

Integrated Analysis of Power and Performance for Pipelined Microprocessors

Victor Zyuban, *Member, IEEE*, David Brooks, *Member, IEEE*, Viji Srinivasan, *Member, IEEE*, Michael Gschwind, *Senior Member, IEEE*, Pradip Bose, *Senior Member, IEEE*, Philip N. Strenski, *Member, IEEE*, and Philip G. Emma, *Fellow, IEEE*

Abstract—Choosing the pipeline depth of a microprocessor is one of the most critical design decisions that an architect must make in the concept phase of a microprocessor design. To be successful in today's cost/performance marketplace, modern CPU designs must effectively balance both performance and power dissipation. The choice of pipeline depth and target clock frequency has a critical impact on both of these metrics. In this paper, we describe an optimization methodology based on both analytical models and detailed simulations for power and performance as a function of pipeline depth. Our results for a set of SPEC2000 applications show that, when both power and performance are considered for optimization, the optimal clock period is around 18 FO4. We also provide a detailed sensitivity analysis of the optimal pipeline depth against key assumptions of our energy models. Finally, we discuss the potential risks in design quality for overly aggressive or conservative choices of pipeline depth.

Index Terms—Low-power design, energy-aware systems, pipeline processors, performance analysis and design aids, microprocessors and microcomputers.

1 INTRODUCTION

CURRENT generation high-end, server-class processors are performance-driven designs. These chips are still somewhat below the power and power density limits afforded by the package/cooling solution of choice in server markets targeted by such processors. In the design of future processors, however, energy efficiency is known to have become one of the primary design constraints [1], [2].

In this paper, we analyze key constraints in choosing the "optimal" pipeline depth (which directly influences the frequency target) of a microprocessor. The choice of pipeline depth is one of the fundamental issues confronting the architect/designer during the early microarchitecture definition phase of high performance, power-efficient processors. Even from a performance-only viewpoint, this issue has been important, if only to understand the limits to which pipelining can scale in the context of real workloads [3], [4], [5], [6], [7]. In certain segments of the market (typically desktops and low-end servers), there is often a market-driven tendency to equate delivered end performance with the frequency of the processor. Enlightened customers do understand the value of net system performance; nonetheless, the instinctive urge of going primarily for the highest frequency processor in a given technology generation is a known weakness even among savvy end users and, therefore, processor design teams. Recent studies [5], [6], [7] seem to suggest that there is still room to increase the pipeline depth with performance optima in the range of

8-11 FO4 inverter delays¹ per stage (consisting of 6-8 FO4 logic delay and 2-3 FO4 latch delay) for current out-of-order superscalar design paradigms. However, even in these performance-centric analysis papers, the authors do point out the practical difficulties of design complexity, verification, and power that must be solved in attaining these idealized limits. Our goal in this paper is to examine the practical, achievable limits when power dissipation constraints are also factored in. We believe that such analysis is needed to realistically bound the scalability limit of future technology generations. In particular, power dissipation must be carefully minimized to avoid design points which on paper promise ever higher performance, yet, under normal operating conditions, with commodity packaging and air cooling, only deliver a fraction of the theoretical peak performance.

In this paper, we first develop an analytical model to understand the power and performance trade offs for superscalar pipelines. From this model, we derive the optimal pipeline depth as a function of both power and performance. Subsequently, the results are validated and further refined using a detailed cycle-accurate simulator of a current generation superscalar processor. The energy model for the core pipeline is based on circuit-extracted power analysis for structures in a current, high-performance PowerPC processor. We then derive a methodology for scaling these energy models to deeper and shallower pipelines. With these performance and power models, we attempt to determine the optimal pipeline depth for a particular power-performance metric. Our results, based on an analysis of the TPC-C transaction processing benchmark

- V. Zyuban, V. Srinivasan, M. Gschwind, P. Bose, P.N. Strenski, and P.G. Emma are with the IBM T.J. Watson Research Center, Yorktown Heights, NY 10598. E-mail: zyuban@us.ibm.com.
- D. Brooks is with Harvard University, Cambridge, MA 02138.

Manuscript received 1 May 2003; revised 16 Dec. 2003; accepted 6 Jan. 2004. For information on obtaining reprints of this article, please send e-mail to: tc@computer.org, and reference IEEECS Log Number TC-0030-0503.

1. Fan-out-of-four (FO4) delay is defined as the delay of one inverter driving four copies of an equally sized inverter. The amount of logic and latch overhead per pipeline stage is often measured in terms of FO4 delay which implies that deeper pipelines have smaller FO4.

and a large set of SPEC2000 programs, indicate that a power-performance optimum is achieved at much shallower pipeline depths than a purely performance-focused evaluation would suggest.

The contributions of this paper are 1) energy models for both dynamic and leakage power that capture the scaling of different power components as a function of pipeline depth; 2) an analytical performance model that can predict the optimal pipeline depth; 3) cycle-accurate, detailed power-performance simulation with a thorough sensitivity analysis of the optimal pipeline depth against key energy model parameters.

This paper is structured as follows: We discuss the prior, related work in Section 2. In Section 3, we describe the proposed analytical model of pipeline depth effects on the performance of superscalar processors. Section 4 describes the experimental setup and derives scaling equations for various components of power with the pipeline depth. Section 5 presents optimum pipeline depth results. In Section 6, we present a detailed sensitivity analysis to understand the effect of variations in key parameters of the derived power models on the optimal pipeline depth. In Section 7, we discuss workload variability across different design points and its relationship to workload properties and, in Section 8, we analyze the implications of our findings on design decisions. We conclude in Section 9 with pointers to future work.

2 RELATED WORK

Previous work has studied the issue of “optimal” pipeline depth exclusively under the constraint of maximizing the performance delivered by the microprocessor. An initial study of optimal pipeline depths was performed by Kunkel and Smith in the context of supercomputers [3]. The authors studied the achievable performance from scalar and vector codes as a function of the number of gate levels per pipeline stage for the Livermore kernels. Subsequently, Dubey and Flynn [4] revisited the topic of optimal pipelining in a more general analytical framework. The authors showed the impact of various workload and design parameters.

More recently, several authors have reexamined this topic in the context of modern superscalar processor microarchitectures. Hartstein and Puzak [5] treat this problem analytically and verify based on detailed simulation of a variety of benchmarks for a 4-issue out-of-order machine with a memory-execute pipeline. They report the optimal logic delay per pipeline stage to be 7.7 FO4 for SPEC2000 and 5.5 FO4 for traditional and Java/C++ workloads. Assuming a latch insertion delay of 3 FO4, this would result in a total delay of 10.7 FO4 and 8.5 FO4 per pipeline stage, respectively.

Hrishikesh et al. [6] treat the question of logic depth per pipeline stage empirically based on simulation of the SPEC2000 benchmarks for an Alpha 21264-like machine. Based on their assumed latch insertion delay of 1.8 FO4, they demonstrate that a performance-optimal point is at logic delay of 6.0 FO4. This would result in a total delay of about 8 FO4 per stage.

Sprangle and Carmean [7] extrapolate from the current performance of the Pentium 4 the IPC degradation factors for adding a cycle to critical processor loops, such as ALU, L1 and L2 cache latencies, and branch miss penalty for a

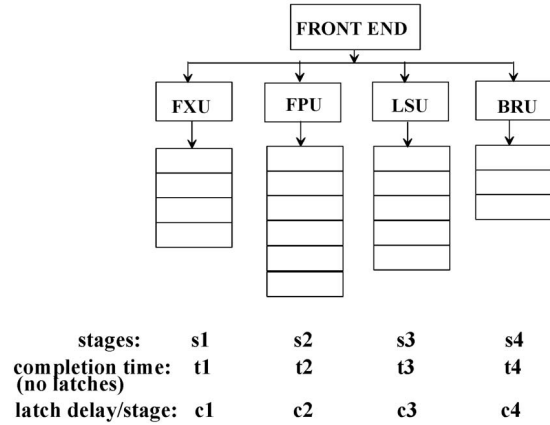


Fig. 1. Pipeline model.

variety of application types. The authors compute an optimal branch misprediction pipeline depth of 52 stages, corresponding to a pipeline stage total delay of 9.9 FO4 (based on a logic depth of 6.3 FO4 and a latch insertion delay of 3.6, of which 3 FO4 are due to latch delay and 0.6 FO4 represent skew and jitter overhead).

All of the above studies (as well as ours) assume that microarchitectural structures can be pipelined without limitation. Collectively, the cited works on optimal pipelining have made a significant contribution to the understanding of workloads and their interaction with pipeline structures by studying the theoretical limits of deep pipelining. However, prior work does not address scalability with respect to the increased *power dissipation* that is associated with deeper pipelines. In this work, we aim to build on this foundation by extending the existing analytical models and by proposing a power modeling methodology that allows us to estimate optimal pipeline depth as a function of *both* power and performance.

3 ANALYTICAL PIPELINE MODEL

In the concept phase definition studies, the exact organization and parameters of the target processor are not known. As such, a custom, cycle-accurate power-performance simulator for the full machine is often not available or relevant. Therefore, the use of analytical reasoning models supplemented by workload characterization and limit studies (obtained from prior generation simulators or trace analysis programs) is common in design groups. We present such an analytical model to understand the power-performance optima and trade offs during the presimulation phase of a design project.

Fig. 1 shows a high-level block diagram of the pipeline model used in our analysis. Our primary goal is to derive the optimum pipeline depth for the various execution units by estimating the various types of stalls in these pipes while using a perfect front-end for the processor. Although Fig. 1 shows only one pipe for each unit (fixed point, floating point, load/store, and branch), the model can be used for a design with multiple pipes per unit as well.

In Fig. 1, let t_i be the latch-free logic time to complete an operation in pipe i and s_i be the number of pipeline stages of pipe i . If c_i is the latch overhead per stage for pipe i , the total time per stage of pipe i is $T_i = ((t_i/s_i) + c_i)$, $\forall i$. The

throughput of the above machine in the absence of stalls is given by $G = \sum(\frac{1}{T_i})$ [8].

We now extend this baseline model to include the effect of data-dependent stalls. Workload analysis is used to derive estimates of the probability that an instruction n depends on another instruction j for all distances $(n - j)$ between them and to estimate the frequency of resulting pipeline stalls. This is illustrated in the example below, where an FXU instruction $(i + 2)$ depends on instruction j , and will be stalled for $(s_1 - 2)$ cycles assuming a bypass to forward the results. Note that, in the above workload analysis, if the source operands of an instruction i are produced by more than one instruction, the largest of the possible stalls is assigned to i .

```
inst (j)      add r1 = r2, r3
inst (j+1)    or r6 = r7, r8
inst (j+2)    and r4 = r1, r5
```

$$T_{fxu} = T_1 + Stall_{fxu-fxu} * T_1 + Stall_{fxu-fpu} * T_2 + Stall_{fxu-lsu} * T_3 + Stall_{fxu-bru} * T_4,$$

where $Stall_{fxu-fxu} = f_1 * (s_1 - 1) + f_2 * (s_1 - 2) + \dots$

The above equation represents the time to complete an FXU operation in the presence of stalls; f_i is the probability that an FXU instruction $(m - i)$ is the producer with the largest stall for instruction m . Similar expressions can be derived for T_{fpu} , T_{lsu} , and T_{bru} , the completion times of an FPU, LSU, and BRU operation, respectively. To account for superscalar (> 1) issue widths, the workload analysis assumes a given issue width along with the number of execution pipes of various types (FXU, FPU, BRU, LSU) and issues independent instructions as an instruction bundle such that the bundle width \leq issue width. Thus, the distance between dependent instructions is the number of instruction bundles issued between them. To account for the dependent instruction stalls due to L1 data cache misses, we use a functional cache simulator to determine cache hits and misses. In addition, we split the load/store pipe into two, namely, load-hit and load-miss pipe, thereby steering all data references that miss in the data cache to the load miss pipeline, which results in longer stall times for the dependent instructions. Since the workload analysis is independent of the machine architecture details and uses only the superscalar issue width to determine the different stalls, it suffices to analyze each application once to derive the stalls.

4 PERFORMANCE AND POWER METHODOLOGY

In this section, we describe the performance simulator used in this study as well as the details of our power modeling toolkit and the methodology that we use to estimate changes in power dissipation as we vary the pipeline depth of the machine.

4.1 Performance Simulation Methodology

We utilize a generic, parameterized, out-of-order 8-way superscalar processor model called Turandot [9], [10] with 32KB I and D-caches and a 2MB L2 cache. The overall pipeline structure (as reported in [9]), is repeated here in Fig. 2. The modeled baseline microarchitecture is similar to a current generation microprocessor. As described in [9],

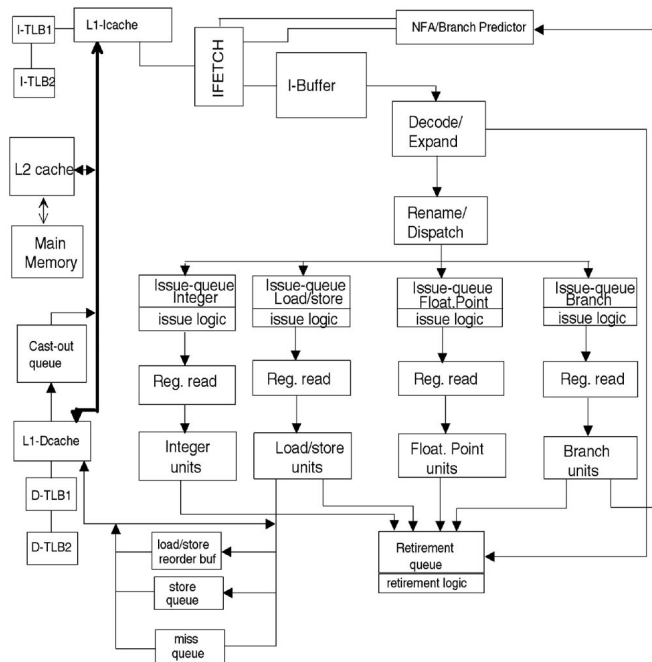


Fig. 2. Modeled processor organization.

this research simulator was calibrated against a pre-RTL, detailed, latch-accurate processor model. Turandot supports a large number of parameters including configurable pipeline latencies discussed below.

Table 1 details the latency values in processor cycles for the 19 FO4 base design point of this study. There is a 2 FO4 latch overhead and 1 FO4 clock skew and jitter overhead. The latency values are then scaled with the FO4-depth. Each latency in Table 1 has two values: The first, labeled STD, is for our detailed simulation model and the second, labeled INF, assumes infinite I-Cache, I-TLB, D-TLB, and a perfect front-end. The INF simulator model is used for validating the analytical model described in Section 3.

4.2 Microprocessor Repipelining

To derive different pipeline design points from the baseline, we consider that the total number of levels of logic—as measured in logic FO4 delays—remains constant regardless of how many pipeline stages the function is partitioned into.² However, as logic is divided into more pipeline stages, the overhead of latch FO4 increases as the number of latches inserted increases with pipeline cut points.

We repipeline the latencies shown in Table 1 using the repipelining equation

$$Latency^{target} = \text{round}\left(Latency^{base} \times FO4_{logic}^{base} / FO4_{logic}^{target}\right).$$

The baseline rounding approach used is “natural rounding,” i.e., $\text{round}(x) = \lfloor x + .5 \rfloor$. All latencies have a minimum of one cycle, regardless of the rounding function used.

Using “natural rounding” for the repipelining assumes that circuits in all pipeline stages have sufficient amount of positive time slack so that, whenever the target latency is rounded down, circuits in all affected pipeline stages can be

2. We note that this is an abstraction for the purpose of this study and an increase in pipeline stages can require changes in the logic design which may impact logic in width and depth.

TABLE 1
Latencies for 19 FO4 Design Point

Fetch Latencies		Decode Latencies	
Latency Parms	STD/INF	Latency Parms	STD/INF
NFA Predictor	1/0	Multiple Decode	2/0
L2 ICache	11/0	Millicode Decode	2/0
L3 (Instruction)	85/0	Expand String	2/0
I-TLB Miss	10/0	Mispredict Cycles	3/0
L2 I-TLB Miss	50/0	Register Read	1/1
Execution Pipe Latencies		Load/Store Latencies	
Latency Parms	STD/INF	Latency Parms	STD/INF
Fix Execute	1/1	L1 D-Load	3/3
Float Execute	4/4	L2 D-Load	9/9
Branch Execute	1/1	L3 (Data)	77/0
Float Divide	12/12	Load Float	2/2
Integer Multiply	7/7	D-TLB Miss	7/0
Integer Divide	35/35	L2 D-TLB Miss	50/0
Retire Delay	2/2	StoreQ Forward	4/0

tuned-up to meet more aggressive timing assertions (constraints) without increasing the energy. On the other hand, whenever the target latency is rounded up, circuits in all affected pipeline stages are assumed not to be tuned down to convert the positive time slack into energy savings.

Under realistic design constraints, restructuring logic or tuning circuits for more aggressive timing constraints to remove the negative timing slack, $\delta D < 0$ inevitably increases the energy dissipation. On the other hand, all macros with a significant positive slack $\delta D > 0$ should be tuned-down or logically restructured to convert the positive slack into energy savings. As a result, the average energy dissipated in the macro changes as $\frac{\delta E}{E} = -\eta \frac{\delta D}{D}$, where η is the hardware intensity³ in the affected macro [11].

To estimate the impact of the selected repipelining method, we have also used strict upper and lower bounds for repipelining latencies by selecting rounding functions $\lceil x \rceil$ and $\lfloor x \rfloor$, corresponding to the two extremes of ample available slack and zero slack, respectively. As shown in Fig. 3, the architectural performance expressed in CPI for these scenarios tracks well, suggesting that the specific repipelining method does not change the trend of the data. In particular, any deviation is not skewed toward a particular pipeline range, but, rather, is periodic in integer multiples of the original design. The CPI impact of rounding is typically within 10 percent for the two extreme scenarios which bound the “natural rounding” scenario.

4.3 Power Modeling Methodology

To estimate power dissipation, we use the PowerTimer toolset developed at the IBM T.J. Watson Research Center

3. We use *hardware intensity* η as a quantitative measure of how aggressively circuits in a hardware macro are tuned to meet the target clock frequency. Hardware intensity shows the energy cost (in percent) required to improve the critical path delay D in the macro by 1 percent by restructuring the logic and retuning the circuits, at a fixed power supply, $\eta = -\frac{\%E}{\%D}$ through retuning.

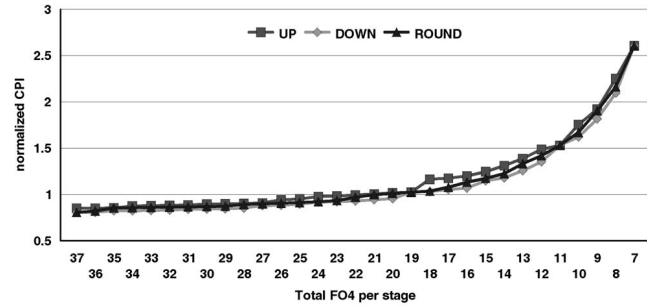


Fig. 3. CPI impact of repipelining design and impact of rounding modes.

[12], [1] as the starting point for the simulator used in this work. PowerTimer is similar to power-performance simulators developed in academia [13], [14], [15], except for the methodology to build energy models.

Fig. 4 depicts the derivation of the energy models in PowerTimer. The energy models are based on circuit-level power analysis that has been performed on structures in a current, high performance PowerPC processor. The power analysis has been performed at the macro level using a circuit-level power analysis tool [16]. Generally, multiple macros combine to form one microarchitectural level structure which we will call a subunit. For example, the fixed-point issue queue (one subunit) might contain separate macros for storage memory, comparison logic, and control. Power analysis has been performed on each macro to determine the macro’s unconstrained (no clock gating) power as a function of the input switching factor. First, the *hold power*, or power when no switching is occurring (SF = 0 percent), is determined. Hold power primarily consists of power dissipated in latches, local clock buffers, the global clock network, and data-independent fraction of the arrays. The *switching power*, which is primarily combinatorial logic and data-dependent array power dissipation, is the additional power that is dissipated when switching factors are applied to the macro’s primary inputs. These two pieces of data allow us to form simple linear equations for each macro’s power. The energy model for a subunit is determined by summing the linear equations for each macro within that subunit. We have generated these power models for all microarchitecture-level structures (subunits) modeled in our research simulator [9], [10]. PowerTimer models over 60 microarchitectural structures which are defined by over 400 macro-level power equations.

PowerTimer uses microarchitectural activity information from the Turandot model to scale down the unconstrained hold and switching power on a per-cycle basis under a

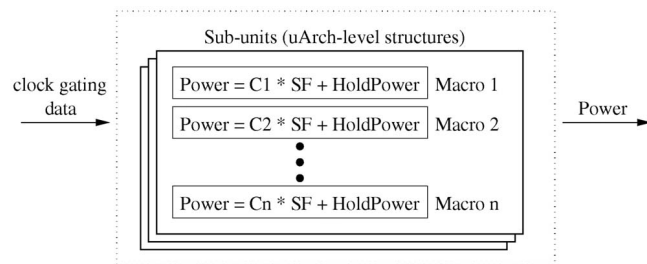


Fig. 4. PowerTimer energy models.

variety of clock gating assumptions. In this study, we use a realistic form of clock gating which considers the applicability of clock gating on a per-macro basis to scale down either the hold power or the combined hold and switching power depending on the microarchitectural event counts. We determine which macros can be clock gated in a fine-grained manner (per-entry or per-stage clock gating) and which can be clock gated in a coarse-grained manner (the entire unit must be idle to be clock gated). For some macros (in particular control logic), we do not apply any clock gating; this corresponds to about 20-25 percent of the unconstrained power dissipation. The overall savings due to clock gating relative to the unconstrained power is roughly 40-45 percent.

In order to quantify the power-performance efficiency of pipelines of a given FO4-depth and to scale the power dissipation from the power models of our base FO4 design point across a range of FO4-depths, we have extended the PowerTimer methodology as discussed below.

Power dissipated by a processor consists of dynamic and leakage components, $P = P_{\text{dynamic}} + P_{\text{leakage}}$. The dynamic power data measured by PowerTimer can be expressed as $P_{\text{dynamic}}^{\text{base}} = CV^2f * (\alpha + \beta) * CGF$, where α is the average “true” switching factor in circuits, i.e., α represents transitions required for the functionality of the circuit and is measured as the switching factor by an RTL-level simulator run under the zero-delay mode. In contrast, β is the average glitching factor that accounts for spurious transitions in circuits due to race conditions. Thus, $(\alpha + \beta)$ is the average number of transitions actually seen inside circuits. Both α and β are averaged over the whole processor over nongated cycles with appropriate energy weights (the higher the capacitance at a particular node, the higher the corresponding energy weight). CGF is the clock gating factor, which is defined as the fraction of cycles where the microarchitectural structures are *not* clock gated. The CGF is measured from our PowerTimer runs for each architectural component at each FO4 design point. The remaining terms C , V , and f are effective switching capacitance, chip supply voltage, and clock frequency, respectively.

Next, we analyze how each of these factors scales with FO4 pipeline depth. To facilitate the explanation, we define the following variables: $FO4_{\text{logic}}$, $FO4_{\text{latch}}$ and $FO4_{\text{pipeline}}$, to designate the depth of the critical path through logic in one pipeline stage, the latch insertion delay including clock skew and jitter, and the sum of the two quantities, respectively, $FO4_{\text{pipeline}} = FO4_{\text{logic}} + FO4_{\text{latch}}$. We use $FO4$ and $FO4_{\text{pipeline}}$ interchangeably. In the remainder of the paper, the qualifier “base” in all quantities designates the value of the quantities measured for the base 19 FO4 design.

Frequency. FreqScale is the scaling factor that is used to account for the changes in the clock frequency with the pipeline depth. This factor applies to both hold power and switching power:

$$\text{FreqScale} = FO4_{\text{pipeline}}^{\text{base}} / FO4_{\text{pipeline}}$$

Latch. With fixed logic hardware for given logic functions, the primary change in the chip effective switching capacitance C with the pipeline depth is due to changes in the latch count. LatchScale is a factor that

appropriately adjusts the hold power dissipation, but does not affect the switching power dissipation.

$$\text{LatchScale} = \text{LatchRatio} * \left(\frac{FO4_{\text{logic}}^{\text{base}}}{FO4_{\text{logic}}} \right)^{LGF},$$

where LatchRatio defines the ratio of hold power to the total power and $\text{LatchGrowthFactor}(LGF)$ captures the growth of latch count due to the logic shape functions. The amount of additional power that is spent in latches and clock in a more deeply pipelined design depends on the *logic shape functions* of the structures that are being pipelined. Logic shape functions describe the number of latches that would need to be inserted at any cut point in a piece of combinatorial logic if it were to be pipelined.

Values of $LGF > 1$ recognize the fact that, for certain hardware structures, the logic shape functions are not flat and, hence, the number of latches in the more deeply pipelined design increases superlinearly with pipeline depth. In our baseline model, we assume an overall LGF of 1.1 and study the sensitivity of the optimal pipeline depth to this parameter in Section 6.1.

Clock Gate Factor. In general, CGF decreases with deeper pipelines because the amount of clock gating potential increases with deeper pipes. This increased clock gating potential is primarily due to the increased number of cycles where pipeline stages are in stall conditions. This in turn leads to an increase in the clock gating potential on a per cycle basis. CGF is workload dependent and is measured directly from simulator runs.

Glitch. The final two factors that must be considered for dynamic power dissipation when migrating to deeper pipelines are α and β , the chip-wide activity and glitching factors.

The “true” switching factor α does not depend on the pipeline depth since it is determined by the functionality of the circuits. The glitching factor at any net, on the other hand, is determined by the difference in delay of paths from the output of a latch that feeds the circuit to the gate that drives that net. Once a glitch is generated at some net, there is a high probability that it will propagate down the circuit until the input of the next latch down the pipeline. Furthermore, the farther the distance from the latch output to the inputs of a gate the higher the probability of the existence of nonequal paths from the output of the latch to the inputs of this gate. Therefore, the average number of spurious transitions grows with $FO4_{\text{logic}}$ —the higher the FO4, the higher the average glitching factor. Experimental data collected by running a dynamic circuit-level simulator on extracted netlists of sample functional units show that the average glitching factor β can be modeled as being linearly dependent on the logic depth [17], $\beta = \beta_{\text{base}} \frac{FO4_{\text{logic}}}{FO4_{\text{logic}}^{\text{base}}}$. To account for the effect of the dependence of β on pipeline depth, we introduce the following factor which applies only to the switching power:

$$\text{GlitchScale} = \frac{1 - \text{LatchRatio}}{1 + \beta_{\text{base}}/\alpha} \left(1 + \frac{\beta_{\text{base}} \text{FO4}_{\text{logic}}}{\alpha \text{FO4}_{\text{logic}}^{\text{base}}} \right).$$

In this formula, β_{base} is the actual glitching factor averaged over the baseline microprocessor for the base FO4 design point. Notice that β_{base} appears in the formula only in the ratio $\beta_{\text{base}}/\alpha$. This is consistent with our experimental results showing that the glitching factor β is roughly proportional to the “true” switching factor α , for the range $0 < \alpha < 0.3$ (for higher values of α , the growth of β typically saturates). For the set of six sample units that we simulated, with the logic depth ranging from 6 FO4 to 20 FO4, the ratio β/α was found to be roughly proportional to the logic depth of the simulated units, $\text{FO4}_{\text{logic}}$, with the coefficient equal to $0.3/\text{FO4}_{\text{logic}}^{\text{base}}$ [17]. Based on these simulation results, we set $\beta_{\text{base}}/\alpha = 0.3$ for the whole microprocessor in the remainder of this section, and study the sensitivity of the results to variations in the glitching factor in Section 6.4.

Leakage Currents. As the technology feature size scales down and the power supply and transistor threshold voltages scale accordingly, the leakage power component becomes more and more significant. Since the magnitude of the leakage power component is affected by the pipeline depth, it is essential to include the effect of the leakage power in the analysis of the optimal pipeline depth. Assuming that the leakage power is proportional to the total channel width of all transistors in the microprocessor, we model the dependence of the leakage power on the depth of the pipeline as follows:

$$P_{\text{leakage}}^{\text{FO4}} = P_{\text{leakage}}^{\text{base}} \left(1 + \frac{w_{\text{latch}}}{w_{\text{total}}} \left[\left(\frac{\text{FO4}_{\text{logic}}^{\text{base}}}{\text{FO4}_{\text{logic}}} \right)^{\text{LGF}} - 1 \right] \right),$$

where LGF is the *LatchGrowthFactor* defined earlier, $w_{\text{latch}}/w_{\text{total}}$ is the ratio of the total channel width of transistors in all pipeline latches (including local clock distribution circuitry) to the total transistor channel width in the base (19 FO4) microprocessor (excluding all low-leakage transistors that might be used in caches or other on-chip memories). If the technology supports multiple thresholds or any of the recently introduced leakage reduction techniques are used on a unit-by-unit basis, such as MTCMOS, back biasing, power-down, or transistor stacking, then the above formula for the leakage power component needs to be modified accordingly and we leave the detailed study of these effects for future work. For the remainder of the study, we empirically set $w_{\text{latch}}/w_{\text{total}}$ to 0.3 for a base 19 FO4 pipeline, assuming that transistors on the scan path and the feedback path in latches, and possibly on the clock distribution path, may use high threshold devices. Even though the portion of the actual area of latches may be larger, when scaling the leakage currents with pipeline depth, we are concerned with the area of low/normal threshold devices. We think that there is higher potential for using high threshold devices in latches than in the logic. Also, rather than giving the absolute value for the leakage current, $P_{\text{leakage}}^{\text{base}}$ in the base microprocessor, we will count it as a fraction of the dynamic power of the base design,

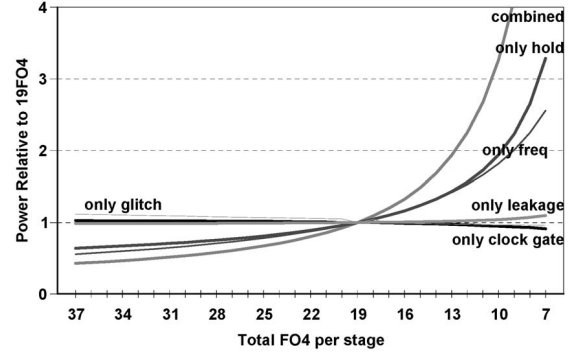


Fig. 5. Power growth breakdown.

$$P_{\text{leakage}}^{\text{base}} = \text{LeakageFactor}^{\text{base}} P_{\text{dynamic}}^{\text{base}}.$$

We set the $\text{LeakageFactor}^{\text{base}}$ to the value of 0.1, typically quoted for state of the art microprocessors, and analyze the sensitivity of the results to LeakageFactor in Section 6.5.

Total Power. The following equation expresses the relationship between the dynamic power for the base FO4 design, $P_{\text{dynamic}}^{\text{base}}$, leakage power, and the scaled power for designs with different depths of the pipeline, $P_{\text{total}}^{\text{FO4}}$, considering all factors above:

$$P_{\text{total}}^{\text{FO4}} = \text{CGF} * \text{FS} * (\text{LS} + \text{GS}) * P_{\text{dynamic}}^{\text{base}} + P_{\text{leakage}}^{\text{FO4}},$$

where FS is *FreqScale*, LS is *LatchScale*, and GS is *GlitchScale*.

Fig. 5 shows contributions of different factors in the above formula, depending on the pipeline depth of the design, $\text{FO4}_{\text{pipeline}}$. The 19 FO4 design was chosen as a base pipeline.

The line labeled “combined” shows the cumulative increase or decrease in power dissipation. The line labeled “only clock gate” quantifies the amount of additional clock gating power savings for deeper pipelines. The relative effect of scaling in clock gating is fairly minor, with slightly more than 10 percent additional power reduction when going from the 19 FO4 to 7 FO4 design points. There are several reasons why the effect of clock gating is not larger. First, the fraction of power dissipation that is not eligible to be clock gated becomes larger with more clock gating, leading to diminishing returns. Second, some of the structures are clock gated in a coarse-grained fashion and, while the average utilization of the structure may decrease, it must become idle in all stages before any additional savings can be realized. Finally, we observe that clock gating is more difficult in deeper pipelined machines because it is harder to deliver cycle-accurate gating signals at lower FO4.

The two lines labeled “only freq” and “only hold” show the power factors due to only frequency and hold power scaling, respectively.⁴ Overall, dynamic power increases more than quadratically with increased pipeline depth.

Fig. 5 shows that the leakage component grows much less rapidly than the dynamic component with the increasing pipeline depth. There are two primary reasons for this. First, the leakage power does not scale with frequency.

4. Although these two factors increase with clock frequency, we are plotting against FO4-depth which is $1/\text{clock frequency}$.

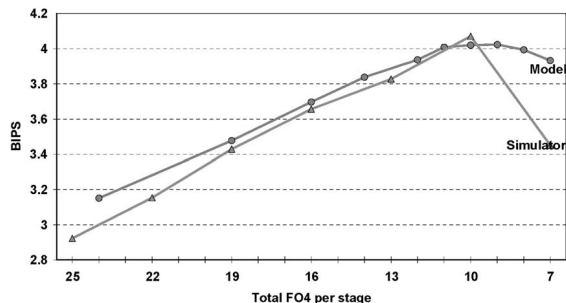


Fig. 6. BIPS: Analytical model versus simulator.

Second, the leakage power growth is proportional to the fraction of channel width of transistors in pipeline latches, whereas the latch dynamic hold power growth is proportional to the fraction of the dynamic power dissipated in pipeline latches. Obviously, the former quantity is much smaller than the latter.

4.4 Workloads and Metrics Used in the Study

In this paper, we report experimental results based on PowerPC traces of a set of 21 SPEC2000 benchmarks, namely, *ammp*, *applu*, *apsi*, *art*, *bzip2*, *crafty*, *equake*, *facerec*, *gap*, *gcc*, *gzip*, *lucas*, *mcf*, *mesa*, *mgrid*, *perl*, *sixtrack*, *swim*, *twolf*, *vpr*, and *wupwise*. We have also used a 172M instruction trace of the TPC-C transaction processing benchmark. The SPEC2000 traces were generated using the tracing facility called *Aria* within the MET toolkit [10]. The particular SPEC2000 trace repository used in this study was created by using the full reference input set. However, sampling was used to reduce the total trace length to 100 million instructions per benchmark program. A systematic validation study to compare the sampled traces against the full traces was done in finalizing the choice of exact sampling parameters [18].

We use $BIPS^3/W$ ($energy * delay^2$) as a basic energy-efficiency metric for comparing different FO4 designs in the power-performance space. The choice of this metric is based on the observation that dynamic power is roughly proportional to the square of supply voltage (V) multiplied by clock frequency and clock frequency is roughly proportional to V . Hence, power is roughly proportional to V^3 , assuming a fixed logic/circuit design. Thus, delay cubed multiplied by power provides a voltage-invariant power-performance characterization metric which we feel is most appropriate for server-class microprocessors (see discussion in [1], [19]). In fact, it was formally shown in [17] that optimizing performance subject to a constant power constraint leads to the $BIPS^3/W$ metric in processors operating at a supply voltage near the maximum allowed value in state of the art CMOS technologies. As a comparative measure, we also consider $BIPS/W$ (energy) and $BIPS^2/W$ (energy-delay [20]) in the baseline power-performance study.

5 ANALYTICAL MODEL AND SIMULATION RESULTS

5.1 Analytical Model Validation

We now present the performance results using the analytical model described in Section 3 and compare it

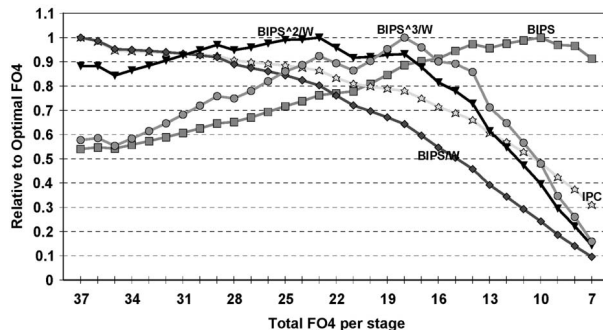


Fig. 7. Simulation results for SPEC2000.

with the results using our detailed cycle-accurate simulator. The results in this section are presented for the average of the SPEC2000 benchmarks described in Section 4. For fair comparison, we have modified the simulator to include a perfect front-end of the machine. The model and the simulator use latencies shown in the column labeled "INF" in Table 1.

Fig. 6 shows BIPS as a function of the FO4 delay per stage of the pipeline. The BIPS for the analytical model was computed after determining the stalls for the different pipelines using independent workload analysis, as explained in Section 3.

From Fig. 6, we observe that the performance optimal pipeline depth is roughly 10 FO4 delay per stage for both the model and simulator. Although the BIPS estimated from the model correlates reasonably well with the BIPS estimated using the simulator, it is true that the analytical model makes many simplifying assumptions while determining the stalls. For example, the analytical model does not account for overlapping data dependencies between instructions, assumes perfect memory disambiguation, uses infinite rename registers, miss queue, and instruction buffer, and does not exploit the out-of-order execution to hide cache miss latencies. For these reasons, the BIPS determined by the model may be higher or lower than the BIPS obtained from the simulator, as seen for the very shallow or deep pipelines. However, the ability of the model to capture the trend in performance when varying the pipeline depth, combined with its generality and simplicity, makes it useful for concept-phase analysis.

5.2 Detailed Power-Performance Simulation

In the remainder of this work, we consider the power and performance results using the detailed performance simulator with parameters corresponding to the STD column in Table 1. Fig. 7 shows the plots for five metrics: BIPS, IPC, $BIPS/W$, $BIPS^2/W$, and $BIPS^3/W$.

Fig. 7 shows that the optimal FO4-depth for performance (defined by BIPS) is 10 FO4, although pipelines of 8 FO4 to 15 FO4 are within 5 percent of the optimal. Because of the superlinear increase in power dissipation and sublinear increases in overall performance, the $BIPS/W$ always decreases with deeper pipelines. $BIPS^3/W$ shows an optimum point at 18 FO4. $BIPS^3/W$ decreases sharply after the optimum and at the performance-optimal pipeline depth of 10 FO4, the $BIPS^3/W$ metric is reduced by 50 percent over the 18 FO4 depth. For metrics that have less emphasis on performance, such as $BIPS^2/W$ and

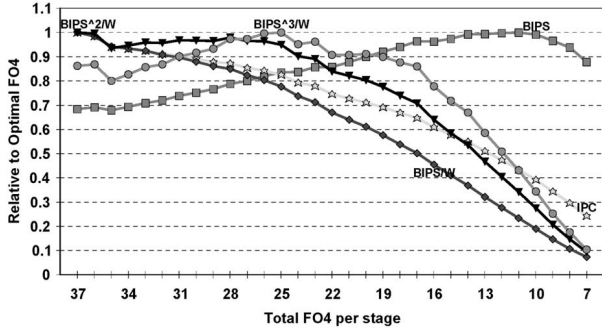


Fig. 8. Simulation results for TPC-C.

BIPS/W, the optimal point shifts toward shallower pipelines, as expected. For the BIPS²/W, the optimal pipeline depth is achieved at 23 FO4.

Fig. 8 presents similar results for the TPC-C trace. The BIPS metric for TPC-C is very flat in the range from 10 FO4 to 14 FO4 around the optimal point (within 1-2 percent for all five design points). The optimum for the BIPS³/W metric shifts over to 25-28 FO4. The main reason why the optimal point is shallower for TPC-C than for SPEC2000 is that TPC-C exhibits high branch misprediction and data cache miss rates, both of which reduce the performance benefits of deeper pipelines. This effect is slightly counterbalanced by a slower power growth for TPC-C with deeper pipes, due to additional opportunities for clock gating resulting from excessive pipeline stalls.

6 SENSITIVITY ANALYSIS

The derived equations that model the dependence of the power dissipation on the pipeline depth depend on several parameters. Some of these parameters, although accurately measured for the baseline microprocessor, are likely to change from one design to another, whereas others are difficult to measure accurately. In this section, we perform sensitivity analysis of the optimal pipeline depth to key parameters of the derived power models such as *LatchGrowthFactor*, *LatchRatio*, latch insertion delay ($FO4_{latch}$), *GlitchRatio*, and *LeakageFactor*.

6.1 Latch Growth Factor

LatchGrowthFactor (*LGF*) is determined by the intrinsic logic shape functions of the structures that are being pipelined. We have analyzed many of the major microarchitectural structures to identify ones that are likely to have *LatchGrowthFactor* greater than 1. One structure that we will highlight is the Booth recoder and Wallace tree, which is common in high-performance floating-point multipliers [21], as shown in Fig. 9. Fig. 9 shows the exponential reduction in the number of result bits as the data passes through the 3-2 and 4-2 compressors. We have estimated the amount of logic that can be inserted between latch cut points for 7, 10, 13, 16, and 19 FO4 designs by assuming 3-4 FO4 delay for 3-2 and 4-2 compressor blocks. As the 7 and 10 FO4 design points require latch insertions just after the Booth multiplexor (where there are 27 partial products), there would be a large increase in the number of latches required for these designs. We note that the 7 FO4 design also requires a latch insertion after the booth recode stage.

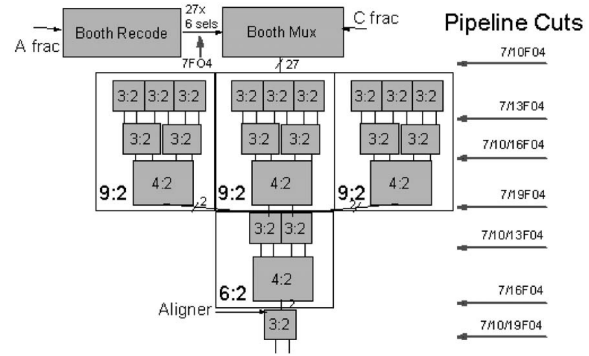


Fig. 9. Wallace Tree diagram and Latch Cut points for 7/10/13/16/19 FO4.

Fig. 10 gives estimates for the cumulative number of latches in the FPU design as a function of the FO4 depth of the FPU. For example, the first stage of the 10 FO4 FPU requires 3x as many latches as the first stage of the 19 FO4 FPU because the first latch cut point of the 19 FO4 FPU is beyond the initial 9:2 compressor tree. Overall, the 10 FO4 FPU requires nearly 3x more latches than the 19 FO4 FPU.

There are many other areas of parallel logic that are likely to see superlinear increases in the number of latches such as structures with decoders, priority encoders, carry look ahead logic, etc. Beyond the pipelining of logic structures, deeper pipelines may require more predecode information to meet aggressive cycle times, which would require more bits to be latched in intermediate stages. On the other hand, the number of latches that comprise storage bits in various on-chip memory arrays (such as register files and queues) does not grow at all with the pipeline depth, meaning that the $LGF = 0$ for those latches. Thus, designs with overall $LGF < 1$ are also possible.

Fig. 11 quantifies the dependence of the optimal pipeline depth on *LGF* using the BIPS³/W metric. It shows that the optimal pipeline FO4 tends to increase as *LGF* increases above the value of 1.1, assumed in the baseline model. As a point of reference, our estimate for the latch growth factor for the 10 FO4 versus 19 FO4 Booth recoder and Wallace tree is $LGF = 1.5$, while, for the entire FPU, *LGF* is slightly smaller.

6.2 Latch Power Ratio

Latch, clock, and array power are the primary components of power dissipation in current generation CPUs. This is especially true in high-performance, superscalar processors

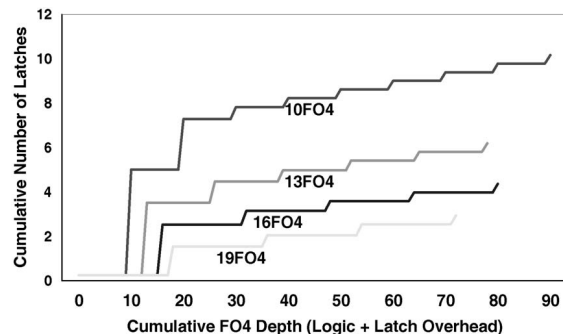


Fig. 10. Cumulative latch count for FPU.

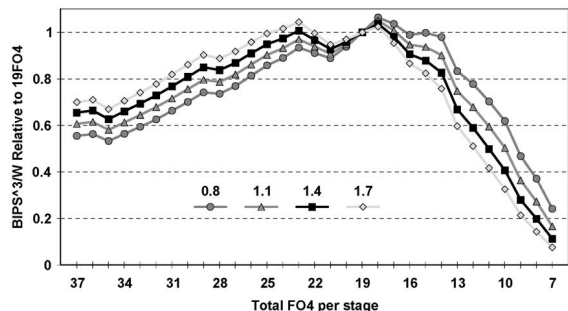


Fig. 11. $BIPS^3/W$ varying *LatchGrowthFactor*.

with speculative execution which require the CPU to maintain an enormous amount of architectural and non-architectural state. One possible reason why the *LatchRatio* could be smaller than the base value of 0.7 chosen in Section 4 is if more energy-efficient SRAM arrays are used in high-power memory structures instead of latches to reduce the data independent array power (which we include as part of hold power).

Fig. 12 shows the optimal FO4 design point while varying the *LatchRatio* of the machine from 80 percent to 40 percent. We see that, while the optimal FO4 point remains 18 FO4, it is less prohibitive to move to deeper pipelines with smaller latch-to-logic ratios. For example, with a *LatchRatio* of 0.4, the 13 FO4 design point is only 19 percent worse than the optimal one, while it is 27 percent worse than optimal with a *LatchRatio* of 0.6.

6.3 Latch Insertion Delay

With the large amount of power spent in latches and clocking, designers may consider the trade off between latch delay and power dissipation as a means to design more energy-efficient CPUs. Researchers have investigated latch power versus delay trade off curves both within a given latch family and across latch families [17], [22], [23]. Table 2, derived from [17], shows latch FO4-delay versus latch energy across several latch styles. The first row of Table 2 shows the latch insertion delay, excluding clock skew and jitter overhead, which is assumed to be constant for all latches. The second row shows the relative latch energy, excluding energy dissipated in the clock distribution. The third row shows the relative energy of the clock system, including both energy of latches (80 percent of the total) and clock distribution system (20 percent). It is assumed that the clock distribution energy cannot be completely scaled with reduced latch load. There is

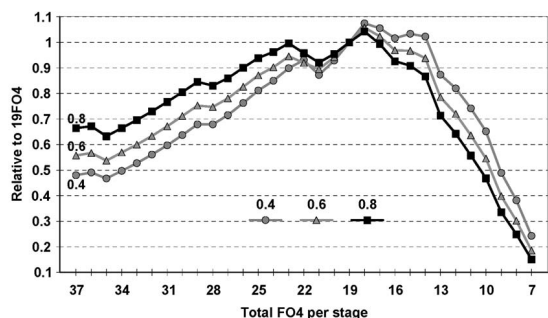


Fig. 12. $BIPS^3/W$ varying *LatchRatio*.

TABLE 2
Latch Insertion Delay (Excluding Skew and Jitter) versus Relative Latch Energy

Latch FO4	2	3	4	5
Relative Latch Energy	1.0	0.53	0.36	0.29
Relative Clocking Energy	1.0	0.62	0.49	0.43

significant overhead simply from driving the wires necessary to distribute the clock over a large area and satisfying the slow constraints.

Replacing the baseline fast 2 FO4 latches with slower, lower power latches increases the latch insertion delay overhead, which impacts both the performance-optimal and power-performance-optimal pipeline depth. Fig. 13 shows the processor performance versus pipeline depth for four latches from Table 2. We see that the performance maxima shift toward shallower pipelines as the latch insertion delay increases. For example, with a 3 FO4 latch, the performance-optimal FO4-depth is 11 FO4 and, with a 4 FO4 latch, it becomes 16 FO4.

Fig. 14 shows the impact of the latch insertion delay on the $BIPS^3/W$ rating of the processor for the same range of pipeline depths. In this figure, all of the data points are shown relative to the 10 FO4 design with the base 5 FO4 latch. Unlike curves on all previous sensitivity graphs, curves in Fig. 14 do not intersect at the base design point because different curves represent different designs, with different power and performance levels.

Fig. 14 shows that using the fastest 2 FO4 latch results in the best $BIPS^3/W$ rating for processors with stages less than 14 FO4. For processors with pipelines ranging from 15 FO4

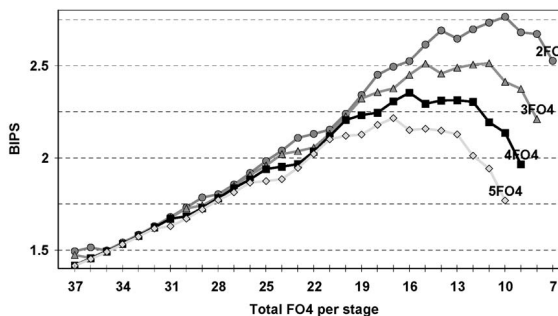


Fig. 13. BIPS varying latch power-delay.

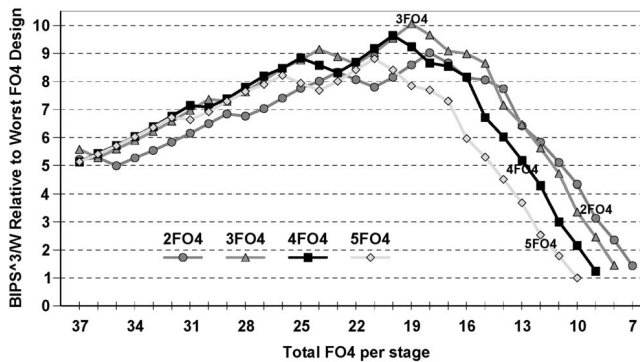
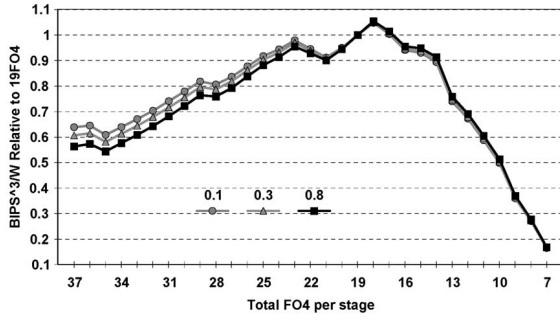


Fig. 14. $BIPS^3/W$ varying latch power-delay.

Fig. 15. $BIPS^3/W$ varying β_{base}/α .

to 24 FO4, a lower power 3 FO4 latch is the most energy efficient, whereas, for shallower pipelines (25 FO4 or more), the highest $BIPS^3/W$ is achieved with even slower 4 FO4 latches.

The use of the 3 FO4 latch, combined with the choice of the pipeline depth in the range from 15 FO4 to 24 FO4 improves the $BIPS^3/W$ rating of the processor by more than 10 percent compared to the base case of 2 FO4 latches. The graph also shows that the optimal $BIPS^3/W$ design point shifts toward shallower pipelines as high-performance latches are replaced with lower power ones. For example, the 18 FO4 design point is optimal for a processor using 2 FO4 latches, whereas the 19 FO4, 20 FO4, and 21 FO4 design points are optimal for processors using 3 FO4 latches, 4 FO4, and 5 FO4 latches, respectively.

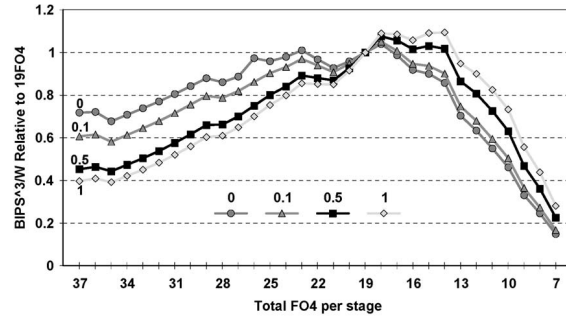
6.4 Glitch Factor

In this section, we quantify the sensitivity of the optimal pipeline depth to the glitching factor. There are no practical means for accurately measuring the actual value of β_{base}/α , averaged over the whole microprocessor. Instead, we measured the glitching factor for a selected set of functional units and used the averaged value of $\beta_{base}/\alpha = 0.3$ throughout the analysis [17]. In this section, we analyze the sensitivity of the optimal pipeline depth to the value of β_{base}/α .

Fig. 15 shows the dependence of the $BIPS^3/W$ rating of the processor on the pipeline depth for three values of β_{base}/α . From this figure, we see that higher glitching factors favor deeper pipelines. However, in the base design, the decrease in power dissipation related to reduced glitching in deeper pipelines did not have a substantial impact on the optimal pipeline depth, primarily because of the relatively small fraction of power dissipated in combinatorial switching. For designs which have smaller *LatchRatio* values, this effect could be more significant.

6.5 Leakage Factor

As explained earlier, the leakage power component grows more slowly with the pipeline depth than the dynamic component. Therefore, the optimum pipeline depth depends on the *LeakageFactor*. Throughout the analysis, we assumed that, for the base 19 FO4 microprocessor, the *LeakageFactor* ($P_{leakage}^{base}/P_{dynamic}^{base}$) to be 0.1. However, as the technology feature size scales down and the power supply and transistor threshold voltages scale accordingly, the leakage power component becomes more and more significant. To study the effect of the growing fraction of the leakage power component, we measured the sensitivity

Fig. 16. $BIPS^3/W$ varying *LeakageFactor*.

of the optimal pipeline depth to the value of the *LeakageFactor*.

Fig. 16 shows the $BIPS^3/W$ rating of the processor versus pipeline depth for three values of the *LeakageFactor*: a value of 0 that represents older CMOS technologies, a value of 0.1, assumed in the current model, and values of 0.5 and 1.0, projected for future generation CMOS technologies (arguably, extreme values). The results in Fig. 16 show that, unless leakage reduction techniques become the standard practice in the design of high-end microprocessors, the high values of the *LeakageFactor* projected for future generations of CMOS technologies may tend to shift the optimum pipeline depth towards deeper pipelines. For current generation technologies, the result for the optimal pipeline depth is sufficiently stable with respect to reasonable variations in the *LeakageFactor*.

Summary of the Sensitivity Analysis. In this section, we considered the sensitivity of optimal pipeline length to five key parameters in the power models using the $BIPS^3/W$ metric. We did not observe a strong dependence of the results on the assumptions and choices of any of these parameters, which demonstrates the stability of the model, and its applicability to a wide range of designs. To summarize the results, higher values of the *LatchGrowthFactor* favor shallower pipelines, lower values of the *LatchRatio* favor deeper pipelines, the use of lower-power latches favors shallower pipelines, higher values of the *GlitchFactor* favor deeper pipelines, and, finally, higher leakage currents favor deeper pipelines.

7 WORKLOAD VARIABILITY

We have studied the distribution of optimal pipeline depth as a function of the workload characteristics. In Fig. 17, we show how the workload optima are distributed.

We note that the optima are clustered at a few preferred configurations. This clustering is a response to the discretization introduced by dividing the data path into several pipeline stages. The (microarchitecturally) optimal points are then those which utilize a given FO4 delay best by having minimal slack, i.e., just before a new pipeline stage must be introduced to achieve the next decrement in FO4 delay per pipeline stage. At this point, the frequency gain for reducing the delay per pipeline stage is $FO4/(FO4-1)$, ranging from $(8/7) = 14$ percent in the very high frequency regime to $(19/18) = 6$ percent in a moderately high frequency design. On the other hand, CPI degrades as a step function response to the increment in latency of a

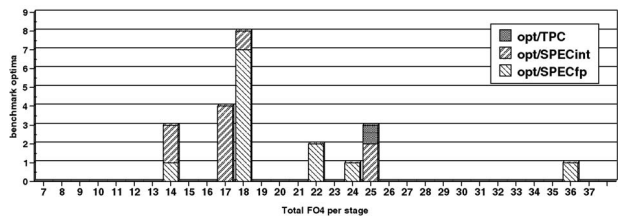


Fig. 17. Distribution of power/performance optimal pipeline depth.

particular unit, making that design point less attractive from a performance perspective.

Examples of this behavior can be seen in Fig. 17 at the transition from 18 to 17 FO4 which increases the latency of the FP unit by one cycle and, hence, degrades CPI for floating-point intensive benchmarks. This causes significant clustering of optima for FP-intensive benchmarks at 18 FO4. Similarly, other key microarchitectural latencies change at the transition from 17 to 16 FO4 and 14 to 13 FO4, leading to clustering of optima of benchmarks with high sensitivity to those particular parameters at 17 and 14 FO4, respectively.

We note that the *lucas* benchmark has a power/performance optimum at a pipeline stage delay of 36 FO4 as its performance curve is essentially flat across any pipeline depth. However, power increases almost as badly as with any other workload due to the increase in latches and clock frequency. While there is some saving due to increased ability to gate the clock at a finer level, this improvement is measured on the order of 20 percent over the polynomial growth in power shown in Fig. 5. The *lucas* benchmark is an extreme case in that it is dominated by L2 misses. As such, more aggressive pipelining has no impact on performance as L2 behavior scales with memory system performance, not core architecture.

8 DESIGN IMPLICATIONS

Fig. 18 gives a graphical interpretation of optimizing the pipeline depth in the power-performance space. Power, plotted on the y-axis, is the average power dissipated on the set of benchmarks, measured in relative units. Delay, plotted on the x-axis, is the average execution time for the same set of benchmarks, measured in relative units. Microarchitectural configurations with different pipeline depths are represented by points on the bold curve marked with with "□." Thus, the bold curve (□) shows how the processor design point moves in the power-performance space as the depth of the pipeline changes, assuming unchanged power supply and circuit tuning point or hardware intensity. The deeper the pipeline (less FO4), the steeper the part of the curve where the corresponding design point sits.

The curves marked with "○" show the circuit energy-delay trade off curves which represent implementations of the microarchitecture with a fixed pipeline depth for varying power supply and circuit hardware intensity such that the condition for the optimal hardware intensity defined in [11] is observed at each point. The base hardware intensity for each circuit energy-delay curve is set to the value of 2, which we found to be typical of modern high-performance microprocessors operating at power supply voltages close to the upper limit of the state of the art CMOS technology. Thus, for any fixed pipeline depth, we can

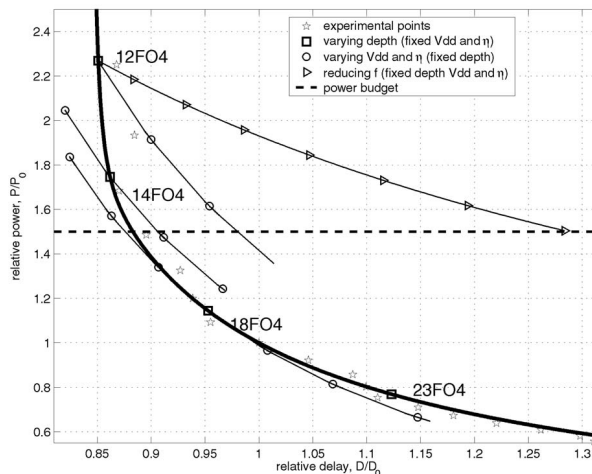


Fig. 18. Cumulative impact of tuning microarchitecture, power supply, and circuits.

move the power-performance point of the processor along one of the hardware tuning curves (○) in the vicinity of the base point. Notice that, on any of these curves (○), every 1 percent in performance (due to frequency) achieved through changing the power supply and retuning the circuits costs 3 percent in power.

The 18FO4 point on the bold curve (□) that minimizes the performance-cube per watt metric is the point where a 1 percent gain in performance through increasing the pipeline depth also costs 3 percent in power. Then, the 18FO4 point is the only point on the bold curve (□) at which the architectural energy-delay trade off curve tangents the circuit energy-delay trade off curve. At this pipeline depth, the processor delivers the required performance while dissipating the minimum power among implementations with any depth of the pipeline. To graphically demonstrate this, consider an initial design with a 14FO4 pipeline. This pipeline depth minimizes the $\frac{\text{BIPS}^{10}}{\text{Watt}}$, which means that, at this point, every percent in performance achieved through increasing the pipeline depth costs 10 percent in power. Thus, the initial design point of 14FO4 sits on a steeper part of the curve than the 18FO4 design point. By reducing the pipeline depth to 18FO4, we can move the processor power-performance point down the bold curve (□), saving (initially) 10 percent in power for every 1 percent loss in performance. Then, by raising the power supply and tuning up the circuits accordingly, we can move the design point up the corresponding hardware tuning curve (○) to recover the performance, spending approximately 3 percent in power for every 1 percent gain in performance. Obviously, the resulting design with 18FO4 is a better design point than the initial 14FO4 because it delivers the same performance at a lower power.

Similarly, if the original pipeline has a depth of 23FO4, which minimizes the BIPS-squared per watt metric, then, by increasing the pipeline depth to 18FO4, we can move the design point up the bold curve (□), spending (initially) approximately 2 percent in power per every percent in gained performance. If the higher performance resulting from the increase in the pipeline depth is not needed, the design point can be moved down the corresponding hardware tuning curve (○) by reducing power supply and

tuning down the circuits, thereby saving 3 percent in power per every percent in performance. Since the circuit energy-delay trade off curve at this point is steeper than the architectural energy-delay trade off curve, the resulting design with 18FO4 will deliver the same performance as the initial one, but at a lower power.

Now, suppose that the design with the optimal 18FO4 pipeline does not deliver the advertised frequency and a business decision is made to go with a deeper than optimum pipeline (say 12FO4). As a result, if the final product still meets the power-dissipating capability of the package (shown by the power budget line in Fig. 18), that company will just end up with an energy-inefficient design (implying that a higher net performance could be delivered at the same power budget). However, if the final product exceeds the power-dissipating capabilities of the package when running at maximum speed (which is a more likely scenario), then stepping down the bold curve (\square) after the processor is in the development (or shipment) phase is not an option. One way to reduce the processor power is to run it at a lower power supply and frequency. This leads to a shift of the power-performance point down the hardware tuning curve (\circ) that originates at the initial 12FO4 power-performance point until it intersects the power budget line. As we can see, eventually such a processor will deliver much lower performance than could be achieved if the depth of the pipeline had been chosen according to the BIPS-cube per watt optimum. Notice that, if such a processor is set to run at a lower than maximum frequency without reducing the power supply, then the power-performance design point will shift down the frequency scaling curve marked with " \triangleright ," which will lead to even worse performance loss in the final product. Also, recently reported execution and fetch throttling techniques to restrain the maximum power, if activated frequently on useful applications, lead to the shift of the design power-performance point down a similar curve.

9 CONCLUSIONS

In this paper, we have demonstrated that it is important to consider both power and performance while optimizing pipelines. For this purpose, we derived detailed energy models using circuit-extracted power analysis for micro-architectural structures. We also developed detailed equations for how the energy functions scale with pipeline depth. Based on the combination of power and performance modeling performed, our results show that a purely performance-driven, power-unaware design may lead to the selection of an overly deep pipelined microprocessor operating at an inherently power-inefficient design point.

As this work is the first quantitative evaluation of power and performance optimal pipelines, we also performed a detailed sensitivity analysis of the optimal pipeline depth against key parameters such as latch growth factor, latch ratio, latch insertion delay, glitch, and leakage currents. Our analysis shows that there is a range of pipeline depth for which performance increases can be achieved at a modest sacrifice in power-performance efficiency. Pipelining beyond that range leads to drastic reduction in power-performance efficiency with little or no further performance improvement.

Our results show that, for a current generation, out-of-order superscalar processor, the optimal delay per stage is about 18 FO4 (consisting of a logic delay of 15 FO4 and 3 FO4 latch insertion delay) when the objective is to maximize net performance subject to a power constraint; this is in contrast to an optimal delay of 10 FO4/stage when considering the BIPS metric alone. We used a broad suite of SPEC2000 benchmarks to arrive at this conclusion.

The optimal pipeline depth depends on a number of parameters in the power models which we have derived from current state-of-the-art microprocessor design methodologies, e.g., higher leakage parameters for a process may favor deeper pipelines as area and, hence, leakage scales linearly with transistor width and not quadratically in the number of latches. Also, as already established through recent prior work, such optimal design points generally depend on the input workload characteristics. Our simulation-based experiments on a typical commercial application (TPC-C) show that, although the optimal pipeline depth is around 10-14 FO4 for performance-only optimization, it increases to 24-28 FO4 when we consider power and performance optimizations.

REFERENCES

- [1] D. Brooks, P. Bose, S.E. Schuster, P.N. Kudva, H. Jacobson, A. Buyuktosunoglu, J.D. Wellman, V. Zyuban, M. Gupta, and P.W. Cook, "Power-Aware Microarchitecture: Design and Modeling Challenges for the Next-Generation Microprocessors," *IEEE Micro*, vol. 20, no. 6, pp. 26-44, Nov./Dec. 2000.
- [2] V. Srinivasan, D. Brooks, M. Gschwind, P. Bose, V. Zyuban, P. Strenski, and P. Emma, "Optimizing Pipelines for Power And Performance," *Proc. ACM/IEEE 35th Int'l Symp. Microarchitecture*, Nov. 2002.
- [3] S.R. Kunkel and J.E. Smith, "Optimal Pipelining in Supercomputers," *Proc. 13th Int'l Symp. Computer Architecture (ISCA-13)*, pp. 404-411, June 1986.
- [4] P. Dubey and M. Flynn, "Optimal Pipelining," *J. Parallel and Distributed Computing*, vol. 8, pp. 10-19, 1990.
- [5] A. Hartstein and T.R. Puzak, "The Optimum Pipeline Depth for a Microprocessor," *Proc. 29th Int'l Symp. Computer Architecture (ISCA-29)*, May 2002.
- [6] M.S. Hrishikesh, K. Farkas, N.P. Jouppi, D.C. Burger, S.W. Keckler, and P. Sivakumar, "The Optimal Logic Depth per Pipeline Stage Is 6 to 8 FO4 Inverter Delays," *Proc. 29th Int'l Symp. Computer Architecture (ISCA-29)*, pp. 14-24, May 2002.
- [7] E. Sprangle and D. Carmean, "Increasing Processor Performance by Implementing Deeper Pipelines," *Proc. 29th Int'l Symp. Computer Architecture (ISCA-29)*, May 2002.
- [8] P. Kogge, *The Architecture of Pipelined Computers*. Hemisphere Publishing, 1981.
- [9] M. Moudgill, P. Bose, and J.H. Moreno, "Validation of Turandot, a Fast Processor Model for Microarchitecture Exploration," *Proc. IEEE Int'l Performance, Computing, and Comm. Conf. (IPCCC)*, pp. 451-457, Feb. 1999.
- [10] M. Moudgill, J.D. Wellman, and J.H. Moreno, "Environment for PowerPC Microarchitecture Exploration," *IEEE Micro*, vol. 19, no. 3, pp. 9-14, May/June 1999.
- [11] V. Zyuban and P. Strenski, "Balancing Hardware Intensity in Microprocessor Pipelines," *IBM J. Research and Development*, vol. 47, nos. 5/6, pp. 585-598, 2003.
- [12] D. Brooks, P. Bose, V. Srinivasan, M. Gschwind, P. Emma, and M. Rosenfield, "Microarchitecture-Level Power-Performance Analysis: The Powertimer Approach," *IBM J. Research and Development*, vol. 47, nos. 5/6, pp. 653-670, 2003.
- [13] D. Brooks, V. Tiwari, and M. Martonosi, "Wattch: A Framework for Architectural-Level Power Analysis and Optimizations," *Proc. 27th Ann. Int'l Symp. Computer Architecture (ISCA-27)*, June 2000.
- [14] N. Vijaykrishnan, M. Kandemir, M.J. Irwin, H.S. Kim, and W. Ye, "Energy-Driven Integrated Hardware-Software Optimizations Using SimplePower," *Proc. 27th Ann. Int'l Symp. Computer Architecture*, June 2000.

- [15] V. Zyuban, "Inherently Lower Power High Performance Superscalar Architectures," PhD thesis, Univ. of Notre Dame, Mar. 2000.
- [16] J.S. Neely, H.H. Chen, S.G. Walker, J. Venuto, and T. Bucelot, "CPAM: A Common Power Analysis Methodology for High-Performance VLSI Design," *Proc. Ninth Topical Meeting Electrical Performance of Electronic Packaging*, pp. 303-306, 2000.
- [17] V. Zyuban, "Optimization of Scannable Latches for Low Energy," *IEEE Trans. VLSI Systems*, vol. 11, pp. 778-787, Oct. 2003.
- [18] V. Iyengar, L.H. Trevillyan, and P. Bose, "Representative Traces for Processor Models with Infinite Cache," *Proc. Second Symp. High Performance Computer Architecture (HPCA-2)*, Feb. 1996.
- [19] M.J. Flynn, P. Hung, and K.W. Rudd, "Deep-Submicron Microprocessor Design Issues," *IEEE Micro*, vol. 19, no. 4, pp. 11-22, July/Aug. 1999.
- [20] R. Gonzalez and M. Horowitz, "Energy Dissipation in General Purpose Microprocessors," *IEEE J. Solid-State Circuits*, vol. 31, no. 9, pp. 1277-1284, Sept. 1996.
- [21] R. Jessani and C. Olson, "The Floating-Point Unit of the PowerPC 603e Microprocessor," *IBM J. Research and Development*, vol. 40, no. 5, pp. 559-566, Sept. 1996.
- [22] S. Heo, R. Krashinsky, and K. Asanovic, "Activity-Sensitive Flip-Flop and Latch Selection for Reduce Energy," *Proc. 19th Conf. Advanced Research in VLSI*, Mar. 2001.
- [23] V.G. Oklobdzija, "Clocking and Storage Elements in a Multi-Gigahertz Environment," *IBM J. Research and Development*, vol. 47, nos. 5/6, pp. 567-584, 2003.



Victor Zyuban received the BS (1993) and MS (1995) degrees from the Moscow Institute of Physics and Technology and the PhD degree in computer science and engineering from the University of Notre Dame in 2000. From 1995 to 1996, he worked in the Moscow Center for SPARC Technologies. He is currently a research staff member at the IBM T.J. Watson Research Center. From 2000 to 2003, he worked on a low-power DSP project in which he was involved in ISA, microarchitecture, and physical design, and was leading the development of a semicustom core. Currently, he is working on a high-performance microprocessor core. His research interests include high-performance microarchitecture, circuits, and methodologies for energy-efficient design. He is a member of the IEEE.



hardware and software for embedded and high-performance computer systems. He is a member of the IEEE.

David Brooks received the BS (1997) degree from the University of Southern California and the MA (1999) and PhD (2001) degrees from Princeton University, all in electrical engineering. He is an assistant professor of computer science at Harvard University. Prior to joining Harvard University, he was a research staff member at the IBM T.J. Watson Research Center. His research interests include architectural-level power-modeling and power-efficient design of

hardware and software for embedded and high-performance computer systems. He is a member of the IEEE.



Viji Srinivasan received the BS degree in physics from the University of Madras in 1990, the MS degree in computer science and engineering from the Indian Institute of Science in 1994, and the PhD degree in computer science and engineering from the University of Michigan in 2001. Currently, she is with the IBM Research Division, T.J. Watson Research Center as a research staff member in the Computer Architecture Department. Her research interests

include computer architecture and performance analysis. She is a member of the IEEE.



Michael Gschwind received the MS and PhD degrees in computer science from the Technische Universität Wien in 1991 and 1996, respectively. He is a research staff member at the IBM T.J. Watson Research Center. He was one of the originators of the supercomputer-on-a-chip "Cell" next-generation high-performance system architecture which is currently being developed in a joint development center by IBM, Sony, and Toshiba in Austin, Texas.

Currently, he leads a project focusing on high-frequency low-power high-performance architectures for media, high-performance, and general purpose computing applications. He is also active in the evolution of the PowerPC architecture, focusing on high performance future PowerPC subsystems. He also holds a visiting faculty appointment at Princeton University. Before joining IBM in 1997, he was a faculty member in the Department of Computer Engineering, Technische Universität Wien, Vienna, Austria. His research interests include compilers, computer architecture, instruction-level parallelism, hardware/software codesign, application-specific processors, and field-programmable gate arrays. He is the author of more than 60 papers and holds several patents on high-performance computer architecture. He is the recipient of eight IBM Invention Plateaus. He is a senior member of the IEEE.



Pradip Bose received the BTech (Hons.) degree from I.I.T Kharagpur, India, in 1977 and the MS and PhD degrees in electrical and computer engineering from the University of Illinois at Urbana-Champaign in 1981 and 1983, respectively. He joined the IBM T.J. Watson Research Center in 1983, where he is currently a research staff member and manager of a group on reliability and power-aware microarchitectures. He has been involved in research and develop-

ment work in support of all the POWER-series, high-end microprocessors marketed by IBM. His current research interests include high performance computer architectures, power, complexity, and reliability-aware design, and power-performance modeling of systems. He is the current editor-in-chief of *IEEE Micro* magazine and is a senior member of the IEEE and the IEEE Computer Society.



ment work in support of all the POWER-series, high-end microprocessors marketed by IBM. His current research interests include high performance computer architectures, power, complexity, and reliability-aware design, and power-performance modeling of systems. He is the current editor-in-chief of *IEEE Micro* magazine and is a senior member of the IEEE and the IEEE Computer Society.

Philip N. Strenski received the AB degree in mathematics in 1979 from Washington University and the PhD degree in physics in 1985 from Stanford University. That same year he joined IBM, engaging in research in design automation, including simulated annealing, parametric extraction, and circuit tuning. In the early 1990s, he switched to microprocessor circuit design, over the years contributing to CPU designs in the mainframe, midrange, and workstation businesses. He is currently involved in the design of the "Cell" system, a collaboration among Sony, Toshiba, and IBM; he manages the Advanced RISC Circuit Design Department; he is also actively involved in issues of power-performance efficiency. He is a member of the IEEE.



Philip G. Emma received the PhD degree in electrical engineering from the University of Illinois in 1983, after which he joined the IBM T.J. Watson Research Center. For 10 years, he did research in architecture, microarchitecture, and systems design, and invented many of the cache and pipeline mechanisms that are in pervasive use today. He did circuit design work for several years and was a design leader on IBM's first CMOS mainframe processors, where

he defined and specified the processor-level RAS features used in the z-series servers. He then led a systems technology group working in the areas of advanced packaging and optical interconnection with a server-level focus. His current group studies advanced technology in light of systems and is working to understand design optimality and to reduce power and increase reliability through microarchitectural and circuit techniques. He holds nearly 100 patents in these areas and received the honorary title of "IBM Master Inventor." He was elected to the IBM Academy of Technology and is a fellow of the IEEE.