but determined by the implementation.

5232        4.   There should be no separate *mem_alloc*( ) or *mem_free*( ) functions. Instead, using
5233             *mmap*( ) on a typed memory object opened with an ''allocate on *mmap*( )'' flag should be
5234             used to force allocation. These are precisely the semantics defined in the current draft.

5235     •  Rationale for no Typed Memory Access Management

5236    The working group had originally defined an additional interface (and an additional kind of
5237    object: typed memory master) to establish and dissolve mappings to typed memory on
5238    behalf of devices or processors which were independent of the operating system and had no
5239    inherent capability to directly establish mappings on their own. This was to have provided
5240    functionality similar to device driver interfaces such as *physio*( ) and their underlying bus-
5241    specific interfaces (for example, *mballoc*( )) which serve to set up and break down DMA
5242    pathways, and derive mapped addresses for use by hardware devices and processor cards.

5243    The ballot group felt that this was beyond the scope of POSIX.1 and its amendments.
5244    Furthermore, the removal of interrupt handling interfaces from a preceding amendment (the
5245    IEEE Std 1003.1d-1999) during its balloting process renders these typed memory access
5246    management interfaces an incomplete solution to portable device management from a user
5247    process; it would be possible to initiate a device transfer to/from typed memory, but
5248    impossible to handle the transfer-complete interrupt in a portable way.

5249    To achieve ballot group consensus, all references to typed memory access management
5250    capabilities were removed. The concept of portable interfaces from a device driver to both
5251    operating system and hardware is being addressed by the Uniform Driver Interface (UDI)
5252    industry forum, with formal standardization deferred until proof of concept and industry-
5253    wide acceptance and implementation.

5254  *B.2.8.4   Process Scheduling*

5255    IEEE PASC Interpretation 1003.1 #96 has been applied, adding the *pthread_setschedprio*( )
5256    function. This was added since previously there was no way for a thread to lower its own
5257    priority without going to the tail of the threads list for its new priority. This capability is
5258    necessary to bound the duration of priority inversion encountered by a thread.

5259    The following portion of the rationale presents models, requirements, and standardization issues
5260    relevant to process scheduling; see also Section B.2.9.4 (on page 167).

5261    In an operating system supporting multiple concurrent processes, the system determines the
5262    order in which processes execute to meet implementation-defined goals. For time-sharing
5263    systems, the goal is to enhance system throughput and promote fairness; the application is
5264    provided with little or no control over this sequencing function. While this is acceptable and
5265    desirable behavior in a time-sharing system, it is inappropriate in a realtime system; realtime
5266    applications must specifically control the execution sequence of their concurrent processes in
5267    order to meet externally defined response requirements.

5268    In IEEE Std 1003.1-2001, the control over process sequencing is provided using a concept of
5269    scheduling policies. These policies, described in detail in this section, define the behavior of the
5270    system whenever processor resources are to be allocated to competing processes. Only the
5271    behavior of the policy is defined; conforming implementations are free to use any mechanism

5272  desired to achieve the described behavior.

5273  • Models

5274  In an operating system supporting multiple concurrent processes, the system determines the
5275  order in which processes execute and might force long-running processes to yield to other
5276  processes at certain intervals. Typically, the scheduling code is executed whenever an event
5277  occurs that might alter the process to be executed next.

5278  The simplest scheduling strategy is a ''first-in, first-out'' (FIFO) dispatcher. Whenever a
5279  process becomes runnable, it is placed on the end of a ready list. The process at the front of
5280  the ready list is executed until it exits or becomes blocked, at which point it is removed from
5281  the list. This scheduling technique is also known as ''run-to-completion'' or ''run-to-block''.

5282  A natural extension to this scheduling technique is the assignment of a ''non-migrating
5283  priority'' to each process. This policy differs from strict FIFO scheduling in only one respect:
5284  whenever a process becomes runnable, it is placed at the end of the list of processes runnable
5285  at that priority level. When selecting a process to run, the system always selects the first
5286  process from the highest priority queue with a runnable process. Thus, when a process
5287  becomes unblocked, it will preempt a running process of lower priority without otherwise
5288  altering the ready list. Further, if a process elects to alter its priority, it is removed from the
5289  ready list and reinserted, using its new priority, according to the policy above.

5290  While the above policy might be considered unfriendly in a time-sharing environment in
5291  which multiple users require more balanced resource allocation, it could be ideal in a
5292  realtime environment for several reasons. The most important of these is that it is
5293  deterministic: the highest-priority process is always run and, among processes of equal
5294  priority, the process that has been runnable for the longest time is executed first. Because of
5295  this determinism, cooperating processes can implement more complex scheduling simply by
5296  altering their priority. For instance, if processes at a single priority were to reschedule
5297  themselves at fixed time intervals, a time-slice policy would result.

5298  In a dedicated operating system in which all processes are well-behaved realtime
5299  applications, non-migrating priority scheduling is sufficient. However, many existing
5300  implementations provide for more complex scheduling policies.

5301  IEEE Std 1003.1-2001 specifies a linear scheduling model. In this model, every process in the
5302  system has a priority. The system scheduler always dispatches a process that has the highest
5303  (generally the most time-critical) priority among all runnable processes in the system. As
5304  long as there is only one such process, the dispatching policy is trivial. When multiple
5305  processes of equal priority are eligible to run, they are ordered according to a strict run-to-
5306  completion (FIFO) policy.

5307  The priority is represented as a positive integer and is inherited from the parent process. For
5308  processes running under a fixed priority scheduling policy, the priority is never altered
5309  except by an explicit function call.

5310  It was determined arbitrarily that larger integers correspond to ''higher priorities''.

5311  Certain implementations might impose restrictions on the priority ranges to which processes
5312  can be assigned. There also can be restrictions on the set of policies to which processes can be
5313  set.

5314  • Requirements

5315  Realtime processes require that scheduling be fast and deterministic, and that it guarantees
5316  to preempt lower priority processes.

5317 Thus, given the linear scheduling model, realtime processes require that they be run at a
5318 priority that is higher than other processes. Within this framework, realtime processes are
5319 free to yield execution resources to each other in a completely portable and implementation-
5320 defined manner.

5321 As there is a generally perceived requirement for processes at the same priority level to share
5322 processor resources more equitably, provisions are made by providing a scheduling policy
5323 (that is, SCHED_RR) intended to provide a timeslice-like facility.

5324 **Note:** The following topics assume that low numeric priority implies low scheduling criticality
5325 and *vice versa*.

5326 • Rationale for New Interface

5327 Realtime applications need to be able to determine when processes will run in relation to
5328 each other. It must be possible to guarantee that a critical process will run whenever it is
5329 runnable; that is, whenever it wants to for as long as it needs. SCHED_FIFO satisfies this
5330 requirement. Additionally, SCHED_RR was defined to meet a realtime requirement for a
5331 well-defined time-sharing policy for processes at the same priority.

5332 It would be possible to use the BSD *setpriority*( ) and *getpriority*( ) functions by redefining the
5333 meaning of the ''nice'' parameter according to the scheduling policy currently in use by the
5334 process. The System V *nice*( ) interface was felt to be undesirable for realtime because it
5335 specifies an adjustment to the ''nice'' value, rather than setting it to an explicit value.
5336 Realtime applications will usually want to set priority to an explicit value. Also, System V
5337 *nice*( ) does not allow for changing the priority of another process.

5338 With the POSIX.1b interfaces, the traditional ''nice'' value does not affect the SCHED_FIFO
5339 or SCHED_RR scheduling policies. If a ''nice'' value is supported, it is implementation-
5340 defined whether it affects the SCHED_OTHER policy.

5341 An important aspect of IEEE Std 1003.1-2001 is the explicit description of the queuing and
5342 preemption rules. It is critical, to achieve deterministic scheduling, that such rules be stated
5343 clearly in IEEE Std 1003.1-2001.

5344 IEEE Std 1003.1-2001 does not address the interaction between priority and swapping. The
5345 issues involved with swapping and virtual memory paging are extremely implementation-
5346 defined and would be nearly impossible to standardize at this point. The proposed
5347 scheduling paradigm, however, fully describes the scheduling behavior of runnable
5348 processes, of which one criterion is that the working set be resident in memory. Assuming
5349 the existence of a portable interface for locking portions of a process in memory, paging
5350 behavior need not affect the scheduling of realtime processes.

5351 IEEE Std 1003.1-2001 also does not address the priorities of ''system'' processes. In general,
5352 these processes should always execute in low-priority ranges to avoid conflict with other
5353 realtime processes. Implementations should document the priority ranges in which system
5354 processes run.

5355 The default scheduling policy is not defined. The effect of I/O interrupts and other system
5356 processing activities is not defined. The temporary lending of priority from one process to
5357 another (such as for the purposes of affecting freeing resources) by the system is not
5358 addressed. Preemption of resources is not addressed. Restrictions on the ability of a process
5359 to affect other processes beyond a certain level (influence levels) is not addressed.

5360 The rationale used to justify the simple time-quantum scheduler is that it is common practice
5361 to depend upon this type of scheduling to ensure ''fair'' distribution of processor resources
5362 among portions of the application that must interoperate in a serial fashion. Note that
5363 IEEE Std 1003.1-2001 is silent with respect to the setting of this time quantum, or whether it is

5364    a system-wide value or a per-process value, although it appears that the prevailing realtime
5365    practice is for it to be a system-wide value.

5366    In a system with *N* processes at a given priority, all processor-bound, in which the time
5367    quantum is equal for all processes at a specific priority level, the following assumptions are
5368    made of such a scheduling policy:

5369    1.  A time quantum *Q* exists and the current process will own control of the processor for
5370        at least a duration of *Q* and will have the processor for a duration of *Q*.

5371    2.  The *N*th process at that priority will control a processor within a duration of $(N-1) \times Q$.

5372    These assumptions are necessary to provide equal access to the processor and bounded
5373    response from the application.

5374    The assumptions hold for the described scheduling policy only if no system overhead, such
5375    as interrupt servicing, is present. If the interrupt servicing load is non-zero, then one of the
5376    two assumptions becomes fallacious, based upon how *Q* is measured by the system.

5377    If *Q* is measured by clock time, then the assumption that the process obtains a duration *Q*
5378    processor time is false if interrupt overhead exists. Indeed, a scenario can be constructed with
5379    *N* processes in which a single process undergoes complete processor starvation if a
5380    peripheral device, such as an analog-to-digital converter, generates significant interrupt
5381    activity periodically with a period of $N \times Q$.

5382    If *Q* is measured as actual processor time, then the assumption that the *N*th process runs in
5383    within the duration $(N-1) \times Q$ is false.

5384    It should be noted that SCHED_FIFO suffers from interrupt-based delay as well. However,
5385    for SCHED_FIFO, the implied response of the system is ''as soon as possible'', so that the
5386    interrupt load for this case is a vendor selection and not a compliance issue.

5387    With this in mind, it is necessary either to complete the definition by including bounds on the
5388    interrupt load, or to modify the assumptions that can be made about the scheduling policy.

5389    Since the motivation of inclusion of the policy is common usage, and since current
5390    applications do not enjoy the luxury of bounded interrupt load, item (2) above is sufficient to
5391    express existing application needs and is less restrictive in the standard definition. No
5392    difference in interface is necessary.

5393    In an implementation in which the time quantum is equal for all processes at a specific
5394    priority, our assumptions can then be restated as:

5395    — A time quantum *Q* exists, and a processor-bound process will be rescheduled after a
5396        duration of, at most, *Q*. Time quantum *Q* may be defined in either wall clock time or
5397        execution time.

5398    — In general, the *N*th process of a priority level should wait no longer than $(N-1) \times Q$ time
5399        to execute, assuming no processes exist at higher priority levels.

5400    — No process should wait indefinitely.

5401    For implementations supporting per-process time quanta, these assumptions can be readily
5402    extended.

**Sporadic Server Scheduling Policy**

The sporadic server is a mechanism defined for scheduling aperiodic activities in time-critical realtime systems. This mechanism reserves a certain bounded amount of execution capacity for processing aperiodic events at a high priority level. Any aperiodic events that cannot be processed within the bounded amount of execution capacity are executed in the background at a low priority level. Thus, a certain amount of execution capacity can be guaranteed to be available for processing periodic tasks, even under burst conditions in the arrival of aperiodic processing requests (that is, a large number of requests in a short time interval). The sporadic server also simplifies the schedulability analysis of the realtime system, because it allows aperiodic processes or threads to be treated as if they were periodic. The sporadic server was first described by Sprunt, et al.

The key concept of the sporadic server is to provide and limit a certain amount of computation capacity for processing aperiodic events at their assigned normal priority, during a time interval called the ''replenishment period''. Once the entity controlled by the sporadic server mechanism is initialized with its period and execution-time budget attributes, it preserves its execution capacity until an aperiodic request arrives. The request will be serviced (if there are no higher priority activities pending) as long as there is execution capacity left. If the request is completed, the actual execution time used to service it is subtracted from the capacity, and a replenishment of this amount of execution time is scheduled to happen one replenishment period after the arrival of the aperiodic request. If the request is not completed, because there is no execution capacity left, then the aperiodic process or thread is assigned a lower background priority. For each portion of consumed execution capacity the execution time used is replenished after one replenishment period. At the time of replenishment, if the sporadic server was executing at a background priority level, its priority is elevated to the normal level. Other similar replenishment policies have been defined, but the one presented here represents a compromise between efficiency and implementation complexity.

The interface that appears in this section defines a new scheduling policy for threads and processes that behaves according to the rules of the sporadic server mechanism. Scheduling attributes are defined and functions are provided to allow the user to set and get the parameters that control the scheduling behavior of this mechanism, namely the normal and low priority, the replenishment period, the maximum number of pending replenishment operations, and the initial execution-time budget.

- Scheduling Aperiodic Activities

  Virtually all realtime applications are required to process aperiodic activities. In many cases, there are tight timing constraints that the response to the aperiodic events must meet. Usual timing requirements imposed on the response to these events are:

  — The effects of an aperiodic activity on the response time of lower priority activities must be controllable and predictable.

  — The system must provide the fastest possible response time to aperiodic events.

  — It must be possible to take advantage of all the available processing bandwidth not needed by time-critical activities to enhance average-case response times to aperiodic events.

  Traditional methods for scheduling aperiodic activities are background processing, polling tasks, and direct event execution:

  — Background processing consists of assigning a very low priority to the processing of aperiodic events. It utilizes all the available bandwidth in the system that has not been consumed by higher priority threads. However, it is very difficult, or impossible, to meet

5450    requirements on average-case response time, because the aperiodic entity has to wait for
5451    the execution of all other entities which have higher priority.

5452    — Polling consists of creating a periodic process or thread for servicing aperiodic requests.
5453    At regular intervals, the polling entity is started and its services accumulated pending
5454    aperiodic requests. If no aperiodic requests are pending, the polling entity suspends itself
5455    until its next period. Polling allows the aperiodic requests to be processed at a higher
5456    priority level. However, worst and average-case response times of polling entities are a
5457    direct function of the polling period, and there is execution overhead for each polling
5458    period, even if no event has arrived. If the deadline of the aperiodic activity is short
5459    compared to the inter-arrival time, the polling frequency must be increased to guarantee
5460    meeting the deadline. For this case, the increase in frequency can dramatically reduce the
5461    efficiency of the system and, therefore, its capacity to meet all deadlines. Yet, polling
5462    represents a good way to handle a large class of practical problems because it preserves
5463    system predictability, and because the amortized overhead drops as load increases.

5464    — Direct event execution consists of executing the aperiodic events at a high fixed-priority
5465    level. Typically, the aperiodic event is processed by an interrupt service routine as soon as
5466    it arrives. This technique provides predictable response times for aperiodic events, but
5467    makes the response times of all lower priority activities completely unpredictable under
5468    burst arrival conditions. Therefore, if the density of aperiodic event arrivals is
5469    unbounded, it may be a dangerous technique for time-critical systems. Yet, for those cases
5470    in which the physics of the system imposes a bound on the event arrival rate, it is
5471    probably the most efficient technique.

5472    — The sporadic server scheduling algorithm combines the predictability of the polling
5473    approach with the short response times of the direct event execution. Thus, it allows
5474    systems to meet an important class of application requirements that cannot be met by
5475    using the traditional approaches. Multiple sporadic servers with different attributes can
5476    be applied to the scheduling of multiple classes of aperiodic events, each with different
5477    kinds of timing requirements, such as individual deadlines, average response times, and
5478    so on. It also has many other interesting applications for realtime, such as scheduling
5479    producer/consumer tasks in time-critical systems, limiting the effects of faults on the
5480    estimation of task execution-time requirements, and so on.

5481    • Existing Practice

5482    The sporadic server has been used in different kinds of applications, including military
5483    avionics, robot control systems, industrial automation systems, and so on. There are
5484    examples of many systems that cannot be successfully scheduled using the classic
5485    approaches, such as direct event execution, or polling, and are schedulable using a sporadic
5486    server scheduler. The sporadic server algorithm itself can successfully schedule all systems
5487    scheduled with direct event execution or polling.

5488    The sporadic server scheduling policy has been implemented as a commercial product in the
5489    run-time system of the Verdix Ada compiler. There are also many applications that have
5490    used a much less efficient application-level sporadic server. These realtime applications
5491    would benefit from a sporadic server scheduler implemented at the scheduler level.

5492    • Library-Level *versus* Kernel-Level Implementation

5493    The sporadic server interface described in this section requires the sporadic server policy to
5494    be implemented at the same level as the scheduler. This means that the process sporadic
5495    server must be implemented at the kernel level and the thread sporadic server policy
5496    implemented at the same level as the thread scheduler; that is, kernel or library level.

5497      In an earlier interface for the sporadic server, this mechanism was implementable at a
5498      different level than the scheduler. This feature allowed the implementor to choose between
5499      an efficient scheduler-level implementation, or a simpler user or library-level
5500      implementation. However, the working group considered that this interface made the use of
5501      sporadic servers more complex, and that library-level implementations would lack some of
5502      the important functionality of the sporadic server, namely the limitation of the actual
5503      execution time of aperiodic activities. The working group also felt that the interface
5504      described in this chapter does not preclude library-level implementations of threads intended
5505      to provide efficient low-overhead scheduling for those threads that are not scheduled under
5506      the sporadic server policy.

5507      • Range of Scheduling Priorities

5508      Each of the scheduling policies supported in IEEE Std 1003.1-2001 has an associated range of
5509      priorities. The priority ranges for each policy might or might not overlap with the priority
5510      ranges of other policies. For time-critical realtime applications it is usual for periodic and
5511      aperiodic activities to be scheduled together in the same processor. Periodic activities will
5512      usually be scheduled using the SCHED_FIFO scheduling policy, while aperiodic activities
5513      may be scheduled using SCHED_SPORADIC. Since the application developer will require
5514      complete control over the relative priorities of these activities in order to meet his timing
5515      requirements, it would be desirable for the priority ranges of SCHED_FIFO and
5516      SCHED_SPORADIC to overlap completely. Therefore, although IEEE Std 1003.1-2001 does
5517      not require any particular relationship between the different priority ranges, it is
5518      recommended that these two ranges should coincide.

5519      • Dynamically Setting the Sporadic Server Policy

5520      Several members of the working group requested that implementations should not be
5521      required to support dynamically setting the sporadic server scheduling policy for a thread.
5522      The reason is that this policy may have a high overhead for library-level implementations of
5523      threads, and if threads are allowed to dynamically set this policy, this overhead can be
5524      experienced even if the thread does not use that policy. By disallowing the dynamic setting of
5525      the sporadic server scheduling policy, these implementations can accomplish efficient
5526      scheduling for threads using other policies. If a strictly conforming application needs to use
5527      the sporadic server policy, and is therefore willing to pay the overhead, it must set this policy
5528      at the time of thread creation.

5529      • Limitation of the Number of Pending Replenishments

5530      The number of simultaneously pending replenishment operations must be limited for each
5531      sporadic server for two reasons: an unlimited number of replenishment operations would
5532      need an unlimited number of system resources to store all the pending replenishment
5533      operations; on the other hand, in some implementations each replenishment operation will
5534      represent a source of priority inversion (just for the duration of the replenishment operation)
5535      and thus, the maximum amount of replenishments must be bounded to guarantee bounded
5536      response times. The way in which the number of replenishments is bounded is by lowering
5537      the priority of the sporadic server to *sched_ss_low_priority* when the number of pending
5538      replenishments has reached its limit. In this way, no new replenishments are scheduled until
5539      the number of pending replenishments decreases.

5540      In the sporadic server scheduling policy defined in IEEE Std 1003.1-2001, the application can
5541      specify the maximum number of pending replenishment operations for a single sporadic
5542      server, by setting the value of the *sched_ss_max_repl* scheduling parameter. This value must
5543      be between one and {SS_REPL_MAX}, which is a maximum limit imposed by the
5544      implementation. The limit {SS_REPL_MAX} must be greater than or equal to
5545      {_POSIX_SS_REPL_MAX}, which is defined to be four in IEEE Std 1003.1-2001. The minimum

5546    limit of four was chosen so that an application can at least guarantee that four different
5547    aperiodic events can be processed during each interval of length equal to the replenishment
5548    period.

5549 *B.2.8.5    Clocks and Timers*

5550        • Clocks

5551    IEEE Std 1003.1-2001 and the ISO C standard both define functions for obtaining system time.
5552    Implicit behind these functions is a mechanism for measuring passage of time. This
5553    specification makes this mechanism explicit and calls it a clock. The CLOCK_REALTIME
5554    clock required by IEEE Std 1003.1-2001 is a higher resolution version of the clock that
5555    maintains POSIX.1 system time. This is a ''system-wide'' clock, in that it is visible to all
5556    processes and, were it possible for multiple processes to all read the clock at the same time,
5557    they would see the same value.

5558    An extensible interface was defined, with the ability for implementations to define additional
5559    clocks. This was done because of the observation that many realtime platforms support
5560    multiple clocks, and it was desired to fit this model within the standard interface. But
5561    implementation-defined clocks need not represent actual hardware devices, nor are they
5562    necessarily system-wide.

5563        • Timers

5564    Two timer types are required for a system to support realtime applications:

5565        1.  One-shot

5566    A one-shot timer is a timer that is armed with an initial expiration time, either relative
5567    to the current time or at an absolute time (based on some timing base, such as time in
5568    seconds and nanoseconds since the Epoch). The timer expires once and then is
5569    disarmed. With the specified facilities, this is accomplished by setting the *it_value*
5570    member of the *value* argument to the desired expiration time and the *it_interval* member
5571    to zero.

5572        2.  Periodic

5573    A periodic timer is a timer that is armed with an initial expiration time, again either
5574    relative or absolute, and a repetition interval. When the initial expiration occurs, the
5575    timer is reloaded with the repetition interval and continues counting. With the
5576    specified facilities, this is accomplished by setting the *it_value* member of the *value*
5577    argument to the desired initial expiration time and the *it_interval* member to the desired
5578    repetition interval.

5579    For both of these types of timers, the time of the initial timer expiration can be specified in
5580    two ways:

5581        1.  Relative (to the current time)

5582        2.  Absolute

5583        • Examples of Using Realtime Timers

5584    In the diagrams below, *S* indicates a program schedule, *R* shows a schedule method request,
5585    and *E* suggests an internal operating system event.

5586        — Periodic Timer: Data Logging

5587    During an experiment, it might be necessary to log realtime data periodically to an
5588    internal buffer or to a mass storage device. With a periodic scheduling method, a logging

5589          module can be started automatically at fixed time intervals to log the data.

5590          Program schedule is requested every 10 seconds.

5591              R           S           S           S           S           S
5592          ----+----+----+----+----+----+----+----+----+----+----+--->
5593              5    10   15   20   25   30   35   40   45   50   55

5594          [Time (in Seconds)]

5595          To achieve this type of scheduling using the specified facilities, one would allocate a per-
5596          process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
5597          a call to *timer_settime*( ) with the TIMER_ABSTIME flag reset, and with an initial
5598          expiration value and a repetition interval of 10 seconds.

5599      — One-shot Timer (Relative Time): Device Initialization

5600          In an emission test environment, large sample bags are used to capture the exhaust from
5601          a vehicle. The exhaust is purged from these bags before each and every test. With a one-
5602          shot timer, a module could initiate the purge function and then suspend itself for a
5603          predetermined period of time while the sample bags are prepared.

5604          Program schedule requested 20 seconds after call is issued.

5605              R                       S
5606          ----+----+----+----+----+----+----+----+----+----+----+--->
5607              5    10   15   20   25   30   35   40   45   50   55

5608          [Time (in Seconds)]

5609          To achieve this type of scheduling using the specified facilities, one would allocate a per-
5610          process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
5611          a call to *timer_settime*( ) with the TIMER_ABSTIME flag reset, and with an initial
5612          expiration value of 20 seconds and a repetition interval of zero.

5613          Note that if the program wishes merely to suspend itself for the specified interval, it
5614          could more easily use *nanosleep*( ).

5615      — One-shot Timer (Absolute Time): Data Transmission

5616          The results from an experiment are often moved to a different system within a network
5617          for postprocessing or archiving. With an absolute one-shot timer, a module that moves
5618          data from a test-cell computer to a host computer can be automatically scheduled on a
5619          daily basis.

5620          Program schedule requested for 2:30 a.m.

5621                      R                                               S
5622          -----+-----+-----+-----+-----+-----+-----+-----+-----+----->
5623              23:00 23:30 24:00 00:30 01:00 01:30 02:00 02:30 03:00

5624          [Time of Day]

5625          To achieve this type of scheduling using the specified facilities, a per-process timer would
5626          be allocated based on clock ID CLOCK_REALTIME. Then the timer would be armed via
5627          a call to *timer_settime*( ) with the TIMER_ABSTIME flag set, and an initial expiration value
5628          equal to 2:30 a.m. of the next day.

5629      — Periodic Timer (Relative Time): Signal Stabilization

5630          Some measurement devices, such as emission analyzers, do not respond instantaneously
5631          to an introduced sample. With a periodic timer with a relative initial expiration time, a

module that introduces a sample and records the average response could suspend itself for a predetermined period of time while the signal is stabilized and then sample at a fixed rate.

Program schedule requested 15 seconds after call is issued and every 2 seconds thereafter.

```
         R                 S S S S S S S S S S S S S S S S S S S S S
    ----+----+----+----+----+----+----+----+----+----+----+--->
         5   10   15   20   25   30   35   40   45   50   55
```

[Time (in Seconds)]

To achieve this type of scheduling using the specified facilities, one would allocate a per-process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via a call to *timer_settime*( ) with TIMER_ABSTIME flag reset, and with an initial expiration value of 15 seconds and a repetition interval of 2 seconds.

— Periodic Timer (Absolute Time): Work Shift-related Processing

Resource utilization data is useful when time to perform experiments is being scheduled at a facility. With a periodic timer with an absolute initial expiration time, a module can be scheduled at the beginning of a work shift to gather resource utilization data throughout the shift. This data can be used to allocate resources effectively to minimize bottlenecks and delays and maximize facility throughput.

Program schedule requested for 2:00 a.m. and every 15 minutes thereafter.

```
           R                                  S   S   S   S   S   S
    -----+-----+-----+-----+-----+-----+-----+-----+-----+----->
         23:00 23:30 24:00 00:30 01:00 01:30 02:00 02:30 03:00
```

[Time of Day]

To achieve this type of scheduling using the specified facilities, one would allocate a per-process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via a call to *timer_settime*( ) with TIMER_ABSTIME flag set, and with an initial expiration value equal to 2:00 a.m. and a repetition interval equal to 15 minutes.

- Relationship of Timers to Clocks

The relationship between clocks and timers armed with an absolute time is straightforward: a timer expiration signal is requested when the associated clock reaches or exceeds the specified time. The relationship between clocks and timers armed with a relative time (an interval) is less obvious, but not unintuitive. In this case, a timer expiration signal is requested when the specified interval, *as measured by the associated clock*, has passed. For the required CLOCK_REALTIME clock, this allows timer expiration signals to be requested at specified ''wall clock'' times (absolute), or when a specified interval of ''realtime'' has passed (relative). For an implementation-defined clock—say, a process virtual time clock—timer expirations could be requested when the process has used a specified total amount of virtual time (absolute), or when it has used a specified *additional* amount of virtual time (relative).

The interfaces also allow flexibility in the implementation of the functions. For example, an implementation could convert all absolute times to intervals by subtracting the clock value at the time of the call from the requested expiration time and ''counting down'' at the supported resolution. Or it could convert all relative times to absolute expiration time by adding in the clock value at the time of the call and comparing the clock value to the expiration time at the supported resolution. Or it might even choose to maintain absolute times as absolute and compare them to the clock value at the supported resolution for absolute timers, and maintain relative times as intervals and count them down at the

5678
5679        resolution supported for relative timers. The choice will be driven by efficiency considerations and the underlying hardware or software clock implementation.

5680     • Data Definitions for Clocks and Timers

5681        IEEE Std 1003.1-2001 uses a time representation capable of supporting nanosecond resolution
5682        timers for the following reasons:

5683        — To enable IEEE Std 1003.1-2001 to represent those computer systems already using
5684            nanosecond or submicrosecond resolution clocks.

5685        — To accommodate those per-process timers that might need nanoseconds to specify an
5686            absolute value of system-wide clocks, even though the resolution of the per-process timer
5687            may only be milliseconds, or *vice versa*.

5688        — Because the number of nanoseconds in a second can be represented in 32 bits.

5689        Time values are represented in the **timespec** structure. The *tv_sec* member is of type **time_t**
5690        so that this member is compatible with time values used by POSIX.1 functions and the ISO C
5691        standard. The *tv_nsec* member is a **signed long** in order to simplify and clarify code that
5692        decrements or finds differences of time values. Note that because 1 billion (number of
5693        nanoseconds per second) is less than half of the value representable by a signed 32-bit value,
5694        it is always possible to add two valid fractional seconds represented as integral nanoseconds
5695        without overflowing the signed 32-bit value.

5696        A maximum allowable resolution for the CLOCK_REALTIME clock of 20 ms (1/50 seconds)
5697        was chosen to allow line frequency clocks in European countries to be conforming. 60 Hz
5698        clocks in the U.S. will also be conforming, as will finer granularity clocks, although a Strictly
5699        Conforming Application cannot assume a granularity of less than 20 ms (1/50 seconds).

5700        The minimum allowable maximum time allowed for the CLOCK_REALTIME clock and the
5701        function *nanosleep*( ), and timers created with *clock_id*=CLOCK_REALTIME, is determined by
5702        the fact that the *tv_sec* member is of type **time_t**.

5703        IEEE Std 1003.1-2001 specifies that timer expirations must not be delivered early, and
5704        *nanosleep*( ) must not return early due to quantization error. IEEE Std 1003.1-2001 discusses
5705        the various implementations of *alarm*( ) in the rationale and states that implementations that
5706        do not allow alarm signals to occur early are the most appropriate, but refrained from
5707        mandating this behavior. Because of the importance of predictability to realtime applications,
5708        IEEE Std 1003.1-2001 takes a stronger stance.

5709        The developers of IEEE Std 1003.1-2001 considered using a time representation that differs
5710        from POSIX.1b in the second 32 bit of the 64-bit value. Whereas POSIX.1b defines this field
5711        as a fractional second in nanoseconds, the other methodology defines this as a binary fraction
5712        of one second, with the radix point assumed before the most significant bit.

5713        POSIX.1b is a software, source-level standard and most of the benefits of the alternate
5714        representation are enjoyed by hardware implementations of clocks and algorithms. It was
5715        felt that mandating this format for POSIX.1b clocks and timers would unnecessarily burden
5716        the application writer with writing, possibly non-portable, multiple precision arithmetic
5717        packages to perform conversion between binary fractions and integral units such as
5718        nanoseconds, milliseconds, and so on.

5719          **Rationale for the Monotonic Clock**

5720          For those applications that use time services to achieve realtime behavior, changing the value of
5721          the clock on which these services rely may cause erroneous timing behavior. For these
5722          applications, it is necessary to have a monotonic clock which cannot run backwards, and which
5723          has a maximum clock jump that is required to be documented by the implementation.
5724          Additionally, it is desirable (but not required by IEEE Std 1003.1-2001) that the monotonic clock
5725          increases its value uniformly. This clock should not be affected by changes to the system time;
5726          for example, to synchronize the clock with an external source or to account for leap seconds.
5727          Such changes would cause errors in the measurement of time intervals for those time services
5728          that use the absolute value of the clock.

5729          One could argue that by defining the behavior of time services when the value of a clock is
5730          changed, deterministic realtime behavior can be achieved. For example, one could specify that
5731          relative time services should be unaffected by changes in the value of a clock. However, there
5732          are time services that are based upon an absolute time, but that are essentially intended as
5733          relative time services. For example, *pthread_cond_timedwait*() uses an absolute time to allow it to
5734          wake up after the required interval despite spurious wakeups. Although sometimes the
5735          *pthread_cond_timedwait*() timeouts are absolute in nature, there are many occasions in which
5736          they are relative, and their absolute value is determined from the current time plus a relative
5737          time interval. In this latter case, if the clock changes while the thread is waiting, the wait interval
5738          will not be the expected length. If a *pthread_cond_timedwait*() function were created that would
5739          take a relative time, it would not solve the problem because to retain the intended ''deadline'' a
5740          thread would need to compensate for latency due to the spurious wakeup, and preemption
5741          between wakeup and the next wait.

5742          The solution is to create a new monotonic clock, whose value does not change except for the
5743          regular ticking of the clock, and use this clock for implementing the various relative timeouts
5744          that appear in the different POSIX interfaces, as well as allow *pthread_cond_timedwait*() to choose
5745          this new clock for its timeout. A new *clock_nanosleep*() function is created to allow an application
5746          to take advantage of this newly defined clock. Notice that the monotonic clock may be
5747          implemented using the same hardware clock as the system clock.

5748          Relative timeouts for *sigtimedwait*() and *aio_suspend*() have been redefined to use the monotonic
5749          clock, if present. The *alarm*() function has not been redefined, because the same effect but with
5750          better resolution can be achieved by creating a timer (for which the appropriate clock may be
5751          chosen).

5752          The *pthread_cond_timedwait*() function has been treated in a different way, compared to other
5753          functions with absolute timeouts, because it is used to wait for an event, and thus it may have a
5754          deadline, while the other timeouts are generally used as an error recovery mechanism, and for
5755          them the use of the monotonic clock is not so important. Since the desired timeout for the
5756          *pthread_cond_timedwait*() function may either be a relative interval or an absolute time of day
5757          deadline, a new initialization attribute has been created for condition variables to specify the
5758          clock that is used for measuring the timeout in a call to *pthread_cond_timedwait*(). In this way, if
5759          a relative timeout is desired, the monotonic clock will be used; if an absolute deadline is required
5760          instead, the CLOCK_REALTIME or another appropriate clock may be used. This capability has
5761          not been added to other functions with absolute timeouts because for those functions the
5762          expected use of the timeout is mostly to prevent errors, and not so often to meet precise
5763          deadlines. As a consequence, the complexity of adding this capability is not justified by its
5764          perceived application usage.

5765          The *nanosleep*() function has not been modified with the introduction of the monotonic clock.
5766          Instead, a new *clock_nanosleep*() function has been created, in which the desired clock may be
5767          specified in the function call.

5768        • History of Resolution Issues

5769        Due to the shift from relative to absolute timeouts in IEEE Std 1003.1d-1999, the amendments
5770        to the *sem_timedwait*( ), *pthread_mutex_timedlock*( ), *mq_timedreceive*( ), and *mq_timedsend*( )
5771        functions of that standard have been removed. Those amendments specified that
5772        CLOCK_MONOTONIC would be used for the (relative) timeouts if the Monotonic Clock
5773        option was supported.

5774        Having these functions continue to be tied solely to CLOCK_MONOTONIC would not
5775        work. Since the absolute value of a time value obtained from CLOCK_MONOTONIC is
5776        unspecified, under the absolute timeouts interface, applications would behave differently
5777        depending on whether the Monotonic Clock option was supported or not (because the
5778        absolute value of the clock would have different meanings in either case).

5779        Two options were considered:

5780        1.  Leave the current behavior unchanged, which specifies the CLOCK_REALTIME clock
5781            for these (absolute) timeouts, to allow portability of applications between
5782            implementations supporting or not the Monotonic Clock option.

5783        2.  Modify these functions in the way that *pthread_cond_timedwait*( ) was modified to allow
5784            a choice of clock, so that an application could use CLOCK_REALTIME when it is trying
5785            to achieve an absolute timeout and CLOCK_MONOTONIC when it is trying to achieve
5786            a relative timeout.

5787        It was decided that the features of CLOCK_MONOTONIC are not as critical to these
5788        functions as they are to *pthread_cond_timedwait*( ). The *pthread_cond_timedwait*( ) function is
5789        given a relative timeout; the timeout may represent a deadline for an event. When these
5790        functions are given relative timeouts, the timeouts are typically for error recovery purposes
5791        and need not be so precise.

5792        Therefore, it was decided that these functions should be tied to CLOCK_REALTIME and not
5793        complicated by being given a choice of clock.

5794    **Execution Time Monitoring**

5795        • Introduction

5796        The main goals of the execution time monitoring facilities defined in this chapter are to
5797        measure the execution time of processes and threads and to allow an application to establish
5798        CPU time limits for these entities.

5799        The analysis phase of time-critical realtime systems often relies on the measurement of
5800        execution times of individual threads or processes to determine whether the timing
5801        requirements will be met. Also, performance analysis techniques for soft deadline realtime
5802        systems rely heavily on the determination of these execution times. The execution time
5803        monitoring functions provide application developers with the ability to measure these
5804        execution times online and open the possibility of dynamic execution-time analysis and
5805        system reconfiguration, if required.

5806        The second goal of allowing an application to establish execution time limits for individual
5807        processes or threads and detecting when they overrun allows program robustness to be
5808        increased by enabling online checking of the execution times.

5809        If errors are detected—possibly because of erroneous program constructs, the existence of
5810        errors in the analysis phase, or a burst of event arrivals—online detection and recovery is
5811        possible in a portable way. This feature can be extremely important for many time-critical
5812        applications. Other applications require trapping CPU-time errors as a normal way to exit an

5813 algorithm; for instance, some realtime artificial intelligence applications trigger a number of
5814 independent inference processes of varying accuracy and speed, limit how long they can run,
5815 and pick the best answer available when time runs out. In many periodic systems, overrun
5816 processes are simply restarted in the next resource period, after necessary end-of-period
5817 actions have been taken. This allows algorithms that are inherently data-dependent to be
5818 made predictable.

5819 The interface that appears in this chapter defines a new type of clock, the CPU-time clock,
5820 which measures execution time. Each process or thread can invoke the clock and timer
5821 functions defined in POSIX.1 to use them. Functions are also provided to access the CPU-
5822 time clock of other processes or threads to enable remote monitoring of these clocks.
5823 Monitoring of threads of other processes is not supported, since these threads are not visible
5824 from outside of their own process with the interfaces defined in POSIX.1.

5825 • Execution Time Monitoring Interface

5826 The clock and timer interface defined in POSIX.1 historically only defined one clock, which
5827 measures wall-clock time. The requirements for measuring execution time of processes and
5828 threads, and setting limits to their execution time by detecting when they overrun, can be
5829 accomplished with that interface if a new kind of clock is defined. These new clocks measure
5830 execution time, and one is associated with each process and with each thread. The clock
5831 functions currently defined in POSIX.1 can be used to read and set these CPU-time clocks,
5832 and timers can be created using these clocks as their timing base. These timers can then be
5833 used to send a signal when some specified execution time has been exceeded. The CPU-time
5834 clocks of each process or thread can be accessed by using the symbols
5835 CLOCK_PROCESS_CPUTIME_ID or CLOCK_THREAD_CPUTIME_ID.

5836 The clock and timer interface defined in POSIX.1 and extended with the new kind of CPU-
5837 time clock would only allow processes or threads to access their own CPU-time clocks.
5838 However, many realtime systems require the possibility of monitoring the execution time of
5839 processes or threads from independent monitoring entities. In order to allow applications to
5840 construct independent monitoring entities that do not require cooperation from or
5841 modification of the monitored entities, two functions have been added: *clock_getcpuclockid*( ),
5842 for accessing CPU-time clocks of other processes, and *pthread_getcpuclockid*( ), for accessing
5843 CPU-time clocks of other threads. These functions return the clock identifier associated with
5844 the process or thread specified in the call. These clock IDs can then be used in the rest of the
5845 clock function calls.

5846 The clocks accessed through these functions could also be used as a timing base for the
5847 creation of timers, thereby allowing independent monitoring entities to limit the CPU time
5848 consumed by other entities. However, this possibility would imply additional complexity
5849 and overhead because of the need to maintain a timer queue for each process or thread, to
5850 store the different expiration times associated with timers created by different processes or
5851 threads. The working group decided this additional overhead was not justified by
5852 application requirements. Therefore, creation of timers attached to the CPU-time clocks of
5853 other processes or threads has been specified as implementation-defined.

5854 • Overhead Considerations

5855 The measurement of execution time may introduce additional overhead in the thread
5856 scheduling, because of the need to keep track of the time consumed by each of these entities.
5857 In library-level implementations of threads, the efficiency of scheduling could be somehow
5858 compromised because of the need to make a kernel call, at each context switch, to read the
5859 process CPU-time clock. Consequently, a thread creation attribute called *cpu-clock-*
5860 *requirement* was defined, to allow threads to disconnect their respective CPU-time clocks.
5861 However, the Ballot Group considered that this attribute itself introduced some overhead,

5862    and that in current implementations it was not worth the effort. Therefore, the attribute was
5863    deleted, and thus thread CPU-time clocks are required for all threads if the Thread CPU-Time
5864    Clocks option is supported.

5865    • Accuracy of CPU-Time Clocks

5866    The mechanism used to measure the execution time of processes and threads is specified in
5867    IEEE Std 1003.1-2001 as implementation-defined. The reason for this is that both the
5868    underlying hardware and the implementation architecture have a very strong influence on
5869    the accuracy achievable for measuring CPU time. For some implementations, the
5870    specification of strict accuracy requirements would represent very large overheads, or even
5871    the impossibility of being implemented.

5872    Since the mechanism for measuring execution time is implementation-defined, realtime
5873    applications will be able to take advantage of accurate implementations using a portable
5874    interface. Of course, strictly conforming applications cannot rely on any particular degree of
5875    accuracy, in the same way as they cannot rely on a very accurate measurement of wall clock
5876    time. There will always exist applications whose accuracy or efficiency requirements on the
5877    implementation are more rigid than the values defined in IEEE Std 1003.1-2001 or any other
5878    standard.

5879    In any case, there is a minimum set of characteristics that realtime applications would expect
5880    from most implementations. One such characteristic is that the sum of all the execution times
5881    of all the threads in a process equals the process execution time, when no CPU-time clocks
5882    are disabled. This need not always be the case because implementations may differ in how
5883    they account for time during context switches. Another characteristic is that the sum of the
5884    execution times of all processes in a system equals the number of processors, multiplied by
5885    the elapsed time, assuming that no processor is idle during that elapsed time. However, in
5886    some implementations it might not be possible to relate CPU time to elapsed time. For
5887    example, in a heterogeneous multi-processor system in which each processor runs at a
5888    different speed, an implementation may choose to define each ''second'' of CPU time to be a
5889    certain number of ''cycles'' that a CPU has executed.

5890    • Existing Practice

5891    Measuring and limiting the execution time of each concurrent activity are common features
5892    of most industrial implementations of realtime systems. Almost all critical realtime systems
5893    are currently built upon a cyclic executive. With this approach, a regular timer interrupt kicks
5894    off the next sequence of computations. It also checks that the current sequence has
5895    completed. If it has not, then some error recovery action can be undertaken (or at least an
5896    overrun is avoided). Current software engineering principles and the increasing complexity
5897    of software are driving application developers to implement these systems on multi-
5898    threaded or multi-process operating systems. Therefore, if a POSIX operating system is to be
5899    used for this type of application, then it must offer the same level of protection.

5900    Execution time clocks are also common in most UNIX implementations, although these
5901    clocks usually have requirements different from those of realtime applications. The POSIX.1
5902    *times*() function supports the measurement of the execution time of the calling process, and
5903    its terminated child processes. This execution time is measured in clock ticks and is supplied
5904    as two different values with the user and system execution times, respectively. BSD supports
5905    the function *getrusage*(), which allows the calling process to get information about the
5906    resources used by itself and/or all of its terminated child processes. The resource usage
5907    includes user and system CPU time. Some UNIX systems have options to specify high
5908    resolution (up to one microsecond) CPU-time clocks using the *times*() or the *getrusage*()
5909    functions.

5910   The *times*( ) and *getrusage*( ) interfaces do not meet important realtime requirements, such as
5911   the possibility of monitoring execution time from a different process or thread, or the
5912   possibility of detecting an execution time overrun. The latter requirement is supported in
5913   some UNIX implementations that are able to send a signal when the execution time of a
5914   process has exceeded some specified value. For example, BSD defines the functions
5915   *getitimer*( ) and *setitimer*( ), which can operate either on a realtime clock (wall-clock), or on
5916   virtual-time or profile-time clocks which measure CPU time in two different ways. These
5917   functions do not support access to the execution time of other processes.

5918   IBM's MVS operating system supports per-process and per-thread execution time clocks. It
5919   also supports limiting the execution time of a given process.

5920   Given all this existing practice, the working group considered that the POSIX.1 clocks and
5921   timers interface was appropriate to meet most of the requirements that realtime applications
5922   have for execution time clocks.  Functions were added to get the CPU time clock IDs, and to
5923   allow/disallow the thread CPU-time clocks (in order to preserve the efficiency of some
5924   implementations of threads).

5925    • Clock Constants

5926   The definition of the manifest constants CLOCK_PROCESS_CPUTIME_ID and
5927   CLOCK_THREAD_CPUTIME_ID allows processes or threads, respectively, to access their
5928   own execution-time clocks. However, given a process or thread, access to its own execution-
5929   time clock is also possible if the clock ID of this clock is obtained through a call to
5930   *clock_getcpuclockid*( ) or *pthread_getcpuclockid*( ).  Therefore, these constants are not necessary
5931   and could be deleted to make the interface simpler. Their existence saves one system call in
5932   the first access to the CPU-time clock of each process or thread. The working group
5933   considered this issue and decided to leave the constants in IEEE Std 1003.1-2001 because they
5934   are closer to the POSIX.1b use of clock identifiers.

5935    • Library Implementations of Threads

5936   In library implementations of threads, kernel entities and library threads can coexist. In this
5937   case, if the CPU-time clocks are supported, most of the clock and timer functions will need to
5938   have two implementations: one in the thread library, and one in the system calls library. The
5939   main difference between these two implementations is that the thread library
5940   implementation will have to deal with clocks and timers that reside in the thread space,
5941   while the kernel implementation will operate on timers and clocks that reside in kernel space.
5942   In the library implementation, if the clock ID refers to a clock that resides in the kernel, a
5943   kernel call will have to be made. The correct version of the function can be chosen by
5944   specifying the appropriate order for the libraries during the link process.

5945    • History of Resolution Issues: Deletion of the *enable* Attribute

5946   In early proposals, consideration was given to inclusion of an attribute called *enable* for CPU-
5947   time clocks. This would allow implementations to avoid the overhead of measuring
5948   execution time for those processes or threads for which this measurement was not required.
5949   However, this is unnecessary since processes are already required to measure execution time
5950   by the POSIX.1 *times*( ) function. Consequently, the *enable* attribute is not present.

5951     **Rationale Relating to Timeouts**

5952     • Requirements for Timeouts

5953     Realtime systems which must operate reliably over extended periods without human
5954     intervention are characteristic in embedded applications such as avionics, machine control,
5955     and space exploration, as well as more mundane applications such as cable TV, security
5956     systems, and plant automation. A multi-tasking paradigm, in which many independent
5957     and/or cooperating software functions relinquish the processor(s) while waiting for a
5958     specific stimulus, resource, condition, or operation completion, is very useful in producing
5959     well engineered programs for such systems. For such systems to be robust and fault-tolerant,
5960     expected occurrences that are unduly delayed or that never occur must be detected so that
5961     appropriate recovery actions may be taken. This is difficult if there is no way for a task to
5962     regain control of a processor once it has relinquished control (blocked) awaiting an
5963     occurrence which, perhaps because of corrupted code, hardware malfunction, or latent
5964     software bugs, will not happen when expected. Therefore, the common practice in realtime
5965     operating systems is to provide a capability to time out such blocking services. Although
5966     there are several methods to achieve this already defined by POSIX, none are as reliable or
5967     efficient as initiating a timeout simultaneously with initiating a blocking service. This is
5968     especially critical in hard-realtime embedded systems because the processors typically have
5969     little time reserve, and allowed fault recovery times are measured in milliseconds rather than
5970     seconds.

5971     The working group largely agreed that such timeouts were necessary and ought to become
5972     part of IEEE Std 1003.1-2001, particularly vendors of realtime operating systems whose
5973     customers had already expressed a strong need for timeouts. There was some resistance to
5974     inclusion of timeouts in IEEE Std 1003.1-2001 because the desired effect, fault tolerance,
5975     could, in theory, be achieved using existing facilities and alternative software designs, but
5976     there was no compelling evidence that realtime system designers would embrace such
5977     designs at the sacrifice of performance and/or simplicity.

5978     • Which Services should be Timed Out?

5979     Originally, the working group considered the prospect of providing timeouts on all blocking
5980     services, including those currently existing in POSIX.1, POSIX.1b, and POSIX.1c, and future
5981     interfaces to be defined by other working groups, as sort of a general policy. This was rather
5982     quickly rejected because of the scope of such a change, and the fact that many of those
5983     services would not normally be used in a realtime context. More traditional timesharing
5984     solutions to timeout would suffice for most of the POSIX.1 interfaces, while others had
5985     asynchronous alternatives which, while more complex to utilize, would be adequate for
5986     some realtime and all non-realtime applications.

5987     The list of potential candidates for timeouts was narrowed to the following for further
5988     consideration:

5989     — POSIX.1b

5990         — *sem_wait*()

5991         — *mq_receive*()

5992         — *mq_send*()

5993         — *lio_listio*()

5994         — *aio_suspend*()

5995         — *sigwait*() (timeout already implemented by *sigtimedwait*())

5996      — POSIX.1c

5997        — *pthread_mutex_lock*( )

5998        — *pthread_join*( )

5999        — *pthread_cond_wait*( ) (timeout already implemented by *pthread_cond_timedwait*( ))

6000      — POSIX.1

6001        — *read*( )

6002        — *write*( )

6003 After further review by the working group, the *lio_listio*( ), *read*( ), and *write*( ) functions (all
6004 forms of blocking synchronous I/O) were eliminated from the list because of the following:

6005 — Asynchronous alternatives exist

6006 — Timeouts can be implemented, albeit non-portably, in device drivers

6007 — A strong desire not to introduce modifications to POSIX.1 interfaces

6008 The working group ultimately rejected *pthread_join*( ) since both that interface and a timed
6009 variant of that interface are non-minimal and may be implemented as a function. See below
6010 for a library implementation of *pthread_join*( ).

6011 Thus, there was a consensus among the working group members to add timeouts to 4 of the
6012 remaining 5 functions (the timeout for *aio_suspend*( ) was ultimately added directly to
6013 POSIX.1b, while the others were added by POSIX.1d). However, *pthread_mutex_lock*( )
6014 remained contentious.

6015 Many feel that *pthread_mutex_lock*( ) falls into the same class as the other functions; that is, it
6016 is desirable to time out a mutex lock because a mutex may fail to be unlocked due to errant or
6017 corrupted code in a critical section (looping or branching outside of the unlock code), and
6018 therefore is equally in need of a reliable, simple, and efficient timeout. In fact, since mutexes
6019 are intended to guard small critical sections, most *pthread_mutex_lock*( ) calls would be
6020 expected to obtain the lock without blocking nor utilizing any kernel service, even in
6021 implementations of threads with global contention scope; the timeout alternative need only
6022 be considered after it is determined that the thread must block.

6023 Those opposed to timing out mutexes feel that the very simplicity of the mutex is
6024 compromised by adding a timeout semantic, and that to do so is senseless. They claim that if
6025 a timed mutex is really deemed useful by a particular application, then it can be constructed
6026 from the facilities already in POSIX.1b and POSIX.1c. The following two C-language library
6027 implementations of mutex locking with timeout represent the solutions offered (in both
6028 implementations, the timeout parameter is specified as absolute time, not relative time as in
6029 the proposed POSIX.1c interfaces).

- Spinlock Implementation

```
6031            #include <pthread.h>
6032            #include <time.h>
6033            #include <errno.h>

6034            int pthread_mutex_timedlock(pthread_mutex_t *mutex,
6035                    const struct timespec *timeout)
6036                {
6037                struct timespec timenow;

6038                while (pthread_mutex_trylock(mutex) == EBUSY)
6039                    {
6040                    clock_gettime(CLOCK_REALTIME, &timenow);
6041                    if (timespec_cmp(&timenow,timeout) >= 0)
6042                        {
6043                        return ETIMEDOUT;
6044                    }
6045                    pthread_yield();
6046                    }
6047                return 0;
6048                }
```

The Spinlock implementation is generally unsuitable for any application using priority-based thread scheduling policies such as SCHED_FIFO or SCHED_RR, since the mutex could currently be held by a thread of lower priority within the same allocation domain, but since the waiting thread never blocks, only threads of equal or higher priority will ever run, and the mutex cannot be unlocked. Setting priority inheritance or priority ceiling protocol on the mutex does not solve this problem, since the priority of a mutex owning thread is only boosted if higher priority threads are blocked waiting for the mutex; clearly not the case for this spinlock.

- Condition Wait Implementation

```
6058            #include <pthread.h>
6059            #include <time.h>
6060            #include <errno.h>

6061            struct timed_mutex
6062                {
6063                int locked;
6064                pthread_mutex_t mutex;
6065                pthread_cond_t cond;
6066                };
6067            typedef struct timed_mutex timed_mutex_t;

6068            int timed_mutex_lock(timed_mutex_t *tm,
6069                    const struct timespec *timeout)
6070                {
6071                int timedout=FALSE;
6072                int error_status;

6073                pthread_mutex_lock(&tm->mutex);

6074                while (tm->locked && !timedout)
6075                    {
6076                    if ((error_status=pthread_cond_timedwait(&tm->cond,
```

```
6077                          &tm->mutex,
6078                          timeout))!=0)
6079                     {
6080                     if (error_status==ETIMEDOUT) timedout = TRUE;
6081                     }
6082                 }
6083             if(timedout)
6084                 {
6085                 pthread_mutex_unlock(&tm->mutex);
6086                 return ETIMEDOUT;
6087                 }
6088             else
6089                 {
6090                 tm->locked = TRUE;
6091                 pthread_mutex_unlock(&tm->mutex);
6092                 return 0;
6093                 }
6094             }
6095     void timed_mutex_unlock(timed_mutex_t *tm)
6096         {
6097         pthread_mutex_lock(&tm->mutex); / for case assignment not atomic /
6098         tm->locked = FALSE;
6099         pthread_mutex_unlock(&tm->mutex);
6100         pthread_cond_signal(&tm->cond);
6101         }
```

6102  The Condition Wait implementation effectively substitutes the *pthread_cond_timedwait*( )
6103  function (which is currently timed out) for the desired *pthread_mutex_timedlock*( ). Since waits
6104  on condition variables currently do not include protocols which avoid priority inversion, this
6105  method is generally unsuitable for realtime applications because it does not provide the same
6106  priority inversion protection as the untimed *pthread_mutex_lock*( ). Also, for any given
6107  implementations of the current mutex and condition variable primitives, this library
6108  implementation has a performance cost at least 2.5 times that of the untimed
6109  *pthread_mutex_lock*( ) even in the case where the timed mutex is readily locked without
6110  blocking (the interfaces required for this case are shown in bold). Even in uniprocessors or
6111  where assignment is atomic, at least an additional *pthread_cond_signal*( ) is required.
6112  *pthread_mutex_timedlock*( ) could be implemented at effectively no performance penalty in
6113  this case because the timeout parameters need only be considered after it is determined that
6114  the mutex cannot be locked immediately.

6115  Thus it has not yet been shown that the full semantics of mutex locking with timeout can be
6116  efficiently and reliably achieved using existing interfaces. Even if the existence of an
6117  acceptable library implementation were proven, it is difficult to justify why the interface
6118  itself should not be made portable, especially considering approval for the other four
6119  timeouts.

6120    • Rationale for Library Implementation of *pthread_timedjoin*()

6121    Library implementation of *pthread_timedjoin*():

```
6122        /*
6123         * Construct a thread variety entirely from existing functions
6124         * with which a join can be done, allowing the join to time out.
6125         */
6126        #include <pthread.h>
6127        #include <time.h>

6128        struct timed_thread {
6129            pthread_t t;
6130            pthread_mutex_t m;
6131            int exiting;
6132            pthread_cond_t exit_c;
6133            void *(*start_routine)(void *arg);
6134            void *arg;
6135            void *status;
6136        };

6137        typedef struct timed_thread *timed_thread_t;
6138        static pthread_key_t timed_thread_key;
6139        static pthread_once_t timed_thread_once = PTHREAD_ONCE_INIT;

6140        static void timed_thread_init()
6141        {
6142            pthread_key_create(&timed_thread_key, NULL);
6143        }

6144        static void *timed_thread_start_routine(void *args)

6145        /*
6146         * Routine to establish thread-specific data value and run the actual
6147         * thread start routine which was supplied to timed_thread_create().
6148         */
6149        {
6150            timed_thread_t tt = (timed_thread_t) args;

6151            pthread_once(&timed_thread_once, timed_thread_init);
6152            pthread_setspecific(timed_thread_key, (void *)tt);
6153            timed_thread_exit((tt->start_routine)(tt->arg));
6154        }

6155        int timed_thread_create(timed_thread_t ttp, const pthread_attr_t *attr,
6156            void *(*start_routine)(void *), void *arg)

6157        /*
6158         * Allocate a thread which can be used with timed_thread_join().
6159         */
6160        {
6161            timed_thread_t tt;
6162            int result;

6163            tt = (timed_thread_t) malloc(sizeof(struct timed_thread));
6164            pthread_mutex_init(&tt->m,NULL);
6165            tt->exiting = FALSE;
6166            pthread_cond_init(&tt->exit_c,NULL);
```

```
6167                    tt->start_routine = start_routine;
6168                    tt->arg = arg;
6169                    tt->status = NULL;

6170                    if ((result = pthread_create(&tt->t, attr,
6171                        timed_thread_start_routine, (void *)tt)) != 0) {
6172                        free(tt);
6173                        return result;
6174                    }

6175                    pthread_detach(tt->t);
6176                    ttp = tt;
6177                    return 0;
6178                }

6179           int timed_thread_join(timed_thread_t tt,
6180                    struct timespec *timeout,
6181                    void **status)
6182                {
6183                    int result;

6184                    pthread_mutex_lock(&tt->m);
6185                    result = 0;
6186                    /*
6187                     * Wait until the thread announces that it is exiting,
6188                     * or until timeout.
6189                     */
6190                    while (result == 0 && ! tt->exiting) {
6191                        result = pthread_cond_timedwait(&tt->exit_c, &tt->m, timeout);
6192                    }
6193                    pthread_mutex_unlock(&tt->m);
6194                    if (result == 0 && tt->exiting) {
6195                        *status = tt->status;
6196                        free((void *)tt);
6197                        return result;
6198                    }
6199                    return result;
6200                }

6201           void timed_thread_exit(void *status)
6202                {
6203                    timed_thread_t tt;
6204                    void *specific;

6205                    if ((specific=pthread_getspecific(timed_thread_key)) == NULL){
6206                        /*
6207                         * Handle cases which won't happen with correct usage.
6208                         */
6209                        pthread_exit( NULL);
6210                    }
6211                    tt = (timed_thread_t) specific;
6212                    pthread_mutex_lock(&tt->m);
6213                    /*
6214                     * Tell a joiner that we're exiting.
6215                     */
```

```
6216            tt->status = status;
6217            tt->exiting = TRUE;
6218            pthread_cond_signal(&tt->exit_c);
6219            pthread_mutex_unlock(&tt->m);
6220            /*
6221             * Call pthread exit() to call destructors and really
6222             * exit the thread.
6223             */
6224            pthread_exit(NULL);
6225        }
```

6226    The *pthread_join*( ) C-language example shown above demonstrates that it is possible, using
6227    existing pthread facilities, to construct a variety of thread which allows for joining such a
6228    thread, but which allows the join operation to time out. It does this by using a
6229    *pthread_cond_timedwait*( ) to wait for the thread to exit. A **timed_thread_t** descriptor structure
6230    is used to pass parameters from the creating thread to the created thread, and from the
6231    exiting thread to the joining thread. This implementation is roughly equivalent to what a
6232    normal *pthread_join*( ) implementation would do, with the single change being that
6233    *pthread_cond_timedwait*( ) is used in place of a simple *pthread_cond_wait*( ).

6234    Since it is possible to implement such a facility entirely from existing pthread interfaces, and
6235    with roughly equal efficiency and complexity to an implementation which would be
6236    provided directly by a pthreads implementation, it was the consensus of the working group
6237    members that any *pthread_timedjoin*( ) facility would be unnecessary, and should not be
6238    provided.

6239    • Form of the Timeout Interfaces

6240    The working group considered a number of alternative ways to add timeouts to blocking
6241    services. At first, a system interface which would specify a one-shot or persistent timeout to
6242    be applied to subsequent blocking services invoked by the calling process or thread was
6243    considered because it allowed all blocking services to be timed out in a uniform manner with
6244    a single additional interface; this was rather quickly rejected because it could easily result in
6245    the wrong services being timed out.

6246    It was suggested that a timeout value might be specified as an attribute of the object
6247    (semaphore, mutex, message queue, and so on), but there was no consensus on this, either on
6248    a case-by-case basis or for all timeouts.

6249    Looking at the two existing timeouts for blocking services indicates that the working group
6250    members favor a separate interface for the timed version of a function. However,
6251    *pthread_cond_timedwait*( ) utilizes an absolute timeout value while *sigtimedwait*( ) uses a
6252    relative timeout value. The working group members agreed that relative timeout values are
6253    appropriate where the timeout mechanism's primary use was to deal with an unexpected or
6254    error situation, but they are inappropriate when the timeout must expire at a particular time,
6255    or before a specific deadline. For the timeouts being introduced in IEEE Std 1003.1-2001, the
6256    working group considered allowing both relative and absolute timeouts as is done with
6257    POSIX.1b timers, but ultimately favored the simpler absolute timeout form.

6258    An absolute time measure can be easily implemented on top of an interface that specifies
6259    relative time, by reading the clock, calculating the difference between the current time and
6260    the desired wake-up time, and issuing a relative timeout call. But there is a race condition
6261    with this approach because the thread could be preempted after reading the clock, but before
6262    making the timed-out call; in this case, the thread would be awakened later than it should
6263    and, thus, if the wake-up time represented a deadline, it would miss it.

6264  There is also a race condition when trying to build a relative timeout on top of an interface
6265  that specifies absolute timeouts. In this case, the clock would have to be read to calculate the
6266  absolute wake-up time as the sum of the current time plus the relative timeout interval. In
6267  this case, if the thread is preempted after reading the clock but before making the timed-out
6268  call, the thread would be awakened earlier than desired.

6269  But the race condition with the absolute timeouts interface is not as bad as the one that
6270  happens with the relative timeout interface, because there are simple workarounds. For the
6271  absolute timeouts interface, if the timing requirement is a deadline, the deadline can still be
6272  met because the thread woke up earlier than the deadline. If the timeout is just used as an
6273  error recovery mechanism, the precision of timing is not really important. If the timing
6274  requirement is that between actions A and B a minimum interval of time must elapse, the
6275  absolute timeout interface can be safely used by reading the clock after action A has been
6276  started. It could be argued that, since the call with the absolute timeout is atomic from the
6277  application point of view, it is not possible to read the clock after action A, if this action is
6278  part of the timed-out call. But looking at the nature of the calls for which timeouts are
6279  specified (locking a mutex, waiting for a semaphore, waiting for a message, or waiting until
6280  there is space in a message queue), the timeouts that an application would build on these
6281  actions would not be triggered by these actions themselves, but by some other external
6282  action. For example, if waiting for a message to arrive to a message queue, and waiting for at
6283  least 20 milliseconds, this time interval would start to be counted from some event that
6284  would trigger both the action that produces the message, as well as the action that waits for
6285  the message to arrive, and not by the wait-for-message operation itself. In this case, the
6286  workaround proposed above could be used.

6287  For these reasons, the absolute timeout is preferred over the relative timeout interface.

6288  **B.2.9    Threads**

6289  Threads will normally be more expensive than subroutines (or functions, routines, and so on) if
6290  specialized hardware support is not provided. Nevertheless, threads should be sufficiently
6291  efficient to encourage their use as a medium to fine-grained structuring mechanism for
6292  parallelism in an application. Structuring an application using threads then allows it to take
6293  immediate advantage of any underlying parallelism available in the host environment. This
6294  means implementors are encouraged to optimize for fast execution at the possible expense of
6295  efficient utilization of storage. For example, a common thread creation technique is to cache
6296  appropriate thread data structures. That is, rather than releasing system resources, the
6297  implementation retains these resources and reuses them when the program next asks to create a
6298  new thread. If this reuse of thread resources is to be possible, there has to be very little unique
6299  state associated with each thread, because any such state has to be reset when the thread is
6300  reused.

6301  **Thread Creation Attributes**

6302  Attributes objects are provided for threads, mutexes, and condition variables as a mechanism to
6303  support probable future standardization in these areas without requiring that the interface itself
6304  be changed.

6305  Attributes objects provide clean isolation of the configurable aspects of threads. For example,
6306  ''stack size'' is an important attribute of a thread, but it cannot be expressed portably. When
6307  porting a threaded program, stack sizes often need to be adjusted. The use of attributes objects
6308  can help by allowing the changes to be isolated in a single place, rather than being spread across
6309  every instance of thread creation.

6310 Attributes objects can be used to set up *classes* of threads with similar attributes; for example,
6311 ''threads with large stacks and high priority'' or ''threads with minimal stacks''. These classes
6312 can be defined in a single place and then referenced wherever threads need to be created.
6313 Changes to ''class'' decisions become straightforward, and detailed analysis of each
6314 *pthread_create*() call is not required.

6315 The attributes objects are defined as opaque types as an aid to extensibility. If these objects had
6316 been specified as structures, adding new attributes would force recompilation of all multi-
6317 threaded programs when the attributes objects are extended; this might not be possible if
6318 different program components were supplied by different vendors.

6319 Additionally, opaque attributes objects present opportunities for improving performance.
6320 Argument validity can be checked once when attributes are set, rather than each time a thread is
6321 created. Implementations will often need to cache kernel objects that are expensive to create.
6322 Opaque attributes objects provide an efficient mechanism to detect when cached objects become
6323 invalid due to attribute changes.

6324 Because assignment is not necessarily defined on a given opaque type, implementation-defined
6325 default values cannot be defined in a portable way. The solution to this problem is to allow
6326 attribute objects to be initialized dynamically by attributes object initialization functions, so that
6327 default values can be supplied automatically by the implementation.

6328 The following proposal was provided as a suggested alternative to the supplied attributes:

6329 1. Maintain the style of passing a parameter formed by the bitwise-inclusive OR of flags to
6330     the initialization routines (*pthread_create*(), *pthread_mutex_init*(), *pthread_cond_init*()). The
6331     parameter containing the flags should be an opaque type for extensibility. If no flags are
6332     set in the parameter, then the objects are created with default characteristics. An
6333     implementation may specify implementation-defined flag values and associated behavior.

6334 2. If further specialization of mutexes and condition variables is necessary, implementations
6335     may specify additional procedures that operate on the **pthread_mutex_t** and
6336     **pthread_cond_t** objects (instead of on attributes objects).

6337 The difficulties with this solution are:

6338 1. A bitmask is not opaque if bits have to be set into bit-vector attributes objects using
6339     explicitly-coded bitwise-inclusive OR operations. If the set of options exceeds an **int**,
6340     application programmers need to know the location of each bit. If bits are set or read by
6341     encapsulation (that is, *get*\*() or *set*\*() functions), then the bitmask is merely an
6342     implementation of attributes objects as currently defined and should not be exposed to the
6343     programmer.

6344 2. Many attributes are not Boolean or very small integral values. For example, scheduling
6345     policy may be placed in 3 bits or 4 bits, but priority requires 5 bits or more, thereby taking
6346     up at least 8 bits out of a possible 16 bits on machines with 16-bit integers. Because of this,
6347     the bitmask can only reasonably control whether particular attributes are set or not, and it
6348     cannot serve as the repository of the value itself. The value needs to be specified as a
6349     function parameter (which is non-extensible), or by setting a structure field (which is non-
6350     opaque), or by *get*\*() and *set*\*() functions (making the bitmask a redundant addition to the
6351     attributes objects).

6352 Stack size is defined as an optional attribute because the very notion of a stack is inherently
6353 machine-dependent. Some implementations may not be able to change the size of the stack, for
6354 example, and others may not need to because stack pages may be discontiguous and can be
6355 allocated and released on demand.

6356　　The attribute mechanism has been designed in large measure for extensibility. Future extensions
6357　　to the attribute mechanism or to any attributes object defined in IEEE Std 1003.1-2001 have to be
6358　　done with care so as not to affect binary-compatibility.

6359　　Attribute objects, even if allocated by means of dynamic allocation functions such as *malloc*(),
6360　　may have their size fixed at compile time. This means, for example, a *pthread_create*() in an
6361　　implementation with extensions to the **pthread_attr_t** cannot look beyond the area that the
6362　　binary application assumes is valid. This suggests that implementations should maintain a size
6363　　field in the attributes object, as well as possibly version information, if extensions in different
6364　　directions (possibly by different vendors) are to be accommodated.

6365　　**Thread Implementation Models**

6366　　There are various thread implementation models. At one end of the spectrum is the ''library-
6367　　thread model''. In such a model, the threads of a process are not visible to the operating system
6368　　kernel, and the threads are not kernel-scheduled entities. The process is the only kernel-
6369　　scheduled entity. The process is scheduled onto the processor by the kernel according to the
6370　　scheduling attributes of the process. The threads are scheduled onto the single kernel-scheduled
6371　　entity (the process) by the runtime library according to the scheduling attributes of the threads.
6372　　A problem with this model is that it constrains concurrency. Since there is only one kernel-
6373　　scheduled entity (namely, the process), only one thread per process can execute at a time. If the
6374　　thread that is executing blocks on I/O, then the whole process blocks.

6375　　At the other end of the spectrum is the ''kernel-thread model''. In this model, all threads are
6376　　visible to the operating system kernel. Thus, all threads are kernel-scheduled entities, and all
6377　　threads can concurrently execute. The threads are scheduled onto processors by the kernel
6378　　according to the scheduling attributes of the threads. The drawback to this model is that the
6379　　creation and management of the threads entails operating system calls, as opposed to subroutine
6380　　calls, which makes kernel threads heavier weight than library threads.

6381　　Hybrids of these two models are common. A hybrid model offers the speed of library threads
6382　　and the concurrency of kernel threads. In hybrid models, a process has some (relatively small)
6383　　number of kernel scheduled entities associated with it. It also has a potentially much larger
6384　　number of library threads associated with it. Some library threads may be bound to kernel-
6385　　scheduled entities, while the other library threads are multiplexed onto the remaining kernel-
6386　　scheduled entities. There are two levels of thread scheduling:

6387　　　1.　The runtime library manages the scheduling of (unbound) library threads onto kernel-
6388　　　　　scheduled entities.

6389　　　2.　The kernel manages the scheduling of kernel-scheduled entities onto processors.

6390　　For this reason, a hybrid model is referred to as a two-level threads scheduling model. In this
6391　　model, the process can have multiple concurrently executing threads; specifically, it can have as
6392　　many concurrently executing threads as it has kernel-scheduled entities.

6393　　**Thread-Specific Data**

6394　　Many applications require that a certain amount of context be maintained on a per-thread basis
6395　　across procedure calls. A common example is a multi-threaded library routine that allocates
6396　　resources from a common pool and maintains an active resource list for each thread. The
6397　　thread-specific data interface provided to meet these needs may be viewed as a two-dimensional
6398　　array of values with keys serving as the row index and thread IDs as the column index (although
6399　　the implementation need not work this way).

6400          • Models

6401          Three possible thread-specific data models were considered:

6402              1.   No Explicit Support

6403                   A standard thread-specific data interface is not strictly necessary to support
6404                   applications that require per-thread context. One could, for example, provide a hash
6405                   function that converted a **pthread_t** into an integer value that could then be used to
6406                   index into a global array of per-thread data pointers. This hash function, in conjunction
6407                   with *pthread_self*(), would be all the interface required to support a mechanism of this
6408                   sort. Unfortunately, this technique is cumbersome. It can lead to duplicated code as
6409                   each set of cooperating modules implements their own per-thread data management
6410                   schemes.

6411              2.   Single (**void** *) Pointer

6412                   Another technique would be to provide a single word of per-thread storage and a pair
6413                   of functions to fetch and store the value of this word. The word could then hold a
6414                   pointer to a block of per-thread memory. The allocation, partitioning, and general use
6415                   of this memory would be entirely up to the application. Although this method is not as
6416                   problematic as technique 1, it suffers from interoperability problems. For example, all
6417                   modules using the per-thread pointer would have to agree on a common usage
6418                   protocol.

6419              3.   Key/Value Mechanism

6420                   This method associates an opaque key (for example, stored in a variable of type
6421                   **pthread_key_t**) with each per-thread datum. These keys play the role of identifiers for
6422                   per-thread data. This technique is the most generic and avoids the problems noted
6423                   above, albeit at the cost of some complexity.

6424          The primary advantage of the third model is its information hiding properties. Modules
6425          using this model are free to create and use their own key(s) independent of all other such
6426          usage, whereas the other models require that all modules that use thread-specific context
6427          explicitly cooperate with all other such modules. The data-independence provided by the
6428          third model is worth the additional interface.

6429          • Requirements

6430          It is important that it be possible to implement the thread-specific data interface without the
6431          use of thread private memory. To do otherwise would increase the weight of each thread,
6432          thereby limiting the range of applications for which the threads interfaces provided by
6433          IEEE Std 1003.1-2001 is appropriate.

6434          The values that one binds to the key via *pthread_setspecific*() may, in fact, be pointers to
6435          shared storage locations available to all threads. It is only the key/value bindings that are
6436          maintained on a per-thread basis, and these can be kept in any portion of the address space
6437          that is reserved for use by the calling thread (for example, on the stack). Thus, no per-thread
6438          MMU state is required to implement the interface. On the other hand, there is nothing in the
6439          interface specification to preclude the use of a per-thread MMU state if it is available (for
6440          example, the key values returned by *pthread_key_create*() could be thread private memory
6441          addresses).

6442          • Standardization Issues

6443          Thread-specific data is a requirement for a usable thread interface.  The binding described in
6444          this section provides a portable thread-specific data mechanism for languages that do not
6445          directly support a thread-specific storage class. A binding to IEEE Std 1003.1-2001 for a

6446      language that does include such a storage class need not provide this specific interface.

6447      If a language were to include the notion of thread-specific storage, it would be desirable (but
6448      *not* required) to provide an implementation of the pthreads thread-specific data interface
6449      based on the language feature. For example, assume that a compiler for a C-like language
6450      supports a *private* storage class that provides thread-specific storage. Something similar to
6451      the following macros might be used to effect a compatible implementation:

```
6452    #define pthread_key_t                 private void *
6453    #define pthread_key_create(key)       /* no-op */
6454    #define pthread_setspecific(key,value) (key)=(value)
6455    #define pthread_getspecific(key)      (key)
```

6456      **Note:**      For the sake of clarity, this example ignores destructor functions. A correct implementation
6457                    would have to support them.

6458      **Barriers**

6459      • Background

6460      Barriers are typically used in parallel DO/FOR loops to ensure that all threads have reached
6461      a particular stage in a parallel computation before allowing any to proceed to the next stage.
6462      Highly efficient implementation is possible on machines which support a ''Fetch and Add''
6463      operation as described in the referenced Almasi and Gottlieb (1989).

6464      The use of return value PTHREAD_BARRIER_SERIAL_THREAD is shown in the following
6465      example:

```
6466    if ( (status=pthread_barrier_wait(&barrier)) ==
6467        PTHREAD_BARRIER_SERIAL_THREAD) {
6468        ...serial section
6469        }
6470            else if (status != 0) {
6471            ...error processing
6472        }
6473    status=pthread_barrier_wait(&barrier);
6474    ...
```

6475      This behavior allows a serial section of code to be executed by one thread as soon as all
6476      threads reach the first barrier. The second barrier prevents the other threads from proceeding
6477      until the serial section being executed by the one thread has completed.

6478      Although barriers can be implemented with mutexes and condition variables, the referenced
6479      Almasi and Gottlieb (1989) provides ample illustration that such implementations are
6480      significantly less efficient than is possible. While the relative efficiency of barriers may well
6481      vary by implementation, it is important that they be recognized in the IEEE Std 1003.1-2001
6482      to facilitate applications portability while providing the necessary freedom to implementors.

6483      • Lack of Timeout Feature

6484      Alternate versions of most blocking routines have been provided to support watchdog
6485      timeouts. No alternate interface of this sort has been provided for barrier waits for the
6486      following reasons:

6487          • Multiple threads may use different timeout values, some of which may be indefinite. It is
6488          not clear which threads should break through the barrier with a timeout error if and when
6489          these timeouts expire.

6490 • The barrier may become unusable once a thread breaks out of a *pthread_barrier_wait*( )
6491 with a timeout error. There is, in general, no way to guarantee the consistency of a
6492 barrier's internal data structures once a thread has timed out of a *pthread_barrier_wait*( ).
6493 Even the inclusion of a special barrier reinitialization function would not help much since
6494 it is not clear how this function would affect the behavior of threads that reach the barrier
6495 between the original timeout and the call to the reinitialization function.

6496 **Spin Locks**

6497 • Background

6498 Spin locks represent an extremely low-level synchronization mechanism suitable primarily
6499 for use on shared memory multi-processors. It is typically an atomically modified Boolean
6500 value that is set to one when the lock is held and to zero when the lock is freed.

6501 When a caller requests a spin lock that is already held, it typically spins in a loop testing
6502 whether the lock has become available. Such spinning wastes processor cycles so the lock
6503 should only be held for short durations and not across sleep/block operations. Callers should
6504 unlock spin locks before calling sleep operations.

6505 Spin locks are available on a variety of systems. The functions included in
6506 IEEE Std 1003.1-2001 are an attempt to standardize that existing practice.

6507 • Lack of Timeout Feature

6508 Alternate versions of most blocking routines have been provided to support watchdog
6509 timeouts. No alternate interface of this sort has been provided for spin locks for the following
6510 reasons:

6511 • It is impossible to determine appropriate timeout intervals for spin locks in a portable
6512 manner. The amount of time one can expect to spend spin-waiting is inversely
6513 proportional to the degree of parallelism provided by the system.

6514 It can vary from a few cycles when each competing thread is running on its own
6515 processor, to an indefinite amount of time when all threads are multiplexed on a single
6516 processor (which is why spin locking is not advisable on uniprocessors).

6517 • When used properly, the amount of time the calling thread spends waiting on a spin lock
6518 should be considerably less than the time required to set up a corresponding watchdog
6519 timer. Since the primary purpose of spin locks is to provide a low-overhead
6520 synchronization mechanism for multi-processors, the overhead of a timeout mechanism
6521 was deemed unacceptable.

6522 It was also suggested that an additional *count* argument be provided (on the
6523 *pthread_spin_lock*( ) call) in *lieu* of a true timeout so that a spin lock call could fail gracefully if
6524 it was unable to apply the lock after *count* attempts. This idea was rejected because it is not
6525 existing practice. Furthermore, the same effect can be obtained with *pthread_spin_trylock*( ),
6526 as illustrated below:

```
6527              int n = MAX_SPIN;

6528              while ( −−n >= 0 )
6529              {
6530                  if ( !pthread_spin_try_lock(...) )
6531                      break;
6532              }
6533              if ( n >= 0 )
6534              {
6535                  /* Successfully acquired the lock */
6536              }
6537              else
6538              {
6539                  /* Unable to acquire the lock */
6540              }
```

6541    • *process-shared* Attribute

6542    The initialization functions associated with most POSIX synchronization objects (for
6543    example, mutexes, barriers, and read-write locks) take an attributes object with a *process-*
6544    *shared* attribute that specifies whether or not the object is to be shared across processes. In the
6545    draft corresponding to the first balloting round, two separate initialization functions are
6546    provided for spin locks, however: one for spin locks that were to be shared across processes
6547    (*spin_init*()), and one for locks that were only used by multiple threads within a single
6548    process (*pthread_spin_init*()). This was done so as to keep the overhead associated with spin
6549    waiting to an absolute minimum. However, the balloting group requested that, since the
6550    overhead associated to a bit check was small, spin locks should be consistent with the rest of
6551    the synchronization primitives, and thus the *process-shared* attribute was introduced for spin
6552    locks.

6553    • Spin Locks *versus* Mutexes

6554    It has been suggested that mutexes are an adequate synchronization mechanism and spin
6555    locks are not necessary. Locking mechanisms typically must trade off the processor resources
6556    consumed while setting up to block the thread and the processor resources consumed by the
6557    thread while it is blocked. Spin locks require very little resources to set up the blocking of a
6558    thread. Existing practice is to simply loop, repeating the atomic locking operation until the
6559    lock is available. While the resources consumed to set up blocking of the thread are low, the
6560    thread continues to consume processor resources while it is waiting.

6561    On the other hand, mutexes may be implemented such that the processor resources
6562    consumed to block the thread are large relative to a spin lock. After detecting that the mutex
6563    lock is not available, the thread must alter its scheduling state, add itself to a set of waiting
6564    threads, and, when the lock becomes available again, undo all of this before taking over
6565    ownership of the mutex. However, while a thread is blocked by a mutex, no processor
6566    resources are consumed.

6567    Therefore, spin locks and mutexes may be implemented to have different characteristics.
6568    Spin locks may have lower overall overhead for very short-term blocking, and mutexes may
6569    have lower overall overhead when a thread will be blocked for longer periods of time. The
6570    presence of both interfaces allows implementations with these two different characteristics,
6571    both of which may be useful to a particular application.

6572    It has also been suggested that applications can build their own spin locks from the
6573    *pthread_mutex_trylock*() function:

```
6574                 while (pthread_mutex_trylock(&mutex));
```

6575    The apparent simplicity of this construct is somewhat deceiving, however. While the actual
6576    wait is quite efficient, various guarantees on the integrity of mutex objects (for example,
6577    priority inheritance rules) may add overhead to the successful path of the trylock operation
6578    that is not required of spin locks. One could, of course, add an attribute to the mutex to
6579    bypass such overhead, but the very act of finding and testing this attribute represents more
6580    overhead than is found in the typical spin lock.

6581    The need to hold spin lock overhead to an absolute minimum also makes it impossible to
6582    provide guarantees against starvation similar to those provided for mutexes or read-write
6583    locks. The overhead required to implement such guarantees (for example, disabling
6584    preemption before spinning) may well exceed the overhead of the spin wait itself by many
6585    orders of magnitude. If a ''safe'' spin wait seems desirable, it can always be provided (albeit
6586    at some performance cost) via appropriate mutex attributes.

6587    **XSI Supported Functions**

6588    On XSI-conformant systems, the following symbolic constants are always defined:

6589    _POSIX_READER_WRITER_LOCKS
6590    _POSIX_THREAD_ATTR_STACKADDR
6591    _POSIX_THREAD_ATTR_STACKSIZE
6592    _POSIX_THREAD_PROCESS_SHARED
6593    _POSIX_THREADS

6594    Therefore, the following threads functions are always supported:

6595    *pthread_atfork*( )                              *pthread_key_create*( )
6596    *pthread_attr_destroy*( )                        *pthread_key_delete*( )
6597    *pthread_attr_getdetachstate*( )                 *pthread_kill*( )
6598    *pthread_attr_getguardsize*( )                   *pthread_mutex_destroy*( )
6599    *pthread_attr_getschedparam*( )                  *pthread_mutex_init*( )
6600    *pthread_attr_getstack*( )                       *pthread_mutex_lock*( )
6601    *pthread_attr_getstackaddr*( )                   *pthread_mutex_trylock*( )
6602    *pthread_attr_getstacksize*( )                   *pthread_mutex_unlock*( )
6603    *pthread_attr_init*( )                           *pthread_mutexattr_destroy*( )
6604    *pthread_attr_setdetachstate*( )                 *pthread_mutexattr_getpshared*( )
6605    *pthread_attr_setguardsize*( )                   *pthread_mutexattr_gettype*( )
6606    *pthread_attr_setschedparam*( )                  *pthread_mutexattr_init*( )
6607    *pthread_attr_setstack*( )                       *pthread_mutexattr_setpshared*( )
6608    *pthread_attr_setstackaddr*( )                   *pthread_mutexattr_settype*( )
6609    *pthread_attr_setstacksize*( )                   *pthread_once*( )
6610    *pthread_cancel*( )                              *pthread_rwlock_destroy*( )
6611    *pthread_cleanup_pop*( )                         *pthread_rwlock_init*( )
6612    *pthread_cleanup_push*( )                        *pthread_rwlock_rdlock*( )
6613    *pthread_cond_broadcast*( )                      *pthread_rwlock_tryrdlock*( )
6614    *pthread_cond_destroy*( )                        *pthread_rwlock_trywrlock*( )
6615    *pthread_cond_init*( )                           *pthread_rwlock_unlock*( )
6616    *pthread_cond_signal*( )                         *pthread_rwlock_wrlock*( )
6617    *pthread_cond_timedwait*( )                      *pthread_rwlockattr_destroy*( )
6618    *pthread_cond_wait*( )                           *pthread_rwlockattr_getpshared*( )
6619    *pthread_condattr_destroy*( )                    *pthread_rwlockattr_init*( )

| | |
|---|---|
| 6620 | *pthread_condattr_getpshared*( ) | *pthread_rwlockattr_setpshared*( ) |
| 6621 | *pthread_condattr_init*( ) | *pthread_self*( ) |
| 6622 | *pthread_condattr_setpshared*( ) | *pthread_setcancelstate*( ) |
| 6623 | *pthread_create*( ) | *pthread_setcanceltype*( ) |
| 6624 | *pthread_detach*( ) | *pthread_setconcurrency*( ) |
| 6625 | *pthread_equal*( ) | *pthread_setspecific*( ) |
| 6626 | *pthread_exit*( ) | *pthread_sigmask*( ) |
| 6627 | *pthread_getconcurrency*( ) | *pthread_testcancel*( ) |
| 6628 | *pthread_getspecific*( ) | *sigwait*( ) |
| 6629 | *pthread_join*( ) | |

6630  On XSI-conformant systems, the symbolic constant _POSIX_THREAD_SAFE_FUNCTIONS is
6631  always defined. Therefore, the following functions are always supported:

| | | |
|---|---|---|
| 6632 | *asctime_r*( ) | *getpwuid_r*( ) |
| 6633 | *ctime_r*( ) | *gmtime_r*( ) |
| 6634 | *flockfile*( ) | *localtime_r*( ) |
| 6635 | *ftrylockfile*( ) | *putc_unlocked*( ) |
| 6636 | *funlockfile*( ) | *putchar_unlocked*( ) |
| 6637 | *getc_unlocked*( ) | *rand_r*( ) |
| 6638 | *getchar_unlocked*( ) | *readdir_r*( ) |
| 6639 | *getgrgid_r*( ) | *strerror_r*( ) |
| 6640 | *getgrnam_r*( ) | *strtok_r*( ) |
| 6641 | *getpwnam_r*( ) | |

6642  The following threads functions are only supported on XSI-conformant systems if the Realtime
6643  Threads Option Group is supported :

| | | |
|---|---|---|
| 6644 | *pthread_attr_getinheritsched*( ) | *pthread_mutex_getprioceiling*( ) |
| 6645 | *pthread_attr_getschedpolicy*( ) | *pthread_mutex_setprioceiling*( ) |
| 6646 | *pthread_attr_getscope*( ) | *pthread_mutexattr_getprioceiling*( ) |
| 6647 | *pthread_attr_setinheritsched*( ) | *pthread_mutexattr_getprotocol*( ) |
| 6648 | *pthread_attr_setschedpolicy*( ) | *pthread_mutexattr_setprioceiling*( ) |
| 6649 | *pthread_attr_setscope*( ) | *pthread_mutexattr_setprotocol*( ) |
| 6650 | *pthread_getschedparam*( ) | *pthread_setschedparam*( ) |

6651  **XSI Threads Extensions**

6652  The following XSI extensions to POSIX.1c are now supported in IEEE Std 1003.1-2001 as part of
6653  the alignment with the Single UNIX Specification:

6654  • Extended mutex attribute types

6655  • Read-write locks and attributes (also introduced by the IEEE Std 1003.1j-2000 amendment)

6656  • Thread concurrency level

6657  • Thread stack guard size

6658  • Parallel I/O

6659  A total of 19 new functions were added.

6660  These extensions carefully follow the threads programming model specified in POSIX.1c. As
6661  with POSIX.1c, all the new functions return zero if successful; otherwise, an error number is

6662        returned to indicate the error.

6663        The concept of attribute objects was introduced in POSIX.1c to allow implementations to extend
6664        IEEE Std 1003.1-2001 without changing the existing interfaces. Attribute objects were defined for
6665        threads, mutexes, and condition variables. Attributes objects are defined as implementation-
6666        defined opaque types to aid extensibility, and functions are defined to allow attributes to be set
6667        or retrieved. This model has been followed when adding the new type attribute of
6668        **pthread_mutexattr_t** or the new read-write lock attributes object **pthread_rwlockattr_t**.

6669        • Extended Mutex Attributes

6670           POSIX.1c defines a mutex attributes object as an implementation-defined opaque object of
6671           type **pthread_mutexattr_t**, and specifies a number of attributes which this object must have
6672           and a number of functions which manipulate these attributes. These attributes include
6673           *detachstate*, *inheritsched*, *schedparm*, *schedpolicy*, *contentionscope*, *stackaddr*, and *stacksize*.

6674           The System Interfaces volume of IEEE Std 1003.1-2001 specifies another mutex attribute
6675           called *type*. The *type* attribute allows applications to specify the behavior of mutex locking
6676           operations in situations where POSIX.1c behavior is undefined. The OSF DCE threads
6677           implementation, based on Draft 4 of POSIX.1c, specified a similar attribute. Note that the
6678           names of the attributes have changed somewhat from the OSF DCE threads implementation.

6679           The System Interfaces volume of IEEE Std 1003.1-2001 also extends the specification of the
6680           following POSIX.1c functions which manipulate mutexes:

6681              *pthread_mutex_lock*( )
6682              *pthread_mutex_trylock*( )
6683              *pthread_mutex_unlock*( )

6684        to take account of the new mutex attribute type and to specify behavior which was declared
6685        as undefined in POSIX.1c. How a calling thread acquires or releases a mutex now depends
6686        upon the mutex *type* attribute.

6687        The *type* attribute can have the following values:

6688        PTHREAD_MUTEX_NORMAL
6689              Basic mutex with no specific error checking built in. Does not report a deadlock error.

6690        PTHREAD_MUTEX_RECURSIVE
6691              Allows any thread to recursively lock a mutex. The mutex must be unlocked an equal
6692              number of times to release the mutex.

6693        PTHREAD_MUTEX_ERRORCHECK
6694              Detects and reports simple usage errors; that is, an attempt to unlock a mutex that is not
6695              locked by the calling thread or that is not locked at all, or an attempt to relock a mutex
6696              the thread already owns.

6697        PTHREAD_MUTEX_DEFAULT
6698              The default mutex type. May be mapped to any of the above mutex types or may be an
6699              implementation-defined type.

6700        *Normal* mutexes do not detect deadlock conditions; for example, a thread will hang if it tries
6701        to relock a normal mutex that it already owns. Attempting to unlock a mutex locked by
6702        another thread, or unlocking an unlocked mutex, results in undefined behavior. Normal
6703        mutexes will usually be the fastest type of mutex available on a platform but provide the
6704        least error checking.

6705        *Recursive* mutexes are useful for converting old code where it is difficult to establish clear
6706        boundaries of synchronization. A thread can relock a recursive mutex without first unlocking

6707   it. The relocking deadlock which can occur with normal mutexes cannot occur with this type
6708   of mutex. However, multiple locks of a recursive mutex require the same number of unlocks
6709   to release the mutex before another thread can acquire the mutex. Furthermore, this type of
6710   mutex maintains the concept of an owner. Thus, a thread attempting to unlock a recursive
6711   mutex which another thread has locked returns with an error. A thread attempting to unlock
6712   a recursive mutex that is not locked returns with an error. Never use a recursive mutex with
6713   condition variables because the implicit unlock performed by *pthread_cond_wait*( ) or
6714   *pthread_cond_timedwait*( ) will not actually release the mutex if it had been locked multiple
6715   times.

6716   *Errorcheck* mutexes provide error checking and are useful primarily as a debugging aid. A
6717   thread attempting to relock an errorcheck mutex without first unlocking it returns with an
6718   error. Again, this type of mutex maintains the concept of an owner. Thus, a thread
6719   attempting to unlock an errorcheck mutex which another thread has locked returns with an
6720   error. A thread attempting to unlock an errorcheck mutex that is not locked also returns with
6721   an error. It should be noted that errorcheck mutexes will almost always be much slower than
6722   normal mutexes due to the extra state checks performed.

6723   The default mutex type provides implementation-defined error checking. The default mutex
6724   may be mapped to one of the other defined types or may be something entirely different.
6725   This enables each vendor to provide the mutex semantics which the vendor feels will be
6726   most useful to their target users. Most vendors will probably choose to make normal
6727   mutexes the default so as to give applications the benefit of the fastest type of mutexes
6728   available on their platform. Check your implementation's documentation.

6729   An application developer can use any of the mutex types almost interchangeably as long as
6730   the application does not depend upon the implementation detecting (or failing to detect) any
6731   particular errors. Note that a recursive mutex can be used with condition variable waits as
6732   long as the application never recursively locks the mutex.

6733   Two functions are provided for manipulating the *type* attribute of a mutex attributes object.
6734   This attribute is set or returned in the *type* parameter of these functions. The
6735   *pthread_mutexattr_settype*( ) function is used to set a specific type value while
6736   *pthread_mutexattr_gettype*( ) is used to return the type of the mutex. Setting the *type* attribute
6737   of a mutex attributes object affects only mutexes initialized using that mutex attributes
6738   object. Changing the *type* attribute does not affect mutexes previously initialized using that
6739   mutex attributes object.

6740   • Read-Write Locks and Attributes

6741   The read-write locks introduced have been harmonized with those in IEEE Std 1003.1j-2000;
6742   see also Section B.2.9.6 (on page 175).

6743   Read-write locks (also known as reader-writer locks) allow a thread to exclusively lock some
6744   shared data while updating that data, or allow any number of threads to have simultaneous
6745   read-only access to the data.

6746   Unlike a mutex, a read-write lock distinguishes between reading data and writing data. A
6747   mutex excludes all other threads. A read-write lock allows other threads access to the data,
6748   providing no thread is modifying the data. Thus, a read-write lock is less primitive than
6749   either a mutex-condition variable pair or a semaphore.

6750   Application developers should consider using a read-write lock rather than a mutex to
6751   protect data that is frequently referenced but seldom modified. Most threads (readers) will be
6752   able to read the data without waiting and will only have to block when some other thread (a
6753   writer) is in the process of modifying the data. Conversely a thread that wants to change the
6754   data is forced to wait until there are no readers. This type of lock is often used to facilitate

6755    parallel access to data on multi-processor platforms or to avoid context switches on single
6756    processor platforms where multiple threads access the same data.

6757    If a read-write lock becomes unlocked and there are multiple threads waiting to acquire the
6758    write lock, the implementation's scheduling policy determines which thread acquires the
6759    read-write lock for writing. If there are multiple threads blocked on a read-write lock for both
6760    read locks and write locks, it is unspecified whether the readers or a writer acquire the lock
6761    first. However, for performance reasons, implementations often favor writers over readers to
6762    avoid potential writer starvation.

6763    A read-write lock object is an implementation-defined opaque object of type
6764    **pthread_rwlock_t** as defined in <**pthread.h**>. There are two different sorts of locks
6765    associated with a read-write lock: a read lock and a write lock.

6766    The *pthread_rwlockattr_init*() function initializes a read-write lock attributes object with the
6767    default value for all the attributes defined in the implementation. After a read-write lock
6768    attributes object has been used to initialize one or more read-write locks, changes to the
6769    read-write lock attributes object, including destruction, do not affect previously initialized
6770    read-write locks.

6771    Implementations must provide at least the read-write lock attribute *process-shared*. This
6772    attribute can have the following values:

6773    PTHREAD_PROCESS_SHARED
6774        Any thread of any process that has access to the memory where the read-write lock
6775        resides can manipulate the read-write lock.

6776    PTHREAD_PROCESS_PRIVATE
6777        Only threads created within the same process as the thread that initialized the read-
6778        write lock can manipulate the read-write lock. This is the default value.

6779    The *pthread_rwlockattr_setpshared*() function is used to set the *process-shared* attribute of an
6780    initialized read-write lock attributes object while the function *pthread_rwlockattr_getpshared*()
6781    obtains the current value of the *process-shared* attribute.

6782    A read-write lock attributes object is destroyed using the *pthread_rwlockattr_destroy*()
6783    function. The effect of subsequent use of the read-write lock attributes object is undefined.

6784    A thread creates a read-write lock using the *pthread_rwlock_init*() function. The attributes of
6785    the read-write lock can be specified by the application developer; otherwise, the default
6786    implementation-defined read-write lock attributes are used if the pointer to the read-write
6787    lock attributes object is NULL. In cases where the default attributes are appropriate, the
6788    PTHREAD_RWLOCK_INITIALIZER macro can be used to initialize statically allocated
6789    read-write locks.

6790    A thread which wants to apply a read lock to the read-write lock can use either
6791    *pthread_rwlock_rdlock*() or *pthread_rwlock_tryrdlock*(). If *pthread_rwlock_rdlock*() is used, the
6792    thread acquires a read lock if a writer does not hold the write lock and there are no writers
6793    blocked on the write lock. If a read lock is not acquired, the calling thread blocks until it can
6794    acquire a lock. However, if *pthread_rwlock_tryrdlock*() is used, the function returns
6795    immediately with the error [EBUSY] if any thread holds a write lock or there are blocked
6796    writers waiting for the write lock.

6797    A thread which wants to apply a write lock to the read-write lock can use either of two
6798    functions: *pthread_rwlock_wrlock*() or *pthread_rwlock_trywrlock*(). If *pthread_rwlock_wrlock*()
6799    is used, the thread acquires the write lock if no other reader or writer threads hold the read-
6800    write lock. If the write lock is not acquired, the thread blocks until it can acquire the write
6801    lock. However, if *pthread_rwlock_trywrlock*() is used, the function returns immediately with

6802          the error [EBUSY] if any thread is holding either a read or a write lock.

6803          The *pthread_rwlock_unlock*( ) function is used to unlock a read-write lock object held by the
6804          calling thread. Results are undefined if the read-write lock is not held by the calling thread. If
6805          there are other read locks currently held on the read-write lock object, the read-write lock
6806          object remains in the read locked state but without the current thread as one of its owners.  If
6807          this function releases the last read lock for this read-write lock object, the read-write lock
6808          object is put in the unlocked read state. If this function is called to release a write lock for this
6809          read-write lock object, the read-write lock object is put in the unlocked state.

6810      • Thread Concurrency Level

6811          On threads implementations that multiplex user threads onto a smaller set of kernel
6812          execution entities, the system attempts to create a reasonable number of kernel execution
6813          entities for the application upon application startup.

6814          On some implementations, these kernel entities are retained by user threads that block in the
6815          kernel. Other implementations do not *timeslice* user threads so that multiple compute-bound
6816          user threads can share a kernel thread. On such implementations, some applications may use
6817          up all the available kernel execution entities before their user-space threads are used up. The
6818          process may be left with user threads capable of doing work for the application but with no
6819          way to schedule them.

6820          The *pthread_setconcurrency*( ) function enables an application to request more kernel entities;
6821          that is, specify a desired concurrency level. However, this function merely provides a hint to
6822          the implementation. The implementation is free to ignore this request or to provide some
6823          other number of kernel entities.  If an implementation does not multiplex user threads onto a
6824          smaller number of kernel execution entities, the *pthread_setconcurrency*( ) function has no
6825          effect.

6826          The *pthread_setconcurrency*( ) function may also have an effect on implementations where the
6827          kernel mode and user mode schedulers cooperate to ensure that ready user threads are not
6828          prevented from running by other threads blocked in the kernel.

6829          The *pthread_getconcurrency*( ) function always returns the value set by a previous call to
6830          *pthread_setconcurrency*( ).  However, if *pthread_setconcurrency*( ) was not previously called, this
6831          function returns zero to indicate that the threads implementation is maintaining the
6832          concurrency level.

6833      • Thread Stack Guard Size

6834          DCE threads introduced the concept of a ''thread stack guard size''.  Most thread
6835          implementations add a region of protected memory to a thread's stack, commonly known as
6836          a ''guard region'', as a safety measure to prevent stack pointer overflow in one thread from
6837          corrupting the contents of another thread's stack. The default size of the guard regions
6838          attribute is {PAGESIZE} bytes and is implementation-defined.

6839          Some application developers may wish to change the stack guard size.  When an application
6840          creates a large number of threads, the extra page allocated for each stack may strain system
6841          resources. In addition to the extra page of memory, the kernel's memory manager has to keep
6842          track of the different protections on adjoining pages. When this is a problem, the application
6843          developer may request a guard size of 0 bytes to conserve system resources by eliminating
6844          stack overflow protection.

6845          Conversely an application that allocates large data structures such as arrays on the stack may
6846          wish to increase the default guard size in order to detect stack overflow. If a thread allocates
6847          two pages for a data array, a single guard page provides little protection against thread stack
6848          overflows since the thread can corrupt adjoining memory beyond the guard page.

6849          The System Interfaces volume of IEEE Std 1003.1-2001 defines a new attribute of a thread
6850          attributes object; that is, the *guardsize* attribute which allows applications to specify the size
6851          of the guard region of a thread's stack.

6852          Two functions are provided for manipulating a thread's stack guard size. The
6853          *pthread_attr_setguardsize*() function sets the thread *guardsize* attribute, and the
6854          *pthread_attr_getguardsize*() function retrieves the current value.

6855          An implementation may round up the requested guard size to a multiple of the configurable
6856          system variable {PAGESIZE}. In this case, *pthread_attr_getguardsize*() returns the guard size
6857          specified by the previous *pthread_attr_setguardsize*() function call and not the rounded up
6858          value.

6859          If an application is managing its own thread stacks using the *stackaddr* attribute, the *guardsize*
6860          attribute is ignored and no stack overflow protection is provided. In this case, it is the
6861          responsibility of the application to manage stack overflow along with stack allocation.

6862     • Parallel I/O

6863          Suppose two or more threads independently issue read requests on the same file. To read
6864          specific data from a file, a thread must first call *lseek*() to seek to the proper offset in the file,
6865          and then call *read*() to retrieve the required data. If more than one thread does this at the
6866          same time, the first thread may complete its seek call, but before it gets a chance to issue its
6867          read call a second thread may complete its seek call, resulting in the first thread accessing
6868          incorrect data when it issues its read call. One workaround is to lock the file descriptor while
6869          seeking and reading or writing, but this reduces parallelism and adds overhead.

6870          Instead, the System Interfaces volume of IEEE Std 1003.1-2001 provides two functions to
6871          make seek/read and seek/write operations atomic. The file descriptor's current offset is
6872          unchanged, thus allowing multiple read and write operations to proceed in parallel. This
6873          improves the I/O performance of threaded applications. The *pread*() function is used to do
6874          an atomic read of data from a file into a buffer. Conversely, the *pwrite*() function does an
6875          atomic write of data from a buffer to a file.

6876  *B.2.9.1    Thread-Safety*

6877          All functions required by IEEE Std 1003.1-2001 need to be thread-safe.  Implementations have to
6878          provide internal synchronization when necessary in order to achieve this goal. In certain
6879          cases—for example, most floating-point implementations—context switch code may have to
6880          manage the writable shared state.

6881          While a read from a pipe of {PIPE_MAX}*2 bytes may not generate a single atomic and thread-
6882          safe stream of bytes, it should generate ''several'' (individually atomic) thread-safe streams of
6883          bytes. Similarly, while reading from a terminal device may not generate a single atomic and
6884          thread-safe stream of bytes, it should generate some finite number of (individually atomic) and
6885          thread-safe streams of bytes. That is, concurrent calls to read for a pipe, FIFO, or terminal device
6886          are not allowed to result in corrupting the stream of bytes or other internal data. However,
6887          *read*(), in these cases, is not required to return a single contiguous and atomic stream of bytes.

6888          It is not required that all functions provided by IEEE Std 1003.1-2001 be either async-cancel-safe
6889          or async-signal-safe.

6890          As it turns out, some functions are inherently not thread-safe; that is, their interface
6891          specifications preclude reentrancy. For example, some functions (such as *asctime*()) return a
6892          pointer to a result stored in memory space allocated by the function on a per-process basis. Such
6893          a function is not thread-safe, because its result can be overwritten by successive invocations.
6894          Other functions, while not inherently non-thread-safe, may be implemented in ways that lead to

6895 them not being thread-safe. For example, some functions (such as *rand*( )) store state information
6896 (such as a seed value, which survives multiple function invocations) in memory space allocated
6897 by the function on a per-process basis. The implementation of such a function is not thread-safe
6898 if the implementation fails to synchronize invocations of the function and thus fails to protect
6899 the state information. The problem is that when the state information is not protected,
6900 concurrent invocations can interfere with one another (for example, applications using *rand*( )
6901 may see the same seed value).

6902 *Thread-Safety and Locking of Existing Functions*

6903 Originally, POSIX.1 was not designed to work in a multi-threaded environment, and some
6904 implementations of some existing functions will not work properly when executed concurrently.
6905 To provide routines that will work correctly in an environment with threads (''thread-safe''), two
6906 problems need to be solved:

6907     1.  Routines that maintain or return pointers to static areas internal to the routine (which may
6908         now be shared) need to be modified. The routines *ttyname*( ) and *localtime*( ) are examples.

6909     2.  Routines that access data space shared by more than one thread need to be modified. The
6910         *malloc*( ) function and the *stdio* family routines are examples.

6911 There are a variety of constraints on these changes. The first is compatibility with the existing
6912 versions of these functions—non-thread-safe functions will continue to be in use for some time,
6913 as the original interfaces are used by existing code.  Another is that the new thread-safe versions
6914 of these functions represent as small a change as possible over the familiar interfaces provided
6915 by the existing non-thread-safe versions. The new interfaces should be independent of any
6916 particular threads implementation. In particular, they should be thread-safe without depending
6917 on explicit thread-specific memory. Finally, there should be minimal performance penalty due to
6918 the changes made to the functions.

6919 It is intended that the list of functions from POSIX.1 that cannot be made thread-safe and for
6920 which corrected versions are provided be complete.

6921 *Thread-Safety and Locking Solutions*

6922 Many of the POSIX.1 functions were thread-safe and did not change at all. However, some
6923 functions (for example, the math functions typically found in **libm**) are not thread-safe because
6924 of writable shared global state. For instance, in IEEE Std 754-1985 floating-point
6925 implementations, the computation modes and flags are global and shared.

6926 Some functions are not thread-safe because a particular implementation is not reentrant,
6927 typically because of a non-essential use of static storage. These require only a new
6928 implementation.

6929 Thread-safe libraries are useful in a wide range of parallel (and asynchronous) programming
6930 environments, not just within pthreads. In order to be used outside the context of pthreads,
6931 however, such libraries still have to use some synchronization method. These could either be
6932 independent of the pthread synchronization operations, or they could be a subset of the pthread
6933 interfaces. Either method results in thread-safe library implementations that can be used without
6934 the rest of pthreads.

6935 Some functions, such as the *stdio* family interface and dynamic memory allocation functions
6936 such as *malloc*( ), are inter-dependent routines that share resources (for example, buffers) across
6937 related calls. These require synchronization to work correctly, but they do not require any
6938 change to their external (user-visible) interfaces.

6939 In some cases, such as *getc*( ) and *putc*( ), adding synchronization is likely to create an
6940 unacceptable performance impact. In this case, slower thread-safe synchronized functions are to

6941 be provided, but the original, faster (but unsafe) functions (which may be implemented as
6942 macros) are retained under new names. Some additional special-purpose synchronization
6943 facilities are necessary for these macros to be usable in multi-threaded programs. This also
6944 requires changes in **<stdio.h>**.

6945 The other common reason that functions are unsafe is that they return a pointer to static storage,
6946 making the functions non-thread-safe. This has to be changed, and there are three natural
6947 choices:

1. Return a pointer to thread-specific storage

6949 This could incur a severe performance penalty on those architectures with a costly
6950 implementation of the thread-specific data interface.

6951 A variation on this technique is to use *malloc*( ) to allocate storage for the function output
6952 and return a pointer to this storage. This technique may also have an undesirable
6953 performance impact, however, and a simplistic implementation requires that the user
6954 program explicitly free the storage object when it is no longer needed. This technique is
6955 used by some existing POSIX.1 functions. With careful implementation for infrequently
6956 used functions, there may be little or no performance or storage penalty, and the
6957 maintenance of already-standardized interfaces is a significant benefit.

2. Return the actual value computed by the function

6959 This technique can only be used with functions that return pointers to structures—routines
6960 that return character strings would have to wrap their output in an enclosing structure in
6961 order to return the output on the stack. There is also a negative performance impact
6962 inherent in this solution in that the output value has to be copied twice before it can be
6963 used by the calling function: once from the called routine's local buffers to the top of the
6964 stack, then from the top of the stack to the assignment target. Finally, many older
6965 compilers cannot support this technique due to a historical tendency to use internal static
6966 buffers to deliver the results of structure-valued functions.

3. Have the caller pass the address of a buffer to contain the computed value

6968 The only disadvantage of this approach is that extra arguments have to be provided by the
6969 calling program. It represents the most efficient solution to the problem, however, and,
6970 unlike the *malloc*( ) technique, it is semantically clear.

6971 There are some routines (often groups of related routines) whose interfaces are inherently non-
6972 thread-safe because they communicate across multiple function invocations by means of static
6973 memory locations. The solution is to redesign the calls so that they are thread-safe, typically by
6974 passing the needed data as extra parameters. Unfortunately, this may require major changes to
6975 the interface as well.

6976 A floating-point implementation using IEEE Std 754-1985 is a case in point. A less problematic
6977 example is the *rand48* family of pseudo-random number generators. The functions *getgrgid*( ),
6978 *getgrnam*( ), *getpwnam*( ), and *getpwuid*( ) are another such case.

6979 The problems with *errno* are discussed in **Alternative Solutions for Per-Thread errno** (on page
6980 92).

6981 Some functions can be thread-safe or not, depending on their arguments. These include the
6982 *tmpnam*( ) and *ctermid*( ) functions. These functions have pointers to character strings as
6983 arguments. If the pointers are not NULL, the functions store their results in the character string;
6984 however, if the pointers are NULL, the functions store their results in an area that may be static
6985 and thus subject to overwriting by successive calls. These should only be called by multi-thread
6986 applications when their arguments are non-NULL.

*Asynchronous Safety and Thread-Safety*

6988    A floating-point implementation has many modes that effect rounding and other aspects of
6989    computation. Functions in some math library implementations may change the computation
6990    modes for the duration of a function call. If such a function call is interrupted by a signal or
6991    cancellation, the floating-point state is not required to be protected.                                    |

6992    There is a significant cost to make floating-point operations async-cancel-safe or async-signal-
6993    safe; accordingly, neither form of async safety is required.

6994    *Functions Returning Pointers to Static Storage*

6995    For those functions that are not thread-safe because they return values in fixed size statically
6996    allocated structures, alternate ''_r'' forms are provided that pass a pointer to an explicit result
6997    structure.  Those that return pointers into library-allocated buffers have forms provided with
6998    explicit buffer and length parameters.

6999    For functions that return pointers to library-allocated buffers, it makes sense to provide ''_r''
7000    versions that allow the application control over allocation of the storage in which results are
7001    returned.  This allows the state used by these functions to be managed on an application-specific
7002    basis, supporting per-thread, per-process, or other application-specific sharing relationships.

7003    Early proposals had provided ''_r'' versions for functions that returned pointers to variable-size
7004    buffers without providing a means for determining the required buffer size. This would have
7005    made using such functions exceedingly clumsy, potentially requiring iteratively calling them
7006    with increasingly larger guesses for the amount of storage required. Hence, *sysconf*( ) variables
7007    have been provided for such functions that return the maximum required buffer size.

7008    Thus, the rule that has been followed by IEEE Std 1003.1-2001 when adapting single-threaded
7009    non-thread-safe functions is as follows: all functions returning pointers to library-allocated
7010    storage should have ''_r'' versions provided, allowing the application control over the storage
7011    allocation. Those with variable-sized return values accept both a buffer address and a length
7012    parameter. The *sysconf*( ) variables are provided to supply the appropriate buffer sizes when
7013    required. Implementors are encouraged to apply the same rule when adapting their own existing
7014    functions to a pthreads environment.

7015  *B.2.9.2    Thread IDs*

7016    Separate applications should communicate through well-defined interfaces and should not
7017    depend on each other's implementation. For example, if a programmer decides to rewrite the *sort*
7018    utility using multiple threads, it should be easy to do this so that the interface to the *sort* utility
7019    does not change. Consider that if the user causes SIGINT to be generated while the *sort* utility is
7020    running, keeping the same interface means that the entire *sort* utility is killed, not just one of its
7021    threads. As another example, consider a realtime application that manages a reactor. Such an
7022    application may wish to allow other applications to control the priority at which it watches the
7023    control rods. One technique to accomplish this is to write the ID of the thread watching the
7024    control rods into a file and allow other programs to change the priority of that thread as they see
7025    fit. A simpler technique is to have the reactor process accept IPCs (Interprocess Communication
7026    messages) from other processes, telling it at a semantic level what priority the program should
7027    assign to watching the control rods. This allows the programmer greater flexibility in the
7028    implementation. For example, the programmer can change the implementation from having one
7029    thread per rod to having one thread watching all of the rods without changing the interface.
7030    Having threads live inside the process means that the implementation of a process is invisible to
7031    outside processes (excepting debuggers and system management tools).

7032    Threads do not provide a protection boundary. Every thread model allows threads to share
7033    memory with other threads and encourages this sharing to be widespread. This means that one

7034      thread can wipe out memory that is needed for the correct functioning of other threads that are
7035      sharing its memory. Consequently, providing each thread with its own user and/or group IDs
7036      would not provide a protection boundary between threads sharing memory.

7037  *B.2.9.3*    *Thread Mutexes*

7038      There is no additional rationale provided for this section.

7039  *B.2.9.4*    *Thread Scheduling*

7040      • Scheduling Implementation Models

7041      The following scheduling implementation models are presented in terms of threads and
7042      ''kernel entities''. This is to simplify exposition of the models, and it does not imply that an
7043      implementation actually has an identifiable ''kernel entity''.

7044      A kernel entity is not defined beyond the fact that it has scheduling attributes that are used to
7045      resolve contention with other kernel entities for execution resources. A kernel entity may be
7046      thought of as an envelope that holds a thread or a separate kernel thread. It is not a
7047      conventional process, although it shares with the process the attribute that it has a single
7048      thread of control; it does not necessarily imply an address space, open files, and so on. It is
7049      better thought of as a primitive facility upon which conventional processes and threads may
7050      be constructed.

7051      — System Thread Scheduling Model

7052          This model consists of one thread per kernel entity. The kernel entity is solely responsible
7053          for scheduling thread execution on one or more processors. This model schedules all
7054          threads against all other threads in the system using the scheduling attributes of the
7055          thread.

7056      — Process Scheduling Model

7057          A generalized process scheduling model consists of two levels of scheduling. A threads
7058          library creates a pool of kernel entities, as required, and schedules threads to run on them
7059          using the scheduling attributes of the threads. Typically, the size of the pool is a function
7060          of the simultaneously runnable threads, not the total number of threads. The kernel then
7061          schedules the kernel entities onto processors according to their scheduling attributes,
7062          which are managed by the threads library. This set model potentially allows a wide range
7063          of mappings between threads and kernel entities.

7064      • System and Process Scheduling Model Performance

7065      There are a number of important implications on the performance of applications using these
7066      scheduling models. The process scheduling model potentially provides lower overhead for
7067      making scheduling decisions, since there is no need to access kernel-level information or
7068      functions and the set of schedulable entities is smaller (only the threads within the process).

7069      On the other hand, since the kernel is also making scheduling decisions regarding the system
7070      resources under its control (for example, CPU(s), I/O devices, memory), decisions that do
7071      not take thread scheduling parameters into account can result in unspecified delays for
7072      realtime application threads, causing them to miss maximum response time limits.

7073      • Rate Monotonic Scheduling

7074      Rate monotonic scheduling was considered, but rejected for standardization in the context of
7075      pthreads. A sporadic server policy is included.

7076 • Scheduling Options

7077 In IEEE Std 1003.1-2001, the basic thread scheduling functions are defined under the Threads
7078 option, so that they are required of all threads implementations. However, there are no
7079 specific scheduling policies required by this option to allow for conforming thread
7080 implementations that are not targeted to realtime applications.

7081 Specific standard scheduling policies are defined to be under the Thread Execution
7082 Scheduling option, and they are specifically designed to support realtime applications by
7083 providing predictable resource-sharing sequences. The name of this option was chosen to
7084 emphasize that this functionality is defined as appropriate for realtime applications that
7085 require simple priority-based scheduling.

7086 It is recognized that these policies are not necessarily satisfactory for some multi-processor
7087 implementations, and work is ongoing to address a wider range of scheduling behaviors. The
7088 interfaces have been chosen to create abundant opportunity for future scheduling policies to
7089 be implemented and standardized based on this interface. In order to standardize a new
7090 scheduling policy, all that is required (from the standpoint of thread scheduling attributes) is
7091 to define a new policy name, new members of the thread attributes object, and functions to
7092 set these members when the scheduling policy is equal to the new value.

7093 **Scheduling Contention Scope**

7094 In order to accommodate the requirement for realtime response, each thread has a scheduling
7095 contention scope attribute. Threads with a system scheduling contention scope have to be
7096 scheduled with respect to all other threads in the system. These threads are usually bound to a
7097 single kernel entity that reflects their scheduling attributes and are directly scheduled by the
7098 kernel.

7099 Threads with a process scheduling contention scope need be scheduled only with respect to the
7100 other threads in the process. These threads may be scheduled within the process onto a pool of
7101 kernel entities. The implementation is also free to bind these threads directly to kernel entities
7102 and let them be scheduled by the kernel. Process scheduling contention scope allows the
7103 implementation the most flexibility and is the default if both contention scopes are supported
7104 and none is specified.

7105 Thus, the choice by implementors to provide one or the other (or both) of these scheduling
7106 models is driven by the need of their supported application domains for worst-case (that is,
7107 realtime) response, or average-case (non-realtime) response.

7108 **Scheduling Allocation Domain**

7109 The SCHED_FIFO and SCHED_RR scheduling policies take on different characteristics on a
7110 multi-processor. Other scheduling policies are also subject to changed behavior when executed
7111 on a multi-processor. The concept of scheduling allocation domain determines the set of
7112 processors on which the threads of an application may run. By considering the application's
7113 processor scheduling allocation domain for its threads, scheduling policies can be defined in
7114 terms of their behavior for varying processor scheduling allocation domain values. It is
7115 conceivable that not all scheduling allocation domain sizes make sense for all scheduling
7116 policies on all implementations. The concept of scheduling allocation domain, however, is a
7117 useful tool for the description of multi-processor scheduling policies.

7118 The ''process control'' approach to scheduling obtains significant performance advantages from
7119 dynamic scheduling allocation domain sizes when it is applicable.

7120 Non-Uniform Memory Access (NUMA) multi-processors may use a system scheduling structure
7121 that involves reassignment of threads among scheduling allocation domains. In NUMA

7122    machines, a natural model of scheduling is to match scheduling allocation domains to clusters of
7123    processors. Load balancing in such an environment requires changing the scheduling allocation
7124    domain to which a thread is assigned.

### Scheduling Documentation

7126    Implementation-provided scheduling policies need to be completely documented in order to be
7127    useful. This documentation includes a description of the attributes required for the policy, the
7128    scheduling interaction of threads running under this policy and all other supported policies, and
7129    the effects of all possible values for processor scheduling allocation domain. Note that for the
7130    implementor wishing to be minimally-compliant, it is (minimally) acceptable to define the
7131    behavior as undefined.

### Scheduling Contention Scope Attribute

7133    The scheduling contention scope defines how threads compete for resources. Within
7134    IEEE Std 1003.1-2001, scheduling contention scope is used to describe only how threads are
7135    scheduled in relation to one another in the system. That is, either they are scheduled against all
7136    other threads in the system (''system scope'') or only against those threads in the process
7137    (''process scope''). In fact, scheduling contention scope may apply to additional resources,
7138    including virtual timers and profiling, which are not currently considered by
7139    IEEE Std 1003.1-2001.

### Mixed Scopes

7141    If only one scheduling contention scope is supported, the scheduling decision is straightforward.
7142    To perform the processor scheduling decision in a mixed scope environment, it is necessary to
7143    map the scheduling attributes of the thread with process-wide contention scope to the same
7144    attribute space as the thread with system-wide contention scope.

7145    Since a conforming implementation has to support one and may support both scopes, it is useful
7146    to discuss the effects of such choices with respect to example applications. If an implementation
7147    supports both scopes, mixing scopes provides a means of better managing system-level (that is,
7148    kernel-level) and library-level resources. In general, threads with system scope will require the
7149    resources of a separate kernel entity in order to guarantee the scheduling semantics. On the
7150    other hand, threads with process scope can share the resources of a kernel entity while
7151    maintaining the scheduling semantics.

7152    The application is free to create threads with dedicated kernel resources, and other threads that
7153    multiplex kernel resources. Consider the example of a window server. The server allocates two
7154    threads per widget: one thread manages the widget user interface (including drawing), while the
7155    other thread takes any required application action. This allows the widget to be ''active'' while
7156    the application is computing. A screen image may be built from thousands of widgets. If each of
7157    these threads had been created with system scope, then most of the kernel-level resources might
7158    be wasted, since only a few widgets are active at any one time. In addition, mixed scope is
7159    particularly useful in a window server where one thread with high priority and system scope
7160    handles the mouse so that it tracks well. As another example, consider a database server. For
7161    each of the hundreds or thousands of clients supported by a large server, an equivalent number
7162    of threads will have to be created. If each of these threads were system scope, the consequences
7163    would be the same as for the window server example above. However, the server could be
7164    constructed so that actual retrieval of data is done by several dedicated threads. Dedicated
7165    threads that do work for all clients frequently justify the added expense of system scope. If it
7166    were not permissible to mix system and process threads in the same process, this type of
7167    solution would not be possible.

**Dynamic Thread Scheduling Parameters Access**

In many time-constrained applications, there is no need to change the scheduling attributes dynamically during thread or process execution, since the general use of these attributes is to reflect directly the time constraints of the application. Since these time constraints are generally imposed to meet higher-level system requirements, such as accuracy or availability, they frequently should remain unchanged during application execution.

However, there are important situations in which the scheduling attributes should be changed. Generally, this will occur when external environmental conditions exist in which the time constraints change. Consider, for example, a space vehicle major mode change, such as the change from ascent to descent mode, or the change from the space environment to the atmospheric environment. In such cases, the frequency with which many of the sensors or actuators need to be read or written will change, which will necessitate a priority change. In other cases, even the existence of a time constraint might be temporary, necessitating not just a priority change, but also a policy change for ongoing threads or processes. For this reason, it is critical that the interface should provide functions to change the scheduling parameters dynamically, but, as with many of the other realtime functions, it is important that applications use them properly to avoid the possibility of unnecessarily degrading performance.

In providing functions for dynamically changing the scheduling behavior of threads, there were two options: provide functions to get and set the individual scheduling parameters of threads, or provide a single interface to get and set all the scheduling parameters for a given thread simultaneously. Both approaches have merit. Access functions for individual parameters allow simpler control of thread scheduling for simple thread scheduling parameters. However, a single function for setting all the parameters for a given scheduling policy is required when first setting that scheduling policy. Since the single all-encompassing functions are required, it was decided to leave the interface as minimal as possible. Note that simpler functions (such as *pthread_setprio*( ) for threads running under the priority-based schedulers) can be easily defined in terms of the all-encompassing functions.

If the *pthread_setschedparam*( ) function executes successfully, it will have set all of the scheduling parameter values indicated in *param*; otherwise, none of the scheduling parameters will have been modified. This is necessary to ensure that the scheduling of this and all other threads continues to be consistent in the presence of an erroneous scheduling parameter.

The [EPERM] error value is included in the list of possible *pthread_setschedparam*( ) error returns as a reflection of the fact that the ability to change scheduling parameters increases risks to the implementation and application performance if the scheduling parameters are changed improperly. For this reason, and based on some existing practice, it was felt that some implementations would probably choose to define specific permissions for changing either a thread's own or another thread's scheduling parameters. IEEE Std 1003.1-2001 does not include portable methods for setting or retrieving permissions, so any such use of permissions is completely unspecified.

**Mutex Initialization Scheduling Attributes**

In a priority-driven environment, a direct use of traditional primitives like mutexes and condition variables can lead to unbounded priority inversion, where a higher priority thread can be blocked by a lower priority thread, or set of threads, for an unbounded duration of time. As a result, it becomes impossible to guarantee thread deadlines. Priority inversion can be bounded and minimized by the use of priority inheritance protocols. This allows thread deadlines to be guaranteed even in the presence of synchronization requirements.

Two useful but simple members of the family of priority inheritance protocols are the basic priority inheritance protocol and the priority ceiling protocol emulation. Under the Basic Priority

7216 Inheritance protocol (governed by the Thread Priority Inheritance option), a thread that is
7217 blocking higher priority threads executes at the priority of the highest priority thread that it
7218 blocks. This simple mechanism allows priority inversion to be bounded by the duration of
7219 critical sections and makes timing analysis possible.

7220 Under the Priority Ceiling Protocol Emulation protocol (governed by the Thread Priority
7221 Protection option), each mutex has a priority ceiling, usually defined as the priority of the
7222 highest priority thread that can lock the mutex. When a thread is executing inside critical
7223 sections, its priority is unconditionally increased to the highest of the priority ceilings of all the
7224 mutexes owned by the thread. This protocol has two very desirable properties in uni-processor
7225 systems. First, a thread can be blocked by a lower priority thread for at most the duration of one
7226 single critical section. Furthermore, when the protocol is correctly used in a single processor, and
7227 if threads do not become blocked while owning mutexes, mutual deadlocks are prevented.

7228 The priority ceiling emulation can be extended to multiple processor environments, in which
7229 case the values of the priority ceilings will be assigned depending on the kind of mutex that is
7230 being used: local to only one processor, or global, shared by several processors. Local priority
7231 ceilings will be assigned the usual way, equal to the priority of the highest priority thread that
7232 may lock that mutex. Global priority ceilings will usually be assigned a priority level higher than
7233 all the priorities assigned to any of the threads that reside in the involved processors to avoid the
7234 effect called remote blocking.

7235 **Change the Priority Ceiling of a Mutex**

7236 In order for the priority protect protocol to exhibit its desired properties of bounding priority
7237 inversion and avoidance of deadlock, it is critical that the ceiling priority of a mutex be the same
7238 as the priority of the highest thread that can ever hold it, or higher. Thus, if the priorities of the
7239 threads using such mutexes never change dynamically, there is no need ever to change the
7240 priority ceiling of a mutex.

7241 However, if a major system mode change results in an altered response time requirement for one
7242 or more application threads, their priority has to change to reflect it. It will occasionally be the
7243 case that the priority ceilings of mutexes held also need to change. While changing priority
7244 ceilings should generally be avoided, it is important that IEEE Std 1003.1-2001 provide these
7245 interfaces for those cases in which it is necessary.

7246 *B.2.9.5* *Thread Cancellation* |

7247 Many existing threads packages have facilities for canceling an operation or canceling a thread. |
7248 These facilities are used for implementing user requests (such as the CANCEL button in a
7249 window-based application), for implementing OR parallelism (for example, telling the other
7250 threads to stop working once one thread has found a forced mate in a parallel chess program), or
7251 for implementing the ABORT mechanism in Ada.

7252 POSIX programs traditionally have used the signal mechanism combined with either *longjmp*( )
7253 or polling to cancel operations. Many POSIX programmers have trouble using these facilities to
7254 solve their problems efficiently in a single-threaded process. With the introduction of threads,
7255 these solutions become even more difficult to use.

7256 The main issues with implementing a cancellation facility are specifying the operation to be |
7257 canceled, cleanly releasing any resources allocated to that operation, controlling when the target
7258 notices that it has been canceled, and defining the interaction between asynchronous signals and |
7259 cancellation. |

**Specifying the Operation to Cancel**

Consider a thread that calls through five distinct levels of program abstraction and then, inside the lowest-level abstraction, calls a function that suspends the thread. (An abstraction boundary is a layer at which the client of the abstraction sees only the service being provided and can remain ignorant of the implementation. Abstractions are often layered, each level of abstraction being a client of the lower-level abstraction and implementing a higher-level abstraction.) Depending on the semantics of each abstraction, one could imagine wanting to cancel only the call that causes suspension, only the bottom two levels, or the operation being done by the entire thread. Canceling operations at a finer grain than the entire thread is difficult because threads are active and they may be run in parallel on a multi-processor. By the time one thread can make a request to cancel an operation, the thread performing the operation may have completed that operation and gone on to start another operation whose cancellation is not desired. Thread IDs are not reused until the thread has exited, and either it was created with the *Attr detachstate* attribute set to PTHREAD_CREATE_DETACHED or the *pthread_join*( ) or *pthread_detach*( ) function has been called for that thread. Consequently, a thread cancellation will never be misdirected when the thread terminates. For these reasons, the canceling of operations is done at the granularity of the thread. Threads are designed to be inexpensive enough so that a separate thread may be created to perform each separately cancelable operation; for example, each possibly long running user request.

For cancellation to be used in existing code, cancellation scopes and handlers will have to be established for code that needs to release resources upon cancellation, so that it follows the programming discipline described in the text.

**A Special Signal Versus a Special Interface**

Two different mechanisms were considered for providing the cancellation interfaces. The first was to provide an interface to direct signals at a thread and then to define a special signal that had the required semantics. The other alternative was to use a special interface that delivered the correct semantics to the target thread.

The solution using signals produced a number of problems. It required the implementation to provide cancellation in terms of signals whereas a perfectly valid (and possibly more efficient) implementation could have both layered on a low-level set of primitives. There were so many exceptions to the special signal (it cannot be used with *kill*( ), no POSIX.1 interfaces can be used with it) that it was clearly not a valid signal. Its semantics on delivery were also completely different from any existing POSIX.1 signal. As such, a special interface that did not mandate the implementation and did not confuse the semantics of signals and cancellation was felt to be the better solution.

**Races Between Cancellation and Resuming Execution**

Due to the nature of cancellation, there is generally no synchronization between the thread requesting the cancellation of a blocked thread and events that may cause that thread to resume execution. For this reason, and because excess serialization hurts performance, when both an event that a thread is waiting for has occurred and a cancellation request has been made and cancellation is enabled, IEEE Std 1003.1-2001 explicitly allows the implementation to choose between returning from the blocking call or acting on the cancellation request.

7302     **Interaction of Cancellation with Asynchronous Signals**

7303     A typical use of cancellation is to acquire a lock on some resource and to establish a cancellation
7304     cleanup handler for releasing the resource when and if the thread is canceled.

7305     A correct and complete implementation of cancellation in the presence of asynchronous signals
7306     requires considerable care. An implementation has to push a cancellation cleanup handler on the
7307     cancellation cleanup stack while maintaining the integrity of the stack data structure. If an
7308     asynchronously-generated signal is posted to the thread during a stack operation, the signal
7309     handler cannot manipulate the cancellation cleanup stack. As a consequence, asynchronous
7310     signal handlers may not cancel threads or otherwise manipulate the cancellation state of a
7311     thread. Threads may, of course, be canceled by another thread that used a *sigwait*() function to
7312     wait synchronously for an asynchronous signal.

7313     In order for cancellation to function correctly, it is required that asynchronous signal handlers
7314     not change the cancellation state. This requires that some elements of existing practice, such as
7315     using *longjmp*() to exit from an asynchronous signal handler implicitly, be prohibited in cases
7316     where the integrity of the cancellation state of the interrupt thread cannot be ensured.

7317     **Thread Cancellation Overview**

7318     • Cancelability States

7319        The three possible cancelability states (disabled, deferred, and asynchronous) are encoded
7320        into two separate bits ((disable, enable) and (deferred, asynchronous)) to allow them to be
7321        changed and restored independently. For instance, short code sequences that will not block
7322        sometimes disable cancelability on entry and restore the previous state upon exit. Likewise,
7323        long or unbounded code sequences containing no convenient explicit cancellation points will
7324        sometimes set the cancelability type to asynchronous on entry and restore the previous value
7325        upon exit.

7326     • Cancellation Points

7327        Cancellation points are points inside of certain functions where a thread has to act on any
7328        pending cancellation request when cancelability is enabled, if the function would block. As
7329        with checking for signals, operations need only check for pending cancellation requests when
7330        the operation is about to block indefinitely.

7331        The idea was considered of allowing implementations to define whether blocking calls such
7332        as *read*() should be cancellation points. It was decided that it would adversely affect the
7333        design of conforming applications if blocking calls were not cancellation points because
7334        threads could be left blocked in an uncancelable state.

7335        There are several important blocking routines that are specifically not made cancellation
7336        points:

7337        — *pthread_mutex_lock*()

7338           If *pthread_mutex_lock*() were a cancellation point, every routine that called it would also
7339           become a cancellation point (that is, any routine that touched shared state would
7340           automatically become a cancellation point). For example, *malloc*(), *free*(), and *rand*()
7341           would become cancellation points under this scheme. Having too many cancellation
7342           points makes programming very difficult, leading to either much disabling and restoring
7343           of cancelability or much difficulty in trying to arrange for reliable cleanup at every
7344           possible place.

7345           Since *pthread_mutex_lock*() is not a cancellation point, threads could result in being
7346           blocked uninterruptibly for long periods of time if mutexes were used as a general

synchronization mechanism. As this is normally not acceptable, mutexes should only be used to protect resources that are held for small fixed lengths of time where not being able to be canceled will not be a problem. Resources that need to be held exclusively for long periods of time should be protected with condition variables.

— *pthread_barrier_wait*()

Canceling a barrier wait will render a barrier unusable. Similar to a barrier timeout (which the standard developers rejected), there is no way to guarantee the consistency of a barrier's internal data structures if a barrier wait is canceled.

— *pthread_spin_lock*()

As with mutexes, spin locks should only be used to protect resources that are held for small fixed lengths of time where not being cancelable will not be a problem.

Every library routine should specify whether or not it includes any cancellation points. Typically, only those routines that may block or compute indefinitely need to include cancellation points.

Correctly coded routines only reach cancellation points after having set up a cancellation cleanup handler to restore invariants if the thread is canceled at that point. Being cancelable only at specified cancellation points allows programmers to keep track of actions needed in a cancellation cleanup handler more easily. A thread should only be made asynchronously cancelable when it is not in the process of acquiring or releasing resources or otherwise in a state from which it would be difficult or impossible to recover.

- Thread Cancellation Cleanup Handlers

The cancellation cleanup handlers provide a portable mechanism, easy to implement, for releasing resources and restoring invariants. They are easier to use than signal handlers because they provide a stack of cancellation cleanup handlers rather than a single handler, and because they have an argument that can be used to pass context information to the handler.

The alternative to providing these simple cancellation cleanup handlers (whose only use is for cleaning up when a thread is canceled) is to define a general exception package that could be used for handling and cleaning up after hardware traps and software-detected errors. This was too far removed from the charter of providing threads to handle asynchrony. However, it is an explicit goal of IEEE Std 1003.1-2001 to be compatible with existing exception facilities and languages having exceptions.

The interaction of this facility and other procedure-based or language-level exception facilities is unspecified in this version of IEEE Std 1003.1-2001. However, it is intended that it be possible for an implementation to define the relationship between these cancellation cleanup handlers and Ada, C++, or other language-level exception handling facilities.

It was suggested that the cancellation cleanup handlers should also be called when the process exits or calls the *exec* function. This was rejected partly due to the performance problem caused by having to call the cancellation cleanup handlers of every thread before the operation could continue. The other reason was that the only state expected to be cleaned up by the cancellation cleanup handlers would be the intraprocess state. Any handlers that are to clean up the interprocess state would be registered with *atexit*(). There is the orthogonal problem that the *exec* functions do not honor the *atexit*() handlers, but resolving this is beyond the scope of IEEE Std 1003.1-2001.

7391          • Async-Cancel Safety

7392          A function is said to be async-cancel-safe if it is written in such a way that entering the
7393          function with asynchronous cancelability enabled will not cause any invariants to be |
7394          violated, even if a cancellation request is delivered at any arbitrary instruction. Functions that |
7395          are async-cancel-safe are often written in such a way that they need to acquire no resources
7396          for their operation and the visible variables that they may write are strictly limited.

7397          Any routine that gets a resource as a side effect cannot be made async-cancel-safe (for
7398          example, *malloc*( )). If such a routine were called with asynchronous cancelability enabled, it
7399          might acquire the resource successfully, but as it was returning to the client, it could act on a |
7400          cancellation request. In such a case, the application would have no way of knowing whether |
7401          the resource was acquired or not.

7402          Indeed, because many interesting routines cannot be made async-cancel-safe, most library
7403          routines in general are not async-cancel-safe. Every library routine should specify whether or
7404          not it is async-cancel safe so that programmers know which routines can be called from code
7405          that is asynchronously cancelable.

7406     IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/8 is applied, adding the *pselect*( ) function |
7407     to the list of functions with cancellation points.                                          |

7408  *B.2.9.6    Thread Read-Write Locks*

7409     **Background**

7410     Read-write locks are often used to allow parallel access to data on multi-processors, to avoid
7411     context switches on uni-processors when multiple threads access the same data, and to protect
7412     data structures that are frequently accessed (that is, read) but rarely updated (that is, written).
7413     The in-core representation of a file system directory is a good example of such a data structure.
7414     One would like to achieve as much concurrency as possible when searching directories, but limit
7415     concurrent access when adding or deleting files.

7416     Although read-write locks can be implemented with mutexes and condition variables, such
7417     implementations are significantly less efficient than is possible. Therefore, this synchronization
7418     primitive is included in IEEE Std 1003.1-2001 for the purpose of allowing more efficient
7419     implementations in multi-processor systems.

7420     **Queuing of Waiting Threads**

7421     The *pthread_rwlock_unlock*( ) function description states that one writer or one or more readers
7422     must acquire the lock if it is no longer held by any thread as a result of the call. However, the
7423     function does not specify which thread(s) acquire the lock, unless the Thread Execution
7424     Scheduling option is supported.

7425     The standard developers considered the issue of scheduling with respect to the queuing of
7426     threads blocked on a read-write lock. The question turned out to be whether
7427     IEEE Std 1003.1-2001 should require priority scheduling of read-write locks for threads whose
7428     execution scheduling policy is priority-based (for example, SCHED_FIFO or SCHED_RR). There
7429     are tradeoffs between priority scheduling, the amount of concurrency achievable among readers,
7430     and the prevention of writer and/or reader starvation.

7431     For example, suppose one or more readers hold a read-write lock and the following threads
7432     request the lock in the listed order:

```
7433          pthread_rwlock_wrlock() - Low priority thread writer_a
7434          pthread_rwlock_rdlock() - High priority thread reader_a
7435          pthread_rwlock_rdlock() - High priority thread reader_b
7436          pthread_rwlock_rdlock() - High priority thread reader_c
```

7437 When the lock becomes available, should *writer_a* block the high priority readers? Or, suppose a
7438 read-write lock becomes available and the following are queued:

```
7439          pthread_rwlock_rdlock() - Low priority thread reader_a
7440          pthread_rwlock_rdlock() - Low priority thread reader_b
7441          pthread_rwlock_rdlock() - Low priority thread reader_c
7442          pthread_rwlock_wrlock() - Medium priority thread writer_a
7443          pthread_rwlock_rdlock() - High priority thread reader_d
```

7444 If priority scheduling is applied then *reader_d* would acquire the lock and *writer_a* would block
7445 the remaining readers. But should the remaining readers also acquire the lock to increase
7446 concurrency? The solution adopted takes into account that when the Thread Execution
7447 Scheduling option is supported, high priority threads may in fact starve low priority threads (the
7448 application developer is responsible in this case for designing the system in such a way that this
7449 starvation is avoided). Therefore, IEEE Std 1003.1-2001 specifies that high priority readers take
7450 precedence over lower priority writers. However, to prevent writer starvation from threads of
7451 the same or lower priority, writers take precedence over readers of the same or lower priority.

7452 Priority inheritance mechanisms are non-trivial in the context of read-write locks. When a high
7453 priority writer is forced to wait for multiple readers, for example, it is not clear which subset of
7454 the readers should inherit the writer's priority. Furthermore, the internal data structures that
7455 record the inheritance must be accessible to all readers, and this implies some sort of
7456 serialization that could negate any gain in parallelism achieved through the use of multiple
7457 readers in the first place. Finally, existing practice does not support the use of priority
7458 inheritance for read-write locks. Therefore, no specification of priority inheritance or priority
7459 ceiling is attempted. If reliable priority-scheduled synchronization is absolutely required, it can
7460 always be obtained through the use of mutexes.

7461 **Comparison to fcntl() Locks**

7462 The read-write locks and the *fcntl*() locks in IEEE Std 1003.1-2001 share a common goal:
7463 increasing concurrency among readers, thus increasing throughput and decreasing delay.

7464 However, the read-write locks have two features not present in the *fcntl*() locks. First, under
7465 priority scheduling, read-write locks are granted in priority order. Second, also under priority
7466 scheduling, writer starvation is prevented by giving writers preference over readers of equal or
7467 lower priority.

7468 Also, read-write locks can be used in systems lacking a file system, such as those conforming to
7469 the minimal realtime system profile of IEEE Std 1003.13-1998.

7470 **History of Resolution Issues**

7471 Based upon some balloting objections, early drafts specified the behavior of threads waiting on a
7472 read-write lock during the execution of a signal handler, as if the thread had not called the lock
7473 operation. However, this specified behavior would require implementations to establish
7474 internal signal handlers even though this situation would be rare, or never happen for many
7475 programs. This would introduce an unacceptable performance hit in comparison to the little
7476 additional functionality gained. Therefore, the behavior of read-write locks and signals was
7477 reverted back to its previous mutex-like specification.

7478   *B.2.9.7   Thread Interactions with Regular File Operations*

7479          There is no additional rationale provided for this section.

7480   **B.2.10   Sockets**

7481          The base document for the sockets interfaces in IEEE Std 1003.1-2001 is the XNS, Issue 5.2
7482          specification. This was primarily chosen as it aligns with IPv6. Additional material has been
7483          added from IEEE Std 1003.1g-2000, notably socket concepts, raw sockets, the *pselect*( ) function,
7484          the *sockatmark*( ) function, and the <**sys/select.h**> header.

7485   *B.2.10.1   Address Families*

7486          There is no additional rationale provided for this section.

7487   *B.2.10.2   Addressing*

7488          There is no additional rationale provided for this section.

7489   *B.2.10.3   Protocols*

7490          There is no additional rationale provided for this section.

7491   *B.2.10.4   Routing*

7492          There is no additional rationale provided for this section.

7493   *B.2.10.5   Interfaces*

7494          There is no additional rationale provided for this section.

7495   *B.2.10.6   Socket Types*

7496          The type **socklen_t** was invented to cover the range of implementations seen in the field. The
7497          intent of **socklen_t** is to be the type for all lengths that are naturally bounded in size; that is, that
7498          they are the length of a buffer which cannot sensibly become of massive size: network addresses,
7499          host names, string representations of these, ancillary data, control messages, and socket options
7500          are examples. Truly boundless sizes are represented by **size_t** as in *read*( ), *write*( ), and so on.

7501          All **socklen_t** types were originally (in BSD UNIX) of type **int**. During the development of
7502          IEEE Std 1003.1-2001, it was decided to change all buffer lengths to **size_t**, which appears at face
7503          value to make sense. When dual mode 32/64-bit systems came along, this choice unnecessarily
7504          complicated system interfaces because **size_t** (with **long**) was a different size under ILP32 and
7505          LP64 models. Reverting to **int** would have happened except that some implementations had
7506          already shipped 64-bit-only interfaces. The compromise was a type which could be defined to be
7507          any size by the implementation: **socklen_t**.

7508   *B.2.10.7   Socket I/O Mode*

7509          There is no additional rationale provided for this section.

7510   *B.2.10.8  Socket Owner*

7511          There is no additional rationale provided for this section.

7512   *B.2.10.9  Socket Queue Limits*

7513          There is no additional rationale provided for this section.

7514   *B.2.10.10 Pending Error*

7515          There is no additional rationale provided for this section.

7516   *B.2.10.11 Socket Receive Queue*

7517          There is no additional rationale provided for this section.

7518   *B.2.10.12 Socket Out-of-Band Data State*

7519          There is no additional rationale provided for this section.

7520   *B.2.10.13 Connection Indication Queue*

7521          There is no additional rationale provided for this section.

7522   *B.2.10.14 Signals*

7523          There is no additional rationale provided for this section.

7524   *B.2.10.15 Asynchronous Errors*

7525          There is no additional rationale provided for this section.

7526   *B.2.10.16 Use of Options*

7527          There is no additional rationale provided for this section.

7528   *B.2.10.17 Use of Sockets for Local UNIX Connections*

7529          There is no additional rationale provided for this section.

7530   *B.2.10.18 Use of Sockets over Internet Protocols*

7531          A raw socket allows privileged users direct access to a protocol; for example, raw access to the
7532          IP and ICMP protocols is possible through raw sockets. Raw sockets are intended for
7533          knowledgeable applications that wish to take advantage of some protocol feature not directly
7534          accessible through the other sockets interfaces.

7535   *B.2.10.19 Use of Sockets over Internet Protocols Based on IPv4*

7536          There is no additional rationale provided for this section.

7537   *B.2.10.20 Use of Sockets over Internet Protocols Based on IPv6*

7538          The Open Group Base Resolution bwg2001-012 is applied, clarifying that IPv6 implementations
7539          are required to support use of AF_INET6 sockets over IPv4.

7540 **B.2.11  Tracing**

7541  The organization of the tracing rationale differs from the traditional rationale in that this tracing
7542  rationale text is written against the trace interface as a whole, rather than against the individual
7543  components of the trace interface or the normative section in which those components are
7544  defined. Therefore the sections below do not parallel the sections of normative text in
7545  IEEE Std 1003.1-2001.

7546  *B.2.11.1  Objectives*

7547  The intended uses of tracing are application-system debugging during system development, as a
7548  ''flight recorder'' for maintenance of fielded systems, and as a performance measurement tool. In
7549  all of these intended uses, the vendor-supplied computer system and its software are, for this
7550  discussion, assumed error-free; the intent being to debug the user-written and/or third-party
7551  application code, and their interactions. Clearly, problems with the vendor-supplied system and
7552  its software will be uncovered from time to time, but this is a byproduct of the primary activity,
7553  debugging user code.

7554  Another need for defining a trace interface in POSIX stems from the objective to provide an
7555  efficient portable way to perform benchmarks. Existing practice shows that such interfaces are
7556  commonly used in a variety of systems but with little commonality. As part of the benchmarking
7557  needs, two aspects within the trace interface must be considered.

7558  The first, and perhaps more important one, is the qualitative aspect.

7559  The second is the quantitative aspect.

7560  • Qualitative Aspect

7561  To better understand this aspect, let us consider an example. Suppose that you want to
7562  organize a number of actions to be performed during the day. Some of these actions are
7563  known at the beginning of the day. Some others, which may be more or less important, will
7564  be triggered by reading your mail. During the day you will make some phone calls and
7565  synchronously receive some more information. Finally you will receive asynchronous phone
7566  calls that also will trigger actions. If you, or somebody else, examines your day at work, you,
7567  or he, can discover that you have not efficiently organized your work. For instance, relative
7568  to the phone calls you made, would it be preferable to make some of these early in the
7569  morning? Or to delay some others until the end of the day? Relative to the phone calls you
7570  have received, you might find that somebody you called in the morning has called you 10
7571  times while you were performing some important work. To examine, afterwards, your day at
7572  work, you record in sequence all the trace events relative to your work. This should give you
7573  a chance of organizing your next day at work.

7574  This is the qualitative aspect of the trace interface. The user of a system needs to keep a trace
7575  of particular points the application passes through, so that he can eventually make some
7576  changes in the application and/or system configuration, to give the application a chance of
7577  running more efficiently.

7578  • Quantitative Aspect

7579  This aspect concerns primarily realtime applications, where missed deadlines can be
7580  undesirable. Although there are, in IEEE Std 1003.1-2001, some interfaces useful for such
7581  applications (timeouts, execution time monitoring, and so on), there are no APIs to aid in the
7582  tuning of a realtime application's behavior (**timespec** in timeouts, length of message queues,
7583  duration of driver interrupt service routine, and so on). The tuning of an application needs a
7584  means of recording timestamped important trace events during execution in order to analyze
7585  offline, and eventually, to tune some realtime features (redesign the system with less

7586            functionalities, readjust timeouts, redesign driver interrupts, and so on).

**Detailed Objectives**

7588   Objectives were defined to build the trace interface and are kept for historical interest. Although
7589   some objectives are not fully respected in this trace interface, the concept of the POSIX trace
7590   interface assumes the following points:

7591       1.   It must be possible to trace both system and user trace events concurrently.

7592       2.   It must be possible to trace per-process trace events and also to trace system trace events
7593            which are unrelated to any particular process. A per-process trace event is either user-
7594            initiated or system-initiated.

7595       3.   It must be possible to control tracing on a per-process basis from either inside or outside
7596            the process.

7597       4.   It must be possible to control tracing on a per-thread basis from inside the enclosing
7598            process.

7599       5.   Trace points must be controllable by trace event type ID from inside and outside of the
7600            process. Multiple trace points can have the same trace event type ID, and will be controlled
7601            jointly.

7602       6.   Recording of trace events is dependent on both trace event type ID and the
7603            process/thread. Both must be enabled in order to record trace events. System trace events
7604            may or may not be handled differently.

7605       7.   The API must not mandate the ability to control tracing for more than one process at the
7606            same time.

7607       8.   There is no objective for trace control on anything bigger than a process; for example,
7608            group or session.

7609       9.   Trace propagation and control:

7610            a.   Trace propagation across *fork*( ) is optional; the default is to not trace a child process.

7611            b.   Trace control must span *pthread_create*( ) operations; that is, if a process is being
7612                 traced, any thread will be traced as well if this thread allows tracing. The default is to
7613                 allow tracing.

7614       10.  Trace control must not span *exec* or *posix_spawn*( ) operations.

7615       11.  A triggering API is not required. The triggering API is the ability to command or stop
7616            tracing   based    on    the    occurrence   of    a    specific    trace    event   other   than   a
7617            POSIX_TRACE_START trace event or a POSIX_TRACE_STOP trace event.

7618       12.  Trace   log   entries   must   have   timestamps   of   implementation-defined   resolution.
7619            Implementations are exhorted to support at least microsecond resolution. When a trace log
7620            entry is retrieved, it must have timestamp, PC address, PID, and TID of the entity that
7621            generated the trace event.

7622       13.  Independently developed code should be able to use trace facilities without coordination
7623            and without conflict.

7624       14.  Even if the trace points in the trace calls are not unique, the trace log entries (after any
7625            processing) must be uniquely identified as to trace point.

7626       15.  There must be a standard API to read the trace stream.

16. The format of the trace stream and the trace log is opaque and unspecified.

17. It must be possible to read a completed trace, if recorded on some suitable non-volatile storage, even subsequent to a power cycle or subsequent cold boot of the system.

18. Support of analysis of a trace log while it is being formed is implementation-defined.

19. The API must allow the application to write trace stream identification information into the trace stream and to be able to retrieve it, without it being overwritten by trace entries, even if the trace stream is full.

20. It must be possible to specify the destination of trace data produced by trace events.

21. It must be possible to have different trace streams, and for the tracing enabled by one trace stream to be completely independent of the tracing of another trace stream.

22. It must be possible to trace events from threads in different CPUs.

23. The API must support one or more trace streams per-system, and one or more trace streams per-process, up to an implementation-defined set of per-system and per-process maximums.

24. It must be possible to determine the order in which the trace events happened, without necessarily depending on the clock, up to an implementation-defined time resolution.

25. For performance reasons, the trace event point call(s) must be implementable as a macro (see the ISO POSIX-1: 1996 standard, 1.3.4, Statement 2).

26. IEEE Std 1003.1-2001 must not define the trace points which a conforming system must implement, except for trace points used in the control of tracing.

27. The APIs must be thread-safe, and trace points should be lock-free (that is, not require a lock to gain exclusive access to some resource).

28. The user-provided information associated with a trace event is variable-sized, up to some maximum size.

29. Bounds on record and trace stream sizes:

    a. The API must permit the application to declare the upper bounds on the length of an application data record. The system must return the limit it used. The limit used may be smaller than requested.

    b. The API must permit the application to declare the upper bounds on the size of trace streams. The system must return the limit it used. The limit used may be different, either larger or smaller, than requested.

30. The API must be able to pass any fundamental data type, and a structured data type composed only of fundamental types. The API must be able to pass data by reference, given only as an address and a length. Fundamental types are the POSIX.1 types (see the **<sys/types.h>** header) plus those defined in the ISO C standard.

31. The API must apply the POSIX notions of ownership and permission to recorded trace data, corresponding to the sources of that data.

**Comments on Objectives**

7664

7665 **Note:** In the following comments, numbers in square brackets refer to the above objectives.

7666 It is necessary to be able to obtain a trace stream for a complete activity. Thus there is a
7667 requirement to be able to trace both application and system trace events. A per-process trace
7668 event is either user-initiated, like the *write*( ) function, or system-initiated, like a timer expiration.
7669 There is also a need to be able to trace an entire process' activity even when it has threads in
7670 multiple CPUs. To avoid excess trace activity, it is necessary to be able to control tracing on a
7671 trace event type basis.
7672 [Objectives 1,2,5,22]

7673 There is a need to be able to control tracing on a per-process basis, both from inside and outside
7674 the process; that is, a process can start a trace activity on itself or any other process. There is also
7675 the perceived need to allow the definition of a maximum number of trace streams per system.
7676 [Objectives 3,23]

7677 From within a process, it is necessary to be able to control tracing on a per-thread basis. This
7678 provides an additional filtering capability to keep the amount of traced data to a minimum. It
7679 also allows for less ambiguity as to the origin of trace events. It is recognized that thread-level
7680 control is only valid from within the process itself. It is also desirable to know the maximum
7681 number of trace streams per process that can be started. The API should not require thread
7682 synchronization or mandate priority inversions that would cause the thread to block. However,
7683 the API must be thread-safe.
7684 [Objectives 4,23,24,27]

7685 There was no perceived objective to control tracing on anything larger than a process; for
7686 example, a group or session. Also, the ability to start or stop a trace activity on multiple
7687 processes atomically may be very difficult or cumbersome in some implementations.
7688 [Objectives 6,8]

7689 It is also necessary to be able to control tracing by trace event type identifier, sometimes called a
7690 trace hook ID. However, there is no mandated set of system trace events, since such trace points
7691 are implementation-defined. The API must not require from the operating system facilities that
7692 are not standard.
7693 [Objectives 6,26]

7694 Trace control must span *fork*( ) and *pthread_create*( ). If not, there will be no way to ensure that an
7695 application's activity is entirely traced. The newly forked child would not be able to turn on its
7696 tracing until after it obtained control after the fork, and trace control externally would be even
7697 more problematic.
7698 [Objective 9]

7699 Since *exec* and *posix_spawn*( ) represent a complete change in the execution of a task (a new
7700 program), trace control need not persist over an *exec* or *posix_spawn*( ).
7701 [Objective 10]

7702 Where trace activities are started on multiple processes, these trace activities should not interfere
7703 with each other.
7704 [Objective 21]

7705 There is no need for a triggering objective, primarily for performance reasons; see also Section
7706 B.2.11.8 (on page 202), rationale on triggering.
7707 [Objective 11]

7708 It must be possible to determine the origin of each traced event. The process and thread
7709 identifiers for each trace event are needed. Also there was a perceived need for a user-specifiable
7710 origin, but it was felt that this would create too much overhead.

7711 [Objectives 12,14]

7712 An allowance must be made for trace points to come embedded in software components from
7713 several different sources and vendors without requiring coordination.
7714 [Objective 13]

7715 There is a requirement to be able to uniquely identify trace points that may have the same trace
7716 stream identifier. This is only necessary when a trace report is produced.
7717 [Objectives 12,14]

7718 Tracing is a very performance-sensitive activity, and will therefore likely be implemented at a
7719 low level within the system. Hence the interface must not mandate any particular buffering or
7720 storage method. Therefore, a standard API is needed to read a trace stream. Also the interface
7721 must not mandate the format of the trace data, and the interface must not assume a trace storage
7722 method. Due to the possibility of a monolithic kernel and the possible presence of multiple
7723 processes capable of running trace activities, the two kinds of trace events may be stored in two
7724 separate streams for performance reasons. A mandatory dump mechanism, common in some
7725 existing practice, has been avoided to allow the implementation of this set of functions on small
7726 realtime profiles for which the concept of a file system is not defined. The trace API calls should
7727 be implemented as macros.
7728 [Objectives 15,16,25,30]

7729 Since a trace facility is a valuable service tool, the output (or log) of a completed trace stream
7730 that is written to permanent storage must be readable on other systems of the type that
7731 produced the trace log. Note that there is no objective to be able to interpret a trace log that was
7732 not successfully completed.
7733 [Objectives 17,18,19]

7734 For trace streams written to permanent storage, a way to specify the destination of the trace
7735 stream is needed.
7736 [Objective 20]

7737 There is a requirement to be able to depend on the ordering of trace events up to some
7738 implementation-defined time interval. For example, there is a need to know the time period
7739 during which, if trace events are closer together, their ordering is unspecified. Events that occur
7740 within an interval smaller than this resolution may or may not be read back in the correct order.
7741 [Objective 24]

7742 The application should be able to know how much data can be traced. When trace event types
7743 can be filtered, the application should be able to specify the approximate maximum amount of
7744 data that will be traced in a trace event so resources can be more efficiently allocated.
7745 [Objectives 28,29]

7746 Users should not be able to trace data to which they would not normally have access. System
7747 trace events corresponding to a process/thread should be associated with the ownership of that
7748 process/thread.
7749 [Objective 31]

7750  *B.2.11.2  Trace Model*

7751  **Introduction**

7752  The model is based on two base entities: the ''Trace Stream'' and the ''Trace Log'', and a recorded
7753  unit called the ''Trace Event''. The possibility of using Trace Streams and Trace Logs separately
7754  gives two use dimensions and solves both the performance issue and the full-information
7755  system issue. In the case of a trace stream without log, specific information, although reduced in
7756  quantity, is required to be registered, in a possibly small realtime system, with as little overhead
7757  as possible. The Trace Log option has been added for small realtime systems. In the case of a
7758  trace stream with log, considerable complex application-specific information needs to be
7759  collected.

7760  **Trace Model Description**

7761  The trace model can be examined for three different subfunctions:  Application Instrumentation,
7762  Trace Operation Control, and Trace Analysis.

7763



7764                    **Figure B-2**  Trace System Overview: for Offline Analysis

7765  Each of these subfunctions requires specific characteristics of the trace mechanism API.

7766  • Application Instrumentation

7767  When instrumenting an application, the programmer is not concerned about the future use of
7768  the trace events in the trace stream or the trace log, the full policy of the trace stream, or the
7769  eventual pre-filtering of trace events. But he is concerned about the correct determination of
7770  the specific trace event type identifier, regardless of how many independent libraries are
7771  used in the same user application; see Figure B-2 and Figure B-3 (on page 185).

7772  This trace API provides the necessary operations to accomplish this subfunction. This is done
7773  by providing functions to associate a programmer-defined name with an implementation-
7774  defined trace event type identifier (see the *posix_trace_eventid_open*( ) function), and to send
7775  this trace event into a potential trace stream (see the *posix_trace_event*( ) function).

7776        • Trace Operation Control

7777        When controlling the recording of trace events in a trace stream, the programmer is
7778        concerned with the correct initialization of the trace mechanism (that is, the sizing of the
7779        trace stream), the correct retention of trace events in a permanent storage, the correct
7780        dynamic recording of trace events, and so on.

7781        This trace API provides the necessary material to permit this efficiently. This is done by
7782        providing functions to initialize a new trace stream, and optionally a trace log:

7783        — Trace Stream Attributes Object Initialization (see *posix_trace_attr_init*())

7784        — Functions    to    Retrieve    or    Set    Information    About    a    Trace    Stream    (see
7785        *posix_trace_attr_getgenversion*())

7786        — Functions    to    Retrieve    or    Set    the    Behavior    of    a    Trace    Stream    (see
7787        *posix_trace_attr_getinherited*())

7788        — Functions    to    Retrieve    or    Set    Trace    Stream    Size    Attributes    (see
7789        *posix_trace_attr_getmaxusereventsize*())

7790        — Trace Stream Initialization, Flush, and Shutdown from a Process (see *posix_trace_create*())

7791        — Clear Trace Stream and Trace Log (see *posix_trace_clear*())

7792        To select the trace event types that are to be traced:

7793        — Manipulate Trace Event Type Identifier (see *posix_trace_trid_eventid_open*())

7794        — Iterate over a Mapping of Trace Event Type (see *posix_trace_eventtypelist_getnext_id*())

7795        — Manipulate Trace Event Type Sets (see *posix_trace_eventset_empty*())

7796        — Set Filter of an Initialized Trace Stream (see *posix_trace_set_filter*())

7797        To control the execution of an active trace stream:

7798        — Trace Start and Stop (see *posix_trace_start*())

7799        — Functions to Retrieve the Trace Attributes or Trace Statuses (see *posix_trace_get_attr*())

7800



7801        **Figure B**-**3**  Trace System Overview: for Online Analysis

7802        • Trace Analysis

7803        Once correctly recorded, on permanent storage or not, an ultimate activity consists of the
7804        analysis of the recorded information. If the recorded data is on permanent storage, a specific
7805        open operation is required to associate a trace stream to a trace log.

7806        The first intent of the group was to request the presence of a system identification structure
7807        in the trace stream attribute. This was, for the application, to allow some portable way to
7808        process the recorded information. However, there is no requirement that the **utsname**
7809        structure, on which this system identification was based, be portable from one machine to
7810        another, so the contents of the attribute cannot be interpreted correctly by an application
7811        conforming to IEEE Std 1003.1-2001.

7812        This modification has been incorporated and requests that some unspecified information be
7813        recorded in the trace log in order to fail opening it if the analysis process and the controller
7814        process were running in different types of machine, but does not request that this
7815        information be accessible to the application. This modification has implied a modification in
7816        the *posix_trace_open*( ) function error code returns.

7817        This trace API provides functions to:

7818        — Extract trace stream identification attributes (see *posix_trace_attr_getgenversion*( ))

7819        — Extract trace stream behavior attributes (see *posix_trace_attr_getinherited*( ))

7820        — Extract     trace     event,     stream,     and     log     size     attributes     (see
7821          *posix_trace_attr_getmaxusereventsize*( ))

7822        — Look up trace event type names (see *posix_trace_eventid_get_name*( ))

7823        — Iterate over trace event type identifiers (see *posix_trace_eventtypelist_getnext_id*( ))

7824        — Open, rewind, and close a trace log (see *posix_trace_open*( ))

7825        — Read trace stream attributes and status (see *posix_trace_get_attr*( ))

7826        — Read trace events (see *posix_trace_getnext_event*( ))

7827     Due to the following two reasons:

7828        1.  The requirement that the trace system must not add unacceptable overhead to the traced
7829            process and so that the trace event point execution must be fast

7830        2.  The traced application does not care about tracing errors

7831     the trace system cannot return any internal error to the application.  Internal error conditions can
7832     range from unrecoverable errors that will force the active trace stream to abort, to small errors
7833     that can affect the quality of tracing without aborting the trace stream. The group decided to
7834     define a system trace event to report to the analysis process such internal errors. It is not the
7835     intention of IEEE Std 1003.1-2001 to require an implementation to report an internal error that
7836     corrupts or terminates tracing operation. The implementor is free to decide which internal
7837     documented errors, if any, the trace system is able to report.

7838          **States of a Trace Stream**

7839



7840                    **Figure B-4**  Trace System Overview: States of a Trace Stream

7841     Figure B-4 shows the different states an active trace stream passes through. After the
7842     *posix_trace_create*( ) function call, a trace stream becomes CREATED and a trace stream is
7843     associated for the future collection of trace events. The status of the trace stream is
7844     POSIX_TRACE_SUSPENDED. The state becomes STARTED after a call to the *posix_trace_start*( )
7845     function, and the status becomes POSIX_TRACE_RUNNING. In this state, all trace events that
7846     are not filtered out will be stored into the trace stream. After a call to *posix_trace_stop*( ), the trace
7847     stream becomes STOPPED (and the status POSIX_TRACE_SUSPENDED). In this state, no new
7848     trace events will be recorded in the trace stream, but previously recorded trace events may
7849     continue to be read.

7850     After a call to *posix_trace_shutdown*( ), the trace stream is in the state COMPLETED. The trace
7851     stream no longer exists but, if the Trace Log option is supported, all the information contained in
7852     it has been logged. If a log object has not been associated with the trace stream at the creation, it
7853     is the responsibility of the trace controller process to not shut the trace stream down while trace
7854     events remain to be read in the stream.

7855          **Tracing All Processes**

7856     Some implementations have a tracing subsystem with the ability to trace all processes. This is
7857     useful to debug some types of device drivers such as those for ATM or X25 adapters. These types
7858     of adapters are used by several independent processes, that are not issued from the same
7859     process.

7860     The POSIX trace interface does not define any constant or option to create a trace stream tracing
7861     all processes. POSIX.1 does not prevent this type of implementation and an implementor is free
7862     to add this capability. Nevertheless, the trace interface allows tracing of all the system trace
7863     events and all the processes issued from the same process.

7864    If such a tracing system capability has to be implemented, when a trace stream is created, it is
7865    recommended that a constant named POSIX_TRACE_ALLPROC be used instead of the process
7866    identifier in the argument of the *posix_trace_create*( ) or *posix_trace_create_withlog*( ) function. A
7867    possible value for POSIX_TRACE_ALLPROC may be −1 instead of a real process identifier.

7868    The implementor has to be aware that there is some impact on the tracing behavior as defined in
7869    the POSIX trace interface. For example:

7870    • If    the    default    value    for    the    inheritance    attribute    is    set    to
7871      POSIX_TRACE_CLOSE_FOR_CHILD, the implementation has to stop tracing for the child
7872      process.

7873    • The trace controller which is creating this type of trace stream must have the appropriate
7874      privilege to trace all the processes.

7875    **Trace Storage**

7876    The model is based on two types of trace events: system trace events and user-defined trace
7877    events. The internal representation of trace events is implementation-defined, and so the
7878    implementor is free to choose the more suitable, practical, and efficient way to design the
7879    internal management of trace events. For the timestamping operation, the model does not
7880    impose the CLOCK_REALTIME or any other clock. The buffering allocation and operation
7881    follow the same principle. The implementor is free to use one or more buffers to record trace
7882    events; the interface assumes only a logical trace stream of sequentially recorded trace events.
7883    Regarding flushing of trace events, the interface allows the definition of a trace log object which
7884    typically can be a file. But the group was also aware of defining functions to permit the use of
7885    this interface in small realtime systems, which may not have general file system capabilities. For
7886    instance,    the    three    functions    *posix_trace_getnext_event*( )    (blocking),
7887    *posix_trace_timedgetnext_event*( ) (blocking with timeout), and *posix_trace_trygetnext_event*( )
7888    (non-blocking) are proposed to read the recorded trace events.

7889    The policy to be used when the trace stream becomes full also relies on common practice:

7890    • For an active trace stream, the POSIX_TRACE_LOOP trace stream policy permits automatic
7891      overrun (overwrite of oldest trace events) while waiting for some user-defined condition to
7892      cause tracing to stop. By contrast, the POSIX_TRACE_UNTIL_FULL trace stream policy
7893      requires the system to stop tracing when the trace stream is full. However, if the trace stream
7894      that is full is at least partially emptied by a call to the *posix_trace_flush*( ) function or by calls
7895      to the *posix_trace_getnext_event*( ) function, the trace system will automatically resume
7896      tracing.

7897      If the Trace Log option is supported, the operation of the POSIX_TRACE_FLUSH policy is an
7898      extension of the POSIX_TRACE_UNTIL_FULL policy. The automatic free operation (by
7899      flushing to the associated trace log) is added.

7900    • If a log is associated with the trace stream and this log is a regular file, these policies also
7901      apply for the log. One more policy, POSIX_TRACE_APPEND, is defined to allow indefinite
7902      extension of the log. Since the log destination can be any device or pseudo-device, the
7903      implementation may not be able to manipulate the destination as required by
7904      IEEE Std 1003.1-2001. For this reason, the behavior of the log full policy may be unspecified
7905      depending on the trace log type.

7906      The current trace interface does not define a service to preallocate space for a trace log file,
7907      because this space can be preallocated by means of a call to the *posix_fallocate*( ) function. This
7908      function could be called after the file has been opened, but before the trace stream is created.
7909      The *posix_fallocate*( ) function ensures that any required storage for regular file data is
7910      allocated on the file system storage media. If *posix_fallocate*( ) returns successfully,

| | |
|---|---|
| 7911 | subsequent writes to the specified file data will not fail due to the lack of free space on the file |
| 7912 | system storage media. Besides trace events, a trace stream also includes trace attributes and |
| 7913 | the mapping from trace event names to trace event type identifiers. The implementor is free |
| 7914 | to choose how to store the trace attributes and the trace event type map, but must ensure that |
| 7915 | this information is not lost when a trace stream overrun occurs. |

7916 *B.2.11.3  Trace Programming Examples*

| | |
|---|---|
| 7917 | Several programming examples are presented to show the code of the different possible |
| 7918 | subfunctions using a trace subsystem. All these programs need to include the **<trace.h>** header. |
| 7919 | In the examples shown, error checking is omitted for more simplicity. |

7920        **Trace Operation Control**

| | |
|---|---|
| 7921 | These examples show the creation of a trace stream for another process; one which is already |
| 7922 | trace instrumented. All the default trace stream attributes are used to simplify programming in |
| 7923 | the first example. The second example shows more possibilities. |

7924        **First Example**

```
7925    /* Caution. Error checks omitted */
7926    {
7927        trace_attr_t attr;
7928        pid_t pid = traced_process_pid;
7929        int fd;
7930        trace_id_t trid;

7931        - - - - - -
7932        /* Initialize trace stream attributes */
7933        posix_trace_attr_init(&attr);
7934        /* Open a trace log */
7935        fd=open("/tmp/mytracelog",...);
7936        /*
7937         * Create a new trace associated with a log
7938         * and with default attributes
7939         */
7940        posix_trace_create_withlog(pid, &attr, fd, &trid);

7941        /* Trace attribute structure can now be destroyed */
7942        posix_trace_attr_destroy(&attr);
7943        /* Start of trace event recording */
7944        posix_trace_start(trid);
7945        - - - - - -
7946        - - - - - -
7947        /* Duration of tracing */
7948        - - - - - -
7949        - - - - - -
7950        /* Stop and shutdown of trace activity */
7951        posix_trace_shutdown(trid);
7952        - - - - - -
7953    }
```

Between the initialization of the trace stream attributes and the creation of the trace stream, these trace stream attributes may be modified; see **Trace Stream Attribute Manipulation** (on page 194) for a specific programming example. Between the creation and the start of the trace stream, the event filter may be set; after the trace stream is started, the event filter may be changed. The setting of an event set and the change of a filter is shown in **Create a Trace Event Type Set and Change the Trace Event Type Filter** (on page 194).

```
/* Caution. Error checks omitted */
{
    trace_attr_t attr;
    pid_t pid = traced_process_pid;
    int fd;
    trace_id_t trid;
    - - - - - -
    /* Initialize trace stream attributes */
    posix_trace_attr_init(&attr);
    /* Attr default may be changed at this place; see example */
    - - - - - -
    /* Create and open a trace log with R/W user access */
    fd=open("/tmp/mytracelog",O_WRONLY|O_CREAT,S_IRUSR|S_IWUSR);
    /* Create a new trace associated with a log */
    posix_trace_create_withlog(pid, &attr, fd, &trid);
    /*
     * If the Trace Filter option is supported
     * trace event type filter default may be changed at this place;
     * see example about changing the trace event type filter
     */
    posix_trace_start(trid);
    - - - - - -

    /*
     * If you have an uninteresting part of the application
     * you can stop temporarily.
     *
     * posix_trace_stop(trid);
     * - - - - - -
     * - - - - - -
     * posix_trace_start(trid);
     */
    - - - - - -
    /*
     * If the Trace Filter option is supported
     * the current trace event type filter can be changed
     * at any time (see example about how to set
     * a trace event type filter)
     */
    - - - - - -

    /* Stop the recording of trace events */
    posix_trace_stop(trid);
    /* Shutdown the trace stream */
    posix_trace_shutdown(trid);
```

```
8004            /*
8005             * Destroy trace stream attributes; attr structure may have
8006             * been used during tracing to fetch the attributes
8007             */
8008            posix_trace_attr_destroy(&attr);
8009            - - - - - -
8010        }
```

**Application Instrumentation**

This example shows an instrumented application. The code is included in a block of instructions, perhaps a function from a library. Possibly in an initialization part of the instrumented application, two user trace events names are mapped to two trace event type identifiers (function *posix_trace_eventid_open*()). Then two trace points are programmed.

```
8016    /* Caution. Error checks omitted */
8017    {
8018        trace_event_id_t eventid1, eventid2;
8019        - - - - - -
8020        /* Initialization of two trace event type ids */
8021        posix_trace_eventid_open("my_first_event",&eventid1);
8022        posix_trace_eventid_open("my_second_event",&eventid2);
8023        - - - - - -
8024        - - - - - -
8025        - - - - - -
8026        /* Trace point */
8027        posix_trace_event(eventid1,NULL,0);
8028        - - - - - -
8029        /* Trace point */
8030        posix_trace_event(eventid2,NULL,0);
8031        - - - - - -
8032    }
```

**Trace Analyzer**

This example shows the manipulation of a trace log resulting from the dumping of a completed trace stream. All the default attributes are used to simplify programming, and data associated with a trace event is not shown in the first example. The second example shows more possibilities.

**First Example**

```
8039    /* Caution. Error checks omitted */
8040    {
8041        int fd;
8042        trace_id_t trid;
8043        posix_trace_event_info trace_event;
8044        char trace_event_name[TRACE_EVENT_NAME_MAX];
8045        int return_value;
8046        size_t returndatasize;
8047        int lost_event_number;

8048        - - - - - -
```

```
8049            /* Open an existing trace log */
8050            fd=open("/tmp/tracelog", O_RDONLY);
8051            /* Open a trace stream on the open log */
8052            posix_trace_open(fd, &trid);
8053            /* Read a trace event */
8054            posix_trace_getnext_event(trid, &trace_event,
8055                NULL, 0, &returndatasize,&return_value);

8056            /* Read and print all trace event names out in a loop */
8057            while (return_value == NULL)
8058            {
8059                /*
8060                 * Get the name of the trace event associated
8061                 * with trid trace ID
8062                 */
8063                posix_trace_eventid_get_name(trid, trace_event.event_id,
8064                    trace_event_name);
8065                /* Print the trace event name out */
8066                printf("%s\n",trace_event_name);
8067                /* Read a trace event */
8068                posix_trace_getnext_event(trid, &trace_event,
8069                    NULL, 0, &returndatasize,&return_value);
8070            }
8071            /* Close the trace stream */
8072            posix_trace_close(trid);
8073            /* Close the trace log */
8074            close(fd);
8075        }
```

8076    **Second Example**

8077    The complete example includes the two other examples in **Retrieve Information from a Trace**
8078    **Log** (on page 195) and in **Retrieve the List of Trace Event Types Used in a Trace Log** (on page
8079    196). For example, the *maxdatasize* variable is set in **Retrieve the List of Trace Event Types Used**
8080    **in a Trace Log** (on page 196).

```
8081        /* Caution. Error checks omitted */
8082        {
8083            int fd;
8084            trace_id_t trid;
8085            posix_trace_event_info trace_event;
8086            char trace_event_name[TRACE_EVENT_NAME_MAX];
8087            char * data;
8088            size_t maxdatasize=1024, returndatasize;
8089            int return_value;
8090            - - - - - -

8091            /* Open an existing trace log */
8092            fd=open("/tmp/tracelog", O_RDONLY);
8093            /* Open a trace stream on the open log */
8094            posix_trace_open( fd, &trid);
8095            /*
8096             * Retrieve information about the trace stream which
```

```
8097                * was dumped in this trace log (see example)
8098                */
8099                - - - - - -

8100                /* Allocate a buffer for trace event data */
8101                data=(char *)malloc(maxdatasize);
8102                /*
8103                 * Retrieve the list of trace events used in this
8104                 * trace log (see example)
8105                 */
8106                - - - - - -

8107                /* Read and print all trace event names and data out in a loop */
8108                while (1)
8109                {
8110                posix_trace_getnext_event(trid, &trace_event,
8111                    data, maxdatasize, &returndatasize,&return_value);
8112                    if (return_value != NULL) break;
8113                    /*
8114                     * Get the name of the trace event type associated
8115                     * with trid trace ID
8116                     */
8117                    posix_trace_eventid_get_name(trid, trace_event.event_id,
8118                        trace_event_name);
8119                    {
8120                    int i;

8121                    /* Print the trace event name out */
8122                    printf("%s: ", trace_event_name);
8123                    /* Print the trace event data out */
8124                    for (i=0; i<returndatasize, i++) printf("%02.2X",
8125                        (unsigned char)data[i]);
8126                    printf("\n");
8127                    }
8128                }

8129            /* Close the trace stream */
8130            posix_trace_close(trid);
8131            /* The buffer data is deallocated */
8132            free(data);
8133            /* Now the file can be closed */
8134            close(fd);
8135        }
```

**Several Programming Manipulations**

The following examples show some typical sets of operations needed in some contexts.

**Trace Stream Attribute Manipulation**

This example shows the manipulation of a trace stream attribute object in order to change the default value provided by a previous *posix_trace_attr_init*() call.

```
/* Caution. Error checks omitted */
{
    trace_attr_t attr;
    size_t logsize=100000;
    - - - - - -
    /* Initialize trace stream attributes */
    posix_trace_attr_init(&attr);
    /* Set the trace name in the attributes structure */
    posix_trace_attr_setname(&attr, "my_trace");
    /* Set the trace full policy */
    posix_trace_attr_setstreamfullpolicy(&attr, POSIX_TRACE_LOOP);
    /* Set the trace log size */
    posix_trace_attr_setlogsize(&attr, logsize);
    - - - - - -
}
```

**Create a Trace Event Type Set and Change the Trace Event Type Filter**

This example is valid only if the Trace Event Filter option is supported. This example shows the manipulation of a trace event type set in order to change the trace event type filter for an existing active trace stream, which may be just-created, running, or suspended. Some sets of trace event types are well-known, such as the set of trace event types not associated with a process, some trace event types are just-built trace event types for this trace stream; one trace event type is the predefined trace event error type which is deleted from the trace event type set.

```
/* Caution. Error checks omitted */
{
    trace_id_t trid = existing_trace;
    trace_event_set_t set;
    trace_event_id_t trace_event1, trace_event2;
    - - - - - -
    /* Initialize to an empty set of trace event types */
    /* (not strictly required because posix_trace_event_set_fill() */
    /* will ignore the prior contents of the event set.) */
    posix_trace_eventset_emptyset(&set);
    /*
     * Fill the set with all system trace events
     * not associated with a process
     */
    posix_trace_eventset_fill(&set, POSIX_TRACE_WOPID_EVENTS);

    /*
     * Get the trace event type identifier of the known trace event name
     * my_first_event for the trid trace stream
     */
    posix_trace_trid_eventid_open(trid, "my_first_event", &trace_event1);
    /* Add the set with this trace event type identifier */
    posix_trace_eventset_add_event(trace_event1, &set);
    /*
```

```
8186              * Get the trace event type identifier of the known trace event name
8187              * my_second_event for the trid trace stream
8188              */
8189             posix_trace_trid_eventid_open(trid, "my_second_event", &trace_event2);
8190             /* Add the set with this trace event type identifier */
8191             posix_trace_eventset_add_event(trace_event2, &set);
8192             - - - - - -
8193             /* Delete the system trace event POSIX_TRACE_ERROR from the set */
8194             posix_trace_eventset_del_event(POSIX_TRACE_ERROR, &set);
8195             - - - - - -
8196             /* Modify the trace stream filter making it equal to the new set */
8197             posix_trace_set_filter(trid, &set, POSIX_TRACE_SET_EVENTSET);
8198             - - - - - -
8199             /*
8200              * Now trace_event1, trace_event2, and all system trace event types
8201              * not associated with a process, except for the POSIX_TRACE_ERROR
8202              * system trace event type, are filtered out of (not recorded in) the
8203              * existing trace stream.
8204              */
8205         }
```

8206     **Retrieve Information from a Trace Log**

8207     This example shows how to extract information from a trace log, the dump of a trace stream.
8208     This code:

8209     • Asks if the trace stream has lost trace events

8210     • Extracts the information about the version of the trace subsystem which generated this trace
8211        log

8212     • Retrieves the maximum size of trace event data; this may be used to dynamically allocate an
8213        array for extracting trace event data from the trace log without overflow

```
8214         /* Caution. Error checks omitted */
8215         {
8216             struct posix_trace_status_info statusinfo;
8217             trace_attr_t attr;
8218             trace_id_t trid = existing_trace;
8219             size_t maxdatasize;
8220             char genversion[TRACE_NAME_MAX];
8221             - - - - - -
8222             /* Get the trace stream status */
8223             posix_trace_get_status(trid, &statusinfo);
8224             /* Detect an overrun condition */
8225             if (statusinfo.posix_stream_overrun_status == POSIX_TRACE_OVERRUN)
8226                 printf("trace events have been lost\n");
8227             /* Get attributes from the trid trace stream */
8228             posix_trace_get_attr(trid, &attr);
8229             /* Get the trace generation version from the attributes */
8230             posix_trace_attr_getgenversion(&attr, genversion);
8231             /* Print the trace generation version out */
8232             printf("Information about Trace Generator:%s\n",genversion);
```

```
8233                /* Get the trace event max data size from the attributes */
8234                posix_trace_attr_getmaxdatasize(&attr, &maxdatasize);
8235                /* Print the trace event max data size out */
8236                printf("Maximum size of associated data:%d\n",maxdatasize);
8237                /* Destroy the trace stream attributes */
8238                posix_trace_attr_destroy(&attr);
8239            }
```

8240            **Retrieve the List of Trace Event Types Used in a Trace Log**

8241            This example shows the retrieval of a trace stream's trace event type list. This operation may be
8242            very useful if you are interested only in tracking the type of trace events in a trace log.

```
8243        /* Caution. Error checks omitted */
8244        {
8245            trace_id_t trid = existing_trace;
8246            trace_event_id_t event_id;
8247            char event_name[TRACE_EVENT_NAME_MAX];
8248            int return_value;
8249            - - - - - -

8250            /*
8251             * In a loop print all existing trace event names out
8252             * for the trid trace stream
8253             */
8254            while (1)
8255            {
8256                posix_trace_eventtypelist_getnext_id(trid, &event_id
8257                    &return_value);
8258                if (return_value != NULL) break;
8259                /*
8260                 * Get the name of the trace event associated
8261                 * with trid trace ID
8262                 */
8263                posix_trace_eventid_get_name(trid, event_id, event_name);
8264                /* Print the name out */
8265                printf("%s\n", event_name);
8266            }
8267        }
```

8268   *B.2.11.4  Rationale on Trace for Debugging*

8269



8270                             **Figure B**-**5**  Trace Another Process

8271   Among the different possibilities offered by the trace interface defined in IEEE Std 1003.1-2001,
8272   the debugging of an application is the most interesting one. Typical operations in the controlling
8273   debugger process are to filter trace event types, to get trace events from the trace stream, to stop
8274   the trace stream when the debugged process is executing uninteresting code, to start the trace
8275   stream when some interesting point is reached, and so on. The interface defined in
8276   IEEE Std 1003.1-2001 should define all the necessary base functions to allow this dynamic debug
8277   handling.

8278   Figure B-5 shows an example in which the trace stream is created after the call to the *fork*( )
8279   function. If the user does not want to lose trace events, some synchronization mechanism
8280   (represented in the figure) may be needed before calling the *exec* function, to give the parent a
8281   chance to create the trace stream before the child begins the execution of its trace points.

8282   *B.2.11.5  Rationale on Trace Event Type Name Space*

8283   At first, the working group was in favor of the representation of a trace event type by an integer
8284   (*event_name*).  It seems that existing practice shows the weakness of such a representation. The
8285   collision of trace event types is the main problem that cannot be simply resolved using this sort
8286   of representation.  Suppose, for example, that a third party designs an instrumented library. The
8287   user does not have the source of this library and wants to trace his application which uses in
8288   some part the third-party library. There is no means for him to know what are the trace event
8289   types used in the instrumented library so he has some chance of duplicating some of them and
8290   thus to obtain a contaminated tracing of his application.

8291



8292                    **Figure B**-6  Trace Name Space Overview: With Third-Party Library

8293   There are requirements to allow program images containing pieces from various vendors to be
8294   traced without also requiring those of any other vendors to coordinate their uses of the trace
8295   facility, and especially the naming of their various trace event types and trace point IDs. The
8296   chosen solution is to provide a very large name space, large enough so that the individual
8297   vendors can give their trace types and tracepoint IDs sufficiently long and descriptive names
8298   making the occurrence of collisions quite unlikely. The probability of collision is thus made
8299   sufficiently low so that the problem may, as a practical matter, be ignored. By requirement, the
8300   consequence of collisions will be a slight ambiguity in the trace streams; tracing will continue in
8301   spite of collisions and ambiguities. ''The show must go on''. The *posix_prog_address* member of
8302   the **posix_trace_event_info** structure is used to allow trace streams to be unambiguously
8303   interpreted, despite the fact that trace event types and trace event names need not be unique.

8304   The *posix_trace_eventid_open*( ) function is required to allow the instrumented third-party library
8305   to get a valid trace event type identifier for its trace event names.  This operation is, somehow,
8306   an allocation, and the group was aware of proposing some deallocation mechanism which the
8307   instrumented application could use to recover the resources used by a trace event type identifier.
8308   This would have given the instrumented application the benefit of being capable of reusing a
8309   possible minimum set of trace event type identifiers, but also the inconvenience to have,
8310   possibly in the same trace stream, one trace event type identifier identifying two different trace
8311   event types. After some discussions the group decided to not define such a function which
8312   would make this API thicker for little benefit, the user having always the possibility of adding
8313   identification information in the *data* member of the trace event structure.

8314   The set of the trace event type identifiers the controlling process wants to filter out is initialized
8315   in the trace mechanism using the function *posix_trace_set_filter*( ), setting the arguments
8316   according to the definitions explained in *posix_trace_set_filter*( ). This operation can be done
8317   statically (when the trace is in the STOPPED state) or dynamically (when the trace is in the
8318   STARTED state). The preparation of the filter is normally done using the function defined in
8319   *posix_trace_eventtypelist_getnext_id*( )          and          eventually          the          function
8320   *posix_trace_eventtypelist_rewind*( ) in order to know (before the recording) the list of the potential

8321    set of trace event types that can be recorded. In the case of an active trace stream, this list may
8322    not be exhaustive. Actually, the target process may not have yet called the function
8323    *posix_trace_eventid_open*( ). But it is a common practice, for a controlling process, to prepare the
8324    filtering of a future trace stream before its start. Therefore the user must have a way to get the
8325    trace event type identifier corresponding to a well-known trace event name before its future
8326    association by the pre-cited function. This is done by calling the *posix_trace_trid_eventid_open*( )
8327    function, given the trace stream identifier and the trace name, and described hereafter. Because
8328    this trace event type identifier is associated with a trace stream identifier, where a unique
8329    process has initialized two or more traces, the implementation is expected to return the same
8330    trace event type identifier for successive calls to *posix_trace_trid_eventid_open*( ) with different
8331    trace stream identifiers. The *posix_trace_eventid_get_name*( ) function is used by the controller
8332    process to identify, by the name, the trace event type returned by a call to the
8333    *posix_trace_eventtypelist_getnext_id*( ) function.

8334    Afterwards, the set of trace event types is constructed using the functions defined in
8335    *posix_trace_eventset_empty*( ),      *posix_trace_eventset_fill*( ),      *posix_trace_eventset_add*( ),      and
8336    *posix_trace_eventset_del*( ).

8337    A set of functions is provided devoted to the manipulation of the trace event type identifier and
8338    names for an active trace stream. All these functions require the trace stream identifier argument
8339    as the first parameter. The opacity of the trace event type identifier implies that the user cannot
8340    associate directly its well-known trace event name with the system-associated trace event type
8341    identifier.

8342    The *posix_trace_trid_eventid_open*( ) function allows the application to get the system trace event
8343    type identifier back from the system, given its well-known trace event name. This function is
8344    useful only when a controlling process needs to specify specific events to be filtered.

8345    The *posix_trace_eventid_get_name*( ) function allows the application to obtain a trace event name
8346    given its trace event type identifier. One possible use of this function is to identify the type of a
8347    trace event retrieved from the trace stream, and print it. The easiest way to implement this
8348    requirement, is to use a single trace event type map for all the processes whose maps are
8349    required to be identical. A more difficult way is to attempt to keep multiple maps identical at
8350    every call to *posix_trace_eventid_open*( ) and *posix_trace_trid_eventid_open*( ).

8351    *B.2.11.6  Rationale on Trace Events Type Filtering*

8352    The most basic rationale for runtime and pre-registration filtering (selection/rejection) of trace
8353    event types is to prevent choking of the trace collection facility, and/or overloading of the
8354    computer system. Any worthwhile trace facility can bring even the largest computer to its
8355    knees. Otherwise, everything would be recorded and filtered after the fact; it would be much
8356    simpler, but impractical.

8357    To achieve debugging, measurement, or whatever the purpose of tracing, the filtering of trace
8358    event types is an important part of trace analysis. Due to the fact that the trace events are put
8359    into a trace stream and probably logged afterwards into a file, different levels of filtering—that
8360    is, rejection of trace event types—are possible.

**Filtering of Trace Event Types Before Tracing**

8361

8362 This function, represented by the *posix_trace_set_filter*( ) function in IEEE Std 1003.1-2001 (see
8363 *posix_trace_set_filter*( )), selects, before or during tracing, the set of trace event types to be filtered
8364 out. It should be possible also (as OSF suggested in their ETAP trace specifications) to select the
8365 kernel trace event types to be traced in a system-wide fashion. These two functionalities are
8366 called the pre-filtering of trace event types.

8367 The restriction on the actual type used for the **trace_event_set_t** type is intended to guarantee
8368 that these objects can always be assigned, have their address taken, and be passed by value as
8369 parameters. It is not intended that this type be a structure including pointers to other data
8370 structures, as that could impact the portability of applications performing such operations. A
8371 reasonable implementation could be a structure containing an array of integer types.

**Filtering of Trace Event Types at Runtime**

8372

8373 It is possible to build this functionality using the *posix_trace_set_filter*( ) function. A privileged
8374 process or a privileged thread can get trace events from the trace stream of another process or
8375 thread, and thus specify the type of trace events to record into a file, using implementation-
8376 defined methods and interfaces. This functionality, called inline filtering of trace event types, is
8377 used for runtime analysis of trace streams.

**Post-Mortem Filtering of Trace Event Types**

8378

8379 The word ''post-mortem'' is used here to indicate that some unanticipated situation occurs
8380 during execution that does not permit a pre or inline filtering of trace events and that it is
8381 necessary to record all trace event types to have a chance to discover the problem afterwards.
8382 When the program stops, all the trace events recorded previously can be analyzed in order to
8383 find the solution. This functionality could be named the post-filtering of trace event types.

**Discussions about Trace Event Type-Filtering**

8384

8385 After long discussions with the parties involved in the process of defining the trace interface, it
8386 seems that the sensitivity to the filtering problem is different, but everybody agrees that the level
8387 of the overhead introduced during the tracing operation depends on the filtering method
8388 elected. If the time that it takes the trace event to be recorded can be neglected, the overhead
8389 introduced by the filtering process can be classified as follows:

8390 Pre-filtering        System and process/thread-level overhead

8391 Inline-filtering     Process/thread-level overhead

8392 Post-filtering       No overhead; done offline

8393 The pre-filtering could be named ''critical realtime'' filtering in the sense that the filtering of
8394 trace event type is manageable at the user level so the user can lower to a minimum the filtering
8395 overhead at some user selected level of priority for the inline filtering, or delay the filtering to
8396 after execution for the post-filtering. The counterpart of this solution is that the size of the trace
8397 stream must be sufficient to record all the trace events. The advantage of the pre-filtering is that
8398 the utilization of the trace stream is optimized.

8399 Only pre-filtering is defined by IEEE Std 1003.1-2001. However, great care must be taken in
8400 specifying pre-filtering, so that it does not impose unacceptable overhead. Moreover, it is
8401 necessary to isolate all the functionality relative to the pre-filtering.

8402 The result of this rationale is to define a new option, the Trace Event Filter option, not
8403 necessarily implemented in small realtime systems, where system overhead is minimized to the
8404 extent possible.

8405 *B.2.11.7 Tracing, pthread API*

8406   The objective to be able to control tracing for individual threads may be in conflict with the
8407   efficiency expected in threads with a *contentionscope* attribute of PTHREAD_SCOPE_PROCESS.
8408   For these threads, context switches from one thread that has tracing enabled to another thread
8409   that has tracing disabled may require a kernel call to inform the kernel whether it has to trace
8410   system events executed by that thread or not. For this reason, it was proposed that the ability to
8411   enable or disable tracing for PTHREAD_SCOPE_PROCESS threads be made optional, through
8412   the introduction of a Trace Scope Process option. A trace implementation which did not
8413   implement the Trace Scope Process option would not honor the tracing-state attribute of a
8414   thread with PTHREAD_SCOPE_PROCESS; it would, however, honor the tracing-state attribute
8415   of a thread with PTHREAD_SCOPE_SYSTEM. This proposal was rejected as:

8416     1.   Removing desired functionality (per-thread trace control)

8417     2.   Introducing counter-intuitive behavior for the tracing-state attribute

8418     3.   Mixing logically orthogonal ideas (thread scheduling and thread tracing)
8419          [Objective 4]

8420   Finally, to solve this complex issue, this API does not provide *pthread_gettracingstate*(),
8421   *pthread_settracingstate*(),     *pthread_attr_gettracingstate*(),     and     *pthread_attr_settracingstate*()
8422   interfaces. These interfaces force the thread implementation to add to the weight of the thread
8423   and cause a revision of the threads libraries, just to support tracing. Worse yet,
8424   *posix_trace_event*() must always test this per-thread variable even in the common case where it is
8425   not used at all. Per-thread tracing is easy to implement using existing interfaces where necessary;
8426   see the following example.

8427   **Example**

```
8428   /* Caution. Error checks omitted */
8429   static pthread_key_t my_key;
8430   static trace_event_id_t my_event_id;
8431   static pthread_once_t my_once = PTHREAD_ONCE_INIT;
8432   void my_init(void)
8433   {
8434       (void) pthread_key_create(&my_key, NULL);
8435       (void) posix_trace_eventid_open("my", &my_event_id);
8436   }
8437   int get_trace_flag(void)
8438   {
8439       pthread_once(&my_once, my_init);
8440       return (pthread_getspecific(my_key) != NULL);
8441   }
8442   void set_trace_flag(int f)
8443   {
8444       pthread_once(&my_once, my_init);
8445       pthread_setspecific(my_key, f? &my_event_id: NULL);
8446   }
8447   fn()
8448   {
8449       if (get_trace_flag())
8450           posix_trace_event(my_event_id, ...)
```

8451  }

8452  The above example does not implement third-party state setting.

8453  Lastly, per-thread tracing works poorly for threads with PTHREAD_SCOPE_PROCESS
8454  contention scope. These ''library'' threads have minimal interaction with the kernel and would
8455  have to explicitly set the attributes whenever they are context switched to a new kernel thread in
8456  order to trace system events. Such state was explicitly avoided in POSIX threads to keep
8457  PTHREAD_SCOPE_PROCESS threads lightweight.

8458  The reason that keeping PTHREAD_SCOPE_PROCESS threads lightweight is important is that
8459  such threads can be used not just for simple multi-processors but also for co-routine style
8460  programming (such as discrete event simulation) without inventing a new threads paradigm.
8461  Adding extra runtime cost to thread context switches will make using POSIX threads less
8462  attractive in these situations.

### B.2.11.8  Rationale on Triggering

8463

8464  The ability to start or stop tracing based on the occurrence of specific trace event types has been
8465  proposed as a parallel to similar functionality appearing in logic analyzers. Such triggering, in
8466  order to be very useful, should be based not only on the trace event type, but on trace event-
8467  specific data, including tests of user-specified fields for matching or threshold values.

8468  Such a facility is unnecessary where the buffering of the stream is not a constraint, since such
8469  checks can be performed offline during post-mortem analysis.

8470  For example, a large system could incorporate a daemon utility to collect the trace records from
8471  memory buffers and spool them to secondary storage for later analysis. In the instances where
8472  resources are truly limited, such as embedded applications, the application incorporation of
8473  application code to test the circumstances of a trace event and call the trace point only if needed
8474  is usually straightforward.

8475  For performance reasons, the *posix_trace_event*( ) function should be implemented using a macro,
8476  so if the trace is inactive, the trace event point calls are latent code and must cost no more than a
8477  scalar test.

8478  The API proposed in IEEE Std 1003.1-2001 does not include any triggering functionality.

### B.2.11.9  Rationale on Timestamp Clock

8479

8480  It has been suggested that the tracing mechanism should include the possibility of specifying the
8481  clock to be used in timestamping the trace events. When application trace events must be
8482  correlated to remote trace events, such a facility could provide a global time reference not
8483  available from a local clock. Further, the application may be driven by timers based on a clock
8484  different from that used for the timestamp, and the correlation of the trace to those untraced
8485  timer activities could be an important part of the analysis of the application.

8486  However, the tracing mechanism needs to be fast and just the provision of such an option can
8487  materially affect its performance. Leaving aside the performance costs of reading some clocks,
8488  this notion is also ill-defined when kernel trace events are to be traced by two applications
8489  making use of different tracing clocks. This can even happen within a single application where
8490  different parts of the application are served by different clocks. Another complication can occur
8491  when a clock is maintained strictly at the user level and is unavailable at the kernel level.

8492  It is felt that the benefits of a selectable trace clock do not match its costs. Applications that wish
8493  to correlate clocks other than the default tracing clock can include trace events with sample
8494  values of those other clocks, allowing correlation of timestamps from the various independent
8495  clocks. In any case, such a technique would be required when applications are sensitive to

8496  multiple clocks.

8498  The analysis of the dynamic behavior of the trace mechanism shows that different overrun
8499  conditions may occur. The API must provide a means to manage such conditions in a portable
8500  way.

8501  **Overrun in Trace Streams Initialized with POSIX_TRACE_LOOP Policy**

8502  In this case, the user of the trace mechanism is interested in using the trace stream with
8503  POSIX_TRACE_LOOP policy to record trace events continuously, but ideally without losing any
8504  trace events. The online analyzer process must get the trace events at a mean speed equivalent to
8505  the recording speed. Should the trace stream become full, a trace stream overrun occurs. This
8506  condition is detected by getting the status of the active trace stream (function
8507  *posix_trace_get_status*( )) and looking at the member *posix_stream_overrun_status* of the read
8508  **posix_stream_status** structure. In addition, two predefined trace event types are defined:

1.  The beginning of a trace overflow, to locate the beginning of an overflow when reading a
    trace stream

2.  The end of a trace overflow, to locate the end of an overflow, when reading a trace stream

8512  As a timestamp is associated with these predefined trace events, it is possible to know the
8513  duration of the overflow.

8514  **Overrun in Dumping Trace Streams into Trace Logs**

8515  The user lets the trace mechanism dump the trace stream initialized with
8516  POSIX_TRACE_FLUSH policy automatically into a trace log. If the dump operation is slower
8517  than the recording of trace events, the trace stream can overrun. This condition is detected by
8518  getting the status of the active trace stream (function *posix_trace_get_status*( )) and looking at the
8519  member *posix_log_overrun_status* of the read **posix_stream_status** structure. This overrun
8520  indicates that the trace mechanism is not able to operate in this mode at this speed. It is the
8521  responsibility of the user to modify one of the trace parameters (the stream size or the trace
8522  event type filter, for instance) to avoid such overrun conditions, if overruns are to be prevented.
8523  The same already predefined trace event types (see **Overrun in Trace Streams Initialized with**
8524  **POSIX_TRACE_LOOP Policy**) are used to detect and to know the duration of an overflow.

8525  **Reading an Active Trace Stream**

8526  Although this trace API allows one to read an active trace stream with log while it is tracing, this
8527  feature can lead to false overflow origin interpretation: the trace log or the reader of the trace
8528  stream. Reading from an active trace stream with log is thus non-portable, and has been left
8529  unspecified.

8530  ## B.2.12  Data Types

8531  The requirement that additional types defined in this section end in ''_t'' was prompted by the
8532  problem of name space pollution. It is difficult to define a type (where that type is not one
8533  defined by IEEE Std 1003.1-2001) in one header file and use it in another without adding symbols
8534  to the name space of the program. To allow implementors to provide their own types, all
8535  conforming applications are required to avoid symbols ending in ''_t'', which permits the
8536  implementor to provide additional types. Because a major use of types is in the definition of
8537  structure members, which can (and in many cases must) be added to the structures defined in
8538  IEEE Std 1003.1-2001, the need for additional types is compelling.

8539    The types, such as **ushort** and **ulong**, which are in common usage, are not defined in
8540    IEEE Std 1003.1-2001 (although **ushort_t** would be permitted as an extension). They can be
8541    added to **<sys/types.h>** using a feature test macro (see Section B.2.2.1 (on page 85)). A suggested
8542    symbol for these is _SYSIII. Similarly, the types like **u_short** would probably be best controlled
8543    by _BSD.

8544    Some of these symbols may appear in other headers; see Section B.2.2.2 (on page 86).

8545    **dev_t**      This type may be made large enough to accommodate host-locality considerations
8546                of networked systems.

8547                This type must be arithmetic. Earlier proposals allowed this to be non-arithmetic
8548                (such as a structure) and provided a *samefile*( ) function for comparison.

8549    **gid_t**      Some implementations had separated **gid_t** from **uid_t** before POSIX.1 was
8550                completed. It would be difficult for them to coalesce them when it was
8551                unnecessary. Additionally, it is quite possible that user IDs might be different from
8552                group IDs because the user ID might wish to span a heterogeneous network,
8553                where the group ID might not.

8554                For current implementations, the cost of having a separate **gid_t** will be only
8555                lexical.

8556    **mode_t**     This type was chosen so that implementations could choose the appropriate
8557                integer type, and for compatibility with the ISO C standard. 4.3 BSD uses
8558                **unsigned short** and the SVID uses **ushort**, which is the same. Historically, only the
8559                low-order sixteen bits are significant.

8560    **nlink_t**     This type was introduced in place of **short** for *st_nlink* (see the **<sys/stat.h>** header)
8561                in response to an objection that **short** was too small.

8562    **off_t**      This type is used only in *lseek*( ), *fcntl*( ), and **<sys/stat.h>**. Many implementations
8563                would have difficulties if it were defined as anything other than **long**. Requiring
8564                an integer type limits the capabilities of *lseek*( ) to four gigabytes. The ISO C
8565                standard supplies routines that use larger types; see *fgetpos*( ) and *fsetpos*( ). XSI-
8566                conformant systems provide the *fseeko*( ) and *ftello*( ) functions that use larger
8567                types.

8568    **pid_t**      The inclusion of this symbol was controversial because it is tied to the issue of the
8569                representation of a process ID as a number. From the point of view of a
8570                conforming application, process IDs should be ‘‘magic cookies’’[1] that are produced
8571                by calls such as *fork*( ), used by calls such as *waitpid*( ) or *kill*( ), and not otherwise
8572                analyzed (except that the sign is used as a flag for certain operations).

8573                The concept of a {PID_MAX} value interacted with this in early proposals. Treating
8574                process IDs as an opaque type both removes the requirement for {PID_MAX} and
8575                allows systems to be more flexible in providing process IDs that span a large range
8576                of values, or a small one.

8577                Since the values in **uid_t**, **gid_t**, and **pid_t** will be numbers generally, and
8578                potentially both large in magnitude and sparse, applications that are based on

8579    _____
8580    1. An historical term meaning: ‘‘An opaque object, or token, of determinate size, whose significance is known only to the entity
8581       which created it. An entity receiving such a token from the generating entity may only make such use of the ‘cookie’ as is defined
8582       and permitted by the supplying entity.’’

8583            arrays of objects of this type are unlikely to be fully portable in any case. Solutions
8584            that treat them as magic cookies will be portable.

8585            {CHILD_MAX} precludes the possibility of a ''toy implementation'', where there
8586            would only be one process.

8587    **ssize_t**    This is intended to be a signed analog of **size_t**. The wording is such that an
8588            implementation may either choose to use a longer type or simply to use the signed
8589            version of the type that underlies **size_t**. All functions that return **ssize_t** (*read*()
8590            and *write*()) describe as ''implementation-defined'' the result of an input exceeding
8591            {SSIZE_MAX}. It is recognized that some implementations might have **int**s that
8592            are smaller than **size_t**. A conforming application would be constrained not to
8593            perform I/O in pieces larger than {SSIZE_MAX}, but a conforming application
8594            using extensions would be able to use the full range if the implementation
8595            provided an extended range, while still having a single type-compatible interface.

8596            The symbols **size_t** and **ssize_t** are also required in **<unistd.h>** to minimize the
8597            changes needed for calls to *read*() and *write*(). Implementors are reminded that it
8598            must be possible to include both **<sys/types.h>** and **<unistd.h>** in the same
8599            program (in either order) without error.

8600    **uid_t**    Before the addition of this type, the data types used to represent these values
8601            varied throughout early proposals. The **<sys/stat.h>** header defined these values as
8602            type **short**, the **<passwd.h>** file (now **<pwd.h>** and **<grp.h>**) used an **int**, and
8603            *getuid*() returned an **int**. In response to a strong objection to the inconsistent
8604            definitions, all the types were switched to **uid_t**.

8605            In practice, those historical implementations that use varying types of this sort can
8606            typedef **uid_t** to **short** with no serious consequences.

8607            The problem associated with this change concerns object compatibility after
8608            structure size changes. Since most implementations will define **uid_t** as a short, the
8609            only substantive change will be a reduction in the size of the **passwd** structure.
8610            Consequently, implementations with an overriding concern for object
8611            compatibility can pad the structure back to its current size. For that reason, this
8612            problem was not considered critical enough to warrant the addition of a separate
8613            type to POSIX.1.

8614            The types **uid_t** and **gid_t** are magic cookies. There is no {UID_MAX} defined by
8615            POSIX.1, and no structure imposed on **uid_t** and **gid_t** other than that they be
8616            positive arithmetic types. (In fact, they could be **unsigned char**.) There is no
8617            maximum or minimum specified for the number of distinct user or group IDs.

8618  **B.3      System Interfaces**

8619      See the RATIONALE sections on the individual reference pages.

8620  **B.3.1    Examples for Spawn**

8621      The following long examples are provided in the Rationale (Informative) volume of
8622      IEEE Std 1003.1-2001 as a supplement to the reference page for *posix_spawn*().

8623      **Example Library Implementation of Spawn**

8624      The *posix_spawn*() or *posix_spawnp*() functions provide the following:

8625  •  Simply start a process executing a process image. This is the simplest application for process
8626         creation, and it may cover most executions of *fork*().

8627  •  Support I/O redirection, including pipes.

8628  •  Run the child under a user and group ID in the domain of the parent.

8629  •  Run the child at any priority in the domain of the parent.

8630      The *posix_spawn*() or *posix_spawnp*() functions do not cover every possible use of the *fork*()
8631      function, but they do span the common applications: typical use by a shell and a login utility.

8632      The price for an application is that before it calls *posix_spawn*() or *posix_spawnp*(), the parent
8633      must adjust to a state that *posix_spawn*() or *posix_spawnp*() can map to the desired state for the
8634      child. Environment changes require the parent to save some of its state and restore it afterwards.
8635      The effective behavior of a successful invocation of *posix_spawn*() is as if the operation were
8636      implemented with POSIX operations as follows:

```
8637      #include <sys/types.h>
8638      #include <stdlib.h>
8639      #include <stdio.h>
8640      #include <unistd.h>
8641      #include <sched.h>
8642      #include <fcntl.h>
8643      #include <signal.h>
8644      #include <errno.h>
8645      #include <string.h>
8646      #include <signal.h>

8647      /* #include <spawn.h> */
8648      /*******************************************/
8649      /* Things that could be defined in spawn.h */
8650      /*******************************************/
8651      typedef struct
8652      {
8653          short posix_attr_flags;
8654      #define POSIX_SPAWN_SETPGROUP        0x1
8655      #define POSIX_SPAWN_SETSIGMASK       0x2
8656      #define POSIX_SPAWN_SETSIGDEF        0x4
8657      #define POSIX_SPAWN_SETSCHEDULER     0x8
8658      #define POSIX_SPAWN_SETSCHEDPARAM    0x10
8659      #define POSIX_SPAWN_RESETIDS         0x20
8660          pid_t posix_attr_pgroup;
8661          sigset_t posix_attr_sigmask;
8662          sigset_t posix_attr_sigdefault;
```

```
8663              int posix_attr_schedpolicy;
8664              struct sched_param posix_attr_schedparam;
8665          }   posix_spawnattr_t;

8666          typedef char *posix_spawn_file_actions_t;

8667          int posix_spawn_file_actions_init(
8668              posix_spawn_file_actions_t *file_actions);
8669          int posix_spawn_file_actions_destroy(
8670              posix_spawn_file_actions_t *file_actions);
8671          int posix_spawn_file_actions_addclose(
8672              posix_spawn_file_actions_t *file_actions, int fildes);
8673          int posix_spawn_file_actions_adddup2(
8674              posix_spawn_file_actions_t *file_actions, int fildes,
8675              int newfildes);
8676          int posix_spawn_file_actions_addopen(
8677              posix_spawn_file_actions_t *file_actions, int fildes,
8678              const char *path, int oflag, mode_t mode);
8679          int posix_spawnattr_init(posix_spawnattr_t *attr);
8680          int posix_spawnattr_destroy(posix_spawnattr_t *attr);
8681          int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8682              short *lags);
8683          int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags);
8684          int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8685              pid_t *pgroup);
8686          int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup);
8687          int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8688              int *schedpolicy);
8689          int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8690              int schedpolicy);
8691          int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8692              struct sched_param *schedparam);
8693          int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8694              const struct sched_param *schedparam);
8695          int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
8696              sigset_t *sigmask);
8697          int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8698              const sigset_t *sigmask);
8699          int posix_spawnattr_getdefault(const posix_spawnattr_t *attr,
8700              sigset_t *sigdefault);
8701          int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8702              const sigset_t *sigdefault);
8703          int posix_spawn(pid_t *pid, const char *path,
8704              const posix_spawn_file_actions_t *file_actions,
8705              const posix_spawnattr_t *attrp, char *const argv[],
8706              char *const envp[]);
8707          int posix_spawnp(pid_t *pid, const char *file,
8708              const posix_spawn_file_actions_t *file_actions,
8709              const posix_spawnattr_t *attrp, char *const argv[],
8710              char *const envp[]);

8711          /*****************************************/
8712          /* Example posix_spawn() library routine */
8713          /*****************************************/
```

```
8714        int posix_spawn(pid_t *pid,
8715            const char *path,
8716            const posix_spawn_file_actions_t *file_actions,
8717            const posix_spawnattr_t *attrp,
8718            char *const argv[],
8719            char *const envp[])
8720        {
8721            /* Create process */
8722            if ((*pid = fork()) == (pid_t) 0)
8723            {
8724                /* This is the child process */
8725                /* Worry about process group */
8726                if (attrp->posix_attr_flags & POSIX_SPAWN_SETPGROUP)
8727                {
8728                    /* Override inherited process group */
8729                    if (setpgid(0, attrp->posix_attr_pgroup) != 0)
8730                    {
8731                        /* Failed */
8732                        exit(127);
8733                    }
8734                }
8735                /* Worry about thread signal mask */                              |
8736                if (attrp->posix_attr_flags & POSIX_SPAWN_SETSIGMASK)
8737                {
8738                    /* Set the signal mask (can't fail) */
8739                    sigprocmask(SIG_SETMASK, &attrp->posix_attr_sigmask, NULL);
8740                }
8741                /* Worry about resetting effective user and group IDs */
8742                if (attrp->posix_attr_flags & POSIX_SPAWN_RESETIDS)
8743                {
8744                    /* None of these can fail for this case. */
8745                    setuid(getuid());
8746                    setgid(getgid());
8747                }
8748                /* Worry about defaulted signals */
8749                if (attrp->posix_attr_flags & POSIX_SPAWN_SETSIGDEF)
8750                {
8751                    struct sigaction deflt;
8752                    sigset_t all_signals;
8753                    int s;
8754                    /* Construct default signal action */
8755                    deflt.sa_handler = SIG_DFL;
8756                    deflt.sa_flags = 0;
8757                    /* Construct the set of all signals */
8758                    sigfillset(&all_signals);
8759                    /* Loop for all signals */
8760                    for (s = 0; sigismember(&all_signals, s); s++)
8761                    {
8762                        /* Signal to be defaulted? */
```

```
8763                        if (sigismember(&attrp->posix_attr_sigdefault, s))
8764                        {
8765                            /* Yes; default this signal */
8766                            if (sigaction(s, &deflt, NULL) == −1)
8767                            {
8768                                /* Failed */
8769                                exit(127);
8770                            }
8771                        }
8772                    }
8773                }
8774            /* Worry about the fds if they are to be mapped */
8775            if (file_actions != NULL)
8776            {
8777                /* Loop for all actions in object file_actions */
8778                /* (implementation dives beneath abstraction) */
8779                char *p = *file_actions;
8780                while (*p != '\0')
8781                {
8782                    if (strncmp(p, "close(", 6) == 0)
8783                    {
8784                        int fd;
8785                        if (sscanf(p + 6, "%d)", &fd) != 1)
8786                        {
8787                            exit(127);
8788                        }
8789                        if (close(fd) == −1)
8790                            exit(127);
8791                    }
8792                    else if (strncmp(p, "dup2(", 5) == 0)
8793                    {
8794                        int fd, newfd;
8795                        if (sscanf(p + 5, "%d,%d)", &fd, &newfd) != 2)
8796                        {
8797                            exit(127);
8798                        }
8799                        if (dup2(fd, newfd) == −1)
8800                            exit(127);
8801                    }
8802                    else if (strncmp(p, "open(", 5) == 0)
8803                    {
8804                        int fd, oflag;
8805                        mode_t mode;
8806                        int tempfd;
8807                        char path[1000];    /* Should be dynamic */
8808                        char *q;
8809                        if (sscanf(p + 5, "%d,", &fd) != 1)
8810                        {
8811                            exit(127);
8812                        }
```

```
8813                            p = strchr(p, ',') + 1;
8814                            q = strchr(p, '*');
8815                            if (q == NULL)
8816                                exit(127);
8817                            strncpy(path, p, q - p);
8818                            path[q - p] = '\0';
8819                            if (sscanf(q + 1, "%o,%o)", &oflag, &mode) != 2)
8820                            {
8821                                exit(127);
8822                            }
8823                            if (close(fd) == −1)
8824                            {
8825                                if (errno != EBADF)
8826                                    exit(127);
8827                            }
8828                            tempfd = open(path, oflag, mode);
8829                            if (tempfd == −1)
8830                                exit(127);
8831                            if (tempfd != fd)
8832                            {
8833                                if (dup2(tempfd, fd) == −1)
8834                                {
8835                                    exit(127);
8836                                }
8837                                if (close(tempfd) == −1)
8838                                {
8839                                    exit(127);
8840                                }
8841                            }
8842                        }
8843                        else
8844                        {
8845                            exit(127);
8846                        }
8847                        p = strchr(p, ')') + 1;
8848                    }
8849                }
8850            /* Worry about setting new scheduling policy and parameters */
8851            if (attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDULER)
8852            {
8853                if (sched_setscheduler(0, attrp->posix_attr_schedpolicy,
8854                    &attrp->posix_attr_schedparam) == −1)
8855                {
8856                    exit(127);
8857                }
8858            }
8859            /* Worry about setting only new scheduling parameters */
8860            if (attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDPARAM)
8861            {
8862                if (sched_setparam(0, &attrp->posix_attr_schedparam) == −1)
8863                {
```

```
8864                    exit(127);
8865                }
8866            }

8867            /* Now execute the program at path */
8868            /* Any fd that still has FD_CLOEXEC set will be closed */
8869            execve(path, argv, envp);
8870            exit(127);                  /* exec failed */
8871        }
8872        else
8873        {
8874            /* This is the parent (calling) process */
8875            if (*pid == (pid_t) - 1)
8876                return errno;
8877            return 0;
8878        }
8879    }
8880    /*****************************************************/
8881    /* Here is a crude but effective implementation of the */
8882    /* file action object operators which store actions as */
8883    /* concatenated token-separated strings.             */
8884    /*****************************************************/
8885    /* Create object with no actions. */
8886    int posix_spawn_file_actions_init(
8887        posix_spawn_file_actions_t *file_actions)
8888    {
8889        *file_actions = malloc(sizeof(char));
8890        if (*file_actions == NULL)
8891            return ENOMEM;
8892        strcpy(*file_actions, "");
8893        return 0;
8894    }

8895    /* Free object storage and make invalid. */
8896    int posix_spawn_file_actions_destroy(
8897        posix_spawn_file_actions_t *file_actions)
8898    {
8899        free(*file_actions);
8900        *file_actions = NULL;
8901        return 0;
8902    }

8903    /* Add a new action string to object. */
8904    static int add_to_file_actions(
8905        posix_spawn_file_actions_t *file_actions, char *new_action)
8906    {
8907        *file_actions = realloc
8908        (*file_actions, strlen(*file_actions) + strlen(new_action) + 1);
8909        if (*file_actions == NULL)
8910            return ENOMEM;
8911        strcat(*file_actions, new_action);
8912        return 0;
8913    }
```

```
8914          /* Add a close action to object. */
8915          int posix_spawn_file_actions_addclose(
8916              posix_spawn_file_actions_t *file_actions, int fildes)
8917          {
8918              char temp[100];

8919              sprintf(temp, "close(%d)", fildes);
8920              return add_to_file_actions(file_actions, temp);
8921          }

8922          /* Add a dup2 action to object. */
8923          int posix_spawn_file_actions_adddup2(
8924              posix_spawn_file_actions_t *file_actions, int fildes,
8925              int newfildes)
8926          {
8927              char temp[100];

8928              sprintf(temp, "dup2(%d,%d)", fildes, newfildes);
8929              return add_to_file_actions(file_actions, temp);
8930          }

8931          /* Add an open action to object. */
8932          int posix_spawn_file_actions_addopen(
8933              posix_spawn_file_actions_t *file_actions, int fildes,
8934              const char *path, int oflag, mode_t mode)
8935          {
8936              char temp[100];

8937              sprintf(temp, "open(%d,%s*%o,%o)", fildes, path, oflag, mode);
8938              return add_to_file_actions(file_actions, temp);
8939          }
8940          /****************************************************/
8941          /* Here is a crude but effective implementation of the */
8942          /* spawn attributes object functions which manipulate  */
8943          /* the individual attributes.                          */
8944          /****************************************************/
8945          /* Initialize object with default values. */
8946          int posix_spawnattr_init(posix_spawnattr_t *attr)
8947          {
8948              attr->posix_attr_flags = 0;
8949              attr->posix_attr_pgroup = 0;
8950              /* Default value of signal mask is the parent's signal mask; */
8951              /* other values are also allowed */
8952              sigprocmask(0, NULL, &attr->posix_attr_sigmask);
8953              sigemptyset(&attr->posix_attr_sigdefault);
8954              /* Default values of scheduling attr inherited from the parent; */
8955              /* other values are also allowed */
8956              attr->posix_attr_schedpolicy = sched_getscheduler(0);
8957              sched_getparam(0, &attr->posix_attr_schedparam);
8958              return 0;
8959          }

8960          int posix_spawnattr_destroy(posix_spawnattr_t *attr)
8961          {
8962              /* No action needed */
```

```
8963                return 0;
8964            }

8965        int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8966            short *flags)
8967        {
8968            *flags = attr->posix_attr_flags;
8969            return 0;
8970        }

8971        int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags)
8972        {
8973            attr->posix_attr_flags = flags;
8974            return 0;
8975        }

8976        int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8977            pid_t *pgroup)
8978        {
8979            *pgroup = attr->posix_attr_pgroup;
8980            return 0;
8981        }

8982        int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup)
8983        {
8984            attr->posix_attr_pgroup = pgroup;
8985            return 0;
8986        }

8987        int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8988            int *schedpolicy)
8989        {
8990            *schedpolicy = attr->posix_attr_schedpolicy;
8991            return 0;
8992        }

8993        int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8994            int schedpolicy)
8995        {
8996            attr->posix_attr_schedpolicy = schedpolicy;
8997            return 0;
8998        }

8999        int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
9000            struct sched_param *schedparam)
9001        {
9002            *schedparam = attr->posix_attr_schedparam;
9003            return 0;
9004        }

9005        int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
9006            const struct sched_param *schedparam)
9007        {
9008            attr->posix_attr_schedparam = *schedparam;
9009            return 0;
9010        }
```

```
9011     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
9012         sigset_t *sigmask)
9013     {
9014         *sigmask = attr->posix_attr_sigmask;
9015         return 0;
9016     }

9017     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
9018         const sigset_t *sigmask)
9019     {
9020         attr->posix_attr_sigmask = *sigmask;
9021         return 0;
9022     }

9023     int posix_spawnattr_getsigdefault(const posix_spawnattr_t *attr,
9024         sigset_t *sigdefault)
9025     {
9026         *sigdefault = attr->posix_attr_sigdefault;
9027         return 0;
9028     }

9029     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
9030         const sigset_t *sigdefault)
9031     {
9032         attr->posix_attr_sigdefault = *sigdefault;
9033         return 0;
9034     }
```

9035     **I/O Redirection with Spawn**

9036     I/O redirection with *posix_spawn*() or *posix_spawnp*() is accomplished by crafting a *file_actions*
9037     argument to effect the desired redirection. Such a redirection follows the general outline of the
9038     following example:

```
9039     /* To redirect new standard output (fd 1) to a file, */
9040     /* and redirect new standard input (fd 0) from my fd socket_pair[1], */
9041     /* and close my fd socket_pair[0] in the new process. */
9042     posix_spawn_file_actions_t file_actions;
9043     posix_spawn_file_actions_init(&file_actions);
9044     posix_spawn_file_actions_addopen(&file_actions, 1, "newout", ...);
9045     posix_spawn_file_actions_dup2(&file_actions, socket_pair[1], 0);
9046     posix_spawn_file_actions_close(&file_actions, socket_pair[0]);
9047     posix_spawn_file_actions_close(&file_actions, socket_pair[1]);
9048     posix_spawn(..., &file_actions, ...);
9049     posix_spawn_file_actions_destroy(&file_actions);
```

9050    **Spawning a Process Under a New User ID**

9051    Spawning a process under a new user ID follows the outline shown in the following example:

```
9052    Save = getuid();
9053    setuid(newid);
9054    posix_spawn(...);
9055    setuid(Save);
```

9056

9057

# Rationale (Informative)

9058 **Part C:**

9059 **Shell and Utilities**

9060 *The Open Group*
9061 *The Institute of Electrical and Electronics Engineers, Inc.*

*Appendix C*

# Rationale for Shell and Utilities

## C.1    Introduction

### C.1.1    Scope

Refer to Section A.1.1 (on page 3).

### C.1.2    Conformance

Refer to Section A.2 (on page 9).

### C.1.3    Normative References

There is no additional rationale provided for this section.

### C.1.4    Change History

The change history is provided as an informative section, to track changes from previous issues of IEEE Std 1003.1-2001.

The following sections describe changes made to the Shell and Utilities volume of IEEE Std 1003.1-2001 since Issue 5 of the base document. The CHANGE HISTORY section for each utility describes technical changes made to that utility from Issue 5. Changes between earlier issues of the base document and Issue 5 are not included.

The change history between Issue 5 and Issue 6 also lists the changes since the ISO POSIX-2:1993 standard.

**Changes from Issue 5 to Issue 6 (IEEE Std 1003.1-2001)**

The following list summarizes the major changes that were made in the Shell and Utilities volume of IEEE Std 1003.1-2001 from Issue 5 to Issue 6:

- This volume of IEEE Std 1003.1-2001 is extensively revised so that it can be both an IEEE POSIX Standard and an Open Group Technical Standard.

- The terminology has been reworked to meet the style requirements.

- Shading notation and margin codes are introduced for identification of options within the volume.

- This volume of IEEE Std 1003.1-2001 is updated to mandate support of FIPS 151-2. The following changes were made:

— Support is mandated for the capabilities associated with the following symbolic constants:

        _POSIX_CHOWN_RESTRICTED
        _POSIX_JOB_CONTROL
        _POSIX_SAVED_IDS

— In the environment for the login shell, the environment variables *LOGNAME* and *HOME* shall be defined and have the properties described in the Base Definitions volume of

9096          IEEE Std 1003.1-2001, Chapter 7, Locale.

9097    • This volume of IEEE Std 1003.1-2001 is updated to align with some features of the Single
9098      UNIX Specification.

9099    • A new section on Utility Limits is added.

9100    • A section on the Relationships to Other Documents is added.

9101    • Concepts and definitions have been moved to a separate volume.

9102    • A RATIONALE section is added to each reference page.

9103    • The *c99* utility is added as a replacement for *c89*, which is withdrawn in this issue.

9104    • IEEE Std 1003.2d-1994 is incorporated, adding the *qalter*, *qdel*, *qhold*, *qmove*, *qmsg*, *qrerun*, *qrls*,
9105      *qselect*, *qsig*, *qstat*, and *qsub* utilities.

9106    • IEEE P1003.2b draft standard is incorporated, making extensive updates and adding the *iconv*
9107      utility.

9108    • IEEE PASC Interpretations are applied.

9109    • The Open Group's corrigenda and resolutions are applied.

9110    **New Features in Issue 6**

9111    The following table lists the new utilities introduced since the ISO POSIX-2: 1993 standard (as
9112    modified by IEEE Std 1003.2d-1994). Apart from the *c99* and *iconv* utilities, these are all part of
9113    the XSI extension.

9114

| New Utilities in Issue 6 | | | | | | | |
|---|---|---|---|---|---|---|---|
| *admin* | *compress* | *gencat* | *ipcrm* | *nl* | *tsort* | *unlink* | *val* |
| *c99* | *cxref* | *get* | *ipcs* | *prs* | *ulimit* | *uucp* | *what* |
| *cal* | *delta* | *hash* | *link* | *sact* | *uncompress* | *uustat* | *zcat* |
| *cflow* | *fuser* | *iconv* | *m4* | *sccs* | *unget* | *uux* | |

9120  **C.1.5    Terminology**

9121    Refer to Section A.1.4 (on page 5).

9122  **C.1.6    Definitions**

9123    Refer to Section A.3 (on page 13).

9124  **C.1.7    Relationship to Other Documents**

9125  *C.1.7.1    System Interfaces*

9126    It has been pointed out that the Shell and Utilities volume of IEEE Std 1003.1-2001 assumes that
9127    a great deal of functionality from the System Interfaces volume of IEEE Std 1003.1-2001 is
9128    present, but never states exactly how much (and strictly does not need to since both are
9129    mandated on a conforming system). This section is an attempt to clarify the assumptions.

9130            **File Removal**

9131            This is intended to be a summary of the *unlink*( ) and *rmdir*( ) requirements. Note that it is
9132            possible using the *unlink*( ) function for item 4. to occur.

9133    *C.1.7.2*   *Concepts Derived from the ISO C Standard*

9134            This section was introduced to address the issue that there was insufficient detail presented by
9135            such utilities as *awk* or *sh* about their procedural control statements and their methods of
9136            performing arithmetic functions.

9137            The ISO C standard was selected as a model because most historical implementations of the
9138            standard utilities were written in C. Thus, it was more likely that they would act in the desired
9139            manner without modification.

9140            Using the ISO C standard is primarily a notational convenience so that the many procedural
9141            languages in the Shell and Utilities volume of IEEE Std 1003.1-2001 would not have to be
9142            rigorously described in every aspect. Its selection does not require that the standard utilities be
9143            written in Standard C; they could be written in Common Usage C, Ada, Pascal, assembler
9144            language, or anything else.

9145            The sizes of the various numeric values refer to C-language data types that are allowed to be
9146            different sizes by the ISO C standard. Thus, like a C-language application, a shell application
9147            cannot rely on their exact size. However, it can rely on their minimum sizes expressed in the
9148            ISO C standard, such as {LONG_MAX} for a **long** type.

9149            The behavior on overflow is undefined for ISO C standard arithmetic. Therefore, the standard
9150            utilities can use ''bignum'' representation for integers so that there is no fixed maximum unless
9151            otherwise stated in the utility description. Similarly, standard utilities can use infinite-precision
9152            representations for floating-point arithmetic, as long as these representations exceed the ISO C
9153            standard requirements.

9154            This section addresses only the issue of semantics; it is not intended to specify syntax. For
9155            example, the ISO C standard requires that 0L be recognized as an integer constant equal to zero,
9156            but utilities such as *awk* and *sh* are not required to recognize 0L (though they are allowed to, as
9157            an extension).

9158            The ISO C standard requires that a C compiler must issue a diagnostic for constants that are too
9159            large to represent. Most standard utilities are not required to issue these diagnostics; for
9160            example, the command:

9161                ```
                diff −C 2147483648 file1 file2
                ```

9162            has undefined behavior, and the *diff* utility is not required to issue a diagnostic even if the
9163            number 2 147 483 648 cannot be represented.

9164    **C.1.8    Portability**

9165            Refer to Section A.1.5 (on page 8).

9166    *C.1.8.1*   *Codes*

9167            Refer to Section A.1.5.1 (on page 8).

9168 **C.1.9      Utility Limits**

9169    This section grew out of an idea that originated with the original POSIX.1, in the tables of system
9170    limits for the *sysconf*( ) and *pathconf*( ) functions. The idea being that a conforming application
9171    can be written to use the most restrictive values that a minimal system can provide, but it should
9172    not have to. The values provided represent compromises so that some vendors can use
9173    historically limited versions of UNIX system utilities. They are the highest values that a strictly
9174    conforming application can assume, given no other information.

9175    However, by using the *getconf* utility or the *sysconf*( ) function, the elegant application can be
9176    tailored to more liberal values on some of the specific instances of specific implementations.

9177    There is no explicitly stated requirement that an implementation provide finite limits for any of
9178    these numeric values; the implementation is free to provide essentially unbounded capabilities
9179    (where it makes sense), stopping only at reasonable points such as {ULONG_MAX} (from the
9180    ISO C standard). Therefore, applications desiring to tailor themselves to the values on a
9181    particular implementation need to be ready for possibly huge values; it may not be a good idea
9182    to allocate blindly a buffer for an input line based on the value of {LINE_MAX}, for instance.
9183    However, unlike the System Interfaces volume of IEEE Std 1003.1-2001, there is no set of limits
9184    that return a special indication meaning ''unbounded''. The implementation should always
9185    return an actual number, even if the number is very large.

9186    The statement:

9187        ''It is not guaranteed that the application ...''

9188    is an indication that many of these limits are designed to ensure that implementors design their
9189    utilities without arbitrary constraints related to unimaginative programming. There are certainly
9190    conditions under which combinations of options can cause failures that would not render an
9191    implementation non-conforming. For example, {EXPR_NEST_MAX} and {ARG_MAX} could
9192    collide when expressions are large; combinations of {BC_SCALE_MAX} and {BC_DIM_MAX}
9193    could exceed virtual memory.

9194    In the Shell and Utilities volume of IEEE Std 1003.1-2001, the notion of a limit being guaranteed
9195    for the process lifetime, as it is in the System Interfaces volume of IEEE Std 1003.1-2001, is not as
9196    useful to a shell script. The *getconf* utility is probably a process itself, so the guarantee would be
9197    without value. Therefore, the Shell and Utilities volume of IEEE Std 1003.1-2001 requires the
9198    guarantee to be for the session lifetime. This will mean that many vendors will either return very
9199    conservative values or possibly implement *getconf* as a built-in.

9200    It may seem confusing to have limits that apply only to a single utility grouped into one global
9201    section. However, the alternative, which would be to disperse them out into their utility
9202    description sections, would cause great difficulty when *sysconf*( ) and *getconf* were described.
9203    Therefore, the standard developers chose the global approach.

9204    Each language binding could provide symbol names that are slightly different from those shown
9205    here. For example, the C-Language Binding option adds a leading underscore to the symbols as a
9206    prefix.

9207    The following comments describe selection criteria for the symbols and their values:

9208    {ARG_MAX}
9209        This is defined by the System Interfaces volume of IEEE Std 1003.1-2001. Unfortunately, it is
9210        very difficult for a conforming application to deal with this value, as it does not know how
9211        much of its argument space is being consumed by the environment variables of the user.

9212    {BC_BASE_MAX}
9213    {BC_DIM_MAX}
9214    {BC_SCALE_MAX}
9215        These were originally one value, {BC_SCALE_MAX}, but it was unreasonable to link all
9216        three concepts into one limit.

9217    {CHILD_MAX}
9218        This is defined by the System Interfaces volume of IEEE Std 1003.1-2001.

9219    {COLL_WEIGHTS_MAX}
9220        The weights assigned to **order** can be considered as ''passes'' through the collation
9221        algorithm.

9222    {EXPR_NEST_MAX}
9223        The value for expression nesting was borrowed from the ISO C standard.

9224    {LINE_MAX}
9225        This is a global limit that affects all utilities, unless otherwise noted. The {MAX_CANON}
9226        value from the System Interfaces volume of IEEE Std 1003.1-2001 may further limit input
9227        lines from terminals. The {LINE_MAX} value was the subject of much debate and is a
9228        compromise between those who wished to have unlimited lines and those who understood
9229        that many historical utilities were written with fixed buffers. Frequently, utility writers
9230        selected the UNIX system constant BUFSIZ to allocate these buffers; therefore, some utilities
9231        were limited to 512 bytes for I/O lines, while others achieved 4 096 bytes or greater.

9232        It should be noted that {LINE_MAX} applies only to input line length; there is no
9233        requirement in IEEE Std 1003.1-2001 that limits the length of output lines. Utilities such as
9234        *awk*, *sed*, and *paste* could theoretically construct lines longer than any of the input lines they
9235        received, depending on the options used or the instructions from the application. They are
9236        not required to truncate their output to {LINE_MAX}. It is the responsibility of the
9237        application to deal with this. If the output of one of those utilities is to be piped into another
9238        of the standard utilities, line length restrictions will have to be considered; the *fold* utility,
9239        among others, could be used to ensure that only reasonable line lengths reach utilities or
9240        applications.

9241    {LINK_MAX}
9242        This is defined by the System Interfaces volume of IEEE Std 1003.1-2001.

9243    {MAX_CANON}
9244    {MAX_INPUT}
9245    {NAME_MAX}
9246    {NGROUPS_MAX}
9247    {OPEN_MAX}
9248    {PATH_MAX}
9249    {PIPE_BUF}
9250        These limits are defined by the System Interfaces volume of IEEE Std 1003.1-2001. Note that
9251        the byte lengths described by some of these values continue to represent bytes, even if the
9252        applicable character set uses a multi-byte encoding.

9253    {RE_DUP_MAX}
9254        The value selected is consistent with historical practice. Although the name implies that it
9255        applies to all REs, only BREs use the interval notation \\{*m,n*\\} addressed by this limit.

9256    {POSIX2_SYMLINKS}
9257        The {POSIX2_SYMLINKS} variable indicates that the underlying operating system supports
9258        the creation of symbolic links in specific directories. Many of the utilities defined in
9259        IEEE Std 1003.1-2001 that deal with symbolic links do not depend on this value. For

9260    example, a utility that follows symbolic links (or does not, as the case may be) will only be
9261    affected by a symbolic link if it encounters one. Presumably, a file system that does not
9262    support symbolic links will not contain any. This variable does affect such utilities as *ln* −**s**
9263    and *pax* that attempt to create symbolic links.

9264    {POSIX2_SYMLINKS} was developed even though there is no comparable configuration
9265    value for the system interfaces.

9266   There are different limits associated with command lines and input to utilities, depending on the
9267   method of invocation. In the case of a C program *exec*-ing a utility, {ARG_MAX} is the
9268   underlying limit. In the case of the shell reading a script and *exec*-ing a utility, {LINE_MAX}
9269   limits the length of lines the shell is required to process, and {ARG_MAX} will still be a limit. If a
9270   user is entering a command on a terminal to the shell, requesting that it invoke the utility,
9271   {MAX_INPUT} may restrict the length of the line that can be given to the shell to a value below
9272   {LINE_MAX}.

9273   When an option is supported, *getconf* returns a value of 1. For example, when C development is
9274   supported:

```
9275    if [ "$(getconf POSIX2_C_DEV)" −eq 1 ]; then
9276        echo C supported
9277    fi
```

9278   The *sysconf*( ) function in the C-Language Binding option would return 1.

9279   The following comments describe selection criteria for the symbols and their values:

9280   POSIX2_C_BIND
9281   POSIX2_C_DEV
9282   POSIX2_FORT_DEV
9283   POSIX2_FORT_RUN
9284   POSIX2_SW_DEV
9285   POSIX2_UPE
9286        It is possible for some (usually privileged) operations to remove utilities that support these
9287        options or otherwise to render these options unsupported. The header files, the *sysconf*( )
9288        function, or the *getconf* utility will not necessarily detect such actions, in which case they
9289        should not be considered as rendering the implementation non-conforming. A test suite
9290        should not attempt tests such as:

```
9291    rm /usr/bin/c99
9292    getconf POSIX2_C_DEV
```

9293   POSIX2_LOCALEDEF
9294        This symbol was introduced to allow implementations to restrict supported locales to only
9295        those supplied by the implementation.

9296   IEEE Std 1003.1-2001/Cor 1-2002, item XCU/TC1/D6/2 is applied, deleting the entry for    |
9297   {POSIX2_VERSION} since it is not a utility limit minimum value.                             |

9298   IEEE Std 1003.1-2001/Cor 1-2002, item XCU/TC1/D6/3 is applied, changing the text in Utility  |
9299   Limits from: ''utility (see *getconf*) through the *sysconf*( ) function defined in the System Interfaces  |
9300   volume of IEEE Std 1003.1-2001. The literal names shown in Table 1-3 apply only to the *getconf*  |
9301   utility; the high-level language binding describes the exact form of each name to be used by the  |
9302   interfaces in that binding.'' to: ''utility (see *getconf*).''.                              |

9303 **C.1.10   Grammar Conventions**

9304   There is no additional rationale provided for this section.

9305 **C.1.11   Utility Description Defaults**

9306   This section is arranged with headings in the same order as all the utility descriptions. It is a
9307   collection of related and unrelated information concerning:

9308   1.   The default actions of utilities

9309   2.   The meanings of notations used in IEEE Std 1003.1-2001 that are specific to individual
9310        utility sections

9311   Although this material may seem out of place here, it is important that this information appear
9312   before any of the utilities to be described later.

9313   **NAME**

9314   There is no additional rationale provided for this section.

9315   **SYNOPSIS**

9316   There is no additional rationale provided for this section.

9317   **DESCRIPTION**

9318   There is no additional rationale provided for this section.

9319   **OPTIONS**

9320   Although it has not always been possible, the standard developers tried to avoid repeating
9321   information to reduce the risk that duplicate explanations could each be modified differently.

9322   The need to recognize −− is required because conforming applications need to shield their
9323   operands from any arbitrary options that the implementation may provide as an extension. For
9324   example, if the standard utility *foo* is listed as taking no options, and the application needed to
9325   give it a pathname with a leading hyphen, it could safely do it as:

9326   ```
   foo −− −myfile
   ```

9327   and avoid any problems with −**m** used as an extension.

9328   **OPERANDS**

9329   The usage of − is never shown in the SYNOPSIS. Similarly, the usage of −− is never shown.

9330   The requirement for processing operands in command-line order is to avoid a ''WeirdNIX''
9331   utility that might choose to sort the input files alphabetically, by size, or by directory order.
9332   Although this might be acceptable for some utilities, in general the programmer has a right to
9333   know exactly what order will be chosen.

9334   Some of the standard utilities take multiple *file* operands and act as if they were processing the
9335   concatenation of those files. For example:

9336   ```
    asa file1 file2
   ```

9337   and:

9338   ```
    cat file1 file2 | asa
   ```

9339     have similar results when questions of file access, errors, and performance are ignored. Other
9340     utilities such as *grep* or *wc* have completely different results in these two cases. This latter type of
9341     utility is always identified in its DESCRIPTION or OPERANDS sections, whereas the former is
9342     not. Although it might be possible to create a general assertion about the former case, the
9343     following points must be addressed:

9344     • Access times for the files might be different in the operand case *versus* the *cat* case.

9345     • The utility may have error messages that are cognizant of the input filename, and this added
9346       value should not be suppressed. (As an example, *awk* sets a variable with the filename at
9347       each file boundary.)

9348     **STDIN**

9349     There is no additional rationale provided for this section.

9350     **INPUT FILES**

9351     A conforming application cannot assume the following three commands are equivalent:

```
9352     tail −n +2 file
9353     (sed −n 1q; cat) < file
9354     cat file | (sed −n 1q; cat)
```

9355     The second command is equivalent to the first only when the file is seekable. In the third
9356     command, if the file offset in the open file description were not unspecified, *sed* would have to be
9357     implemented so that it read from the pipe 1 byte at a time or it would have to employ some
9358     method to seek backwards on the pipe. Such functionality is not defined currently in POSIX.1
9359     and does not exist on all historical systems. Other utilities, such as *head*, *read*, and *sh*, have similar
9360     properties, so the restriction is described globally in this section.

9361     The definition of ''text file'' is strictly enforced for input to the standard utilities; very few of
9362     them list exceptions to the undefined results called for here. (Of course, ''undefined'' here does
9363     not mean that historical implementations necessarily have to change to start indicating error
9364     conditions. Conforming applications cannot rely on implementations succeeding or failing when
9365     non-text files are used.)

9366     The utilities that allow line continuation are generally those that accept input languages, rather
9367     than pure data. It would be unusual for an input line of this type to exceed {LINE_MAX} bytes
9368     and unreasonable to require that the implementation allow unlimited accumulation of multiple
9369     lines, each of which could reach {LINE_MAX}. Thus, for a conforming application the total of all
9370     the continued lines in a set cannot exceed {LINE_MAX}.

9371     The format description is intended to be sufficiently rigorous to allow other applications to
9372     generate these input files. However, since <blank>s can legitimately be included in some of the
9373     fields described by the standard utilities, particularly in locales other than the POSIX locale, this
9374     intent is not always realized.

9375     **ENVIRONMENT VARIABLES**

9376     There is no additional rationale provided for this section.

**ASYNCHRONOUS EVENTS**

Because there is no language prohibiting it, a utility is permitted to catch a signal, perform some additional processing (such as deleting temporary files), restore the default signal action (or action inherited from the parent process), and resignal itself.

**STDOUT**

The format description is intended to be sufficiently rigorous to allow post-processing of output by other programs, particularly by an *awk* or *lex* parser.

**STDERR**

This section does not describe error messages that refer to incorrect operation of the utility. Consider a utility that processes program source code as its input. This section is used to describe messages produced by a correctly operating utility that encounters an error in the program source code on which it is processing. However, a message indicating that the utility had insufficient memory in which to operate would not be described.

Some utilities have traditionally produced warning messages without returning a non-zero exit status; these are specifically noted in their sections. Other utilities shall not write to standard error if they complete successfully, unless the implementation provides some sort of extension to increase the verbosity or debugging level.

The format descriptions are intended to be sufficiently rigorous to allow post-processing of output by other programs.

**OUTPUT FILES**

The format description is intended to be sufficiently rigorous to allow post-processing of output by other programs, particularly by an *awk* or *lex* parser.

Receipt of the SIGQUIT signal should generally cause termination (unless in some debugging mode) that would bypass any attempted recovery actions.

**EXTENDED DESCRIPTION**

There is no additional rationale provided for this section.

**EXIT STATUS**

Note the additional discussion of exit values in *Exit Status for Commands* in the *sh* utility. It describes requirements for returning exit values greater than 125.

A utility may list zero as a successful return, 1 as a failure for a specific reason, and greater than 1 as ''an error occurred''. In this case, unspecified conditions may cause a 2 or 3, or other value, to be returned. A strictly conforming application should be written so that it tests for successful exit status values (zero in this case), rather than relying upon the single specific error value listed in IEEE Std 1003.1-2001. In that way, it will have maximum portability, even on implementations with extensions.

The standard developers are aware that the general non-enumeration of errors makes it difficult to write test suites that test the *incorrect* operation of utilities. There are some historical implementations that have expended effort to provide detailed status messages and a helpful environment to bypass or explain errors, such as prompting, retrying, or ignoring unimportant syntax errors; other implementations have not. Since there is no realistic way to mandate system behavior in cases of undefined application actions or system problems—in a manner acceptable to all cultures and environments—attention has been limited to the correct operation of utilities

9419   by the conforming application. Furthermore, the conforming application does not need detailed
9420   information concerning errors that it caused through incorrect usage or that it cannot correct.

9421   There is no description of defaults for this section because all of the standard utilities specify
9422   something (or explicitly state ''Unspecified'') for exit status.

9423   **CONSEQUENCES OF ERRORS**

9424   Several actions are possible when a utility encounters an error condition, depending on the
9425   severity of the error and the state of the utility. Included in the possible actions of various
9426   utilities are: deletion of temporary or intermediate work files; deletion of incomplete files; and
9427   validity checking of the file system or directory.

9428   The text about recursive traversing is meant to ensure that utilities such as *find* process as many
9429   files in the hierarchy as they can. They should not abandon all of the hierarchy at the first error
9430   and resume with the next command-line operand, but should attempt to keep going.

9431   **APPLICATION USAGE**

9432   This section provides additional caveats, issues, and recommendations to the developer.

9433   **EXAMPLES**

9434   This section provides sample usage.

9435   **RATIONALE**

9436   There is no additional rationale provided for this section.

9437   **FUTURE DIRECTIONS**

9438   FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
9439   the future, and often cautions the developer to architect the code to account for a change in this
9440   area. Note that a future directions statement should not be taken as a commitment to adopt a
9441   feature or interface in the future.

9442   **SEE ALSO**

9443   There is no additional rationale provided for this section.

9444   **CHANGE HISTORY**

9445   There is no additional rationale provided for this section.

9446 **C.1.12   Considerations for Utilities in Support of Files of Arbitrary Size**

9447   This section is intended to clarify the requirements for utilities in support of large files.

9448   The utilities listed in this section are utilities which are used to perform administrative tasks
9449   such as to create, move, copy, remove, change the permissions, or measure the resources of a
9450   file. They are useful both as end-user tools and as utilities invoked by applications during
9451   software installation and operation.

9452   The *chgrp*, *chmod*, *chown*, *ln*, and *rm* utilities probably require use of large file-capable versions of
9453   *stat*( ), *lstat*( ), *ftw*( ), and the **stat** structure.

9454   The *cat*, *cksum*, *cmp*, *cp*, *dd*, *mv*, *sum*, and *touch* utilities probably require use of large file-capable
9455   versions of *creat*( ), *open*( ), and *fopen*( ).

9456  The *cat*, *cksum*, *cmp*, *dd*, *df*, *du*, *ls*, and *sum* utilities may require writing large integer values. For
9457  example:

9458      • The *cat* utility might have a –**n** option which counts <newline>s.

9459      • The *cksum* and *ls* utilities report file sizes.

9460      • The *cmp* utility reports the line number at which the first difference occurs, and also has a –**l**
9461        option which reports file offsets.

9462      • The *dd*, *df*, *du*, *ls*, and *sum* utilities report block counts.

9463  The *dd*, *find*, and *test* utilities may need to interpret command arguments that contain 64-bit
9464  values. For *dd*, the arguments include *skip=n*, *seek=n*, and *count=n*. For *find*, the arguments
9465  include –**size***n*. For *test*, the arguments are those associated with algebraic comparisons.

9466  The *df* utility might need to access large file systems with *statvfs*( ).

9467  The *ulimit* utility will need to use large file-capable versions of *getrlimit*( ) and *setrlimit*( ) and be
9468  able to read and write large integer values.

## 9469  C.1.13  Built-In Utilities

9470  All of these utilities can be *exec*-ed. There is no requirement that these utilities are actually built
9471  into the shell itself, but many shells need the capability to do so because the Shell and Utilities
9472  volume of IEEE Std 1003.1-2001, Section 2.9.1.1, Command Search and Execution requires that
9473  they be found prior to the *PATH* search. The shell could satisfy its requirements by keeping a list
9474  of the names and directly accessing the file-system versions regardless of *PATH*. Providing all of
9475  the required functionality for those such as *cd* or *read* would be more difficult.

9476  There were originally three justifications for allowing the omission of *exec*-able versions:

9477  1.  It would require wasting space in the file system, at the expense of very small systems.
9478      However, it has been pointed out that all 16 utilities in the table can be provided with 16
9479      links to a single-line shell script:

9480      ```
      $0 "$@"
      ```

9481  2.  It is not logical to require invocation of utilities such as *cd* because they have no value
9482      outside the shell environment or cannot be useful in a child process. However, counter-
9483      examples always seemed to be available for even the most unusual cases:

9484      ```
      find . –type d –exec cd {} \; –exec foo {} \;
      ```
9485              (which invokes ''foo'' on accessible directories)

9486      ```
      ps ... | sed ... | xargs kill
      ```

9487      ```
      find . –exec true \; –a ...
      ```
9488              (where ''true'' is used for temporary debugging)

9489  3.  It is confusing to have a utility such as *kill* that can easily be in the file system in the base
9490      standard, but that requires built-in status for the User Portability Utilities option (for the %
9491      job control job ID notation). It was decided that it was more appropriate to describe the
9492      required functionality (rather than the implementation) to the system implementors and
9493      let them decide how to satisfy it.

9494  On the other hand, it was realized that any distinction like this between utilities was not useful
9495  to applications, and that the cost to correct it was small. These arguments were ultimately the
9496  most effective.

9497       There were varying reasons for including utilities in the table of built-ins:

9498       *alias*, *fc*, *unalias*
9499            The functionality of these utilities is performed more simply within the shell itself and that
9500            is the model most historical implementations have used.

9501       *bg*, *fg*, *jobs*
9502            All of the job control-related utilities are eligible for built-in status because that is the model
9503            most historical implementations have used.

9504       *cd*, *getopts*, *newgrp*, *read*, *umask*, *wait*
9505            The functionality of these utilities is performed more simply within the context of the
9506            current process. An example can be taken from the usage of the *cd* utility. The purpose of
9507            the *cd* utility is to change the working directory for subsequent operations.  The actions of *cd*
9508            affect the process in which *cd* is executed and all subsequent child processes of that process.
9509            Based on the POSIX standard process model, changes in the process environment of a child
9510            process have no effect on the parent process. If the *cd* utility were executed from a child
9511            process, the working directory change would be effective only in the child process. Child
9512            processes initiated subsequent to the child process that executed the *cd* utility would not
9513            have a changed working directory relative to the parent process.

9514       *command*
9515            This utility was placed in the table primarily to protect scripts that are concerned about
9516            their *PATH* being manipulated. The ''secure'' shell script example in the *command* utility in
9517            the Shell and Utilities volume of IEEE Std 1003.1-2001 would not be possible if a *PATH*
9518            change retrieved an alien version of *command*. (An alternative would have been to
9519            implement *getconf* as a built-in, but the standard developers considered that it carried too
9520            many changing configuration strings to require in the shell.)

9521       *kill*  Since *kill* provides optional job control functionality using shell notation (`%1`, `%2`, and so on),
9522            some implementations would find it extremely difficult to provide this outside the shell.

9523       *true*, *false*
9524            These are in the table as a courtesy to programmers who wish to use the `"while true"`
9525            shell construct without protecting *true* from *PATH* searches. (It is acknowledged that
9526            `"while :"` also works, but the idiom with *true* is historically pervasive.)

9527       All utilities, including those in the table, are accessible via the *system*( ) and *popen*( ) functions in
9528       the System Interfaces volume of IEEE Std 1003.1-2001. There are situations where the return
9529       functionality of *system*( ) and *popen*( ) is not desirable. Applications that require the exit status of
9530       the invoked utility will not be able to use *system*( ) or *popen*( ), since the exit status returned is
9531       that of the command language interpreter rather than that of the invoked utility. The alternative
9532       for such applications is the use of the *exec* family.

9533 **C.2      Shell Command Language**

9534 **C.2.1      Shell Introduction**

9535    The System V shell was selected as the starting point for the Shell and Utilities volume of
9536    IEEE Std 1003.1-2001. The BSD C shell was excluded from consideration for the following
9537    reasons:

9538        • Most historically portable shell scripts assume the Version 7 Bourne shell, from which the
9539          System V shell is derived.

9540        • The majority of tutorial materials on shell programming assume the System V shell.

9541    The construct "#!" is reserved for implementations wishing to provide that extension. If it were
9542    not reserved, the Shell and Utilities volume of IEEE Std 1003.1-2001 would disallow it by forcing
9543    it to be a comment. As it stands, a strictly conforming application must not use "#!" as the first
9544    two characters of the file.

9545 **C.2.2      Quoting**

9546    There is no additional rationale provided for this section.

9547 *C.2.2.1    Escape Character (Backslash)*

9548    There is no additional rationale provided for this section.

9549 *C.2.2.2    Single-Quotes*

9550    A backslash cannot be used to escape a single-quote in a single-quoted string. An embedded
9551    quote can be created by writing, for example: "'a'\''b'", which yields "a'b". (See the Shell
9552    and Utilities volume of IEEE Std 1003.1-2001, Section 2.6.5, Field Splitting for a better
9553    understanding of how portions of words are either split into fields or remain concatenated.) A
9554    single token can be made up of concatenated partial strings containing all three kinds of quoting
9555    or escaping, thus permitting any combination of characters.

9556 *C.2.2.3    Double-Quotes*

9557    The escaped <newline> used for line continuation is removed entirely from the input and is not
9558    replaced by any white space. Therefore, it cannot serve as a token separator.

9559    In double-quoting, if a backslash is immediately followed by a character that would be
9560    interpreted as having a special meaning, the backslash is deleted and the subsequent character is
9561    taken literally. If a backslash does not precede a character that would have a special meaning, it
9562    is left in place unmodified and the character immediately following it is also left unmodified.
9563    Thus, for example:

9564        "\$"   →   $

9565        "\a"   →   \a

9566    It would be desirable to include the statement ''The characters from an enclosed "${" to the
9567    matching '}' shall not be affected by the double quotes'', similar to the one for "$()".
9568    However, historical practice in the System V shell prevents this.

9569    The requirement that double-quotes be matched inside "${...}" within double-quotes and the
9570    rule for finding the matching '}' in the Shell and Utilities volume of IEEE Std 1003.1-2001,
9571    Section 2.6.2, Parameter Expansion eliminate several subtle inconsistencies in expansion for
9572    historical shells in rare cases; for example:

9573
```
"${foo-bar"}
```

9574 yields **bar** when **foo** is not defined, and is an invalid substitution when **foo** is defined, in many
9575 historical shells. The differences in processing the "${...}" form have led to inconsistencies
9576 between historical systems. A consequence of this rule is that single-quotes cannot be used to
9577 quote the '}' within "${...}"; for example:

9578
9579
```
unset bar
foo="${bar-'}'}"
```

9580 is invalid because the "${...}" substitution contains an unpaired unescaped single-quote. The
9581 backslash can be used to escape the '}' in this example to achieve the desired result:

9582
9583
```
unset bar
foo="${bar-\}}"
```

9584 The differences in processing the "${...}" form have led to inconsistencies between the
9585 historical System V shell, BSD, and KornShells, and the text in the Shell and Utilities volume of
9586 IEEE Std 1003.1-2001 is an attempt to converge them without breaking too many applications.
9587 The only alternative to this compromise between shells would be to make the behavior
9588 unspecified whenever the literal characters ''', '{', '}', and '"' appear within "${...}".
9589 To write a portable script that uses these values, a user would have to assign variables; for
9590 example:

9591
9592
```
squote=\' dquote=\" lbrace='{' rbrace='}'
${foo-$squote$rbrace$squote}
```

9593 rather than:

9594
```
${foo-"'}'"}
```

9595 Some implementations have allowed the end of the word to terminate the backquoted command
9596 substitution, such as in:

9597
```
"`echo hello"
```

9598 This usage is undefined; the matching backquote is required by the Shell and Utilities volume of
9599 IEEE Std 1003.1-2001. The other undefined usage can be illustrated by the example:

9600
```
sh -c '' echo "foo''
```

9601 The description of the recursive actions involving command substitution can be illustrated with
9602 an example. Upon recognizing the introduction of command substitution, the shell parses input
9603 (in a new context), gathering the source for the command substitution until an unbalanced ')'
9604 or '`' is located. For example, in the following:

9605
9606
```
echo "$(date; echo "
    one" )"
```

9607 the double-quote following the *echo* does not terminate the first double-quote; it is part of the
9608 command substitution script. Similarly, in:

9609
```
echo "$(echo *)"
```

9610 the asterisk is not quoted since it is inside command substitution; however:

9611
```
echo "$(echo "*")"
```

9612 is quoted (and represents the asterisk character itself).

9613 **C.2.3    Token Recognition**

9614    The "`((`" and "`))`" symbols are control operators in the KornShell, used for an alternative
9615    syntax of an arithmetic expression command. A conforming application cannot use "`((`" as a
9616    single token (with the exception of the "`$((`" form for shell arithmetic).

9617    On some implementations, the symbol "`((`" is a control operator; its use produces unspecified
9618    results.  Applications that wish to have nested subshells, such as:

9619        `((echo Hello);(echo World))`

9620    must separate the "`((`" characters into two tokens by including white space between them.
9621    Some systems may treat these as invalid arithmetic expressions instead of subshells.

9622    Certain combinations of characters are invalid in portable scripts, as shown in the grammar.
9623    Implementations may use these combinations (such as "`|&`") as valid control operators. Portable
9624    scripts cannot rely on receiving errors in all cases where this volume of IEEE Std 1003.1-2001
9625    indicates that a syntax is invalid.

9626    The (3) rule about combining characters to form operators is not meant to preclude systems from
9627    extending the shell language when characters are combined in otherwise invalid ways.
9628    Conforming applications cannot use invalid combinations, and test suites should not penalize
9629    systems that take advantage of this fact. For example, the unquoted combination "`|&`" is not
9630    valid in a POSIX script, but has a specific KornShell meaning.

9631    The (10) rule about '`#`' as the current character is the first in the sequence in which a new token
9632    is being assembled. The '`#`' starts a comment only when it is at the beginning of a token. This
9633    rule is also written to indicate that the search for the end-of-comment does not consider escaped
9634    <newline> specially, so that a comment cannot be continued to the next line.

9635 *C.2.3.1   Alias Substitution*

9636    The alias capability was added in the User Portability Utilities option because it is widely used in
9637    historical implementations by interactive users.

9638    The definition of ''alias name'' precludes an alias name containing a slash character. Since the
9639    text applies to the command words of simple commands, reserved words (in their proper
9640    places) cannot be confused with aliases.

9641    The placement of alias substitution in token recognition makes it clear that it precedes all of the
9642    word expansion steps.

9643    An example concerning trailing <blank>s and reserved words follows. If the user types:

9644        **$** `alias foo="/bin/ls "`
9645        **$** `alias while="/"`

9646    The effect of executing:

9647        **$** `while true`
9648        **>** `do`
9649        **>** `echo "Hello, World"`
9650        **>** `done`

9651    is a never-ending sequence of "`Hello, World`" strings to the screen. However, if the user
9652    types:

9653        **$** `foo while`

9654    the result is an *ls* listing of /.  Since the alias substitution for **foo** ends in a <space>, the next word
9655    is checked for alias substitution. The next word, **while**, has also been aliased, so it is substituted

9656   as well. Since it is not in the proper position as a command word, it is not recognized as a
9657   reserved word.

9658   If the user types:

9659       **$** foo; while

9660   **while** retains its normal reserved-word properties.


## C.2.4   Reserved Words

9661

9662   All reserved words are recognized syntactically as such in the contexts described. However, note
9663   that **in** is the only meaningful reserved word after a **case** or **for**; similarly, **in** is not meaningful as
9664   the first word of a simple command.

9665   Reserved words are recognized only when they are delimited (that is, meet the definition of the
9666   Base Definitions volume of IEEE Std 1003.1-2001, Section 3.435, Word), whereas operators are
9667   themselves delimiters. For instance, ′(′ and ′)′ are control operators, so that no <space> is
9668   needed in (*list*). However, ′{′ and ′}′ are reserved words in { *list*;}, so that in this case the
9669   leading <space> and semicolon are required.

9670   The list of unspecified reserved words is from the KornShell, so conforming applications cannot
9671   use them in places a reserved word would be recognized. This list contained **time** in early
9672   proposals, but it was removed when the *time* utility was selected for the Shell and Utilities
9673   volume of IEEE Std 1003.1-2001.

9674   There was a strong argument for promoting braces to operators (instead of reserved words), so
9675   they would be syntactically equivalent to subshell operators. Concerns about compatibility
9676   outweighed the advantages of this approach. Nevertheless, conforming applications should
9677   consider quoting ′{′ and ′}′ when they represent themselves.

9678   The restriction on ending a name with a colon is to allow future implementations that support
9679   named labels for flow control; see the RATIONALE for the *break* built-in utility.

9680   It is possible that a future version of the Shell and Utilities volume of IEEE Std 1003.1-2001 may
9681   require that ′{′ and ′}′ be treated individually as control operators, although the token "{ }"
9682   will probably be a special-case exemption from this because of the often-used *find*{} construct.


## C.2.5   Parameters and Variables

9683

### C.2.5.1   Positional Parameters

9684

9685   There is no additional rationale provided for this section.


### C.2.5.2   Special Parameters

9686

9687   Most historical implementations implement subshells by forking; thus, the special parameter
9688   ′$′ does not necessarily represent the process ID of the shell process executing the commands
9689   since the subshell execution environment preserves the value of ′$′.

9690   If a subshell were to execute a background command, the value of "$!" for the parent would
9691   not change. For example:

9692       (
9693       date &
9694       echo $!
9695       )
9696       echo $!

9697          would echo two different values for `"$!"`.

9698          The `"$-"` special parameter can be used to save and restore *set* options:

```
9699          Save=$(echo $- | sed 's/[ics]//g')
9700          ...
9701          set +aCefnuvx
9702          if [ -n "$Save" ]; then
9703              set -$Save
9704          fi
```

9705          The three options are removed using *sed* in the example because they may appear in the value of
9706          `"$-"` (from the *sh* command line), but are not valid options to *set*.

9707          The descriptions of parameters `'*'` and `'@'` assume the reader is familiar with the field splitting
9708          discussion in the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.6.5, Field Splitting
9709          and understands that portions of the word remain concatenated unless there is some reason to
9710          split them into separate fields.

9711          Some examples of the `'*'` and `'@'` properties, including the concatenation aspects:

```
9712          set "abc" "def ghi" "jkl"

9713          echo $*    => "abc" "def" "ghi" "jkl"
9714          echo "$*"  => "abc def ghi jkl"
9715          echo $@    => "abc" "def" "ghi" "jkl"
```

9716          but:

```
9717          echo "$@"       => "abc" "def ghi" "jkl"
9718          echo "xx$@yy"   => "xxabc" "def ghi" "jklyy"
9719          echo "$@$@"     => "abc" "def ghi" "jklabc" "def ghi" "jkl"
```

9720          In the preceding examples, the double-quote characters that appear after the `"=>"` do not appear
9721          in the output and are used only to illustrate word boundaries.

9722          The following example illustrates the effect of setting *IFS* to a null string:

```
9723          $ IFS=''
9724          $ set foo bar bam
9725          $ echo "$@"
9726          foo bar bam
9727          $ echo "$*"
9728          foobarbam
9729          $ unset IFS
9730          $ echo "$*"
9731          foo bar bam
```

9732  *C.2.5.3    Shell Variables*

9733          See the discussion of *IFS* in Section C.2.6.5 (on page 241) and the RATIONALE for the *sh* utility.

9734          The prohibition on *LC_CTYPE* changes affecting lexical processing protects the shell
9735          implementor (and the shell programmer) from the ill effects of changing the definition of
9736          <blank> or the set of alphabetic characters in the current environment. It would probably not be
9737          feasible to write a compiled version of a shell script without this rule. The rule applies only to
9738          the current invocation of the shell and its subshells—invoking a shell script or performing *exec sh*
9739          would subject the new shell to the changes in *LC_CTYPE*.

9740    Other common environment variables used by historical shells are not specified by the Shell and
9741    Utilities volume of IEEE Std 1003.1-2001, but they should be reserved for the historical uses.

9742    Tilde expansion for components of *PATH* in an assignment such as:

9743        `PATH=˜hlj/bin:˜dwc/bin:$PATH`

9744    is a feature of some historical shells and is allowed by the wording of the Shell and Utilities
9745    volume of IEEE Std 1003.1-2001, Section 2.6.1, Tilde Expansion. Note that the tildes are expanded
9746    during the assignment to *PATH*, not when *PATH* is accessed during command search.

9747    The following entries represent additional information about variables included in the Shell and
9748    Utilities volume of IEEE Std 1003.1-2001, or rationale for common variables in use by shells that
9749    have been excluded:

9750    _              (Underscore.) While underscore is historical practice, its overloaded usage in
9751                      the KornShell is confusing, and it has been omitted from the Shell and Utilities
9752                      volume of IEEE Std 1003.1-2001.

9753    *ENV*       This variable can be used to set aliases and other items local to the invocation
9754                      of a shell. The file referred to by *ENV* differs from **$HOME/.profile** in that
9755                      **.profile** is typically executed at session start-up, whereas the *ENV* file is
9756                      executed at the beginning of each shell invocation. The *ENV* value is
9757                      interpreted in a manner similar to a dot script, in that the commands are
9758                      executed in the current environment and the file needs to be readable, but not
9759                      executable. However, unlike dot scripts, no *PATH* searching is performed.
9760                      This is used as a guard against Trojan Horse security breaches.

9761    *ERRNO*   This variable was omitted from the Shell and Utilities volume of
9762                      IEEE Std 1003.1-2001 because the values of error numbers are not defined in
9763                      IEEE Std 1003.1-2001 in a portable manner.

9764    *FCEDIT*   Since this variable affects only the *fc* utility, it has been omitted from this more
9765                      global place. The value of *FCEDIT* does not affect the command-line editing
9766                      mode in the shell; see the description of *set* −**o** *vi* in the *set* built-in utility.

9767    *PS1*       This variable is used for interactive prompts. Historically, the ''superuser''
9768                      has had a prompt of ′#′. Since privileges are not required to be monolithic, it
9769                      is difficult to define which privileges should cause the alternate prompt.
9770                      However, a sufficiently powerful user should be reminded of that power by
9771                      having an alternate prompt.

9772    *PS3*       This variable is used by the KornShell for the *select* command. Since the POSIX
9773                      shell does not include *select*, *PS3* was omitted.

9774    *PS4*       This variable is used for shell debugging. For example, the following script:

9775                      `PS4=′[${LINENO}]+ ′`
9776                      `set −x`
9777                      `echo Hello`

9778                      writes the following to standard error:

9779                      `[3]+ echo Hello`

9780    *RANDOM*  This pseudo-random number generator was not seen as being useful to
9781                      interactive users.

9782    *SECONDS*  Although this variable is sometimes used with *PS1* to allow the display of the
9783                      current time in the prompt of the user, it is not one that would be manipulated

9784                                  frequently enough by an interactive user to include in the Shell and Utilities
9785                                  volume of IEEE Std 1003.1-2001.

9786 **C.2.6       Word Expansions**

9787          Step (2) refers to the ''portions of fields generated by step (1)''. For example, if the word being
9788          expanded were `"$x+$y"` and *IFS*=+, the word would be split only if `"$x"` or `"$y"` contained
9789          ´+´; the ´+´ in the original word was not generated by step (1).

9790          *IFS* is used for performing field splitting on the results of parameter and command substitution;
9791          it is not used for splitting all fields. Previous versions of the shell used it for splitting all fields
9792          during field splitting, but this has severe problems because the shell can no longer parse its own
9793          script. There are also important security implications caused by this behavior. All useful
9794          applications of *IFS* use it for parsing input of the *read* utility and for splitting the results of
9795          parameter and command substitution.

9796          The rule concerning expansion to a single field requires that if **foo**=**abc** and **bar**=**def**, that:

9797              `"$foo""$bar"`

9798          expands to the single field:

9799              `abcdef`

9800          The rule concerning empty fields can be illustrated by:

9801          `$      unset foo`
9802          `$      set $foo bar '' xyz "$foo" abc`
9803          `$      for i`
9804          `>      do`
9805          `>          echo "−$i−"`
9806          `>      done`
9807          −**bar**−
9808          −−
9809          −**xyz**−
9810          −−
9811          −**abc**−

9812          Step (1) indicates that parameter expansion, command substitution, and arithmetic expansion
9813          are all processed simultaneously as they are scanned. For example, the following is valid
9814          arithmetic:

9815              `x=1`
9816              `echo $(( $(echo 3)+$x ))`

9817          An early proposal stated that tilde expansion preceded the other steps, but this is not the case in
9818          known historical implementations; if it were, and if a referenced home directory contained a ´$´
9819          character, expansions would result within the directory name.

9820 *C.2.6.1   Tilde Expansion*

9821          Tilde expansion generally occurs only at the beginning of words, but an exception based on
9822          historical practice has been included:

9823              `PATH=/posix/bin:~dgk/bin`

9824          This is eligible for tilde expansion because tilde follows a colon and none of the relevant
9825          characters is quoted. Consideration was given to prohibiting this behavior because any of the
9826          following are reasonable substitutes:

```
9827          PATH=$(printf %s ˜karels/bin : ˜bostic/bin)

9828          for Dir in ˜maart/bin ˜srb/bin ...
9829          do
9830              PATH=${PATH:+$PATH:}$Dir
9831          done
```

9832    In the first command, explicit colons are used for each directory. In all cases, the shell performs
9833    tilde expansion on each directory because all are separate words to the shell.

9834    Note that expressions in operands such as:

9835          `make −k mumble LIBDIR=˜chet/lib`

9836    do not qualify as shell variable assignments, and tilde expansion is not performed (unless the
9837    command does so itself, which *make* does not).

9838    Because of the requirement that the word is not quoted, the following are not equivalent; only
9839    the last causes tilde expansion:

9840          `\˜hlj/    ~h\lj/    ~"hlj"/    ~hlj\/    ~hlj/`

9841    In an early proposal, tilde expansion occurred following any unquoted equals sign or colon, but
9842    this was removed because of its complexity and to avoid breaking commands such as:

9843          `rcp hostname:˜marc/.profile .`

9844    A suggestion was made that the special sequence "$˜" should be allowed to force tilde
9845    expansion anywhere. Since this is not historical practice, it has been left for future
9846    implementations to evaluate. (The description in the Shell and Utilities volume of
9847    IEEE Std 1003.1-2001, Section 2.2, Quoting requires that a dollar sign be quoted to represent
9848    itself, so the "$˜" combination is already unspecified.)

9849    The results of giving tilde with an unknown login name are undefined because the KornShell
9850    "˜+" and "˜−" constructs make use of this condition, but in general it is an error to give an
9851    incorrect login name with tilde. The results of having *HOME* unset are unspecified because some
9852    historical shells treat this as an error.

9853    *C.2.6.2    Parameter Expansion*

9854    The rule for finding the closing '}' in "${...}" is the one used in the KornShell and is
9855    upwardly-compatible with the Bourne shell, which does not determine the closing '}' until the
9856    word is expanded. The advantage of this is that incomplete expansions, such as:

9857          `${foo`

9858    can be determined during tokenization, rather than during expansion.

9859    The string length and substring capabilities were included because of the demonstrated need for
9860    them, based on their usage in other shells, such as C shell and KornShell.

9861    Historical versions of the KornShell have not performed tilde expansion on the word part of
9862    parameter expansion; however, it is more consistent to do so.

9863    *C.2.6.3    Command Substitution*

9864    The "`$()`" form of command substitution solves a problem of inconsistent behavior when using
9865    backquotes. For example:

9866
| Command | Output |
|---|---|
| `echo '\$x'` | `\$x` |
| `echo `echo '\$x'`` | `$x` |
| `echo $(echo '\$x')` | `\$x` |

9870    Additionally, the backquoted syntax has historical restrictions on the contents of the embedded
9871    command. While the newer "`$()`" form can process any kind of valid embedded script, the
9872    backquoted form cannot handle some valid scripts that include backquotes. For example, these
9873    otherwise valid embedded scripts do not work in the left column, but do work on the right:

```
9874    echo `                          echo $(
9875    cat <<\eof                      cat <<\eof
9876    a here-doc with `               a here-doc with )
9877    eof                             eof
9878    `                               )

9879    echo `                          echo $(
9880    echo abc # a comment with `     echo abc # a comment with )
9881    `                               )

9882    echo `                          echo $(
9883    echo '`'                        echo ')'
9884    `                               )
```

9885    Because of these inconsistent behaviors, the backquoted variety of command substitution is not
9886    recommended for new applications that nest command substitutions or attempt to embed
9887    complex scripts.

9888    The KornShell feature:

9889        If *command* is of the form <*word*, *word* is expanded to generate a pathname, and the value of
9890        the command substitution is the contents of this file with any trailing <newline>s deleted.

9891    was omitted from the Shell and Utilities volume of IEEE Std 1003.1-2001 because $(*cat word*) is
9892    an appropriate substitute. However, to prevent breaking numerous scripts relying on this
9893    feature, it is unspecified to have a script within "`$()`" that has only redirections.

9894    The requirement to separate "`$(`" and '`(`' when a single subshell is command-substituted is to
9895    avoid any ambiguities with arithmetic expansion.

9896    IEEE Std 1003.1-2001/Cor 1-2002, item XCU/TC1/D6/4 is applied, changing the text from: ''If a        |
9897    command substitution occurs inside double-quotes, it shall not be performed on the results of        |
9898    the substitution.'' to: ''If a command substitution occurs inside double-quotes, field splitting and   |
9899    pathname expansion shall not be performed on the results of the substitution.''. The              |
9900    replacement text taken from the ISO POSIX-2:1993 standard is clearer about the items that are        |
9901    not performed.                                                                                |

9902 *C.2.6.4*   *Arithmetic Expansion*

9903   The "`(())`" form of KornShell arithmetic in early proposals was omitted. The standard
9904   developers concluded that there was a strong desire for some kind of arithmetic evaluator to
9905   replace *expr*, and that relating it to '`$`' makes it work well with the standard shell language, and
9906   it provides access to arithmetic evaluation in places where accessing a utility would be
9907   inconvenient.

9908   The syntax and semantics for arithmetic were changed for the ISO/IEC 9945-2:1993 standard.
9909   The language is essentially a pure arithmetic evaluator of constants and operators (excluding
9910   assignment) and represents a simple subset of the previous arithmetic language (which was
9911   derived from the KornShell "`(())`" construct). The syntax was changed from that of a
9912   command denoted by ((*expression*)) to an expansion denoted by $((*expression*)). The new form is
9913   a dollar expansion ('`$`') that evaluates the expression and substitutes the resulting value.
9914   Objections to the previous style of arithmetic included that it was too complicated, did not fit in
9915   well with the use of variables in the shell, and its syntax conflicted with subshells. The
9916   justification for the new syntax is that the shell is traditionally a macro language, and if a new
9917   feature is to be added, it should be accomplished by extending the capabilities presented by the
9918   current model of the shell, rather than by inventing a new one outside the model; adding a new
9919   dollar expansion was perceived to be the most intuitive and least destructive way to add such a
9920   new capability.

9921   In early proposals, a form $[*expression*] was used. It was functionally equivalent to the "`$(())`"
9922   of the current text, but objections were lodged that the 1988 KornShell had already implemented
9923   "`$(())`" and there was no compelling reason to invent yet another syntax. Furthermore, the
9924   "`$[]`" syntax had a minor incompatibility involving the patterns in **case** statements.

9925   The portion of the ISO C standard arithmetic operations selected corresponds to the operations
9926   historically supported in the KornShell.

9927   It was concluded that the *test* command (**[**) was sufficient for the majority of relational arithmetic
9928   tests, and that tests involving complicated relational expressions within the shell are rare, yet
9929   could still be accommodated by testing the value of "`$(())`" itself. For example:

```
9930      # a complicated relational expression
9931      while [ $(( (($x + $y)/($a * $b)) < ($foo*$bar) )) −ne 0 ]
```

9932   or better yet, the rare script that has many complex relational expressions could define a
9933   function like this:

```
9934      val() {
9935          return $((!$1))
9936      }
```

9937   and complicated tests would be less intimidating:

```
9938      while val $(( (($x + $y)/($a * $b)) < ($foo*$bar) ))
9939      do
9940          # some calculations
9941      done
```

9942   A suggestion that was not adopted was to modify *true* and *false* to take an optional argument,
9943   and *true* would exit true only if the argument was non-zero, and *false* would exit false only if the
9944   argument was non-zero:

```
9945      while true $(($x > 5 && $y <= 25))
```

9946   There is a minor portability concern with the new syntax. The example "`$((2+2))`" could have
9947   been intended to mean a command substitution of a utility named "`2+2`" in a subshell. The

9948    standard developers considered this to be obscure and isolated to some KornShell scripts
9949    (because "$()" command substitution existed previously only in the KornShell). The text on
9950    command substitution requires that the "$(" and '(' be separate tokens if this usage is needed.

9951    An example such as:

9952        echo $((echo hi);(echo there))

9953    should not be misinterpreted by the shell as arithmetic because attempts to balance the
9954    parentheses pairs would indicate that they are subshells. However, as indicated by the Base
9955    Definitions volume of IEEE Std 1003.1-2001, Section 3.112, Control Operator, a conforming
9956    application must separate two adjacent parentheses with white space to indicate nested
9957    subshells.

9958    Although the ISO/IEC 9899:1999 standard now requires support for **long long** and allows
9959    extended integer types with higher ranks, IEEE Std 1003.1-2001 only requires arithmetic
9960    expansions to support **signed long** integer arithmetic. Implementations are encouraged to
9961    support signed integer values at least as large as the size of the largest file allowed on the
9962    implementation.

9963    Implementations are also allowed to perform floating-point evaluations as long as an
9964    application won't see different results for expressions that would not overflow **signed long**
9965    integer expression evaluation. (This includes appropriate truncation of results to integer values.)

9966    Changes made in response to IEEE PASC Interpretation 1003.2 #208 removed the requirement
9967    that the integer constant suffixes l and L had to be recognized. The ISO POSIX-2:1993 standard
9968    did not require the u, ul, uL, U, Ul, UL, lu, lU, Lu, and LU suffixes since only signed integer
9969    arithmetic was required. Since all arithmetic expressions were treated as handling **signed long**
9970    integer types anyway, the l and L suffixes were redundant. No known scripts used them and
9971    some historic shells did not support them. When the ISO/IEC 9899:1999 standard was used as
9972    the basis for the description of arithmetic processing, the ll and LL suffixes and combinations
9973    were also not required. Implementations are still free to accept any or all of these suffixes, but
9974    are not required to do so.

9975    There was also some confusion as to whether the shell was required to recognize character
9976    constants. Syntactically, character constants were required to be recognized, but the
9977    requirements for the handling of backslash ('\') and quote (''') characters (needed to specify
9978    character constants) within an arithmetic expansion were ambiguous. Furthermore, no known
9979    shells supported them. Changes made in response to IEEE PASC Interpretation 1003.2 #208
9980    removed the requirement to support them (if they were indeed required before).
9981    IEEE Std 1003.1-2001 clearly does not require support for character constants.

9982    *C.2.6.5   Field Splitting*

9983    The operation of field splitting using *IFS*, as described in early proposals, was based on the way
9984    the KornShell splits words, but it is incompatible with other common versions of the shell.
9985    However, each has merit, and so a decision was made to allow both. If the *IFS* variable is unset
9986    or is <space><tab><newline>, the operation is equivalent to the way the System V shell splits
9987    words. Using characters outside the <space><tab><newline> set yields the KornShell behavior,
9988    where each of the non-<space><tab><newline>s is significant. This behavior, which affords the
9989    most flexibility, was taken from the way the original *awk* handled field splitting.

9990    Rule (3) can be summarized as a pseudo-ERE:

9991        (*s*\**ns*\* | *s*+)

9992    where *s* is an *IFS* white space character and *n* is a character in the *IFS* that is not white space.
9993    Any string matching that ERE delimits a field, except that the *s*+ form does not delimit fields at

9994    the beginning or the end of a line. For example, if *IFS* is <space>/<comma>/<tab>, the string:

9995        `<space><space>red<space><space>,<space>white<space>blue`

9996    yields the three colors as the delimited fields.

9997  *C.2.6.6  Pathname Expansion*

9998    There is no additional rationale provided for this section.

9999  *C.2.6.7  Quote Removal*

10000    There is no additional rationale provided for this section.

## 10001 C.2.7   Redirection

10002    In the System Interfaces volume of IEEE Std 1003.1-2001, file descriptors are integers in the range
10003    0–({OPEN_MAX}–1). The file descriptors discussed in the Shell and Utilities volume of
10004    IEEE Std 1003.1-2001, Section 2.7, Redirection are that same set of small integers.

10005    Having multi-digit file descriptor numbers for I/O redirection can cause some obscure
10006    compatibility problems. Specifically, scripts that depend on an example command:

10007        `echo 22>/dev/null`

10008    echoing "2" to standard error or "22" to standard output are no longer portable. However, the
10009    file descriptor number must still be delimited from the preceding text. For example:

10010        `cat file2>foo`

10011    writes the contents of **file2**, not the contents of **file**.

10012    The "`>|`" format of output redirection was adopted from the KornShell. Along with the
10013    *noclobber* option, *set* –**C**, it provides a safety feature to prevent inadvertent overwriting of
10014    existing files. (See the RATIONALE for the *pathchk* utility for why this step was taken.) The
10015    restriction on regular files is historical practice.

10016    The System V shell and the KornShell have differed historically on pathname expansion of *word*;
10017    the former never performed it, the latter only when the result was a single field (file). As a
10018    compromise, it was decided that the KornShell functionality was useful, but only as a shorthand
10019    device for interactive users. No reasonable shell script would be written with a command such
10020    as:

10021        `cat foo > a*`

10022    Thus, shell scripts are prohibited from doing it, while interactive users can select the shell with
10023    which they are most comfortable.

10024    The construct "`2>&1`" is often used to redirect standard error to the same file as standard
10025    output. Since the redirections take place beginning to end, the order of redirections is significant.
10026    For example:

10027        `ls > foo 2>&1`

10028    directs both standard output and standard error to file **foo**. However:

10029        `ls 2>&1 > foo`

10030    only directs standard output to file **foo** because standard error was duplicated as standard
10031    output before standard output was directed to file **foo**.

10032    The "`<>`" operator could be useful in writing an application that worked with several terminals,
10033    and occasionally wanted to start up a shell. That shell would in turn be unable to run
10034    applications that run from an ordinary controlling terminal unless it could make use of "`<>`"
10035    redirection. The specific example is a historical version of the pager *more*, which reads from
10036    standard error to get its commands, so standard input and standard output are both available
10037    for their usual usage. There is no way of saying the following in the shell without "`<>`":

10038
```
cat food | more − >/dev/tty03 2<>/dev/tty03
```

10039    Another example of "`<>`" is one that opens **/dev/tty** on file descriptor 3 for reading and writing:

10040
```
exec 3<> /dev/tty
```

10041    An example of creating a lock file for a critical code region:

10042
10043
10044
10045
10046
10047
10048
```
set −C
until    2> /dev/null > lockfile
do       sleep 30
done
set +C
perform critical function
rm lockfile
```

10049    Since **/dev/null** is not a regular file, no error is generated by redirecting to it in *noclobber* mode.

10050    Tilde expansion is not performed on a here-document because the data is treated as if it were
10051    enclosed in double quotes.

10052 *C.2.7.1*    *Redirecting Input*

10053    There is no additional rationale provided for this section.

10054 *C.2.7.2*    *Redirecting Output*

10055    There is no additional rationale provided for this section.

10056 *C.2.7.3*    *Appending Redirected Output*

10057    Note that when a file is opened (even with the O_APPEND flag set), the initial file offset for that
10058    file is set to the beginning of the file. Some historic shells set the file offset to the current end-of-
10059    file when **append** mode shell redirection was used, but this is not allowed by
10060    IEEE Std 1003.1-2001.

10061 *C.2.7.4*    *Here-Document*

10062    There is no additional rationale provided for this section.

10063 *C.2.7.5*    *Duplicating an Input File Descriptor*

10064    There is no additional rationale provided for this section.

10065 *C.2.7.6*    *Duplicating an Output File Descriptor*

10066    There is no additional rationale provided for this section.

10068          There is no additional rationale provided for this section.

10069 **C.2.8      Exit Status and Errors**

10071          There is no additional rationale provided for this section.

10073          There is a historical difference in *sh* and *ksh* non-interactive error behavior. When a command
10074          named in a script is not found, some implementations of *sh* exit immediately, but *ksh* continues
10075          with the next command. Thus, the Shell and Utilities volume of IEEE Std 1003.1-2001 says that
10076          the shell ''may'' exit in this case. This puts a small burden on the programmer, who has to test
10077          for successful completion following a command if it is important that the next command not be
10078          executed if the previous command was not found. If it is important for the command to have
10079          been found, it was probably also important for it to complete successfully. The test for successful
10080          completion would not need to change.

10081          Historically, shells have returned an exit status of 128+$n$, where $n$ represents the signal number.
10082          Since signal numbers are not standardized, there is no portable way to determine which signal
10083          caused the termination. Also, it is possible for a command to exit with a status in the same range
10084          of numbers that the shell would use to report that the command was terminated by a signal.
10085          Implementations are encouraged to choose exit values greater than 256 to indicate programs
10086          that terminate by a signal so that the exit status cannot be confused with an exit status generated
10087          by a normal termination.

10088          Historical shells make the distinction between ''utility not found'' and ''utility found but cannot
10089          execute'' in their error messages. By specifying two seldomly used exit status values for these
10090          cases, 127 and 126 respectively, this gives an application the opportunity to make use of this
10091          distinction without having to parse an error message that would probably change from locale to
10092          locale. The *command*, *env*, *nohup*, and *xargs* utilities in the Shell and Utilities volume of
10093          IEEE Std 1003.1-2001 have also been specified to use this convention.

10094          When a command fails during word expansion or redirection, most historical implementations
10095          exit with a status of 1. However, there was some sentiment that this value should probably be
10096          much higher so that an application could distinguish this case from the more normal exit status
10097          values. Thus, the language ''greater than zero'' was selected to allow either method to be
10098          implemented.

10099 **C.2.9      Shell Commands**

10100          A description of an ''empty command'' was removed from an early proposal because it is only
10101          relevant in the cases of *sh* −**c** " ", *system*(" "), or an empty shell-script file (such as the
10102          implementation of *true* on some historical systems). Since it is no longer mentioned in the Shell
10103          and Utilities volume of IEEE Std 1003.1-2001, it falls into the silently unspecified category of
10104          behavior where implementations can continue to operate as they have historically, but
10105          conforming applications do not construct empty commands. (However, note that *sh* does
10106          explicitly state an exit status for an empty string or file.) In an interactive session or a script with
10107          other commands, extra <newline>s or semicolons, such as:

```
10108        $ false
10109        $
10110        $ echo $?
10111        1
```

10112 would not qualify as the empty command described here because they would be consumed by
10113 other parts of the grammar.

## 10114 C.2.9.1    *Simple Commands*

10115 The enumerated list is used only when the command is actually going to be executed. For
10116 example, in:

```
10117        true || $foo *
```

10118 no expansions are performed.

10119 The following example illustrates both how a variable assignment without a command name
10120 affects the current execution environment, and how an assignment with a command name only
10121 affects the execution environment of the command:

```
10122        $ x=red
10123        $ echo $x
10124        red
10125        $ export x
10126        $ sh −c 'echo $x'
10127        red
10128        $ x=blue sh −c 'echo $x'
10129        blue
10130        $ echo $x
10131        red
```

10132 This next example illustrates that redirections without a command name are still performed:

```
10133        $ ls foo
10134        ls: foo: no such file or directory
10135        $ > foo
10136        $ ls foo
10137        foo
```

10138 A command without a command name, but one that includes a command substitution, has an
10139 exit status of the last command substitution that the shell performed. For example:

```
10140        if       x=$(command)
10141        then     ...
10142        fi
```

10143 An example of redirections without a command name being performed in a subshell shows that
10144 the here-document does not disrupt the standard input of the **while** loop:

```
10145        IFS=:
10146        while    read a b
10147        do       echo $a
10148                 <<−eof
10149                 Hello
10150                 eof
10151        done </etc/passwd
```

10152   Following are examples of commands without command names in AND-OR lists:

```
10153       > foo || {
10154           echo "error: foo cannot be created" >&2
10155           exit 1
10156       }

10157       # set saved if /vmunix.save exists
10158       test −f /vmunix.save && saved=1
```

10159   Command substitution and redirections without command names both occur in subshells, but
10160   they are not necessarily the same ones. For example, in:

```
10161       exec 3> file
10162       var=$(echo foo >&3) 3>&1
```

10163   it is unspecified whether **foo** is echoed to the file or to standard output.

10164   **Command Search and Execution**

10165   This description requires that the shell can execute shell scripts directly, even if the underlying
10166   system does not support the common "#!" interpreter convention. That is, if file **foo** contains
10167   shell commands and is executable, the following executes **foo**:

```
10168       ./foo
```

10169   The command search shown here does not match all historical implementations. A more typical
10170   sequence has been:

10171   • Any built-in (special or regular)

10172   • Functions

10173   • Path search for executable files

10174   But there are problems with this sequence. Since the programmer has no idea in advance which
10175   utilities might have been built into the shell, a function cannot be used to override portably a
10176   utility of the same name. (For example, a function named *cd* cannot be written for many
10177   historical systems.) Furthermore, the *PATH* variable is partially ineffective in this case, and only
10178   a pathname with a slash can be used to ensure a specific executable file is invoked.

10179   After the *execve*( ) failure described, the shell normally executes the file as a shell script. Some
10180   implementations, however, attempt to detect whether the file is actually a script and not an
10181   executable from some other architecture. The method used by the KornShell is allowed by the
10182   text that indicates non-text files may be bypassed.

10183   The sequence selected for the Shell and Utilities volume of IEEE Std 1003.1-2001 acknowledges
10184   that special built-ins cannot be overridden, but gives the programmer full control over which
10185   versions of other utilities are executed. It provides a means of suppressing function lookup (via
10186   the *command* utility) for the user's own functions and ensures that any regular built-ins or
10187   functions provided by the implementation are under the control of the path search. The
10188   mechanisms for associating built-ins or functions with executable files in the path are not
10189   specified by the Shell and Utilities volume of IEEE Std 1003.1-2001, but the wording requires that
10190   if either is implemented, the application is not able to distinguish a function or built-in from an
10191   executable (other than in terms of performance, presumably). The implementation ensures that
10192   all effects specified by the Shell and Utilities volume of IEEE Std 1003.1-2001 resulting from the
10193   invocation of the regular built-in or function (interaction with the environment, variables, traps,
10194   and so on) are identical to those resulting from the invocation of an executable file.

10195     **Examples**

10196     Consider three versions of the *ls* utility:

10197        1.   The application includes a shell function named *ls*.

10198        2.   The user writes a utility named *ls* and puts it in /**fred/bin**.

10199        3.   The example implementation provides *ls* as a regular shell built-in that is invoked (either
10200             by the shell or directly by *exec*) when the path search reaches the directory /**posix/bin**.

10201     If *PATH*=/**posix/bin**, various invocations yield different versions of *ls*:

10202
10203

| Invocation | Version of *ls* |
|---|---|
| *ls* (from within application script) | (1) function |
| *command ls* (from within application script) | (3) built-in |
| *ls* (from within makefile called by application) | (3) built-in |
| *system*("*ls*") | (3) built-in |
| *PATH*="/**fred/bin**:$*PATH*" *ls* | (2) user's version |

10204
10205
10206
10207
10208

10209  *C.2.9.2   Pipelines*

10210     Because pipeline assignment of standard input or standard output or both takes place before
10211     redirection, it can be modified by redirection. For example:

10212        **$** `command1 2>&1 | command2`

10213     sends both the standard output and standard error of *command1* to the standard input of
10214     *command2*.

10215     The reserved word **!** allows more flexible testing using AND and OR lists.

10216     It was suggested that it would be better to return a non-zero value if any command in the
10217     pipeline terminates with non-zero status (perhaps the bitwise-inclusive OR of all return values).
10218     However, the choice of the last-specified command semantics are historical practice and would
10219     cause applications to break if changed. An example of historical behavior:

10220        ```
           $ sleep 5 | (exit 4)
```
10221        ```
           $ echo $?
```
10222        ```
           4
```
10223        ```
           $ (exit 4) | sleep 5
```
10224        ```
           $ echo $?
```
10225        ```
           0
```

10226  *C.2.9.3   Lists*

10227     The equal precedence of "&&" and "||" is historical practice. The standard developers
10228     evaluated the model used more frequently in high-level programming languages, such as C, to
10229     allow the shell logical operators to be used for complex expressions in an unambiguous way, but
10230     they could not allow historical scripts to break in the subtle way unequal precedence might
10231     cause. Some arguments were posed concerning the "{}" or "()" groupings that are required
10232     historically. There are some disadvantages to these groupings:

10233        • The "()" can be expensive, as they spawn other processes on some implementations. This
10234          performance concern is primarily an implementation issue.

10235        • The "{ }" braces are not operators (they are reserved words) and require a trailing space
10236          after each '{', and a semicolon before each '}'. Most programmers (and certainly

10237    interactive users) have avoided braces as grouping constructs because of the problematic
10238    syntax required. Braces were not changed to operators because that would generate
10239    compatibility issues even greater than the precedence question; braces appear outside the
10240    context of a keyword in many shell scripts.

10241    IEEE PASC Interpretation 1003.2 #204 is applied, clarifying that the operators "`&&`" and "`||`"
10242    are evaluated with left associativity.

**Asynchronous Lists**

10244    The grammar treats a construct such as:

10245
```
foo & bar & bam &
```

10246    as one ''asynchronous list'', but since the status of each element is tracked by the shell, the term
10247    ''element of an asynchronous list'' was introduced to identify just one of the **foo**, **bar**, or **bam**
10248    portions of the overall list.

10249    Unless the implementation has an internal limit, such as {CHILD_MAX}, on the retained process
10250    IDs, it would require unbounded memory for the following example:

10251
10252
10253
```
while true
do      foo & echo $!
done
```

10254    The treatment of the signals SIGINT and SIGQUIT with asynchronous lists is described in the
10255    Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.11, Signals and Error Handling.

10256    Since the connection of the input to the equivalent of **/dev/null** is considered to occur before
10257    redirections, the following script would produce no output:

10258
10259
10260
```
exec < /etc/passwd
cat <&0 &
wait
```

**Sequential Lists**

10262    There is no additional rationale provided for this section.

**AND Lists**

10264    There is no additional rationale provided for this section.

**OR Lists**

10266    There is no additional rationale provided for this section.

10267  *C.2.9.4  Compound Commands*

**Grouping Commands**

10268

10269  The semicolon shown in {*compound-list*;} is an example of a control operator delimiting the **}**
10270  reserved word. Other delimiters are possible, as shown in the Shell and Utilities volume of
10271  IEEE Std 1003.1-2001, Section 2.10, Shell Grammar; <newline> is frequently used.

10272  A proposal was made to use the **<do-done>** construct in all cases where command grouping in
10273  the current process environment is performed, identifying it as a construct for the grouping
10274  commands, as well as for shell functions. This was not included because the shell already has a
10275  grouping construct for this purpose ("{ }"), and changing it would have been counter-
10276  productive.

**For Loop**

10277

10278  The format is shown with generous usage of <newline>s. See the grammar in the Shell and
10279  Utilities volume of IEEE Std 1003.1-2001, Section 2.10, Shell Grammar for a precise description of
10280  where <newline>s and semicolons can be interchanged.

10281  Some historical implementations support ′{′ and ′}′ as substitutes for **do** and **done**. The
10282  standard developers chose to omit them, even as an obsolescent feature. (Note that these
10283  substitutes were only for the **for** command; the **while** and **until** commands could not use them
10284  historically because they are followed by compound-lists that may contain "{ . . . }" grouping
10285  commands themselves.)

10286  The reserved word pair **do** … **done** was selected rather than **do** … **od** (which would have
10287  matched the spirit of **if** … **fi** and **case** … **esac**) because *od* is already the name of a standard
10288  utility.

10289  PASC Interpretation 1003.2 #169 has been applied changing the grammar.

**Case Conditional Construct**

10290

10291  An optional left parenthesis before *pattern* was added to allow numerous historical KornShell
10292  scripts to conform. At one time, using the leading parenthesis was required if the **case** statement
10293  was to be embedded within a "$( )" command substitution; this is no longer the case with the
10294  POSIX shell. Nevertheless, many historical scripts use the left parenthesis, if only because it
10295  makes matching-parenthesis searching easier in *vi* and other editors. This is a relatively simple
10296  implementation change that is upwards-compatible for all scripts.

10297  Consideration was given to requiring *break* inside the *compound-list* to prevent falling through to
10298  the next pattern action list. This was rejected as being nonexisting practice. An interesting
10299  undocumented feature of the KornShell is that using ";&" instead of ";;" as a terminator
10300  causes the exact opposite behavior—the flow of control continues with the next *compound-list*.

10301  The pattern ′*′, given as the last pattern in a **case** construct, is equivalent to the default case in
10302  a C-language **switch** statement.

10303  The grammar shows that reserved words can be used as patterns, even if one is the first word on
10304  a line. Obviously, the reserved word **esac** cannot be used in this manner.

10305       **If Conditional Construct**

10306       The precise format for the command syntax is described in the Shell and Utilities volume of
10307       IEEE Std 1003.1-2001, Section 2.10, Shell Grammar.

10308       **While Loop**

10309       The precise format for the command syntax is described in the Shell and Utilities volume of
10310       IEEE Std 1003.1-2001, Section 2.10, Shell Grammar.

10311       **Until Loop**

10312       The precise format for the command syntax is described in the Shell and Utilities volume of
10313       IEEE Std 1003.1-2001, Section 2.10, Shell Grammar.

10314  *C.2.9.5   Function Definition Command*

10315       The description of functions in an early proposal was based on the notion that functions should
10316       behave like miniature shell scripts; that is, except for sharing variables, most elements of an
10317       execution environment should behave as if they were a new execution environment, and
10318       changes to these should be local to the function. For example, traps and options should be reset
10319       on entry to the function, and any changes to them do not affect the traps or options of the caller.
10320       There were numerous objections to this basic idea, and the opponents asserted that functions
10321       were intended to be a convenient mechanism for grouping common commands that were to be
10322       executed in the current execution environment, similar to the execution of the *dot* special
10323       built-in.

10324       It was also pointed out that the functions described in that early proposal did not provide a local
10325       scope for everything a new shell script would, such as the current working directory, or *umask*,
10326       but instead provided a local scope for only a few select properties.  The basic argument was that
10327       if a local scope is needed for the execution environment, the mechanism already existed: the
10328       application can put the commands in a new shell script and call that script. All historical shells
10329       that implemented functions, other than the KornShell, have implemented functions that operate
10330       in the current execution environment. Because of this, traps and options have a global scope
10331       within a shell script. Local variables within a function were considered and included in another
10332       early proposal (controlled by the special built-in *local*), but were removed because they do not fit
10333       the simple model developed for functions and because there was some opposition to adding yet
10334       another new special built-in that was not part of historical practice.  Implementations should
10335       reserve the identifier *local* (as well as *typeset*, as used in the KornShell) in case this local variable
10336       mechanism is adopted in a future version of IEEE Std 1003.1-2001.

10337       A separate issue from the execution environment of a function is the availability of that function
10338       to child shells. A few objectors maintained that just as a variable can be shared with child shells
10339       by exporting it, so should a function. In early proposals, the *export* command therefore had a –**f**
10340       flag for exporting functions. Functions that were exported were to be put into the environment
10341       as *name*()=*value* pairs, and upon invocation, the shell would scan the environment for these and
10342       automatically define these functions. This facility was strongly opposed and was omitted. Some
10343       of the arguments against exportable functions were as follows:

10344       • There was little historical practice. The Ninth Edition shell provided them, but there was
10345         controversy over how well it worked.

10346       • There are numerous security problems associated with functions appearing in the
10347         environment of a user and overriding standard utilities or the utilities owned by the
10348         application.

10349    • There was controversy over requiring *make* to import functions, where it has historically used
10350      an *exec* function for many of its command line executions.

10351    • Functions can be big and the environment is of a limited size. (The counter-argument was
10352      that functions are no different from variables in terms of size: there can be big ones, and there
10353      can be small ones—and just as one does not export huge variables, one does not export huge
10354      functions. However, this might not apply to the average shell-function writer, who typically
10355      writes much larger functions than variables.)

10356    As far as can be determined, the functions in the Shell and Utilities volume of
10357    IEEE Std 1003.1-2001 match those in System V. Earlier versions of the KornShell had two
10358    methods of defining functions:

10359
```
function fname { compound-list }
```

10360    and:

10361
```
fname() { compound-list }
```

10362    The latter used the same definition as the Shell and Utilities volume of IEEE Std 1003.1-2001, but
10363    differed in semantics, as described previously. The current edition of the KornShell aligns the
10364    latter syntax with the Shell and Utilities volume of IEEE Std 1003.1-2001 and keeps the former as
10365    is.

10366    The name space for functions is limited to that of a *name* because of historical practice.
10367    Complications in defining the syntactic rules for the function definition command and in dealing
10368    with known extensions such as the "@()" usage in the KornShell prevented the name space
10369    from being widened to a *word*.  Using functions to support synonyms such as the "!!" and '%'
10370    usage in the C shell is thus disallowed to conforming applications, but acceptable as an
10371    extension. For interactive users, the aliasing facilities in the Shell and Utilities volume of
10372    IEEE Std 1003.1-2001 should be adequate for this purpose. It is recognized that the name space
10373    for utilities in the file system is wider than that currently supported for functions, if the portable
10374    filename character set guidelines are ignored, but it did not seem useful to mandate extensions
10375    in systems for so little benefit to conforming applications.

10376    The "()" in the function definition command consists of two operators.  Therefore, intermixing
10377    <blank>s with the *fname*, '(', and ')' is allowed, but unnecessary.

10378    An example of how a function definition can be used wherever a simple command is allowed:

10379
```
# If variable i is equal to "yes",
```
10380
```
# define function foo to be ls −l
```
10381
```
#
```
10382
```
[ "$i" = yes ] && foo() {
```
10383
```
    ls −l
```
10384
```
}
```

10385 ## C.2.10   Shell Grammar

10386    There are several subtle aspects of this grammar where conventional usage implies rules about
10387    the grammar that in fact are not true.

10388    For *compound_list*, only the forms that end in a *separator* allow a reserved word to be recognized,
10389    so usually only a *separator* can be used where a compound list precedes a reserved word (such as
10390    **Then**, **Else**, **Do**, and **Rbrace**).  Explicitly requiring a separator would disallow such valid (if rare)
10391    statements as:

10392
```
if (false) then (echo x) else (echo y) fi
```

10393          See the Note under special grammar rule (1).

10394          Concerning the third sentence of rule (1) (''Also, if the parser ...''):

10395          • This sentence applies rather narrowly: when a compound list is terminated by some clear
10396            delimiter (such as the closing **fi** of an inner **if_clause**) then it would apply; where the
10397            compound list might continue (as in after a ' ; '), rule (7a) (and consequently the first
10398            sentence of rule (1)) would apply. In many instances the two conditions are identical, but this
10399            part of rule (1) does not give license to treating a **WORD** as a reserved word unless it is in a
10400            place where a reserved word has to appear.

10401          • The statement is equivalent to requiring that when the LR(1) lookahead set contains exactly
10402            one reserved word, it must be recognized if it is present. (Here ''LR(1)'' refers to the
10403            theoretical concepts, not to any real parser generator.)

10404          For example, in the construct below, and when the parser is at the point marked with ' ^ ',
10405          the only next legal token is **then** (this follows directly from the grammar rules):

10406          ```
               if if...fi then ... fi
10407                        ^
               ```

10408          At that point, the **then** must be recognized as a reserved word.

10409          (Depending on the parser generator actually used, ''extra'' reserved words may be in some
10410          lookahead sets. It does not really matter if they are recognized, or even if any possible
10411          reserved word is recognized in that state, because if it is recognized and is not in the
10412          (theoretical) LR(1) lookahead set, an error is ultimately detected. In the example above, if
10413          some other reserved word (for example, **while**) is also recognized, an error occurs later.

10414          This is approximately equivalent to saying that reserved words are recognized after other
10415          reserved words (because it is after a reserved word that this condition occurs), but avoids the
10416          ''except for ...'' list that would be required for **case**, **for**, and so on. (Reserved words are of
10417          course recognized anywhere a *simple_command* can appear, as well. Other rules take care of
10418          the special cases of non-recognition, such as rule (4) for **case** statements.)

10419          Note that the body of here-documents are handled by token recognition (see the Shell and
10420          Utilities volume of IEEE Std 1003.1-2001, Section 2.3, Token Recognition) and do not appear in
10421          the grammar directly. (However, the here-document I/O redirection operator is handled as part
10422          of the grammar.)

10423          The start symbol of the grammar (**complete_command**) represents either input from the
10424          command line or a shell script. It is repeatedly applied by the interpreter to its input and
10425          represents a single ''chunk'' of that input as seen by the interpreter.

10426 *C.2.10.1  Shell Grammar Lexical Conventions*

10427          There is no additional rationale provided for this section.

10428 *C.2.10.2  Shell Grammar Rules*

10429          There is no additional rationale provided for this section.

10430 **C.2.11  Signals and Error Handling**

10431    There is no additional rationale provided for this section.

10432 **C.2.12  Shell Execution Environment**

10433    Some implementations have implemented the last stage of a pipeline in the current environment
10434    so that commands such as:

10435    `   command | read foo`

10436    set variable **foo** in the current environment. This extension is allowed, but not required;
10437    therefore, a shell programmer should consider a pipeline to be in a subshell environment, but
10438    not depend on it.

10439    In early proposals, the description of execution environment failed to mention that each
10440    command in a multiple command pipeline could be in a subshell execution environment. For
10441    compatibility with some historical shells, the wording was phrased to allow an implementation
10442    to place any or all commands of a pipeline in the current environment. However, this means that
10443    a POSIX application must assume each command is in a subshell environment, but not depend
10444    on it.

10445    The wording about shell scripts is meant to convey the fact that describing ''trap actions'' can
10446    only be understood in the context of the shell command language. Outside of this context, such
10447    as in a C-language program, signals are the operative condition, not traps.

10448 **C.2.13  Pattern Matching Notation**

10449    Pattern matching is a simpler concept and has a simpler syntax than REs, as the former is
10450    generally used for the manipulation of filenames, which are relatively simple collections of
10451    characters, while the latter is generally used to manipulate arbitrary text strings of potentially
10452    greater complexity. However, some of the basic concepts are the same, so this section points
10453    liberally to the detailed descriptions in the Base Definitions volume of IEEE Std 1003.1-2001,
10454    Chapter 9, Regular Expressions.

10455 *C.2.13.1  Patterns Matching a Single Character*

10456    Both quoting and escaping are described here because pattern matching must work in three
10457    separate circumstances:

10458    1.  Calling directly upon the shell, such as in pathname expansion or in a **case** statement. All
10459        of the following match the string or file **abc**:

10460        `    abc "abc" a"b"c a\bc a[b]c a["b"]c a[\b]c a["\b"]c a?c a*c`

10461        The following do not:

10462        `    "a?c" a\*c a\[b]c`

10463    2.  Calling a utility or function without going through a shell, as described for *find* and the
10464        *fnmatch*() function defined in the System Interfaces volume of IEEE Std 1003.1-2001.

10465    3.  Calling utilities such as *find*, *cpio*, *tar*, or *pax* through the shell command line. In this case,
10466        shell quote removal is performed before the utility sees the argument. For example, in:

10467        `    find /bin −name "e\c[\h]o" −print`

10468        after quote removal, the backslashes are presented to *find* and it treats them as escape
10469        characters. Both precede ordinary characters, so the *c* and *h* represent themselves and *echo*
10470        would be found on many historical systems (that have it in **/bin**). To find a filename that

10471    contained shell special characters or pattern characters, both quoting and escaping are
10472    required, such as:

10473        pax −r ... "*a\(\?"

10474    to extract a filename ending with "a(?".

10475    Conforming applications are required to quote or escape the shell special characters (sometimes
10476    called metacharacters). If used without this protection, syntax errors can result or
10477    implementation extensions can be triggered. For example, the KornShell supports a series of
10478    extensions based on parentheses in patterns.

10479    The restriction on a circumflex in a bracket expression is to allow implementations that support
10480    pattern matching using the circumflex as the negation character in addition to the exclamation
10481    mark. A conforming application must use something like "[\^!]" to match either character.

10482 *C.2.13.2  Patterns Matching Multiple Characters*

10483    Since each asterisk matches zero or more occurrences, the patterns "a*b" and "a**b" have
10484    identical functionality.

10485    **Examples**

10486    a[bc]       Matches the strings "ab" and "ac".

10487    a*d         Matches the strings "ad", "abd", and "abcd", but not the string "abc".

10488    a*d*        Matches the strings "ad", "abcd", "abcdef", "aaaad", and "adddd".

10489    *a*d        Matches the strings "ad", "abcd", "efabcd", "aaaad", and "adddd".

10490 *C.2.13.3  Patterns Used for Filename Expansion*

10491    The caveat about a slash within a bracket expression is derived from historical practice. The
10492    pattern "a[b/c]d" does not match such pathnames as **abd** or **a/d**. On some implementations
10493    (including those conforming to the Single UNIX Specification), it matched a pathname of
10494    literally "a[b/c]d". On other systems, it produced an undefined condition (an unescaped '['
10495    used outside a bracket expression). In this version, the XSI behavior is now required.

10496    Filenames beginning with a period historically have been specially protected from view on
10497    UNIX systems. A proposal to allow an explicit period in a bracket expression to match a leading
10498    period was considered; it is allowed as an implementation extension, but a conforming
10499    application cannot make use of it. If this extension becomes popular in the future, it will be
10500    considered for a future version of the Shell and Utilities volume of IEEE Std 1003.1-2001.

10501    Historical systems have varied in their permissions requirements. To match **f*/bar** has required
10502    read permissions on the **f*** directories in the System V shell, but the Shell and Utilities volume of
10503    IEEE Std 1003.1-2001, the C shell, and KornShell require only search permissions.

10504 **C.2.14   Special Built-In Utilities**

10505       See the RATIONALE sections on the individual reference pages.

10506 **C.3      Batch Environment Services and Utilities**

10507       **Scope of the Batch Environment Services and Utilities Option**                    |

10508       This section summarizes the deliberations of the IEEE P1003.15 (Batch Environment) working    |
10509       group in the development of the Batch Environment Services and Utilities option, which covers  |
10510       a set of services and utilities defining a batch processing system.                            |

10511       This informative section contains historical information concerning the contents of the
10512       amendment and describes why features were included or discarded by the working group.

10513       **History of Batch Systems**

10514       The supercomputing technical committee began as a ''Birds Of a Feather'' (BOF) at the January
10515       1987 Usenix meeting. There was enough general interest to form a supercomputing attachment
10516       to the /usr/group working groups. Several subgroups rapidly formed. Of those subgroups, the
10517       batch group was the most ambitious. The first early meetings were spent evaluating user needs
10518       and existing batch implementations.

10519       To evaluate user needs, individuals from the supercomputing community came and presented
10520       their needs. Common requests were flexibility, interoperability, control of resources, and ease-
10521       of-use. Backward-compatibility was not an issue. The working group then evaluated some
10522       existing systems. The following different systems were evaluated:

10523         • PROD

10524         • Convex Distributed Batch

10525         • NQS

10526         • CTSS

10527         • MDQS from Ballistics Research Laboratory (BRL)

10528       Finally, NQS was chosen as a model because it satisfied not only the most user requirements, but
10529       because it was public domain, already implemented on a variety of hardware platforms, and
10530       network-based.

10531       **Historical Implementations of Batch Systems**

10532       Deferred processing of work under the control of a scheduler has been a feature of most
10533       proprietary operating systems from the earliest days of multi-user systems in order to maximize
10534       utilization of the computer.

10535       The arrival of UNIX systems proved to be a dilemma to many hardware providers and users
10536       because it did not include the sophisticated batch facilities offered by the proprietary systems.
10537       This omission was rectified in 1986 by NASA Ames Research Center who developed the
10538       Network Queuing System (NQS) as a portable UNIX application that allowed the routing and
10539       processing of batch ''jobs'' in a network. To encourage its usage, the product was later put into
10540       the public domain. It was promptly picked up by UNIX hardware providers, and ported and
10541       developed for their respective hardware and UNIX implementations.

10542  Many major vendors, who traditionally offer a batch-dominated environment, ported the
10543  public-domain product to their systems, customized it to support the capabilities of their
10544  systems, and added many customer-requested features.

10545  Due to the strong hardware provider and customer acceptance of NQS, it was decided to use
10546  NQS as the basis for the POSIX Batch Environment amendment in 1987. Other batch systems
10547  considered at the time included CTSS, MDQS (a forerunner of NQS from the Ballistics Research
10548  Laboratory), and PROD (a Los Alamos Labs development). None were thought to have both the
10549  functionality and acceptability of NQS.

10550  **NQS Differences from the at utility**

10551  The base standard *at* and *batch* utilities are not sufficient to meet the batch processing needs in a
10552  supercomputing environment and additional functionality in the areas of resource management,
10553  job scheduling, system management, and control of output is required.

10554  **Batch Environment Services and Utilities Option Definitions**                    |

10555  The concept of a batch job is closely related to a session with a session leader. The main  |
10556  difference is that a batch job does not have a controlling terminal. There has been much debate
10557  over whether to use the term ''request'' or ''job''. Job was the final choice because of the
10558  historical use of this term in the batch environment.

10559  The current definition for job identifiers is not sufficient with the model of destinations. The
10560  current definition is:

10561      ```
      sequence_number.originating_host
      ```

10562  Using the model of destination, a host may include multiple batch nodes, the location of which is
10563  identified uniquely by a name or directory service. If the current definition is used, batch nodes
10564  running on the same host would have to coordinate their use of sequence numbers, as sequence
10565  numbers are assigned by the originating host. The alternative is to use the originating batch node
10566  name instead of the originating host name.

10567  The reasons for wishing to run more than one batch system per host could be the following.

10568  A test and production batch system are maintained on a single host.  This is most likely in a
10569  development facility, but could also arise when a site is moving from one version to another.
10570  The new batch system could be installed as a test version that is completely separate from the
10571  production batch system, so that problems can be isolated to the test system. Requiring the batch
10572  nodes to coordinate their use of sequence numbers creates a dependency between the two
10573  nodes, and that defeats the purpose of running two nodes.

10574  A site has multiple departments using a single host, with different management policies. An
10575  example of contention might be in job selection algorithms. One group might want a FIFO type
10576  of selection, while another group wishes to use a more complex algorithm based on resource
10577  availability. Again, requiring the batch nodes to coordinate is an unnecessary binding.

10578  The proposal eventually accepted was to replace originating host with originating batch node.
10579  This supplies sufficient granularity to ensure unique job identifiers. If more than one batch node
10580  is on a particular host, they each have their own unique name.

10581  The queue portion of a destination is not part of the job identifier as these are not required to be
10582  unique between batch nodes. For instance, two batch nodes may both have queues called small,
10583  medium, and large.  It is only the batch node name that is uniquely identifiable throughout the
10584  batch system. The queue name has no additional function in this context.

10585  Assume there are three batch nodes, each of which has its own name server. On batch node one,
10586  there are no queues. On batch node two, there are fifty queues. On batch node three, there are
10587  forty queues.  The system administrator for batch node one does not have to configure queues,
10588  because there are none implemented. However, if a user wishes to send a job to either batch
10589  node two or three, the system administrator for batch node one must configure a destination
10590  that maps to the appropriate batch node and queue. If every queue is to be made accessible from
10591  batch node one, the system administrator has to configure ninety destinations.

10592  To avoid requiring this, there should be a mechanism to allow a user to separate the destination
10593  into a batch node name and a queue name. Then, an implementation that is configured to get to
10594  all the batch nodes does not need any more configuration to allow a user to get to all of the
10595  queues on all of the batch nodes. The node name is used to locate the batch node, while the
10596  queue name is sent unchanged to that batch node.

10597  The following are requirements that a destination identifier must be capable of providing:

10598  - The ability to direct a job to a queue in a particular batch node.

10599  - The ability to direct a job to a particular batch node.

10600  - The ability to group at a higher level than just one queue. This includes grouping similar
10601    queues across multiple batch nodes (this is a pipe queue).

10602  - The ability to group batch nodes. This allows a user to submit a job to a group name with no
10603    knowledge of the batch node configuration. This also provides aliasing as a special case.
10604    Aliasing is a group containing only one batch node name. The group name is the alias.

10605  In addition, the administrator has the following requirements:

10606  - The ability to control access to the queues.

10607  - The ability to control access to the batch nodes.

10608  - The ability to control access to groups of queues (pipe queues).

10609  - The ability to configure retry time intervals and durations.

10610  The requirements of the user are met by destination as explained in the following.

10611  The user has the ability to specify a queue name, which is known only to the batch node
10612  specified. There is no configuration of these queues required on the submitting node.

10613  The user has the ability to specify a batch node whose name is network-unique. The
10614  configuration required is that the batch node be defined as an application, just as other
10615  applications such as FTP are configured.

10616  Once a job reaches a queue, it can again become a user of the batch system. The batch node can
10617  choose to send the job to another batch node or queue or both. In other words, the routing is at
10618  an application level, and it is up to the batch system to choose where the job will be sent.
10619  Configuration is up to the batch node where the queue resides.  This provides grouping of
10620  queues across batch nodes or within a batch node. The user submits the job to a queue, which by
10621  definition routes the job to other queues or nodes or both.

10622  A node name may be given to a naming service, which returns multiple addresses as opposed to
10623  just one. This provides grouping at a batch node level. This is a local issue, meaning that the
10624  batch node must choose only one of these addresses. The list of addresses is not sent with the
10625  job, and once the job is accepted on another node, there is no connection between the list and the
10626  job. The requirements of the administrator are met by destination as explained in the following.

10627  The control of queues is a batch system issue, and will be done using the batch administrative
10628  utilities.

10629    The control of nodes is a network issue, and will be done through whatever network facilities
10630    are available.

10631    The control of access to groups of queues (pipe queues) is covered by the control of any other
10632    queue. The fact that the job may then be sent to another destination is not relevant.

10633    The propagation of a job across more than one point-to-point connection was dropped because
10634    of its complexity and because all of the issues arising from this capability could not be resolved.
10635    It could be provided as additional functionality at some time in the future.

10636    The addition of *network* as a defined term was done to clarify the difference between a network
10637    of batch nodes as opposed to a network of hosts. A network of batch nodes is referred to as a
10638    batch system. The network refers to the actual host configuration. A single host may have
10639    multiple batch nodes.

10640    In the absence of a standard network naming convention, this option establishes its own
10641    convention for the sake of consistency and expediency. This is subject to change, should a future
10642    working group develop a standard naming convention for network pathnames.

10643 **C.3.1    Batch General Concepts**

10644    During the development of the Batch Environment Services and Utilities option, a number of    |
10645    topics were discussed at length which influenced the wording of the normative text but could    |
10646    not be included in the final text. The following items are some of the most significant terms and    |
10647    concepts of those discussed:    |

10648    • Small and Consistent Command Set

10649    Often, conventional utilities from UNIX systems have a very complicated utility syntax and
10650    usage. This can often result in confusion and errors when trying to use them. The Batch    |
10651    Environment Services and Utilities option utility set, on the other hand, has been paired to a    |
10652    small set of robust utilities with an orthogonal calling sequence.    |

10653    • Checkpoint/Restart

10654    This feature permits an already executing process to checkpoint or save its contents. Some
10655    implementations permit this at both the batch utility level (for example, checkpointing this
10656    job upon its abnormal termination) or from within the job itself via a system call. Support of
10657    checkpoint/restart is optional. A conscious, careful effort was made to make the *qsub* utility
10658    consistently refer to checkpoint/restart as optional functionality.

10659    • Rerunability

10660    When a user submits a job for batch processing, they can designate it ''rerunnable'' in that it
10661    will automatically resume execution from the start of the job if the machine on which it was
10662    executing crashes for some reason. The decision on whether the job will be rerun or not is
10663    entirely up to the submitter of the job and no decisions will be made within the batch system.
10664    A job that is rerunnable and has been submitted with the proper checkpoint/restart switch
10665    will first be checkpointed and execution begun from that point. Furthermore, use of the
10666    implementation-defined checkpoint/restart feature will not be defined in this context.

10667    • Error Codes

10668    All utilities exit with error status zero (0) if successful, one (1) if a user error occurred, and
10669    two (2) for an internal Batch Environment Services and Utilities option error.    |

10670    • Level of Portability

10671    Portability is specified at both the user, operator, and administrator levels. A conforming
10672    batch implementation prevents identical functionality and behavior at all these levels.

| 10673 | Additionally, portable batch shell scripts with embedded Batch Environment Services and |
| 10674 | Utilities option utilities add an additional level of portability. |

10675 • Resource Specification

| 10676 | A small set of globally understood resources, such as memory and CPU time, is specified. All |
| 10677 | conforming batch implementations are able to process them in a manner consistent with the |
| 10678 | yet-to-be-developed resource management model. Resources not in this amendment set are |
| 10679 | ignored and passed along as part of the argument stream of the utility. |

10680 • Queue Position

| 10681 | Queue position is the place a job occupies in a queue. It is dependent on a variety of factors |
| 10682 | such as submission time and priority. Since priority may be affected by the implementation |
| 10683 | of fair share scheduling, the definition of queue position is implementation-defined. |

10684 • Queue ID

| 10685 | A numerical queue ID is an external requirement for purposes of accounting. The |
| 10686 | identification number was chosen over queue name for processing convenience. |

10687 • Job ID

| 10688 | A common notion of ''jobs'' is a collection of processes whose process group cannot be |
| 10689 | altered and is used for resource management and accounting. This concept is |
| 10690 | implementation-defined and, as such, has been omitted from the batch amendment. |

10691 • Bytes *versus* Words

| 10692 | Except for one case, bytes are used as the standard unit for memory size. Furthermore, the |
| 10693 | definition of a word varies from machine to machine. Therefore, bytes will be the default unit |
| 10694 | of memory size. |

10695 • Regular Expressions

| 10696 | The standard definition of regular expressions is much too broad to be used in the batch |
| 10697 | utility syntax. All that is needed is a simple concept of ''all''; for example, delete all my jobs |
| 10698 | from the named queue. For this reason, regular expressions have been eliminated from the |
| 10699 | batch amendment. |

10700 • Display Privacy

| 10701 | How much data should be displayed locally through functions? Local policy dictates the |
| 10702 | amount of privacy. Library functions must be used to create and enforce local policy. |
| 10703 | Network and local *qstat*s must reflect the policy of the server machine. |

10704 • Remote Host Naming Convention

| 10705 | It was decided that host names would be a maximum of 255 characters in length, with at |
| 10706 | most 15 characters being shown in displays. The 255 character limit was chosen because it is |
| 10707 | consistent with BSD. The 15-character limit was an arbitrary decision. |

10708 • Network Administration

| 10709 | Network administration is important, but is outside the scope of the batch amendment. |
| 10710 | Network administration could be done with *rsh*. However, authentication becomes two- |
| 10711 | sided. |

10712 • Network Administration Philosophy

| 10713 | Keep it simple. Centralized management should be possible. For example, Los Alamos needs |
| 10714 | a dumb set of CPUs to be managed by a central system *versus* several independently- |

10715    managed systems as is the general case for the Batch Environment Services and Utilities    |
10716    option.    |

10717    • Operator Utility Defaults (that is, Default Host, User, Account, and so on)

10718    It was decided that usability would override orthogonality and syntactic consistency.

10719    • The Batch System Manager and Operator Distinction

10720    The distinction between manager and operator is that operators can only control the flow of
10721    jobs. A manager can alter the batch system configuration in addition to job flow. POSIX
10722    makes a distinction between user and system administrator but goes no further. The
10723    concepts of manager and operator privileges fall under local policy. The distinction between
10724    manager and operator is historical in batch environments, and the Batch Environment    |
10725    Services and Utilities option has continued that distinction.    |

10726    • The Batch System Administrator

10727    An administrator is equivalent to a batch system manager.

## 10728   C.3.2    Batch Services

10729    This rationale is provided as informative rather than normative text, to avoid placing
10730    requirements on implementors regarding the use of symbolic constants, but at the same time to
10731    give implementors a preferred practice for assigning values to these constants to promote
10732    interoperability.

10733    The *Checkpoint* and *Minimum_Cpu_Interval* attributes induce a variety of behavior depending
10734    upon their values. Some jobs cannot or should not be checkpointed. Other users will simply
10735    need to ensure job continuation across planned downtimes; for example, scheduled preventive
10736    maintenance. For users consuming expensive resources, or for jobs that run longer than the
10737    mean time between failures, however, periodic checkpointing may be essential. However,
10738    system administrators must be able to set minimum checkpoint intervals on a queue-by-queue
10739    basis to guard against, for example, naive users specifying interval values too small on
10740    memory-intensive jobs. Otherwise, system overhead would adversely affect performance.

10741    The use of symbolic constants, such as NO_CHECKPOINT, was introduced to lend a degree of
10742    formalism and portability to this option.

10743    Support for checkpointing is optional for servers. However, clients must provide for the −**c**
10744    option, since in a distributed environment the job may run on a server that does provide such
10745    support, even if the host of the client does not support the checkpoint feature.

10746    If the user does not specify the −**c** option, the default action is left unspecified by this option.
10747    Some implementations may wish to do checkpointing by default; others may wish to checkpoint
10748    only under an explicit request from the user.

10749    The *Priority* attribute has been made non-optional. All clients already had been required to
10750    support the −**p** option. The concept of prioritization is common in historical implementations.
10751    The default priority is left to the server to establish.

10752    The *Hold_Types* attribute has been modified to allow for implementation-defined hold types to
10753    be passed to a batch server.

10754    It was the intent of the IEEE P1003.15 working group to mandate the support for the
10755    *Resource_List* attribute in this option by referring to another amendment, specifically the
10756    IEEE P1003.1a draft standard. However, during the development of the IEEE P1003.1a draft
10757    standard this was excluded. As such this requirement has been removed from the normative
10758    text.

10759    The *Shell_Path* attribute has been modified to accept a list of shell paths that are associated with
10760    a host. The name of the attribute has been changed to *Shell_Path_List*.

10761 **C.3.3    Common Behavior for Batch Environment Utilities**

10762    This section was defined to meet the goal of a ''Small and Consistent Command Set'' for this
10763    option.

# 10764 C.4    Utilities

10765    For the utilities included in IEEE Std 1003.1-2001, see the RATIONALE sections on the individual
10766    reference pages.

**Exclusion of Utilities**

10768    The set of utilities contained in IEEE Std 1003.1-2001 is drawn from the base documents, with
10769    one addition: the *c99* utility. This section contains rationale for some of the deliberations that led
10770    to this set of utilities, and why certain utilities were excluded.

10771    Many utilities were evaluated by the standard developers; more historical utilities were
10772    excluded from the base documents than included. The following list contains many common
10773    UNIX system utilities that were not included as mandatory utilities, in the User Portability
10774    Utilities option, in the XSI extension, or in one of the software development groups. It is
10775    logistically difficult for this rationale to distribute correctly the reasons for not including a utility
10776    among the various utility options. Therefore, this section covers the reasons for all utilities not
10777    included in IEEE Std 1003.1-2001.

10778    This rationale is limited to a discussion of only those utilities actively or indirectly evaluated by
10779    the standard developers of the base documents, rather than the list of all known UNIX utilities
10780    from all its variants.

10781    *adb*        The intent of the various software development utilities was to assist in the
10782                 installation (rather than the actual development and debugging) of applications.
10783                 This utility is primarily a debugging tool. Furthermore, many useful aspects of *adb*
10784                 are very hardware-specific.

10785    *as*         Assemblers are hardware-specific and are included implicitly as part of the
10786                 compilers in IEEE Std 1003.1-2001.

10787    *banner*     The only known use of this command is as part of the *lp* printer header pages. It
10788                 was decided that the format of the header is implementation-defined, so this utility
10789                 is superfluous to application portability.

10790    *calendar*   This reminder service program is not useful to conforming applications.

10791    *cancel*     The *lp* (line printer spooling) system specified is the most basic possible and did
10792                 not need this level of application control.

10793    *chroot*     This is primarily of administrative use, requiring superuser privileges.

10794    *col*        No utilities defined in IEEE Std 1003.1-2001 produce output requiring such a filter.
10795                 The *nroff* text formatter is present on many historical systems and will continue to
10796                 remain as an extension; *col* is expected to be shipped by all the systems that ship
10797                 *nroff*.

10798    *cpio*       This has been replaced by *pax*, for reasons explained in the rationale for that utility.

| | | |
|---|---|---|
| 10799 | *cpp* | This is subsumed by *c99*. |
| 10800 10801 | *cu* | This utility is terminal-oriented and is not useful from shell scripts or typical application programs. |
| 10802 10803 10804 10805 | *dc* | The functionality of this utility can be provided by the *bc* utility; *bc* was selected because it was easier to use and had superior functionality. Although the historical versions of *bc* are implemented using *dc* as a base, IEEE Std 1003.1-2001 prescribes the interface and not the underlying mechanism used to implement it. |
| 10806 10807 10808 | *dircmp* | Although a useful concept, the historical output of this directory comparison program is not suitable for processing in application programs. Also, the *diff* –**r** command gives equivalent functionality. |
| 10809 | *dis* | Disassemblers are hardware-specific. |
| 10810 10811 10812 10813 10814 10815 10816 10817 | *emacs* | The community of *emacs* editing enthusiasts was adamant that the full *emacs* editor not be included in the base documents because they were concerned that an attempt to standardize this very powerful environment would encourage vendors to ship versions conforming strictly to the standard, but lacking the extensibility required by the community. The author of the original *emacs* program also expressed his desire to omit the program. Furthermore, there were a number of historical UNIX systems that did not include *emacs*, or included it without supporting it, but there were very few that did not include and support *vi*. |
| 10818 | *ld* | This is subsumed by *c99*. |
| 10819 | *line* | The functionality of *line* can be provided with *read*. |
| 10820 10821 10822 | *lint* | This technology is partially subsumed by *c99*. It is also hard to specify the degree of checking for possible error conditions in programs in any compiler, and specifying what *lint* would do in these cases is equally difficult. |
| 10823 10824 10825 10826 10827 10828 | | It is fairly easy to specify what a compiler does. It requires specifying the language, what it does with that language, and stating that the interpretation of any incorrect program is unspecified. Unfortunately, any description of *lint* is required to specify what to do with erroneous programs. Since the number of possible errors and questionable programming practices is infinite, one cannot require *lint* to detect all errors of any given class. |
| 10829 10830 10831 10832 10833 10834 10835 | | Additionally, some vendors complained that since many compilers are distributed in a binary form without a *lint* facility (because the ISO C standard does not require one), implementing the standard as a stand-alone product will be much harder. Rather than being able to build upon a standard compiler component (simply by providing *c99* as an interface), source to that compiler would most likely need to be modified to provide the *lint* functionality. This was considered a major burden on system providers for a very small gain to developers (users). |
| 10836 10837 | *login* | This utility is terminal-oriented and is not useful from shell scripts or typical application programs. |
| 10838 10839 | *lorder* | This utility is an aid in creating an implementation-defined detail of object libraries that the standard developers did not feel required standardization. |
| 10840 10841 | *lpstat* | The *lp* system specified is the most basic possible and did not need this level of application control. |
| 10842 10843 | *mail* | This utility was omitted in favor of *mailx* because there was a considerable functionality overlap between the two. |

| | | |
|---|---|---|
| 10844<br>10845 | *mknod* | This was omitted in favor of *mkfifo*, as *mknod* has too many implementation-defined functions. |
| 10846<br>10847 | *news* | This utility is terminal-oriented and is not useful from shell scripts or typical application programs. |
| 10848 | *pack* | This compression program was considered inferior to *compress.* |
| 10849<br>10850 | *passwd* | This utility was proposed in a historical draft of the base documents but met with too many objections to be included. There were various reasons: |

<div style="margin-left:2em">

- Changing a password should not be viewed as a command, but as part of the login sequence. Changing a password should only be done while a trusted path is in effect.

- Even though the text in early drafts was intended to allow a variety of implementations to conform, the security policy for one site may differ from another site running with identical hardware and software. One site might use password authentication while the other did not. Vendors could not supply a *passwd* utility that would conform to IEEE Std 1003.1-2001 for all sites using their system.

- This is really a subject for a system administration working group or a security working group.

</div>

| | | |
|---|---|---|
| 10862 | *pcat* | This compression program was considered inferior to *zcat*. |
| 10863<br>10864 | *pg* | This duplicated many of the features of the *more* pager, which was preferred by the standard developers. |
| 10865<br>10866<br>10867 | *prof* | The intent of the various software development utilities was to assist in the installation (rather than the actual development and debugging) of applications. This utility is primarily a debugging tool. |
| 10868<br>10869<br>10870 | RCS | RCS was originally considered as part of a version control utilities portion of the scope. However, this aspect was abandoned by the standard developers. SCCS is now included as an optional part of the XSI extension. |
| 10871<br>10872 | *red* | Restricted editor. This was not considered by the standard developers because it never provided the level of security restriction required. |
| 10873<br>10874<br>10875 | *rsh* | Restricted shell. This was not considered by the standard developers because it does not provide the level of security restriction that is implied by historical documentation. |
| 10876<br>10877<br>10878<br>10879 | *sdb* | The intent of the various software development utilities was to assist in the installation (rather than the actual development and debugging) of applications. This utility is primarily a debugging tool. Furthermore, some useful aspects of *sdb* are very hardware-specific. |
| 10880<br>10881<br>10882 | *sdiff* | The ''side-by-side *diff*'' utility from System V was omitted because it is used infrequently, and even less so by conforming applications. Despite being in System V, it is not in the SVID or XPG. |
| 10883<br>10884 | *shar* | Any of the numerous ''shell archivers'' were excluded because they did not meet the requirement of existing practice. |
| 10885<br>10886<br>10887 | *shl* | This utility is terminal-oriented and is not useful from shell scripts or typical application programs. The job control aspects of the shell command language are generally more useful. |

| | | |
|---|---|---|
| 10888<br>10889<br>10890 | *size* | The intent of the various software development utilities was to assist in the installation (rather than the actual development and debugging) of applications. This utility is primarily a debugging tool. |
| 10891<br>10892<br>10893<br>10894 | *spell* | This utility is not useful from shell scripts or typical application programs. The *spell* utility was considered, but was omitted because there is no known technology that can be used to make it recognize general language for user-specified input without providing a complete dictionary along with the input file. |
| 10895<br>10896 | *su* | This utility is not useful from shell scripts or typical application programs. (There was also sentiment to avoid security-related utilities.) |
| 10897 | *sum* | This utility was renamed *cksum*. |
| 10898 | *tar* | This has been replaced by *pax*, for reasons explained in the rationale for that utility.   \| |
| 10899 | *unpack* | This compression program was considered inferior to *uncompress*. |
| 10900<br>10901 | *wall* | This utility is terminal-oriented and is not useful in shell scripts or typical applications. It is generally used only by system administrators. |

10902　*Rationale (Informative)*

10903　**Part D:**

10904　**Portability Considerations**

10905　*The Open Group*
10906　*The Institute of Electrical and Electronics Engineers, Inc.*

*Appendix D*

# *Portability Considerations (Informative)*

10909    This section contains information to satisfy various international requirements:

10910      • Section D.1 describes perceived user requirements.

10911      • Section D.2 (on page 270) indicates how the facilities of IEEE Std 1003.1-2001 satisfy those
10912        requirements.

10913      • Section D.3 (on page 277) offers guidance to writers of profiles on how the configurable
10914        options, limits, and optional behavior of IEEE Std 1003.1-2001 should be cited in profiles.

## 10915 D.1    User Requirements

10916    This section describes the user requirements that were perceived by the developers of
10917    IEEE Std 1003.1-2001. The primary source for these requirements was an analysis of historical
10918    practice in widespread use, as typified by the base documents listed in Section A.1.1 (on page 3).

10919    IEEE Std 1003.1-2001 addresses the needs of users requiring open systems solutions for source
10920    code portability of applications. It currently addresses users requiring open systems solutions
10921    for source-code portability of applications involving multi-programming and process
10922    management (creating processes, signaling, and so on); access to files and directories in a
10923    hierarchy of file systems (opening, reading, writing, deleting files, and so on); access to
10924    asynchronous communications ports and other special devices; access to information about
10925    other users of the system; facilities supporting applications requiring bounded (realtime)
10926    response.

10927    The following users are identified for IEEE Std 1003.1-2001:

10928      • Those employing applications written in high-level languages, such as C, Ada, or FORTRAN.

10929      • Users who desire conforming applications that do not necessarily require the characteristics
10930        of high-level languages (for example, the speed of execution of compiled languages or the
10931        relative security of source code intellectual property inherent in the compilation process).

10932      • Users who desire conforming applications that can be developed quickly and can be
10933        modified readily without the use of compilers and other system components that may be
10934        unavailable on small systems or those without special application development capabilities.

10935      • Users who interact with a system to achieve general-purpose time-sharing capabilities
10936        common to most business or government offices or academic environments: editing, filing,
10937        inter-user communications, printing, and so on.

10938      • Users who develop applications for POSIX-conformant systems.

10939      • Users who develop applications for UNIX systems.

10940    An acknowledged restriction on applicable users is that they are limited to the group of
10941    individuals who are familiar with the style of interaction characteristic of historically-derived
10942    systems based on one of the UNIX operating systems (as opposed to other historical systems
10943    with different models, such as MS/DOS, Macintosh, VMS, MVS, and so on). Typical users
10944    would include program developers, engineers, or general-purpose time-sharing users.

10945    The requirements of users of IEEE Std 1003.1-2001 can be summarized as a single goal:
10946    *application source portability*. The requirements of the user are stated in terms of the requirements

10947    of portability of applications. This in turn becomes a requirement for a standardized set of
10948    syntax and semantics for operations commonly found on many operating systems.

10949    The following sections list the perceived requirements for application portability.

### 10950   D.1.1    Configuration Interrogation

10951    An application must be able to determine whether and how certain optional features are
10952    provided and to identify the system upon which it is running, so that it may appropriately adapt
10953    to its environment.

10954    Applications must have sufficient information to adapt to varying behaviors of the system.

### 10955   D.1.2    Process Management

10956    An application must be able to manage itself, either as a single process or as multiple processes.
10957    Applications must be able to manage other processes when appropriate.

10958    Applications must be able to identify, control, create, and delete processes, and there must be
10959    communication of information between processes and to and from the system.

10960    Applications must be able to use multiple flows of control with a process (threads) and
10961    synchronize operations between these flows of control.

### 10962   D.1.3    Access to Data

10963    Applications must be able to operate on the data stored on the system, access it, and transmit it
10964    to other applications. Information must have protection from unauthorized or accidental access
10965    or modification.

### 10966   D.1.4    Access to the Environment

10967    Applications must be able to access the external environment to communicate their input and
10968    results.

### 10969   D.1.5    Access to Determinism and Performance Enhancements

10970    Applications must have sufficient control of resource allocation to ensure the timeliness of
10971    interactions with external objects.

### 10972   D.1.6    Operating System-Dependent Profile

10973    The capabilities of the operating system may make certain optional characteristics of the base
10974    language in effect no longer optional, and this should be specified.

### 10975   D.1.7    I/O Interaction

10976    The interaction between the C language I/O subsystem (*stdio*) and the I/O subsystem of
10977    IEEE Std 1003.1-2001 must be specified.

10978 **D.1.8    Internationalization Interaction**

10979   The effects of the environment of IEEE Std 1003.1-2001 on the internationalization facilities of the
10980   C language must be specified.

10981 **D.1.9    C-Language Extensions**

10982   Certain functions in the C language must be extended to support the additional capabilities
10983   provided by IEEE Std 1003.1-2001.

10984 **D.1.10   Command Language**

10985   Users should be able to define procedures that combine simple tools and/or applications into
10986   higher-level components that perform to the specific needs of the user. The user should be able
10987   to store, recall, use, and modify these procedures. These procedures should employ a powerful
10988   command language that is used for recurring tasks in conforming applications (scripts) in the
10989   same way that it is used interactively to accomplish one-time tasks. The language and the
10990   utilities that it uses must be consistent between systems to reduce errors and retraining.

10991 **D.1.11   Interactive Facilities**

10992   Use the system to accomplish individual tasks at an interactive terminal. The interface should be
10993   consistent, intuitive, and offer usability enhancements to increase the productivity of terminal
10994   users, reduce errors, and minimize retraining costs. Online documentation or usage assistance
10995   should be available.

10996 **D.1.12   Accomplish Multiple Tasks Simultaneously**

10997   Access applications and interactive facilities from a single terminal without requiring serial
10998   execution: switch between multiple interactive tasks; schedule one-time or periodic background
10999   work; display the status of all work in progress or scheduled; influence the priority scheduling of
11000   work, when authorized.

11001 **D.1.13   Complex Data Manipulation**

11002   Manipulate data in files in complex ways: sort, merge, compare, translate, edit, format, pattern
11003   match, select subsets (strings, columns, fields, rows, and so on). These facilities should be
11004   available to both conforming applications and interactive users.

11005 **D.1.14   File Hierarchy Manipulation**

11006   Create, delete, move/rename, copy, backup/archive, and display files and directories. These
11007   facilities should be available to both conforming applications and interactive users.

11008 **D.1.15   Locale Configuration**

11009   Customize applications and interactive sessions for the cultural and language conventions of the
11010   user. Employ a wide variety of standard character encodings. These facilities should be available
11011   to both conforming applications and interactive users.

11012 **D.1.16  Inter-User Communication**

11013       Send messages or transfer files to other users on the same system or other systems on a network.
11014       These facilities should be available to both conforming applications and interactive users.

11015 **D.1.17  System Environment**

11016       Display information about the status of the system (activities of users and their interactive and
11017       background work, file system utilization, system time, configuration, and presence of optional
11018       facilities) and the environment of the user (terminal characteristics, and so on). Inform the
11019       system operator/administrator of problems. Control access to user files and other resources.

11020 **D.1.18  Printing**

11021       Output files on a variety of output device classes, accessing devices on local or network-
11022       connected systems. Control (or influence) the formatting, priority scheduling, and output
11023       distribution of work. These facilities should be available to both conforming applications and
11024       interactive users.

11025 **D.1.19  Software Development**

11026       Develop (create and manage source files, compile/interpret, debug) portable open systems
11027       applications and package them for distribution to, and updating of, other systems.

11028 **D.2  Portability Capabilities**

11029       This section describes the significant portability capabilities of IEEE Std 1003.1-2001 and
11030       indicates how the user requirements listed in Section D.1 (on page 267) are addressed. The
11031       capabilities are listed in the same format as the preceding user requirements; they are
11032       summarized below:

11033       • Configuration Interrogation

11034       • Process Management

11035       • Access to Data

11036       • Access to the Environment

11037       • Access to Determinism and Performance Enhancements

11038       • Operating System-Dependent Profile

11039       • I/O Interaction

11040       • Internationalization Interaction

11041       • C-Language Extensions

11042       • Command Language

11043       • Interactive Facilities

11044       • Accomplish Multiple Tasks Simultaneously

11045       • Complex Data Manipulation

11046       • File Hierarchy Manipulation

<div style="margin-left:2em">

11047       • Locale Configuration

11048       • Inter-User Communication

11049       • System Environment

11050       • Printing

11051       • Software Development

</div>

## 11052 D.2.1   Configuration Interrogation

11053 The *uname*( ) operation provides basic identification of the system. The *sysconf*( ), *pathconf*( ), and
11054 *fpathconf*( ) functions and the *getconf* utility provide means to interrogate the implementation to
11055 determine how to adapt to the environment in which it is running. These values can be either
11056 static (indicating that all instances of the implementation have the same value) or dynamic
11057 (indicating that different instances of the implementation have the different values, or that the
11058 value may vary for other reasons, such as reconfiguration).

### 11059 Unsatisfied Requirements

11060 None directly. However, as new areas are added, there will be a need for additional capability in
11061 this area.

## 11062 D.2.2   Process Management

11063 The *fork*( ), *exec* family, *posix_spawn*( ), and *posix_spawnp*( ) functions provide for the creation of
11064 new processes or the insertion of new applications into existing processes. The *_Exit*( ), *_exit*( ),
11065 *exit*( ), and *abort*( ) functions allow for the termination of a process by itself. The *wait*( ) and
11066 *waitpid*( ) functions allow one process to deal with the termination of another.

11067 The *times*( ) function allows for basic measurement of times used by a process. Various
11068 functions, including *fstat*( ), *getegid*( ), *geteuid*( ), *getgid*( ), *getgrgid*( ), *getgrnam*( ), *getlogin*( ),
11069 *getpid*( ), *getppid*( ), *getpwnam*( ), *getpwuid*( ), *getuid*( ), *lstat*( ), and *stat*( ), provide for access to the
11070 identifiers of processes and the identifiers and names of owners of processes (and files).

11071 The various functions operating on environment variables provide for communication of
11072 information (primarily user-configurable defaults) from a parent to child processes.

11073 The operations on the current working directory control and interrogate the directory from
11074 which relative filename searches start. The *umask*( ) function controls the default protections
11075 applied to files created by the process.

11076 The *alarm*( ), *pause*( ), *sleep*( ), *ualarm*( ), and *usleep*( ) operations allow the process to suspend until
11077 a timer has expired or to be notified when a period of time has elapsed. The *time*( ) operation
11078 interrogates the current time and date.

11079 The signal mechanism provides for communication of events either from other processes or
11080 from the environment to the application, and the means for the application to control the effect
11081 of these events. The mechanism provides for external termination of a process and for a process
11082 to suspend until an event occurs. The mechanism also provides for a value to be associated with
11083 an event.

11084 Job control provides a means to group processes and control them as groups, and to control their
11085 access to the function between the user and the system (the ''controlling terminal''). It also
11086 provides the means to suspend and resume processes.

11087 The Process Scheduling option provides control of the scheduling and priority of a process.

11088   The Message Passing option provides a means for interprocess communication involving small
11089   amounts of data.

11090   The Memory Management facilities provide control of memory resources and for the sharing of
11091   memory. This functionality is mandatory on XSI-conformant systems.

11092   The Threads facilities provide multiple flows of control with a process (threads),
11093   synchronization between threads, association of data with threads, and controlled cancellation   |
11094   of threads.

11095   The XSI interprocess communications functionality provide an alternate set of facilities to
11096   manipulate semaphores, message queues, and shared memory. These are provided on XSI-
11097   conformant systems to support conforming applications developed to run on UNIX systems.

11098 ## D.2.3   Access to Data

11099   The *open*( ), *close*( ), *fclose*( ), *fopen*( ), and *pipe*( ) functions provide for access to files and data.
11100   Such files may be regular files, interprocess data channels (pipes), or devices.  Additional types
11101   of objects in the file system are permitted and are being contemplated for standardization.

11102   The *access*( ), *chmod*( ), *chown*( ), *dup*( ), *dup2*( ), *fchmod*( ), *fcntl*( ), *fstat*( ), *ftruncate*( ), *lstat*( ),
11103   *readlink*( ), *realpath*( ), *stat*( ), and *utime*( ) functions allow for control and interrogation of file and
11104   file-related objects (including symbolic links), and their ownership, protections, and timestamps.

11105   The *fgetc*( ), *fputc*( ), *fread*( ), *fseek*( ), *fsetpos*( ), *fwrite*( ), *getc*( ), *getchar*( ), *lseek*( ), *putchar*( ), *putc*( ),
11106   *read*( ), and *write*( ) functions provide for data transfer from the application to files (in all their
11107   forms).

11108   The *closedir*( ), *link*( ), *mkdir*( ), *opendir*( ), *readdir*( ), *rename*( ), *rmdir*( ), *rewinddir*( ), and *unlink*( )
11109   functions provide for a complete set of operations on directories.  Directories can arbitrarily
11110   contain other directories, and a single file can be mentioned in more than one directory.

11111   The file-locking mechanism provides for advisory locking (protection during transactions) of
11112   ranges of bytes (in effect, records) in a file.

11113   The *confstr*( ), *fpathconf*( ), *pathconf*( ), and *sysconf*( ) functions provide for enquiry as to the
11114   behavior of the system where variability is permitted.

11115   The Synchronized Input and Output option provides for assured commitment of data to media.

11116   The Asynchronous Input and Output option provides for initiation and control of asynchronous
11117   data transfers.

11118 ## D.2.4   Access to the Environment

11119   The operations and types in the Base Definitions volume of IEEE Std 1003.1-2001, Chapter 11,
11120   General Terminal Interface are provided for access to asynchronous serial devices. The primary
11121   intended use for these is the controlling terminal for the application (the interaction point
11122   between the user and the system). They are general enough to be used to control any
11123   asynchronous serial device. The functions are also general enough to be used with many other
11124   device types as a user interface when some emulation is provided.

11125   Less detailed access is provided for other device types, but in many instances an application
11126   need not know whether an object in the file system is a device or a regular file to operate
11127   correctly.

11128      **Unsatisfied Requirements**

11129      Detailed control of common device classes, specifically magnetic tape, is not provided.

## 11130 D.2.5   Bounded (Realtime) Response

11131      The Realtime Signals Extension provides queued signals and the prioritization of the handling of
11132      signals. The SCHED_FIFO, SCHED_SPORADIC, and SCHED_RR scheduling policies provide
11133      control over processor allocation. The Semaphores option provides high-performance
11134      synchronization. The Memory Management functions provide memory locking for control of
11135      memory allocation, file mapping for high-performance, and shared memory for high-
11136      performance interprocess communication. The Message Passing option provides for interprocess
11137      communication without being dependent on shared memory.

11138      The Timers option provides a high resolution function called *nanosleep*( ) with a finer resolution
11139      than the *sleep*( ) function.

11140      The Typed Memory Objects option, the Monotonic Clock option, and the Timeouts option
11141      provide further facilities for applications to use to obtain predictable bounded response.

## 11142 D.2.6   Operating System-Dependent Profile

11143      IEEE Std 1003.1-2001 makes no distinction between text and binary files. The values of
11144      EXIT_SUCCESS and EXIT_FAILURE are further defined.

11145      **Unsatisfied Requirements**

11146      None known, but the ISO C standard may contain some additional options that could be
11147      specified.

## 11148 D.2.7   I/O Interaction

11149      IEEE Std 1003.1-2001 defines how each of the ISO C standard *stdio* functions interact with the
11150      POSIX.1 operations, typically specifying the behavior in terms of POSIX.1 operations.

11151      **Unsatisfied Requirements**

11152      None.

## 11153 D.2.8   Internationalization Interaction

11154      The IEEE Std 1003.1-2001 environment operations provide a means to define the environment
11155      for *setlocale*( ) and time functions such as *ctime*( ). The *tzset*( ) function is provided to set time
11156      conversion information.

11157      The *nl_langinfo*( ) function is provided as an XSI extension to query locale-specific cultural
11158      settings.

11159      **Unsatisfied Requirements**

11160      None.

### 11161 D.2.9    C-Language Extensions

11162 The *setjmp*( ) and *longjmp*( ) functions are not defined to be cognizant of the signal masks defined
11163 for POSIX.1. The *sigsetjmp*( ) and *siglongjmp*( ) functions are provided to fill this gap.

11164 The *_setjmp*( ) and *_longjmp*( ) functions are provided as XSI extensions to support historic
11165 practice.

11166 **Unsatisfied Requirements**

11167 None.

### 11168 D.2.10    Command Language

11169 The shell command language, as described in the Shell and Utilities volume of
11170 IEEE Std 1003.1-2001, Chapter 2, Shell Command Language, is a common language useful in
11171 batch scripts, through an API to high-level languages (for the C-Language Binding option,
11172 *system*( ) and *popen*( )) and through an interactive terminal (see the *sh* utility). The shell language
11173 has many of the characteristics of a high-level language, but it has been designed to be more
11174 suitable for user terminal entry and includes interactive debugging facilities.  Through the use of
11175 pipelining, many complex commands can be constructed from combinations of data filters and
11176 other common components. Shell scripts can be created, stored, recalled, and modified by the
11177 user with simple editors.

11178 In addition to the basic shell language, the following utilities offer features that simplify and
11179 enhance programmatic access to the utilities and provide features normally found only in high-
11180 level languages: *basename*, *bc*, *command*, *dirname*, *echo*, *env*, *expr*, *false*, *printf*, *read*, *sleep*, *tee*, *test*,
11181 *time*\*,[2] *true*, *wait*, *xargs*, and all of the special built-in utilities in the Shell and Utilities volume of
11182 IEEE Std 1003.1-2001, Section 2.14, Special Built-In Utilities.

11183 **Unsatisfied Requirements**

11184 None.

### 11185 D.2.11    Interactive Facilities

11186 The utilities offer a common style of command-line interface through conformance to the Utility
11187 Syntax Guidelines (see the Base Definitions volume of IEEE Std 1003.1-2001, Section 12.2, Utility
11188 Syntax Guidelines) and the common utility defaults (see the Shell and Utilities volume of
11189 IEEE Std 1003.1-2001, Section 1.11, Utility Description Defaults). The *sh* utility offers an
11190 interactive command-line history and editing facility. The following utilities in the User
11191 Portability Utilities option have been customized for interactive use: *alias*, *ex*, *fc*, *mailx*, *more*, *talk*,
11192 *vi*, *unalias*, and *write*; the *man* utility offers online access to system documentation.

11193 _____
11194 2.  The utilities listed with an asterisk here and later in this section are present only on systems which support the User Portability
11195     Utilities option. There may be further restrictions on the utilities offered with various configuration option combinations; see the
11196     individual utility descriptions.

11197 **Unsatisfied Requirements**

11198 The command line interface to individual utilities is as intuitive and consistent as historical
11199 practice allows. Work underway based on graphical user interfaces may be more suitable for
11200 novice or occasional users of the system.

11201 ## D.2.12 Accomplish Multiple Tasks Simultaneously

11202 The shell command language offers background processing through the asynchronous list
11203 command form; see the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.9, Shell
11204 Commands. The *nohup* utility makes background processing more robust and usable. The *kill*
11205 utility can terminate background jobs. When the User Portability Utilities option is supported,
11206 the following utilities allow manipulation of jobs: *bg, fg,* and *jobs.* Also, if the User Portability
11207 Utilities option is supported, the following can support periodic job scheduling, control, and
11208 display: *at*, *batch*, *crontab*, *nice*, *ps*, and *renice*.

11209 **Unsatisfied Requirements**

11210 Terminals with multiple windows may be more suitable for some multi-tasking interactive uses
11211 than the job control approach in IEEE Std 1003.1-2001. See the comments on graphical user
11212 interfaces in Section D.2.11 (on page 274). The *nice* and *renice* utilities do not necessarily take
11213 advantage of complex system scheduling algorithms that are supported by the realtime options
11214 within IEEE Std 1003.1-2001.

11215 ## D.2.13 Complex Data Manipulation

11216 The following utilities address user requirements in this area: *asa, awk, bc, cmp, comm, csplit\*, cut,*
11217 *dd, diff, ed, ex\*, expand\*, expr, find, fold, grep, head, join, od, paste, pr, printf, sed, sort, split\*, tabs\*, tail,*
11218 *tr, unexpand\*, uniq, uudecode\*, uuencode\*,* and *wc.*

11219 **Unsatisfied Requirements**

11220 Sophisticated text formatting utilities, such as *troff* or *TeX*, are not included. Standards work in
11221 the area of SGML may satisfy this.

11222 ## D.2.14 File Hierarchy Manipulation

11223 The following utilities address user requirements in this area: *basename, cd, chgrp, chmod, chown,*
11224 *cksum, cp, dd, df\*, diff, dirname, du\*, find, ls, ln, mkdir, mkfifo, mv, patch\*, pathchk, pax, pwd, rm, rmdir,*
11225 *test,* and *touch.*

11226 **Unsatisfied Requirements**

11227 Some graphical user interfaces offer more intuitive file manager components that allow file
11228 manipulation through the use of icons for novice users.

11229 **D.2.15  Locale Configuration**

11230 The standard utilities are affected by the various *LC_* variables to achieve locale-dependent
11231 operation: character classification, collation sequences, regular expressions and shell pattern
11232 matching, date and time formats, numeric formatting, and monetary formatting. When the
11233 POSIX2_LOCALEDEF option is supported, applications can provide their own locale definition
11234 files. The following utilities address user requirements in this area: *date*, *ed*, *ex**, *find*, *grep*, *locale*,
11235 *localedef*, *more**, *sed*, *sh*, *sort*, *tr*, *uniq*, and *vi**.

11236 The *iconv*( ), *iconv_close*( ), and *iconv_open*( ) functions are available to allow an application to
11237 convert character data between supported character sets.

11238 The *gencat* utility and the *catopen*( ), *catclose*( ), and *catgets*( ) functions for message catalog
11239 manipulation are available on XSI-conformant systems.

11240 **Unsatisfied Requirements**

11241 Some aspects of multi-byte character and state-encoded character encodings have not yet been
11242 addressed. The C-language functions, such as *getopt*( ), are generally limited to single-byte
11243 characters. The effect of the *LC_MESSAGES* variable on message formats is only suggested at
11244 this time.

11245 **D.2.16  Inter-User Communication**

11246 The following utilities address user requirements in this area: *cksum*, *mailx**, *mesg**, *patch**, *pax*,
11247 *talk**, *uudecode**, *uuencode**, *who**, and *write**.

11248 The historical UUCP utilities are included on XSI-conformant systems.

11249 **Unsatisfied Requirements**

11250 None.

11251 **D.2.17  System Environment**

11252 The following utilities address user requirements in this area: *chgrp*, *chmod*, *chown*, *df**, *du**, *env*,
11253 *getconf*, *id*, *logger*, *logname*, *mesg**, *newgrp**, *ps**, *stty*, *tput**, *tty*, *umask*, *uname*, and *who**.

11254 The *closelog*( ), *openlog*( ), *setlogmask*( ), and *syslog*( ) functions provide System Logging facilities
11255 on XSI-conformant systems; these are analogous to the *logger* utility.

11256 **Unsatisfied Requirements**

11257 None.

11258 **D.2.18  Printing**

11259 The following utilities address user requirements in this area: *pr* and *lp*.

11260 **Unsatisfied Requirements**

11261 There are no features to control the formatting or scheduling of the print jobs.

11262 **D.2.19  Software Development**

11263    The following utilities address user requirements in this area: *ar, asa, awk, c99, ctags\*, fort77,*
11264    *getconf, getopts, lex, localedef, make, nm\*, od, patch\*, pax, strings\*, strip, time\*,* and *yacc.*

11265    The *system*(), *popen*(), *pclose*(), *regcomp*(), *regexec*(), *regerror*(), *regfree*(), *fnmatch*(), *getopt*(),
11266    *glob*(), *globfree*(), *wordexp*(), and *wordfree*() functions allow C-language programmers to access
11267    some of the interfaces used by the utilities, such as argument processing, regular expressions,
11268    and pattern matching.

11269    The SCCS source-code control system utilities are available on systems supporting the XSI
11270    Development option.

11271    **Unsatisfied Requirements**

11272    There are no language-specific development tools related to languages other than C and
11273    FORTRAN. The C tools are more complete and varied than the FORTRAN tools. There is no
11274    data dictionary or other CASE-like development tools.

11275 **D.2.20  Future Growth**

11276    It is arguable whether or not all functionality to support applications is potentially within the
11277    scope of IEEE Std 1003.1-2001. As a simple matter of practicality, it cannot be. Areas such as
11278    graphics, application domain-specific functionality, windowing, and so on, should be in unique
11279    standards. As such, they are properly ''Unsatisfied Requirements'' in terms of providing fully
11280    conforming applications, but ones which are outside the scope of IEEE Std 1003.1-2001.

11281    However, as the standards evolve, certain functionality once considered ''exotic'' enough to be
11282    part of a separate standard become common enough to be included in a core standard such as
11283    this. Realtime and networking, for example, have both moved from separate standards (with
11284    much difficult cross-referencing) into IEEE Std 1003.1 over time, and although no specific areas
11285    have been identified for inclusion in future revisions, such inclusions seem likely.

11286 **D.3      Profiling Considerations**

11287    This section offers guidance to writers of profiles on how the configurable options, limits, and
11288    optional behavior of IEEE Std 1003.1-2001 should be cited in profiles. Profile writers should
11289    consult the general guidance in POSIX.0 when writing POSIX Standardized Profiles.

11290    The information in this section is an inclusive list of features that should be considered by profile
11291    writers. Subsetting of IEEE Std 1003.1-2001 should follow the Base Definitions volume of
11292    IEEE Std 1003.1-2001, Section 2.1.5.1, Subprofiling Considerations. A set of profiling options is
11293    described in Appendix E (on page 291).

11294 **D.3.1    Configuration Options**

11295    There are two set of options suggested by IEEE Std 1003.1-2001: those for POSIX-conforming
11296    systems and those for X/Open System Interface (XSI) conformance. The requirements for XSI
11297    conformance are documented in the Base Definitions volume of IEEE Std 1003.1-2001 and not
11298    discussed further here, as they superset the POSIX conformance requirements.

## D.3.2    Configuration Options (Shell and Utilities)

11299

11300 There are three broad optional configurations for the Shell and Utilities volume of
11301 IEEE Std 1003.1-2001: basic execution system, development system, and user portability
11302 interactive system. The options to support these, and other minor configuration options, are
11303 listed in the Base Definitions volume of IEEE Std 1003.1-2001, Chapter 2, Conformance. Profile
11304 writers should consult the following list and the comments concerning user requirements
11305 addressed by various components in Section D.2 (on page 270).

11306 POSIX2_UPE
11307 The system supports the User Portability Utilities option.

11308 This option is a requirement for a user portability interactive system. It is required
11309 frequently except for those systems, such as embedded realtime or dedicated application
11310 systems, that support little or no interactive time-sharing work by users or operators. XSI-
11311 conformant systems support this option.

11312 POSIX2_SW_DEV
11313 The system supports the Software Development Utilities option.

11314 This option is required by many systems, even those in which actual software development
11315 does not occur. The *make* utility, in particular, is required by many application software
11316 packages as they are installed onto the system. If POSIX2_C_DEV is supported,
11317 POSIX2_SW_DEV is almost a mandatory requirement because of *ar* and *make*.

11318 POSIX2_C_BIND
11319 The system supports the C-Language Bindings option.

11320 This option is required on some implementations developing complex C applications or on
11321 any system installing C applications in source form that require the functions in this option.
11322 The *system*() and *popen*() functions, in particular, are widely used by applications; the
11323 others are rather more specialized.

11324 POSIX2_C_DEV
11325 The system supports the C-Language Development Utilities option.

11326 This option is required by many systems, even those in which actual C-language software
11327 development does not occur. The *c99* utility, in particular, is required by many application
11328 software packages as they are installed onto the system. The *lex* and *yacc* utilities are used
11329 less frequently.

11330 POSIX2_FORT_DEV
11331 The system supports the FORTRAN Development Utilities option

11332 As with C, this option is needed on any system developing or installing FORTRAN
11333 applications in source form.

11334 POSIX2_FORT_RUN
11335 The system supports the FORTRAN Runtime Utilities option.

11336 This option is required for some FORTRAN applications that need the *asa* utility to convert
11337 Hollerith printing statement output. It is unknown how frequently this occurs.

11338 POSIX2_LOCALEDEF
11339 The system supports the creation of locales.

11340 This option is needed if applications require their own customized locale definitions to
11341 operate. It is presently unknown whether many applications are dependent on this.
11342 However, the option is virtually mandatory for systems in which internationalized
11343 applications are developed.

11344            XSI-conformant systems support this option.

11345      POSIX2_PBS
11346            The system supports the Batch Environment Services and Utilities option.                    |

11347      POSIX2_PBS_ACCOUNTING
11348            The system supports the optional feature of accounting within the Batch Environment        |
11349            Services and Utilities option. It will be required in servers that implement the optional  |
11350            feature of accounting.                                                                      |

11351      POSIX2_PBS_CHECKPOINT
11352            The system supports the optional feature of checkpoint/restart within the Batch            |
11353            Environment Services and Utilities option.                                                  |

11354      POSIX2_PBS_LOCATE
11355            The system supports the optional feature of locating batch jobs within the Batch           |
11356            Environment Services and Utilities option.                                                  |

11357      POSIX2_PBS_MESSAGE
11358            The system supports the optional feature of sending messages to batch jobs within the      |
11359            Batch Environment Services and Utilities option.                                            |

11360      POSIX2_PBS_TRACK
11361            The system supports the optional feature of tracking batch jobs within the Batch           |
11362            Environment Services and Utilities option.                                                  |

11363      POSIX2_CHAR_TERM
11364            The system supports at least one terminal type capable of all operations described in
11365            IEEE Std 1003.1-2001.

11366            On systems with POSIX2_UPE, this option is almost always required. It was developed
11367            solely to allow certain specialized vendors and user applications to bypass the requirement
11368            for general-purpose asynchronous terminal support. For example, an application and
11369            system that was suitable for block-mode terminals, such as IBM 3270s, would not need this
11370            option.

11371            XSI-conformant systems support this option.

## 11372 D.3.3    Configurable Limits

11373      Very few of the limits need to be increased for profiles. No profile can cite lower values.

11374      {POSIX2_BC_BASE_MAX}
11375      {POSIX2_BC_DIM_MAX}
11376      {POSIX2_BC_SCALE_MAX}
11377      {POSIX2_BC_STRING_MAX}
11378            No increase is anticipated for any of these *bc* values, except for very specialized applications
11379            involving huge numbers.

11380      {POSIX2_COLL_WEIGHTS_MAX}
11381            Some natural languages with complex collation requirements require an increase from the
11382            default 2 to 4; no higher numbers are anticipated.

11383      {POSIX2_EXPR_NEST_MAX}
11384            No increase is anticipated.

11385      {POSIX2_LINE_MAX}
11386            This number is much larger than most historical applications have been able to use. At some
11387            future time, applications may be rewritten to take advantage of even larger values.

11388   {POSIX2_RE_DUP_MAX}
11389           No increase is anticipated.

11390   {POSIX2_VERSION}
11391           This is actually not a limit, but a standard version stamp. Generally, a profile should specify
11392           the Shell and Utilities volume of IEEE Std 1003.1-2001, Chapter 2, Shell Command Language
11393           by name in the normative references section, not this value.

## 11394  D.3.4    Configuration Options (System Interfaces)

11395   {NGROUPS_MAX}
11396           A non-zero value indicates that the implementation supports supplementary groups.

11397           This option is needed where there is a large amount of shared use of files, but where a
11398           certain amount of protection is needed. Many profiles[3] are known to require this option; it
11399           should only be required if needed, but it should never be prohibited.

11400   _POSIX_ADVISORY_INFO
11401           The system provides advisory information for file management.

11402           This option allows the application to specify advisory information that can be used to
11403           achieve better or even deterministic response time in file manager or input and output
11404           operations.

11405   _POSIX_ASYNCHRONOUS_IO
11406           The system provides concurrent process execution and input and output transfers.

11407           This option was created to support historical systems that did not provide the feature. It
11408           should only be required if needed, but it should never be prohibited.

11409   _POSIX_BARRIERS
11410           The system supports barrier synchronization.

11411           This option was created to allow efficient synchronization of multiple parallel threads in
11412           multi-processor systems in which the operation is supported in part by the hardware
11413           architecture.

11414   _POSIX_CHOWN_RESTRICTED
11415           The system restricts the right to ''give away'' files to other users.

11416           This option should be carefully investigated before it is required. Some applications expect
11417           that they can change the ownership of files in this way. It is provided where either security
11418           or system account requirements cause this ability to be a problem. It is also known to be
11419           specified in many profiles.

11420   _POSIX_CLOCK_SELECTION
11421           The system supports the Clock Selection option.

11422           This option allows applications to request a high resolution sleep in order to suspend a
11423           thread during a relative time interval, or until an absolute time value, using the desired
11424           clock. It also allows the application to select the clock used in a *pthread_cond_timedwait*( )
11425           function call.

11426   _____
11427   3.  There are no formally approved profiles of IEEE Std 1003.1-2001 at the time of publication; the reference here is to various
11428       profiles generated by private bodies or governments.

11429      _POSIX_CPUTIME
11430          The system supports the Process CPU-Time Clocks option.

11431          This option allows applications to use a new clock that measures the execution times of
11432          processes or threads, and the possibility to create timers based upon these clocks, for
11433          runtime detection (and treatment) of execution time overruns.

11434      _POSIX_FSYNC
11435          The system supports file synchronization requests.

11436          This option was created to support historical systems that did not provide the feature.
11437          Applications that are expecting guaranteed completion of their input and output operations
11438          should require the _POSIX_SYNC_IO option. This option should never be prohibited.

11439          XSI-conformant systems support this option.

11440      _POSIX_IPV6
11441          The system supports facilities related to Internet Protocol Version 6 (IPv6).

11442          This option was created to allow systems to transition to IPv6.

11443      _POSIX_JOB_CONTROL
11444          Job control facilities are mandatory in IEEE Std 1003.1-2001.

11445          The option was created primarily to support historical systems that did not provide the
11446          feature. Many existing profiles now require it; it should only be required if needed, but it
11447          should never be prohibited. Most applications that use it can run when it is not present,
11448          although with a degraded level of user convenience.

11449      _POSIX_MAPPED_FILES
11450          The system supports the mapping of regular files into the process address space.

11451          XSI-conformant systems support this option.

11452          Both this option and the Shared Memory Objects option provide shared access to memory
11453          objects in the process address space. The functions defined under this option provide the
11454          functionality of existing practice for mapping regular files. This functionality was deemed
11455          unnecessary, if not inappropriate, for embedded systems applications and, hence, is
11456          provided under this option. It should only be required if needed, but it should never be
11457          prohibited.

11458      _POSIX_MEMLOCK
11459          The system supports the locking of the address space.

11460          This option was created to support historical systems that did not provide the feature. It
11461          should only be required if needed, but it should never be prohibited.

11462      _POSIX_MEMLOCK_RANGE
11463          The system supports the locking of specific ranges of the address space.

11464          For applications that have well-defined sections that need to be locked and others that do
11465          not, IEEE Std 1003.1-2001 supports an optional set of functions to lock or unlock a range of
11466          process addresses. The following are two reasons for having a means to lock down a
11467          specific range:

11468          1.  An asynchronous event handler function that must respond to external events in a
11469              deterministic manner such that page faults cannot be tolerated

11470          2.  An input/output ''buffer'' area that is the target for direct-to-process I/O, and the
11471              overhead of implicit locking and unlocking for each I/O call cannot be tolerated

11472          It should only be required if needed, but it should never be prohibited.

11473    _POSIX_MEMORY_PROTECTION
11474          The system supports memory protection.

11475          XSI-conformant systems support this option.

11476          The provision of this option typically imposes additional hardware requirements. It should
11477          never be prohibited.

11478    _POSIX_PRIORITIZED_IO
11479          The system provides prioritization for input and output operations.

11480          The use of this option may interfere with the ability of the system to optimize input and
11481          output throughput. It should only be required if needed, but it should never be prohibited.

11482    _POSIX_MESSAGE_PASSING
11483          The system supports the passing of messages between processes.

11484          This option was created to support historical systems that did not provide the feature. The
11485          functionality adds a high-performance XSI interprocess communication facility for local
11486          communication. It should only be required if needed, but it should never be prohibited.

11487    _POSIX_MONOTONIC_CLOCK
11488          The system supports the Monotonic Clock option.

11489          This option allows realtime applications to rely on a monotonically increasing clock that
11490          does not jump backwards, and whose value does not change except for the regular ticking
11491          of the clock.

11492    _POSIX_PRIORITY_SCHEDULING
11493          The system provides priority-based process scheduling.

11494          Support of this option provides predictable scheduling behavior, allowing applications to
11495          determine the order in which processes that are ready to run are granted access to a
11496          processor. It should only be required if needed, but it should never be prohibited.

11497    _POSIX_REALTIME_SIGNALS
11498          The system provides prioritized, queued signals with associated data values.

11499          This option was created to support historical systems that did not provide the features. It
11500          should only be required if needed, but it should never be prohibited.

11501    _POSIX_REGEXP
11502          Support for regular expression facilities is mandatory in IEEE Std 1003.1-2001.

11503    _POSIX_SAVED_IDS
11504          Support for this feature is mandatory in IEEE Std 1003.1-2001.

11505          Certain classes of applications rely on it for proper operation, and there is no alternative
11506          short of giving the application root privileges on most implementations that did not provide
11507          _POSIX_SAVED_IDS.

11508    _POSIX_SEMAPHORES
11509          The system provides counting semaphores.

11510          This option was created to support historical systems that did not provide the feature. It
11511          should only be required if needed, but it should never be prohibited.

11512    _POSIX_SHARED_MEMORY_OBJECTS
11513          The system supports the mapping of shared memory objects into the process address space.

| | |
|---|---|
| 11514 | Both this option and the Memory Mapped Files option provide shared access to memory |
| 11515 | objects in the process address space. The functions defined under this option provide the |
| 11516 | functionality of existing practice for shared memory objects. This functionality was deemed |
| 11517 | appropriate for embedded systems applications and, hence, is provided under this option. It |
| 11518 | should only be required if needed, but it should never be prohibited. |

11519 _POSIX_SHELL

11520    Support for the *sh* utility command line interpreter is mandatory in IEEE Std 1003.1-2001.

11521 _POSIX_SPAWN

11522    The system supports the spawn option.

11523    This option provides applications with an efficient mechanism to spawn execution of a new
11524    process.

11525 _POSIX_SPINLOCKS

11526    The system supports spin locks.

11527    This option was created to support a simple and efficient synchronization mechanism for
11528    threads executing in multi-processor systems.

11529 _POSIX_SPORADIC_SERVER

11530    The system supports the sporadic server scheduling policy.

11531    This option provides applications with a new scheduling policy for scheduling aperiodic
11532    processes or threads in hard realtime applications.

11533 _POSIX_SYNCHRONIZED_IO

11534    The system supports guaranteed file synchronization.

11535    This option was created to support historical systems that did not provide the feature.
11536    Applications that are expecting guaranteed completion of their input and output operations
11537    should require this option, rather than the File Synchronization option. It should only be
11538    required if needed, but it should never be prohibited.

11539 _POSIX_THREADS

11540    The system supports multiple threads of control within a single process.

11541    This option was created to support historical systems that did not provide the feature.
11542    Applications written assuming a multi-threaded environment would be expected to require
11543    this option. It should only be required if needed, but it should never be prohibited.

11544    XSI-conformant systems support this option.

11545 _POSIX_THREAD_ATTR_STACKADDR

11546    The system supports specification of the stack address for a created thread.

11547    Applications may take advantage of support of this option for performance benefits, but
11548    dependence on this feature should be minimized. This option should never be prohibited.

11549    XSI-conformant systems support this option.

11550 _POSIX_THREAD_ATTR_STACKSIZE

11551    The system supports specification of the stack size for a created thread.

11552    Applications may require this option in order to ensure proper execution, but such usage
11553    limits portability and dependence on this feature should be minimized. It should only be
11554    required if needed, but it should never be prohibited.

11555    XSI-conformant systems support this option.

11556    _POSIX_THREAD_PRIORITY_SCHEDULING
11557         The system provides priority-based thread scheduling.

11558         Support of this option provides predictable scheduling behavior, allowing applications to
11559         determine the order in which threads that are ready to run are granted access to a processor.
11560         It should only be required if needed, but it should never be prohibited.

11561    _POSIX_THREAD_PRIO_INHERIT
11562         The system provides mutual-exclusion operations with priority inheritance.

11563         Support of this option provides predictable scheduling behavior, allowing applications to
11564         determine the order in which threads that are ready to run are granted access to a processor.
11565         It should only be required if needed, but it should never be prohibited.

11566    _POSIX_THREAD_PRIO_PROTECT
11567         The system supports a priority ceiling emulation protocol for mutual-exclusion operations.

11568         Support of this option provides predictable scheduling behavior, allowing applications to
11569         determine the order in which threads that are ready to run are granted access to a processor.
11570         It should only be required if needed, but it should never be prohibited.

11571    _POSIX_THREAD_PROCESS_SHARED
11572         The system provides shared access among multiple processes to synchronization objects.

11573         This option was created to support historical systems that did not provide the feature. It
11574         should only be required if needed, but it should never be prohibited.

11575         XSI-conformant systems support this option.

11576    _POSIX_THREAD_SAFE_FUNCTIONS
11577         The system provides thread-safe versions of all of the POSIX.1 functions.

11578         This option is required if the Threads option is supported. This is a separate option because
11579         thread-safe functions are useful in implementations providing other mechanisms for
11580         concurrency. It should only be required if needed, but it should never be prohibited.

11581         XSI-conformant systems support this option.

11582    _POSIX_THREAD_SPORADIC_SERVER
11583         The system supports the thread sporadic server scheduling policy.

11584         Support for this option provides applications with a new scheduling policy for scheduling
11585         aperiodic threads in hard realtime applications.

11586    _POSIX_TIMEOUTS
11587         The system provides timeouts for some blocking services.

11588         This option was created to provide a timeout capability to system services, thus allowing
11589         applications to include better error detection, and recovery capabilities.

11590    _POSIX_TIMERS
11591         The system provides higher resolution clocks with multiple timers per process.

11592         This option was created to support historical systems that did not provide the features. This
11593         option is appropriate for applications requiring higher resolution timestamps or needing to
11594         control the timing of multiple activities. It should only be required if needed, but it should
11595         never be prohibited.

11596    _POSIX_TRACE
11597         The system supports the Trace option.

11598          This option was created to allow applications to perform tracing.

11599     _POSIX_TRACE_EVENT_FILTER
11600          The system supports the Trace Event Filter option.

11601          This option is dependent on support of the Trace option.

11602     _POSIX_TRACE_INHERIT
11603          The system supports the Trace Inherit option.

11604          This option is dependent on support of the Trace option.

11605     _POSIX_TRACE_LOG
11606          The system supports the Trace Log option.

11607          This option is dependent on support of the Trace option.

11608     _POSIX_TYPED_MEMORY_OBJECTS
11609          The system supports the Typed Memory Objects option.

11610          This option was created to allow realtime applications to access different kinds of physical
11611          memory, and allow processes in these applications to share portions of this memory.

## 11612 D.3.5  Configurable Limits

11613     In general, the configurable limits in the **<limits.h>** header defined in the Base Definitions
11614     volume of IEEE Std 1003.1-2001 have been set to minimal values; many applications or
11615     implementations may require larger values. No profile can cite lower values.

11616     {AIO_LISTIO_MAX}
11617          The current minimum is likely to be inadequate for most applications. It is expected that
11618          this value will be increased by profiles requiring support for list input and output
11619          operations.

11620     {AIO_MAX}
11621          The current minimum is likely to be inadequate for most applications. It is expected that
11622          this value will be increased by profiles requiring support for asynchronous input and
11623          output operations.

11624     {AIO_PRIO_DELTA_MAX}
11625          The functionality associated with this limit is needed only by sophisticated applications. It
11626          is not expected that this limit would need to be increased under a general-purpose profile.

11627     {ARG_MAX}
11628          The current minimum is likely to need to be increased for profiles, particularly as larger
11629          amounts of information are passed through the environment. Many implementations are
11630          believed to support larger values.

11631     {CHILD_MAX}
11632          The current minimum is suitable only for systems where a single user is not running
11633          applications in parallel. It is significantly too low for any system also requiring windows,
11634          and if _POSIX_JOB_CONTROL is specified, it should be raised.

11635     {CLOCKRES_MIN}
11636          It is expected that profiles will require a finer granularity clock, perhaps as fine as 1 $\mu$s,
11637          represented by a value of 1 000 for this limit.

11638     {DELAYTIMER_MAX}
11639          It is believed that most implementations will provide larger values.

11640          {LINK_MAX}
11641                  For most applications and usage, the current minimum is adequate. Many implementations
11642                  have a much larger value, but this should not be used as a basis for raising the value unless
11643                  the applications to be used require it.

11644          {LOGIN_NAME_MAX}
11645                  This is not actually a limit, but an implementation parameter. No profile should impose a
11646                  requirement on this value.

11647          {MAX_CANON}
11648                  For most purposes, the current minimum is adequate. Unless high-speed burst serial
11649                  devices are used, it should be left as is.

11650          {MAX_INPUT}
11651                  See {MAX_CANON}.

11652          {MQ_OPEN_MAX}
11653                  The current minimum should be adequate for most profiles.

11654          {MQ_PRIO_MAX}
11655                  The current minimum corresponds to the required number of process scheduling priorities.
11656                  Many realtime practitioners believe that the number of message priority levels ought to be
11657                  the same as the number of execution scheduling priorities.

11658          {NAME_MAX}
11659                  Many implementations now support larger values, and many applications and users
11660                  assume that larger names can be used. Many existing profiles also specify a larger value.
11661                  Specifying this value will reduce the number of conforming implementations, although this
11662                  might not be a significant consideration over time. Values greater than 255 should not be
11663                  required.

11664          {NGROUPS_MAX}
11665                  The value selected will typically be 8 or larger.

11666          {OPEN_MAX}
11667                  The historically common value for this has been 20. Many implementations support larger
11668                  values. If applications that use larger values are anticipated, an appropriate value should be
11669                  specified.

11670          {PAGESIZE}
11671                  This is not actually a limit, but an implementation parameter. No profile should impose a
11672                  requirement on this value.

11673          {PATH_MAX}
11674                  Historically, the minimum has been either 1 024 or indefinite, depending on the
11675                  implementation. Few applications actually require values larger than 256, but some users
11676                  may create file hierarchies that must be accessed with longer paths. This value should only
11677                  be changed if there is a clear requirement.

11678          {PIPE_BUF}
11679                  The current minimum is adequate for most applications. Historically, it has been larger. If
11680                  applications that write single transactions larger than this are anticipated, it should be
11681                  increased. Applications that write lines of text larger than this probably do not need it
11682                  increased, as the text line is delimited by a <newline>.

11683          {POSIX_VERSION}
11684                  This is actually not a limit, but a standard version stamp. Generally, a profile should specify
11685                  IEEE Std 1003.1-2001 by a name in the normative references section, not this value.

11686    {PTHREAD_DESTRUCTOR_ITERATIONS}
11687        It is unlikely that applications will need larger values to avoid loss of memory resources.

11688    {PTHREAD_KEYS_MAX}
11689        The current value should be adequate for most profiles.

11690    {PTHREAD_STACK_MIN}
11691        This should not be treated as an actual limit, but as an implementation parameter. No
11692        profile should impose a requirement on this value.

11693    {PTHREAD_THREADS_MAX}
11694        It is believed that most implementations will provide larger values.

11695    {RTSIG_MAX}
11696        The current limit was chosen so that the set of POSIX.1 signal numbers can fit within a 32-
11697        bit field. It is recognized that most existing implementations define many more signals than
11698        are specified in POSIX.1 and, in fact, many implementations have already exceeded 32
11699        signals (including the ''null signal''). Support of {_POSIX_RTSIG_MAX} additional signals
11700        may push some implementations over the single 32-bit word line, but is unlikely to push
11701        any implementations that are already over that line beyond the 64 signal line.

11702    {SEM_NSEMS_MAX}
11703        The current value should be adequate for most profiles.

11704    {SEM_VALUE_MAX}
11705        The current value should be adequate for most profiles.

11706    {SSIZE_MAX}
11707        This limit reflects fundamental hardware characteristics (the size of an integer), and should
11708        not be specified unless it is clearly required. Extreme care should be taken to assure that
11709        any value that might be specified does not unnecessarily eliminate implementations
11710        because of accidents of hardware design.

11711    {STREAM_MAX}
11712        This limit is very closely related to {OPEN_MAX}. It should never be larger than
11713        {OPEN_MAX}, but could reasonably be smaller for application areas where most files are
11714        not accessed through *stdio*. Some implementations may limit {STREAM_MAX} to 20 but
11715        allow {OPEN_MAX} to be considerably larger. Such implementations should be allowed for
11716        if the applications permit.

11717    {TIMER_MAX}
11718        The current limit should be adequate for most profiles, but it may need to be larger for
11719        applications with a large number of asynchronous operations.

11720    {TTY_NAME_MAX}
11721        This is not actually a limit, but an implementation parameter. No profile should impose a
11722        requirement on this value.

11723    {TZNAME_MAX}
11724        The minimum has been historically adequate, but if longer timezone names are anticipated
11725        (particularly such values as UTC−1), this should be increased.

11726 **D.3.6    Optional Behavior**

11727    In IEEE Std 1003.1-2001, there are no instances of the terms unspecified, undefined,
11728    implementation-defined, or with the verbs ''may'' or ''need not'', that the developers of
11729    IEEE Std 1003.1-2001 anticipate or sanction as suitable for profile or test method citation. All of
11730    these are merely warnings to conforming applications to avoid certain areas that can vary from
11731    system to system, and even over time on the same system. In many cases, these terms are used
11732    explicitly to support extensions, but profiles should not anticipate and require such extensions;
11733    future versions of IEEE Std 1003.1 may do so.

# Rationale (Informative)

11734

## Part E:

## Subprofiling Considerations

11735

11736

*The Open Group*
*The Institute of Electrical and Electronics Engineers, Inc.*

11737
11738

11739

*Appendix E*

# Subprofiling Considerations (Informative)

11741     This section contains further information to satisfy the requirement that the project scope enable
11742     subprofiling of IEEE Std 1003.1-2001. The original intent was to have included a set of options
11743     similar to the ''Units of Functionality'' contained in IEEE Std 1003.13-1998. However, as the
11744     development of IEEE Std 1003.1-2001 continued, the standard developers felt it premature to fix
11745     these in normative text. The approach instead has been to include a general requirement in
11746     normative text regarding subprofiling and to include an informative section (here) containing a
11747     proposed set of subprofiling options.

## 11748   E.1    Subprofiling Option Groups

11749     The following Option Groups[4] are defined to support profiling. Systems claiming support to
11750     IEEE Std 1003.1-2001 need not implement these options apart from the requirements stated in
11751     the Base Definitions volume of IEEE Std 1003.1-2001, Section 2.1.3, POSIX Conformance. These
11752     Option Groups allow profiles to subset the System Interfaces volume of IEEE Std 1003.1-2001 by
11753     collecting sets of related functions.

11754     POSIX_C_LANG_JUMP: Jump Functions
11755        *longjmp*(), *setjmp*()

11756     POSIX_C_LANG_MATH: Maths Library
11757        *acos*(), *acosf*(), *acosh*(), *acoshf*(), *acoshl*(), *acosl*(), *asin*(), *asinf*(), *asinh*(), *asinhf*(), *asinhl*(),
11758        *asinl*(), *atan*(), *atan2*(), *atan2f*(), *atan2l*(), *atanf*(), *atanh*(), *atanhf*(), *atanhl*(), *atanl*(), *cabs*(),
11759        *cabsf*(), *cabsl*(), *cacos*(), *cacosf*(), *cacosh*(), *cacoshf*(), *cacoshl*(), *cacosl*(), *carg*(), *cargf*(), *cargl*(),
11760        *casin*(), *casinf*(), *casinh*(), *casinhf*(), *casinhl*(), *casinl*(), *catan*(), *catanf*(), *catanh*(), *catanhf*(),
11761        *catanhl*(), *catanl*(), *cbrt*(), *cbrtf*(), *cbrtl*(), *ccos*(), *ccosf*(), *ccosh*(), *ccoshf*(), *ccoshl*(), *ccosl*(),
11762        *ceil*(), *ceilf*(), *ceill*(), *cexp*(), *cexpf*(), *cexpl*(), *cimag*(), *cimagf*(), *cimagl*(), *clog*(), *clogf*(), *clogl*(),
11763        *conj*(), *conjf*(), *conjl*(), *copysign*(), *copysignf*(), *copysignl*(), *cos*(), *cosf*(), *cosh*(), *coshf*(),
11764        *coshl*(), *cosl*(), *cpow*(), *cpowf*(), *cpowl*(), *cproj*(), *cprojf*(), *cprojl*(), *creal*(), *crealf*(), *creall*(),
11765        *csin*(), *csinf*(), *csinh*(), *csinhf*(), *csinhl*(), *csinl*(), *csqrt*(), *csqrtf*(), *csqrtl*(), *ctan*(), *ctanf*(),
11766        *ctanh*(), *ctanhf*(), *ctanhl*(), *ctanl*(), *erf*(), *erfc*(), *erfcf*(), *erfcl*(), *erff*(), *erfl*(), *exp*(), *exp2*(),
11767        *exp2f*(), *exp2l*(), *expf*(), *expl*(), *expm1*(), *expm1f*(), *expm1l*(), *fabs*(), *fabsf*(), *fabsl*(), *fdim*(),
11768        *fdimf*(), *fdiml*(), *floor*(), *floorf*(), *floorl*(), *fma*(), *fmaf*(), *fmal*(), *fmax*(), *fmaxf*(), *fmaxl*(), *fmin*(),
11769        *fminf*(), *fminl*(), *fmod*(), *fmodf*(), *fmodl*(), *fpclassify*(), *frexp*(), *frexpf*(), *frexpl*(), *hypot*(),
11770        *hypotf*(), *hypotl*(), *ilogb*(), *ilogbf*(), *ilogbl*(), *isfinite*(), *isgreater*(), *isgreaterequal*(), *isinf*(),
11771        *isless*(), *islessequal*(), *islessgreater*(), *isnan*(), *isnormal*(), *isunordered*(), *ldexp*(), *ldexpf*(),
11772        *ldexpl*(), *lgamma*(), *lgammaf*(), *lgammal*(), *llrint*(), *llrintf*(), *llrintl*(), *llround*(), *llroundf*(),
11773        *llroundl*(), *log*(), *log10*(), *log10f*(), *log10l*(), *log1p*(), *log1pf*(), *log1pl*(), *log2*(), *log2f*(), *log2l*(),
11774        *logb*(), *logbf*(), *logbl*(), *logf*(), *logl*(), *lrint*(), *lrintf*(), *lrintl*(), *lround*(), *lroundf*(), *lroundl*(),
11775        *modf*(), *modff*(), *modfl*(), *nan*(), *nanf*(), *nanl*(), *nearbyint*(), *nearbyintf*(), *nearbyintl*(),
11776        *nextafter*(), *nextafterf*(), *nextafterl*(), *nexttoward*(), *nexttowardf*(), *nexttowardl*(), *pow*(), *powf*(),
11777        *powl*(), *remainder*(), *remainderf*(), *remainderl*(), *remquo*(), *remquof*(), *remquol*(), *rint*(), *rintf*(),
11778        *rintl*(), *round*(), *roundf*(), *roundl*(), *scalbln*(), *scalblnf*(), *scalblnl*(), *scalbn*(), *scalbnf*(), *scalbnl*(),
11779        *signbit*(), *sin*(), *sinf*(), *sinh*(), *sinhf*(), *sinhl*(), *sinl*(), *sqrt*(), *sqrtf*(), *sqrtl*(), *tan*(), *tanf*(),

11780 _____

11781   4.   These are equivalent to the Units of Functionality from IEEE Std 1003.13-1998.

| | |
|---|---|
| 11782 | *tanh*( ), *tanhf*( ), *tanhl*( ), *tanl*( ), *tgamma*( ), *tgammaf*( ), *tgammal*( ), *trunc*( ), *truncf*( ), *truncl*( ) |

11783    POSIX_C_LANG_SUPPORT: General ISO C Library
11784    *abs*( ), *asctime*( ), *atof*( ), *atoi*( ), *atol*( ), *atoll*( ), *bsearch*( ), *calloc*( ), *ctime*( ), *difftime*( ), *div*( ),
11785    *feclearexcept*( ), *fegetenv*( ), *fegetexceptflag*( ), *fegetround*( ), *feholdexcept*( ), *feraiseexcept*( ),
11786    *fesetenv*( ), *fesetexceptflag*( ), *fesetround*( ), *fetestexcept*( ), *feupdateenv*( ), *free*( ), *gmtime*( ),
11787    *imaxabs*( ), *imaxdiv*( ), *isalnum*( ), *isalpha*( ), *isblank*( ), *iscntrl*( ), *isdigit*( ), *isgraph*( ), *islower*( ),
11788    *isprint*( ), *ispunct*( ), *isspace*( ), *isupper*( ), *isxdigit*( ), *labs*( ), *ldiv*( ), *llabs*( ), *lldiv*( ), *localeconv*( ),
11789    *localtime*( ), *malloc*( ), *memchr*( ), *memcmp*( ), *memcpy*( ), *memmove*( ), *memset*( ), *mktime*( ),
11790    *qsort*( ), *rand*( ), *realloc*( ), *setlocale*( ), *snprintf*( ), *sprintf*( ), *srand*( ), *sscanf*( ), *strcat*( ), *strchr*( ),
11791    *strcmp*( ), *strcoll*( ), *strcpy*( ), *strcspn*( ), *strerror*( ), *strftime*( ), *strlen*( ), *strncat*( ), *strncmp*( ),
11792    *strncpy*( ), *strpbrk*( ), *strrchr*( ), *strspn*( ), *strstr*( ), *strtod*( ), *strtof*( ), *strtoimax*( ), *strtok*( ), *strtol*( ),
11793    *strtold*( ), *strtoll*( ), *strtoul*( ), *strtoull*( ), *strtoumax*( ), *strxfrm*( ), *time*( ), *tolower*( ), *toupper*( ),
11794    *tzname*, *tzset*( ), *va_arg*( ), *va_copy*( ), *va_end*( ), *va_start*( ), *vsnprintf*( ), *vsprintf*( ), *vsscanf*( )

11795    POSIX_C_LANG_SUPPORT_R: Thread-Safe General ISO C Library
11796    *asctime_r*( ), *ctime_r*( ), *gmtime_r*( ), *localtime_r*( ), *rand_r*( ), *strerror_r*( ), *strtok_r*( )

11797    POSIX_C_LANG_WIDE_CHAR: Wide-Character ISO C Library
11798    *btowc*( ), *iswalnum*( ), *iswalpha*( ), *iswblank*( ), *iswcntrl*( ), *iswctype*( ), *iswdigit*( ), *iswgraph*( ),
11799    *iswlower*( ), *iswprint*( ), *iswpunct*( ), *iswspace*( ), *iswupper*( ), *iswxdigit*( ), *mblen*( ), *mbrlen*( ),
11800    *mbrtowc*( ), *mbsinit*( ), *mbsrtowcs*( ), *mbstowcs*( ), *mbtowc*( ), *swprintf*( ), *swscanf*( ), *towctrans*( ),
11801    *towlower*( ), *towupper*( ), *vswprintf*( ), *vswscanf*( ), *wcrtomb*( ), *wcscat*( ), *wcschr*( ), *wcscmp*( ),
11802    *wcscoll*( ), *wcscpy*( ), *wcscspn*( ), *wcsftime*( ), *wcslen*( ), *wcsncat*( ), *wcsncmp*( ), *wcsncpy*( ),
11803    *wcspbrk*( ), *wcsrchr*( ), *wcsrtombs*( ), *wcsspn*( ), *wcsstr*( ), *wcstod*( ), *wcstof*( ), *wcstoimax*( ),
11804    *wcstok*( ), *wcstol*( ), *wcstold*( ), *wcstoll*( ), *wcstombs*( ), *wcstoul*( ), *wcstoull*( ), *wcstoumax*( ),
11805    *wcsxfrm*( ), *wctob*( ), *wctomb*( ), *wctrans*( ), *wctype*( ), *wmemchr*( ), *wmemcmp*( ), *wmemcpy*( ),
11806    *wmemmove*( ), *wmemset*( )

11807    POSIX_C_LIB_EXT: General C Library Extension
11808    *fnmatch*( ), *getopt*( ), *optarg*, *opterr*, *optind*, *optopt*

11809    POSIX_DEVICE_IO: Device Input and Output
11810    *FD_CLR*( ), *FD_ISSET*( ), *FD_SET*( ), *FD_ZERO*( ), *clearerr*( ), *close*( ), *fclose*( ), *fdopen*( ), *feof*( ),
11811    *ferror*( ), *fflush*( ), *fgetc*( ), *fgets*( ), *fileno*( ), *fopen*( ), *fprintf*( ), *fputc*( ), *fputs*( ), *fread*( ), *freopen*( ),
11812    *fscanf*( ), *fwrite*( ), *getc*( ), *getchar*( ), *gets*( ), *open*( ), *perror*( ), *printf*( ), *pselect*( ), *putc*( ), *putchar*( ),
11813    *puts*( ), *read*( ), *scanf*( ), *select*( ), *setbuf*( ), *setvbuf*( ), *stderr*, *stdin*, *stdout*, *ungetc*( ), *vfprintf*( ),
11814    *vfscanf*( ), *vprintf*( ), *vscanf*( ), *write*( )

11815    POSIX_DEVICE_SPECIFIC: General Terminal
11816    *cfgetispeed*( ), *cfgetospeed*( ), *cfsetispeed*( ), *cfsetospeed*( ), *ctermid*( ), *isatty*( ), *tcdrain*( ), *tcflow*( ),
11817    *tcflush*( ), *tcgetattr*( ), *tcsendbreak*( ), *tcsetattr*( ), *ttyname*( )

11818    POSIX_DEVICE_SPECIFIC_R: Thread-Safe General Terminal
11819    *ttyname_r*( )

11820    POSIX_FD_MGMT: File Descriptor Management
11821    *dup*( ), *dup2*( ), *fcntl*( ), *fgetpos*( ), *fseek*( ), *fseeko*( ), *fsetpos*( ), *ftell*( ), *ftello*( ), *ftruncate*( ), *lseek*( ),
11822    *rewind*( )

11823    POSIX_FIFO: FIFO
11824    *mkfifo*( )

11825    POSIX_FILE_ATTRIBUTES: File Attributes
11826    *chmod*( ), *chown*( ), *fchmod*( ), *fchown*( ), *umask*( )

11827    POSIX_FILE_LOCKING: Thread-Safe Stdio Locking
11828    *flockfile*( ), *ftrylockfile*( ), *funlockfile*( ), *getc_unlocked*( ), *getchar_unlocked*( ), *putc_unlocked*( ),

11829          *putchar_unlocked*()

11830     POSIX_FILE_SYSTEM: File System
11831          *access*(), *chdir*(), *closedir*(), *creat*(), *fpathconf*(), *fstat*(), *getcwd*(), *link*(), *mkdir*(), *opendir*(),
11832          *pathconf*(), *readdir*(), *remove*(), *rename*(), *rewinddir*(), *rmdir*(), *stat*(), *tmpfile*(), *tmpnam*(),
11833          *unlink*(), *utime*()

11834     POSIX_FILE_SYSTEM_EXT: File System Extensions
11835          *glob*(), *globfree*()

11836     POSIX_FILE_SYSTEM_R: Thread-Safe File System
11837          *readdir_r*()

11838     POSIX_JOB_CONTROL: Job Control
11839          *setpgid*(), *tcgetpgrp*(), *tcsetpgrp*()

11840     POSIX_MULTI_PROCESS: Multiple Processes
11841          *_Exit*(), *_exit*(), *assert*(), *atexit*(), *clock*(), *execl*(), *execle*(), *execlp*(), *execv*(), *execve*(), *execvp*(),
11842          *exit*(), *fork*(), *getpgrp*(), *getpid*(), *getppid*(), *setsid*(), *sleep*(), *times*(), *wait*(), *waitpid*()

11843     POSIX_NETWORKING: Networking
11844          *accept*(), *bind*(), *connect*(), *endhostent*(), *endnetent*(), *endprotoent*(), *endservent*(),
11845          *freeaddrinfo*(), *gai_strerror*(), *getaddrinfo*(), *gethostbyaddr*(), *gethostbyname*(), *gethostent*(),
11846          *gethostname*(), *getnameinfo*(), *getnetbyaddr*(), *getnetbyname*(), *getnetent*(), *getpeername*(),
11847          *getprotobyname*(), *getprotobynumber*(), *getprotoent*(), *getservbyname*(), *getservbyport*(),
11848          *getservent*(), *getsockname*(), *getsockopt*(), *h_errno*, *htonl*(), *htons*(), *if_freenameindex*(),
11849          *if_indextoname*(), *if_nameindex*(), *if_nametoindex*(), *inet_addr*(), *inet_ntoa*(), *inet_ntop*(),
11850          *inet_pton*(), *listen*(), *ntohl*(), *ntohs*(), *recv*(), *recvfrom*(), *recvmsg*(), *send*(), *sendmsg*(), *sendto*(),
11851          *sethostent*(), *setnetent*(), *setprotoent*(), *setservent*(), *setsockopt*(), *shutdown*(), *socket*(),
11852          *sockatmark*(), *socketpair*()

11853     POSIX_PIPE: Pipe
11854          *pipe*()

11855     POSIX_REGEXP: Regular Expressions
11856          *regcomp*(), *regerror*(), *regexec*(), *regfree*()

11857     POSIX_SHELL_FUNC: Shell and Utilities
11858          *pclose*(), *popen*(), *system*(), *wordexp*(), *wordfree*()

11859     POSIX_SIGNALS: Signal
11860          *abort*(), *alarm*(), *kill*(), *pause*(), *raise*(), *sigaction*(), *sigaddset*(), *sigdelset*(), *sigemptyset*(),
11861          *sigfillset*(), *sigismember*(), *signal*(), *sigpending*(), *sigprocmask*(), *sigsuspend*(), *sigwait*()

11862     POSIX_SIGNAL_JUMP: Signal Jump Functions
11863          *siglongjmp*(), *sigsetjmp*()

11864     POSIX_SINGLE_PROCESS: Single Process
11865          *confstr*(), *environ*, *errno*, *getenv*(), *setenv*(), *sysconf*(), *uname*(), *unsetenv*()

11866     POSIX_SYMBOLIC_LINKS: Symbolic Links
11867          *lstat*(), *readlink*(), *symlink*()

11868     POSIX_SYSTEM_DATABASE: System Database
11869          *getgrgid*(), *getgrnam*(), *getpwnam*(), *getpwuid*()

11870     POSIX_SYSTEM_DATABASE_R: Thread-Safe System Database
11871          *getgrgid_r*(), *getgrnam_r*(), *getpwnam_r*(), *getpwuid_r*()

| | |
|---|---|
| 11872 | POSIX_USER_GROUPS: User and Group |
| 11873 | *getegid*( ), *geteuid*( ), *getgid*( ), *getgroups*( ), *getlogin*( ), *getuid*( ), *setegid*( ), *seteuid*( ), *setgid*( ), |
| 11874 | *setuid*( ) |
| 11875 | POSIX_USER_GROUPS_R: Thread-Safe User and Group |
| 11876 | *getlogin_r*( ) |
| 11877 | POSIX_WIDE_CHAR_DEVICE_IO: Device Input and Output |
| 11878 | *fgetwc*( ), *fgetws*( ), *fputwc*( ), *fputws*( ), *fwide*( ), *fwprintf*( ), *fwscanf*( ), *getwc*( ), *getwchar*( ), |
| 11879 | *putwc*( ), *putwchar*( ), *ungetwc*( ), *vfwprintf*( ), *vfwscanf*( ), *vwprintf*( ), *vwscanf*( ), *wprintf*( ), |
| 11880 | *wscanf*( ) |
| 11881 | XSI_C_LANG_SUPPORT: XSI General C Library |
| 11882 | *_tolower*( ), *_toupper*( ), *a64l*( ), *daylight*( ), *drand48*( ), *erand48*( ), *ffs*( ), *getcontext*( ), *getdate*( ), |
| 11883 | *getsubopt*( ), *hcreate*( ), *hdestroy*( ), *hsearch*( ), *iconv*( ), *iconv_close*( ), *iconv_open*( ), *initstate*( ), |
| 11884 | *insque*( ), *isascii*( ), *jrand48*( ), *l64a*( ), *lcong48*( ), *lfind*( ), *lrand48*( ), *lsearch*( ), *makecontext*( ), |
| 11885 | *memccpy*( ), *mrand48*( ), *nrand48*( ), *random*( ), *remque*( ), *seed48*( ), *setcontext*( ), *setstate*( ), |
| 11886 | *signgam*, *srand48*( ), *srandom*( ), *strcasecmp*( ), *strdup*( ), *strfmon*( ), *strncasecmp*( ), *strptime*( ), |
| 11887 | *swab*( ), *swapcontext*( ), *tdelete*( ), *tfind*( ), *timezone*( ), *toascii*( ), *tsearch*( ), *twalk*( ) |
| 11888 | XSI_DBM: XSI Database Management |
| 11889 | *dbm_clearerr*( ), *dbm_close*( ), *dbm_delete*( ), *dbm_error*( ), *dbm_fetch*( ), *dbm_firstkey*( ), |
| 11890 | *dbm_nextkey*( ), *dbm_open*( ), *dbm_store*( ) |
| 11891 | XSI_DEVICE_IO: XSI Device Input and Output |
| 11892 | *fmtmsg*( ), *poll*( ), *pread*( ), *pwrite*( ), *readv*( ), *writev*( ) |
| 11893 | XSI_DEVICE_SPECIFIC: XSI General Terminal |
| 11894 | *grantpt*( ), *posix_openpt*( ), *ptsname*( ), *unlockpt*( ) |
| 11895 | XSI_DYNAMIC_LINKING: XSI Dynamic Linking |
| 11896 | *dlclose*( ), *dlerror*( ), *dlopen*( ), *dlsym*( ) |
| 11897 | XSI_FD_MGMT: XSI File Descriptor Management |
| 11898 | *truncate*( ) |
| 11899 | XSI_FILE_SYSTEM: XSI File System |
| 11900 | *basename*( ), *dirname*( ), *fchdir*( ), *fstatvfs*( ), *ftw*( ), *lchown*( ), *lockf*( ), *mknod*( ), *mkstemp*( ), *nftw*( ), |
| 11901 | *realpath*( ), *seekdir*( ), *statvfs*( ), *sync*( ), *telldir*( ), *tempnam*( ) |
| 11902 | XSI_I18N: XSI Internationalization |
| 11903 | *catclose*( ), *catgets*( ), *catopen*( ), *nl_langinfo*( ) |
| 11904 | XSI_IPC: XSI Interprocess Communication |
| 11905 | *ftok*( ), *msgctl*( ), *msgget*( ), *msgrcv*( ), *msgsnd*( ), *semctl*( ), *semget*( ), *semop*( ), *shmat*( ), *shmctl*( ), |
| 11906 | *shmdt*( ), *shmget*( ) |
| 11907 | XSI_JOB_CONTROL: XSI Job Control |
| 11908 | *tcgetsid*( ) |
| 11909 | XSI_JUMP: XSI Jump Functions |
| 11910 | *_longjmp*( ), *_setjmp*( ) |
| 11911 | XSI_MATH: XSI Maths Library |
| 11912 | *j0*( ), *j1*( ), *jn*( ), *scalb*( ), *y0*( ), *y1*( ), *yn*( ) |
| 11913 | XSI_MULTI_PROCESS: XSI Multiple Process |
| 11914 | *getpgid*( ), *getpriority*( ), *getrlimit*( ), *getrusage*( ), *getsid*( ), *nice*( ), *setpgrp*( ), *setpriority*( ), |
| 11915 | *setrlimit*( ), *ulimit*( ), *usleep*( ), *vfork*( ), *waitid*( ) |

11916    XSI_SIGNALS: XSI Signal
11917    *bsd_signal*( ), *killpg*( ), *sigaltstack*( ), *sighold*( ), *sigignore*( ), *siginterrupt*( ), *sigpause*( ), *sigrelse*( ),
11918    *sigset*( ), *ualarm*( )

11919    XSI_SINGLE_PROCESS: XSI Single Process
11920    *gethostid*( ), *gettimeofday*( ), *putenv*( )

11921    XSI_SYSTEM_DATABASE: XSI System Database
11922    *endpwent*( ), *getpwent*( ), *setpwent*( )

11923    XSI_SYSTEM_LOGGING: XSI System Logging
11924    *closelog*( ), *openlog*( ), *setlogmask*( ), *syslog*( )

11925    XSI_THREAD_MUTEX_EXT: XSI Thread Mutex Extensions
11926    *pthread_mutexattr_gettype*( ), *pthread_mutexattr_settype*( )

11927    XSI_THREADS_EXT: XSI Threads Extensions
11928    *pthread_attr_getguardsize*( ), *pthread_attr_getstack*( ), *pthread_attr_setguardsize*( ),
11929    *pthread_attr_setstack*( ), *pthread_getconcurrency*( ), *pthread_setconcurrency*( )

11930    XSI_TIMERS: XSI Timers
11931    *getitimer*( ), *setitimer*( )

11932    XSI_USER_GROUPS: XSI User and Group
11933    *endgrent*( ), *endutxent*( ), *getgrent*( ), *getutxent*( ), *getutxid*( ), *getutxline*( ), *pututxline*( ),
11934    *setgrent*( ), *setregid*( ), *setreuid*( ), *setutxent*( )

11935    XSI_WIDE_CHAR: XSI Wide-Character Library
11936    *wcswidth*( ), *wcwidth*( )

# *Index*

# Index