

Low Power VLSI: High End Design Techniques

P.Bujjibabu, V.Satyanarayana, G.Jyothirmai, GNPK Mahalakshmi.E

Abstract — Low power has emerged as a principal argument in today's electronics diligence. The need for low power has caused a major hypothesis shift where power dissipation has become as important a consideration as performance and area. This object reviews various strategies and methodologies for designing low power circuits and systems. It describes the many issues facing designers at architectural, logic, circuit and device levels and presents some of the techniques that have been proposed to overcome these complications. The paper concludes with the future challenges that must be met to design low power, high performance systems.

Keywords — Architecture driven voltage scaling, Complex functionality, Glitching, Monte-Carlo Simulation, Power estimation techniques, Output entropy

1. INTRODUCTION

IN the past, the major concerns of the VLSI designer were area, performance, cost and reliability; power consideration was mostly of only inferior importance. In recent years, however, this has begun to change and, increasingly, power is being given similar weight to area and speed considerations. Several factors have contributed to this trend. In high-speed computation and complex functionality applications with low power consumption, average power consumption is a critical design concern. The projected power budget for a battery-powered, A4 format, portable multimedia terminal, when implemented using off-the-shelf components not optimized for low-power operation, is about 40 W. With advanced Nickel-Metal-Hydride (secondary) battery technologies offering around 65watt-hours/kilogram [42], this terminal would require an unacceptable 6 kilograms of batteries for 10 hours of operation between recharges. In the absence of low-power design techniques then, current and future portable devices will suffer from either very short battery life or very heavy battery pack. There also exists a strong pressure for producers of high-end products to reduce their power consumption. The cost associated with packaging and cooling such devices is prohibitive. Since core power consumption must be dissipated through the packaging, increasingly expensive packaging and cooling strategies are required as chip power consumption increases. Consequently, there is a clear financial advantage to reducing the power consumed in high performance systems.

In addition to cost, there is the issue of reliability. High power systems often run hot and high temperature to

exacerbate several silicon failures mechanisms. Every 10 °C increase in operating temperature roughly doubles a component's failure rate [48]. In this context, peak power (maximum possible power dissipation) is a critical design factor as it determines the thermal and electrical limits of designs, impacts the system cost, size and weight, dictates specific battery type, component and system packaging and heat sinks, and aggravates the resistive and inductive voltages drop problems. Another crucial driving factor is that excessive power consumption is becoming the limiting factor in integrating more transistors on a single chip or on a multiple-chip module. From the environmental viewpoint, the smaller the power dissipation of electronic systems, the lower the heat pumped into the rooms, the lower the electricity consumed and hence the lower the impact on global environment, the less the office noise (e.g., due to elimination of a fan from the desktop), and the less stringent the environment/office power delivery or heat removal requirements. In the class of micro-powered battery-operated, portable applications, such as cellular phones and personal digital assistants, the goal is to keep the battery lifetime and weight reasonable and the packaging cost low. For high performance, portable computers, such as laptop and notebook computers, the goal is to reduce the power dissipation of the system to a point which is about half of the total power dissipation (including that of display and hard disk). Finally, for high performance, non-battery operated systems, such as workstations, set-top computers and multimedia digital signal processors, the overall goal of power minimization is to reduce system cost (cooling, packaging and energy bill) while ensuring long-term device reliability. If extending the battery life is the only concern, then the energy (that is, the power-delay product) should be minimized. In this case the battery consumption is minimized even though an operation may take a very long time. On the other hand, if both the battery life and the circuit delay are important, then the energy-delay product must be minimized [21]. In most design scenarios, the circuit delay is set based on system-level considerations, and hence during circuit

- P.Bujjibabu, Assistant Professor, Aditya Engineering College, Andhra Pradesh, India, 9949370089. bujjibabuforu@gmail.com
- V.Satyanarayana, Sr. Assistant Professor, Aditya Engineering College, Andhra Pradesh, India, 9908173673. vasece453@gmail.com
- G.Jyothirmai, Assistant Professor, Aditya Engineering College, Andhra Pradesh, India, 9346332714. sivajyothi1427@gmail.com
- GNPK Mahalakshmi.E, Assistant Professor, Al Ameer college of engineering & Technology, Andhra Pradesh, India, 8885685765.

optimization, one minimizes power under user-specified timing constraints.

2. SOURCES OF POWER DISSIPATION

Power dissipation in CMOS circuits is caused by three sources: 1) *the leakage current* which is primarily determined by the fabrication technology, consists of reverse bias current in the parasitic diodes formed between source and drain diffusions and the bulk region in a MOS transistor as well as the sub threshold current that arises from the inversion charge that exists at the gate voltages below the threshold voltage, 2) *the short-circuit (rush-through) current* which is due to the DC path between the supply rails during output transitions and 3) *the charging and discharging of capacitive loads* during logic changes. The resulting current due to transistor ON/OFF is proportional to the area of the drain diffusion and the leakage current density. The sub threshold leakage current for long channel devices increases linearly with the ratio of the channel width over channel length and decreases exponentially with $V_{gs}-V_t$ where V_{gs} is the gate bias and V_t is the threshold voltage. With reduced power supply and device threshold voltages, the sub threshold current will however become more pronounced. In addition, at short channel lengths, the sub threshold current also becomes exponentially dependent on drain voltage instead of being independent of V_{DS} (see [18] for a recent analysis). The short-circuit (crowbar current) power consumption for an inverter gate is proportional to the gain of the inverter, the cubic power of supply voltage minus device threshold, the input rise/fall time, and the operating frequency. The maximum short circuit current flows when there is no load; this current decreases with the load. If gate sizes are selected so that the input and output rise/fall times are about equal, the short-circuit power consumption will be less than 15% of the dynamic power consumption. If, however, design for high performance is taken to the extreme where large gates are used to drive relatively small loads, then there will be a stiff penalty in terms of short-circuit power consumption. The dominant source of power dissipation is thus the charging and discharging of the node capacitances (also referred to as the dynamic power dissipation) and is given by:

$$P = 0.5CV_{dd}^2 E(sw) f_{clk}$$

where C is the physical capacitance of the circuit, V is the supply voltage, $E(sw)$ (referred as the switching activity) is the average number of transitions in the circuit per $1/f_{clk}$ time, and f_{clk} is the clock frequency.

3. LOW POWER DESIGN SPACE

The previous section revealed the three degrees of freedom inherent in the low-power design space: voltage, physical

capacitance, and data activity. Optimizing for power entails an attempt to reduce one or more of these factors.

3.1. Voltage

Because of its quadratic relationship to power, voltage reduction offers the most effective means of minimizing power consumption. Without requiring any special circuits or technologies, a factor of two reduction in supply voltage yields a factor of four decrease in power consumption. Furthermore, this power reduction is a global effect, experienced not only in one sub-circuit or block of the chip, but throughout the entire design. Unfortunately, we pay a speed penalty for supply voltage reduction, with delays drastically increasing as V_{dd} approaches the threshold voltage V_t of the devices. This tends to limit the useful range of V_{dd} to a minimum of about 2-3 V_t . In [9], an architecture driven voltage scaling strategy is presented in which parallel and pipelined architectures are used to compensate for the increased gate delays at reduced supply voltages and meet throughput constraints. Reducing the V_t allows the supply voltage to be scaled down without loss in speed. The limit of how low the V_t can go is set by the requirement to set adequate noise margins and control the increase in sub threshold leakage currents. The optimum V_t must be determined based on the current drives at low supply voltage operation and control of the leakage currents. Since the inverse threshold slope (S) of a MOSFET is invariant with scaling, for every 80-100 mV (based on the operating temperature) reduction in V_t , the standby current will be increased by one order of magnitude. This tends to limit V_t to about 0.3 V for room temperature operation of CMOS circuits. Another important concern in the low V_{dd} - low V_t regime is the fluctuation in V_t . Basically, delay increases by 3x for a delta V_{dd} of plus/minus 0.15 V_{at} V_{dd} of 1 V.

3.2. Physical Capacitance

Dynamic power consumption depends linearly on the physical capacitance being switched. So, in addition to operating at low voltages, minimizing capacitances offers another technique for minimizing power consumption. In order to consider this possibility we must first understand what factors contribute to the physical capacitance of a circuit. Power dissipation is dependent on the physical capacitances seen by individual gates in the circuit. Capacitance calculation can however be done easily after technology mapping by using the logic and delay information from the library. Interconnect plays an increasing role in determining the total chip area, delay and power dissipation, and hence, must be accounted for as early as possible during the design process. Approximate interconnect capacitance estimates can be obtained by using information derived from a companion placement solution

^[39] or by using stochastic / procedural interconnect models ^[40]. Interconnect capacitance estimation after layout is straight-forward and in general accurate. It is observed that capacitances can be kept at a minimum by using less logic, smaller devices, fewer and shorter wires. Example techniques for reducing the active area include *resource sharing, logic minimization and gate sizing*. Example techniques for reducing interconnect include *register sharing, common sub-function extraction, placement and routing*. But reducing device sizes reduces physical capacitance, and also reduces the current drive of the transistors making the circuit operate more slowly. This loss in performance might prevent us from lowering V_{dd} as much as we might otherwise be able to do.

3.3. Switching Activity

In addition to voltage and physical capacitance, switching activity also influences dynamic power consumption. The data activity determines how often this switching occurs. There are two components to switching activity: fclk which determines the average periodicity of data arrivals and E(sw) which determines how many transitions each arrival will generate. For circuits that do not experience glitching, E(sw) can be interpreted as the probability that a power consuming transition will occur during a single data period. For certain logic styles, however, glitching can be an important source of signal activity and, therefore, deserves some mention here. Glitching refers to spurious and unwanted transitions that occur before a node settles down to its final steady-state value. Glitching often arises when paths with unbalanced propagation delays converge at the same point in the circuit. The data activity E(sw) can be combined with the physical capacitance C to obtain switched capacitance, C_{sw}=C.E(sw), which describes the average capacitance charged during each data period 1/fclk. It should be noted that it is the switched capacitance that determines the power consumed by a CMOS circuit.

Calculation of Switching Activity

Calculation of the switching activity in a logic circuit is difficult as it depends on a number of circuit parameters and technology-dependent factors which are not readily available or precisely characterized. Some of these factors are described next.

Input Pattern Dependence:

Switching activity at the output of a gate depends not only on the switching activities at the inputs and the logic function of the gate, but also on the spatial and temporal dependencies among the gate inputs. The straight-forward approach of estimating power by using a simulator is greatly complicated by this pattern dependence problem. It is clearly infeasible to estimate the power by exhaustive simulation of the circuit. Recent techniques overcome this

difficulty by using probabilities that describe the set of possible logic values at the circuit inputs and developing mechanisms to calculate these probabilities for gates inside the circuit. Alternatively, exhaustive simulation may be replaced by Monte-Carlo simulation with well-defined stopping criterion for specified relative or absolute error in power estimates for a given confidence level ^[7].

Delay Model:

Based on the delay model used, the power estimation techniques could account for steady-state transitions (which consume power, but are necessary to perform a computational task) and/or hazards and glitches (which dissipate power without doing any useful computation). Sometimes, the first component of power consumption is referred as the functional activity while the latter is referred as the spurious activity. It is shown in ^[3] that the mean value of the ratio of hazardous component to the total power dissipation varies significantly with the considered circuits (from 9% to 38% in random logic circuits) and that the spurious power dissipation cannot be neglected in CMOS circuits. The spurious activity is much higher in certain data path modules (such as adders and multipliers); Indeed, in a 32-bit pipelined multiplier, the power dissipation due to hazard activity is 3-4 times higher than that due to functional activity!

Current power estimation techniques often handle both zero-delay (non-glitch) and real delay models. In the first model, it is assumed that all changes at the circuit inputs propagate through the internal gates of the circuits instantaneously. The latter model assigns each gate in the circuit a finite delay and can thus account for the hazards in the circuit. A real delay model significantly increases the computational requirements of the power estimation techniques while improving the accuracy of the estimates

Logic Function:

Switching activity at the output of a logic gate is also strongly dependent on the Boolean function of the gate itself. This is because the logic function of a gate determines the probability that the present value of the gate output is different from its previous value. For example, under the assumption that the input signals are uncorrelated, switching activity at the output of a (static) two-input NAND or NOR gate is 3/8 while that at the output of a two-input XOR gate is 1/2. Indeed, switching activity at the output of a K-input NAND or NOR gate approaches 1/2K-1 for large K whereas that for a K-input XOR gate remains at 1/2.

Logic Style switching activity:

Logic Style switching activity of the circuits is also a function of the logic style used to implement the circuit.

The functional activity in dynamic circuits is always higher than that in static implementation of the same circuit as all nodes are pre charged to some value (one in N-type dynamic and zero in P-type dynamic) before the new input data arrives. This effectively increases the number of power consuming transitions. For example, under pseudo-random input signals, switching activities of two-input N-type dynamic NAND, NOR and XOR gates are 3/2, 1/2 and 1, respectively and those of the P-type version of these same gates are 1/2, 3/2 and 1, respectively. These values should be compared to the switching activities of these gates in static CMOS which are 3/8, 3/8 and 1/2, respectively. Dynamic circuits are free from glitches

Circuit Structure:

The major difficulty in computing the switching activities is the re-convergent nodes. Indeed, if a network consists of simple gates and has no re convergent fan-out nodes (that is, circuit nodes that receive inputs from two paths that fan-out from some other circuit node), then the exact switching activities can be computed during a single post-order traversal of the network. For networks with re convergent fan-out, the problem is much more challenging as internal signals may become strongly correlated and exact consideration of these correlations cannot be performed with reasonable computational effort or memory usage. Current power estimation techniques either ignore these correlations or approximate them, thereby improving the accuracy at the expense of longer run times.

Statistical Variation of Circuit Parameters:

In real networks, statistical perturbations of circuit parameters may change the propagation delays and produce changes in the number of transitions because of the appearance or disappearance of hazards. It is therefore useful to resolve the change in the signal transition count as a function of these statistical perturbations. Variation of gate delay parameters may change the number of hazards occurring during a transition as well as their duration.

4. POWER ESTIMATION TECHNIQUES

The design for low power problem cannot be achieved without accurate power prediction and optimization tools or without power efficient gate and module libraries. Therefore, there is a critical need for CAD tools to estimate power dissipation during the design process to meet the power budget without having to go through a costly redesign effort and enable efficient design and characterization of the design libraries. In the following section, various techniques for power estimation at the circuit, logic and behavioral levels will be reviewed. These techniques are divided into two general categories: simulation based and non simulation based.

4.1. SIMULATIVE APPROACHES

4.1.1 Brute-force simulation

Circuit simulation based techniques [27] simulate the circuit with a representative set of input vectors. They are accurate and capable of handling various device models, different circuit design styles, single and multi-phase clocking methodologies, tri-state drives, etc. However, they suffer from memory and execution time constraints and are not suitable for large, cell-based designs. In addition, it is difficult to generate a compact stimulus vector set to calculate accurate activity factors at the circuit nodes. The size of such a vector set is dependent on the application and the system environment [44]. Power Mill [16] is a transistor-level power simulator and analyzer which applies an event-driven timing simulation algorithm (based on simplified table-driven device models, circuit partitioning and single-step nonlinear iteration) to increase the speed by two to three orders of magnitude over SPICE. Switch-level simulation techniques are in general much faster than circuit-level simulation techniques, but are not as accurate or versatile. Verilog-XL logic simulator is a Verilog-based gate-level simulation program that relies on the accuracy of the macro models built for the gates in the ASIC library as well as gate-level timing analysis to produce fast and accurate power estimates. The accuracy depends heavily on the quality of the macro models, the glitch filtering scheme used and the accuracy of physical capacitances provided at the gate level. The speed is 3-4 orders of magnitude faster than SPICE. Important statistics include the number of operations of a given type, the number of bus, register and memory accesses and the number of I/O operations executed within a given period [10] [24]. Instruction level simulation or behavioral simulators are easily (and have indeed been) adapted to produce information.

4.1.2 Hierarchical simulation

A simulation method based on a hierarchy of simulators is presented also. The idea is to use a hierarchy of power simulators (for example, at architectural, gate-level and circuit-level) to achieve a reasonable accuracy and efficiency trade off. Another good example is Entice-Aspen [9]. This power analysis system consists of two components: Aspen which computes the circuit activity information and Entice which computes the power characterization data. A stimulus file is to be supplied to Entice where power and timing delay vectors are specified. The set of power vectors discretizes all possible events in which power can be dissipated by the cell. With the relevant parameters set according to the user's specs, a SPICE circuit simulation is invoked to accurately obtain the power dissipation of each vector. During logic simulation, Aspen monitors the transition count of each cell and computes the total power

consumption as the sum of the power dissipation for all cells in the power vector path.

4.1.3 Monte Carlo simulation

A Monte Carlo simulation approach for power estimation which alleviates the input pattern dependence problem has been proposed in [8]. This approach consists of applying randomly generated input patterns at the circuit inputs and monitoring the power dissipation per time interval T using a simulator. Based on the assumption that the power consumed by the circuit over any period T has a normal distribution, and for a desired percentage error in the power estimate and a given confidence level, the number of required power samples is estimated. The designer can use an existing simulator (circuit-level, gate-level or behavioral) in the inner loop of the Monte-Carlo program, thus trading accuracy for higher efficiency. The convergence time for this approach is fast when estimating the total power consumption of the circuit. However, when signal probability (or power consumption) values on individual lines of the circuit are required, the convergence rate is very slow.

4.2. Non-simulative Approaches

4.2.1 Behavioural Level

For functional units (adders, multipliers and registers) or for memories, power estimates are directly obtained from the design library whereby each functional unit has been simulated using pseudo-random white noise data and the average switched capacitance per clock cycle has been calculated and stored in the library.

The power model for a functional unit may be parametrized in terms of its input bit width. The library thus contains interface descriptions of each module, description of its parameters, its area, delay and internal power dissipation (assuming pseudo-random white noise data inputs). The latter is determined by extracting a circuit or logic level model from the layout or logic level descriptions of the module, simulating it using a long stream of randomly generated input patterns and calculating the average power dissipation per pattern. These characteristics are available in terms of the parameter values (i.e., equations) or in the form of tables. The power model thus generated and stored for each module in the library has to be "conditioned" or "modulated" by the actual input switching activities in order to provide power estimates which are sensitive to the input activities. In [41] and [28], the model consists of a single physical capacitance value and a single switching activity value which represents the average switching activity on each input bit. In [29], a more detailed model is presented where it is projected that data in the data path of a digital system can be divided into two regions: the Least Significant Bits (LSB)

which act as uncorrelated white noise and the Most Significant Bits (MSB) which correspond to sign bits and exhibit strong temporal dependence. The power model thus uses two capacitance values and requires two input switching activity values corresponding to the LSB and MSB regions. Another parametric model is described in [49], where the power dissipation of the various components of typical processor architecture is expressed as a function of set of primary parameters. Word-level behavior of a data input can be properly captured by its probability density function (pdf). Similarly, spatial correlation between two data inputs can be captured by their joint pdf. This observation is used in [11],[12] to develop a probabilistic technique for behavioral level power prediction which consists of four steps:

- 1) Building the joint pdf of the input variables of a data flow graph (DFG) based on the given input vectors
- 2) Computing the joint pdf for any combination of internal arcs in the DFG
- 3) Calculating the switching activity at the inputs of each functional block or register in the DFG using the joint pdf of the inputs and the data representation format which determines the (bit-level) Hamming distances of (word-level) data values
- 4) Estimating the power dissipation of each functional block using the input statistics obtained in step 3 and the library characterization data that gives the physical capacitance information for each module in the library. This method is very robust, but suffers from the worst-case complexity of joint pdf computation and inaccuracies associated with the library characterization data. An information theoretic approach is described in [33] and [37] which relies on information theoretic measures of activity (for example, entropy) to devise fast, yet accurate, power estimation at the algorithmic and structural behavioral levels. In the following, the approach presented in [33] will be summarized. Entropy characterizes the uncertainty of a sequence of applied vectors and thus, intuitively, is related to switching activity. Indeed, it is shown that an upper bound on the average switching activity of a bit is half of its entropy. Knowing the statistics of the input stream and having some information about the structure (or functionality) of the circuit, the input and output entropies per bit are calculated using a closed form expression that gives the output entropy per bit as a function of the input entropy per bit, a structure-dependent information scaling factor, and the distribution of gates as a function of logic depth in the circuit (or using a compositional technique which has a linear complexity in terms of the circuit size). Next the average entropy per circuit line is calculated and used as an estimate of the average switching activity per signal line. This is then used to estimate the power dissipation of the module. A major advantage of this

technique is that it is not simulation based and is thus very fast, yet it produces accurate power estimates. The above techniques apply to data paths. Behavioral power prediction models have also been proposed for the controller circuitry in [30],[28]. These techniques provide quick estimation of the power dissipation in a controller based on the knowledge of its target implementation style (that is, pre-charged pseudo-NMOS or dynamic PLA), the number of inputs, outputs, states, and so on.

4.2.2 Logic Level Estimation under a Zero Delay Model

Most of the power in CMOS circuits is consumed during charging and discharging of the load capacitance. To estimate the power consumption, one has to calculate the (switching) activity factors of the internal nodes of the circuit. Methods of estimating the activity factor $E(sw)$ at a circuit node n involve estimation of signal probability $prob(n)$, which is the probability that the signal value at the node is one. Under the assumption that the values applied to each circuit input are temporally independent (that is, value of any input signal at time t is independent of its value at time $t-1$), we can write:

$$E(sw) = 2 prob(n) (1 - prob(n))$$

In the recent years, a computational procedure based on Ordered Binary-Decision Diagrams (OBDDs) [5] has become widespread. In this method, which is known as the OBDD-based method, the signal probability at the output of a node is calculated by first building a BDD corresponding to the global function of the node (i.e., function of the node in terms of the circuit inputs) and then performing a post order traversal of the OBDD using equation:

$$prob(y) = prob(x)prob(f_x) + \overline{prob(x)}prob(\overline{f_x})$$

This leads to a very efficient computational procedure for signal probability estimation. In [17], a procedure for propagating signal probabilities from the circuit inputs toward the circuit outputs using only pair wise correlations between circuit lines and ignoring higher order correlation terms is described. In [47] and [32], the temporal correlation between values of some signal x in two successive clock cycles is modeled by a time-homogeneous Markov chain which has two states 0 and 1 and four edges where each edge ij ($i, j = 0, 1$) is annotated with the conditional probability $prob_{ijx}$ that x will go to state j at time $t+1$ if it is in state i at time t .

The various transition probabilities can be computed exactly using the OBDD representation of the logic function of x in terms of the circuit inputs. The authors of [32] also describe a mechanism for propagating the transition probabilities through the circuit which is more efficient as there is no need to build the global function of each node in

terms of the circuit inputs. The loss in accuracy is often small while the computational saving is significant. They then extend the model to account for spatio-temporal correlations.

Estimation under a Real Delay Model

The above methods only account for steady-state behavior of the circuit and thus ignore hazards and glitches. This section reviews some techniques that examine the dynamic behavior of the circuit and thus estimate the power dissipation due to hazards and glitches. In [20], the exact power estimation of a given combinational logic circuit is carried out by creating a set of symbolic functions such that summing the signal probabilities of the functions corresponds to the average switching activity at a circuit line x in the original combinational circuit (this method is known as the symbolic simulation method). The inputs to the created symbolic functions are the circuit input lines at time instances 0-and \bullet . Each function is the exclusive or of the characteristic functions describing the logic values of x at two consecutive instances. The concept of a probability waveform is introduced in [6]. This waveform consists of a sequence of transition edges or events over time from the initial steady state (time 0-) to the final steady state (time 8) where each event is annotated with an occurrence probability. The probability waveform of a node is a compact representation of the set of all possible logical waveforms at that node. Given such waveforms at the circuit inputs and with some convenient partitioning of the circuit, the authors examine every sub-circuit and derive the corresponding waveforms at the internal circuit nodes. In [36], an efficient probabilistic simulation technique is described that propagates transition waveforms at the circuit primary inputs up in the circuit and thus estimates the total power consumption (ignoring signal correlations due to the re-convergent fan out nodes). A tagged probabilistic simulation approach is described in that correctly accounts for re-convergent fan-out and glitches. The key idea is to break the set of possible logical waveforms at a node n into four groups, each group being characterized by its steady state values (i.e., values at time instance 0-and \bullet).Next, each group is combined into a probability waveform with the appropriate steady-state tag. Given the tagged probability waveforms at the input of a simple gate, it is then possible to compute the tagged probability waveforms at the output of the gate. The correlation between probability waveforms at the inputs is approximated by the correlation between the steady state values of these lines. This is much more efficient than trying to estimate the dynamic correlations between each pair of events. This approach requires significantly less memory and runs much faster than symbolic simulation, yet achieves very high accuracy.

4.2.3 Sequential Circuits

Recently developed methods for power estimation have primarily focused on combinational logic circuits. The estimates produced by purely combinational methods can greatly differ from those produced by the exact method. Indeed, accurate average switching activity estimation for finite state machines (FSMs) is considerably more difficult than that for combinational circuits for two reasons:

- 1) The probability of the circuit being in each of its possible states has to be calculated;
- 2) The present state line inputs of the FSM are strongly correlated (that is, they are temporally correlated due to the machine behavior as represented in its State Transition Graph description and they are spatially correlated because of the given state encoding). A first attempt at estimating switching activity in FSMs has been presented in [24]. The idea is to unroll the next state logic once (thus capturing the temporal correlations of present state lines) and then perform symbolic simulation on the resulting circuit (which is hence treated as a combinational circuit). The authors of and [35] also describe a method for approximate switching activity estimation of sequential circuits. The basic computation step is the solution of a non-linear system of equations in terms of the present state bit probabilities and signal probabilities for the combinational inputs of the FSM. The fixed point (or zero) of this system of equations can be found using the Picard-Peano (or Newton-Raphson) iteration [31]. Increasing the number of variables or the number of equations in the above system results in increased

5. POWER MINIMIZATION TECHNIQUES

To address the challenge to reduce power, the semiconductor industry has adopted a multifaceted approach, attacking the problem on four fronts:

1. *Reducing chip and package capacitance*: This can be achieved through process development such as SOI with partially or fully depleted wells, CMOS scaling to submicron device sizes, and advanced interconnect substrates such as Multi-Chip Modules (MCM). This approach can be very effective but is also very expensive and has its own pace of development and introduction to the market.
2. *Scaling the supply voltage*: This approach can be very effective in reducing the power dissipation, but often requires new IC fabrication processing. Supply voltage scaling also requires support circuitry for low-voltage operation including level-converters and DC/DC converters as well as detailed consideration of issues such as signal-to-noise.
3. *Employing better design techniques*: This approach promises to be very successful because the investment to reduce power by design is relatively small in comparison to the

other three approaches and because it is relatively untapped in potential.

4. *Using power management strategies*: The power savings that can be achieved by various static and dynamic power management techniques are very application dependent, but can be significant. In the following we will discuss these strategies in some depth. The various approaches interact with one another, for example CMOS device scaling, supply voltage scaling, and choice of circuit architecture must be done judiciously and carefully in order to find an optimum power-area-delay trade-off.

5.1. CMOS Device and Voltage Scaling

The main force behind this voltage scaling drive is the ability to produce complex, high performance systems on a chip. This is further exacerbated by the projected explosion in demand for portable and wireless systems with very low power consumption. It is also expected that various memory and ASIC's will also switch to lower supply voltages to maintain manageable power densities. A key concern is the availability of the complete chip set to make up systems at reduced supply voltages. However, most of the difficulties can be circumvented by techniques to mix and match different supply voltages on board or on the chip. In [15], two CMOS device and voltage scaling scenarios are described, one optimized for the highest speed and one trading off high performance for significantly lower power (the speed of the low power case in one generation is about the same as the speed of the high-performance case of the previous generation, with greatly reduced power consumption). This paper also presents a discussion of how high the electric field in a transistor channel can go without impacting the long term device reliability, while at the same time achieving high performance and low power. Next the speed/standby current trade-off is addressed, dealing with the issue of non scalability of the threshold voltage. The status of silicon-on-insulator (SOI) approach to scaled CMOS is also reviewed, showing that the potential for 3x savings in power compared to the bulk case at the same speed. The performance improvement of SOI compared to bulk CMOS is mainly due to the reduction of parasitic capacitances and body effect. Also, in partially depleted device designs, the floating body effect can give rise to a sharper sub threshold slope (< 60 mV/dec) at high drain bias, which effectively reduces the threshold voltage and can actually improve the performance at a given standby current. In addition, CMOS on SOI offers significant reduction in soft error rate, latch-up elimination, and simpler isolation which results in reduced wafer fabrication steps. The main challenges are the availability of low cost wafers with low defect density at high volumes, floating body effects on the device and circuit operation, and heat dissipation through the buried oxide.

5.2. CAD Methodologies and Techniques

Low power VLSI design can be achieved at various levels of the design abstraction from algorithmic and system levels down to layout and circuit levels. In the following, some of these optimization techniques will be briefly mentioned.

5.2.1 System Design

At the system level, inactive hardware modules may be automatically turned off to save power; Modules may be provided with the optimum supply voltage and interfaced by means of level converters; Some of the energy that is delivered from the power supply may be cycled back to the power supply; A given task may be partitioned between various hardware modules or programmable processors or both so as to reduce the system-level power consumption.

5.2.2 Behavioral Synthesis

Behavioral synthesis is the process of generating a register-transfer level (RTL) design from an algorithmic behavioral specification. In particular, it constructs a structural view of the data path and a logical view of the control unit of a circuit. The data path consists of a set of interconnected functional units (arithmetic, logic, memory and registers) and steering units (multiplexers and busses) while the control unit sends signals to the data path to schedule the appropriate sequence of operations in time. The behavioral synthesis process consists of three steps: allocation, assignment and scheduling. A wide class of transformations can be done at the behavioral level and most of them are typically aimed at either reducing the number of cycles in a computation or reducing the number of resources used in the computation. One interesting approach is to introduce more concurrency in a circuit to speed it up and then to reduce the voltage until it realizes its originally required speed. Another interesting approach is to reduce the supply voltage of each functional unit (thus reducing the power consumption, but increasing the delay of the unit) in the data path as much as possible while satisfying the timing requirements in terms of the cycle-time or throughput (in the case of pipelined circuits). These transformations include concurrency increasing transformations such as (time) loop unrolling and control flow optimizations and critical path reducing transformations such as retiming and pipelining. At the early stages of the behavioral design process, concurrency increasing transformations such as loop unrolling, pipelining and control flow optimization as well as critical path reducing transformations such as height minimization, retiming and pipelining may be used to allow a reduction in supply voltage without degrading system throughput; Algorithm-specific instruction sets may be utilized that boost code density and minimize switching.

Consider a module M in an RTL circuit that performs two operations A and B . The switching activity at the inputs of M , is determined by the number of bit flips between the values taken on by the variables that are inputs to the two operations, which in turn depend on the bit-level statistical characteristics of the variables. Hence, the power dissipation depends on the module binding. Similarly, consider a register R that is shared between two data values X and Y . The switching activity of R depends on the correlations between these two variables X and Y . Hence, the power dissipation depends on the register binding as well. These observations form the basis for power optimization during module and register allocation and binding in [14], [13], [45] and [46]. In [43], an exact (graph-theoretic) algorithm for minimizing the system power through variable-voltage scheduling is presented. The idea is to establish a supply voltage level for each of the operations in a data flow graph, thereby fixing the latency of that operation, such that the system timing constraint is met while power is minimized (because each operation will be executed using minimum possible supply voltage).

5.2.3 Logic Synthesis

Logic synthesis fits between the register transfer level and the net list of gates specification. It provides the automatic synthesis of net lists minimizing some objective function subject to various constraints. Example inputs to a logic synthesis system include two-level logic representation, multi-level Boolean net works, finite state machines and technology mapped circuits. Depending on the input specification (combinational versus sequential, synchronous versus asynchronous), the target implementation (two-level versus multi-level, unmapped versus mapped, ASICs versus FPGAs), the objective function (area, delay, power, testability) and the delay models used (zero-delay, unit-delay, unit-fanout delay, or library delay models), different techniques are applied to transform and optimize the original RTL description. The strategy for synthesizing circuits for low power consumption will be to restructure or optimize the circuit to obtain low switching activity factors at nodes which drive large capacitive loads. At the register-transfer (RT) and logic levels, symbolic states of a finite state machine (FSM) can be assigned binary codes to minimize the number of bit changes in the combinational logic for the most likely state transitions¹; Latches in a pipelined design can be repositioned to eliminate hazardous activity in the circuit [34]; Parts of the circuit that do not contribute to the present computation may be shut off completely; Output logic values of a circuit may be pre-computed one cycle before they are required and then used to reduce the internal switching activity of the circuit in the succeeding clock cycle [1]; Common sub-expressions with low transition

probability values can be extracted [25]; Network don't cares can be used to modify the input variable support and thus the local expression of a node so as to reduce the bit switching in the transitive fan-out of the node [28]; Nodes with high switching activity may be hidden inside CMOS gates where they drive smaller physical capacitances, Hazards/glitches in the circuit can be reduced by appropriate use of selective collapse, logic decomposition or delay insertion which lead to path balanced circuit structures; Circuit depth and power dissipation may be simultaneously minimized using a node clustering approach ; PLAs can be implemented to reduce static or dynamic power dissipation in pseudo-NMOS or dynamic NOR-NOR implementations [26]. Power dissipation may be further reduced by gate resizing [4], signal-to-pin assignment and I/O encoding.

5.2.4 Physical Design

Physical design fits between the netlist of gates specification and the geometric (mask) representation known as the layout. It provides the automatic layout of circuits minimizing some objective function subject to given constraints. Depending on the target design style (full-custom, standard-cell, gate arrays, FPGAs), the packaging technology (printed circuit boards, multi-chip modules, wafer-scale integration) and the objective function (area, delay, power, reliability), various optimization techniques are used to partition, place, resize and route gates. Under a zero-delay model, the switching activity of gates remains unchanged during layout optimization, and hence, the only way to reduce power dissipation is to decrease the load on high switching activity gates by proper net list partitioning and gate placement, gate and wire sizing, transistor reordering, and routing. At the same time, if a real-delay model is used, various layout optimization operations influence the hazard activity in the circuit. It should be noted that by applying post-layout optimization techniques (such as buffer and wire sizing, local restructuring and re-mapping, etc.), power can be further reduced. Under a zero-delay model, the switching activity of gates remains unchanged during layout optimization, and hence, the only way to reduce power dissipation is to decrease the load on high switching activity gates by proper netlist partitioning and gate placement, gate and wire sizing, transistor reordering, and routing. At the physical design level, power may be reduced by using appropriate net weights during netlist partitioning, floor planning, placement and routing; Individual transistors may be sized down to reduce the

power dissipation along the non-critical paths in a circuit; Large capacitive loads can be buffered using optimally sized inverter chains so as to minimize the power dissipation subject to a given delay constraint, Wire and driver sizing may be combined to reduce the interconnect delay with only a small increase in the power dissipation [13]; Clock trees may be constructed that minimize the load on the clock drivers subject to meeting a tolerable clock skew [14] [23].

5.2.5 Circuit Design

At the circuit level, power savings techniques that recycle the signal energies using the adiabatic switching principles rather than dissipating them as heat are promising in certain applications where speed can be traded for lower power [2]. Similarly, techniques based on combining self-timed circuits with a mechanism for selective adjustment of the supply voltage that minimizes the power while satisfying the performance constraints [38], those based on partial transfer of the energy stored on a capacitance to some charge sharing capacitance and then reusing this energy at a later time [21], and those based on electronic compensation for variations in VT thus making it possible to scale power supply voltages down to very low levels [8], show good signs . Design of energy efficient level-converters and DC/DC converters is also essential to the success of adaptive supply voltage strategies.

5.3. Power Management Strategies

In many synchronous applications a lot of power is dissipated by the clock. The clock is the only signal that switches all the time and it usually has to drive a very large clock tree. The circuit itself is partitioned in different blocks and each block is clocked with its own (derived) clock. The power savings that can be achieved this way are very application dependent, but can be significant. Power savings techniques that recycle the signal energies using the adiabatic switching principles rather than dissipating them as heat are promising in certain applications where speed can be traded for lower power. Similarly, techniques based on combining self-timed circuits with a mechanism for selective adjustment of the supply voltage that minimizes the power while satisfying the performance constraints show good signs.

6. CHALLENGES AHEAD

The need for lower power systems is being driven by many market segments. There are several approaches to reducing power, however the highest Return On Investment approach is through designing for low power. Unfortunately designing for low power adds another dimension to the already complex design problem; the

design has to be optimized for Power as well as Performance and Area. The successful development of new power conscious tools and methodologies requires a clear and measurable goal. In this context the research work should strive to reduce power by 5-10x in three years through design and tool development. To conclude this introduction, it is worthwhile to summarize the major challenges that, to our belief, have to be addressed if we want to keep power dissipation within bounds in the next generations of digital integrated circuits.

- A low voltage/low threshold technology and circuit design approach, targeting supply voltages around 1 Volt and operating with reduced thresholds.
- Low power interconnect, using advanced technology, reduced swing or reduced activity approaches.
- Dynamic power management techniques, varying supply voltage and execution speed according to activity measurements. This can be achieved by partitioning the design into sub-circuits whose energy levels can be independently controlled and by powering down sub-circuits which are not in use.
- System performance can be improved by moving the work to less energy constrained parts of the system, for example, by performing the task on fixed stations rather than mobile sites, by using asymmetric communication protocols, or unbalanced data compression schemes.
- Application specific processing. This might rely on the increased use of application specific circuits or application or domain specific processors.
- Move toward self-adjusting and adaptive circuit architectures that can quickly and efficiently respond to the environmental change as well as varying data statistics.
- An integrated design methodology - including synthesis and compilation tools. This might require the progression to higher level programming and specification paradigms (e.g. data flow or object oriented programming).
- Development of power conscious techniques and tools for behavioral synthesis, logic synthesis and layout optimization. The key requirements for these techniques are accurate and efficient estimation of the power cost of alternative organizations and / or implementations and the ability to minimize the power dissipation subject to given performance (or throughput in case of pipelined designs) constraints and supply voltage levels.
- Power savings techniques that recycle the signal energies using the adiabatic switching principles rather than dissipating them as heat are promising in certain applications where speed can be traded for lower power.

References

- [1] M. Alidina, J. Monteiro, S. Devadas, A. Ghosh, and M. Papaefthymiou. " Precomputation- based sequential logic optimization for low power." In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 57-62, April 1994.
- [2] W.C.Athas, L. J. Svensson, J.G.Koller, Thartzanis and E. Chou. " Low-Power Digital Systems Based on Adiabatic-Switching Principles. " *IEEE Transactions on VLSI Systems*, 2(4):398-407, December 1994
- [3] L. Benini, M. Favalli, and B. Ricco. " Analysis of hazard contribution to power dissipation in CMOS IC's. " In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 27-32, April 1994.
- [4] M. Berkelaar and J. Jess. " Gate sizing in MOS digital circuits with linear programming. " In *Proceedings of the European Design Automation Conference*, pages 217-221, 1990.
- [5] R. Bryant. " Graph-based algorithms for Boolean function manipulation. " *IEEE Transactions on Computers*, volume C-35, pages 677-691, August 1986.
- [6] R. Burch, F. Najm, P. Yang, and D. Hocevar. " Pattern independent current estimation for reliability analysis of CMOS circuits. " In *Proceedings of the 25th Design Automation Conference*, pages 294-299, June 1988.
- [7] R. Burch, F. N. Najm, P. Yang, and T. Trick. " A Monte Carlo approach for power estimation. " *IEEE Transactions on VLSI Systems*, 1(1):63-71, March 1993.
- [8] L. Richard Carley and Ihor Lys. " QuadRail: A design methodology for low power ICs. " *IEEE Transactions on VLSI Systems*, 2(4):383-390, December 1994.
- [9] A. Chandrakasan, S. Sheng, and R. W. Brodersen, "Low-power CMOS design. " *IEEE Journal of Solid-State Circuits*, pages 472-484, April 1992.
- [10] A. Chandrakasan, M. Potkonjak, J. Rabaey and R. W. Brodersen, " HYPER-LP: A System for Power Minimization Using Architectural Transformation. " In *Proceedings of the IEEE International Conference on Computer Aided Design*, pages 300-303, November 1992.
- [11] J-M. Chang and M. Pedram. " Low power register allocation and binding. " In *Proceedings of the 32nd Design Automation Conference*, pages 29-35, June 1995.
- [12] J-M. Chang and M. Pedram. " Power efficient module allocation and binding. " *CENG Technical Report 95-16*, University of Southern California, June 1995.
- [13] J. Cong, C-K. Koh and K-S. Leung. " Simultaneous driver and wire sizing for performance and power optimization. " *IEEE Transactions on VLSI Systems*, 2(4):408-425, December 1994.
- [14] J. Cong and C-K. Koh. " Minimum-cost bounded-skew clock routing. " In *Proceedings of the International Symposium on Circuits and Systems*, pages 215-218, 1995.
- [15] B. Davari, R. H. Dennard and G. G. Shahidi. " CMOS scaling for high performance and low power.. " *Proceedings of IEEE*, 83(4):408-425, April 1995.
- [16] C. Deng. " Power analysis for CMOS/BiCMOS circuits. " In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 3-8, April 1994.
- [17] S. Ercolani, M. Favalli, M. Damiani, P. Olivo, and B. Ricco. "Estimate of signal probability in combinational logic networks. " In *First European Test Conference*, pages 132-138, 1989.
- [18] T. A. Fjeldly and M. Shur. " Threshold voltage modeling and the subthreshold regime of operation of short-channel MOSFET's. " *IEEE Transactions on Electron Devices*, 40(1):137-145, Jan. 1993.
- [19] B. J. George, D. Gossain, S. C. Tyler, M. G. Wloka, and G. K. H. Yeap. " Power analysis and characterization for semi-custom design. " In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 215-218, April 1994.
- [20] A. Ghosh, S. Devadas, K. Keutzer, and J. White. " Estimation of average switching activity in combinational and sequential circuits. " In *Proceedings of the 29th Design Automation Conference*, pages 253-259, June 1992.

- [21] M. Hahm. " Modest power savings for applications dominated by switching of large capacitive loads. " In *Proceedings of the 1995 IEEE Symposium on Low Power Electronics*, pages 60-61, October 1995
- [22] M. Horowitz, T. Indermaur and R. Gonzalez. " Low-power digital design. " In *Proceedings of the 1995 IEEE Symposium on Low Power Electronics*, pages 8-11, October 1995.
- [23] D. J. Huang, A. B. Kahng and C. W. Tsao. " On the bounded-skew clock and Steiner tree problems. " In *Proceedings of the 32nd Design Automation Conference*, pages 508-513, June 1995.
- [24] S. Iman and M. Pedram. " Multi-level network optimization for low power. " In *Proceedings of the IEEE International Conference on Computer Aided Design*, pages 372-377, November 1994.
- [25] S. Iman and M. Pedram. " Logic extraction and decomposition for low power. " In *Proceedings of the 32nd Design Automation Conference*, June 1995.
- [26] S. Iman, C. Y. Tsui and M. Pedram. " PLA minimization for low power VLSI designs. " *CENG Technical Report*, Dept. of EE-Systems, University of Southern California, April 1995.
- [27] S. M. Kang. " Accurate simulation of power dissipation in VLSI circuits. " *IEEE Journal of Solid State Circuits*, 21(5):889-891, Oct. 1986.
- [28] N. Kumar, S. Katkooi, L. Rader and R. Vemuri. " Profile-driven behavioral synthesis for low power VLSI systems. " *To appear*, 1995.
- [29] P.E. Landman and J. Rabaey. " Power estimation for high level synthesis. " In *Proceedings of the European Conference on Design Automation*, pages 361-366, February 1993.
- [30] P.E. Landman and J. Rabaey. " Activity sensitive architectural power analysis for control path. " In *Proceedings of the 1995 International Symposium on Low Power Design*, pages 93-98, April 1995.
- [31] H. M. Lieberstein. " *A Course in Numerical Analysis*. " Harper & Row Publishers, 1968.
- [32] R. Marculescu, D. Marculescu, and M. Pedram. " Logic level power estimation considering spatiotemporal correlations. " In *Proceedings of the IEEE International Conference on Computer Aided Design*, pages 294-299, November 1994.
- [33] D. Marculescu, R. Marculescu, and M. Pedram. " Information theoretic measures for energy consumption at register transfer level. " In *Proceedings of the 1995 International Symposium on Low Power Design*, pages 81-86, April 1995.
- [34] J. Monteiro, S. Devadas, and A. Ghosh. " Retiming sequential circuits for low power. " In *Proceedings of the IEEE International Conference on Computer Aided Design*, pages 398-402, November 1993.
- [35] J. Monteiro, S. Devadas, and A. Ghosh. " Estimation of switching activity in sequential logic circuits with applications to synthesis for low power. " In *Proceedings of the 31st Design Automation Conference*, pages 12-17, June 1994.
- [36] F. N. Najm, R. Burch, P. Yang, and I. Hajj. " Probabilistic simulation for reliability analysis of CMOS VLSI circuits. " *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 9(4):439-450, April 1990.
- [37] F. N. Najm. " Towards a high-level power estimation capability. " In *Proceedings of the 1995 International Symposium on Low Power Design*, pages 87-92, April 1995.
- [38] L.S. Nielsen, C. Niessen, J. Sparso and C.H. van Berke. " Low-power operation using self-timed circuits and adaptive scaling of the supply voltage. " *IEEE Transactions on VLSI Systems*, 2(4):391-397, December 1994.
- [39] M. Pedram and N. Bhat. " Layout driven technology mapping. " In *Proceedings of the 28th Design Automation Conference*, pages 99-105, June 1991.
- [40] M. Pedram, B. T. Preas. " Interconnection length estimation for optimized standard cell layouts. " In *Proceedings of the IEEE International Conference on Computer Aided Design*, pages 390-393, November 1989.
- [41] S. R. Powell and P. M. Chau. " A model for estimating power dissipation in a class of DSP VLSI chips. " *IEEE Transactions on Circuits and Systems*, 36(6):646-650, June 1995.
- [42] R. A. Powers. " Batteries for low power electronics. " *Proceedings of IEEE*, 38(4):687-693, April 1995.
- [43] S. Raje and M. Sarrafzadeh. " Variable voltage scheduling. " In *Proceedings of the 1995 International Symposium on Low Power Design*, pages 9-13, April 1995.
- [44] S. Rajgopal and G. Mehta. " Experiences with simulation-based schematic level current estimation. " In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 9-14, April 1994.
- [45] A. Raghunathan and N. K. Jha. " Behavioral synthesis for low power. " In *Proceedings of the IEEE International Conference on Computer Design*, pages 318-322, October 1994.
- [46] A. Raghunathan and N. K. Jha. " An ILP formulation for low power based on minimizing switched capacitance during data path allocation. " *To appear*, 1995.
- [47] P. Schneider and U. Schlichtmann. " Decomposition of boolean functions for low power based on a new power estimation technique. " In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 123-128, April 1994.
- [48] C. Small, " Shrinking devices put the squeeze on system packaging. " *EDN*, vol. 39, no. 4, pages 41-46, Feb. 17, 1994.
- [49] C. Svensson and D. Liu. " A power estimation tool and prospects of power savings in CMOS VLSI chips. " In *Proceedings of the 1994 International Workshop on Low Power Design*, pages 171-176, April 1994.