

Precomputation-Based Sequential Logic Optimization for Low Power

Mazhar Alidina, José Monteiro, Srinivas Devadas, *Member, IEEE*,
Abhijit Ghosh, *Member, IEEE*, and Marios Papaefthymiou

Abstract—We address the problem of optimizing logic-level sequential circuits for low power. We present a powerful sequential logic optimization method that is based on selectively *precomputing* the output logic values of the circuit one clock cycle before they are required, and using the precomputed values to reduce internal switching activity in the succeeding clock cycle. We present two different precomputation architectures which exploit this observation. The primary optimization step is the synthesis of the precomputation logic, which computes the output values for a *subset* of input conditions. If the output values can be precomputed, the original logic circuit can be “turned off” in the next clock cycle and will have substantially reduced switching activity. The size of the precomputation logic determines the power dissipation reduction, area increase and delay increase relative to the original circuit. Given a logic-level sequential circuit, we present an automatic method of synthesizing precomputation logic so as to achieve maximal reductions in power dissipation. We present experimental results on various sequential circuits. Up to 75% reductions in average switching activity and power dissipation are possible with marginal increases in circuit area and delay.

I. INTRODUCTION

AVERAGE POWER DISSIPATION has recently emerged as an important parameter in the design of general-purpose and application-specific integrated circuits. Optimization for low power can be applied at many different levels of the design hierarchy. For instance, algorithmic and architectural transformations can trade off throughput, circuit area, and power dissipation [4], and logic optimization methods have been shown to have a significant impact on the power dissipation of combinational logic circuits [12].

In CMOS circuits, the probabilistic average switching activity of the circuit is a good measure of the average power

dissipation of the circuit. Average power dissipation can thus be computed by estimating the average switching activity. Several methods to estimate power dissipation for CMOS combinational circuits have been developed (e.g., [6], [10]). More recently, efficient and accurate methods of power dissipation estimation for sequential circuits have been developed [9], [13].

In this paper, we are concerned with the problem of optimizing logic-level sequential circuits for low power. Previous work in the area of sequential logic synthesis for low power has focused on state encoding [11] and retiming [8] algorithms. We present a powerful sequential logic optimization method that is based on selectively *precomputing* the output logic values of the circuit one clock cycle before they are required, and using the precomputed values to reduce internal switching activity in the succeeding clock cycle.

The primary optimization step is the synthesis of precomputation logic, which computes the output values for a *subset* of an input conditions. If the output values can be precomputed, the original logic circuit can be “turned off” in the next clock cycle and will not have any switching activity. Since the savings in the power dissipation of the original circuit is offset by the power dissipated in the precomputation phase, the selection of the subset of input conditions for which the output is precomputed is critical. The precomputation logic adds to the circuit area and can also result in an increased clock period.

Given a logic-level sequential circuit, we present an automatic method of synthesizing the precomputation logic so as to achieve a maximal reduction in power dissipation. We present experimental results on various sequential circuits. For some circuits, 75% reductions in average switching activity are possible with marginal increases in circuit area and delay.

In Section II, we briefly describe our model for power dissipation. In Section III we describe two different precomputation architectures. Algorithms that synthesize precomputation logic so as to achieve power dissipation reduction are presented in Section IV. In Section V we describe a method for multiple-cycle precomputation. Experimental results are presented in Section VI. In Section VII we describe additional precomputation architectures which are the subject of ongoing research.

II. A POWER DISSIPATION MODEL

Under a simplified model, the energy dissipation of a CMOS circuit is directly related to the switching activity.

Manuscript received June 15, 1994; revised August 23, 1994. The work of M. Alidina and S. Devadas was supported in part by the Defense Advanced Research Projects Agency under Contract N00014-91-J-1698 and by a NSF Young Investigator Award with matching funds from Mitsubishi and IBM Corporation. The work of J. Monteiro was supported by the Portuguese “Junta Nacional de Investigação Científica e Tecnológica” under project “Ciência”.

M. Alidina was with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA. He is now with AT&T Bell Laboratories, Allentown, PA 18103 USA.

J. Monteiro and S. Devadas are with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139 USA.

A. Ghosh is with Mitsubishi Electric Research Laboratories, Sunnyvale, CA 94086 USA.

M. Papaefthymiou is with the Department of Electrical Engineering, Yale University, New Haven, CT 06520 USA.

IEEE Log Number 9406366.

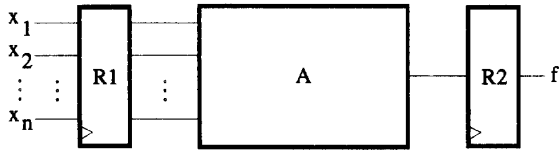


Fig. 1. Original Circuit.

In particular, the three simplifying assumptions are:

- The only capacitance in a CMOS logic-gate is at the output node of the gate.
- Either current is flowing through some path from V_{DD} to the output capacitor, or current is flowing from the output capacitor to ground.
- Any change in a logic-gate output voltage is a change from V_{DD} to ground or vice-versa.

All of these are reasonably accurate assumptions for well-designed CMOS gates [7], and when combined imply that the energy dissipated by a CMOS logic gate each time its output changes is roughly equal to the change in energy stored in the gate's output capacitance. If the gate is part of a synchronous digital system controlled by a global clock, it follows that the average power dissipated by the gate is given by:

$$P_{avg} = 0.5 \times C_{load} \times (V_{dd}^2 / T_{cyc}) \times E(transitions) \quad (1)$$

where P_{avg} denotes the average power, C_{load} is the load capacitance, V_{dd} is the supply voltage, T_{cyc} is the global clock period, and $E(transitions)$ is the *expected value* of the number of gate output transitions per global clock cycle [10], or equivalently the average number of gate output transitions per clock cycle. All of the parameters in (1) can be determined from technology or circuit layout information except $E(transitions)$, which depends on both the logic function being performed and the statistical properties of the primary inputs.

Equation (1) is used by the power estimation techniques such as [6], [10] to relate switching activity to power dissipation.

III. PRECOMPUTATION ARCHITECTURES

We describe two different precomputation architectures and discuss their characteristics in terms of their impact on power dissipation, circuit area and circuit delay.

A. First Precomputation Architecture

Consider the circuit of Fig. 1. We have a combinational logic block **A** that is separated by registers R_1 and R_2 . While R_1 and R_2 are shown as distinct registers in Fig. 1 they could, in fact, be the same register. We will first assume that block **A** has a single output and that it implements the Boolean function f .

In Fig. 2 the first precomputation architecture is shown. Two Boolean functions g_1 and g_2 are the *predictor* functions. We require:

$$g_1 = 1 \Rightarrow f = 1 \quad (2)$$

$$g_2 = 1 \Rightarrow f = 0 \quad (3)$$

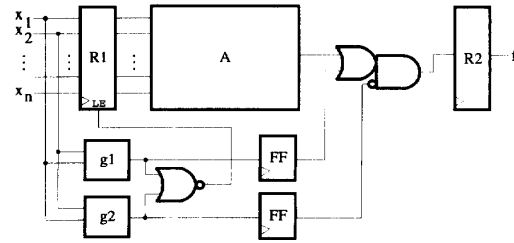


Fig. 2. First precomputation architecture.

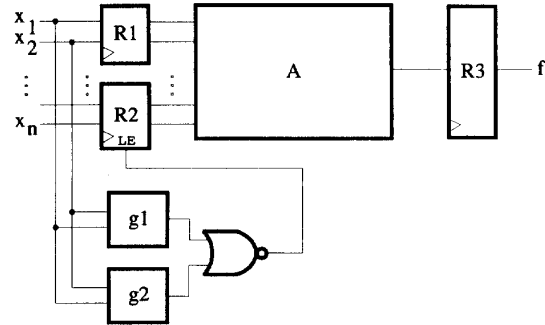


Fig. 3. Second precomputation architecture.

Therefore, during clock cycle t if either g_1 or g_2 evaluates to a 1, we set the load enable signal of the register R_1 to be 0. This means that in clock cycle $t+1$ the inputs to the combinational logic block **A** do not change. If g_1 evaluates to a 1 in clock cycle t , the input to register R_2 is a 1 in clock cycle $t+1$, and if g_2 evaluates to a 1, then the input to register R_2 is a 0. Note that g_1 and g_2 cannot both be 1 during the same clock cycle due to the conditions imposed by (2) and (3).

A power reduction in block **A** is obtained because for a subset of input conditions corresponding to $g_1 + g_2$ the inputs to **A** do not change implying zero switching activity. However, the area of the circuit has increased due to additional logic corresponding to g_1 , g_2 , the two additional gates shown in the figure, and the two flip-flops marked **FF**. The delay between R_1 and R_2 has increased due to the addition of the AND-NOR gate. Note also that g_1 and g_2 add to the delay of paths that originally ended at R_1 but now pass through g_1 or g_2 and the NOR gate before ending at the load enable signal of the register R_1 . Therefore, we would like to apply this transformation on noncritical logic blocks.

The choice of g_1 and g_2 is critical. We wish to include as many input conditions as we can in g_1 and g_2 . In other words, we wish to maximize the probability of g_1 or g_2 evaluating to a 1. In the extreme case this probability can be made unity if $g_1 = f$ and $g_2 = \bar{f}$. However, this would imply a duplication of the logic block **A** and no reduction in power with a twofold increase in area! To obtain reduction in power with marginal increases in circuit area and delay, g_1 and g_2 have to be significantly less complex than f . One way of ensuring this is to make g_1 and g_2 depend on significantly fewer inputs than f . This leads us to the second precomputation architecture of Fig. 3.

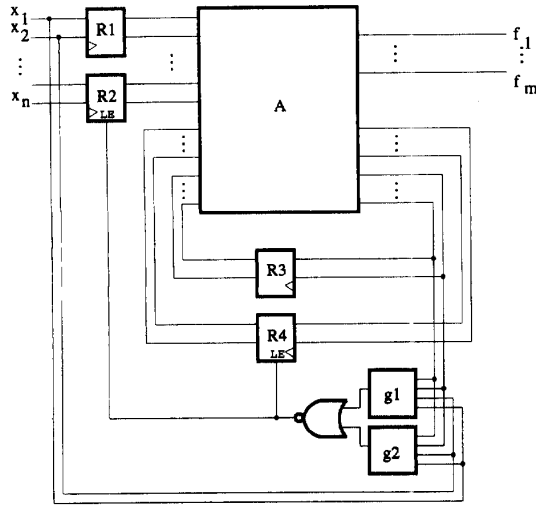


Fig. 4. Second precomputation architecture for a finite state machine.

B. Second Precomputation Architecture

In the architecture of Fig. 3, the inputs to the block **A** have been partitioned into two sets, corresponding to the registers R_1 and R_2 . The output of the logic block **A** feeds the register R_3 . The functions g_1 and g_2 satisfy the conditions of (2) and (3) as before, but g_1 and g_2 only depend on a subset of the inputs to f . If g_1 or g_2 evaluates to a 1 during clock cycle t , the load enable signal to the register R_2 is turned off. This implies that the outputs of R_2 during clock cycle $t+1$ do not change. However, since the outputs of register R_1 are updated, the function f will evaluate to the correct logical value. A power reduction is achieved because only a subset of the inputs to block **A** change implying reduced switching activity.

As before, g_1 and g_2 have to be significantly less complex than f and the probability of $g_1 + g_2$ begin a 1 should be high in order to achieve substantial power gains. The delay of the circuit between R_1/R_2 and R_3 unchanged, allowing precomputation of logic that is on the critical path. However, the delay of paths that originally ended at R_1 has increased.

The choice of inputs to g_1 and g_2 has to be made first, and then the particular functions that satisfy (2) and (3) have to be selected. Methods to perform this selection for this second precomputation architecture are described in Section IV.

C. Precomputation for Finite State Machines

As mentioned in Section III-A, these precomputation architectures are not restricted to pipeline circuits. We present in Fig. 4 an example of precomputation for a finite state machine using this second precomputation architecture.

D. An Example

We give an example that illustrates the fact that substantial power gains can be achieved with marginal increases in circuit area and delay. The circuit we are considering is a n -bit comparator that compares two n -bit numbers C and D and computes the function $C > D$. The optimized circuit with

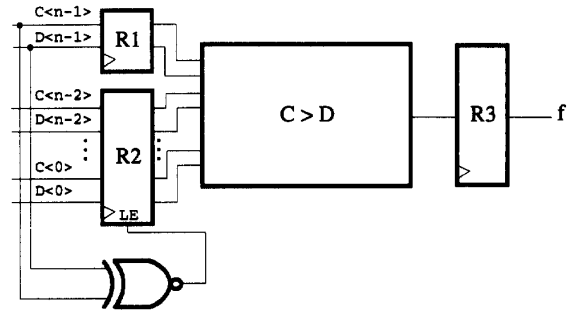


Fig. 5. A comparator example.

precomputation logic is shown in Fig. 5. The precomputation logic is as follows:

$$g_1 = C\langle n-1 \rangle \cdot \overline{D\langle n-1 \rangle}$$

$$g_2 = \overline{C\langle n-1 \rangle} \cdot D\langle n-1 \rangle$$

Clearly, when $g_1 = 1$, C is greater than D , and when $g_2 = 1$, C is less than D . We have to implement

$$\overline{g_1 + g_2} = C\langle n-1 \rangle \otimes D\langle n-1 \rangle$$

where \otimes stands for the exclusive-nor operator.

Assuming a uniform probability for the inputs¹, the probability that the XNOR gate evaluates to a 1 is 0.5, regardless of n . For large n , we can neglect the power dissipation in the XNOR gate, and therefore, we can achieve a power reduction of close to 50%. The reduction will depend upon the relative power dissipated by the vector pairs with $C\langle n-1 \rangle \otimes D\langle n-1 \rangle = 1$ and the vector pairs with $C\langle n-1 \rangle \otimes D\langle n-1 \rangle = 0$. If we add the inputs $C\langle n-2 \rangle$ and $D\langle n-2 \rangle$ to g_1 and g_2 it is possible to achieve a power reduction close to 75%.

IV. SYNTHESIS OF PRECOMPUTATION LOGIC

A. Introduction

In this section, we will describe methods to determine the functionality of the precomputation logic, and then describe methods to efficiently implement the logic.

We will focus on the second precomputation architecture (Section III-B) illustrated in Fig. 3. In order to ensure that the precomputation logic is significantly less complex than the combinational logic in the original circuit, we will restrict ourselves to identifying g_1 and g_2 such that they depend on a relatively small subset of the inputs to the logic block **A**.

B. Precomputation and Observability Don't-Cares

Assume that we have a logic function $f(X)$, with $X = \{x_1, \dots, x_n\}$, corresponding to block **A** of Fig. 1. Given that the logic function implemented by block **A** is f , then the *observability don't-care set* for input x_i is given by:

$$\text{ODC}_i = f_{x_i} \cdot \overline{f_{\overline{x_i}}} + \overline{f_{x_i}} \cdot f_{\overline{x_i}}$$

¹ The assumption here is that each $C\langle i \rangle$ and $D\langle i \rangle$ has a 0.5 static probability of being a 0 or a 1.

where f_{x_i} and $f_{\bar{x}_i}$ are the cofactors of f with respect to x_i , and similarly for \bar{f} . (The cofactors can be obtained simply by setting x_i to a 1 or 0 in f or \bar{f} .)

If we determine that a given input combination is in ODC_i then we can disable the loading of x_i into the register since that means that we do not need the value of x_i in order to know what the value of f is. If we wish to disable the loading of registers x_{m+1}, \dots, x_n , we will have to implement the function

$$g = \prod_{i=m+1}^n ODC_i$$

and use \bar{g} as the load enable signal for the registers corresponding to x_{m+1}, \dots, x_n .

For more details regarding observability don't-cares the reader is referred to [5] (p. 179).

C. Precomputation Logic

Let us now consider the architecture of Fig. 3. Assume that the inputs x_1, \dots, x_m , with $m < n$ have been selected as the variables that g_1 and g_2 depend on. We have to find g_1 and g_2 such that they satisfy the constraints of (2) and (3), respectively, and such that $\text{prob}(g_1 + g_2 = 1)$ is maximum.

We can determine g_1 and g_2 using universal quantification on f . The *universal quantification* of a function f with respect to a variable x_i is defined as:

$$U_{x_i} f = f_{x_i} \cdot f_{\bar{x}_i}$$

This gives all the combinations over the inputs $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$, that result in $f = 1$ such that the value of x_i does not matter.

Given a subset of inputs $S = \{x_1, \dots, x_m\}$, set $D = X - S$. We can define:

$$U_D f = U_{x_{m+1}} \dots U_{x_n} f$$

Theorem 4.1: $g_1 = U_D f$ satisfies (2). Further, no function $h(x_1, \dots, x_m)$ exists such that $\text{prob}(h = 1) > \text{prob}(g_1 = 1)$ and such that $h = 1 \Rightarrow f = 1$.

Proof: By construction, if for some input combination a_1, \dots, a_m causes $g_1(a_1, \dots, a_m) = 1$, then for that combination of x_1, \dots, x_m and all possible combinations of variables in x_{m+1}, \dots, x_n $f(a_1, \dots, a_m, x_{m+1}, \dots, x_n) = 1$.

We cannot add any minterm over x_1, \dots, x_m to g_1 because for any minterm that is added, there will be some combination of x_{m+1}, \dots, x_n for which $f(x_1, \dots, x_n)$ will evaluate to a 0. Therefore, we cannot find any function h that satisfies (2) and such that $\text{prob}(h = 1) > \text{prob}(g_1 = 1)$. \square

Similarly, given a subset of inputs S , we can obtain a maximal g_2 by:

$$g_2 = U_D \bar{f} = U_{x_{m+1}} \dots U_{x_n} \bar{f}$$

We can compute the functionality of the precomputation logic as $g_1 + g_2$.

Selecting a Subset of Inputs: Exact Method: Given a function f we wish to select the "best" subset of inputs S of cardinality k . Given S , we have $D = X - S$ and we compute $g_1 = U_D f$, $g_2 = U_D \bar{f}$. In the sequel, we assume that the best set of inputs corresponds to the inputs which result in

```

SELECT_INPUTS( f, k ):
{
  /* f = function to precompute */
  /* k = # of inputs to precompute with */
  BEST_PROB = 0 ;
  SELECTED_SET =  $\phi$  ;
  SELECT_RECUR( f,  $\bar{f}$ ,  $\phi$ , X, |X| - k ) ;
  return( SELECTED_SET ) ;
}

SELECT_RECUR( f_a, f_b, D, Q, l ):
{
  if( |D| + |Q| < l )
    return ;
  pr = prob( f_a = 1 ) + prob( f_b = 1 ) ;
  if( pr  $\leq$  BEST_PROB )
    return ;
  else if( |D| == l ) {
    BEST_PROB = pr ;
    SELECTED_SET = X - D ;
    return ;
  }
  choose  $x_i \in Q$  such that  $i$  is minimum ;
  SELECT_RECUR( U_{x_i} f_a, U_{x_i} f_b, D  $\cup$   $x_i$ , Q -  $x_i$ , l ) ;
  SELECT_RECUR( f_a, f_b, D, Q -  $x_i$ , l ) ;
  return ;
}

```

Fig. 6. Procedure to determine the optimal set of inputs.

$\text{prob}(g_1 + g_2 = 1)$ being maximum for a given k . We know that $\text{prob}(g_1 + g_2 = 1) = \text{prob}(g_1 = 1) + \text{prob}(g_2 = 1)$ since g_1 and g_2 cannot both be 1 on the same input vector. The above cost function ignores the power dissipated in the precomputation logic, but since the number of inputs to the precomputation logic is significantly smaller than the total number of inputs this is a good approximation.

In the sequel we describe a branching algorithm that determines the optimal set of inputs maximizing the probability of the g_1 and g_2 functions. This algorithm is shown in pseudocode in Fig. 6.

The procedure **SELECT_INPUTS** receives as arguments the function f and the desired number of inputs k to the precomputation logic. **SELECT_INPUTS** calls the recursive procedure **SELECT_RECUR** with five arguments. The first two arguments corresponds to the g_1 and g_2 functions, which are initially f and \bar{f} . A variable is selected within the recursive procedure and the two functions are universally quantified with respect to the selected variable. The third argument D corresponds to the set of variables that g_1 and g_2 do not depend on. The fourth argument Q corresponds to the set of "active" variables, which may be selected or discarded. Finally, the argument l correspond to the number of variables that have to be universally quantified in order to obtain g_1 and g_2 with k or fewer inputs.

If $|D| + |Q| < l$ it means that we have selected too many variables in the earlier recursions and we will not be able to quantify with respect to enough input variables. The functions g_1 and g_2 will depend on too many variables ($> k$).

We calculate the probability of $g_1 + g_2$. If this probability is less than the maximum probability we have encountered thus far, we can immediately return because of the following invariant which is true because f contains $U_{x_i}f$.

$$\text{prob}(U_{x_i}f) = \text{prob}(f_{x_i} \cdot \overline{f_{x_i}}) \leq \text{prob}(f) \quad \forall x_i, f$$

Therefore as we universally quantify variables from a given f_a and f_b function pair the pr quantity monotonically decreases.

We store the selected set corresponding to the maximum encountered probability.

Selecting a Subset of Inputs: Approximate Method: The worst-case running time of the exact method is exponential in the number of input variables and, although we have a nice pruning condition, there are many examples for which we cannot use it. Thus we have also implemented an approximate algorithm that looks at each input individually and chooses the k most promising inputs.

For each input we calculate:

$$p_i = \text{prob}(U_{x_i}f) + \text{prob}(U_{x_i}\overline{f})$$

p_i is the probability that we know the value of f without knowing the value of x_i . If p_i is high then most of the time we do not need x_i to compute f . Therefore we select the k inputs corresponding to smaller values of p_i .

Implementing the Logic: The Boolean operations of OR and universal quantification required in the input selection procedure can be carried out efficiently using reduced, ordered Binary Decision Diagrams (ROBDDs) [3]. We obtain a ROBDD for the $g_1 + g_2$ function. A ROBDD can be converted into a multiplexor-based network (see [1]) or into a sum-of-products cover. The network or cover can be optimized using standard combinational logic optimization methods that reduce area [2] or those that target low power dissipation [12].

D. Multiple-Output Functions

In general, we have a multiple-output function f_1, \dots, f_m that corresponds to the logic block **A** in Fig. 1. All the procedures described thus far can be generalized to the multiple-output case.

The functions g_{1i} and g_{2i} are obtained using the equations below.

$$\begin{aligned} g_{1i} &= U_D f_i \\ g_{2i} &= U_D \overline{f_i} \end{aligned}$$

where $D = X - S$ as before. The function g whose complement drives the load enable signal is obtained as:

$$g = \prod_{i=1}^m (g_{1i} + g_{2i})$$

The function g corresponds to the set of input conditions where the variables in S control the values of *all* the f_i 's regardless of the values of variables in $D = X - S$.

Selecting a Subset of Outputs: Exact Method: We describe an algorithm, which given a multiple-output function, selects a subset of outputs *and* a subset of inputs so as to maximize a given cost function that is dependent on the probability of

```

SELECT_OUTPUTS( F = {f1, ..., fm}, k ):
{
  /* F = multi-output function to precompute */
  /* k = # of inputs to precompute with */
  BEST_COST = 0 ;
  SEL_OP_SET = φ ;
  SELECT_ORECUR( φ, F, 1, k ) ;
  return( SEL_OP_SET ) ;
}

SELECT_ORECUR( G, H, proldG, k ):
{
  lf = gates(F - H)/total_gates × proldG ;
  if( lf ≤ BEST_COST )
    return ;
  BEST_PROB = total_gates/gates(F - H) × BEST_COST ;
  if( G ≠ φ )
    if( SELECT_INPUTS( G, k ) == φ )
      return ;
  prG = BEST_PROB ; /* BEST_PROB is set in SELECT_INPUTS */
  cost = prG × gates(G)/total_gates ;
  if( cost > BEST_COST ) {
    BEST_COST = cost ;
    SEL_OP_SET = G ;
  }
  choose fi ∈ H such that i is minimum ;
  SELECT_ORECUR( G ∪ fi, H - fi, prG, k ) ;
  SELECT_ORECUR( G, H - fi, prG, k ) ;

  return ;
}

```

Fig. 7. Procedure to determine the optimal set of outputs.

the precomputation logic and the number of selected outputs. This algorithm is described in pseudocode in Fig. 7.

The inputs to procedure **SELECT_OUTPUTS** are the multiple-output function F , and a number k corresponding to the number of inputs to the precomputation logic.

The procedure **SELECT_ORECUR** receives as inputs two sets G and H , which correspond to the current set of outputs that have been selected and the set of outputs which can be added to the selected set, respectively. Initially, $G = \phi$ and $H = F$. The cost of a particular selection of outputs, namely G , is given by $\text{pr}G \times \text{gates}(F - H)/\text{total_gates}$, where $\text{pr}G$ corresponds to the signal probability of the precomputation logic, $\text{gates}(F - H)$ corresponds to the number of gates in the logic corresponding to the outputs in G and not shared by any output in H , and total_gates corresponds to the total number of gates in the network (across all outputs of F).

There are two pruning conditions that are checked for in the procedure **SELECT_ORECUR**. The first corresponds to assuming that all the outputs in H can be added to G without decreasing the probability of the precomputation logic. This is a valid condition because the quantity $\text{prold}G$ in each recursive call can only decrease with the addition of outputs of G . We then set a lower bound on the probability of the precomputation logic prior to calling the input selection procedure. Optimistically assuming that all the outputs in H can be added to G without lowering the precomputation logic probability, we are not interested in a precomputation logic probability for G that would result in a cost that is equal to or lower than **BEST_COST**.

Logic Duplication: Since we are only precomputing a subset of outputs, we may incorrectly evaluate the outputs that we are *not* precomputing as we disable certain inputs during particular clock cycles. If an output that is not being precom-

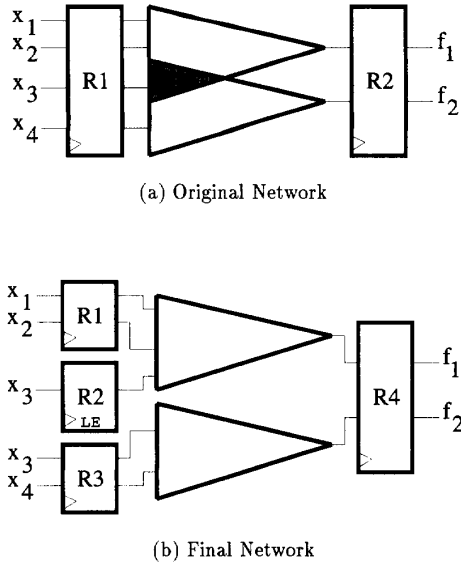


Fig. 8. Logic duplication in a multiple-output function.

puted depends on an input that is being disabled, then the output will be incorrect.

The **support** of f , denoted as $support(f)$, is the set of all variables x_i that occur in f as x_i or \bar{x}_i . Once a set of outputs $G \subset F$ and a set of precomputation logic inputs $S \subset X$ have been selected, we need to duplicate the registers corresponding to $(support(G) - S) \cap support(F - G)$. The inputs that are being disabled are in $support(G) - S$. Logic in the $F - G$ outputs that depends on the set of duplicated inputs has to be duplicated as well. It is precisely for this reason that we maximize $prG \times gates(F - H)/total_gates$ rather than prG in the output-selection algorithm. This way we are maximizing the number of gates (logic corresponding to the outputs in G) that will not switch when precomputation is possible but not taking into account gates that are shared by the outputs in H , thus reducing the amount of duplication as much as possible.

An example of a multiple-output function where the registers and logic need to be duplicated is shown in Fig. 8.

The original network has outputs f_1 and f_2 and inputs x_1, \dots, x_4 . The function f_1 depends on inputs x_1, \dots, x_3 and the function f_2 depends on inputs x_3 and x_4 . Hence, the two outputs are sharing the input x_3 . Suppose that the output-selection procedure determines that f_1 is the best output to precompute and that inputs x_1 and x_2 are the best inputs to the precomputation logic. Therefore, just as in the case of a single-output function, the inputs x_1 and x_2 feed the input register, whereas, x_3 feeds the register with the load-enable signal. However, since f_2 depends on x_3 and the register with the load-enable signal contains stale values in some clock cycles, we need to duplicate the register for x_3 and the logic from x_3 to f_2 .

Selecting a Subset of Outputs: Approximate Method: Again the exact algorithm is worst-case exponential in the number of inputs plus number of outputs, thus we need an approximate method to handle larger circuits. We designed an approximate algorithm which is presented in pseudocode in Fig. 9.

```

SELECT_OUTPUTS_APPROX(  $F = \{f_1, \dots, f_m\}, k$  ):
{
  BEST_COST = 0 ;
  foreach  $x_i \in X$  { /* Output selection */
    foreach  $f_j \in F$  {
       $g_j = U_{x_i}f_j + U_{x_i}\bar{f}_j$  ;
    }
    foreach  $f_j \in F$  {
       $G = \{f_j\}$  ;
       $H = F - \{f_j\}$  ;
       $probG = prob(g_j)$  ;
       $curr\_cost = probG \times gates(F - H)/total\_gates$  ;

      /* Add any outputs that make the cost increase */
       $g = g_j$  ;
      foreach  $f_i \in F$  {
         $G = G \cup \{f_i\}$  ;
         $probG = prob(g \cdot g_i)$  ;
         $cost = probG \times gates(F - H)/total\_gates$  ;
        if(  $cost > curr\_cost$  ) {
           $curr\_cost = cost$  ;
           $g = g \cdot g_i$  ;
        } else
           $G = G - \{f_i\}$  ;
      }
    }
  }
  if(  $curr\_cost > BEST\_COST$  ) {
    BEST_COST =  $curr\_cost$  ;
    SEL.OP.SET =  $G$  ;
  }
}
foreach  $x_i \in X$  { /* Input selection */
   $g = 1$  ;
  foreach  $f_j \in SEL.OP.SET$ 
     $g = g \cdot (U_{x_i}f_j + U_{x_i}\bar{f}_j)$  ;
   $p_i = prob(g)$  ;
}
select the  $k$   $x_i$ 's corresponding to smaller  $p_i$ 's
}

```

Fig. 9. Procedure to determine a good set of outputs.

In this algorithm we first select the set of outputs that will be precomputed and then select the inputs that we are going to precompute those outputs with. When we are selecting the outputs we still do not know which inputs are going to be selected, thus we select those outputs that seem to be the *most precomputable*. Universally quantifying just one of the inputs, we start with one output and compute the same cost function as in the exact method, $prG \times gates(F - H)/total_gates$. Then we add outputs that make the cost function increase. We repeat this process for each input. At the end we keep the set of outputs corresponding to the maximum cost.

Once we have a set of promising outputs to precompute we can use the approximate algorithm described in Section IV-C-2 to select the inputs. This algorithm runs in polynomial time in the number outputs times the number of inputs.

V. MULTIPLE CYCLE PRECOMPUTATION

A. Basic Strategy

It is possible to precompute output values that are not required in the succeeding clock cycle, but required 2 or more clock cycles later.

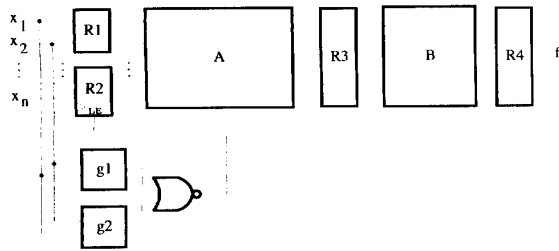


Fig. 10. Multiple cycle precomputation.

Consider the topology of Fig. 10. If the register outputs of R_3 are not used except to compute f , then we can precompute the value of the function f using a selected set of inputs, namely those corresponding to register R_1 . If f can be precomputed to a 1 or a 0 for a set of input conditions, then for these inputs we can turn off the load enable signal to R_2 . This will reduce switching activity not only in logic block A , but also in logic block B , because there will be reduced switching activity at the outputs of R_3 in the clock cycle following the one where the outputs of R_2 do not change.

B. Examples

We give examples illustrating multiple-cycle precomputation.

Consider the circuit of Fig. 11. The function f computes $(C + D) > (X + Y)$ in two clock cycles². Attempting to precompute $C + D$ or $X + Y$ using the methods of the previous section do not result in any savings because there are too many outputs to consider. However, 2-cycle precomputation can reduce switching activity by close to 12.5% if the functions below are used.

$$g_1 = C\langle n-1 \rangle \cdot D\langle n-1 \rangle \cdot \overline{X\langle n-1 \rangle} \cdot \overline{Y\langle n-1 \rangle}$$

$$g_2 = \overline{C\langle n-1 \rangle} \cdot \overline{D\langle n-1 \rangle} \cdot X\langle n-1 \rangle \cdot Y\langle n-1 \rangle$$

where g_1 and g_2 satisfy the constraints of (2) and (3), respectively. Since $\text{prob}(g_1 + g_2) = \frac{2}{16} = 0.125$, we can disable the loading of registers $C\langle n-2:0 \rangle$, $D\langle n-2:0 \rangle$, $X\langle n-2:0 \rangle$, and $Y\langle n-2:0 \rangle$ 12.5% of the time, which results in switching activity reduction. This percentage can be increased to over 45% by using $C\langle n-2 \rangle$ through $Y\langle n-2 \rangle$. We can additionally use single-cycle precomputation logic (as illustrated in Fig. 5) to further reduce switching activity in the $>$ comparator of Fig. 11.

Next, consider the circuit of Fig. 12. The multiple-output function f computes $\text{MAX}(C + D, X + Y)$ in two clock cycles. We can use exactly the same g_1 and g_2 functions as those immediately above, but g_1 is used to disable the loading of registers $X\langle n-2:0 \rangle$ and $Y\langle n-2:0 \rangle$, and g_2 is used to disable the loading of $C\langle n-2:0 \rangle$ and $D\langle n-2:0 \rangle$. We exploit the fact that if we know that $C + D > X + Y$, there is no need to compute $X + Y$, and vice versa. Finally, we can implement the MAX function as shown in Fig. 13, duplicate registers and use single-cycle precomputation on the $>$ operator (as illustrated in Fig. 5) to achieve switching activity reduction.

²+ in the figure stands for addition.

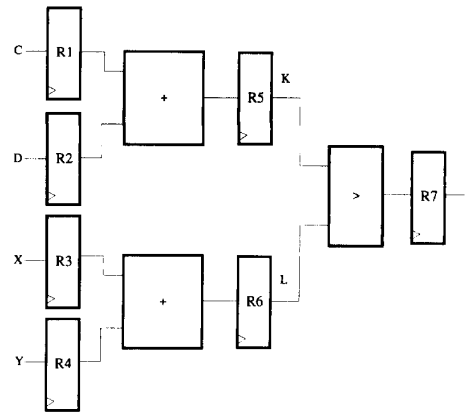


Fig. 11. Adder-comparator circuit.

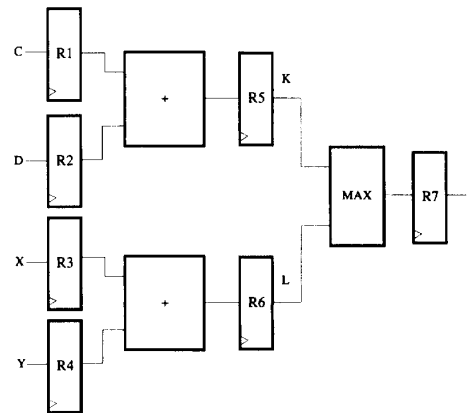


Fig. 12. Adder-maximum circuit.

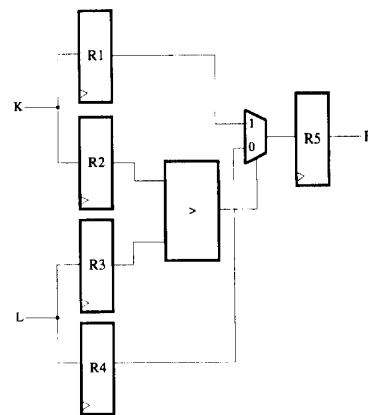


Fig. 13. Precomputation applied to maximum circuit.

VI. EXPERIMENTAL RESULTS

We first present results on datapath circuits such as carry-select adders, comparators, and interconnections of adders and comparators in Table I. In all examples the precomputation architecture of Fig. 3 was used and all the outputs of each circuit were precomputed. For each circuit, the number of

TABLE I
POWER REDUCTIONS FOR DATAPATH CIRCUITS

CKT	Original			Precompute Logic				Optimized	
	Literals	Levels	Power	Bits	Literals	Levels	Power	% Red	
comp16	286	7	1281	2	4	2	965	25	
				4	8	2	683	47	
				6	12	2	550	57	
				8	16	2	518	60	
				10	20	2	538	58	
add_comp16	3026	8	6941	4/0	8	2	6346	9	
				4/8	24	4	5711	18	
				8/0	51	4	4781	31	
				8/8	67	6	3933	43	
max16	350	9	1744	8	16	2	1281	27	
add_max16	3090	9	7370	4/0	8	2	7174	3	
				4/8	24	4	6751	8	
				8/0	51	4	6624	10	
				8/8	67	6	6116	17	
csa16	975	10	2945	2	4	2	2958	0	
				4	11	4	2775	6	
				6	18	4	2676	9	
				8	25	5	2644	10	

literals, levels of logic and power of the original circuit, the number of inputs, literals and levels of the precompute logic, the final power and the percent reduction in power are shown. All power estimates are in micro-Watts and are computed using the techniques described in [6], [9]. A zero delay model and a clock frequency of 20 MHz was assumed. The rugged script of *sis* was used to optimize the precompute logic.

Power dissipation decreases for almost all cases. For circuit **comp16**, a 16-bit parallel comparator, the power decreases by as much as 60% when 8 inputs are used for precomputation. Multiple-cycle precomputation results are given for circuits **add_comp16** and **add_max16**, shown in Figs. 11 and 12, respectively. For circuit **add_comp16**, for instance, the numbers 4/8 under the fifth column indicates that four inputs are used to precompute the adders in the first cycle and eight inputs are used to precompute the comparator in the next cycle.

The number of levels of the precompute logic is an indication of the performance penalty in using precomputation. The logic that is driving the input flip-flops to the original circuit is increased in depth by the number of levels of the precompute logic. In most cases, the increase in the number of levels is small.

Results on random logic circuits are presented in Table II. The random logic circuits are taken from the MCNC combinational benchmark sets. We have presented results for those examples where significant savings in power was obtained. Again, the second precomputation architecture was used and the input and output selection algorithms described in Section IV were used. Due to the size of the circuits, on most examples the approximate algorithm was used. Circuits for which the exact algorithm was used are marked with *. The columns in this table have the same meaning as in Table I, except for the second and third columns which show the number of inputs and outputs of each circuit, and the eighth column which shows the number of outputs that are being precomputed. It

TABLE II
POWER REDUCTIONS FOR RANDOM LOGIC CIRCUITS

Circuit	Original				Precompute Logic				Optimized		
	I	O	Lits	Levels	Power	I	O	Lits	Levels	Power	% Red
apex2	39	3	395	11	2387	4	3	4	1	1378	42
cht	47	36	167	3	1835	1	35	1	1	1537	16
cm138*	6	8	35	2	286	3	8	3	1	153	47
cm150*	21	1	61	4	744	1	1	1	1	574	23
cmb*	16	4	62	5	620	5	4	10	1	353	43
comp	32	3	185	6	1352	6	3	13	2	627	54
cordic*	23	2	194	13	1049	10	2	18	2	645	39
cps	24	109	1203	9	3726	7	101	26	3	2191	41
dalu	75	16	3067	24	11048	5	16	12	2	7344	34
duke2	22	29	424	7	1732	9	29	24	3	1328	23
e64	65	65	253	32	2039	5	65	5	1	513	75
i2	201	1	230	3	5606	17	1	42	5	1943	65
majority*	5	1	12	3	173	1	1	1	1	141	19
misex2	25	18	113	5	976	8	18	16	3	828	15
misex3	25	18	626	14	2350	2	14	2	1	1903	19
mux*	21	1	54	5	715	1	1	0	0	557	22
pcie	19	9	71	7	692	3	9	3	1	486	30
pcler8	27	17	95	8	917	3	17	3	1	571	38
sao2*	10	4	270	17	1191	2	4	2	1	422	65
seq	42	35	1724	11	6112	2	35	1	1	2134	65
spla	16	46	634	9	2267	4	46	6	1	1340	41
term1	34	10	625	9	3605	8	10	14	3	2133	41
too_large	38	3	491	11	2718	1	3	1	1	1756	35
unreg	36	16	144	2	1499	2	15	2	1	1234	18

*Precompute logic calculated using the exact algorithm.

is noteworthy that in some cases, as much as 75% reduction in power dissipation is obtained.

The area penalty incurred is indicated by the number of literals in the precomputation logic and is 3% on the average. The extra delay incurred is proportional to the number of levels in the precomputation logic and is quite small in most cases. It should be noted that it may be possible to use the other precomputation architectures for all of the examples presented here. Some of these examples are perhaps better suited to other architectures than the one we used to derive the results, and therefore larger savings in power may be possible. Secondly, the inputs and outputs to be selected and the precomputation logic are determined automatically, making this approach suitable for automatic logic synthesis systems. Finally, the significant power savings obtained for random logic circuits indicate that this approach is not restricted only to datapath circuits.

VII. OTHER PRECOMPUTATION ARCHITECTURES

In this section, we describe additional precomputation architectures. We first present an architecture that is applicable to *all* logic circuits and does not require, for instance, that the inputs should be in the observability don't-care set in order to be disabled, which was the case for the architectures shown in Section III. We also extend precomputation so that it can be used in combinational logic circuits.

A. Multiplexor-Based Precomputation

All logic functions can be written in a Shannon expansion. For the function f with inputs $X = \{x_1, \dots, x_n\}$ we can

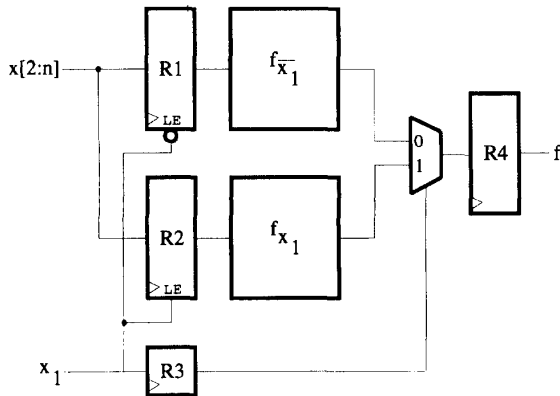


Fig. 14. Precomputation using the Shannon expansion.

write:

$$f = x_1 \cdot f_{x_1} + \bar{x}_1 \cdot f_{\bar{x}_1} \quad (4)$$

where f_{x_1} and $f_{\bar{x}_1}$ are the cofactors of f with respect to x_1 .

Fig. 14 shows an architecture based on (4). We implement the functions f_{x_1} and $f_{\bar{x}_1}$. Depending on the value of x_1 , only one of the cofactors is computed while the other is disabled by setting the load-enable signal of its input register. The input x_1 drives the select line of a multiplexer which chooses the correct cofactor.

The main advantage of this architecture is that it applies to *all* logic functions. The input x_1 in the example was chosen for the purpose illustration. In fact, any input x_1, \dots, x_n could have been selected. Unlike the architectures described earlier, we do not require that the inputs being disabled should be don't-cares for the input conditions which we are precomputing. In other words, the inputs being disabled do not have to be in the observability don't-care set. A disadvantage of this architecture is that we need to duplicate the registers for the inputs not being used to turn off part of the logic. On the other hand, no precomputation logic functions have been added to the circuit.

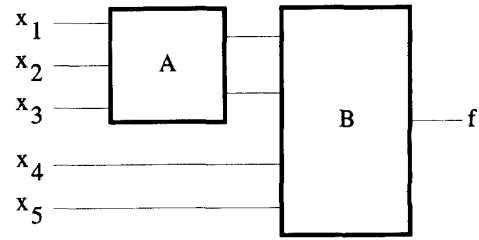
The algorithm to select the best input for this architecture is also quite different. We will not discuss this algorithm in detail, except to mention that in this case, we are interested in finding the input that yields the most area efficient f_{x_1} and $f_{\bar{x}_1}$ functions.

B. Combinational Logic Precomputation

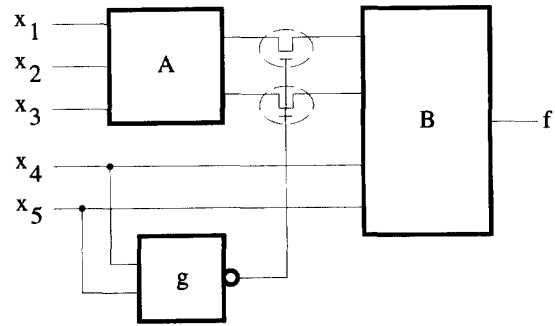
The architectures described so far apply only to sequential circuits. We now describe precomputation of combinational circuits.

Suppose we have some combinational logic function f composed of two subfunctions **A** and **B** as shown in Fig. 15(a). Suppose we also want to precompute this function with the inputs x_4 and x_5 . Fig. 15(b) shows how this can be accomplished. For simplicity, pass transistors are shown, however, we have several choices as to what to use within the dotted circles instead of the pass transistors.

Transmission Gates: Assume that transmission gates are used in place of the pass transistors in Fig. 15(b). The function



(a) Original Network



(b) Final Network

Fig. 15. Combinational logic precomputation.

g with inputs x_4 and x_5 drives the transmission gates. As in the previous architectures, $g = g_1 + g_2$. Hence, when g is a 0, the transmission gates are turned off and the new values of logic block **A** are prevented from propagating into logic block **B**. The inputs x_4 and x_5 are also inputs to the logic block **B** just as in the original network in order to ensure that the output is set correctly.

For the combinational architecture, there is an implied delay constraint, i.e., the transmission gates should be off *before* the new values of **A** are computed. In the example shown, the worst-case delay of the g block plus the arrival time of inputs x_4 or x_5 should be less than the best-case delay of logic block **A** plus the arrival time of the inputs x_1, x_2 , or x_3 . The *arrival time* of an input is defined as the time at which the input settles to its steady state value [5]. If the delay constraint is not met, then it may be necessary to delay the x_1, x_2 and x_3 inputs with respect to the x_4 and x_5 inputs in order to get the switching activity reduction in logic block **B**.

Transparent Latches: A violation of the delay constraint described immediately above can result in nodes in the circuit being stuck at metastable states (halfway between the supply voltages) causing excessive power dissipation. In order to ensure that this does not occur, transparent latches can be used instead of transmission gates. This results in increased overhead for precomputation. Note that a violation of the delay constraint may cause glitching in the circuit, but the nodes will settle to the supply voltages.

AND Gates: One can also replace the pass transistor with an AND gate. This will reduce switching activity, though not

as much as in the transparent latch case. This is because g may make a $0 \rightarrow 1$ transition during a clock cycle, possible causing unnecessary $1 \rightarrow 0$ transitions at the outputs of the AND gates. This option works best for percharged logic.

VIII. CONCLUSION AND ONGOING WORK

We have presented a method of precomputing the output response of a sequential circuit one clock cycle before the output is required, and exploited this knowledge to reduce power dissipation in the succeeding clock cycle. Several different architectures that utilize precomputation logic were presented.

In a finite state machine there is typically a single register, whose inputs are combinational functions of the register outputs. The precomputation architectures make no assumptions regarding feedback. For instance, R_1 and R_2 in Fig. 2 can be the same register.

Precomputation increases circuit area and can adversely impact circuit performance. In order to keep area and delay increases small, it is best to synthesize precomputation logic which depends on a small set of inputs.

Precomputation works best when there are a small number of complex functions corresponding to the logic block **A** of Figs. 2 and 3. If the logic block has a large number of outputs, then it may be worthwhile to selectively apply precomputation-based power optimization to a small number of complex outputs. This selective partitioning will entail a duplication of combinational logic and registers, and the savings in power is offset by this duplication.

Other precomputation architectures are being explored, including the architectures of Section VII, and those that rely on a history of previous input vectors. More work is required in the automation of a logic design methodology that exploits multiplexor-based, combinational and multiple-cycle precomputation.

ACKNOWLEDGMENT

The authors would like to thank A. Chandrakasan for providing us with information regarding power dissipation in registers and P. Vanbekbergen for pointing out that transparent latches should be used in Fig. 15(b).

REFERENCES

- [1] P. Ashar, S. Devadas, and K. Keutzer, "Path-delay-fault testability properties of multiplexor-based networks," *Integration, the VLSI J.*, vol. 15, no. 1, pp. 1–23, July 1993.
- [2] R. Brayton, R. Rudell, A. Sangiovanni-Vincentelli, and A. Wang, "MIS: A multiple-level logic optimization system," in *IEEE Trans. Computer-Aided Design*, vol. CAD-6, pp. 1062–1081, Nov. 1987.
- [3] R. Bryant, "Graph-based algorithms for Boolean function manipulation," *IEEE Trans. Comput.*, vol. C-35, no. 8, pp. 677–691, Aug. 1986.
- [4] A. Chandrakasan, T. Sheng, and R. W. Brodersen, "Low power CMOS digital design," in *IEEE J. Solid-State Circ.*, pp. 473–484, Apr. 1992.
- [5] S. Devadas, A. Ghosh, and K. Keutzer, *Logic Synthesis*. New York: McGraw Hill, 1994.
- [6] A. Ghosh, S. Devadas, K. Keutzer, and J. White, "Estimation of average switching activity in combinational and sequential circuits," in *Proc. 29th Design Automation Conf.*, June 1992, pp. 253–259.
- [7] L. Glasser and D. Dobberpuhl, *The Design and Analysis of VLSI Circuits*. Reading, MA: Addison-Wesley, 1985.
- [8] J. Monteiro, S. Devadas, and A. Ghosh, "Retiming sequential circuits for low power," in *Proc. Int. Conf. Computer-Aided Design*, Nov. 1993, pp. 398–402.
- [9] J. Monteiro, S. Devadas, and B. Lin, "A methodology for efficient estimation of switching activity in sequential logic circuits," in *Proc. 31st Design Automation Conf.*, June 1994, pp. 12–17.
- [10] F. Najm, "Transition density, a stochastic measure of activity in digital circuits," in *Proc. 28th Design Automation Conf.*, June 1991, pp. 644–649.
- [11] K. Roy and S. Prasad, "SYCLOP: Synthesis of CMOS logic for low power applications," in *Proc. Int. Conf. Computer Design: VLSI in Computers and Processors*, Oct. 1992, pp. 464–467.
- [12] A. Shen, A. Devadas, A. Ghosh, and K. Keutzer, "On average power dissipation and random pattern testability of combinational logic circuits," in *Proc. Int. Conf. Computer-Aided Design*, Nov. 1992, pp. 402–407.
- [13] C.-Y. Tsui, M. Pedram, and A. Despain, "Exact and approximate methods for switching activity estimation in sequential logic circuits," in *Proc. 31st Design Automation Conf.*, June 1994, pp. 18–23.



Mazhar Alidina received the B.S. degree from Lehigh University in 1992 and the S.M. degree from the Massachusetts Institute of Technology in 1994, both in electrical engineering.

He is currently a member of the Technical Staff with the Signal Processing and Integrated Circuit Design group of AT&T Bell Laboratories. His research interests are in low power design, CAD for low power, and VLSI design.



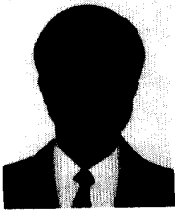
José Monteiro was born in Lisbon, Portugal. He received the Engineer's and Master's degrees in electrical and computer engineering in 1989 and 1992 respectively, from Instituto Superior Técnico at the Technical University of Lisbon.

He is currently working on the Ph.D. degree at the Massachusetts Institute of Technology in the area of power estimation and synthesis for low power of VLSI circuits.

Srinivas Devadas (M'88) received the B. Tech degree in electrical engineering from the Indian Institute of Technology, Madras in 1985 and the M.S. and Ph.D. degrees in electrical engineering from the University of California, Berkeley, in 1986 and 1988, respectively.

Since August 1988, he has been at the Massachusetts Institute of Technology, Cambridge, and is currently an Associate Professor of Electrical Engineering and Computer Science. His research interests span all aspects of synthesis of VLSI circuits, with emphasis on optimization techniques for synthesis at the logic, layout and architectural levels, testing of VLSI circuits, formal verification, hardware/software co-design, design-for-testability methods and interactions between synthesis and testability of VLSI systems.

Dr. Devadas held the Analog Devices Career Development Chair of Electrical Engineering from 1989 to 1991. He has received five Best Paper awards at CAD conferences and journals, including the 1990 IEEE TRANSACTIONS ON CIRCUITS AND DEVICES Best Paper award. In 1992, he received a NSF Young Investigator Award. He has served on the technical program committees of several conferences and workshops including the International Conference on Computer Design, and the International Conference on Computer-Aided Design. He is a member of the ACM.



Abhijit Ghosh (M'91) received the B.Tech degree in electrical and electronics engineering from the Indian Institute of Technology, Kharagpur, in 1986. He received the M.S. and Ph.D. degrees in electrical engineering and computer science from the University of California at Berkeley in 1988 and 1991, respectively.

Since 1991, he has been a Senior Engineer at Mitsubishi Electric Research Laboratories, Inc., Sunnyvale, CA, conducting research in CAD for VLSI and system design. His research interests include all aspects of CAD for VLSI with special emphasis on logic synthesis, formal verification, testing, parallel processing, fault-tolerant computing, compiler optimization, low power design and synthesis, and electronic system design automation.

Dr. Ghosh received the Best Paper award at the 27th IEEE Design Automation Conference, 1990. He is a member of ACM.



Marios Papaefthymiou received the B.S. degree in electrical engineering from the California Institute of Technology in 1988 and the S.M. and Ph.D. degrees in computer science from the Massachusetts Institute of Technology in 1990 and 1993, respectively.

Currently, he is an Assistant Professor of Electrical Engineering and Computer Science at Yale University, New Haven, CT. His research interests include algorithms, parallel computation, and VLSI design.