

A Combined Gate Replacement and Input Vector Control Approach for Leakage Current Reduction

Lin Yuan and Gang Qu

Abstract—Input vector control (IVC) is a popular technique for leakage power reduction. It utilizes the transistor stack effect in CMOS gates by applying a minimum leakage vector (MLV) to the primary inputs of combinational circuits during the standby mode. However, the IVC technique becomes less effective for circuits of large logic depth because the input vector at primary inputs has little impact on leakage of internal gates at high logic levels. In this paper, we propose a technique to overcome this limitation by replacing those internal gates in their worst leakage states by other library gates while maintaining the circuit's correct functionality during the active mode. This modification of the circuit does not require changes of the design flow, but it opens the door for further leakage reduction when the MLV is not effective. We then present a divide-and-conquer approach that integrates gate replacement, an optimal MLV searching algorithm for tree circuits, and a genetic algorithm to connect the tree circuits. Our experimental results on all the MCNC91 benchmark circuits reveal that 1) the gate replacement technique alone can achieve 10% leakage current reduction over the best known IVC methods with no delay penalty and little area increase; 2) the divide-and-conquer approach outperforms the best pure IVC method by 24% and the existing control point insertion method by 12%; and 3) compared with the leakage achieved by optimal MLV in small circuits, the gate replacement heuristic and the divide-and-conquer approach can reduce on average 13% and 17% leakage, respectively.

Index Terms—Gate replacement, leakage reduction, minimum leakage vector (MLV).

I. INTRODUCTION

AS THE VLSI technology and supply/threshold voltage continue scaling down, leakage power has become more and more significant in the power dissipation of today's CMOS circuits. For example, it is projected that subthreshold leakage power can contribute as much as 42% of the total power in the 90-nm process generation [11]. Many techniques thus have been proposed recently to reduce the leakage power consumption. Dual threshold voltage process uses devices with higher threshold voltage along noncritical paths to reduce leakage current while maintaining the performance [16]. Multiple-threshold CMOS (MTCMOS) technique places a high V_{th} device in series with low V_{th} circuitry, creating a sleep transistor [13]. In dynamic threshold MOS (DTMOS) [3], the gate and body are tied together and the threshold voltage is altered dynamically to suit the operating state of the circuit. Another technique to dynamically adjust threshold voltages is

(a)		(c)	
INPUT	Leakage(nA)	INPUT	Leakage (nA)
0	best:100.3	000	best: 22.84
1	worst: 227.2	001	37.84
		010	37.84
		011	2nd worst: 100.30
		100	37.01
		101	95.17
		110	94.87
		111	worst: 852.40

(b)	
INPUT	Leakage(nA)
00	best: 37.84
01	2nd worst: 100.30
10	95.17
11	worst: 454.50

Fig. 1. Leakage current of (a) INVERTER, (b) NAND2, and (c) NAND3. Data obtained by simulation in cadence spectre using 0.18- μ m process.

the variable threshold CMOS (VTCMOS) [14]. All of these approaches require the process technology support.

The input vector control (IVC) technique is applied to reduce leakage current at circuit level with little or no performance overhead [7]. It is based on the well-known transistor stack effect: a CMOS gate's subthreshold leakage current varies dramatically with the input vector applied to the gate [10]. Recently, Lee *et al.* observed that gate oxide leakage is also dependent on the input vectors to a CMOS gate [12]. Besides, the maximal and minimal leakage vectors are the same for both subthreshold leakage and gate leakage. In our study, we use Cadence Spectre to measure the overall leakage current in a CMOS gate that includes both subthreshold leakage and gate leakage. Fig. 1 lists the overall leakage current in INVERTER, NAND2 and NAND3 gates under all the possible input combinations. We see that the worst case leakage (marked in bold) is much higher than the other cases. The idea of IVC technique is to manipulate the input vector with the help of a sleep signal to reduce the leakage when the circuit is at the standby mode [9]. The associated minimum leakage vector (MLV) problem seeks to find a primary input vector that minimizes the total leakage current in a given circuit. [1], [4] [6], [8]–[10], [15]. The MLV problem is NP-complete and both exact and heuristic approaches have been proposed to search for the MLV. A detailed survey is given in Section II.

In this paper, we consider how to enhance IVC technique with little or no re-design effort. In particular, we study the **MLV+ problem** that seeks to modify a given circuit and determine an input vector such that the circuit's functionality is maintained at the active mode and the circuit leakage is minimized when the circuit is at standby mode. Our solution to this problem is based on the concept of gate replacement that is motivated by the large discrepancy between the worst leakage and the other cases (see Fig. 1). The essence of gate replacement is to replace a logic gate that is at its worst leakage state (WLS) by another library gate. This is illustrated by the following example.

Manuscript received May 28, 2005; revised October 15, 2005.

The authors are with the Electrical and Computer Engineering Department and Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742 USA (e-mail: yuanl@eng.umd.edu).

Digital Object Identifier 10.1109/TVLSI.2005.863747

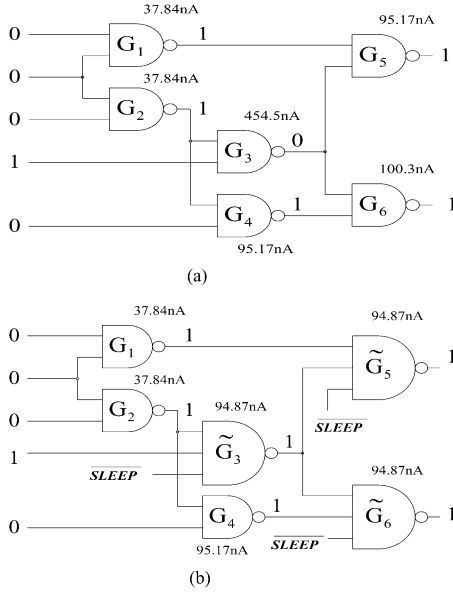


Fig. 2. Motivation example for gate replacement. (a) Original MCNC benchmark circuit C17 with total leakage 831.08 nA under the optimal MLV. (b) New circuit C17 with three gates replaced and total leakage 476.88 nA under the same MLV.

Consider circuit C17 from the MCNC91 benchmark suite [21] [Fig. 2(a)]. An exhaustive search finds the MLV $\{0, 0, 0, 1, 0\}$, with the corresponding minimum total leakage current of 831.08 nA. Note that gate G_3 has its worst leakage current (454.5 nA) with input $\{1, 1\}$, which contributes more than half of the total leakage. In fact, we have observed that a significant portion of the total leakage is often caused by the gates that are in their WLS (see Table II in Section V).

Instead of controlling the primary inputs, we consider replacing these leakage-intensive gates. In particular, we replace the NAND2 gate G_3 by a NAND3 \tilde{G}_3 , where the third input \overline{SLEEP} is the complement of the $SLEEP$ signal [Fig. 2(b)]. At active mode, $\overline{SLEEP} = 1$ and \tilde{G}_3 produces the same output as G_3 . But at the standby mode, $\overline{SLEEP} = 0$ and \tilde{G}_3 has a leakage of 94.87 nA [Fig. 1(b)], which is much smaller than G_3 's 454.5 nA.

However, this replacement also changes the output of this gate at the sleep mode and affects the leakage on gates G_5 and G_6 . In this case, we replace them in a similar fashion. As a result, the new circuit's total leakage becomes 476.88 nA, a 43% reduction from the original 831.08 nA in Fig. 2(a).

The proposed gate replacement technique is conceptually different from the existing IVC methods. In fact, they are complementary to each other. Specifically, IVC method considers the entire circuit and searches for an appropriate input vector in favor of small leakage. The gate replacement technique targets directly at the logic gates that are in their WLS under a specific input vector and replace them to reduce leakage. This paper has the following contributions.

- 1) We examine the effectiveness of IVC methods¹ in multilevel circuits. For all the 69 MCNC91 benchmarks, we

¹IVC-based approaches such as internal control point insertion [1] will be discussed in Section II.

obtain the optimal MLV for small circuits and the best over 10 000 random input vectors for large circuits. The number of gates in their WLS are on average 15% and 17%, respectively, but they contribute more than 40% of the circuit's total leakage.

- 2) Motivated by the above observation, we propose the technique to replace gates that are in their WLS by other library gates that will generate less leakage current at those states. Unlike other leakage reduction techniques such as MTCMOS and DTMOS, this modification of the circuit does not require changes of process technology in the design flow. Hence, it will not increase the design complexity or the leakage sensitivity.
- 3) We implement a fast gate replacement algorithm that gives an average of 10% leakage reduction for a fixed input vector. This algorithm's run time complexity is linear to the number of gates in the circuit in average cases and quadratic in the worst case.
- 4) We develop a divide-and-conquer approach to combine gate replacement and IVC. It reduces the leakage by 17% and 24% over the optimal/suboptimal MLV mentioned in 1) with little area and delay overhead. The number of gates in their WLS is dropped to 4% and 9%, respectively.

II. RELATED WORK

In this section, we mainly survey the efforts on IVC-based leakage reduction techniques. A survey on other leakage minimization techniques can be found in [7].

The effect of circuit input logic values on leakage current was observed by Halter and Najm [9]. The underlying reason of this effect was explained by Johnson *et al.* [10] as the transistor stack effect. Authors in [9] proposed a technique to insert a set of latches with MLV stored in to the primary inputs of a circuit, forcing the combinational logic into a low-leakage state when the circuit is idle. Many algorithms have been proposed to find such MLV. Based on the nature of these algorithms, they can be classified into the following groups:

Heuristic Algorithms: These include the random search algorithm developed by Halter and Najm [9] and the genetic algorithm proposed by Chen *et al.* [5].

Johnson *et al.* [10] defined *leakage observability* for each primary input as the degree to which the value of a particular input is observable in the magnitude of leakage current. They iteratively chose the input with the largest leakage observability and assigned it with a value that results in the smallest leakage. The input combination constructed in this greedy fashion was taken as the MLV.

In [15], Rao *et al.* introduced the concept of *node controllability*, which is defined as the minimum number of inputs that have to be assigned to particular values to ensure that a node (or gate) is in a specific state. Based on this, they proposed a fast greedy heuristic to determine the values of the primary inputs that minimize the node's leakage.

Exact Algorithms: The MLV problem can be modeled as a pseudo-Boolean satisfiability (SAT) problem. This formulation allows us to apply the off-the-shelf SAT solvers to find the MLV for leakage reduction [1], [2].

Gao and Hayes [8] formulated the MLV problem as an integer linear programming (ILP) problem. They first use pseudo-Boolean functions to represent leakage current in different types of cells with the general sum-of-products form. Then they apply the well-known Boole–Shannon expansion [19] to linearize the objective function and constraints. At last, they use an off-the-shelf ILP solver to solve the ILP optimization. For large circuits, the authors proposed a simplified mixed-ILP formulation that uses selective variable-type relaxation to reduce the runtime.

Based on the pseudo-Boolean formulation of the leakage in CMOS gates, two implicit pseudo-Boolean enumeration algorithms are presented in [6]. The input space enumeration method leverages integer valued decision diagrams and works well for small circuits. The hyper-graph partitioning based recursive algorithm represents a given circuit as a hyper-graph, partitions it, and uses divide-and-conquer to solve the problem. The trade-off between dynamic and leakage power in choosing the MLV has also been discussed.

Internal Point Control: Due to the ineffectiveness of IVC technique for circuits with large logic levels, Abdollahi *et al.* proposed a technique to directly control the value of internal pins to reduce leakage [1]. Their first approach inserts multiplexers at the input pins of each gate. The *SLEEP* signal selects the correct input in active mode and chooses the input values that produce low leakage current in standby mode. This approach can reduce leakage in the CMOS gates significantly; however, the inserted multiplexers will also generate leakage current and introduce extra delay and area. In their second approach, they modify the library gates by adding *SLEEP* signal-controlled transistors in the gate to select the low-leakage inputs for its fanout gates. However, since the structure of the gates is changed, a new set of library gates are needed.

Our gate replacement technique belongs to the class of internal point control, but is conceptually different from [1] in the following aspects.

- 1) They treat each input pin of the gates as potential places to insert multiplexers, while we consider only roots of each tree. The search space is reduced substantially.
- 2) Their purpose of modifying a gate G is to produce the low-leakage input for G 's fanout gate while we aim to reduce leakage current at G itself.
- 3) They modify gates whenever necessary while we restrict our algorithm to replace gates only by the available gates in the library, and, hence, do not require gate structure modification.

However, these two approaches can be combined as we will discuss in more details in Sections III and IV.

III. LEAKAGE REDUCTION BY GATE REPLACEMENT

A logic gate is at its WLS when its input yields the largest leakage current. Regardless of the primary input vector, a large number of gates are at WLS, particularly when the circuit has high logic depth. Take the 69 MCNC91 benchmarks for example. For each of the 69 circuits, when we apply the optimal (or suboptimal) MLVs to these circuits, 16% of the gates on average remain at WLS, producing more than 40% of the circuit's

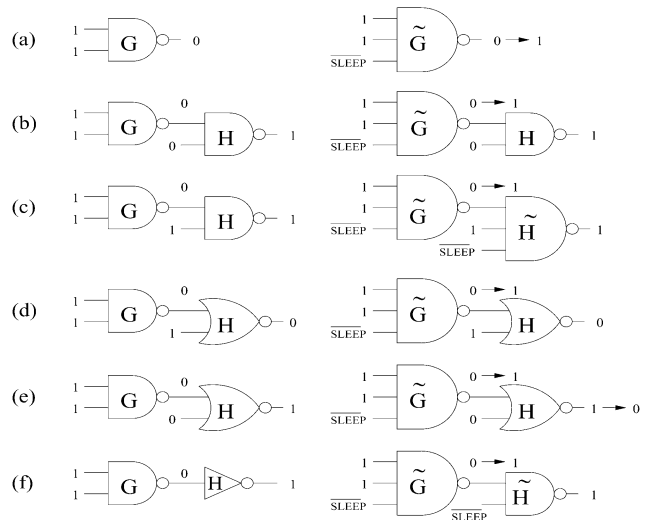


Fig. 3. Gate replacement and the consequence to its fanout gate.

total leakage. A detailed report can be found in Section V. In this section, we describe the gate replacement technique that targets directly the leakage reduction in WLS gates.

A. Basic Gate Replacement Technique

As we have shown in the motivation example in Section I, the proposed gate replacement technique replaces a gate $G(\vec{x})$ by another library gate $\tilde{G}(\vec{x}, SLEEP)$, where \vec{x} is the input vector at G , such that

- 1) $\tilde{G}(\vec{x}, 0) = G(\vec{x})$ when the circuit is active ($SLEEP = 0$);
- 2) $\tilde{G}(\vec{x}, 1)$ has smaller leakage than $G(\vec{x})$ when the circuit is in standby ($SLEEP = 1$).

The first condition guarantees the correct functionality of the circuit at active mode. The second condition reduces the leakage on gate G at the standby mode, but it may change the output of this gate. Note that, although we do not need to maintain the circuit's functionality at the standby mode, this change may affect the leakage of other gates and should be carefully considered.

Fig. 3(a) shows that the replacement of G by \tilde{G} changes the output from 0 to 1. For simplicity, we assume that G 's fanout only goes to gate H which can be either a NAND or a NOR or an INVERTER. In Fig. 3(b) and (d), we see that such change does not affect the output of gate H and, therefore, it will not affect any other gates in the circuit. Let $L(G(11))$ be the leakage of gate G with input 11, we can conveniently compute the leakage reduction by this replacement, which is $L(G(11)) + L(H(00)) - L(\tilde{G}(110)) - L(H(10))$ in the case of (b) for example.

In Fig. 3(c), the replacement at gate G not only changes the output of gate H , it also puts H at its WLS. Our solution is to replace the NAND2 gate H by a NAND3 \tilde{H} . This preserves the output of H and the leakage change will be $L(G(11)) + L(H(01)) - L(\tilde{G}(110)) - L(\tilde{H}(110))$. Similarly, in Fig. 3(f), we replace the INVERTER by a NAND2 gate. Finally, in Fig. 3(e), the replacement of G moves both gates G and H away from their WLS. It also changes the output of the NOR gate H , which we can conduct similar analysis.

Remarks:

- **General Fanout:** The above analysis is applicable to G 's fanout gate H of any type. The change of G 's output either does not affect H 's output [Fig. 3(b) and (d)] or changes H 's output. In the latter case, we either change H 's output back (Fig. 3(c) and (f)) or continue the analysis starting from H [Fig. 3(e)].
- **Beyond library gates:** If the library does not have a replacement for G , we can add one transistor into the N or P sections of G to meet conditions 1 and 2. This is similar to the gate modification method proposed in [1]. However, they attempt to control the output of the modified gate in order to reduce the leakage in its fanout gate by producing the desirable signal. Our gate replacement targets directly at the leakage reduction of the current gate.
- **Multiple fanouts:** When gate G has multiple fanouts, we analyze each of them and then consider their total leakage when we compute the leakage change due to the replacement of gate G .
- **Compatibility:** The gate replacement technique does not change the primary input vector of the circuit. This implies that we can combine it with existing MLV searching strategies to further reduce leakage. The MLV+ problem is based on this observation and is discussed in details in Section III-B.
- **Power overhead:** There is not much dynamic power overhead because the SLEEP signal remains constant at active mode and will not cause any additional switching activities. The leakage in gates \tilde{G} and G may be different at active mode. Such difference becomes negligible when the circuit stays at standby mode long enough [1].
- **Other overhead:** Gate replacement may introduce delay and area overhead. This overhead can be controlled by restricting the replacement off critical path and transistor re-sizing. Gate replacement does not add new logic gates and thus requires little or no effort to redo the place-and-route.

B. Fast Gate Replacement Algorithm

Based on the above gate replacement technique, we propose a fast algorithm that selectively replaces gates to reduce the circuit's total leakage for a given input vector. Fig. 4 gives the pseudocode of this algorithm.

We visit the gates in the circuit by the topological order. We skip all the gates that are not at WLS and the gates that have already been visited or marked (line 16) until we find a new gate G_i at WLS (line 2). Lines 3–9 find a subset of gates \mathcal{S} and temporarily replace them. \mathcal{S} includes all the unmarked gates whose leakage and/or output is affected by the replacement we attempt to do on gate G_i and other gates in \mathcal{S} . We then compute the total leakage change caused by the replacement of gates in \mathcal{S} (line 10) and adopt these replacements if there is a leakage reduction (lines 11–13). Otherwise, we simply mark gate G_i as visited and do not make any replacement (line 14). We then look for the next unmarked gate at WLS and this procedure stops when all the gates in the circuits are marked.

Correctness: The topological order guarantees that when we find a gate at its WLS, all its predecessors have already been

Input: $\{G_1, G_2, \dots\}$: gates in a circuit sorted topologically,
 $\{x_1, x_2, \dots\}$: an input vector,
 SLEEP: the sleep signal.

Output: a circuit of the same functionality when SLEEP = 0 and with less leakage when SLEEP = 1.

Gate Replacement Algorithm:

1. **for** each gate $G_i \in \{G_1, G_2, \dots\}$
2. **if** (G_i is at WLS and not marked)
3. include G_i in the selection \mathcal{S} ;
4. **while** (there is new addition to \mathcal{S})
5. **for** each newly selected gate G in \mathcal{S}
6. **if** (there exists library gate \tilde{G} meets the conditions in Section III-A)
7. temporarily replace G by \tilde{G} ;
8. **if** (output of G is changed due to this replacement)
9. include G 's unmarked fanout gate G_j in \mathcal{S} ;
10. compute the total leakage change of gates in \mathcal{S} ;
11. **if** (there is leakage reduction)
12. mark all gates G_j in the selection \mathcal{S} ;
13. make the replacements in lines 7,9,or 10 permanent;
14. **else** mark gate G_i only;
15. empty the selection \mathcal{S} ;
16. **else** mark G_i if it has not been marked yet;

Fig. 4. Pseudocode of the gate replacement algorithm.

considered. The replacement at line 7 ensures that the functionality will not change at the active mode. The subset \mathcal{S} constructed in the **while** loop (lines 4–9) is the *transitive closure* of gates that are affected by the replacement action at gate G_i . Therefore, we only need to compute the leakage change on gates within \mathcal{S} (line 10). We make the replacement only when this leakage change is in favor of us, so the new circuit will have less leakage in standby mode.

Complexity: Let n be the number of gates in the circuit. The for loop is linear to n . Inside the for loop, the computation of leakage change and the marking of all gates in \mathcal{S} (line 10–15) is linear to $|\mathcal{S}|$, the number of gates in \mathcal{S} . The **while** loop (lines 3–9) stops when there is no new addition to \mathcal{S} and this will be executed no more than $|\mathcal{S}|$ times. As we have discussed in Section III-A (see Fig. 3), in most cases, \mathcal{S} includes only G and its fanout gates. However, it may include all the gates of the circuit in cases similar to Fig. 3(e) and so $|\mathcal{S}|$ cannot be bounded by any constant. That is, $|\mathcal{S}|$ is $O(n)$ in the worst case and $O(k)$ on average, where k is the maximal fanout of the gates in the circuit. Consequently, the complexity of this gate replacement algorithm is $O(n^2)$ in the worst case and $O(kn)$ on average.

Improvement: There are several ways to improve the leakage reduction performance of the above gate replacement heuristic. The tradeoff will be either increased design complexity, or reduced circuit performance, or both. First, one can consider gates that are not in the library as we have commented in the second remark in Section III-A (line 6). However, this requires the measurement of leakage current, area and delay in these new gates as they are not available in the library. A second alternative is to insert control point at one of G 's fanins. For example, one can find the fanin y such that replacing y by its complement y' gives G the largest leakage reduction. If $y = 0$, replace it by $\text{OR}(y, \text{SLEEP})$; if $y = 1$, replace it by $\text{AND}(y, \text{SLEEP})$. However, the addition of new gates may require the repeat of placement and routing and will incur more area and delay penalty in general. Third, one may also consider both the library gate replacement and control point insertion at the same time and choose the one that gives more

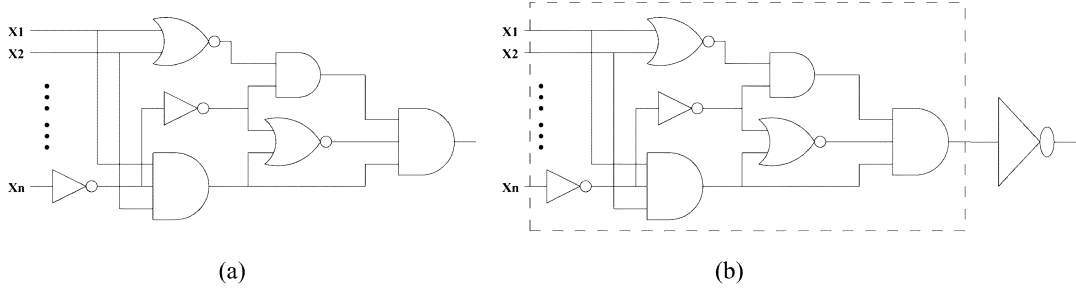


Fig. 5. Illustration for the proof of the NP-completeness of the MLV problem. (a) Circuit for satisfiability test. (b) Reducing the satisfiability test to MLV.

leakage reduction. Finally, whenever we replace gate G_i , we also make the replacement for all the other gates in the selection \mathcal{S} permanent (line 13). We have tested a couple of alternatives and they give limited improvement in leakage reduction at very high cost of run time complexity.

The incentive to keep the run time complexity of this gate replacement algorithm low is that it will be combined with IVC technique under the following divide-and-conquer approach to solve the MLV+ problem.

IV. SOLVING THE MLV+ PROBLEM

Recall that the MLV problem seeks the input vector that minimizes the circuit's total leakage. It has been claimed that this problem is NP-complete for general circuits [1], [6], [10], [15]. But no formal proof has been given to our knowledge. In this section, we first give a brief proof of the NP-completeness of the MLV problem and then define the MLV+ problem, an extension of the MLV problem. Our main focus will be on the divide-and-conquer approach that solves the MLV+ problem.

A. NP-Completeness of the MLV Problem

The MLV problem can be defined as follows: given a combinational circuit consisting of primary inputs (PIs), primary outputs (POs), internal logic gates connected by nets/wires, and the leakage current of each gate under different input combinations, determine an input vector at the PIs such that the total leakage current of all the gates in the circuit is minimized.

Theorem: The MLV problem is NP-complete.

Proof: On one side, we have already mentioned a couple of exact algorithms that solve the MLV problem by reducing it to NP-complete problems such as pseudo-Boolean satisfiability and ILP.

On the other side, we show that the NP-complete CIRCUIT-SAT problem [18] can be reduced to the MLV problem. Consider an arbitrary circuit shown in Fig. 5(a), to test whether the circuit is satisfiable (i.e., producing a logic "1" at its output), we construct a new circuit by adding a big inverter at its output [Fig. 5(b)]. The inverter is big in the sense that it has a huge leakage value L when its input is "0" and a small leakage ϵ when its input is "1." Actually, we can set L to be the sum of ϵ and the leakage of each gate in the circuit when it is in its WLS. Now we solve the MLV problem for this modified circuit. If the total leakage is less than L , clearly the original circuit is satisfiable and the MLV is one input vector that makes the circuit output logic "1." Otherwise, because that the only

way for the total leakage to be larger than L is when the input to the big inverter is "0," the original circuit is not satisfiable. ■

B. MLV+ Problem and Outline of the Divide-and-Conquer Approach

In the previous section, we have seen that leakage current can be further reduced over the MLV by the proposed gate replacement technique. We have also mentioned that this technique is independent of the input vector and can be combined with the MLV method. We, hence, formulate the following **MLV+ problem**.

Given a combinational circuit with PIs, POs, the internal logic gates that implement the PI-PO functionality, and the leakage current of each library gate under its different input patterns, determine a gate level implementation of the same PI-PO functionality without changing the place-and-route and an input vector at the PIs that minimizes the total leakage.

Apparently, this is an extension of the MLV problem with the relaxation of modifying circuit by gate replacement. It enlarges the search space of MLV and provides us with the opportunity of finding better solutions. For a circuit of k PIs and n internal logic gates, the search space for the original MLV problem is the 2^k different input combinations. Under the above MLV+ formulation, the search space becomes $2^k \cdot \prod_{i=1}^n l_i$, where l_i is the number of library gates that can replace gate i , including gate i itself. Assuming that half of the gates have one replacement, then the solution space for MLV+ problem will be $2^{n/2}$ times larger than the solution space for the MLV problem. Even when we restrict the gate replacement technique only to gates that are at their WLS, this will be significant because 1) a circuit normally has more gates than PIs ($n \gg k$) and 2) the percentage of gates in WLS is considerably high (16% on the MCNC91 benchmark when MLV is applied, and will be higher as the logic depth of the circuit increases).

As we have analyzed in the previous section, the MLV+ problem not only enlarges the solution space for the IVC method, it also has the great potential in improving the solution quality (in terms of leakage reduction) because of the stack effect. However, one challenge is how to explore such enormous solution space for better solutions. Given the NP-completeness of the MLV problem, we consider special circuits where the MLV+ can be solved optimally and develop heuristics for the general case. In the rest of this section, we describe details of our proposed divide-and-conquer approach that consists of the following phases:

- 1) decompose a general circuit into tree circuits.
- 2) find the MLV for each tree circuit optimally by dynamic programming.
- 3) apply the gate replacement technique to the MLV for each tree to further reduce leakage.
- 4) connect the tree circuits by a genetic algorithm.

C. Finding the Optimal MLV for Tree Circuits

A tree circuit is a single output circuit in which each gate, except the primary output, feeds exactly one other gate. A general combinational circuit can be trivially decomposed into nonoverlapping tree circuits [19]. (This is illustrated in Fig. 7.) The circuit in (a) is not a tree because gate G_3 has two fan-out gates G_5 and G_6 . By splitting at the fanout of G_3 , we get three trees with G_3 , G_5 and G_6 being the root of each tree, respectively.

We consider a tree circuit with gates $\{G_1, G_2, \dots, G_n\}$ sorted in the topological order, which is preserved by the tree decomposition.

Let $L(G_i(\vec{x}))$ be the leakage current in the gate G_i when vector \vec{x} is applied at G_i 's fanins. Each gate G_i can be treated as the root of a subtree circuit. Let $LK(i, z)$ be the minimum total leakage of the tree circuit when it outputs logic value z at root G_i and $\vec{V}(i, z)$ be the input vector to the tree circuit that achieves $LK(i, z)$. We develop a dynamic programming approach to compute the pairs $(LK(i, 0), \vec{V}(i, 0))$ and $(LK(i, 1), \vec{V}(i, 1))$ for each gate G_i . The MLV for the tree circuit rooted at gate G_n , with gates $\{G_1, G_2, \dots, G_n\}$ sorted in the topological order, can then be determined conveniently.

- 1) For each input signal to the tree, define

$$LK(0, z) = 0, \quad \vec{V}(0, z) = z. \quad (1)$$

- 2) For each gate $G_i (i = 1, 2, \dots, n)$, let

$$LK(i, z) = \min_{\forall \vec{x}, s.t. G_i \text{ outputs } z} \left(L(G_i(\vec{x})) + \sum_{j=1}^t LK(i_j, x_{i_j}) \right) \quad (2)$$

$$\vec{V}(i, z) = \cup_{j=1}^t \vec{V}(i_j, x_{i_j}^0) \quad (3)$$

where $\{x_{i_1}, x_{i_2}, \dots, x_{i_t}\}$ are the fanins of G_i from gates $\{G_{i_1}, G_{i_2}, \dots, G_{i_t}\}$, respectively, and the input combination $\{x_{i_1}^0, \dots, x_{i_t}^0\}$ achieves $LK(i, z)$.

- 3) The minimum leakage of the tree circuit with gates $\{G_1, \dots, G_n\}$ is given by

$$\min\{LK(n, 0), LK(n, 1)\} \quad (4)$$

and the MLV will be either $\vec{V}(n, 0)$ or $\vec{V}(n, 1)$ accordingly.

A step-by-step illustration of the dynamic programming can be found in [17].

Correctness: We show the correctness of the recursive formula in (2) and (3). To compute $LK(i, z)$, we need to consider all the possible combination of fanins $\{x_{i_1}, \dots, x_{i_t}\}$ that

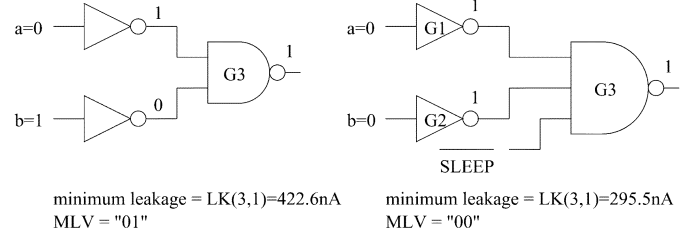


Fig. 6. MLV in a circuit before and after gate replacement.

produces output z at gate G_i . For each such combination, the minimum leakage in the subtree rooted at G_i is the sum of leakage at gate G_i and the minimum leakage at each of its fan-in gate G_{i_j} with output x_{i_j} , $LK(i_j, x_{i_j})$. Equation (2) takes the overall minimum leakage and gives the correct $LK(i, z)$. Assume that this minimum leakage is achieved when G_i has fanins $x_{i_1} = x_{i_1}^0, \dots, x_{i_t} = x_{i_t}^0$. Note that $\vec{V}(i_j, x_{i_j}^0)$ is the input vector for the subtree circuit rooted at G_j to produce $x_{i_j}^0$ with the minimum leakage $LK(i_j, x_{i_j})$. The tree structure of the circuit guarantees that the subtrees rooted at $\{G_{i_1}, \dots, G_{i_t}\}$ will not share any common inputs. Therefore, $\vec{V}(i, z)$ is the simple concatenation of $\vec{V}(i_j, x_{i_j}^0)$ as given in (3).

Complexity: Equations (1) and (4) take constant time. For each gate G_i , we need to compute $(LK(i, 0), \vec{V}(i, 0))$ and $(LK(i, 1), \vec{V}(i, 1))$ by (2) and (3). This requires the enumeration of all the 2^t different combinations of G_i 's t fanins. For the first time, we need to perform t additions in (2). If we enumerate the rest $2^t - 1$ cases following a Gray code, we only need to update $L(G_i(\vec{x}))$ (two operations), replace one $LK(i_j, x_{i_j})$ (two operations) and compare the result with the current minimum leakage, a total of five operations. Therefore, we need $t + 5 \cdot (2^t - 1)$ operations for each G_i and this gives a complexity of $O(K \cdot n)$, where K is a constant depending on the largest number of fanins in the circuit.

After obtaining the MLV for the tree circuit, we perform the gate replacement algorithm proposed in Section III to further reduce leakage. Note that, although the MLV is optimal, this does not guarantee us an optimal solution for the MLV+ problem on the tree circuit. For example, consider the circuit in Fig. 6, the algorithm finds the optimal MLV $\{a = 0, b = 1\}$ with leakage 422 nA. Gate 2 is at its WLS and the gate replacement algorithm does not give any improvement. The input vector $\{0, 0\}$ gives the maximum leakage 654 nA; however, when we apply gate replacement technique and replace G_3 , the leakage is reduced to 295 nA. In fact, $\{0, 0\}$ is the optimal solution for the MLV+ problem.²

D. Connecting the Tree Circuits

In the previous phase, we have determined the output and required input for each individual tree circuit to yield the minimum leakage. The goal of this phase is to combine all the tree circuits to solve the MLV+ problem for the original circuit. The root of each tree circuit may have multiple fanouts that go to

²We conjecture that the MLV+ problem remains NP-hard for tree circuit. Because we have already lost the optimality when we do the tree decomposition, we will not discuss in details on how to find better solutions to MLV+ on tree circuits. For the same reason, we did not focus on how to improve the fast gate replacement algorithm in Section III-B.

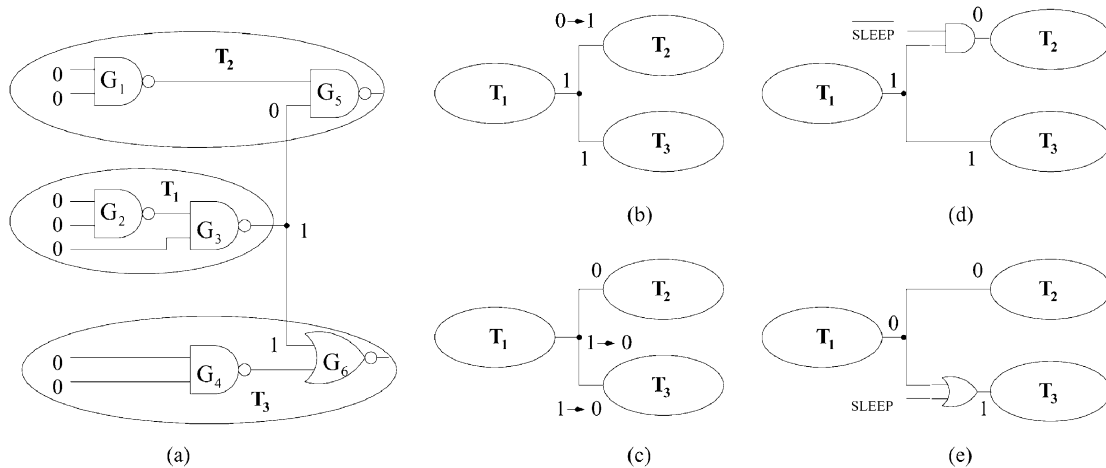


Fig. 7. Resolving the conflict in connecting tree circuits.

other tree circuits as input. Since we treat the tree circuits independently, conflict occurs if the output of a tree circuit and the value required by its fanout gates are not consistent. For example, in Fig. 7(a), the circuit is decomposed into three tree circuits T_1 , T_2 and T_3 . T_1 outputs “1” when its MLV is applied, while T_2 and T_3 require “0” and “1” from T_1 in their respective MLVs. So we have a conflict.

There are basically three ways to resolve this conflict:

- (I) enforcing T_1 's output at all the fanout gates [Fig. 7(b)];
- (II) changing T_1 's output and enforcing this new value at all the fanout gates [Fig. 7(c)];
- (III) inserting an AND gate to allow them to be inconsistent [Fig. 7(d)]. Similarly, if T_1 output “0” and some of its fanouts require “1,” we can add an OR gate [as shown in Fig. 7(e)].

To decide which one we should use to resolve the conflict, we apply each of them and re-evaluate the circuit's total leakage. In (I), this requires the recomputing of the minimum leakage and the MLV for tree circuit T_2 under the condition that its input from T_1 is logic “1.” The dynamic programming algorithm in Section IV-B can be trivially modified for this purpose. In (II), we need to do the same procedure for tree circuit T_3 . Besides, we have to replace the pair $\{LK(n, 1), \vec{V}(n, 1)\}$ for tree circuit T_1 by $\{LK(n, 0), \vec{V}(n, 0)\}$.

Both (I) and (II) resolve the conflict by sacrificing the minimum leakage of tree circuits under the provably optimal MLV. In (III), we successfully connect the tree circuits while preserving the minimum leakage and MLV for each tree with the help of the *SLEEP* signal-controlled AND or OR gates. The cost is that we have to add the leakage of the inserted AND or OR gate into the total leakage. We mention that this gate addition also preserves the correctness of the circuit at active mode when *SLEEP* = 0.

It is now easy to make a decision on which method to adopt to resolve a single conflict: use the one that gives the minimum leakage. However, the decision at one conflict may affect the existence of conflict at other places in the circuit. For example, method (I) in Fig. 7(b) could change the output of tree T_2 and directly affect whether there is a conflict at the root of T_2 .

We use a genetic algorithm (GA) to resolve the conflicts and connect all the tree circuits. A solution by the GA is in the form of a binary bit stream, each bit indicates whether there is a conflict at the root of a tree and which method to use to resolve it. In particular, a “1” means there is a conflict and method (III) should be used; a “0” means that there is either no conflict or we should use the better one of methods (I) and (II) to resolve the conflict.

The GA follows a standard routine where we start with a population of N random bit streams (referred to as *chromosomes*). Based on each bit stream, we resolve the conflict, apply the dynamic programming algorithm in Section IV-B to re-compute the minimum leakage of a tree circuit when methods (I) and (II) are used, run the gate replacement algorithm in Fig. 4 on the entire circuit, and compute the circuit's total leakage. The *fitness* for a bit stream is calculated from the leakage value. The smaller the leakage, the larger the *fitness*. We sort all the chromosomes according to their fitness and create the next generation by the *roulette wheel* method. In this method, the probability that a *chromosome* is selected as one of the two parents is proportional to its fitness. *Crossover*, which refers to the exchange of substrings in two chromosomes, is performed among parents to produce children. A simple *mutation* operation, which flips a bit in the chromosome at the *bit mutation rate*, is also used. The GA continues to generate a total of N new chromosomes and starts for the next generation. This process repeats for certain number of times (50 in our simulation) and the best chromosome is returned as the optimal solution.

E. Overhead Analysis

As the control gates are introduced in the tree-connecting stage of the algorithm, they also require sleep signal to control. Hence, we need to consider the extra power these control gates and sleep signal may consume, and their effect on the overall power saving. In this subsection, we will discuss the power overheads.

1) *Control Gates*: The control gates will consume extra dynamic power and leakage power. In this paper, we only consider the leakage power overhead of the inserted gates and ignore their dynamic power due to the following reasons. First, the number

of inserted control gates only accounts for 5% to 6% of the total number of gates in the circuit. Second, they are simple 2-input AND and OR gates, which have a relatively small intrinsic capacitance at the node compared to other gates. Third, the switching activities in these control gates are very limited because one of the two inputs is the sleep signal, which changes only at the moment when the circuit switches between active mode and sleep mode. As dynamic power is dependent on physical capacitance and switching activities, we consider this dynamic power overhead is negligible.

As for leakage power, we measured the average leakage current in control gates over all possible inputs. In our algorithm, we add this extra leakage current to the objective function, i.e., the overall leakage current to be minimized. Therefore, the leakage saving achieved in our algorithm has already considered this overhead.

2) *Sleep Signal*: Both the gate replacement and the control gates require the sleep signal to drive them during active and sleep mode. The generation of the sleep signal may consume extra power. However, due to the fact that our experiment was conducted at the logic synthesis level before placement and routing, it is not practical to obtain such power data quantitatively. On the other hand, the sleep signal is required by many other leakage minimization techniques, such as [1], [3], and [13]. Hence, in this paper, we expect the generation of the sleep signal to be similar to those approaches and we believe this problem can be better solved at the physical level of circuit design.

V. EXPERIMENTAL RESULTS

We implemented the gate replacement and divide-and-conquer techniques in SIS environment [20] and applied them on 69 MCNC91 benchmark circuits. Each circuit is synthesized and mapped to a 0.18- μm technology library. We use Cadence Spectre to simulate the leakage current for all the library gates under every possible input vector. The supply voltage and threshold voltage are 1.5 and 0.2 V, respectively. The measured leakage current includes both subthreshold and gate leakage. The simulations are conducted on an Ultra SPARC SUN workstation.

Our results are compared with traditional IVC methods in terms of leakage saving, run time, area and delay penalty. The 69 benchmarks include 26 small circuits with 22 or fewer primary inputs (Table I) and 43 large circuits (Table II). For each small circuit, we find the optimal MLV by exhaustive search. For each large circuit, we choose the best MLV from 10 000 distinct random input vectors. It is reported that this will give us a 99% confidence that the vectors with less leakage is less than 0.5% of the entire vector population [9], [15]. To have a fair comparison with [1], we also collect the average leakage of 1000 random input vectors for each large circuit.

Table I reports the results for the 26 small circuits. Column 4 lists the leakage current for each circuit when the best MLV is applied. Even in this case, an average of 15% of the gates are at WLS as shown in column 5. The fast gate replacement algorithm is able to move about half of these gates from their WLS (column 7). This results in a 13% leakage reduction with

TABLE I
RESULTS ON 26 SMALL CIRCUITS WITH 22 OR LESS PRIMARY INPUTS

circuit	pi #	gate #	exhaustive			gate replace			divide-and-conquer			
			leak(nA)	wls	imprv	wls	ar inc	imprv	wls	# tr	cg	ar inc
b1	3	13	2195	23%	2%	15%	5%	2%	10%	5	0%	5%
cm42a	4	25	2941	0%	0%	0%	0%	8%	0%	18	4%	8%
C17	5	6	831	17%	43%	0%	17%	43%	0%	4	0%	17%
cm82a	5	28	5017	21%	29%	4%	12%	40%	1%	10	4%	18%
decod	5	22	1921	0%	0%	0%	0%	8%	0%	21	5%	3%
cm138a	6	19	1760	0%	0%	0%	0%	1%	0%	12	5%	5%
z4ml	7	66	12246	24%	25%	11%	11%	37%	4%	20	5%	17%
f51m	8	136	26038	26%	37%	7%	12%	48%	4%	25	3%	14%
9symml	9	166	34018	26%	20%	17%	5%	38%	8%	18	8%	14%
alu2	10	356	64153	21%	2%	20%	0%	21%	5%	89	7%	11%
x2	10	44	6159	9%	15%	2%	3%	12%	2%	18	9%	10%
cm85a	11	38	4925	8%	14%	3%	3%	13%	3%	16	0%	3%
cm151a	12	34	5745	24%	9%	18%	4%	3%	18%	5	3%	5%
alu4	14	728	133127	25%	1%	21%	1%	15%	4%	166	7%	10%
cm162a	14	45	6947	18%	2%	9%	3%	0%	9%	13	4%	12%
cu	14	49	6182	12%	16%	6%	2%	9%	5%	21	6%	7%
cm163a	16	43	6376	19%	2%	9%	3%	1%	9%	11	5%	13%
cmb	16	42	5405	10%	11%	5%	2%	4%	4%	8	2%	6%
parity	16	75	12764	20%	11%	15%	5%	15%	7%	15	7%	20%
pm1	16	39	3474	3%	0%	0%	1%	-2%	0%	16	3%	3%
t481	16	1945	251184	2%	1%	1%	0%	26%	0%	17	2%	1%
tcon	17	41	6491	20%	43%	0%	14%	41%	0%	9	2%	17%
pcl	19	74	12594	20%	32%	4%	6%	32%	4%	22	0%	6%
sct	19	92	11811	18%	14%	9%	4%	10%	6%	24	4%	6%
cc	21	48	5823	13%	6%	10%	1%	6%	9%	22	0%	1%
cm150a	21	72	12270	15%	4%	14%	1%	1%	10%	9	7%	10%
Average				15%	13%	8%	4%	17%	4%		4%	9%

only 4% area increase (columns 6 and 8). We mention that we restrict ourselves to replace only gates off critical paths. This leaves 8% of the gates in the circuits at their WLS, but it also guarantees us that there is no delay overhead.

The last four columns show that the divide-and-conquer algorithm gives a 17% leakage reduction over the best MLV at the cost of 9% more area. We incorporate delay constraints in the genetic algorithm to ensure that the delay overhead to be within 5%. The two columns in the middle are the number of tree circuits in each case and the number of control gates we have used to connect these trees. Only in three cases, we have inserted more than five control gates. Note that the addition of control gates may decrease the delay because it reduces the fanouts of the gate. The area increase comes from the addition of control gates and the replacement of “smaller” gates by “bigger” library gates.

Fig. 8 reports the leakage and wls gates reduction in the 43 large circuits (x -axis) with 22 PIs or more. We replace the infeasible exhaustive search by the best solution from a random search of 10 K input vectors. The fast gate replacement algorithm are restricted only on gates off critical paths; for the divide-and-conquer approach, we set the maximal delay increase to be 5%.

The benchmarks are sorted by the total leakage achieved by the divide-and-conquer method normalized to the best over 10 K random search, which is shown one of the two curves at the top part of the figure. The average leakage reductions are 10% by gate replacement only (leakage G.R.) and 24% by divide-and-conquer method (leakage D.C.). The maximal leakage reductions are 46.4% and 60%, respectively. The three curves at the bottom give the ratio of WLS gates. On average, the 10 K random search has 17% gates at WLS(orig, wls); the proposed fast gate replacement and divide-and-conquer techniques reduce this ratio to 11%(G.R. wls) and 9%(D.C. wls), respectively.

TABLE II
RESULTS ON 43 LARGE CIRCUITS WITH PRIMARY INPUTS MORE THAN 22

circuit	pi #	gate #	random search (10k)			gate replacement (G.R.)				divide-and-conquer (D.C.)							over 1K average	
			leak(nA)	time(s)	wls(%)	imprv(%)	time(s)	wls(%)	area(%)	imprv(%)	time(s)	wls(%)	#tree	#gates/tree	cg(%)	area(%)	G.R.(%)	D.C.(%)
cordic	23	102	18434.0	9.9	21.6	15.1	0.01	11.8	5.7	27.4	10.1	7.8	52	3.1	7	9.3	28.4	38.8
ttt2	24	207	33801.5	22.7	18.4	9.5	0.02	17.4	4.4	18.4	72.6	14.5	43	4.7	6	9.6	30.9	37.7
il	25	39	5250.6	5.4	7.7	27.7	0	0.0	4.3	26.3	6.0	0.0	16	2.1	3	5.1	45.5	44.4
pcler8	27	90	14670.1	10.0	16.7	11.1	0.01	11.1	4.0	27.0	14.9	10.3	31	2.4	0	4.0	35.2	46.8
c8	28	164	26083.0	17.4	19.5	19.0	0.01	4.3	8.4	14.4	21.5	0.0	38	5.9	8	6.9	31.7	27.8
C6288	32	2400	480084.2	222.0	29.0	2.9	0.11	27.7	1.9	8.8	398.7	11.7	1424	1.7	29	27.3	7.0	12.6
comp	32	163	28322.3	15.2	22.1	5.6	0.01	11.7	2.4	13.2	85.4	9.7	77	3.4	2	5.4	34.1	39.4
C1908	33	615	117029.6	57.2	20.5	2.5	0.02	17.1	0.9	31.0	66.0	13.4	218	2.9	10	10.1	6.4	33.7
my_adder	33	225	40842.1	21.0	22.2	2.0	0.02	20.0	1.5	31.1	32.1	18.2	95	2.8	7	6.4	8.9	36.0
term1	34	363	60460.5	37.3	18.5	11.7	0.02	9.6	4.0	15.4	160.0	8.8	75	6.8	5	8.8	23.9	27.0
count	35	144	22445.4	15.2	17.4	0.0	0.01	17.4	0.0	3.4	14.2	16.7	37	4.2	2	2.4	0.0	15.4
C432	36	200	38101.4	20.1	15.0	11.2	0.01	9.0	3.3	37.5	24.7	8.0	79	4.1	6	8.9	21.6	44.8
unreg	36	113	18188.4	12.7	19.5	4.6	0.01	5.3	3.1	17.3	84.4	5.3	18	6.3	2	4.9	20.1	30.7
too_large	38	582	107888.1	61.4	17.4	12.5	0.05	9.6	2.2	37.1	80.1	9.6	113	5.2	7	10.9	24.5	45.7
b9	41	111	16100.3	12.8	11.7	8.6	0.01	8.1	2.0	19.7	68.0	7.9	34	3.3	4	8.7	30.1	38.5
C1355	41	517	91739.0	50.7	22.1	4.5	0.02	13.0	1.4	19.1	95.0	6.9	265	2.0	15	13.1	12.1	25.4
C499	41	532	95292.0	48.3	20.3	5.0	0.05	13.3	2.2	18.2	84.5	8.8	197	2.7	7	5.8	16.8	28.4
cht	47	232	38560.8	25.3	16.8	4.5	0.02	11.6	3.7	14.7	22.8	10.1	66	3.5	2	3.3	18.4	27.1
apex7	49	239	41955.1	26.0	20.1	19.3	0.02	8.4	5.8	30.3	25.6	7.4	82	2.9	3	11.1	26.9	36.9
C3540	50	1136	218977.1	115.0	18.2	2.9	0.08	15.3	1.3	21.3	133.8	7.7	381	3.0	15	2.1	11.5	28.2
x1	51	295	45351.2	32.8	16.3	17.7	0.02	4.7	4.8	25.0	105.9	4.0	61	4.8	7	11.9	32.1	38.2
C880	60	354	61978.8	35.8	18.9	12.6	0.04	11.6	4.1	25.8	39.9	11.6	115	3.1	13	10.8	21.7	33.5
dalu	75	1865	349299.8	187.5	25.6	3.8	0.15	23.2	1.4	23.2	194.9	17.9	321	5.8	8	14.2	29.1	43.5
example2	85	286	51036.6	32.6	17.5	4.3	0.02	15.0	1.4	41.5	28.9	13.2	110	2.6	2	9.8	11.3	45.7
i9	88	510	88469.6	63.9	1.0	0.0	0.04	1.0	0.0	17.3	156.0	1.0	113	4.5	3	2.1	0.0	50.1
x4	94	378	61336.3	46.4	18.3	28.2	0.03	4.5	5.3	33.6	206.5	4.5	110	3.4	11	8.6	40.1	44.7
i3	132	92	16166.9	14.9	21.7	0.0	0.00	21.7	0.0	18.5	0.0	20.7	6	15.3	0	0.0	0.0	27.2
i5	133	269	44848.1	34.3	12.6	19.9	0.02	4.8	2.9	42.0	45.6	4.0	68	4.0	2	7.8	35.8	53.5
i8	133	1898	305924.5	224.4	14.2	9.1	0.15	11.4	0.8	39.4	7591.3	4.0	259	7.3	6	6.3	43.5	62.3
apex6	135	710	126523.6	86.1	20.8	3.9	0.06	5.9	2.1	26.8	399.5	3.0	215	3.3	10	5.7	11.4	32.6
rot	135	601	109944.1	67.1	20.0	17.5	0.06	13.8	5.5	23.1	403.3	12.0	208	2.9	10	12.7	23.5	28.7
x3	135	742	116641.0	89.5	14.3	15.6	0.07	9.0	3.2	20.4	384.4	5.6	192	3.9	8	10.0	29.7	33.7
i6	138	340	47021.1	47.3	9.1	46.4	0.03	0.6	2.1	59.0	89.8	0.0	71	4.8	1	3.0	68.9	76.2
frg2	143	1030	165090.4	136.0	16.1	12.9	0.11	7.4	3.2	28.4	176.5	6.8	244	4.2	5	7.4	28.0	40.8
pair	173	1538	270729.8	160.9	18.9	7.6	0.14	13.2	2.4	17.5	366.0	5.4	434	3.5	12	12.0	14.9	24.0
C5315	178	1777	343295.9	188.3	18.7	6.0	0.15	15.0	2.0	11.5	534.5	9.9	532	3.3	12	15.1	11.6	16.8
i4	192	136	22699.8	22.8	8.8	3.1	0.01	8.8	0.4	27.8	34.6	8.8	6	22.7	0	4.6	28.3	46.6
i7	199	405	58431.5	58.4	6.2	1.2	0.04	5.7	0.2	13.5	117.9	5.7	76	5.3	2	1.1	37.7	45.5
i2	201	109	13174.8	22.1	4.6	19.7	0.01	0.0	2.2	36.8	36.1	0.0	12	9.1	4	3.6	36.1	49.7
C7552	207	2801	515320.2	293.3	20.8	0.6	0.18	15.3	0.3	5.9	726.0	6.9	908	3.1	15	16.1	20.6	24.8
C2670	233	807	155992.3	94.5	18.1	0.8	0.09	17.8	0.2	11.9	98.6	14.6	235	3.4	11	9.9	5.4	16.0
des	256	3995	931447.4	471.2	23.6	7.2	0.24	18.5	2.5	45.7	8502.6	7.3	847	4.7	11	14.2	17.6	51.8
ii0	257	2281	440552.2	261.6	20.4	6.7	0.2	19.2	1.9	14.3	162.8	4.5	695	3.3	14	6.1	11.7	18.8
Average				80.9	17%	10%	0.05	11%	2%	24%	510.2	9%			6%	7%	23%	37%

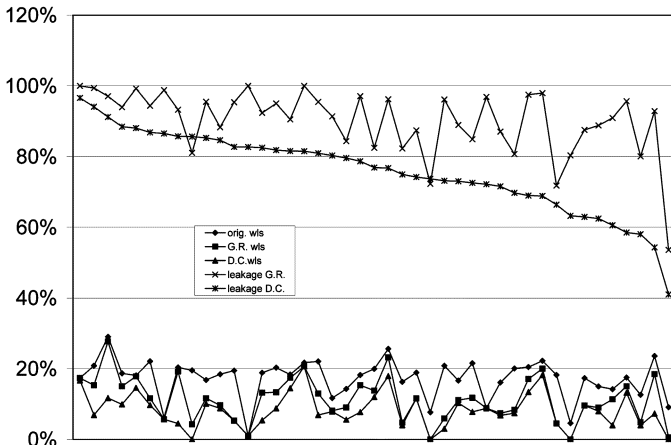


Fig. 8. Leakage and WLS percentage on 43 large circuits with 22 PIs or more. X-axis lists benchmarks sorted by leakage current in divide-and-conquer approach; Y-axis shows percentage of leakage and WLS gates.

More detailed results for these 43 circuits are shown in Table II. Columns 4–6 list the leakage current, runtime, and

percentage of gates at WLS when the best MLV from 10000 random vectors is applied to each circuit. The next four columns show the results when the fast gate replacement algorithm is applied to such best MLV. The average run time is only 0.05 s and increases linearly to the number of gates in the circuit. There is no delay overhead and the area increase is only 2%.

The next seven columns show results by the divide-and-conquer approach where we set a 5% maximum delay constraint. In the genetic algorithm, we start with a population size of $N = 150$ and it converges after 50 generations. We are able to achieve, over the best MLV from 10000 random vectors, 24% leakage saving with 7% area penalty on average. Although the average run time is $6\times$ of the random search, we mention that this is mainly caused by the two circuits, *i8* and *des*. They have a couple of large tree circuits and, therefore, the frequently called dynamic programming takes considerably long time. Excluding these two circuits, the average run time for random search and the divide-and-conquer algorithm drop to 64.7s and 143s, respectively. More importantly, we see clearly the run time for random search increases exponentially to the number of primary input and linearly to the number of gates (columns

TABLE III
AVERAGE PERFORMANCE COMPARISON WITH ALGORITHM

	algorithm in [1]	gate replacement	divide-and-conquer
leakage reduction	25%	23%	37%
delay penalty	$\leq 5\%$	0%	$\leq 5\%$
area penalty	$\leq 15\%$	2%	7%

2,3,5). However, the run time for the divide-and-conquer approach grows at a much slower pace (column 12).

Finally, the last two columns compare our results with those reported in [1]. Because their detailed results are not available, we can only compare the average performance. In their experimental setup, the leakage reduction is compared with the average value among 1000 random vectors. For a fair comparison, we also report in the last two columns the improvement of our approaches over the same baseline. Table III summarizes the performance improvement in the control point insertion approach [1], our gate replacement algorithm, and the divide-and-conquer approach.

VI. CONCLUSION

We study the MLV+ problem which seeks to modify a given circuit and determine an input vector such that the correct functionality is maintained when the circuit is active and the leakage is minimized under the determined input vector when the circuit is at stand-by mode. The relaxation of circuit modification with changing its functionality enlarges the solution space of the IVC method. We show that MLV (and, hence, MLV+) problem is a hard problem and propose low-complexity heuristics to solve the MLV+ problem. The proposed algorithms are practical and effective in the sense that we do not need to change the design flow and re-do place-and-route. The experimental results show that this technique improves significantly the performance of IVC in leakage reduction at gate level with little area and delay overhead.

ACKNOWLEDGMENT

The authors would like to thank the Editor-in-Chief, the Associate Editor, and the reviewers for their valuable comments. A full version of this paper can be found in [17].

REFERENCES

- [1] A. Abdollahi, F. Fallah, and M. Pedram, "Leakage current reduction in CMOS VLSI circuits by input vector control," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 12, no. 2, pp. 140–154, Feb. 2004.
- [2] F. Aloul, S. Hassoun, K. Sakallah, and D. Blaauw, "Robust SAT-based search algorithm for leakage power reduction," in *Proc. Int. Workshop Integr. Circuit Des.*, 2002, pp. 167–177.
- [3] F. Assaderaghi, D. Sinitzky, S. A. Parke, J. Bokor, P. K. Ko, and C. Hu, "Dynamic threshold-voltage MOSFET (DTMOS) for ultra-low voltage VLSI," *IEEE Trans. Electron Devices*, vol. 44, no. 3, pp. 414–422, Mar. 1997.
- [4] S. Bobba and I. N. Hajj, "Maximum leakage power estimation for CMOS circuits," in *Proc. IEEE Alessandro Volta Memorial Workshop Low Power Des.*, 1999, pp. 116–116.
- [5] Z. Chen, M. Johnson, L. Wei, and K. Roy, "Estimation of standby leakage power in CMOS circuits considering accurate modeling of transistor stacks," in *Proc. ISLPED*, 1998, pp. 239–244.
- [6] K. Chopra and S. B. K. Vrudhula, "Implicit pseudo-Boolean enumeration algorithms for input vector control," in *Proc. DAC*, 2004, pp. 767–772.
- [7] D. Duarte, Y. Tsai, N. Vijaykrishnan, and M. Irwin, "Evaluating run-time techniques for leakage power reduction," in *Proc. VLSI Des.*, 2002, pp. 31–38.
- [8] F. Gao and J. P. Hayes, "Exact and heuristic approaches to input vector control for leakage power reduction," in *Proc. ICCAD*, 2004, pp. 527–532.
- [9] J. Halter and F. Najm, "A gate-level leakage power reduction method for ultra low power CMOS circuits," in *Proc. CICC*, 1997, pp. 475–478.
- [10] M. C. Johnson, D. Somasekhar, and K. Roy, "Models and algorithms for bounds on leakage in CMOS circuits," *IEEE Trans. Comput.-Aided Des. Integr. Circuits Syst.*, vol. 18, no. 6, pp. 714–725, Jun. 1999.
- [11] J. Kao, S. Narendra, and A. Chandrakasan, "Subthreshold leakage modeling and reduction techniques," in *Proc. ICCAD*, 2002, pp. 141–148.
- [12] D. Lee, W. Kwong, D. Blaauw, and D. Sylvester, "Analysis and minimization techniques for total leakage considering gate oxide leakage," in *Proc. DAC*, 2003, pp. 175–180.
- [13] V. Khandelwal and A. Srinivastava, "Leakage control through fine-grained placement and sizing of sleep transistors," in *Proc. IEEE/ACM Int. Conf. Comput.-Aided Des.*, Nov. 2004, pp. 533–536.
- [14] T. Kuroda *et al.*, "A 0.9 V 150 MHz 10 mW 4 mm² 2-D discrete cosine transform core processor with variable threshold-voltage (VT) scheme," *IEEE J. Solid-State Circuits*, vol. 31, no. 11, pp. 1770–1779, Nov. 1996.
- [15] R. M. Rao, F. Liu, J. L. Burns, and R. B. Brown, "A heuristic to determine low leakage sleep state vectors for CMOS combinational circuits," in *Proc. ICCAD*, 2003, pp. 689–692.
- [16] V. Khandelwal, A. Davoodi, and A. Srinivastava, "Simultaneous V_t selection and assignment for leakage optimization," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 13, no. 6, pp. 762–765, Jun. 2005.
- [17] L. Yuan and G. Qu, "A combined gate replacement and input vector control approaches for leakage current reduction," Inst. Adv. Comput. Studies (UMIACS), Univ. Maryland, College Park, MD, Tech. Rep. TR 2005-63, 2005.
- [18] M. R. Garey and D. S. Johnson, *Computers and Intractability, a Guide to the Theory of NP-Completeness*. San Francisco, CA: Freeman, 2001.
- [19] G. D. Hachtel and F. Somenzi, *Logic Synthesis and Verification Algorithms*. Norwell, MA: Kluwer, 1996.
- [20] E. Sentovich *et al.*, "SIS: A system for sequential circuit synthesis," Univ. California, Electron. Res. Lab. Memorandum, Berkeley, CA, no. UCB/ERL M92/41, 1992.
- [21] *Synthesis and Optimization Benchmarks User Guide*, Microelectronic Center, Triangle Park, NC, 1991.



Lin Yuan (M'03) received the B.S. degree in information engineering from Xi'an Jiaotong University, China, in 2001. He is currently working toward the Ph.D. degree in the Department of Electrical and Computer Engineering, University of Maryland, College Park.

His research interests include low-power design, VLSI design automation, and wireless sensor networks.



Gang Qu (S'98–A'00–M'03) received the B.S. and M.S. degrees in mathematics from the University of Science and Technology, China, in 1992 and 1994, respectively, and the Ph.D. degree in computer science from the University of California, Los Angeles, in 2000.

In 2000, he joined the Electrical and Computer Engineering Department, University of Maryland, College Park. In 2001, he became a member of the University of Maryland Institute of Advanced Computer Studies. His research interests include

low-power system design, computer-aided synthesis, sensor network, and intellectual property reuse and protection.

Dr. Qu has received many awards for his academic achievements and service and has served on the Technical Program Committee for many conferences. Currently, he is the General Co-Chair of the 16th ACM Great Lakes Symposium on VLSI.