

ECO Timing Optimization Using Spare Cells and Technology Remapping

Yen-Pin Chen¹, Jia-Wei Fang², and Yao-Wen Chang^{1,2}

¹Department of Electrical Engineering, National Taiwan University, Taipei 106, Taiwan

²Graduate Institute of Electronics Engineering, National Taiwan University, Taipei 106, Taiwan

Abstract—Spare cells rewiring is a technique used to fix defects or deficiencies after the placement stage. It is traditionally done by manual work but becomes extremely hard nowadays. In this paper, we propose a spare cells selection algorithm consisting of two phases to optimize timing of the circuit by rewiring spare cells. In the first phase, we apply gate sizing and buffer insertion to all timing violated paths to fix timing violations. In the second phase we further fix timing violations by extracting timing critical parts and applies technology remapping to them. Experimental results based on five industrial benchmarks show that our algorithm can fix up to 99.82% of the total negative slack. The runtime is very short and linear to the gate count. The performance proves that our algorithm is efficient and effective.

I. INTRODUCTION

A. Spare Cells for Timing Optimization

ECO (Engineering Change Order) is usually performed during the chip implementation cycle. If engineers need to change only a small portion of the netlist in a very short time, running the traditional back-end design flow to the whole netlist is very time-consuming. The most efficient way is to change the netlist locally without affecting other parts of the chip. Using spare cells is a good choice for this purpose because rewiring the circuit by spare cells can change the netlist without changing the chip placement. Engineers do not need to run placement tools to place the netlist after the rewiring process. Since timing closure is hard to be achieved in today's nanometer designs and engineers have to run the back-end design flow many times to meet timing constraints, using spare cells to do netlist changes can save a lot of time and effort. Besides, if masks are already produced before netlist change, rewiring the netlist using spare cells only needs the masks of the routing layers to be re-produced. This will save a large amount of production cost because masks are quite expensive in the nanometer process.

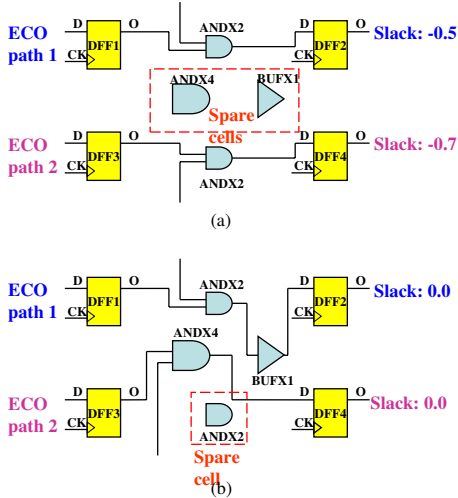


Fig. 1. (a) ECO paths before rewiring. (b) ECO paths after rewiring.

Although spare cells rewiring is a very effective ECO technique, using it to fix timing deficiencies is getting tougher and tougher nowadays. This is because the local netlist change cannot consider its effect on the circuit timing and makes the circuit fail to achieve timing closure. Additionally, increasing of the gate count of chip designs also makes the problem

substantially harder. Thus we need an efficient algorithm to deal with the problem of timing optimization by spare cells.

Figure 1 shows an instance of timing optimization by rewiring spare cells. The AND gate ANDX4 and BUFFER gate BUF1 are spare cells and not connected to any path. Gate DFF1, gate DFF2, gate DFF3, and gate DFF4 are D flip-flops. They are start points and end points of path 1 and path 2. Arrival times of DFF2 and DFF4 are larger than the clock cycle, and the timing of path 1 and path 2 needs to be fixed by ECO. Thus we call these two timing violated paths ECO paths. We can improve the timing of ECO path 1 by inserting the adjacent BUFFER gate BUF1 in that path to help driving the load. To fix the timing of ECO path 2, we can use the AND gate ANDX4 instead of the AND gate ANDX2 on ECO path 2 because ANDX4 has a larger driving capability. After the sizing and buffering operations, both ECO paths meet the timing constraints. The AND gate ANDX2 is now released from the netlist and becomes a spare cell.

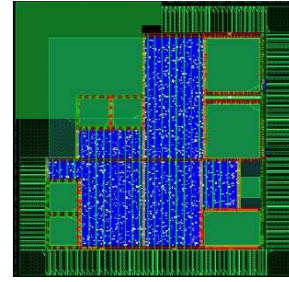


Fig. 2. Spare cells placement.

Spare cells are designed for further design changes, and are evenly placed on the chip layout. The type and number of spare cells vary from different chip characteristics, and are usually determined empirically. The number of spare cells is usually small compared to other standard cells. Thus using spare cells to perform ECO operations needs to consider the resource sharing problem. Figure 2 by [4] shows the spare cells placement. The spare cells are plotted as points, and they are spread throughout the standard cell area for good performance.

B. Previous Work

Spare cells selection is a very common technique in industrial designs, and it was traditionally done by manual labor. Since the gate count of a chip design is increasing and timing closure is more and more difficult to be achieved, this problem becomes so tough that it can no longer be solved by manual ways. To the best knowledge of the authors, there is still no published literature for the ECO timing optimization by spare cells. There are two topics related to this problem: (1) buffer insertion and gate sizing problem, and (2) physical and logic co-synthesis.

Buffer insertion is a well-known technique for timing optimization. It can not only reduce the path delay but also eliminate signal noise. Authors of [5] proposed a dynamic programming method for the buffer insertion problem. When the candidate buffer locations of a signal net are known, the dynamic programming method (VGDP) can find the maximum timing slack solution in quadratic time. Many works are proposed based on this method, and they can be grouped into three categories: (1) net based buffer insertion algorithms [27] [6], (2) network flow based buffer insertion algorithms [13] [14], and (3) path based buffer insertion algorithms [28].

Logic synthesis and placement are two important stages of the IC design flow. Logic synthesis tools translate functions into logic gates, and placement tools place them on the chip layout. Since the logic synthesis tools do not have the information of exact positions of gates, they cannot fully optimize the circuits in the right way. On the other hand, placement tools have limited flexibility to optimize because they cannot change the

circuit netlist. Thus combining the two stages has been an attractive topic in recent years. Previous works can be grouped into two categories: (1) layout driven logic synthesis [10], [11], [23], and [25], and (2) local netlist transformation [17] [21].

C. Our Contributions

To our best knowledge, this paper is the first work for the ECO timing optimization by spare cells rewiring. The major difference between this problem and the traditional buffer insertion problem is the cost metric. The traditional buffer insertion problem is to insert buffers at some candidate locations. These candidate locations are not related to the placement and the placement after buffering may overlap. The cost of buffering is known, such as area overhead or buffer delay. The ECO timing optimization problem considers rewiring the netlist with spare cells, and all spare cells are the candidate buffering/sizing locations for each net/gate. Since every spare cell becomes a normal standard cell if it is rewired to the netlist, the candidate buffering locations of each net vary along the optimization process. Thus the buffering/sizing cost is dynamic, making this problem much harder than the traditional buffer insertion problems. We propose a dynamic programming method considering such dynamic cost, and so we call our algorithm “Dynamic Cost Programming (DCP)”.

Our spare cells selection algorithm consists of two phases. The first phase is buffer insertion and gate sizing. We iteratively perform buffer insertion and gate sizing simultaneously to the ECO paths. This loop terminates until all timing violations of ECO paths are fixed or every timing violated ECO path cannot be further optimized. We also propose several heuristics to reduce the solution size during dynamic cost programming.

The second phase is technology remapping. We extract timing critical parts of the ECO paths and remap it using spare cells. From our proposed optimization flow, our method can be smoothly integrated into commercial a design flow. Experimental results show that our algorithm can fix almost all of the timing deficiencies in a short CPU time.

II. PROBLEM FORMULATION

In this chapter, we introduce the notations used in this paper and the problem formulation. A timing path is defined as (1) a path from one primary input to one primary output, (2) a path from one primary input to one D flip-flop input pin, (3) a path from one D flip-flop output pin to one primary output, and (4) a path between one D flip-flop output pin and one D flip-flop input pin. An ECO path is a timing path which violates the timing constraints and we are going to fix it by spare cell rewiring. We denote the start point of a ECO path as the D flip-flop or the primary input at the beginning of the ECO path. We also denote the end point of a ECO path as the D flip-flop or the primary output at the end of the ECO path. Figure 3 (a) shows the modeling of the ECO path. Let N be the set of nets of the netlist and N^E be set of nets of the ECO paths. Let G be the set of all standard cells. We denote G^E as the cells on the ECO paths, and G^S as the spare cells. The coordinate of a gate g_i is denoted as $p_i(x_i, y_i)$. We define the buffering and sizing operations below.

Definition 1: A buffering operation is to insert a buffer type spare cell g_i^S into a net n_j^E in the ECO paths. A gate sizing operation is to exchange a spare cell g_i^S with a gate g_j^E in the ECO paths by rewiring.

Definition 2: The delay of a gate g_i is $\text{delay}(g_i)$ while the delay of g_i after sizing or buffering operations is $\text{delay}'(g_i)$.

During our optimization process, we generate many solutions. A solution can be formally defined as follows:

Definition 3: The target gate of a solution is a gate on the ECO path and this gate is being considered to be sized or buffered.

Definition 4: The scope of a solution is a sub-path between the ECO path end point and the target gate. The delay of the scope is the sum of delays of the gates in the scope.

Definition 5: A solution is a set of gate sizing and buffering operations to gates and nets in its scope. The cost of a solution is the sum of delays of gates in the scope if operations of the solution are performed.

Figure 3 (b) shows a solution S_1 corresponding to the ECO path shown in Figure 3 (a). Solution S_1 consists of one buffering operation and one sizing operation. The scope of S_1 is the sub-path between g_2^E and g_5^E , and the cost of S_1 is the sum of delays of g_2^E , g_1^E , g_4^E , g_3^E , and g_5^E . Another solution S_2 is shown in Figure 3 (c) with the same scope, the cost of S_2 is the sum of delays of gates in the scope after inserting two buffers.

At the end of the optimization process, we get a set of solutions which scope is the whole ECO path and the cost is the ECO path delay. We choose a solution meeting timing constraints that uses minimum number of buffers as our final solution. Operations of the final solution are performed to the ECO path and STA is re-run to update the timing information.

Based on the definition above, the ECO timing optimization problem can be formally defined as the follows:

Problem 1: Given a netlist after ECO process, the ECO timing optimization problem is to re-wire the netlist using spare cells G^S so that the netlist meets the timing constraints, and the functionality of the netlist cannot be unchanged.

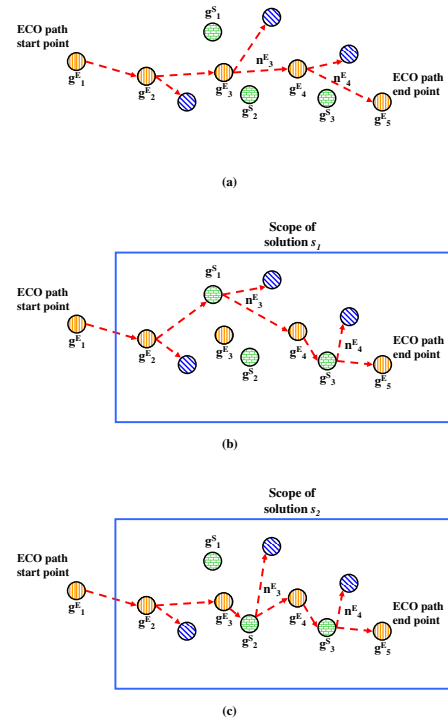


Fig. 3. (a) The model of ECO paths. (b) Solution S_1 . (c) Solution S_2 .

III. THE SPARE CELLS SELECTION ALGORITHM

In this chapter, we present our spare cells selection algorithm for the ECO timing optimization problem. First we give the overview of our algorithm and the optimization flow. Then we detail the methods used in each phase.

A. Algorithm Overview

We propose a two-phase flow to solve the ECO timing optimization problem. The input information is a placed netlist with some timing violations. We use static timing analysis (STA) to identify those timing violated paths as ECO paths. Then we enter the first phase.

In the first phase we iteratively choose the ECO path with the largest negative slack and optimize it by gate sizing and buffering operations. Since ECO paths are usually overlapped, we update netlist and timing information of all ECO paths after optimization to one of them. If the processed path is not improved after optimization, we put it into a denied list. It means that this path can not be improved under the current netlist and the status of spare cells. Then we choose the ECO path whose negative slack is the largest among all ECO paths not in the denied list and optimize it. Whenever optimization to an ECO path is effective and the timing of that path is improved, we clear the denied list to make paths in the denied list be able to be optimized again. This is because operations to one ECO path may change the status of spare cells or modify the structure of some paths in the denied list. This criterion prevents us from continuously trying to optimize an ECO path which can not be further improved under current situation. This loop continues until all ECO paths meet the timing constraints or all remained timing violated ECO paths can not be improved any more.

After the first phase, we fix most timing violations. If there are still violations left, we forward those information to the second phase. In the second phase, we identify timing critical parts of the netlist. Those parts are extracted from the netlist and remapped using spare cells. The remapping process stopped when all timing violations are fixed or no more paths can be optimized. After the second phase, we write the optimized netlist to a DEF file, and the program terminates. The optimization flow is shown in Figure 4.

B. The Timing Model

We apply the Synopsys Liberty library format as our timing model to evaluate the circuit timing. Although the wirelength is the commonest cost metric during placement, we can use a more accurate timing model because the placement is fixed and all cells' locations are known. The timing model combines wire loading and gate input capacitance with gate driving loading. The coordinate of a gate g_i is denoted as $p_i(x_i, y_i)$. The relation between

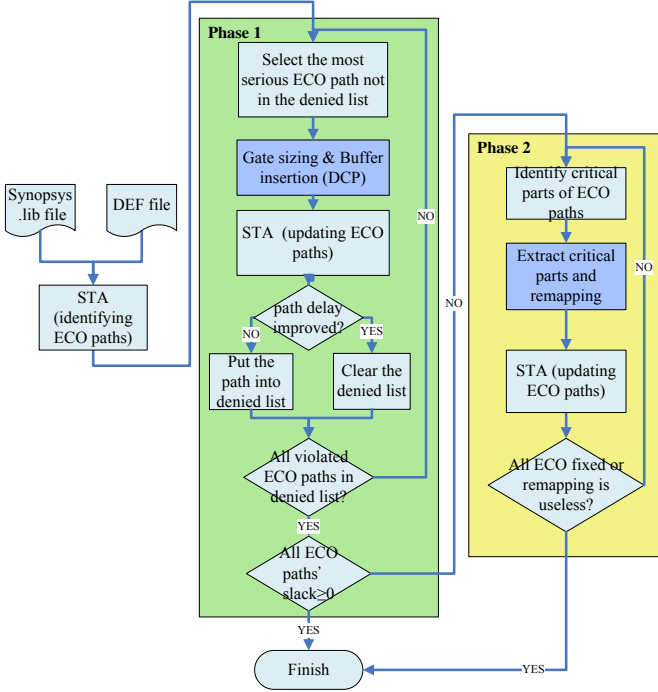


Fig. 4. The ECO timing optimization flow.

fanout wirelength of a gate g_i and its corresponding capacitance loading value C_{Wi} is defined as:

$$C_{Wi} = \sum_{g_j \in \text{fanouts of } g_i} (|x_i - x_j| + |y_i - y_j|) \times \phi, \quad (1)$$

where ϕ is the amount of capacitance from per unit wirelength. Then the capacitance loading of gate g_i can be defined as:

$$C_i = C_{Wi} + C_{Oi} + \sum_{g_j \in \text{fanout of } g_i} C_{Ij}, \quad (2)$$

where C_{Oi} is the output pin capacitance of gate g_i , and C_{Ij} is the input pin capacitance of gate g_j . The delay and the output transition time of a gate are functions of its input transition time and output driving capacitance, and the functions are characterized by lookup tables. An example of the lookup table for a gate's delay when the output signal is falling is shown in Figure 5.

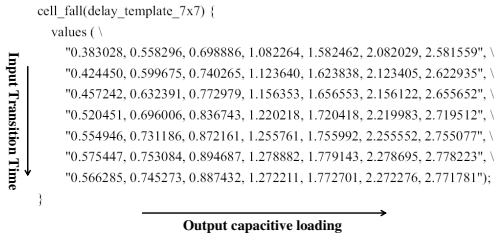


Fig. 5. Lookup table example.

It is obvious that if a gate and its fanouts are far away, the gate must have a large delay due to its large capacitive loading. The timing path delay is the sum of delays of all gates on the timing path.

There are several properties of this timing model:

- 1) Loading dominance: For a gate g_i , the effect of the output loading to the gate delay is much larger than that of the input transition time.
- 2) Shielding: If two gates, g_i and g_j , on the same ECO path are not directly connected, and g_i is the former gate while g_j is the latter

gate in the path, then gate sizing to g_j does not affect the gate delay of g_i . On the other hand, gate sizing to g_i has only a little effect to g_j because of the loading dominance property.

The first property is summarized from the technology data empirically. This property is important because during our optimization process, we only know one factor among the output loading and the input slope when calculating the gate delay. We assume a value to the input slope while we have an exact output loading value to get an approximated gate delay. The delay calculated in this way is more accurate than the case that the output loading is an assumed value and the input slope is known.

The second property means that changing one part of the ECO path by gate sizing and buffering does little effect on timing of the unchanged part. It facilitates us to judge one operation by its local delay effect other than its global delay effect. This will speed up the algorithm greatly. Figure 6 shows a conceptual example.

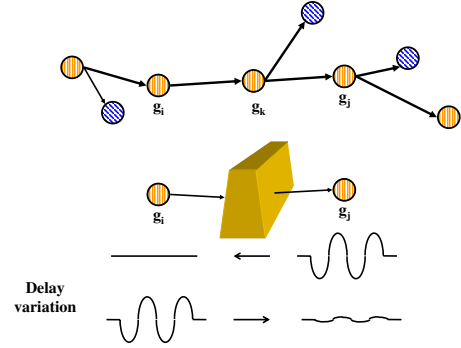


Fig. 6. The gate g_k serves as a barrier and mitigates the delay variation of g_i and g_j due to changes of the other gate.

C. Gate Sizing and Buffer Insertion

In this section, we present the Dynamic Cost Programming (DCP) algorithm which uses sizing/buffering operations to optimize the timing of a path. This algorithm is named DCP because it is based on the dynamic programming framework and considers dynamic cost. Figure 7 shows an overview of the DCP algorithm.

1) *Delay Cost Calculation*: In our optimization process, we have two timing-calculation operations. The first one is the static timing analysis (STA), which is a well known technique. The second one is applied during dynamic cost programming and different from STA because it only calculates the timing of a small region. We detail the second timing-calculation operation in the following paragraphs.

From the Synopsys Liberty timing model, we have two observations for our ECO timing optimization problem:

- 1) The buffering to one net n_i^E changes the delay of the source of n_i^E , g_i^E , while other gates are little affected or not affected. Thus the effect of buffering to the timing is the delay change of g_i^E and the delay increase of the inserted buffer.
- 2) The sizing to the gate g_i^E changes the delay of the fanin gates of g_i^E , G_j^E , while other gates are little affected or not affected. Thus the effect of sizing g_i^E to the timing is the delay change of G_j^E and the sized gate.

Based on the above description, we can evaluate the effect of buffering and sizing by calculating delay of the changed part of the path without applying STA to the whole path. An illustration of the above observations is shown in Figure 8.

It is important that the delay value calculated in the second operation is an approximated value. As described in Chapter III-B, the gate delay is a function of the input transition time and the output loading. Since the dynamic programming method optimizes the path along one direction, we cannot know both factors of a gate at the same time. Thus we apply dynamic programming method from the end point of the ECO path to the start point of the ECO path. During sizing and buffering we calculate the gate delay as costs with a known output loading value, while the input slope is assumed to be the input slope calculated by STA before the optimization to this path.

On the other hand, if we apply dynamic programming from the start point to the end point of the ECO path, the gate delay is calculated with a known input slope value while the output loading is unknown. Gate delays calculated in this way have larger error because estimation error of the output loading causes larger delay error than that of the input slope due to loading dominance property.

Algorithm: Dynamic Cost Programming(P, G^S, N^E, G^E)

P : the ECO path to be optimized;
 G^S : set of all spare cells;
 $N^E \{n_1^E, \dots, n_{M-1}^E\}$: set of all nets on the target ECO path;
 $G^E \{g_1^E, \dots, g_M^E\}$: set of all cells on the target ECO path;
 M : size of G^E
 S_i^N : set of solutions stored for buffering net n_i^E
 S_i^G : set of solutions stored for sizing gate g_i^E
1 begin
2 Merge G^E and fanouts of $\{g_1^E, \dots, g_{M-1}^E\}$ into a routing tree
3 for $i = M - 1 \rightarrow 2$
4 for all buffer-type spare cells $g_j^S \in G^S$ in bounding box of n_i^E
5 apply buffer insertion to n_i^E using g_j^S based on S_{i+1}^G
6 store the assignment in S_i^N if $\text{delay}'(g_i^E) + \text{delay}'(g_j^S) < \text{delay}(g_i^E)$
7 $\text{delay}(g_i^E)$
8 prune the solutions in S_i^N
9 for all spare cells $g_j^S \in G^S$ with type same as g_i^E in the
10 bounding polygon of g_i^E
11 apply gate sizing to g_i^E using g_j^S based on S_i^N
12 store the assignment in S_i^G if $\text{delay}'(g_{i-1}^E) < \text{delay}(g_{i-1}^E)$
13 & $\text{delay}'(g_{i-1}^E) + \text{delay}'(g_j^S) < \text{delay}(g_{i-1}^E) + \text{delay}(g_i^E)$
14 prune the solutions in S_i^G
15 for all buffer-type spare cells $g_j^S \in G^S$ in bounding box of n_1^E
16 apply buffer insertion to n_1^E using g_j^S based on S_2^G
17 store assignment in S_1^N if $\text{delay}'(g_1^E) + \text{delay}'(g_j^S) < \text{delay}(g_1^E)$
18 choose the solution in S_1^N that meets the timing constraints and
19 uses fewest buffers, and rewire the netlist according to the solution
20 end

Fig. 7. Overview of the DCP algorithm.

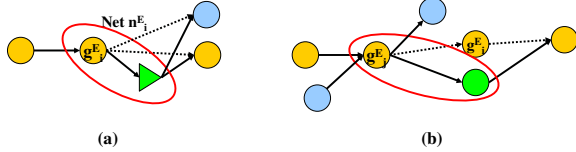


Fig. 8. (a) The delay effect of a buffering operation. (b) The delay effect of a gate sizing operation.

2) Sizing and Buffering Operations: Here we describe the overall process of the DCP algorithm.

Given an ECO path to be optimized, we merge the nets of the ECO path into a big routing tree. Fanouts of gates along the path are also merged into the tree because they affect the loading of the ECO path. M is the number of gates along the ECO path. From the start point to the end point, gates on the ECO path are numbered as g_i^E , $i = 1, \dots, M$, and nets of the path are numbered as n_i^E , $i = 1, \dots, M - 1$. We perform one gate sizing operation to each gate g_i^E , and one buffering operation for each net n_i^E . Although a net may need more than one buffer to fix the timing, we only insert one buffer into the net at one time. This is because the structure of the routing tree is undetermined and the number of possible buffer assignments is too large to be considered. So we insert one buffer into a net in one iteration and insert more buffers into the net in latter iterations if needed.

We start applying buffering to net n_{M-1}^E . We try to insert one buffer in the neighborhood into n_{M-1}^E and calculate the approximated delay of the driving gate, g_{M-1}^E , of n_{M-1}^E . Each possible buffering assignment is a solution whose scope is the sub-path consisting of g_{M-1}^E and g_M^E (the end point), and we store it if sum of delay of the gates in its scope is smaller than the case without buffering. At this time we only estimate the effect of buffering n_{M-1}^E without actually inserting a buffer to the net. Figure 9 shows the result of buffering n_{M-1}^E ($M = 5$).

Since we only store the buffering solutions that reduces the delay of g_{M-1}^E , timing of paths that are overlapped with the ECO path by g_{M-1}^E

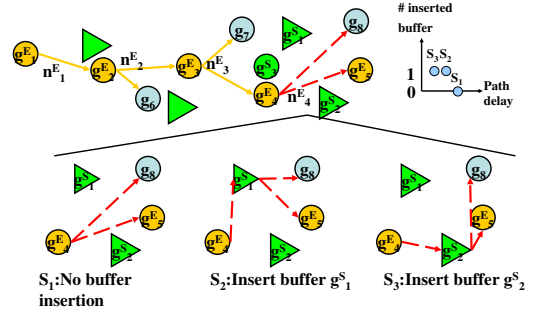


Fig. 9. The solutions for buffering n_4^E .

will not be worse than the case without buffering. This property is very important because we must guarantee no timing change to unprocessed paths (ex: $\{g_1^E, g_2^E, g_3^E, g_4^E, g_5^E\}$ in Figure 9) when optimizing ECO paths.

Additionally, since spare cells assignment at early stages will affect the assignment at latter stages, the cost of buffering and sizing is dynamic. Thus we must store a set of solutions instead of only the best one at every operation.

After buffering n_{M-1}^E , we apply gate sizing to g_{M-1}^E using nearby spare cells with the same type as g_{M-1}^E . Timing and loading information of every solution stored in n_{M-1}^E is considered to generate new solutions for g_{M-1}^E . For one sizing solution, if (1) $\text{delay}'(g_{M-2}^E) + \text{delay}'(\text{the sizing spare cell}) < \text{delay}(g_{M-2}^E) + \text{delay}(g_{M-1}^E)$, and (2) $\text{delay}'(g_{M-2}^E) < \text{delay}(g_{M-2}^E)$, we store this solution for g_{M-1}^E .

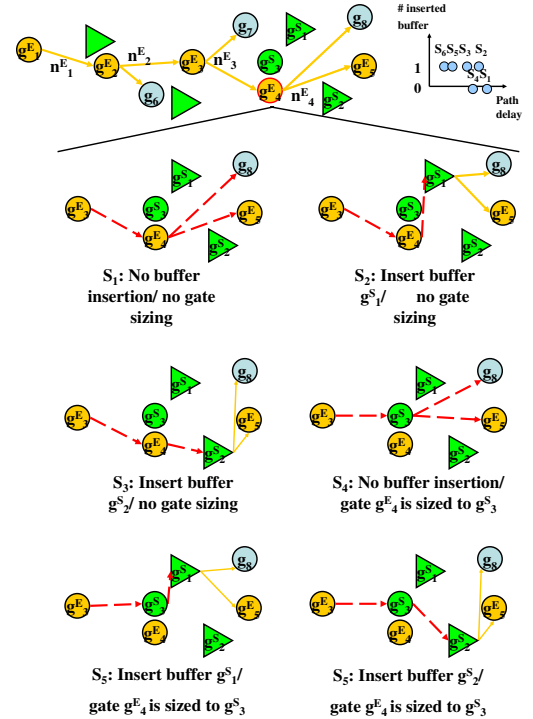


Fig. 10. The solutions for sizing g_4^E .

With the second criterion like buffering n_{M-1}^E , timing of paths overlapped with the ECO path by g_{M-2}^E is no worse than the case with no sizing. Figure 10 shows the solutions of sizing g_{M-1}^E ($M = 5$).

Then we apply buffer insertion to n_{M-2}^E . After buffering we size g_{M-2}^E . We recursively apply buffer insertion to n_{M-i}^E , $i = 1, \dots, M - 1$, and gate

sizing to g_{M-i}^E , $i = 1, \dots, M-2$, one after the other until the start point is reached. During buffering n_i^E , we consider sizing solutions of the driven gate, g_{i+1}^E , to generate new solutions. We also calculate solutions of sizing g_i^E based on buffering solutions of the driven net, n_{M-1}^E .

The DCP algorithm starts from the buffering to n_{M-1}^E and stops at the buffering to n_1^E . Sizing operations to g_1^E and g_M^E are not considered because they will influence the timing of non-ECO paths. Among the buffering solutions stored in n_1^E , we choose the one that makes the ECO path meet the timing constraints and uses the fewest buffers as the final solution of the ECO path. We rewire the netlist according to the operations of the final solution and run STA to update the circuit timing.

3) *Bounding Box for Choosing Spare Cells*: When buffering a net and sizing to a gate, we need to use spare cells as resources. Since the amount of spare cells is large, exhaustive search for every possible assignment is not efficient. We propose a heuristic to greatly reduce the number of assignments during gate sizing and buffer insertion.

For the case of buffer insertion, we consider inserting a buffer into the net n_i^E along the ECO path with a source g_i^E and sinks G_j^E . We generate a square bounding box for selection. The box is centered at g_i^E and the width of the box is defined as below:

$$width = 2 * \left(\sum_{g_j^E \in \text{fanouts of } g_i^E} \left(distance(g_i^E, g_j^E) + \frac{C_{Ij}^E}{\phi} \right) \right), \quad (3)$$

where C_{Ij}^E is the input pin capacitance of gate g_j^E .

We choose buffer-type spare cells in the bounding box as candidate buffer locations for net n_i^E . Buffer-type spare cells outside the box is not considered because they can not improve the delay of the net. The reason is described below.

We assume that the effect of input transition time to the gate delay is much less than that of the output loading and can be neglected. Then if a buffer-type spare cell g_i^S outside the bounding box is inserted into n_i^E , g_i^E must drive a larger loading than the case with no buffering. Thus gate delay of g_i^E increases because of its larger loading. Since gate delay of g_i^S is larger than zero, the sub-path delay ($delay'(g_i^E) + delay'(g_i^S)$) is larger than the original case that no buffer is inserted. This means the buffer insertion does not help the path delay. Based on the discussion, we have the following theorem.

Theorem 1: Suppose that the gate delay is independent of the input transition time. Given a net n_i^E with source g_i^E and sinks G_j^E to be buffered, let the bounding box for buffering n_i^E be square and centered at g_i^E with $width = 2 * (\sum_{g_j^E \in \text{fanouts of } g_i^E} (distance(g_i^E, g_j^E) + \frac{C_{Ij}^E}{\phi}))$, where C_{Ij}^E is the input pin capacitance of gate g_j^E . Then inserting any buffer-type spare cell outside the bounding box into the net increases the path delay.

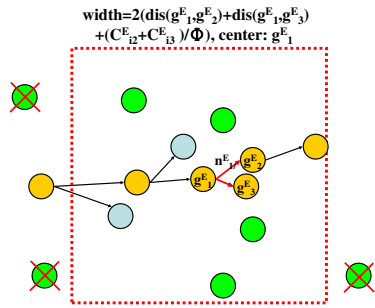


Fig. 11. The bounding box for n_1^E to reduce buffer assignment.

For the case of gate sizing, we have a similar conclusion for reducing spare cell assignment. Figure 12 shows an example of bounding polygon.

Theorem 2: Suppose that the gate delay is independent of the input transition time, and the driving capabilities of the sized gate and the sizing spare cell are the same. Given a gate g_i^E with fanins G_j^E and fanouts G_k^E to be sized, let the bounding polygon be the union of a set of square bounding boxes of G_j^E and G_k^E . A bounding box of g_j^E ($g_j^E \in G_j^E$) is centered at

g_j^E with $width = 2 * (\frac{C_{Ii}^E}{\phi} + distance(g_i^E, g_j^E))$, where C_{Ii}^E is the input pin capacitance of gate g_i^E . A bounding box of g_k^E ($g_k^E \in G_k^E$) is centered at g_k^E with $width = 2 * (\frac{C_{Oj}^E}{\phi} + \sum_{g_k^E \in \text{fanouts of } g_j^E} distance(g_i^E, g_k^E))$,

where C_{Oj}^E is the output pin capacitance of gate g_j^E . Then sizing g_i^E with any spare cell with the same type as g_i^E and located outside the bounding polygon increases the path delay.

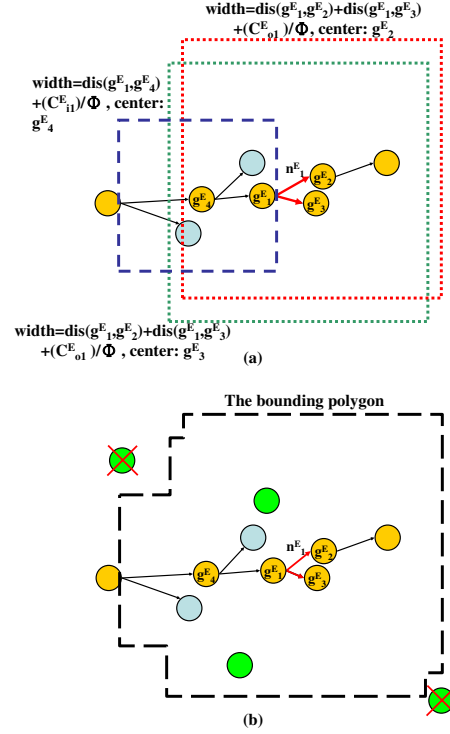


Fig. 12. (a) Bounding boxes of fanins and fanouts. (b) The union of bounding boxes. Spare cells outside the polygon are unconsidered to size g_i^E .

We compare the results of using the two theorems with that of not using in section IV.

4) *Solution Control*: Although we carefully delete many redundant spare cells assignments during DCP by the bounding box theorem and the bounding polygon theorem, number of feasible spare cells assignments is still too large. To speed up the algorithm, we propose two heuristics to control the number of solutions during DCP.

When performing sizing/buffering operation, we generate a set of solutions based on solutions of previous net/gate. The generated solutions can be pruned by the two criteria:

- 1) The number of used buffers.
- 2) The delay of the scope.

Number of sized gates is not counted in the first criterion because gate sizing operation changes a cell in G^E with a spare cell in G^S and does not reduce the number of available spare cells. Thus we prefer to use gate sizing operations to fix timing rather than buffer insertion operations. If a solution s_i uses more buffers than another solution s_j but delay of s_i is larger than s_j , we can delete s_i because it is dominated by s_j .

After deleting dominated solutions, left solutions are grouped into classes according to the number of used buffers. We keep at most $\frac{k}{\# \text{ classes}}$ solutions for each class, where k is a user-defined parameter and can be modified to trade solution quality with runtime.

It is important that both heuristics can not guarantee that the optimum solution will never be deleted, but in general we can delete a lot of sub-optimal solutions and keep the optimum one. Figure 13 illustrates the pruning idea.

D. Technology Remapping

After DCP, we fix timing violations by technology remapping method. This method is to deal with the case that cannot be solved by gate sizing and buffer insertion. For the example shown in Figure 14, there is an AND gate g_1^E driving a large loading but there is no BUFFER and AND-type spare cells near g_1^E . We can use a NAND-type spare cell g_1^S and an INVERTER-type spare cell g_2^S to replace g_1^E and separate the loading.

The placement driven technology mapping methods in Section I-B first place the base gates. Then they map the base gates according to the coordinates of the initial placement. Similar to the methods above, we

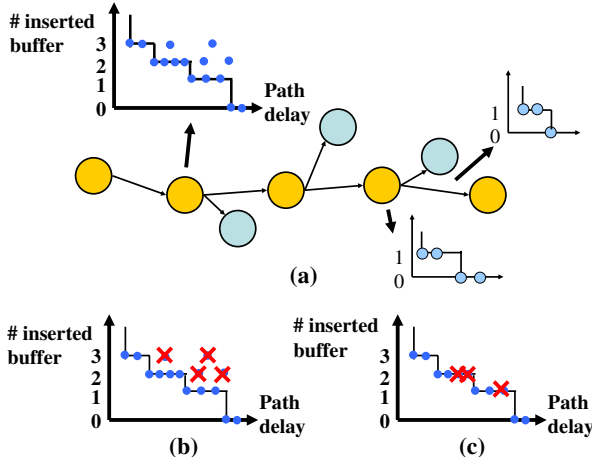


Fig. 13. (a) A set of solutions. (b) Delete the dominated solutions. (c) Keep at most k solutions for every solution class. ($k=8$)

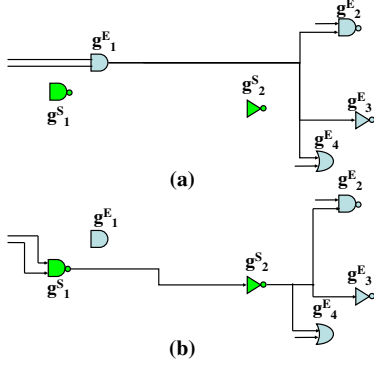


Fig. 14. (a) An AND gate driving a large loading. (b) Map the AND gate to a NAND gate and an INVERTER gate.

first calculate ideal coordinates of the base gates, and map them using this information. The remapping method has following four steps:

- 1) Identify critical parts of the netlist and extract them from the netlist. We denote the extracted gates as G^M .
- 2) Decompose gates in G^M into base gates G^B (NANDs and INVERTERS) by MVSIS [22].
- 3) Calculate ideal locations of base gates of G^B .
- 4) Map G^B . The mapping cost is related to their ideal locations.

We will detail the methods of calculating ideal locations and technology mapping in the following subsections.

1) *Ideal locations:* We know from [1] that optimal buffering to a line is to insert buffers with equal distance, and the distance is $\sqrt{\frac{2R_b C_b}{RC}}$. If we want to map a path that locations of the input pins and output pins are known and fixed, it is intuitive that we map the gates along the path in a way that they are evenly located between input pins and output pins. Since the wire delay is proportional to square of the wirelength, distributing wire loading evenly between the gates reduces the total delay along the path. Furthermore, the buffering after mapping is easier because no gate drives a large loading. Figure 15 shows this concept.

Given a part of netlist to be mapped and locations of input pins and output pins, we calculate the ideal mapping locations of base gates by:

- 1) For every paths from one input pin to one output pin, calculate the candidate locations of base gates along the path as equal distance between the input pin and the output pin.
- 2) If a base gate has more than one candidate locations, average these values to get the final ideal location of the base gate.

An example of calculating ideal locations is shown in Figure 16.

2) *Mapping:* Our mapping algorithm uses dynamic programming method [8]. After decomposing the extracted netlist and calculating ideal locations of the base gates, we cut the network into a forest. Then we map each tree by the following cost function:

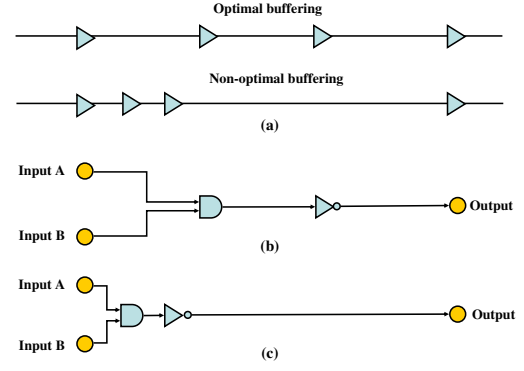


Fig. 15. (a) Optimal buffering and Non-optimal buffering. (b) Distribute the gates between input pins and output pins evenly. (c) Gates are not placed evenly. The inverter needs to drive a large loading.

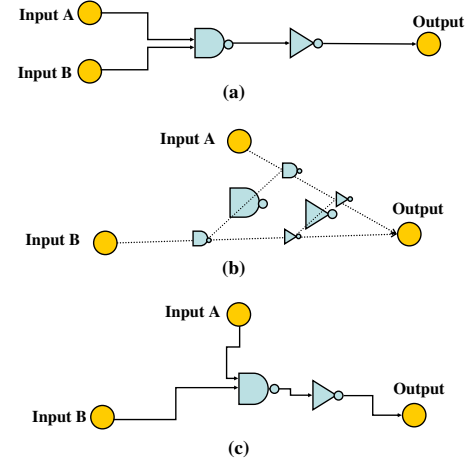


Fig. 16. (a) Subject graph of the netlist without location information. (b) When locations of inputs and outputs are known, calculate the ideal locations of the base gates. (c) The resulted placement if the base gates are placed at their ideal locations.

$$Cost(g_i) = \sum_{g_j \in \text{fanins of } g_i} (Cost(g_j) - d(g_i, g_j)) + d(g_i, \text{fanout of } g_i) \quad (4)$$

where $d(g_i, g_j)$ denotes the distance between g_i and g_j .

The locations of already mapped gates are actual locations of spare cells, while locations of the unmapped base gates are their ideal locations. Thus we can perform technology mapping with the placement information.

After remapping a part of the netlist, we apply STA to update the circuit timing. We apply the remapping process until no timing violations left. An example of technology remapping is shown in Figure 17.

E. Time Complexity Analysis

In this section we analyze the timing complexity of phase 1 of our spare cells selection algorithm.

There are P ECO paths of the netlist. Total gate count is $|G|$, and the number of spare cells is $|G^S|$. If we can finish phase 1 in L iterations, and we keep at most k solutions during each sizing and buffering operation, then the timing complexity of each sizing and buffering operation is $O(k|G^S|)$. If an ECO path has at most M gates, the complexity of the DCP algorithm is $O(kM|G^S|)$. We apply DCP and STA once in an iteration, and complexity of STA is $O(|G|)$. Hence the timing complexity of phase 1 is $O((kM|G^S| + |G|)L)$.

F. Summary

We propose an algorithm which consists of dynamic cost programming and technology remapping to optimize the circuit timing. The whole spare

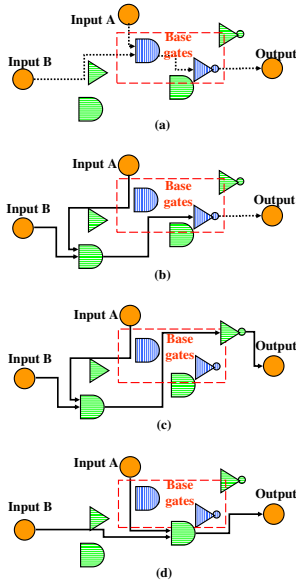


Fig. 17. (a) The base gates located at ideal locations and neighbor spare cells. (b) A match of the NAND base gate. (c) A match of the NAND base gate and the INV base gate. (d) Another match of the NAND base gate and the INV base gate.

cells selection algorithm is illustrated in Figure 18.

IV. EXPERIMENTAL RESULTS

We implemented our algorithm in the C++ programming language on a 3.2GHz Linux workstation with 3 GB memory. The benchmark circuits Case 1, Case 2, Case 3, Case 4, and Case 5, are real industry designs.

In Table I, “Case name” denotes the names of circuits, “Gate count” denotes the number of gates of the circuit, “# Spare cells” denotes number of spare cells in the circuit, “# ECO path” denotes the number of timing violated paths in the unoptimized circuit, “TNS” denotes the total negative slack, and “Max # gates” denotes the largest number of gates along the path among all ECO paths.

Case name	Gate count	# Spare cells	# ECO path	TNS (ns)	Max # gates
Case 1	28927	860	16	9.8	164
Case 2	200504	860	80	312	178
Case 3	91107	860	27	319	173
Case 4	18932	860	22	70	85
Case 5	38011	8600	137	161	72

TABLE I
STATISTICS OF THE TEST CASES.

The experimental results are shown in Table II. We compare three versions of our algorithm: (1) Dynamic cost programming using the bounding box/polygon theorem. (2) Dynamic cost programming and technology remapping using the bounding box/polygon theorem. (3) Dynamic cost programming and technology remapping without the bounding box/polygon theorem. We also compare our algorithm with a greedy wire-cost heuristic algorithm in this table. This heuristic algorithm performs buffer insertion and gate sizing to nets/gates on the ECO paths to reduce the wirelength.

We report the number of left ECO paths, the total negative slack after optimization, the CPU time, and the memory usage. The timing of the circuit is checked by PrimeTime. The third version has similar TNS as the second version and we do not report it. It is seen that using the bounding box/polygon theorem can greatly reduce the runtime while the solution quality is preserved. Besides, technology remapping can further improve the timing but it costs more CPU time than DCP. Since the technology remapping method tends to fix timing deficiencies of misplaced gates by reducing the wirelength, it improves only a little to Case 4 when most misplaced gates are fixed by DCP in the earlier stage. From the comparison

Algorithm: Spare Cells Selection(G, G^S, N)

N : set of all nets of the netlist;
 G : set of all cells;
 P : set of all ECO paths;
 L : denied list;
 G^S : set of all spare cells;
 N^E : set of all nets on the ECO paths;
 G^E : set of all cells on the ECO paths;

```

1 begin
2   apply STA to identify all ECO paths  $P$  and  $G^E, N^E$ 
3   While  $P$  is not empty (phase 1)
4      $p_i$  = the most critical path in  $P$ ;
5     apply DCP to  $p_i$ ;
6     apply STA to update all ECO paths  $P$  and the circuit timing;
7     if path delay of  $p_i \leq$  clock cycle
8       delete  $p_i$  from  $P$ ;
9     if path delay of  $p_i$  is improved
10      move all paths in  $L$  to  $P$  and clear  $L$ ;
11    else
12      move  $p_i$  from  $P$  to  $L$ ;
13  identify critical parts of the paths in  $L$  and remap them (phase 2);
14  apply STA to update the circuit timing;
15 end

```

Fig. 18. Overview of the spare cells selection algorithm.

with the greedy heuristic algorithm, our algorithm is much more effective with a reasonable runtime cost.

We plot the chip layouts of Case 2 for better visualization. Figure 19 (a) shows all timing violated paths (ECO paths) in Case 2, while spare cells are plotted as points. Several gates on the ECO paths are misplaced and cause long wires. This is identical to our analysis that a large wire loading results in a large gate delay. The ECO paths after optimization by our algorithm are shown in Figure 19 (b). Those paths are more compact because misplaced gates on the paths are rewired by spare cells. Additionally, since we prefer gate sizing operations to buffer insertion operations when optimizing timing, the number of gate sizing operations are larger than that of buffer insertion operations.

In Chapter III-E, the timing complexity of the DCP algorithm is $O((kM|G^S| + |G|)L)$, where the term $kM|G^S|$ is dominated by $|G|$ if the number of spare cells and the number of gates along ECO paths are much smaller than the gate count. Since STA is applied to the whole netlist in every iteration, the runtime is proportional to the gate count. By the CPU time in Table II, results of Case 1, Case 2, and Case 3 confirm our analysis. Case 4 and Case 5 are special because the ECO paths of Case 4 are seriously misplaced, and we need a longer time to fix them. The number of spare cells of Case 5 is much larger than other cases and makes the term $kM|G^S|$ dominate the term $|G|$.

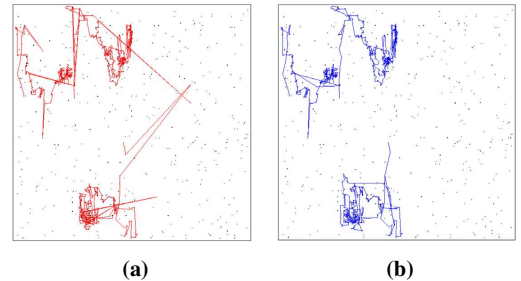


Fig. 19. (a) ECO paths of Case 2 before optimization. (b) ECO paths of Case 2 after optimization.

Based on the comparison with the heuristic algorithm, experimental results demonstrate the effectiveness of our spare cells selection algorithm for the ECO timing optimization problem.

	(1) DCP					(2) DCP+Tech. remapping					(3) DCP (no bounding box/polygon)+Tech. remapping	Greedy wire-cost heuristic				
Case name	# ECO paths	TNS (ns)	Run time (s)	Memory (MB)	Comp. rate	# ECO paths	TNS (ns)	Run time (s)	Memory (MB)	Comp. rate	Run time (s)	# ECO paths	TNS (ns)	Run time (s)	Memory (MB)	Comp. rate
Case 1	0	0	7.7	36	100%	0	0	7.7	36	100%	14.6	6	0.76	2.8	27	92.20%
Case 2	0	0	37	177	100%	0	0	37	177	100%	131	0	0	9.6	174	100%
Case 3	0	0	17.4	78	100%	0	0	17.4	78	100%	502	14	26.96	4.2	78	91.55%
Case 4	3	0.57	30	33	99.20%	3	0.48	75	93	99.31%	180.8	8	10.71	1.5	18	84.59%
Case 5	0	0	2125.5	84	100%	0	0	2125.5	84	100%	12816.8	18	3.34	4.3	33	97.93%
avg			0.88x		99.84%			1.00		99.86%	8.55x			0.18x		93.25%

TABLE II
COMPARISON OF THE THREE VERSIONS OF OUR TOOL WITH A HEURISTIC ALGORITHM.

V. CONCLUSION

In this paper, we propose a spare cells selection algorithm to fix circuit timing after placement. The algorithm consists of gate sizing, buffer insertion, and technology remapping operations. The gate sizing and buffer insertion operations change the netlist gate by gate, while the technology remapping fix timing violations in a more global view. Experimental results show that our algorithm can fix almost timing violations in a much shorter time than the manual method.

We leave the gates not on ECO paths unchanged in our method to reduce the problem size. A more general way of timing optimization using spare cells is to replace gates on non-ECO paths with spare cells and use those replaced gates as spare cells to optimize ECO paths. We will extend our work in this way later.

Functional change is a more difficult work than the timing optimization. Since the changed functions may be very complex and timing issues of the changed netlist also need to be considered, we have to carefully cope with logic and physical co-optimization. This is the direction we plan to advance our research in the future.

REFERENCES

- [1] C. J. Alpert, J. Hu, S. S. Sapatnekar, and C. N. Sze, "Accurate Estimation of Global Buffer Delay within a Floorplan," in *Proceeding of International Conference on Computer Aided Design*, pp. 1140-1146, 2004.
- [2] K. Chaudhary and M. Pedram, "A Near Optimal Algorithm for Technology Mapping Minimizing Area under Delay Constraints," in *Proceeding of Design Automation Conference*, pp. 492-498, 1992.
- [3] J. Cong and Y. Ding, "An Optimal Technology Mapping Algorithm for Delay Optimization in Lookup-table Based FPGA Designs," in *Proceeding of International Conference on Computer Aided Design*, pp. 48-53, 1992.
- [4] Faraday Technology Corporation, <http://www.faraday-tech.com/index.html>.
- [5] L. P. P. van Ginneken, "Buffer Placement in Distributed RC-tree Networks for Minimal Elmore Delay," in *Proceeding of International Symposium on Circuits and Systems*, pp. 865-868, 1990.
- [6] Z. Li, W. Shi, "An $O(mn)$ Time Algorithm for Optimal Buffer Insertion of Nets With m Sinks," in *Proceeding of Asia and South Pacific Design Automation Conference*, pp. 320-325, 2006.
- [7] D.-J. Jongeneel, Y. Watanabe, R. K. Brayton, and R. Otten, "Area and Search Space Control for Technology Mapping," in *Proceeding of Design Automation Conference*, pp. 86-91, 2000.
- [8] K. Keutzer, "DAGON: Technology Binding and Local Optimization by DAG Matching," in *Proceeding of Design Automation Conference*, pp. 617-623, 1987.
- [9] Y. Kukimoto, R. K. Brayton, and P. Sawkar, "Delay-optimal Technology Mapping by DAG Covering," in *Proceeding of Design Automation Conference*, pp. 348-351, 1998.
- [10] T. Kutzschebauch and L. Stok, "Congestion Aware Layout Driven Logic Synthesis," in *Proceeding of International Conference on Computer Aided Design*, pp. 216-223, 2001.
- [11] T. Kutzschebauch and L. Stok, "Layout Driven Decomposition with Congestion Consideration," in *Proceeding of Design Automation and Test in Europe*, pp. 672-676, 2002.
- [12] E. Lehman, Y. Watanabe, J. Grodstein, and H. Harkness, "Logic Decomposition during Technology Mapping," in *Proceeding of International Conference on Computer Aided Design*, pp. 264-271, 1995.
- [13] I.-M. Liu, A. Aziz, D.F. Wong, and H. Zhou, "An Efficient Buffer Insertion Algorithm for Large Networks Based on Lagrangian Relaxation," in *Proceeding of International Conference on Computer Design*, pp. 614-621, 1999.
- [14] I.-M. Liu, A. Aziz, and D.F. Wong, "Meeting Delay Constraints in DSM by Minimal Repeater Insertion," in *Proceeding of Design Automation and Test in Europe*, pp. 436-440, 2000.
- [15] Q. Liu and M. Marek-Sadowska, "Pre-layout Wire Length and Congestion Estimation," in *Proceeding of Design Automation Conference*, pp. 582-587, 2004.
- [16] Q. Liu and M. Marek-Sadowska, "Technology Mapping: Wire Length Prediction-based Technology Mapping and Fanout Optimization," in *Proceeding of International Symposium on Physical Design*, pp. 145-151, 2005.
- [17] J. Lou, W. Chen, and M. Pedram, "Concurrent Logic Restructuring and Placement for Timing Closure," in *Proceeding of International Conference on Computer Aided Design*, pp. 31-36, 1999.
- [18] A. Lu, G. Stenz, and F. M. Johannes, "Technology Mapping for Minimizing Gate and Routing Area," in *Proceeding of Design Automation and Test in Europe*, pp. 664-669, 1998.
- [19] Y. Matsunaga, "On Accelerating Pattern Matching for Technology Mapping," in *Proceeding of International Conference on Computer Aided Design*, pp. 118-122, 1998.
- [20] A. Mishchenko, X. Wang, and T. Kam, "A New Enhanced Constructive Decomposition and Mapping Algorithm," in *Proceeding of Design Automation Conference*, pp. 143-148, 2003.
- [21] M. Murofushi, T. Ishioka, M. Murakata and T. Mitsuhashi, "Layout Driven Re-synthesis for Low Power Consumption LSIs," in *Proceeding of Design Automation Conference*, pp. 666-669, 1997.
- [22] MVSIS: Logic Synthesis and Verification, <http://embedded.eecs.berkeley.edu/Respep/Research/mvsis>.
- [23] D. Pandini, L. T. Pileggi, and A. J. Strojwas, "Understanding and Addressing the Impact of Wiring Congestion during Technology Mapping," in *Proceeding of International Symposium on Physical Design*, pp. 131-136, 2002.
- [24] D. Pandini, L. Pileggi, and A. Strojwas, "Congestion-aware Logic Synthesis," in *Proceeding of Design Automation and Test in Europe*, pp. 664-671, 2002.
- [25] M. Pedram and N. Bhat, "Layout Driven Technology Mapping," in *Proceeding of Design Automation Conference*, pp. 99-105, 1991.
- [26] R. S. Shelar, P. Saxena, X. Wang, and S. S. Sapatnekar, "Technology Mapping: An Efficient Technology Mapping Algorithm Targeting Routing Congestion under Delay Constraints," in *Proceeding of International Symposium on Physical Design*, pp. 137-144, 2005.
- [27] W. Shi and Z. Li, "An $O(n \log n)$ Time Algorithm for Optimal Buffer Insertion," in *Proceeding of Design Automation Conference*, pp. 580-585, 2003.
- [28] C. N. Sze, C. J. Alpert, J. Hu, and W. Shi, "Path Based Buffer Insertion," in *Proceeding of Design Automation Conference*, pp. 509-514, 2005.