

Figure 10: Optimum voltage of operation.

Voltage	Parallel area/power	Parallel- Pipeline area/power	IIR area/power
5	1 / 1	1 / 1	1 / 1
2	6 / 0.19	3.7 / 0.2	2.6 / 0.23
1.5	11 / 0.13	7 / 0.12	7 / 0.14
1.4	15 / 0.14	10 / 0.11	Recursive bottleneck reached

Table 3: Normalized Area/Power for various supply voltages for Plots 2,3,4 in Figure 10.

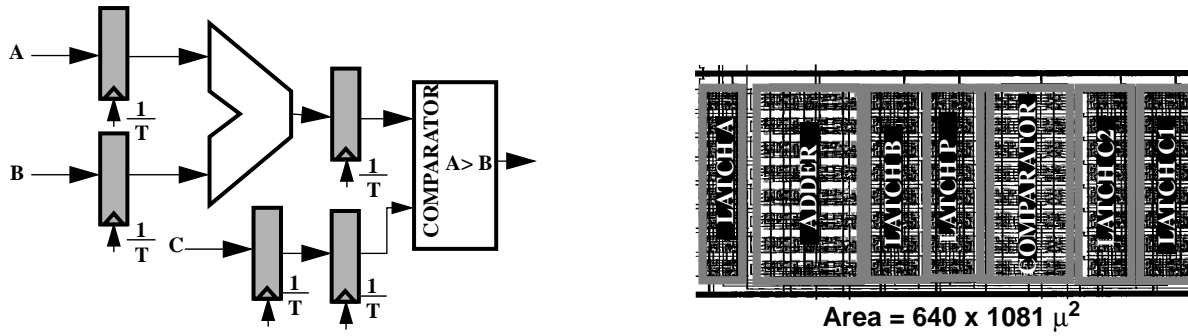


Figure 9: Pipelined implementation of the simple datapath.

Architecture type	Voltage	Area	Power
Simple datapath (no pipelining or parallelism)	5V	1	1
Pipelined datapath	2.9V	1.3	0.39
Parallel datapath	2.9V	3.4	0.36
Pipeline-Parallel	2.0V	3.7	0.2

Table 2: Architecture summary.

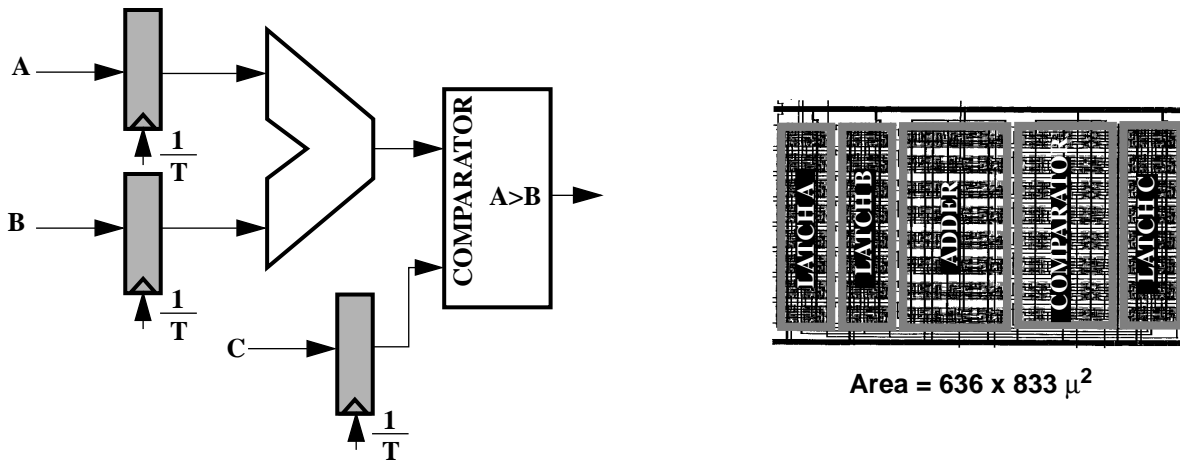


Figure 7: A simple datapath with corresponding layout.

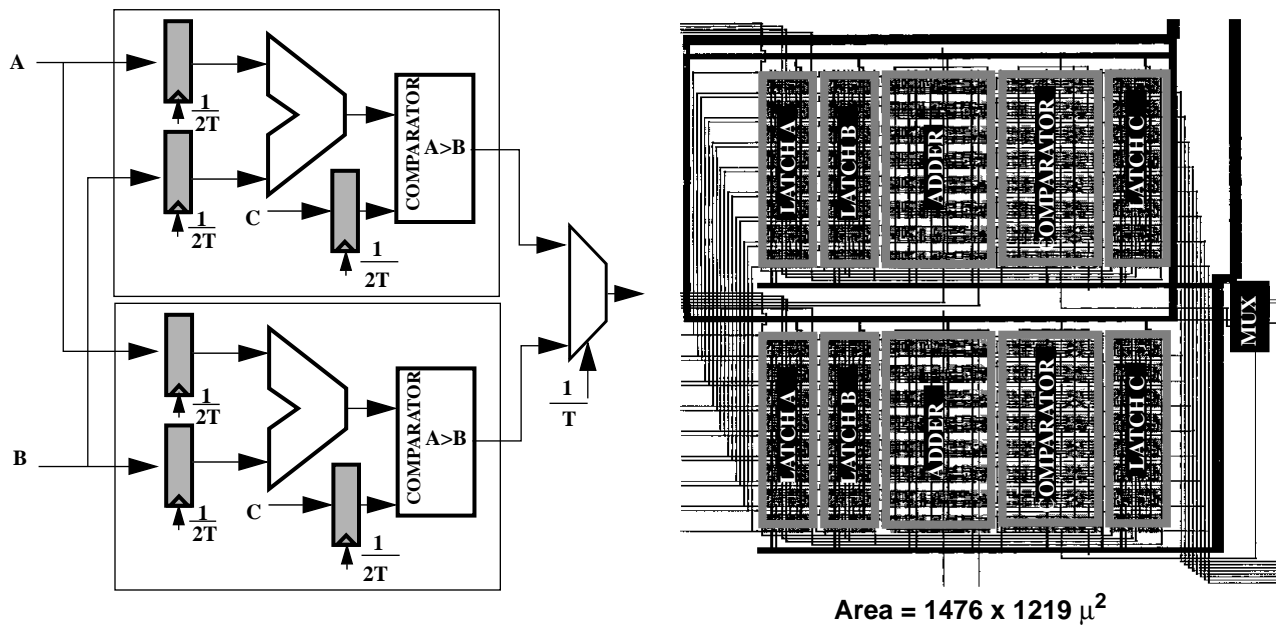


Figure 8: Parallel implementation of the simple datapath.

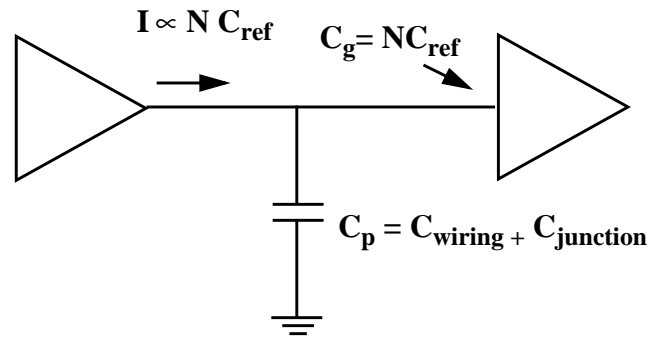


Figure 5: Circuit model for analyzing the effect of transistor sizing.

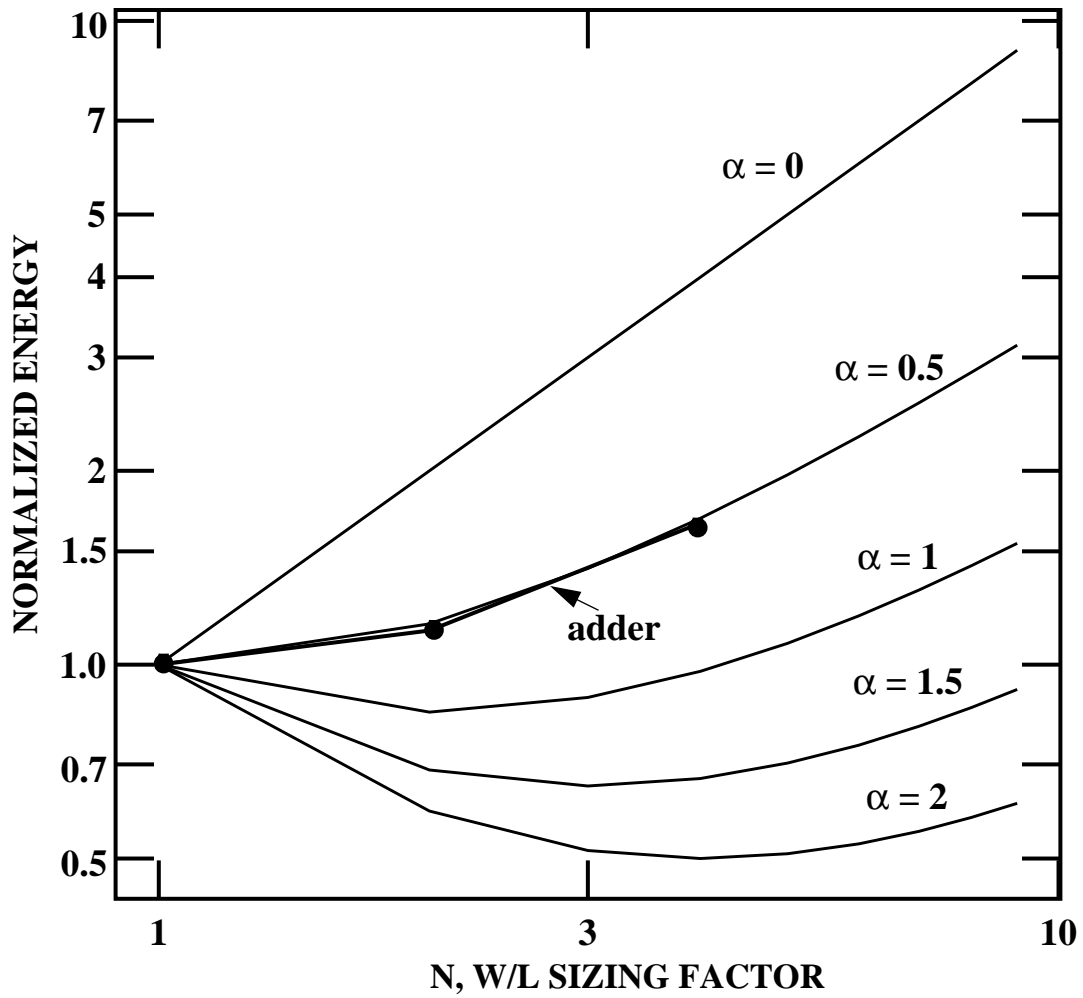


Figure 6: Plot of energy vs. transistor sizing factor for various parasitic contributions.

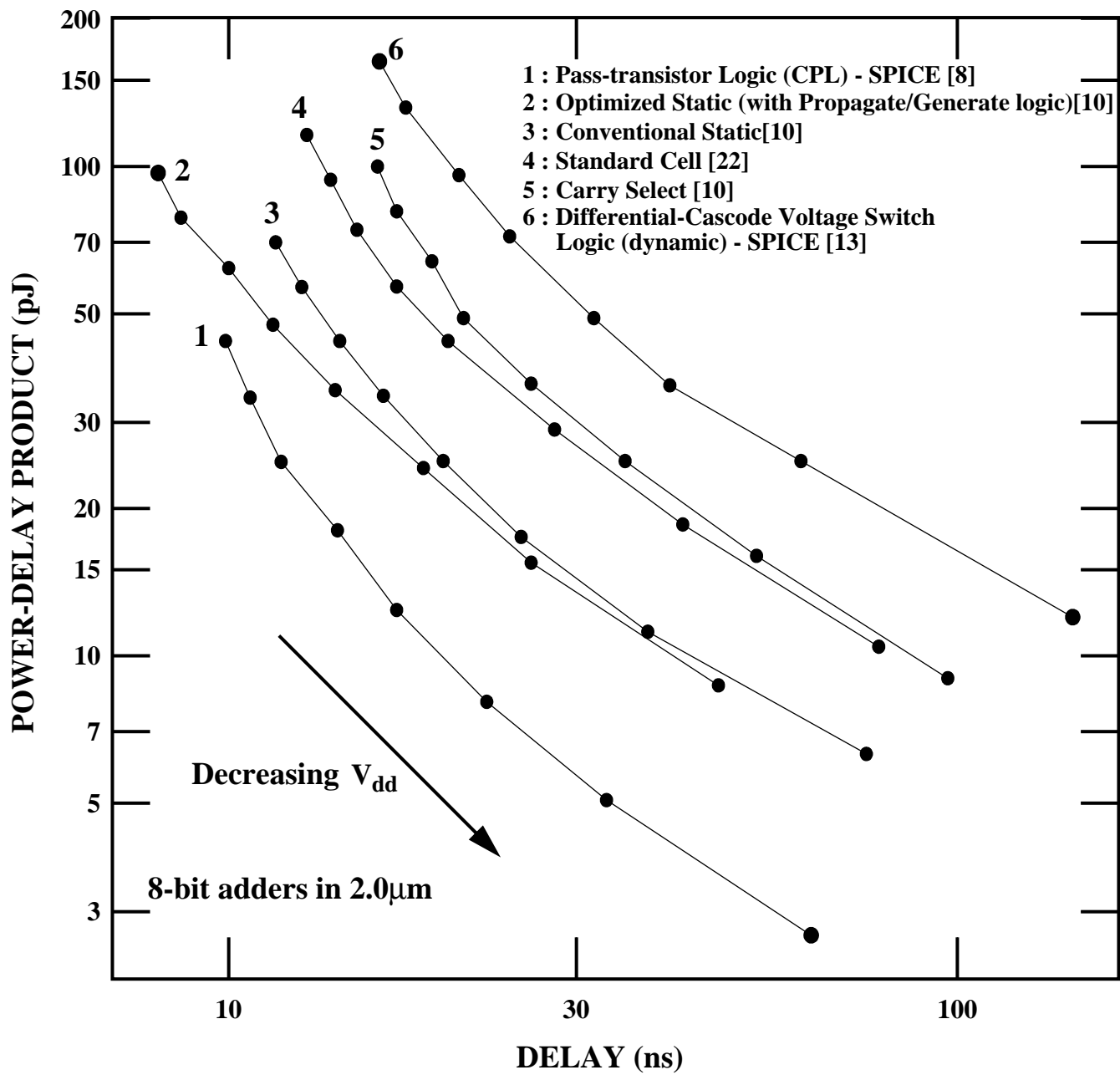


Figure 4: Data showing improvement in power-delay product at the cost of speed for various circuit approaches.

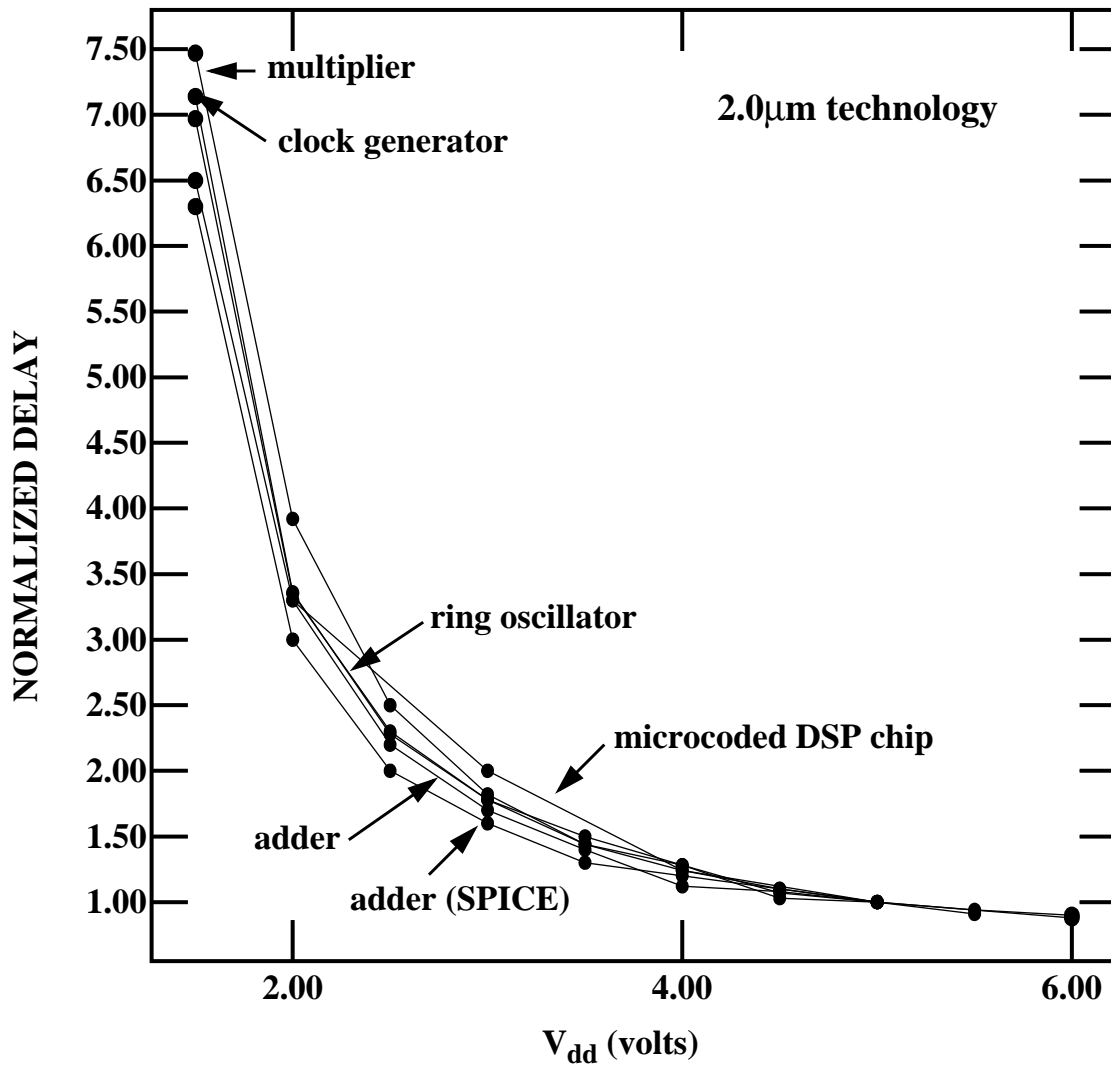


Figure 3: Data demonstrating delay characteristics follow simple first order theory.

Component (all in 2µm)	# of transistors	Area	Comments
Microcoded DSP Chip [21]	44802	94mm ²	20-bit datapath
Multiplier	20432	12.2mm ²	24x24 bits
Adder	256	0.083mm ²	conventional static
Ring Oscillator	102	0.055mm ²	51-stages
Clock Generator	56	0.04mm ²	cross-coupled NOR

Table 1: Details of components used for the study in Figure 3.

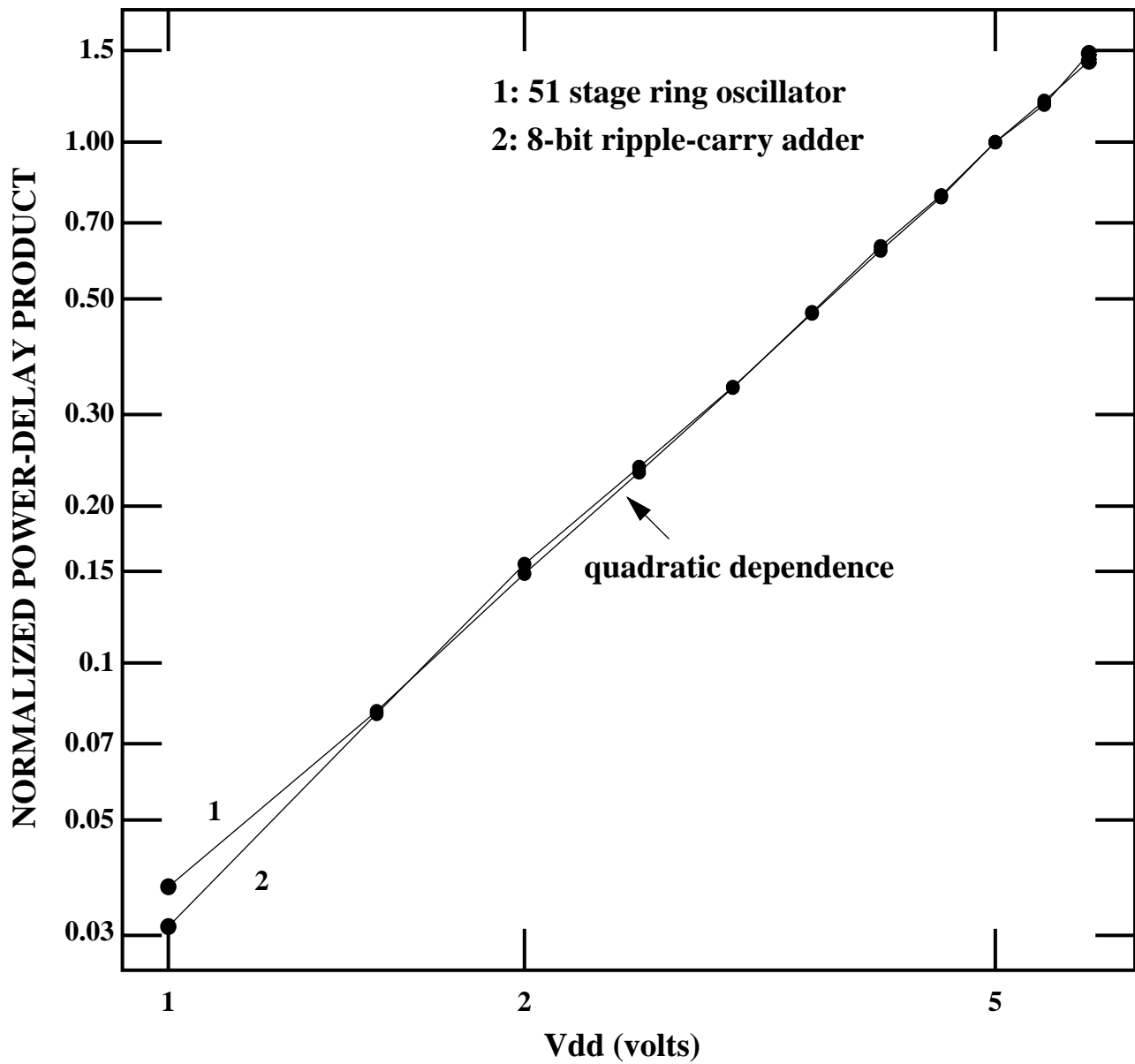
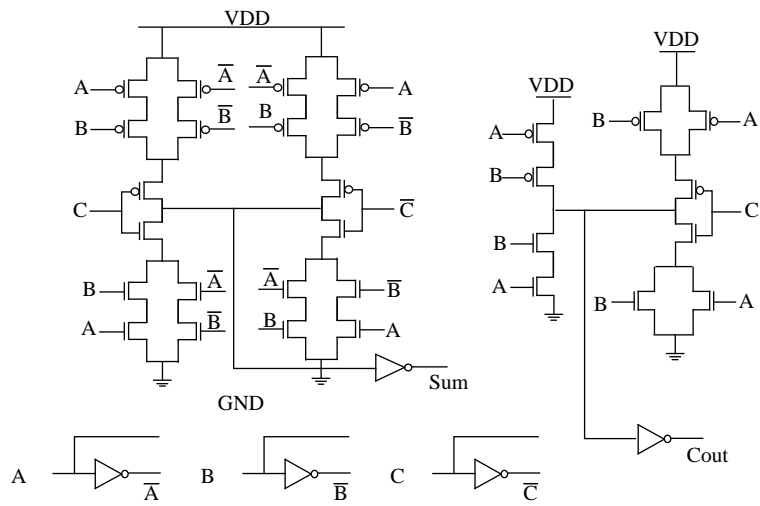
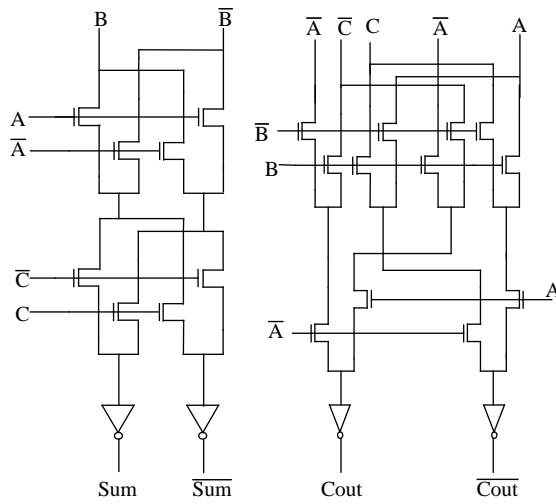


Figure 2: Power-delay product exhibiting square law dependence for two different circuits.



Transistor count (conventional CMOS) : 40



Transistor count (CPL) : 28

Figure 1: Comparison of a Conventional CMOS and CPL adders [8].

Conclusions

- [6] G. Geschwind and R.M. Clary, "Multichip Modules- An Overview", Expo SMT/HiDEP 1990 Technical Proceedings, San Jose, CA, 1990.
- [7] D.A. Patterson and C.H. Sequin, "Design Considerations for Single-Chip Computers of the Future", Joint Special Issue, IEEE Journal of Solid-State Circuits SC-15, No.1 and IEEE Transactions on Computers C-29, No.2, pp. 108-116, Feb. 1980.
- [8] K. Yano, et al., "A 3.8ns CMOS 16x16 Multiplier Using Complementary Pass Transistor Logic", IEEE Journal of Solid-State Circuits, pp. 388-395, April 1990.
- [9] H.J.M. Veendrick, "Short-Circuit Dissipation of Static CMOS Circuitry and Its Impact on the Design of Buffer Circuits", IEEE Journal of Solid-State Circuits, Vol. SC-19, pp. 468-473, August 1984.
- [10] N. Weste and K. Eshragian, Principles of CMOS VLSI Design: A Systems Perspective, Addison-Wesley, MA, 1988.
- [11] R.K. Watts (ed.), Submicron Integrated Circuits, John Wiley & Sons, NY, 1989.
- [12] D. Green, Modern Logic Design, Addison-Wesley, pp. 15-17, 1986.
- [13] G. Jacobs and R.W. Brodersen, "A Fully Asynchronous Digital Signal Processor Using Self-Timed Circuits", IEEE Journal of Solid-State Circuits, pp. 1526-1537, December 1990.
- [14] D. Hodges and H. Jackson, Analysis and Design of Digital Integrated Circuits, McGraw-Hill, Inc., 1988.
- [15] M. Shoji, CMOS Digital Circuit Technology, Prentice-Hall, 1988.
- [16] S. Sze, Physics of Semiconductor Devices, John Wiley & Sons, 1981.
- [17] M. Aoki et al., "0.1 μ m CMOS Devices Using Low-Impurity-Channel Transistors (LICT)", IEDM, pp. 939-941, 1990.
- [18] M. Nagata, "Limitations, Innovations, and Challenges of Circuits & Devices into Half-Micron and Beyond", Symposium on VLSI circuits, pp. 39-42, 1991.
- [19] S.C. Ellis, "Power Management in Notebook PC's", Silicon Valley Personal Computer Conference, pp. 749-754, 1991.
- [20] K. Chu and D. Pulfrey, "A Comparison of CMOS Circuit Techniques: Differential Cascode Voltage Switch Logic Versus Conventional Logic", IEEE Journal of Solid-State Circuits, pp. 528-532, August 1987.
- [21] R.W. Brodersen, Lager: A Silicon Compiler, to be published, Kluwer.
- [22] R.W. Brodersen et al., LagerIV Cell Library Documentation, Electronics Research Laboratory, University of California, Berkeley, June 23, 1988.
- [23] B. Davari et al., "A High Performance 0.25 μ m CMOS Technology", IEDM, pp. 56-59, 1988.
- [24] M. Kakumu and M. Kinugawa, "Power-Supply Voltage Impact on Circuit Performance for Half and Lower Submicrometer CMOS LSI", IEEE Transactions on Electron Devices, Vol 37, No. 8, pp. 1902-1908, August 1990.
- [25] D. Dahle, "Designing High Performance Systems to Run from 3.3V or Lower Sources", Silicon Valley Personal Computer Conference, pp. 685-691, 1991.
- [26] R. Swanson and J. Meindl, "Ion-Implanted Complementary MOS Transistors in Low-Voltage Circuits", IEEE Journal of Solid-State Circuits, pp. 146-153, April 1972.
- [27] D. Messerschmitt: "Breaking the Recursive Bottleneck", Skwirzynski (editor), Performance Limits in Communication Theory and Practice, pp. 3-19, 1988.
- [28] M. Potkonjak and J. Rabaey, "Optimizing Resource Utilization using Transformations", ICCAD, pp. 88-91, November 1991.

With the continuing trend of denser technology through scaling and the development of advanced packaging techniques, a new degree of freedom in architectural design has been made possible in which silicon area can be traded off against power consumption. Parallel architectures, utilizing pipelining or hardware replication, provide the mechanism for this trade-off by maintaining throughput while using slower device speeds and thus allowing reduced voltage operation. The well-behaved nature of the dependencies of power dissipation and delay as a function of supply voltage over a wide variety of situations allows optimizations of the architecture. In this way, for a wide variety of situations, the optimum voltage was found to be less than 1.5 volts, below which the overhead associated with the increased parallelism becomes prohibitive.

There are other limitations which may not allow the optimum supply voltage to be achieved. The algorithm that is being implemented may be sequential in nature and/or have feedback which will limit the degree of parallelism which can be exploited. Another possibility is that the optimum degree of parallelism may be so large that the number of transistors may be inordinately large, thus making the optimum solution unreasonable. However, in any case, the goal in minimizing power consumption is clear; operate the circuit as slowly as possible, with the lowest possible supply voltage.

Acknowledgments

This project was sponsored by DARPA. We wish to thank Andy Burstein, Miodrag Potkonjak, Mani Srivastava, Anton Stoelzle, and Lars Thon for providing us with examples. We also would like to thank professors Ping Ko, Teresa Meng, Jan Rabaey, and Charles Sodini for their invaluable feedback. Lastly, we would also wish to acknowledge the support of Sam Sheng by the Fannie and John Hertz Foundation.

References

- [1] T. Bell, "Incredible shrinking computers", IEEE Spectrum, pp. 37-43, May 1991.
- [2] J. Eager, "Advances in Rechargeable Batteries Pace Portable Computer Growth", Silicon Valley Personal Computer Conference, pp. 693-697, 1991.
- [3] A. Chandrakasan, S. Sheng, R.W. Brodersen, "Design Considerations for a Future Multimedia Terminal", 1990 WIN-LAB Workshop, Rutgers University, New Jersey, October, 1990.
- [4] H.B. Bakoglu, Circuits, Interconnections, and Packaging for VLSI, Addison-Wesley, Menlo Park, CA, 1990.
- [5] D. Benson, Y. Bobra, B. McWilliams, et al., "Silicon Multichip Modules", Hot Chips Symposium III, Santa Clara, CA, August 1990.

figure represents the power dissipation which would be achieved if there were no overhead associated with increased parallelism. For this case, the power is a strictly decreasing function of V_{dd} and the optimum voltage would be set by the minimum value allowed from noise margin constraints (assuming that no recursive bottleneck was reached). Curve 5 assumes that the inter-processor capacitance has an N^2 dependence while Curve 6 assumes an N^3 dependence. It is expected that in most practical cases the dependence is actually less than N^2 , but even with the extremely strong N^3 dependence an optimal value around 2V is found.

Curves 2 and 3 are obtained from data from actual layouts[22], and exhibit a dependence of the interface capacitance which lies between linear and quadratic on the degree of parallelism, N . For these cases, there was no interprocessor communication. Curves 2 and 3 are extensions of the example described in Section 4.5 in which the parallel and parallel-pipeline implementations of the simple datapath were duplicated N times. Curve 4 is for a much more complex example, a 7th order IIR filter, also obtained from actual layout data. The overhead in this case arose primarily from interprocessor communication. This curve terminates around 1.4V, because at that point the algorithm has been made maximally parallel, reaching a recursive bottleneck. For this case, at a supply voltage of 5V, the architecture is basically a single hardware unit that is being time-multiplexed and requires about 7 times more power than the optimal parallel case which is achieved with a supply of around 1.5V. In Table 3 is a summary of the power reduction and normalized areas that were obtained from layouts. The increase in areas gives an indication of the amount of parallelism being exploited. The key point is that the optimal voltage was found to be relatively independent over all the cases considered, and occurred around 1.5V for the 2.0μ technology; a similar analysis using a 0.5V threshold, 0.8μ process (with an L_{eff} of 0.5μ) resulted in optimal voltages around 1V, with power reductions in excess of a factor of 10. Further scaling of the threshold would allow even lower voltage operation, and hence even greater power savings.

6.0 Conclusions

There are a variety of considerations that must be taken into account in low power design which include the style of logic, the technology used and the logic implemented. Factors that were shown to contribute to power dissipation included spurious transitions due to hazards and critical race conditions, leakage and direct path currents, pre-charge transitions and power consuming transitions in unused circuitry. A passgate logic family with modified threshold voltages was found to be the best performer for low-power designs, due to the minimal number of transistors required to implement the important logic functions. An analysis of transistor sizing has been shown that minimum sized transistors should be used if the parasitic capacitances are less than the active gate capacitances in a cascade of logic gates.

$$P_{normalized} = \left(1 + \frac{C_{ip}(N)}{NC_{ref}} + \frac{C_{interface}(N)}{C_{ref}} \right) \left(\frac{V}{V_{ref}} \right)^2 \quad (\text{EQ 12})$$

At very low supply voltages (near the device thresholds), the number of processors (and hence the corresponding overhead in the above equation) typically increases at a faster rate than the V^2 term decreases, resulting in a power increase with further reduction in voltage.

Reduced threshold devices tends to lower the optimal voltage; however, as seen in Section 3.3, at thresholds below 0.2V, power dissipation due to the subthreshold current will soon start to dominate and limit further power improvement. An even lower bound on the power supply voltage for a CMOS inverter with “correct” functionality was found to be 0.2V (8kT/q) [26]. This gives a limit on the power-delay product that can be achieved with CMOS digital circuits; however, the amount of parallelism to retain throughput at this voltage level would no doubt be prohibitive for any practical situation.

So far, we have seen that parallel and pipelined architectures can allow for a reduction in supply voltages to the “optimal” level; this will indeed be the case if the algorithm being implemented does not display any recursion (feedback). However, there are a wide class of applications inherently recursive in nature, ranging from the simple ones, such as infinite impulse response and adaptive filters, to more complex cases such as systems solving non-linear equations and adaptive compression algorithms. There is therefore also an algorithmic bound on the level to which pipelining and parallelism can be exploited for voltage reduction. Although the application of data control flowgraph transformations can alleviate this bottleneck to some extent, both the constraints on latency and the structure of computation of some algorithms can prevent voltage reduction to the optimal voltage levels discussed above[27,28].

Another constraint on the lowest allowable supply voltages is set by system noise margin constraints ($V_{noise\ margin}$). Thus, we must lower-bound the “optimal” voltage by:

$$V_{noise\ margin} \leq V_{optimal} \leq V_{critical} \quad (\text{EQ 13})$$

with $V_{critical}$ defined in Section 4.4. Hence, the “optimal” supply voltage (for a fixed technology) will lie somewhere between the voltage set by noise margin constraints and the critical voltage.

Figure 10 shows power (normalized to 1 at $V_{dd}=5V$) as a function of V_{dd} for a variety of cases for a 2.0 μ technology. As will be shown there is a wide variety of assumptions in these various cases and it is important to note that the all have roughly the same optimum value of supply voltage, approximately 1.5V. Curve 1 in this

the level of pipelining also has the effect of reducing logic depth and hence power contributed due to hazards and critical races (see section 3.1.1).

Furthermore, an obvious extension is to utilize a combination of pipelining and parallelism. Since this architecture reduces the critical path and hence speed requirement by a factor of 4, the voltage can be dropped until the delay increases by a factor of 4. The power consumption in this case is:

$$P_{parpipe} = C_{parpipe} V_{parpipe}^2 f_{parpipe} = (2.5C_{ref})(0.4V_{ref})^2 \left(\frac{f_{ref}}{2}\right) \approx 0.2P_{ref} \quad (\text{EQ 10})$$

The parallel-pipeline implementation results in a 5 times reduction in power. Table 2 shows a comparative summary of the various architectures described for the simple adder-comparator datapath.

From the above examples, it is clear that the traditional time-multiplexed architectures, as used in general purpose microprocessors and DSP chips are the least desirable for low-power applications. This follows since time multiplexing actually increases the speed requirements on the logic circuitry, thus not allowing reduction in the supply voltage.

5.0 Optimal Supply Voltage

In the previous section, we saw that the delay increase due to reduced supply voltages below the critical voltage can be compensated by exploiting parallel architectures. However, as seen in Figure 3 and Equation 3, as supply voltages approach the device thresholds, the gate delays increase rapidly. Correspondingly, the amount of parallelism and overhead circuitry increases to a point where the added overhead dominates any gains in power reduction from further voltage reduction, leading to the existence of an “optimal” voltage from an architectural point of view. To determine the value of this voltage, the following model is used for the power for a fixed system throughput as a function of voltage (and hence degree of parallelism):

$$Power(N) = NC_{ref} V^2 \frac{f_{ref}}{N} + C_{ip} V^2 \frac{f_{ref}}{N} + C_{interface} V^2 f_{ref} \quad (\text{EQ 11})$$

where N is the number of parallel processors, C_{ref} is the capacitance of a single processor, C_{ip} is the inter-processor communication overhead introduced due to the parallelism (due to control and routing), and $C_{interface}$ is the overhead introduced at the interface which is not decreased in speed as the architecture is made more parallel. In general, C_{ip} and $C_{interface}$ are functions of N , and the power improvement over the reference case (i.e. without parallelism) can be expressed as:

One way to maintain throughput while reducing the supply voltage is to utilize a parallel architecture. As shown in Figure 8, two identical adder-comparator datapaths are used, allowing each unit to work at half the original rate while maintaining the original throughput. Since the speed requirements for the adder, comparator, and latch have decreased from 25ns to 50ns, the voltage can be dropped from 5V to 2.9V (the voltage at which the delay doubled, from Figure 3). While the datapath capacitance has increased by a factor of 2, the operating frequency has correspondingly decreased by a factor of 2. Unfortunately, there is also a slight increase in the total “effective” capacitance introduced due to the extra routing, resulting in an increased capacitance by a factor of 2.15. Thus the power for the parallel datapath is given by:

$$P_{par} = C_{par} V_{par}^2 f_{par} = (2.15C_{ref})(0.58V_{ref})^2 \left(\frac{f_{ref}}{2}\right) \approx 0.36P_{ref} \quad (\text{EQ 8})$$

This method of reducing power by using parallelism has the overhead of increased area, and would not be suitable for area-constrained designs. In general, parallelism will have the overhead of extra routing (and hence extra power), and careful optimization must be performed to minimize this overhead (for example, partitioning techniques for minimal overhead). Interconnect capacitance will especially play a very important role in deep sub-micron implementations, since the fringing capacitance of the interconnect capacitance ($C_{wiring} = C_{area} + C_{fringing} + C_{wiring}$) can become a dominant part of the total capacitance (equal to $C_{gate} + C_{junction} + C_{wiring}$) and cease to scale [4].

Another possible approach is to apply pipelining to the architecture, as shown in Figure 9. With the additional pipeline latch, the critical path becomes the $\max[T_{adder}, T_{comparator}]$, allowing the adder and the comparator to operate at a slower rate. For this example, the two delays are equal, allowing the supply voltage to again be reduced from 5V used in the reference datapath to 2.9V (the voltage at which the delay doubles) with no loss in throughput. However, there is a much lower area overhead incurred by this technique, as we only need to add pipeline registers. Note that there is again a slight increase in hardware due to the extra latches, increasing the “effective” capacitance by approximately a factor of 1.15. The power consumed by the pipelined datapath is:

$$P_{pipe} = C_{pipe} V_{pipe}^2 f_{pipe} = (1.15C_{ref})(0.58V_{ref})^2 f_{ref} \approx 0.39 P_{ref} \quad (\text{EQ 9})$$

With this architecture, the power reduces by a factor of approximately 2.5, providing approximately the same power reduction as the parallel case with the advantage of lower area overhead. As an added bonus, increasing

very little penalty in speed performance. This implies that there is little advantage to operating above a certain voltage. This idea has been formalized by Kakumu et. al., yielding the concept of a “critical voltage” which provides a lower limit on the supply voltage [24]. The critical voltage is defined as $V_c = 1.1E_cL_{eff}$, where E_c is the critical electric field causing velocity saturation; this is the voltage at which the delay vs. V_{dd} curve approaches a $\sqrt{V_{dd}}$ dependence. For 0.3μ technology, the proposed lower limit on supply voltage (or the critical voltage) was found to be 2.43V.

Because of this effect, there is some movement to a 3.3V industrial voltage standard since at this level of voltage reduction there is not a significant loss of circuit speed[1,25]. This was found to achieve a 60% reduction in power when compared to a 5 volt operation[25].

4.5 Architecture-Driven Voltage Scaling

The above mentioned “technology” based approaches are focusing on reducing the voltage while maintaining device speed, and are not attempting to achieve the minimum possible power. As shown in Figures 2 and 4, CMOS logic gates achieve lower power-delay products (energy per computation) as the supply voltages are reduced. In fact, once a device is in velocity saturation there is a further degradation in the energy per computation, so in minimizing the energy required for computation, Kakumu’s critical voltage provides an *upper* bound on the supply voltage (whereas for his analysis it provided a *lower* bound!). It now will be the task of the architecture to compensate for the reduced circuit speed, that comes with operating below the critical voltage.

To illustrate how architectural techniques can be used to compensate for reduced speeds, a simple 8-bit datapath consisting of an adder and a comparator is analyzed assuming a $2.0\mu\text{m}$ technology. As shown in Figure 7, inputs A and B are added, and the result compared to input C. Assuming the worst-case delay through the adder, comparator and latch is approximately 25ns at a supply voltage of 5V, the system in the best case can be clocked with a clock period of $T = 25\text{ns}$. When required to run at this maximum possible throughput, it is clear that the operating voltage cannot be reduced any further since no extra delay can be tolerated, hence yielding no reduction in power. We will use this as the reference datapath for our architectural study and present power improvement numbers with respect to this reference. The power for the reference datapath is given by:

$$P_{ref} = C_{ref} V_{ref}^2 f_{ref} \quad (\text{EQ 7})$$

where C_{ref} is the total effective capacitance being switched per clock cycle. The effective capacitance was determined by averaging the energy over a sequence of input patterns with a uniform distribution.

Also plotted in Figure 6 are simulation results from extracted layouts of an 8-bit adder carry chain for three different device (W/L) ratios (N=1, N=2, and N=4). The curve follows the simple first-order model derived very well, and suggests that this example is dominated more by the effect of gate capacitance rather than parasitics. In this case, increasing devices (W/L)'s does not help, and the solution using the smallest possible (W/L) ratios results in the best sizing.

From this section, it is clear that the determination of an “optimal” supply voltage is the key to minimizing power consumption; hence we focus on this issue in the following sections. First, we will review the previous work dealing with choice of supply voltage which were based on reliability and speed considerations[23,24], followed by an architecturally driven approach to supply voltage scaling.

4.3 Reliability-Driven Voltage Scaling

One approach to the selection of an optimal power supply voltage for deep-submicron technologies is based on optimizing the trade-off between speed and reliability[23]. Constant-voltage scaling - the most commonly used technique - results in higher electric fields that create hot carriers. As a result of this, the devices degrade with time (including changes in threshold voltages, degradation of transconductance, and increase in sub-threshold currents), leading to eventual breakdown[11]. One solution to reducing the number of hot carriers is to change the physical device structure, such as the use of LDD (lightly doped drain), usually at the cost of decreased performance. Assuming the use of an LDD structure and a constant hot carrier margin, an optimal voltage of 2.5V was found for a 0.25 μ technology by choosing the minimum point on the delay vs. V_{dd} curve [23]. For voltages above this minimum point, the delay was found to increase with increasing V_{dd} , since the LDD structure used for the purposes of reliability resulted in increased parasitic resistances.

4.4 Technology-Driven Voltage Scaling

The simple first order delay analysis presented in Section 4.1 is reasonably accurate for long channel devices. However, as feature sizes shrink below 1.0 μ , the delay characteristics as a function of lowering the supply voltage deviate from the first order theory presented since it does not consider carrier velocity saturation under high electric fields[11]. As a result of velocity saturation, the current is no longer a quadratic function of the voltage but linear; hence, the current drive is significantly reduced and is approximately given by $I = WC_{ox}(V_{dd} - V_t)v_{max}$ [4]. Given this and the equation for delay in Equation 3, we see that the delay for submicron circuits is relatively independent of supply voltages at high electric fields.

A “technology” based approach proposes choosing the power supply voltage based on maintaining the speed performance for a given submicron technology[24]. By exploiting the relative independence of delay on supply voltage at high electric fields, the voltage can be dropped to some extent for a velocity-saturated device with

$$T_N = K \frac{(C_p + NC_{ref})}{(NC_{ref})} \frac{V_{ref}}{(V_{ref} - V_t)^2} = K(1 + \alpha/N) \frac{V_{ref}}{(V_{ref} - V_t)^2} \quad (\text{EQ 4})$$

where α is defined as the ratio of C_p to C_{ref} , and K represents terms independent of device width and voltage. For a given supply voltage V_{ref} , the speed up of a circuit whose W/L ratios are sized up by a factor of N over a reference circuit using minimum size transistors ($N=1$) is given by $(1 + \alpha/N) / (1 + \alpha)$. In order to evaluate the energy performance of the two designs at the same speed, the voltage of the scaled solution is allowed to vary as to keep delay constant. Assuming that the delay scales as $1/V_{dd}$ (ignoring threshold voltage reductions in signal swings) the supply voltage, V_N , where the delay of the scaled design and the reference design are equal is given by:

$$V_N = \frac{(1 + \alpha/N)}{(1 + \alpha)} V_{ref} \quad (\text{EQ 5})$$

Under these conditions, the energy consumed by the first stage as a function of N is given by:

$$\text{Energy}(N) = (C_p + NC_{ref})V_N^2 = \frac{NC_{ref}(1 + \alpha/N)^3 V_{ref}^2}{(1 + \alpha)^2} \quad (\text{EQ 6})$$

After normalizing against E_{ref} (the energy for the minimum size case), Figure 6 shows a plot of $\text{Energy}(N) / \text{Energy}(1)$ vs. N for various values of α . When there is no parasitic capacitance contribution (i.e. $\alpha = 0$), the energy increases linearly with respect to N , and the solution utilizing devices with the smallest (W/L) ratios results in the lowest power. At high values of α , when parasitic capacitances begin to dominate over the gate capacitances, the power decreases temporarily with increasing device sizes and then starts to increase, resulting in a optimal value for N . The initial decrease in supply voltage achieved from the reduction in delays more than compensates the increase in capacitance due to increasing N . However, after some point the increase in capacitance dominates the achievable reduction in voltage, since the incremental speed increase with transistor sizing is very small (this can be seen in Equation 4, with the delay becoming independent of α as N goes to infinity). Throughout the analysis we have assumed that the parasitic capacitance is independent of device sizing. However, the drain and source diffusion and perimeter capacitances actually increase with increasing area, favoring smaller size devices and making the above a worst-case analysis.

We also evaluated (through experimental measurements and SPICE simulations) the energy and delay performance for several different logic styles and topologies using an 8-bit adder as a reference; the results are shown on a log-log plot in Figure 4. We see that the power-delay product improves as delays increase (through reduction of the supply voltage), and therefore it is desirable to operate at the *slowest* possible speed. Since the objective is to reduce power consumption while maintaining the overall system throughput, compensation for these increased delays at low voltages is required. Of particular interest in this figure is the range of energies required for a transition at a given amount of delay. The best logic family we analyzed (over 10 times better than the worst that we investigated) was the passgate family, CPL, (see section 3.2) if a reduced value for the threshold is assumed[8].

Figures 2, 3, and 4 suggest that the delay and energy behavior as a function of V_{dd} scaling for a given technology is “well-behaved” and relatively independent of logic style and circuit complexity. We will use this result during our optimization of architecture for low-power by treating V_{dd} as a free variable and by allowing the architectures to vary to retain constant throughput. By exploiting the monotonic dependencies of delay and energy versus supply voltage that hold over wide circuit variations, it is possible to make relatively strong predictions about the types of architectures that are best for low power design. Of course, as mentioned previously, there are some logic styles such as NMOS pass-transistor logic without reduced thresholds whose delay and energy characteristics would deviate from the ones presented above, but even for these cases, though the quantitative results will be different, the basic conclusions will still hold.

4.2 Optimal Transistor Sizing with Voltage Scaling

Independent of the choice of logic family or topology, optimized transistor sizing will play an important role in reducing power consumption. For low power, as is true for high speed design, it is important to equalize all delay paths so that a single critical path does not unnecessarily limit the performance of the entire circuit. However, beyond this constraint, there is the issue of what extent the (W/L) ratios should be uniformly raised for all the devices, yielding a uniform decrease in the gate delay and hence allowing for a corresponding reduction in voltage and power. It is shown in this section, that if voltage is allowed to vary, that the optimal sizing for low power operation is quite different from that required for high speed.

In Figure 5, a simple two-gate circuit is shown, with the first stage driving the gate capacitance of the second, in addition to the parasitic capacitance C_p due to substrate coupling and interconnect. Assuming that the input gate capacitance of both stages is given by NC_{ref} , where C_{ref} represents the gate capacitance of a MOS device with the smallest allowable (W/L), then the delay through the first gate at a supply voltage V_{ref} is given by:

rail DCVSL family consumes at least two times more in energy per input transition than a conventional static family. Hence self-timed implementations can prove to be expensive in terms of energy for datapaths that are continuously computing.

4.0 Voltage Scaling

Thus far, we have been primarily concerned with the contributions of capacitance to the power expression CV^2f . Clearly, though, the reduction of V should yield even greater benefits; indeed, reducing the supply voltage is the key to low-power operation, even after taking into account the modifications to the system architecture which is required to maintain the computational throughput. First, a review of circuit behavior (delay and energy characteristics) as a function of scaling supply voltage and feature sizes will be presented. By comparison with experimental data, it is found that simple first order theory yields an amazingly accurate representation of the various dependencies over a wide variety of circuit styles and architectures. A survey of two previous approaches to supply-voltage scaling is then presented, which were focused on maintaining reliability and performance. This is followed by our architecture-driven approach, from which an “optimal” supply voltage based on technology, architecture and noise margin constraints is derived.

4.1 Impact on Delay and Power-Delay Product

As noted in Equation 2, the energy per transition or equivalently the power-delay product in “properly designed” CMOS circuits (as discussed in Section 2), is proportional to V^2 . This is seen from Figure 2, which is a plot of two experimental circuits which exhibit the expected V^2 dependence. Therefore, it is only necessary to reduce the supply voltage for a *quadratic* improvement in the power-delay product of a logic family.

Unfortunately, this simple solution to low power design comes at a cost. As shown in Figure 3, the effect of reducing V_{dd} on the delay is shown for a variety of different logic circuits, that range in size from 56 to 44,000 transistors spanning a variety of functions, all exhibit essentially the same dependence (see Table I). Clearly, we pay a speed penalty for a V_{dd} reduction, with the delays drastically increasing as V_{dd} approaches the sum of the threshold voltages of the devices. Even though the exact analysis of the delay is quite complex if the non-linear characteristic of a CMOS gate are taken into account, it is found that a simple first-order derivation adequately predicts the experimentally determined dependence and is given by:

$$T_d = \frac{C_L \times V_{dd}}{I} = \frac{C_L \times V_{dd}}{\mu C_{ox} (W/L)(V_{dd} - V_t)^2} \quad (\text{EQ 3})$$

with 60 mV/dec being the lower limit. Clearly, the lower S_{th} is, the better, since it is desirable to have the device “turn-off” as close to V_t as possible. As a reference, for an $L=1.5\mu$, $W=70\mu$ NMOS device, at the point where V_{gs} equals V_t , with V_t defined as where the surface inversion charge density is equal to the bulk doping, approximately $1\mu A$ of leakage current is exhibited, or $.014\mu A/\text{micron}$ of gate width[16]. The issue is whether this extra current is negligible in comparison to the time-average current during switching. For a CMOS inverter (PMOS: $W=4\mu$, NMOS: $W=8\mu$), the current was measured to be $64\mu A$ over 3.7nsec at a supply voltage of 2V. This implies that there would be a 100% power penalty for subthreshold leakage if the device were operating at a clock speed of 25 MHz with an activity factor of $p_t = 1/6th$, i.e. the devices were left idle and leaking current 83% of the time. It is not advisable, therefore, to use a true zero threshold device, but instead to use thresholds of at least 0.2V, which provides for at least two orders of magnitude of reduction of subthreshold current. This provides a good compromise between improvement of current drive at low supply voltage operation and keeping subthreshold power dissipation to a negligible level. This value may have to be higher in dynamic circuits to prevent accidental discharge during the evaluation phase [11]. Fortunately, device technologists are addressing the problem of subthreshold currents in future scaled technologies, and reducing the supply voltages also serves to reduce the current by reducing the maximum allowable drain-source voltage [17,18]. The design of future circuits for lowest power operation should therefore explicitly take into account the effect of subthreshold currents.

3.4 Power-down strategies

In synchronous designs, the logic between registers is continuously computing every clock cycle based on its new inputs. To reduce the power in synchronous designs, it is important to minimize switching activity by powering down execution units when they are not performing “useful” operations. This is an important concern since logic modules can be switching and consuming power even when they are not being actively utilized [19].

While the design of synchronous circuits requires special design effort and power-down circuitry to detect and shut down unused units, self-timed logic has inherent power-down of unused modules, since transitions occur only when requested. However, since self-timed implementations require the generation of a completion signal indicating the outputs of the logic module are valid, there is additional overhead circuitry. There are several circuit approaches to generate the requisite completion signal. One method is to use dual-rail coding, which is implicit in certain logic families such as the DCVSL[13,20]. The completion signal in a combinational macrocell made up of cascading DCVSL gates consists of simply ORing the outputs of only the last gate in the chain, leading to small overhead requirements. However, for each computation, dual-rail coding guarantees a switching event will occur since at least one of the outputs must evaluate to zero. We found that the dual

3.2 Conventional Static vs. Pass-gate Logic

A more clear situation exists in the use of transfer gates to implement logic functions, as is used in the CPL (Complementary Passgate Logic) family [8,10]. In Figure 1, the schematic of a typical static CMOS logic circuit for a full adder is shown along with a static CPL version [8]. The passgate design uses only a single transmission NMOS gate, instead of a full complementary passgate to reduce node capacitance. Passgate logic is attractive as fewer transistors are required to implement important logic functions, such as XOR's which only require 2 pass transistors in a CPL implementation. This particularly efficient implementation of an XOR is important since it is key to most arithmetic functions, permitting adders and multipliers to be created using a minimal number of devices. Likewise, multiplexers, registers, and other key building blocks are simplified using passgate designs.

However, a CPL implementation as shown in Figure 1 has two basic problems. First, the threshold drop across the single channel pass transistors results in reduced current drive and hence slower operation at reduced supply voltages; this is important for low-power design since it is desirable to operate at the lowest possible voltages levels. Second, since the "high" input voltage levels at the regenerative inverters is not V_{dd} , the PMOS device in the inverter is not fully turned off, and hence direct-path static power dissipation could be significant. To solve these problems, reduction of the threshold voltage has proven effective, although if taken too far will incur a cost in dissipation due to subthreshold leakage (see section 3.3) and reduced noise margins. The power dissipation for a passgate family adder with zero-threshold pass transistors at a supply voltage of 4V was reported to be 30% lower than a conventional static design, with the difference being even more significant at lower supply voltages [8].

3.3 Threshold voltage scaling

Since a significant power improvement can be gained through the use of low-threshold MOS devices, the question of how low the thresholds can be reduced must be addressed. The limit is set by the requirement to retain adequate noise margins and the increase in subthreshold currents. Noise margins will be relaxed in low power designs because of the reduced currents being switched, however, the subthreshold currents can result in significant static power dissipation. Essentially, subthreshold leakage occurs due to carrier diffusion between the source and the drain when the gate-source voltage, V_{gs} , has exceeded the weak inversion point, but is still below the threshold voltage V_t , where carrier drift is dominant. In this regime, the MOSFET behaves similarly to a bipolar transistor, and the subthreshold current is exponentially dependent on the gate-source voltage V_{gs} , and approximately independent of the drain-source voltage V_{ds} , for V_{ds} approximately larger than 0.1V. Associated with this is the subthreshold slope S_{th} , which is the amount of voltage required to drop the subthreshold current by one decade. At room temperature, typical values for S_{th} lie between 60 to 90 mV/(decade current),

3.1.2 Short-circuit currents

Short circuit (direct-path) currents, I_{sc} in Equation 1, are found in static CMOS circuits. However, by sizing transistors for equal rise and fall times, the short-circuit component of the total power dissipated can be kept to less than 20% [9] (typically $< 5-10\%$) of the dynamic switching component. Dynamic logic does not exhibit this problem, except for those cases in which static pull-up devices are used to control charge sharing[13] or when clock skew is significant.

3.1.3 Parasitic capacitance

Dynamic logic typically uses fewer transistors to implement a given logic function, which directly reduces the amount of capacitance being switched and thus has a direct impact on the power-delay product[14,15]. However, extra transistors may be required to insure that charge-sharing does not result in incorrect evaluation.

3.1.4 Switching activity

The one area in which dynamic logic is at a distinct disadvantage is in its necessity for a precharge operation. Since in dynamic logic every node must be precharged every clock cycle, this means that some nodes are precharged only to be immediately discharged again as the node is evaluated, leading to a higher activity factor. If a two-input N-tree (precharged high) dynamic NOR gate has a uniform input distribution of high and low levels, then the four possible input combinations (00,01,10,11) will be equally likely. There is then a 75% probability that the output node will discharge immediately after the precharge phase, implying that the activity for such a gate is 0.75 (i.e, $P_{NOR} = 0.75 C_L V_{dd}^2 f_{clk}$). On the other hand, the activity factor for the static NOR counterpart will be only 3/16, excluding the component due to the spurious transitions mentioned in section 3.1.1 [Power is only drawn on a 0 to 1 transition, so $p_{0 \rightarrow 1} = p(0)p(1) = p(0) (1-p(0))$]. In general, gate activities will be different for static and dynamic logic and will depend on the type of operation being performed and the input signal probabilities. In addition, the clock buffers to drive the precharge transistors will also require power that is not needed in a static implementation.

3.1.5 Power-down modes

Lastly, power-down techniques achieved by disabling the clock signal have been used effectively in static circuits, but are not as well-suited for dynamic techniques. If the logic state is to be preserved during shut-down, a relatively small amount of extra circuitry must be added to the dynamic circuits to preserve the state, resulting in a slight increase in parasitic capacitance and slower speeds.

3.0 Circuit Design and Technology Considerations

There are a number of options available in choosing the basic circuit approach and topology for implementing various logic and arithmetic functions. Choices between static *vs.* dynamic implementations, passgate *vs.* conventional CMOS logic styles, and synchronous *vs.* asynchronous timing are just some of the options open to the system designer. At another level, there are also various architectural/structural choices for implementing a given logic function; for example, to implement an adder module one can utilize a ripple-carry, carry-select, or carry-lookahead topology. In this section, the trade-offs with respect to low power design between a selected set of circuit approaches will be discussed, followed by an discussion of some general issues and factors affecting the choice of logic family.

3.1 Dynamic *vs.* Static Logic

The choice of using static or dynamic logic is dependent on many other criteria than just its low power performance, e.g. testability and ease of design. However, if only the low power performance is analyzed it would appear that dynamic logic has some inherent advantages in a number of areas including reduced switching activity due to hazards, elimination of short-circuit dissipation and reduced parasitic node capacitances. Static logic has advantages since there is no pre-charge operation and charge-sharing does not exist. Below, each of these considerations will be discussed in more detail.

3.1.1 Spurious Transitions

Static designs can exhibit spurious transitions due to finite propagation delays from one logic block to the next (also called critical races and dynamic hazards[12]), i.e. a node can have multiple transitions in a single clock cycle before settling to the correct logic level. For example, consider a static N-bit adder, with all bits of the summands going from “zero” to “one”, with the carry input set to “zero”. For all bits, the resultant sum should be zero; however, the propagation of the carry signal causes a “one” to appear briefly at most of the outputs. These spurious transitions dissipate extra power over that strictly required to perform the computation. The number of these extra transitions is a function of input patterns, internal state assignment in the logic design, delay skew, and logic depth. To be specific about the magnitude of this problem, an 8-bit ripple-carry adder with an uniformly distributed set of random input patterns, will typically consume an extra 30% in energy. Though it is possible with careful logic design to eliminate these transitions, dynamic logic intrinsically does not have this problem, since any node can undergo at most one power-consuming transition per clock cycle.

Another important consideration, particularly in portable applications, is that many computation tasks are likely to be real-time; the radio modem, speech and video compression, and speech recognition all require computation that is always at near-peak rates. Conventional schemes for conserving power in laptops, which are generally based on power-down schemes, are not appropriate for these continually active computations. On the other hand, there is a degree of freedom in design that is available in implementing these functions, in that once the real-time requirements of these applications are met, there is no advantage in increasing the computational throughput. This fact, along with the availability of almost “limitless” numbers of transistors, allows a strategy to be developed for architecture design, which if can be followed, will be shown to provide significant power savings.

2.0 Sources of Power Dissipation

There are three major sources of power dissipation in digital CMOS circuits which are summarized in the following equation:

$$P_{\text{total}} = p_t (C_L * V * V_{\text{dd}} * f_{\text{clk}}) + I_{\text{sc}} * V_{\text{dd}} + I_{\text{leakage}} * V_{\text{dd}} \quad (\text{EQ 1})$$

The first term represents the switching component of power, where C_L is the loading capacitance, f_{clk} is the clock frequency and p_t is the probability that a power consuming transition occurs (the activity factor). In most cases, the voltage swing, V , is the same as the supply voltage, V_{dd} ; however, some logic circuits, such as in single-gate pass-transistor implementations, the voltage swing on some internal nodes may be slightly less [8]. The second term is due to the direct-path short circuit current, I_{sc} , which arises when both the NMOS and PMOS transistors are simultaneously active, conducting current directly from supply to ground [9,10]. Finally, leakage current, I_{leakage} , which can arise from substrate injection and sub-threshold effects, is primarily determined by fabrication technology considerations [11] (see section 3.3). The dominant term in a “well-designed” circuit is the switching component, and low power design thus becomes the task of minimizing p_t , C_L , V_{dd} and f_{clk} , while retaining the required functionality.

Power-delay product can be interpreted as the amount of energy expended in each switching event (or transition) and is thus particularly useful in comparing the power dissipation of various circuit styles. If it is assumed that only the switching component of the power dissipation is important then it is given by:

$$\text{Energy per transition} = P_{\text{total}} / f_{\text{clk}} = C_{\text{effective}} V_{\text{dd}}^2 \quad (\text{EQ 2})$$

where $C_{\text{effective}}$ is the effective capacitance being switched to perform a computation and is given by $C_{\text{effective}} = p_t * C_L$.

Although the traditional mainstay of portable digital applications has been in low-power, low-throughput uses such as wristwatches and pocket calculators, there are an ever-increasing number of portable applications requiring low-power and high-throughput. For example, notebook and laptop computers, representing the fastest growing segment of the computer industry, are demanding the same computation capabilities as found in desktop machines. Equally demanding are developments in personal communications services (PCS), such as the current generation of digital cellular telephony networks which employ complex speech compression algorithms and sophisticated radio modems in a pocket sized device. Even more dramatic are the proposed future PCS applications, with universal, portable multimedia access supporting full-motion digital video and control via speech recognition[3]. In these applications, not only will voice be transmitted via wireless links, but data as well. This will facilitate new services such as multimedia database access (video and audio in addition to text) and supercomputing for simulation and design, through an intelligent network which allows communication with these services or other people at any place and time. Power for video compression and decompression, and speech recognition must be added to the portable unit to support these services - on top of the already-lean power budget for the analog transceiver and speech encoding. Indeed, it is apparent that portability can no longer be associated with low-throughput; instead, vastly increased capabilities, actually in excess of that demanded of fixed workstations, must be placed in a low-power, portable environment.

Even when power is available in non-portable applications, the issue of low power design is becoming critical. Up until now, this power consumption has not been of great concern, since large packages, cooling fins and fans have been capable of dissipating the generated heat. However, as the density and size of the chips and systems continues to increase, the difficulty in providing adequate cooling might either add significant cost to the system or provide a limit on the amount of functionality that can be provided.

Thus, it is evident that methodologies for the design of high-throughput, low-power digital systems are needed. Fortunately, there are clear technological trends that give us a new degree of freedom, so that it may be possible to satisfy these seemingly contradictory requirements. Scaling of device feature sizes, along with the development of high density, low-parasitic packaging, such as multi-chip modules[4,5,6], will alleviate the overriding concern with the numbers of transistors being used. When MOS technology has scaled to 0.2 μ m minimum feature size it will be possible to place from 1-10 \times 10⁹ transistors in an area of 8" by 10" if a high-density packaging technology is used. The question then becomes how can this increased capability be used to meet a goal of low power operation. Previous analyses on the question of how to best utilize increased transistor density at the chip level, concluded that for high-performance microprocessors the best use is to provide increasing amounts of on-chip memory[7]. It will be shown here that for computationally intensive functions that the best use is to provide additional circuitry to parallelize the computation.

Low Power CMOS Digital Design

Anantha P. Chandrakasan

Samuel Sheng

Robert W. Brodersen[†]

EECS Department,
University of California at Berkeley

Abstract: Motivated by emerging battery operated applications that demand intensive computation in portable environments, techniques are investigated which reduce power consumption in CMOS digital circuits while maintaining computational throughput. Techniques for low power operation are shown which use the lowest possible supply voltage coupled with architectural, logic style, circuit and technology optimizations. An architectural based scaling strategy is presented which indicates that the optimum voltage is much lower than that determined by other scaling considerations. This optimum is achieved by trading increased silicon area for reduced power consumption.

1.0 Introduction

With much of research efforts of the past ten years directed toward increasing the speed of digital systems, present-day technologies possess computing capabilities which make possible powerful personal workstations, sophisticated computer graphics, and multi-media capabilities such as real-time speech recognition and real-time video. High-speed computation has thus become the expected norm from the average user, instead of being the province of the few with access to a powerful mainframe. Likewise, another significant change in the attitude of users is the desire to have access to this computation at any location, without the need to be physically tethered to a wired network. The requirement of portability thus places severe restrictions on size, weight, and power. Power is particularly important since conventional nickel-cadmium battery technology only provides 20 W-hrs of energy for each pound of weight [1]. Improvements in battery technology are being made, but it is unlikely that a dramatic solution to the power problem is forthcoming; it is projected that only a 30% improvement in battery performance will be obtained over the next 5 years[2].

[†] Please mail all correspondence to: c/o Prof. Robert W. Brodersen,
EECS Department, Cory Hall, University of California, Berkeley, CA94720