

Fault-Tolerant Resynthesis with Dual-Output LUTs

Ju-Yueh Lee¹, Yu Hu², Rupak Majumdar³, Lei He¹ and Minming Li⁴

1. Electrical Engineering Department, University of California, Los Angeles
2. Electrical and Computer Engineering Department, University of Alberta, Edmonton Canada
3. Computer Science Department, University of California, Los Angeles
4. Computer Science Department, City University of Hong Kong

ABSTRACT

We present a fault-tolerant post-mapping resynthesis for FPGA-based designs that exploits the dual-output feature of modern FPGA architectures to improve the reliability of a mapped circuit against faults. Emerging FPGA architectures, such as 6-LUTs in Xilinx Virtex-5 and 8-input ALMs in Altera Stratix-III, have a secondary LUT output that allows access to non-occupied SRAM bits. We show that this architectural feature can be used to build redundancy for fault masking with limited area and performance overhead. Our algorithm improves reliability of a mapping by performing two basic operations: *duplication* (in which free configuration bits are used to duplicate a logic function whose value is obtained at the secondary output) and *encoding* (in which two copies of the same logic function are ANDed or ORed together in the fanout of the duplicated logic). The problem of fault tolerant post-mapping resynthesis is then formulated as the optimal duplication and encoding scheme that ensures the minimal circuit fault rate w.r.t. a stochastic single fault model. We present an ILP formulation of this problem and an efficient algorithm based on generalized network flow. On MCNC benchmarks, experimental results show that for combinational circuits the proposed approach improves mean-time-to-failure(MTTF) by 27% with 4% area overhead, and the proposed approach with explicit area redundancy improves MTTF by 113% with 36% area overhead, compared to the baseline mapping by ABC. This provides a viable fault tolerance solution for non-mission critical applications compared to TMR (triple modular redundancy) which has a $5 \times 6 \times$ area overhead.

1. INTRODUCTION

Faults are becoming increasingly pronounced in emerging applications and technologies, from permanent faults arising from circuit processing at nanometer scales to transient soft errors arising from high-energy particle hits. This has spurred much research on ways to improve fault tolerance without incurring substantial area, power, or performance penalties [1].

In this paper, we explore a new opportunity for achieving fault tolerance with respect to a stochastic fault model by exploiting architectural features of emerging FPGA architectures. For example, in Xilinx's Virtex-5 FPGA [2], the 6-input LUT has two usable output pins (see Figure 1). It can implement two independent LUTs if the total number of unique pins in each does not exceed five. Similarly, in Altera's Stratix II FPGA [3], an ALM can implement two independent LUTs if the total number of input pins does not exceed 8 and constraints on input sharing and LUT sizes are met. One use of dual-output LUTs is to map circuits that contain many small logic functions, for example, with 2 or 3 inputs, to produce higher logic density and reduced circuit power [4]. Unfortunately, it is not easy to pack logic functions with 4 or 5 inputs in one dual-output LUT. In fact, empirical observation of a wide collection of benchmark sets shows that less than 50% of LUTs in Virtex-5 FPGAs are fully occupied, even with sophisticated "dual-output-aware" technology mapping and an LUT merging procedure [5].

The low logic utilization rate in real designs motivates us to utilize non-occupied SRAM bits of dual-output LUTs for fault masking. Specifically, redundant logics can be implemented in non-occupied SRAM bits of these LUTs and can be accessed via the

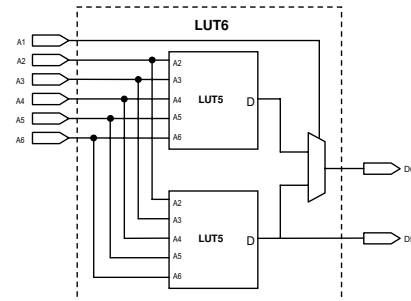


Figure 1: Xilinx Virtex-5 dual-output LUT

secondary LUT outputs. Logic duplication can effectively mask $0 \rightarrow 1$ (resp. $1 \rightarrow 0$) single event upsets (SEUs) by ANDing (resp. ORing) the original and the copy of a duplicated LUT. The additional AND and OR logic can be encoded into LUTs at the next logic stage. We consider two versions of our algorithm. In the *fully masked* version (FMD) we assume the duplicated LUT is encoded in the same way (AND or OR) in all fanouts. In the *partially masked* version (PMD) a duplicated LUT may be encoded differently in different fanouts.

We make the following contributions in this paper. Assuming a stochastic single fault model[6], we formulate the fault masking problems for dual-output LUT-based circuits using duplication and (full or partial) encoding as an ILP optimization problem. We show how the structure of the ILP problem allows an efficient generalized network flow-based algorithm. We test our algorithm on MCNC benchmarks under two different CAD flows with and without explicit area overhead, respectively. Experimental results show that for combinational circuits the proposed algorithm improves mean time to failure (MTTF) by 27% with 4% area overhead. The proposed approach with explicit area overhead improves MTTF by 113% with 36% area overhead, when compared to the technology mapping obtained by [5].

In contrast to related fault-tolerance techniques proposed for FPGAs, such as triple-modular redundancy (TMR) which has a $5 \times$ to $6 \times$ area/power overhead [7], our approach is lightweight and incurs minor area and performance overhead. In contrast to chip-wise synthesis [8, 9], our generic solution has much lower test complexity. In contrast to a recent stochastic synthesis algorithm (ROSE) [10], our algorithm is much faster on wide input LUTs (by exploiting architectural features) and it obtains more MTTF improvement with slightly higher overhead. Clearly, this paper provides a new and viable solution for the spectrum of techniques for reliability and overhead trade off.

The rest of this paper is organized as follows. Section 2 presents the problem formulation. Section 3 proposes our primary algorithm for dual-output LUT-based fault masking, and this algorithm is generalized in Section 4. The experimental results are given in Section 5, and the paper is concluded by Section 6. To the best of our knowledge, this paper is the first systematic study on stochastic fault tolerance using the dual-output architectural feature of modern FPGAs.

2. PROBLEM FORMULATION

2.1 Fault Modeling

In this paper, we assume a *stochastic single fault* model for faults¹, i.e., there is at most one fault occurring at a time for the entire circuit with the identical fault rate for each SRAM bit.

To quantitatively measure the sensitivity of an SRAM bit with respect to an SEU, we first define the *criticality* of SRAM bit b as

$$c_b = \frac{1}{2^n} |\{x \mid C(x) \neq C(\bar{b})(x)\}|, \quad (1)$$

where b can be an SRAM bit in an LUT or a routing switch². C is the circuit and n is the number of primary inputs of C . x is the primary input vector $x \in \{0, 1\}^n$. $C(x)$ is the value of the primary outputs when the input vector x is applied to C . $C(\bar{b})(x)$ is the circuit which is identical to C except that the SRAM bit b is flipped. For simplicity, a routing switch is a programmable connection in connection boxes or switch boxes. In general, the criticality of an SRAM bit is the percentage of the primary input vectors which cause observable erroneous circuit outputs due to the flip of SRAM bit b . Accordingly, the criticality c_L of an LUT L is the average criticality of its every SRAM bit:

$$c_L = \frac{1}{2^K} \sum_i c_{L_i}, \quad (2)$$

where L_i the i th SRAM bit of LUT L .

The *full-chip fault rate* of a circuit is the percentage of the primary input vectors which cause observable erroneous outputs due to faults. Based on the single fault assumption, the full-chip fault rate of a circuit C is calculated by the average criticality of all SRAM bits in C , and it is defined by the following equation.

$$\text{fault_rate}(C) = \frac{\sum_{L \in \text{Luts}(C)} 2^K \cdot c_L + \sum_{b \in \text{Rout}(C)} c_b}{2^K \cdot |\text{Luts}(C)| + |\text{Rout}(C)|} \cdot P_F,$$

where $\text{Luts}(C)$ is the set of LUTs in C and $\text{Rout}(C)$ is the set of routing SRAM bits in C . P_F is the probability that an SEU occurs in an occupied SRAM bit. Therefore we have

$$P_F = \frac{2^K \cdot |\text{Luts}(C)| + |\text{Rout}(C)|}{\mathcal{A}} \cdot \gamma,$$

where \mathcal{A} is the total number of SRAM bits of the underlying FPGA device and γ is a constant denoting the single SRAM bit sensitivity to an SEU strike. Thus, the full chip fault rate can be expressed as

$$\text{fault_rate}(C) = \left(\sum_{L \in \text{Luts}(C)} 2^K \cdot c_L + \sum_{b \in \text{Rout}(C)} c_b \right) \cdot \frac{\gamma}{\mathcal{A}}. \quad (3)$$

The criticality of each SRAM bit in a mapped circuit can be obtained by random simulation or other more efficient approaches, such as the analytical models based on signal masking probability [11] or the hardware-based emulator, which is used in this paper. Specifically for sequential circuits, multiple clock cycles are simulated in order to capture the propagation of faults in registers. The computation of criticality needs only to be performed once. During the course of the optimization, criticality values can be updated efficiently, as will be described later in the paper.

2.2 Design Freedom

In our optimization, we shall use two atomic operations: *duplication* and *encoding*. Duplication computes the same function twice in an LUT, such that the two outputs produce two independently copied versions of the function. Encoding takes the two outputs (representing two independently computed versions of the same logic function) from a dual-output LUT and computes their AND or OR in the fanout nodes. One of the two different encoding schemes, i.e., AND-encoding and OR-encoding, is applied to the fanout LUTs of a duplicated LUT based on the logic masking effectiveness. As

¹ As shown in [6], simultaneous multiple SRAM flips almost never occur, and therefore the single fault is a valid assumption in practice.

² An SRAM bit flips in a routing switch may cause a short/open/bridging fault. In this paper, we model these faults as the logic value flipping. In addition, we do not consider multiple errors caused by the single SEU-induced bridging fault.

shown in Figure 2, any 0→1 fault that occurs in LUT A can be masked. In this way, any 0→1 fault that occurs in intermediate wire z_{org} or z_{cpy} can also be masked. It is important to note that the proposed duplication/encoding scheme can be performed spanning through the register boundary. Consider Figure 2 again, and image z is now a register driven by LUT A . The duplication of both LUT A and register z prevent the propagation of 0→1 fault in the next clock cycle, and therefore the fault will not be observable before it reaches primary outputs. A nice feature of this approach is that it preserves the logic depth as well as the number of LUTs, which controls the area and performance overhead due to the fault masking insertion.

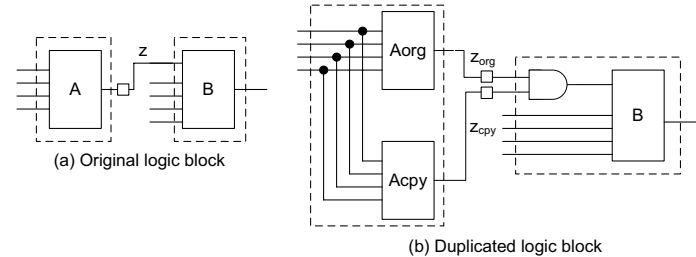


Figure 2: LUT duplication (LUT A) and AND-encoding of its fanout (LUT B)

Given an LUT L with duplication, if the same encoding (AND or OR) is applied to encode all its fanout LUTs, LUT L is called *fully masked*. Otherwise, LUT L is called *partially masked*. In this paper, the duplication scheme in which all duplicated LUTs are fully masked is called *fully masking-based duplication* (FMD); otherwise it is called *partially masking-based duplication* (PMD). FMD is a special case of PMD, and therefore FMD can be solved more efficiently but will result in a lower quality of fault tolerance, as will be shown later. In general, the duplication-based fault tolerance problem can be formulated as follows.

FORMULATION 1. *Given a circuit, perform duplication and encoding on LUTs (if possible), such that the full-chip fault rate (defined in (3)) is minimized.*

3. FMD ALGORITHM

3.1 Criticality Update in FMD

After the FMD of an LUT, the following lemmas show that the criticalities of this LUT, its fanout SRAM bits in the routing switches, and its fanout LUTs can be updated efficiently. In addition, the update of the criticality of one duplication is independent of that of the other duplication. Therefore the FMD of the full-chip can be performed optimally via a mathematical formulation.

LEMMA 1. *The criticality of LUT L after a duplication is $2 \cdot (\sum_{L_i \in \text{OnSet}} c_{L_i})/2^K$ (resp. $2 \cdot (\sum_{L_i \in \text{OffSet}} c_{L_i})/2^K$) if all fanouts of LUT L are AND-encoded (resp. OR-encoded) assuming the single fault model.*

PROOF. Suppose LUT L is AND (resp. OR)-encoded, all 0 → 1 (resp. 1 → 0) flips due to SEU can be masked by its fanout AND (resp. OR) gate, and therefore the criticality of $L_i \in \text{OffSet}$ (resp. OnSet) is $c_{L_i} = 0$, where $L_i \in \text{OffSet}$ (resp. OnSet) means $L_i = 0$ (resp. $L_i = 1$). \square

LEMMA 2. *The encoding operation preserves the criticality of LUT L assuming the single fault model.*

PROOF. Suppose the Boolean function implemented by LUT L is $f(i_1, \dots, i_m)$, where m is the number of occupied inputs of LUT L before encoding, and the criticality of configuration bit L_i is c_{L_i} . Without loss of the generality, the Boolean function after encoding can be expressed as $f(i_1, \dots, i_m \odot i_{m+1})$, where i_{m+1} is a newly

occupied input pin due to the addition of \odot logic (AND or OR). The criticality of the encoded LUT, L^e , is

$$\frac{1}{2K} \cdot \left(\sum_{L_i^e \in \{i_m \neq i_{m+1}\}} c_{L_i^e} + \sum_{L_i^e \in \{i_m \equiv i_{m+1}\}} c_{L_i^e} \right),$$

where $\{i_m \neq i_{m+1}\}$ (resp. $\{i_m \equiv i_{m+1}\}$) is the set of min-terms where LUT pins i_m and i_{m+1} have the different (resp. same) logic values. Assuming the single fault, any fault occurring in LUT L indicates that no faults occur in its fanin LUTs and therefore we always have $i_m \equiv i_{m+1}$. As a result, we have $c_{L_i} = 0$ for all $L_i \in \{i_m \neq i_{m+1}\}$, and therefore the criticality of the encoded LUT L is

$$\frac{1}{2K} \cdot \sum_{L_i^e \in \{i_m \equiv i_{m+1}\}} c_{L_i^e} = \frac{1}{2K} \cdot \sum_{\forall L_i} c_{L_i} = c_L.$$

This completes the proof. \square

The criticality of an SRAM bit in routing switches can be divided into two components: $c_b^{1 \rightarrow 0}$ and $c_b^{0 \rightarrow 1}$, i.e., the percentage of input vectors that make the $0 \rightarrow 1$ and $1 \rightarrow 0$ fault observable at the primary outputs, respectively. In addition, $c_b^{1 \rightarrow 0}$ and $c_b^{0 \rightarrow 1}$ are bounded by their driver LUT, i.e., $c_b^{1 \rightarrow 0}$ (resp. $c_b^{0 \rightarrow 1}$) equals to the total OnSet (resp. OffSet) criticality of its driver LUT.

Based on the above lemmas, after the duplication of LUT L , the total criticality reduction under the AND-encoding scheme is

$$\Delta c_{\text{AND}} = \sum_{L_i \in \text{OffSet}} c_{L_i} - \sum_{L_i \in \text{OnSet}} c_{L_i} + \delta_b \cdot (c_b^{0 \rightarrow 1} - c_b^{1 \rightarrow 0}),$$

and the total criticality reduction under OR-encoding scheme is

$$\Delta c_{\text{OR}} = \sum_{L_i \in \text{OnSet}} c_{L_i} - \sum_{L_i \in \text{OffSet}} c_{L_i} + \delta_b \cdot (c_b^{1 \rightarrow 0} - c_b^{0 \rightarrow 1}),$$

where δ_b is the number of SRAM bits in the corresponding routing switches.

3.2 Algorithm and Complexity

The possibility of the LUT duplication and encoding is determined by the available LUT resource. An LUT cannot be duplicated more than once since it can only use up to two outputs. On the other hand, an LUT can be encoded more than once. In addition, duplication and encoding can be performed simultaneously on a LUT. The possibility of duplication and encoding is constrained by the number of spare (non-occupied) input pins of an LUT. Particularly, we have the following lemma to quantitatively express such constraints.

LEMMA 3. *If a K -LUT has p occupied input pins and $p < K$, the total number of atomic operations (duplication and encoding) that can be applied to this LUT cannot exceed $(K - p)$, and the duplication can only be done at most once.*

PROOF. We show that each atomic operation occupies one spare input pin. It is clear that each encoding operation occupies exactly one original spare input pin. A duplication operation effectively occupies one input pin. In other words, A_1 must be set as constant 0 in a Xilinx Virtex-5 LUT, as shown in Figure 1. \square

Based on Lemma 3, the FMD problem can be formulated as the following ILP problem:

$$\begin{aligned} & \text{Maximize} && \sum_{L \in \text{Luts}(C)} w_L \cdot d_L \\ & \text{Subject to} && d_L + \sum_{f \in \text{fanin}(L)} d_f \leq S_L, \forall L \in \text{Luts}(C), \\ & && d_L \in \{0, 1\}, \end{aligned} \quad (4)$$

where for each LUT $L \in \text{Luts}(C)$, we introduce a 0-1 variable d_L where $d_L = 1$ if LUT L is duplicated, and $d_L = 0$ otherwise. S_L is the number of spare input pins of LUT L , and $w_L = \max(\Delta c_{\text{AND}}, \Delta c_{\text{OR}})$. Formulation (4) guarantees that the encoding scheme which leads to the maximal criticality reduction between AND-encoding and OR-encoding is selected.

THEOREM 1. *The decision version of problem (4) is NP-hard.*

PROOF. We prove the NP-hardness of Problem (4) by a reduction from 3-SAT. Given a 3-SAT instance with variables x_1, \dots, x_n and clauses C_1, \dots, C_m , we construct Problem (4) as follows. We construct DAG $G = (V, E)$ where $V = \{S, T\} \cup \{x_1, \dots, x_n, \bar{x}_1, \dots, \bar{x}_n\} \cup \{C_1, \dots, C_m\}$. For the edge set E , we have $(S, x_i) \in E, (T, x_i) \in E, (\bar{x}_i, x_i) \in E$ for $1 \leq i \leq n$ and $(x_i, C_j) \in E$ if C_j contains literal \bar{x}_i ($(\bar{x}_i, C_j) \in E$ if C_j contains literal x_i). The weight for each node is defined as follows. $w_S = w_T = (m + 1)(2n + 2)$, $w_{x_i} = w_{\bar{x}_i} = m + 1$ and $w_{C_j} = 1$. We finally set $K = 6$. Then the constraints in the ILP can be expressed as follows.

$$\begin{aligned} d_T + d_S + d_{x_i} + d_{\bar{x}_i} &\leq 3, 1 \leq i \leq n \\ d_{\bar{z}_1} + d_{z_2} + d_{z_3} + d_{C_j} &\leq 3, 1 \leq j \leq m, z_1, z_2, z_3 \text{ are literals in } C_j. \end{aligned}$$

Based on the weight selection, we can see that the optimal solution of the constructed Problem (4) will duplicate both S and T ; furthermore, it will duplicate exactly one out of each pair x_i and \bar{x}_i ; finally it will try to duplicate as many C_j 's as possible. By the second set of ILP constraints, a C_j can be duplicated if and only if one of its literals is duplicated. Hence, the given 3-SAT instance is satisfiable if and only if the constructed Problem (4) can achieve objective value $2(m + 1)(2n + 2) + (m + 1)n + m$, and a literal is set true in the truth assignment if and only if its corresponding LUT is duplicated. \square

Although Theorem 1 proves that Problem (4) is NP-Hard under arbitrary weights w_i , the complexity of FMD remains unknown as it is not clear if there exists a configuration which results in weight assignments for w_i used in the proof.

Problem (4) exhibits a nice structure of the min-cost discrete generalized network flow (GNF) as shown in Figure 4(a). Although the worst case complexity of a discrete GNF problem remains NP-Hard [12], there exist efficient combinatorial algorithms to solve it in polynomial behavior [13].

3.3 Example

Let us illustrate the proposed FMD algorithm using the circuit shown in Figure 3 assuming $K = 4$. The FMD problem can be transformed to a min-cost discrete GNF shown in Figure 4(a) as follows. Each level-1 node d_L in the flow network represents the duplicability of LUT L , and each level-2 node S_L expresses the spare pin number of LUT L . Note that only d_A, d_B, d_D are shown in level-1 because LUT C and LUT E are in the last logic level, and therefore cannot be duplicated. The tuple in each arc denotes (capacity, cost). The capacity of a (S_L, t) -arc is S_L , which serves as the constraints in the above ILP formulation. The weight of a (s, d_L) -arc is $-w_L$ and the supply flow in source node s is ∞ . The value r associated with each node is the gain factor, which means that flow f becomes $r \cdot f$ when it exits from this node. In FMD the gain factor for node d_L is equal to the number of fanins of LUT L . It is easy to verify that this GNF is equivalent to the integer programming-based Formulation (4). The solution of the FMD problem can be obtained by the flow amount in (s, d_L) -arcs after solving this min-cost GNF problem with integer flow values, i.e., discrete GNF. One feasible solution of this GNF problem is shown by the bold arcs in Figure 4(a) and the corresponding FMD result is shown in Figure 4(b). For instance, LUT D is encoded due to flow in (d_A, S_D) -arc, and it is also duplicated due to flows in (s, d_D) -arc and (d_D, S_D) -arc.

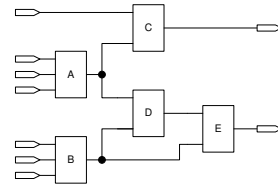
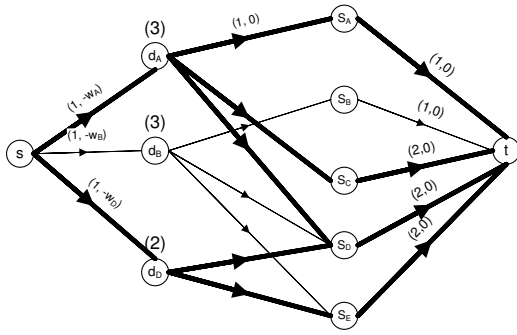
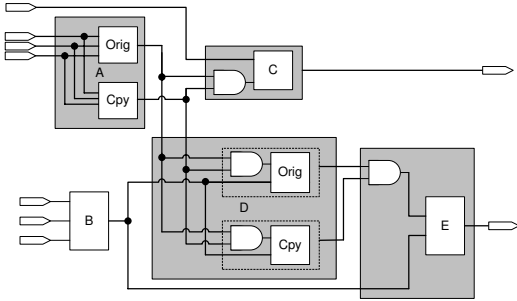


Figure 3: A sub-circuit example



(a) Min-cost generalized flow network for FMD



(b) A FMD solution corresponding to the flow (bold arcs) in Figure 4(a) for the circuit in Figure 3

Figure 4: FMD example

4. PMD ALGORITHM

As a generalization of FMD, PMD allows partial encoding of the fanout LUTs of a duplicated LUT, and therefore it is more flexible. It can be shown that Lemma 2 remains true under PMD but Lemma 1 no longer holds under PMD. Essentially the PMD of LUT L does not mask faults observable along paths starting from its non-encoded fanouts, and therefore the update of criticality cannot be simply obtained based on Lemma 1.

Given an LUT L with non-saturated³ fanouts f_1, \dots, f_p , the duplication of LUT L requires the encoding of at least one of its non-saturated fanouts. In general, each fanout can be AND-encoded, OR-encoded or non-encoded for LUT L , and therefore there are $(3^p - 1)$ fanout encoding options, which are denoted as $\Phi(L) = \{\phi(L, 1), \dots, \phi(L, 3^p - 1)\}$. For each option $\phi(L, j)$, we can pre-compute the modified criticality of LUT L after applying such an encoding, and denote the reduction of criticality as w_L^j . Since AND/OR-encoding only changes the criticalities of 0(1) configuration bits of an LUT, we only need to compute the modified criticality of LUT L by applying pure AND/OR-encoding, and the criticality caused by the mixed encoding can be calculated by merging the corresponding pure encodings as shown in Figure 5. As a result, we only need to perform 2^{p+1} computations for the criticality instead of 3^p .

Note that w_L^j can be computed by performing $O(N \cdot \max(|\Phi(L)|))$ iterations of full-chip random simulations, where N is the number of LUTs and $\max(|\Phi(L)|)$ is bounded by the maximal fanout number, which is a small constant in practice. In addition, only those options that have positive w_L^j are kept since our objective is to maximize the criticality reduction. As shown in Figure 6, most LUTs have less than 9 fanouts and a threshold 512 is used to control the size of $\Phi(L)$ taken into account. We use an emulator-based approach to perform logic simulation and therefore can pre-compute w_L^j very efficiently.

Similar to FMD, the update of the criticality for one LUT after PMD is independent of that for another LUT under the single fault assumption, and therefore PMD can be solved optimally by

³ A fanout is saturated if the corresponding fanout LUT has no spare pins for duplication or encoding. Otherwise, it is non-saturated.

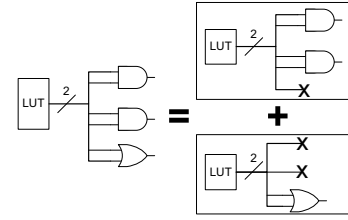


Figure 5: The modified criticality after two AND-encoding and one OR-encoding can be calculated by merging the modified bit criticalities in OffSet due to two AND-encodings and those in OnSet due to one OR-encoding.

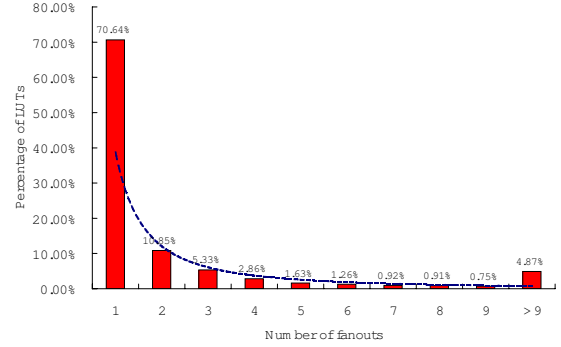


Figure 6: Fanout distribution of 20 MCNC benchmarks (mapped by 6-LUTs)

a mathematical formulation. Essentially, PMD selects a subset of LUT to duplicate. For each of those duplicated LUTs, the encoding schemes for its fanout LUTs are decided. We formulate the PMD problem as an ILP shown in (5) by introducing a 0-1 variable d_L with the same meaning as it is in Formulation (4), and a 0-1 variable e_L^j , where $e_L^j = 1$ if option $\phi(L, j)$ is used for encoding of the fanouts for LUT L .

$$\begin{aligned}
 & \text{Maximize} && \sum_{L \in Luts(C)} \sum_{j=1}^{|\Phi(L)|} w_L^j \cdot e_L^j \\
 & \text{Such that} && d_L = \sum_{j=1}^{|\Phi(L)|} e_L^j, && \forall L \in Luts(C) \\
 & && d_L + \sum_{f \in fanin(L)} \sum_{j=1}^{|\Phi(L)|} g_{f,j}^L \cdot e_f^j \leq S_L, && \forall L \in Luts(C) \\
 & && d_L \in \{0, 1\} \quad e_L^j \in \{0, 1\}
 \end{aligned} \tag{5}$$

where the first set of constraints guarantee that one and only one encoding scheme is used if an LUT L is duplicated, and the second set of constraints guarantee that the required resource (duplication and encoding) for an LUT does not exceed the available amount (similar to the FMD constraints), and constant coefficient $g_{f,j}^L$ is defined as

$$g_{f,j}^L = \begin{cases} 1, & \text{if } L \in \phi(f, j) \\ 0, & \text{otherwise.} \end{cases}$$

It is easy to show that Formulation (5) is a generalized case of Formulation (4). Particularly, for a circuit with a tree structure, where each LUT has only one fanout, we have $|\Phi(L)| \equiv 1$ and thereby $d_L \equiv e_L$ for all LUT L . In this case, d_L and e_L can be merged as one variable, and therefore Formulation (5) is reduced to Formulation (4). Since we have shown that Formulation (4) is NP-hard in Theorem 1, we have

THEOREM 2. *The decision version of problem (5) is NP-hard.*

Again we show that Formulation (5) has a nice generalized network flow structure, which means it can be solved efficiently.

4.1 Example

Consider the circuit example shown in Figure 3. Only AND-encoding is assumed for the sake of simplicity. The corresponding GNF for PMD has the similar structure as its counterpart for FMD, except that for each node d_L corresponding to LUT L we introduce a set of new nodes e_L^j to represent encoding options $\phi(L, j)$ for L 's fanouts. An arc (e_L^j, s_K) is added if LUT K is one of the fanouts of L and it is encoded in option $\phi(L, j)$. A gain factor is associated with each node e_L^j , and it is equal to the number of encoded fanouts in option $\phi(L, j)$. For instance, there are 2 fanouts of LUT A and therefore $|\Phi(A)| = 3$, i.e., $\phi(A, 1) = \{C\}$ (only encoding LUT C), $\phi(A, 2) = \{D\}$ (only encoding LUT D) and $\phi(A, 3) = \{C, D\}$ (encoding both LUT C and LUT D), if LUT A is duplicated. It is easy to verify that this GNF is equivalent to Formulation (5). A feasible solution of the GNF in Figure 7(a) is shown in bold arcs and the corresponding PMD result is shown in Figure 7(b), where three LUTs (A, B, D) are duplicated and three LUTs (C, D, E) are encoded.

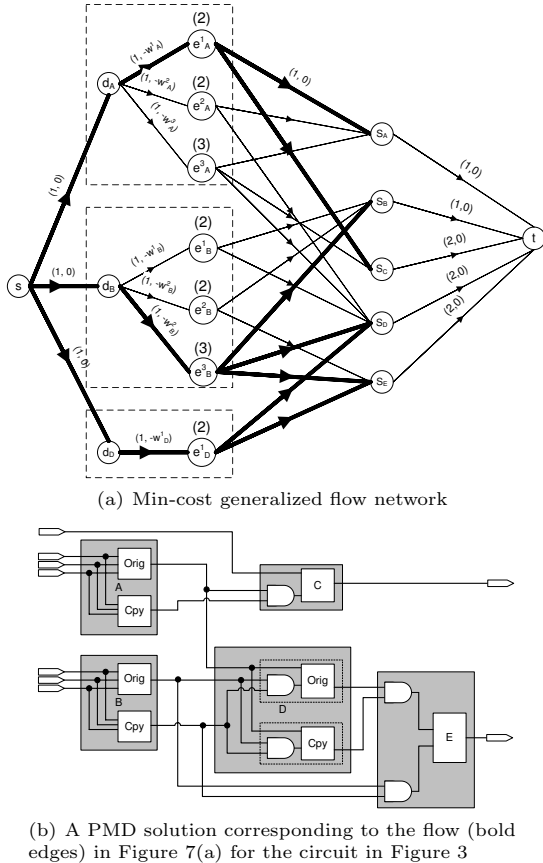


Figure 7: PMD example

5. EXPERIMENTAL RESULTS

The proposed FMD and PMD algorithms are implemented in C++ with the mosek [14] solver on a Ubuntu server with Xeon 2.4GHz CPU and 2Gb of memory. The 20 biggest MCNC benchmarks are tested. Throughout the experiments, the 6-input 2-output LUT structure and a cluster size of 4 in Xilinx Virtex-5 FPGA architecture is assumed. All benchmarks are first mapped by Berkeley ABC technology mapper [15] with edge flow optimization, which has a special property that the mapped circuits are more suitable to be packed into dual-output LUTs [5]. The original circuits are clustered by TV-Pack packing tool, and placed and routed by VPR [16]. Then, the routing resource utilization (including the configuration bits in the switch boxes and the connection

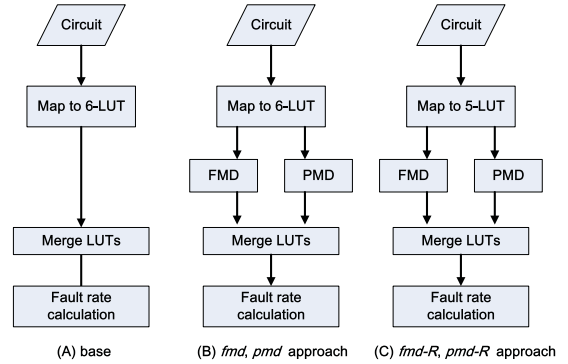


Figure 8: Experimental flows

boxes along the path of an interconnect) is extracted. As the dedicated clustering, placement, and routing that consider dual-output LUTs are not available in the current TV-Pack and VPR tool set, we use the post-routing results produced by the original circuit to estimate the routing resource utilization for the optimized (by FMD/PMD) circuits with dual-output LUTs⁴. Since the number of routing configuration bits of a net may increase after duplication (due to the extra interconnect of the duplicating output), we estimate the number of configuration bits after routing by the following empirical rule: the wire length (proportional to the number of configuration bits) is proportional to the square root of the area (proportional to the number of LUTs). Once the number of configuration bits in the interconnects is obtained, we perform the full-chip logic simulation to calculate the fault rate, considering the single fault in both LUTs and interconnects. The criticality and the fault rate calculation can be done by an FPGA-based emulator, which takes less than 1s for runtime per circuit. In the set of experiments, we assume $\gamma = 1.0$ in (3) for all FPGA devices. In addition, we assume the same FPGA device for each circuit, i.e., the same A in eq. 3.

Three sets of CAD flows (shown in Figure 8) are compared in our experiments. In the baseline algorithm (“base” in Figure 8 (a)), benchmarks are first mapped to 6-LUTs and then the LUT merge algorithm in [4] is used to merge pairs of small LUTs (<4 inputs) into dual-output LUTs. In the second set of CAD flows (“*fmd*” and “*pmd*” in Figure 8 (b)), the proposed resynthesis (FMD or PMD) is performed on the 6-LUT-mapped circuits, and then the aforementioned LUT merge algorithm is applied. Note that both *fmd* and *pmd* preserve logic depths of the baseline flow. To further explore the potential of the proposed algorithms, the third set of CAD flows (“*fmd-R*” and “*pmd-R*” in Figure 8 (c)) adds explicit area redundancy by first mapping the benchmarks with 5-LUTs. This reserves at least one spare pin for each LUT to increase the opportunities of duplication by FMD and PMD.

The experimental results for the baseline algorithm and the proposed resynthesis flows (*fmd*, *pmd*, *fmd-R* and *pmd-R*) are summarized in Table 1. It shows that the four proposed resynthesis flows achieve different tradeoff among fault tolerance, area and performance. For combinational circuits as shown in the upper part of Table 1, *fmd* increase MTTF⁵ by 14% without area overhead, and *pmd* increases MTTF by 27% with 4% area overhead, while both preserving the logic depth, compared to the baseline algorithm. For sequential circuits, both *fmd* and *pmd* increases MTTF by about 23% without area overhead. We found that the majority of the duplication and encoding for sequential circuits occur in those LUTs lying in the boundary defined by FFs, and those LUTs that drive FFs mostly do not have multiple fanouts. Therefore, *fmd* and *pmd* produce very similar results for sequential circuits.

With explicit area redundancy, more significant improvement of MTTF is obtained by both *fmd-R* and *pmd-R*. Specifically, *pmd-R*

⁴ Note that the area and performance of the duplicated circuits can be further improved using intelligent clustering, placement and routing algorithms, which absorb the duplicated interconnects into CLBs and therefore minimize the impact to the global routing. These algorithms will be our future work.

⁵ MTTF ratio is calculated based on the reciprocal of the fault rate for each benchmark circuit [17].

| combinational circuits | | | | | | | | | | | | | | | | |
|-------------------------|-----|-----|------|-----------|--------|--------|--------|--------|------------|--------|--------|--------|--------|--------------|-------------|--|
| circuit characteristics | | | | # of LUTs | | | | | fault rate | | | | | logic depth | | |
| circuit | PI# | PO# | reg# | base | fmd | pmd | fmd-R | pmd-R | base | fmd | pmd | fmd-R | pmd-R | base/fmd/pmd | fmd-R/pmd-R | |
| alu4 | 14 | 8 | 0 | 442 | 442 | 471 | 537 | 594 | 0.51% | 0.48% | 0.41% | 0.33% | 0.25% | 5 | 6 | |
| apex2 | 39 | 3 | 0 | 600 | 602 | 634 | 741 | 779 | 0.43% | 0.36% | 0.35% | 0.23% | 0.20% | 6 | 7 | |
| apex4 | 9 | 19 | 0 | 533 | 537 | 550 | 651 | 664 | 1.67% | 1.40% | 1.23% | 1.09% | 0.71% | 5 | 6 | |
| des | 256 | 245 | 0 | 531 | 531 | 531 | 949 | 950 | 1.94% | 1.94% | 1.94% | 1.28% | 1.07% | 4 | 4 | |
| ex1010 | 10 | 10 | 0 | 642 | 647 | 652 | 735 | 781 | 1.66% | 1.25% | 1.22% | 1.14% | 0.95% | 5 | 6 | |
| exp5p | 8 | 63 | 0 | 337 | 340 | 348 | 436 | 437 | 1.07% | 0.98% | 0.76% | 0.71% | 0.50% | 4 | 5 | |
| misex3 | 14 | 14 | 0 | 425 | 428 | 463 | 506 | 553 | 0.82% | 0.71% | 0.58% | 0.50% | 0.33% | 5 | 6 | |
| pdc | 16 | 40 | 0 | 1377 | 1385 | 1397 | 1910 | 1911 | 1.29% | 1.15% | 1.04% | 0.77% | 0.60% | 7 | 8 | |
| seq | 41 | 35 | 0 | 629 | 631 | 661 | 804 | 850 | 0.90% | 0.79% | 0.70% | 0.53% | 0.42% | 5 | 5 | |
| spla | 16 | 46 | 0 | 1296 | 1300 | 1312 | 1836 | 1842 | 1.63% | 1.42% | 1.39% | 0.99% | 0.85% | 7 | 8 | |
| GeoMean | 21 | 24 | 0 | 615 | 618 | 638 | 804 | 835 | 1.07% | 0.94% | 0.85% | 0.67% | 0.51% | 5.21 | 5.98 | |
| Ratio | | | | 1 | 1.0047 | 1.0366 | 1.3066 | 1.3577 | 1 | 0.8769 | 0.7904 | 0.6232 | 0.4801 | 1.00 | 1.15 | |
| MTTF Ratio | | | | | | | | | 1 | 1.1410 | 1.2756 | 1.6476 | 2.1254 | | | |
| sequential circuits | | | | | | | | | | | | | | | | |
| bigkey | 263 | 197 | 224 | 518 | 518 | 518 | 798 | 742 | 1.56% | 1.33% | 1.26% | 1.16% | 1.12% | 3 | 3 | |
| clma | 383 | 82 | 33 | 2868 | 2868 | 2869 | 3242 | 3234 | 0.10% | 0.08% | 0.08% | 0.07% | 0.07% | 9 | 11 | |
| diffeq | 28 | 3 | 305 | 534 | 536 | 537 | 573 | 566 | 1.20% | 0.98% | 0.96% | 0.75% | 0.68% | 8 | 10 | |
| dsip | 228 | 197 | 224 | 665 | 665 | 665 | 782 | 782 | 1.70% | 1.73% | 1.73% | 1.46% | 1.46% | 3 | 3 | |
| elliptic | 19 | 2 | 194 | 322 | 323 | 323 | 271 | 270 | 1.18% | 0.95% | 0.94% | 0.90% | 0.88% | 6 | 8 | |
| frisc | 20 | 116 | 886 | 1868 | 1868 | 1872 | 2233 | 2228 | 1.38% | 1.15% | 1.14% | 0.66% | 0.64% | 14 | 17 | |
| s298 | 3 | 6 | 14 | 20 | 20 | 20 | 24 | 23 | 2.03% | 1.50% | 1.50% | 1.25% | 1.19% | 2 | 2 | |
| s38417 | 29 | 106 | 1462 | 1991 | 1993 | 2005 | 2338 | 2314 | 1.70% | 1.43% | 1.39% | 1.20% | 1.05% | 6 | 8 | |
| s38584.1 | 39 | 304 | 1260 | 1972 | 1976 | 1985 | 2408 | 2378 | 1.51% | 1.27% | 1.24% | 1.01% | 0.95% | 7 | 8 | |
| tseng | 52 | 122 | 385 | 640 | 664 | 667 | 743 | 741 | 1.51% | 1.16% | 1.13% | 0.86% | 0.82% | 7 | 10 | |
| GeoMean | 46 | 47 | 247 | 661 | 664 | 665 | 767 | 755 | 1.15% | 0.95% | 0.94% | 0.77% | 0.73% | 5.63 | 6.66 | |
| Ratio | | | | 1 | 1.0047 | 1.0066 | 1.1601 | 1.1415 | 1 | 0.8287 | 0.8145 | 0.6649 | 0.6338 | 1.00 | 1.18 | |
| MTTF Ratio | | | | | | | | | 1 | 1.2337 | 1.2426 | 1.4700 | 1.5023 | | | |

Table 1: Summary of experimental results

increases MTTF by 113% with 36% area overhead and 15% logic depth increase for combinational circuits, and it increases MTTF by 50% with 14% area overhead and 18% logic depth increase for sequential circuits. Remarkably, *pmd-R* reduces fault rate by about 2 \times with only 24% area overhead and 20% logic depth increase for “apex4”. Note that the effectiveness of *fmd-R* and *pmd-R* for sequential circuits is less significant than that for combinational circuits. In fact, optimization such as technology mapping and retiming can be performed before or simultaneously with *fmd-R* and *pmd-R* to create more opportunities for the duplication and encoding. This will be an interesting future research topic.

The runtime for our proposed algorithms are less than 2 minutes for all cases. In addition, it is interesting that although *fmd-R* duplicates fewer LUTs than *pmd-R*, it results in larger area for sequential benchmarks (767 vs. 755 on average). This is because more small LUTs (<4 inputs) are used for duplication and encoding in *fmd-R* due to the fully-masking constraint, and as a result the LUT merge algorithm becomes less effective for *fmd-R*.

6. CONCLUSIONS AND FUTURE WORK

We have presented a fault-tolerant post-mapping resynthesis which explicitly exploits the dual-output LUT architecture. The LUT duplication problem is formulated as a generalized network flow problem, which can be solved very efficiently. Our experimental results are encouraging and show 27% MTTF improvement with 4% area overhead and 2 \times MTTF improvement with 36% explicit area overhead, compared to ABC technology mapping.

In the future, the following research directions can be explored. As our experiments present the area-robustness-performance trade-off of the two extreme cases (mapping with 6-LUTs and 5-LUTs), simultaneous technology mapping and resynthesis should be studied in order to traverse the Pareto points in the solution space and to maximize the robustness under area/performance constraints. In addition, besides AND-Encoding and OR-Encoding, other logics (e.g., NAND and NOR, or other more complicated logics) can also be used to mask faults in their upstream circuit. Furthermore, while the proposed technique places masking logic immediately after a duplicated LUT, the area overhead may be reduced by placing masking logics at the end of a series of duplications. Finally, we will study the dual-output-aware physical synthesis, including clustering, placement and routing, to mitigate the impact of the proposed approach on performance.

7. REFERENCES

- [1] A. Djupdal and P. C. Haddow, “Yield enhancing defect tolerance techniques for FPGAs,” in *MAPLD*, 2006.
- [2] A. Cosoroaba and F. Rivoallon, “Achieving higher system performance with the virtex-5 family of FPGAs,”
- [3] “Altera stratix II device handbook,” 2007.
- [4] T. Ahmed, P. D. Kundarewich, J. H. Anderson, B. L. Taylor, and R. Aggarwal, “Architecture-specific packing for virtex-5 FPGAs,” in *FPGA*, 2007.
- [5] S. Jang, B. Chan, K. Chung, and A. Mishchenko, “Wiremap: FPGA technology mapping for improved routability,” in *FPGA*, 2008.
- [6] K. Chapman and L. Jones, “SEU Strategies for Virtex-5 Devices,” 2009.
- [7] R. Lyons and W. Vanderkulk, “The use of triple-modular redundancy to improve computer reliability,” *IBM Journal of Research and Development*, 1962.
- [8] H. Naeimi, “A greedy algorithm for tolerating defective crosspoints in NanoPLA design,” in *Master Thesis, California Institute of Technology*, 2005.
- [9] M. Joshi and W. Al-Assadi, *Development and Analysis of Defect Tolerant Bipartite Mapping Techniques for Programmable cross-points in Nanofabric Architecture*. Springer Netherlands, 2007.
- [10] Y. Hu, Z. Feng, R. Majumdar, and L. He, “Robust FPGA resynthesis based on fault tolerant boolean matching,” in *ICCAD*, 2008.
- [11] S. Krishnaswamy, S. Plaza, I. L. Markov, and J. P. Hayes, “Signature-based SER analysis and design of logic circuits,” in *TCAD*, 2009.
- [12] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.
- [13] T.-Y. Hsu and T.-C. Wang, “A generalized network flow based algorithm for power-aware FPGA memory mapping,” in *DAC*, 2008.
- [14] “mosek,” in <http://www.mosek.com>.
- [15] “ABC: A system for sequential synthesis and verification.”
- [16] V. Betz and J. Rose, “VPR: A new packing, placement and routing tool for fpga research,” pp. 213–222, 1997.
- [17] S. Mukherjee, J. Emer, and S. K. Reinhardt, “Radiation-induced soft errors: An architectural perspective,” in *HPCA*, 2005.