

Throughput-Driven IC Communication Fabric Synthesis[†]

Tao Lin[‡]

Lawrence T. Pileggi

Department of Electrical and Computer Engineering, Carnegie Mellon University
5000 Forbes Avenue, Pittsburgh, PA, 15213
{tl, pileggi}@ece.cmu.edu

ABSTRACT

As the scale of system integration continues to grow, the on-chip communication becomes the ultimate bottleneck of system performance and the primary determinant of system architecture. In this paper we propose a *throughput-driven* synthesis methodology for on-chip communication fabrics based on optimized bus models. Compared with traditional *delay-driven*, *wire-by-wire* planning methods, the *throughput-driven* methodology provides a feasible and accurate system-level solution to address delay and congestion problems simultaneously during early-phase design planning. Unlike the conventional methods which are based on rather inaccurate RC models and simplistic delay metrics, in our methodology the communication fabrics are characterized in terms of realistic Partial Element Equivalent Circuits (PEEC) extracted from the multi-layer interconnects and transistor level transient analysis via SPICE-like tools. The characterized models facilitate a flexible interconnect fabric optimization engine that can be embedded into a system planner for *throughput-driven* synthesis. Furthermore, engineering trade-offs considering repeater area and interconnect power consumption are further considered as part of this methodology.

1. Motivation

On-chip communication is widely accepted as a key performance bottleneck for cutting-edge deep sub-micron (DSM) ICs ([1][2][3][4]). With the advent of the Giga-Scale Integration (GSI) era, it is foreseeable that the design of interconnects that can support the communication between a billion gates operating at multi-gigahertz frequency is going to be a daunting task. It follows that system level planning of the interconnect architecture is becoming more and more important ([15][16][17][18]).

Classical interconnect planning methodologies, however, are limited to very simplified models([13][14]), and do not properly consider the cost of the routing and area resources. Due to the inaccuracy in planning and prediction, time-consuming post-routing analyses are frequently required to identify and correct hidden interconnect problems. In many cases the entire design process has to be restarted in order to overcome problems that resulted as a consequence of bad planning.

Traditionally, the focus of interconnect planning has been on the interconnect delay problem. In order to achieve timing closure, global interconnects are often “reverse-scaled” to reduce the wiring delay due to interconnect resistance. Although timing can often be solved by this means, other problems can be created in the process.

Firstly, the risk of increasing routing congestion is greatly increased by the use of “fat” wires and large spacing on a growing number of global interconnects. It was projected in [5] that this technique alone would require an unrealistic number of metal layers as early as in 2005. Ultimately, the number of routing layers required will be nearly an order of magnitude larger than the number of layers prescribed by the International Technology Roadmap for Semiconductors (ITRS).

In contrast, the authors of [5] proposed that the congestion problem can be mitigated by the use of repeaters. Unfortunately, in order to achieve the reduction in number of layers, the projected repeater area is as much as 40% of the total area beyond 2005. Note that such a large amount of repeater area can completely change the overall floorplan or design plan.

Finally, but importantly, the power consumption of the global interconnect is becoming a more and more important concern. Wide wires and large repeaters will increase the power consumption. It was projected by the same study that the power dissipation by the repeaters could account for 20%-30% of the total power. Quite a few works proposed to reduce the bus power consumption by encoding techniques. However, the extra data bits required for encoding would in turn potentially increase the congestion.

In addition to the failure of considering the cost of available resources, wire planning is also based on very simplistic RC models. With the trend of rapidly shrinking feature size of CMOS and interconnect technology, more and more gates and wires are fabricated inside the same chip area than ever before. The parasitics among the massive amount of metal wires are beginning to play a significant role in circuit response. Not only is the delay of interconnects often found dominating that of the gates, but the actual delay also becomes extremely uncertain and hard to bound due to the strong crosstalk noises abundant in DSM chips.

Moreover, the simplistic RC delay models are unable to capture the emerging on-chip inductance effects [6][7][8]. With increasing signal frequency, the inductance reactance ($j\omega L$) is not necessarily a negligible part of the total interconnect impedance. Both signal delay and signal integrity can deteriorate with increasing inductance effects. One sub-problem is to determine just how wide wires should be in order to balance the trade-offs of increasing signal speed vs. increasing inductance.

The classic wire planning methodology ([13][14]) uses simple RC models obtained via empirical formulas for the wire and a constant linear resistor model for the transistors. The delay calculation is based on first order approximations, such as Elmore delay or Sakurai's expression. These models are of limited accuracy even when only RC delay is concerned, but are futile for RCL interconnects and crosstalk, delay uncertainty, and signal integrity problems.

Therefore, a system planner based on the conventional *delay-driven*, *wire-by-wire* planning paradigm can not guarantee the reliability or feasibility of the synthesized communication links. It also lacks the

[†] This work was supported in part by Semiconductor Research Corporation contract 2000-TJ-778 and a grant from Intel Corporation

[‡] Tao Lin was with Carnegie Mellon University, Dept. of Electrical and Computer Engineering. He is now with Monterey Design Systems, Sunnyvale, CA.

ability to achieve a satisfactory performance cost trade-off for the aforementioned scenario.

In this paper, we propose a new methodology for interconnect system planning that utilizes the most accurate models and analysis methods and simultaneously integrates the ability for an aggressive throughput optimization. Importantly, these models provide for the opportunity to make engineering trade-off decisions for area and power consumption as part of a floorplanning or design planning process.

The remainder of this paper is organized as follows. Section 2 provides the preliminaries for our approach, including the description of the fundamental bus style interconnect fabric and its performance models. Section 3 describes the primary throughput-driven optimization formulation based on the models in Section 2. The optimization formulation and its variations are cast in an application framework in Section 4. Section 5 and Section 6 discuss possible trade-off methodologies for repeater area and power consumption during synthesis, followed by our conclusions in Section 7.

2. Preliminaries

2.1 Bus Fabrics

Point-to-point or *shared* bus networks are commonly seen in communication architectures of SOCs. As the backbones of on-chip communication, busses usually consist of a number of long parallel wires in the same metal layer and/or in the neighboring orthogonal layer dedicated for global routing. The wire widths and spacings are often sized differently than the minimum sizes allowed by the technology in order to minimize the signal delay. Moreover, repeater and shield insertion techniques are often employed as well to further reduce delay and crosstalk noise.

We consider the busses to be uniform and periodical interconnect fabrics as shown in Fig. 1 and Fig. 2 with the following structure and regularity:

- The wire width and spacing are uniform for every wire. All the signal wires have the same width, denoted by the variable, W_{si} . The shield wires may have another width, W_{sh} . The spacing between two neighboring wires are identical, S_p .
- The repeater insertion is uniform. All the repeaters inserted along the signal wires, including the driver and the receiver, have the same gate size (equal driving strength), S_{gate} . The wire length between a repeater to the next, L_{seg} , is the same for all the segments.
- The shield wires are inserted periodically. We refer to the number of the signal wires between a pair of VDD/GND shields as the shielding period, denoted by the variable N .

Note that the group of parameters $(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N)$ uniquely defines a uniform and periodical interconnect fabric (pattern). Actual on-chip busses can be built via replication of this basic fabric.

2.2 Performance Models

In order to measure the quality of the bus fabrics, we define their performance and cost characteristics as functions of the fabric parameters. In this paper, we consider a scalable model containing the following set of characteristics for each fabric design:

- *Normalized Throughput, TH_N* . We define the throughput of a bus with unit length and unit total width as the normalized throughput of the bus fabric $(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N)$. For a bus of given length and width, its actual throughput has the following relation with the normalized throughput of the bus fabric with which it is constructed.

$$TH(W, L) = TH_N \times \frac{W}{L} \quad (1)$$

- *Signal Speed, SS* . Instead of delay, the average signal speed is used to denote the timing performance of a bus fabric. Due to the uniformity of the fabrics, the worst case flight time for signal to pass from one repeater to the next is a constant, WD_{seg} . The signal speed in a bus fabric is L_{seg}/WD_{seg} . The worst case signal delay of a bus with length L can be computed from the signal speed in the fabric.

$$Delay(L) = \frac{L}{SS} \quad (2)$$

- *Energy Consumption, EC* . This parameter represents a measure of power consumption for busses constructed with a given type of interconnect fabric. It is the energy consumed to transfer one bit active signal “0->1” over a unit distance by a bus of the given fabric. Assuming the switching behavior of each data line is independent and irrelevant to the activity on other lines of the bus, the actual power consumption of a bus of bit-width M and length L is estimated by,

$$P(M, L) = EC \times L \times M \times \alpha \times F \quad (3)$$

where M is the bit width, F is the operating frequency, and α is the switching (0->1) probability of the signal that can be computed via a stochastic analysis of the data on the bus.

Given a bus fabric’s design parameters $(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N)$ and its performance parameters (TH_N, SS, EC) , the performance and cost of any bus built with it can be easily computed.

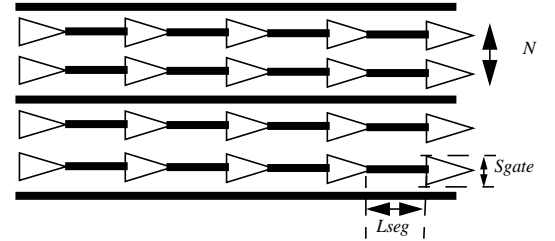


Figure 1. A uniform on-chip bus fabric.

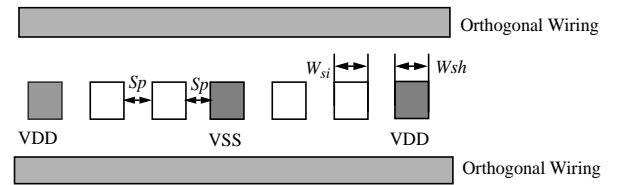


Figure 2. Cross section view of on-chip interconnect.

2.3 Model Generation

Generation of the scalable models involves the following two major steps.

1. Extraction: Input - $(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N)$, Output - (R, L, C) .

During this step, the layout of the bus fabric is generated from the set of design parameters $(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N)$. Each conductor between the repeaters is further discretized into several segments to capture the transmission line effect. Field solvers ([10][11]) are then applied to generate the PEEC models ([9]) in form of R, L, C matrices.

Since inductive coupling is a long range effect, the mutual coupling inductance between far away conductors can not be simply truncated without loss of accuracy and stability (as in the capacitive coupling case). It has been shown, however, that the forward coupling inductance among the segments of parallel bus structures can be either ignored or

reallocated uniformly between parallel segments [12]. Based on this observation, the inductance matrix can be sparsified into a block diagonal matrix via the combined 2D inductance models proposed in [12]. The sparsified inductance model greatly reduces the complexity of subsequent analyses while maintaining the required modeling accuracy and stability.

2. Simulation: Input - (R, L, C) , Output - (TH_N, SS, EC)

Instead of using a simplistic metric such as the Elmore delay, our models are characterized via transistor level transient simulation. The circuit models of the interconnect parasitics are combined with the SPICE models of the given CMOS technology. The combined circuits are then simulated with a set of input patterns to uncover the worst case behaviors.

This differs from traditional RC based analyses in two respects. First, inductive coupling exists beyond immediate neighbors. For RC wires, the capacitive couplings between the victim and the wires beyond the neighbors are negligible due to the electrostatic shielding effect. However, with inductive coupling, the faraway aggressors can still contribute a significant, albeit collective, impact on the victim. Second, the opposite switching case is no longer guaranteed to be the worst case. Instead, aggressors switching in the same direction can aggravate the inductive effect on the victim by superposition. It can be observed that in extreme cases, all the aggressors switching in the same direction can cause a worse delay than the opposite switching case.

In our implementation, the effective crosstalk window is determined by greedy heuristics. An empirical set of switching patterns, including the opposite switching and same direction switching, as well as the “W” switching pattern (i.e., the immediate neighbors switch in the opposite direction while the other aggressors switch in the same direction with the victim), are applied in circuit simulation.

With worst case signal delay, noise peak, and total current measured from simulation output. The performance parameters (TH_N, SS, EC) are computed via the relations in (1), (2) and (3).

2.4 Interconnect Fabric Optimization

Based on the performance parameters (TH_N, SS, EC) , for a given technology and specific application requirements, the bus fabric is optimized to achieve designated features with minimum cost. A standard optimization package is used to optimize our carefully formulated optimal performance functions in terms of the set of design variables $(W_{si}, W_{shr}, S_p, S_{gate}, L_{seg}, N)$ via either of the following multi-variable constrained nonlinear programs:

$$\begin{aligned} & \text{maximize } Performance_Func(W_{si}, W_{shr}, S_p, S_{gate}, L_{seg}, N) \\ & \text{subject to } Constraint_Func(W_{si}, W_{shr}, S_p, S_{gate}, L_{seg}, N) \geq 0 \quad (4) \\ & \text{lower and upper bounds on } (W_{si}, W_{shr}, S_p, S_{gate}, L_{seg}, N) \end{aligned}$$

or

$$\begin{aligned} & \text{minimize } Cost_Func(W_{si}, W_{shr}, S_p, S_{gate}, L_{seg}, N) \\ & \text{subject to } Constraint_Func(W_{si}, W_{shr}, S_p, S_{gate}, L_{seg}, N) \geq 0 \quad (5) \\ & \text{and lower and upper bounds} \end{aligned}$$

Where the performance, cost, and constraints are combinations or simple functions of the performance characteristics (TH_N, SS, EC) defined in Section 2.2.

The above optimization problems are solved by the following algorithm:

1. Solve the subproblem of $(W_{si}, W_{shr}, S_p, S_{gate}, L_{seg})$ using a Sequential Quadratic Programming (SQP) solver.

2. Downhill transverse the possible values of N and repeater step 1 at each N . The best of the solutions to the subproblems is the optimal solution to the original problem.

The overall flow of model generation and fabric optimization is summarized in Fig. 3.

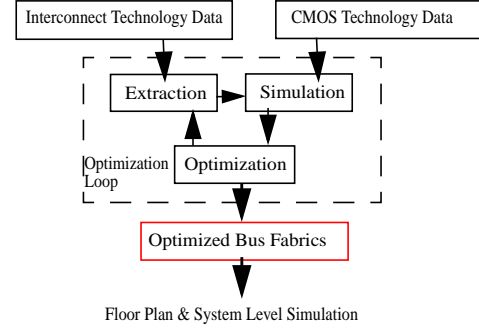


Figure 3. Communication Fabric Optimization.

3. Throughput-Driven Optimization

As discussed in the previous sections, the interconnect delay problem and the routing congestion problem often render the traditional delay-driven, wire-by-wire planner ineffective. In this section we show that both problems can be modeled and solved via a single metric, namely, throughput.

Throughput is the common performance measure for both processing units and communication links. It is the ultimate index of the overall system performance. For communication links, throughput (a.k.a. bandwidth) is the amount of data that can be transmitted through the bus over a fixed time period. The throughput of the communication links must be as much as the slowest of the associated processing unit in order not to drag down the system performance. Similar to the timing closure, the closure of throughput is a more fundamental requirement of high performance designs.

The throughput of a bus is determined by two factors, the operating frequency and the bit-width.

$$Throughput = Frequency \times Bitwidth \quad (6)$$

The frequency is related to the signal delay by the following equation.

$$MaximalFrequency = \frac{1}{\beta \times Worst\ Case\ Delay} \quad (7)$$

where β is a positive constant. β is larger than 1 to guarantee the signal can be correctly locked in the receiver in one clock period.

Traditional delay-driven wire planning methodology improves the communication throughput solely by reducing the signal delay, i.e., increasing the maximal operating frequency. The bit-width of on-chip busses is left to be determined by the I/O width of IP blocks they are connecting to. However, according to today's interconnect/communication centric SOC design methodology, the restriction of the bit-width of the busses is no longer as important. The bit-width of the busses can be varied to further increase the performance or lower the power consumption with little interfacing overhead ([15][16]). Arbitrating, encoding and decoding, and interfacing circuits may be required for bus controlling and data caching in this case. Nevertheless, the overall performance of the system can be boosted due to the fact that the effective throughput of the busses is improved.

For example, as shown in Fig. 4, for the fixed channel width shown, design (a) contains wider wires with larger spacing than design (b). It is possible that the signal delay of design (a) is smaller than (b), thus design (a) could function at a high frequency than (b). However, by shrinking the wire size and spacing, design (b) is able to route more wires in the given width. It is important to note that with a larger bit

width, the “slower” design (b) could achieve higher throughput than the “faster” design (a).

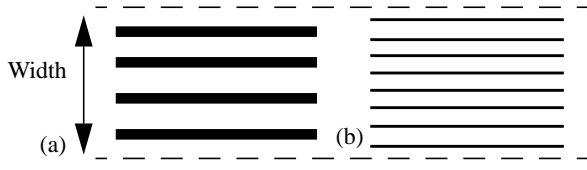


Figure 4. Bus designs with different bit-width.

According to (1), the optimal bus design that allows the maximal throughput in a given channel width must be constructed with the interconnect fabric that has the maximum normalized throughput.

$$TH_N = TH \times \frac{L}{W} = \frac{L}{\beta \times \text{Worst Case Delay}} \times \frac{\text{Bit-Width}}{W}, \quad (8)$$

The optimal fabric can be synthesized via solving the following formulation using the fabric optimization engine described in the previous section:

$$\begin{aligned} & \text{maximize } TH_N(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N) \\ & \text{subject to:} \\ & \text{lower bounds and upper bounds} \end{aligned} \quad (9)$$

Similarly, according to (2), the bus design with minimum signal delay must be constructed with the bus fabric that allows the maximal signal speed.

$$SS = L \times \text{Frequency} \times \beta = \frac{L}{\text{Worst Case Delay}}. \quad (10)$$

For comparison, we show the corresponding optimization formulation below for delay-driven (speed-driven) planning.

$$\begin{aligned} & \text{maximize } SS(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N) \\ & \text{subject to:} \\ & \text{lower bounds and upper bounds} \end{aligned} \quad (11)$$

A 0.18 μm ASIC technology is used here to demonstrate these performance measures. Suppose the global interconnects are routed using metal 5. The corresponding technology parameters are listed in Table 1.

Table 1: Metal 5 technology parameters

thickness	min width	min spacing	material	k
0.4 μm	0.3 μm	0.3 μm	Al	3.9

The maximal normalized¹ throughput fabric and the maximal signal speed fabric are generated respectively using the optimization engine described in the previous section. The minimum feature sizes are applied as lower bound constraints, and an upper bound constraint of 4 μm is added for the wire size and spacing to represent a reasonable design area. An additional constraint is applied to strictly limit the crosstalk noise to below 0.2volts. This tight constraint caused both solutions to be *fully shielded*. The maximal signal speed fabric solution corresponds to the widest wires and largest spacings allowed by the upper bound. The performance of both solutions are compared in Table 4. Although Fabric II is much faster than solution I, Fabric I produces a significantly larger throughput within an equivalent bus width.

Table 2: Maximal normalized throughput solution - I

W_{si}	W_{sh}	S_p	L_{seg}	S_{gate}^a	N
0.66 μm	0.3 μm	0.3 μm	1561 μm	29x	1

a. 1x is a basic gate size with NMOS W/L=4

Table 3: Maximal signal speed solution - II

W_{si}	W_{sh}	S_p	L_{seg}	S_{gate}^a	N
4 μm	4 μm	4 μm	4010 μm	47x	1

Table 4: Performance Comparison

Fabrics	SS (m/s)	TH_N (bps)
Solution I	1.34e7	8.57e12
Solution II	4.71e7	2.94e12

Given a throughput requirement for a bus of certain length, it can be determined by (1) that by using the maximal normalized throughput solution, the bus width is minimized. Therefore, the throughput-driven solution naturally reduces the congestion problem while maintaining the performance requirement. During early design planning, this would also determine the number of I/O ports required for the IP blocks to which these interconnect fabrics are connected.

4. Communication Fabric Synthesis

The fabric optimization formulations can be easily integrated into a floorplanner or design planner to provide for synthesis and prediction of on-chip communication channels. In this section, we demonstrated the application of these models for three design scenarios.

4.1 Globally-Asynchronous Locally Synchronous

In the first scenario, we consider the application of the *throughput-driven* communication fabric synthesis for Globally-Asynchronous Locally-Synchronous (GALS) designs[19][20]. A GALS design methodology allows the IP cores to work at different local clock frequencies. The communication links between the IPs operate asynchronously or at independent clock frequencies. In addition to the potential benefits of lower power consumption and higher system performance, GALS brings a new degree of freedom to the system level interconnect design by removing the hard timing constraints on the busses, since the operating frequency of each bus can be tuned solely by the interconnect delay.

Since the interconnect delay is not limited by any global clock cycle, the signal speed requirement on the bus fabrics does not exist. Therefore, the unconstrained *throughput-driven* optimization solution given by the previous section is always a feasible golden fabric for all the busses. Every bus should simply be implemented with the same optimal fabric. For example, using the ASIC technology in Section 3, the bus fabrics should all be configured following the design parameters in Table 2.

During floor planning of a GALS design, the bus width of a link connecting two IPs with distance L and throughput requirement TH_{req} can be easily computed as follows:

$$W = \frac{TH_{req}}{TH_N^*} \times L \quad (12)$$

$$\text{BitWidth} = \frac{W}{(N^* \times W_{si}^* + W_{sh}^* + (N^* + 1) \times S_p^*) / N^*}$$

where $(W_{si}^*, W_{sh}^*, S_p^*, S_{gate}^*, L_{seg}^*, N^*)$ is the optimal solution of the throughput driven optimization in Section 3.

For example, given the technology used in the previous section and two IPs located 2cms apart that require 64Gbps communication in-

1. Note that the normalized throughput values are several order of magnitude larger than the bus throughput we commonly see. This is because although the *normalized throughput* has the same units as throughput, it is not the actual throughput of any realistic bus. Numerically it is equivalent to the throughput of an imaginary bus with equal width and length.

between, it is straightforward to determine from Table 2 and equation (12) that a channel width of $150\mu m$ will be required.

4.2 Synchronous Communication

In the second scenario, everything on the chip operates synchronously at a common clock frequency. The bus signals are required to arrive at the destination in one clock cycle (or a given number of cycles for pipelined signals). In this case, specific timing constraints are applied to the communication channels, given the bus length, L_{latch} , between the latch repeaters. These length constraints are translated into the requirement on the signal speed of the bus fabric by (10). The optimal bus fabric for these designs can be synthesized from a constrained optimization problem as below:

$$\begin{aligned} & \text{maximize } TH_N(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N) \\ & \text{s.t. } SS(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N) \geq L_{latch} \times F \times \beta \quad (13) \\ & \text{and lower and upper bounds} \end{aligned}$$

For the same bus used in Section 4.1, assume that it is unpipelined and required to operate at 1Ghz. The signal speed is now constrained to be larger than $2e7m/s$ (assume $\beta=1$). Solving the above optimization problem we find that the optimal fabric under this condition has a normalized throughput of $8e12bps$.

Similar to (12), the required channel width can be easily computed as a function of the bus length, given the required throughput of the communication path and the normalized throughput of the fabric.

$$W = \frac{TH_{req}}{TH_N} \times L. \quad (14)$$

Here, TH_N is the normalized throughput of the solution to (13). Thus the channel width required to achieve 64Gbps throughput is $160\mu m$.

4.3 Throughput VS Latency

Similar to the synchronous case in the previous section, our third scenario considers the case of a global clock, however, instead of having a fixed number of latch repeaters in the bus, the planner is allowed to choose to the number of clock cycles a signal should take to reach the destination.

One simple approach for this synthesis would be to simply use the maximal normalized throughput fabric for all the busses, and insert the latch repeaters whenever necessary so that the timing constraints are satisfied. This becomes exactly the same situation as Section 4.1, such that the latch repeater distance is determined by

$$L_{latch} = \frac{SS^*}{Frequency \times \beta} \quad (15)$$

where SS^* is the signal speed in the maximal normalized throughput fabric.

However, the more latch repeaters that are introduced, the larger the latency of the bus. Clearly, this is also part of the architectural and early design planning problem. For real time applications, where the response time is of concern, the number of latch repeaters should be limited. Since the number of latch repeaters required is inversely proportional to the distance the signal can propagate in a clock cycle, the objective in our synthesis is to find an interconnect fabric that provides acceptable trade-off between signal speed and throughput. We generate such an optimal bus fabric via the following process:

1. Start with the maximal speed solution, determine the latency and throughput
2. Gradually relax the signal speed constraint in (13) and solve the problem repeatedly

3. Choose an acceptable point before the throughput stops increasing

We generated such a speed vs. normalized throughput trade-off curve for the above process. Given a clock frequency, one can easily convert this Throughput vs. Signal Speed into a Throughput vs. Latency plot.

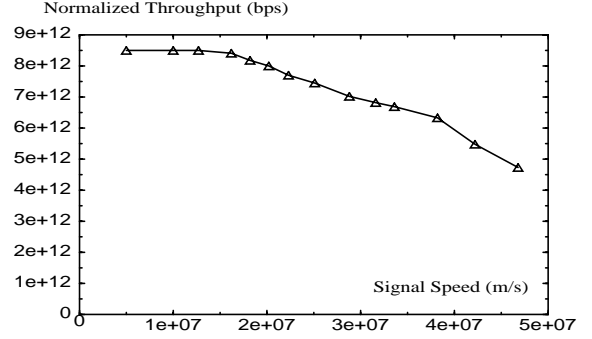


Figure 5. Throughput and Speed Trade-off.

5. Repeater Resource Planning

As shown in previous sections, throughput driven communication fabric synthesis and its variations facilitate the simultaneously consideration of interconnect delay and wiring resources during design planning. However, other resources, such as repeater area, are not reflected in our metric of *normalized-throughput*. From Table 2 and Table 3 we observe that the repeater sizes of the optimal solution can be an order of magnitude larger than the common standard cell sizes. According to the projection in [5], repeater usage will take a significant portion of the total chip area in future technology generations in the very near future. Therefore, it is important for floorplanners and design planners to have the ability to select an optimal, yet feasible, bus design based on the area and location of the repeater blocks as part of that early design planning process.

The uniform and periodical bus fabric make it relatively easy to realize such a function. With the model parameters L_{seg} and S_g of a fabric, the ultimate repeater usage can be computed by bus length and bit-width:

$$RepeaterArea = A \times S_g \times bitwidth \times \frac{L}{L_{seg}} \quad (16)$$

Approaches similar to those applied in throughput-latency trade-off in the previous section can be applied to select an optimal interconnect fabric with controlled repeater usage. However, perhaps more importantly, the floorplanner should be able to allocate the desired chip area for repeaters and include the repeater blocks as part of the complete floorplan using this model.

6. Low Power Communication Synthesis

We further extend the interconnect fabric synthesis methodology to include power considerations. Following the discussions of design trade-offs with respect to throughput-latency and throughput-repeater area, the solution to power aware interconnect fabric synthesis follows straightforwardly:

1. Start with the maximal normalized throughput solution, evaluate the normalized throughput and energy efficiency (EC).
2. Solve the following optimization formulation repeatedly and gradually relax the throughput constraint TH_x :

$$\begin{aligned}
& \text{minimize } EC(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N) \\
& \text{s.t. } TH_N(W_{si}, W_{sh}, S_p, S_{gate}, L_{seg}, N) \geq TH_x \quad (17) \\
& \text{other constraints and bounds}
\end{aligned}$$

3. Choose an acceptable point before the energy consumption stops decreasing

The resulting energy-throughput trade-off curve is shown in Fig. 6 for the technology parameters of Section 3. Note that the energy it takes to transfer one switching bit over a unit distance drops sharply at the initial stage of the curve when approximately only 10% normalized throughput is sacrificed. After that, the energy consumption gradually flattens out when more throughput is yielded. It follows that an appropriate design trade-off point might be one which would save as much as 45% of the power while only sacrificing 10% of the throughput.

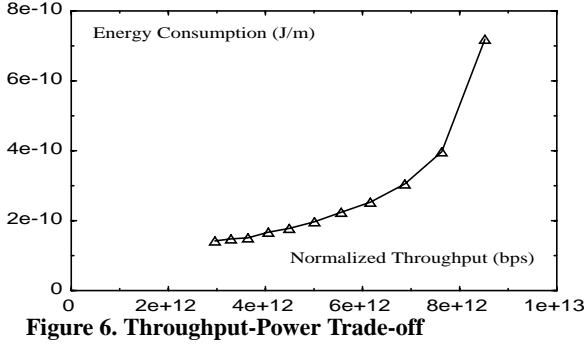


Figure 6. Throughput-Power Trade-off

7. Conclusions and Future Works

In this paper we propose a *throughput-driven* on-chip communication fabric synthesis methodology for system level interconnect and communication planning of SOCs. Unlike the classic *wire planning* paradigm, the *throughput-driven* methodology provides a feasible and accurate system-level solution to the interconnect bottleneck. In order to accurately capture the pervasive deep submicron effects, the communication fabrics in the new methodology are modeled by realistic Partial Element Equivalent Circuits (PEEC) extracted from the multi-tier interconnects via field solvers and simulated using transistor level simulators. A flexible optimization engine is used to generate optimal interconnect fabric designs for various design scenarios. Engineering trade-offs between throughput, latency, repeater area, and power consumption can be achieved conveniently via the proposed communication fabric synthesis procedure.

REFERENCES

[1] M. T. Bohr, "Interconnect Scaling-the real limiter to high performance ULSI", *Proc. IEDM*, pp. 241-244, 1995

[2] D. Sylvester and K. Keutzer, "Getting to the bottom of deep sub-micron", *Proc. ICCAD*, pp. 203-211, Nov. 1998

[3] R. Otten and Paul Stravers, "Challenges in physical chip design", *Proc. ICCAD*, pp. 84-91, Nov. 2000

[4] J. A. Davis and J. D. Meindl, "Is interconnect the weak link?", *Circuit and Device Magazine*, pp. 30-36, March 1998

[5] J. A. Davis, et. al., "Gigascale Integration (GSI) Interconnect Limits and N-Tier Multilevel Interconnect Architectural Solutions", *Proc. SLIP*, pp. 147-148 April 2001

[6] P. J. Restle, et. al., "Measurement and Modeling of On-Chip Transmission-Line Effects in a 400 Mhz Microprocessor", *IEEE Journal of Solid-State Circuits*, Vol. 33 No. 4, pp. 662-665, Apr. 1998

[7] A. Deutsch, et. al., "When are transmission-line effects important for on chip interconnect?", *IEEE Trans. Microwave Theory* Vol. 45, No. 10, pp. 1836-1846, Oct. 1997

[8] S. Morton, "On chip inductance issues in multiconductor systems", *Proc. DAC*, pp. 921-926, June 1999

[9] A. Ruehli, "Equivalent Circuit Models for Three-Dimensional Multiconductor Systems", *IEEE Trans. Microwave Theory and Techniques*, MTT-22, No. 3, pp. 216-221, March 1974.

[10] K. Nabors and J. White, "FastCap: A Multipole Accelerated 3D Capacitance Extraction Program", *IEEE Trans. CAD*, 10, pp. 1447-1459, November 1991.

[11] M. Kamon, M. Tsuk, and J. White, "FastHenry: A Multipole Accelerated 3D Inductance Extraction Program", *IEEE Trans. Microwave Theory and Techniques*, 42, pp. 1750-1758, September 1994.

[12] T. Lin and L. T. Pileggi, "On the efficacy of simplified 2D on-chip inductance models", *Proc. DAC*, June 2002

[13] H. B. Bakoglu, *Circuit, Interconnections, and Packaging for VLSI*, Addison-Wesley, Pub. Co. 1988

[14] R. Otten and R. K. Brayton, "Planning for performance", *Proc. DAC*, pp. 122-127, June 1998

[15] R. Ho, K. W. Mai, M. Horowitz, "The future of wires", *Proceedings of the IEEE*, vol. 89, pp. 490-504, April 2001

[16] W. Dally, "Interconnect-Limited VLSI Architecture", *Proceedings of the International Interconnect Technology Conference*, pp.15-17, May 1999

[17] M. Drinic, "Latency-guided on-chip bus network design", *Proc. ICCAD*, pp.420-423, Nov. 2000

[18] S. Meguerdichian, et. al., "Latency-driven design of multi-purpose system-on-chip", *Proc. DAC*, pp. 27-30, June 2001

[19] A. Hemani, et. al., "Lowering power consumption in clock by using globally asynchronous locally synchronous design style", *IEEE/ACM DAC*, pp. 873-878, 1999

[20] S. W. Moore, et. al., "Self calibrating clocks for globally asynchronous locally synchronous systems" *Proc. ICCD*, pp. 73-78, 2000