

# An Efficient Technology Mapping Algorithm Targeting Routing Congestion under Delay Constraints\*

Rupesh S. Shelar  
Desktop Products Group  
Intel Corporation  
Hillsboro, OR

rupesh.s.shelar@intel.com

Xinning Wang  
Strategic CAD Labs.  
Intel Corporation  
Hillsboro, OR

xinning.wang@intel.com

Prashant Saxena  
Advanced Technology Group  
Synopsys Inc.  
Hillsboro, OR

saxena@synopsys.com

Sachin S. Sapatnekar  
Department of ECE  
University of Minnesota  
Minneapolis, MN

sachin@ece.umn.edu

## ABSTRACT

Routing congestion has become a serious concern in today's VLSI designs. In this paper, we propose a technology mapping algorithm that minimizes routing congestion under delay constraints. The algorithm employs a dynamic programming framework in the matching phase to generate probabilistic congestion maps for all the matches. These congestion maps are then utilized to minimize routing congestion during the covering, which preserves the delay-optimality of the solution using the notion of slack. Experimental results on benchmark circuits in a 100 nm technology show that the algorithm can improve track overflows by 44%, on an average, as compared to the conventional technology mapping while satisfying delay constraints.

## Categories and Subject Descriptors

B.6.3 [Design Aids]: Automatic synthesis; Optimization

## General Terms

Algorithms, Design

## Keywords

Routing congestion, technology mapping, logic synthesis

## 1. INTRODUCTION

Following Moore's law [1], the number of on-chip transistors are doubling every two years, while the number of wires are growing

\*Partially, this work was performed, when Rupesh Shelar and Prashant Saxena were with the University of Minnesota and Strategic CAD Labs., Intel, respectively.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ISPD'05, April 3–6, 2005, San Francisco, California, USA.  
Copyright 2005 ACM 1-59593-021-3/05/0005 ...\$5.00.

almost linearly with the number of gates. This increasing design complexity results in circuits that face the problem of routing congestion, which can be described as the unavailability of a sufficient number of tracks to route wires. Moreover, wires are becoming increasingly resistive with each technology generation in spite of the advances in manufacturing techniques [2, 3], and therefore, interconnect delays have been seen to dominate gate delays since the last couple of technology generations. Together, routing congestion and the dominance of interconnect delays make timing closure difficult: if wires are detoured to avoid congested regions, they may incur larger delays and thus violate timing constraints.

The flexibilities offered by the placement and routing stages to alleviate congestion are often insufficient, and addressing the problem of routing congestion only during these stages is known to result in a number of design iterations. On the other hand, logic synthesis offers a large degree of freedom in handling the routing congestion problem, but it may suffer from inaccurate estimates as it operates at a higher level of abstraction than the placement and routing stages. In the synthesis domain, technology mapping is a powerful transformation which makes decisions about wires and therefore, affects congestion. Consequently, it would be an excellent stage during which one could try to alleviate congestion, provided it were possible to obtain reasonably accurate congestion estimates.

## 1.1 Previous Work

Several technology mapping algorithms that target traditional objectives such as area, delay, or power exist in the literature [4–6]. Recently, there have been attempts to consider congestion during the mapping. These approaches include the following: placement driven mapping for FPGA's [7], methods employing a cost function that involved wirelength as a metric for routing congestion [8–10], a procedure based on predictive probabilistic congestion estimates [11], and a method based on pre-layout wirelength prediction [12]. In [7], Cong *et al.* present an iterative congestion-aware mapping and placement procedure for FPGA's; however, in this work, the congestion metric used by them pertains not to routing but to cells, being defined as the number of cells placed in a given location. The approaches due to Pandini *et al.* [8,9] and Stok *et al.* [10] rely on the total wirelength, which, being a global metric, fails to capture the locality property of the routing congestion. The work in [11] em-

loys predictive probabilistic congestion estimates, and therefore, suffers from the inaccuracies inherent in any predictive scheme. Furthermore, it focuses solely on the circuit area. Other related work lies in the domain of structural logic synthesis [13, 14], where metrics for routing congestion are proposed to guide the logic synthesis. An *adhesion* is presented as a metric for routing congestion in [13], while structural pin density is shown to correlate well with congestion in [14]. The adhesion, being computationally expensive, may not be suitable for technology mapping purposes. The structural pin density, on the other hand, ignores congestion contribution of wires passing over a given region and therefore, may not be accurate.

## 1.2 Our Contributions

Considering routing congestion during the mapping is more complex than traditional objectives such as area or delay due to the following reasons.

- Unlike conventional objectives, routing congestion, being locality dependent, cannot be captured using a single number at the technology mapping stage [9].
- Even with the application of a probabilistic congestion map, there is a “chicken-and-egg” problem between mapping and placement stages, since such a congestion map is required before mapping, but cannot be created until after the placement of a mapped netlist.

To overcome this “chicken-and-egg” problem, previous approaches have either used predictive congestion maps, as in [11] or have employed metrics such as wirelength or mutual contraction, as in [9, 10]. The limitation of the former approach has been the reliance on empirical data both to justify the heuristic objective function driving the mapping and to predict congestion maps for mapped netlists based solely on subject graphs, while the latter approaches attempt minimizing wirelength, assuming that it correlates well with congestion. In contrast, our current work provides a sound theoretical basis for the mapping procedure that guarantees optimal delays as well as allows the use of accurate congestion maps that are created as the mapping proceeds. The contributions of this work are as follows:

1. We formulate the technology mapping problem targeting routing congestion as that of minimizing the total track overflow under the specified delay constraints. Exploiting the dynamic programming framework, we provide a delay-optimal solution to the problem under the assumption that the placement assigned to the cells during the mapping is preserved.
2. Instead of predicting congestion from a generic netlist, such as a subject graph, and justifying its use empirically to overcome the cyclic dependence between the mapping and placement stages, we propose a matching procedure to generate two-dimensional congestion maps for all delay-optimal mapping solutions in a bottom-up manner. The procedure is general enough and can be applied not only to optimize different cost functions, such as maximum congestion or track overflow, defined over the congestion map, but also to optimize other physical properties that can be accurately captured using two-dimensional maps, for instance, temperature or power density maps.
3. In the covering phase, where the matches are selected from among the stored choices, we employ an explicit notion of the slack to further optimize the design unlike the classical covering approach [15, 16], which does not explore this

potential. Our covering technique chooses the congestion-optimal matches that minimize the total track overflow and also satisfy the slack constraints. This technique can be easily extended to optimize even traditional objectives, such area or power under delay constraints without introducing any sub-optimality in delays. Experimental results on an entire ISCAS’85 benchmark suite confirm that delay constraints are always satisfied while improving track overflow by 44%, on an average.

The rest of the paper is organized as follows. Section 2 introduces formal definitions and the background for the technology mapping problem, while Section 3 describes the generation of congestion maps during the matching phase. Section 4 illustrates the slack-constrained congestion-aware covering algorithm and Section 5 discusses the extensions to the algorithm. Section 6 presents experimental results followed by the conclusion in Section 7.

## 2. PRELIMINARIES

The following terminology is used in this paper. A Boolean network is a directed acyclic graph (DAG), in which a node denotes a Boolean function,  $f : B^n \rightarrow B$ , where  $B = \{0, 1\}$ , and  $n$  is the number of inputs to the node. Traditional technology mapping is usually preceded by a decomposition of this abstract network into one that contains primitive gates, such as 2-input NAND’s and inverters. The decomposed network is referred to as a subject graph or a premapped netlist. The subject graph is mapped on to a set of cells in the library during technology mapping; the resulting network is known as a mapped netlist, which is placed in a given block area and routed. The block area is divided into bins for congestion analysis purposes or for global routing. Each bin contains a limited number of horizontal and vertical tracks. The track overflow and congestion can be defined for every bin as follows.

**DEFINITION 2.1.** *The horizontal (vertical) track overflow for a given bin, ( $T_{h(v)}^{bin}$ ), is defined as the difference between the number of horizontal (vertical) tracks required to route the nets through the bin and the available number of horizontal (vertical) tracks.*

**DEFINITION 2.2.** *The horizontal (vertical) congestion for a given bin,  $C_{h(v)}^{bin}$ , is the ratio of number of horizontal (vertical) tracks required to route the nets through the bin to the number of horizontal (vertical) tracks available.*

In this paper, when we use the terms “track overflow” or “congestion” without specifying a horizontal or vertical direction, we mean that the terms are equally applicable to both horizontal and vertical directions.

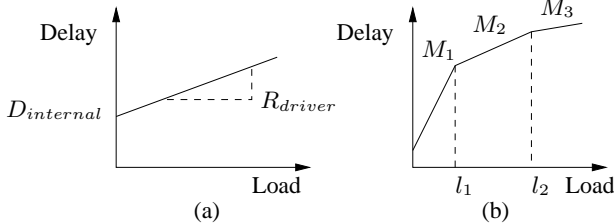
A positive track overflow or a congestion of more than 1.0 means that sufficient tracks are unavailable for the routing, while a negative value of the overflow or a congestion smaller than 1.0 indicates the availability of tracks. The total track overflow ( $OF$ ) is the sum of positive track overflows over all the bins, as shown in the following equation

$$OF = \sum_{\forall bins: C^{bin} > 1.0} T^{bin} \quad (1)$$

This overflow can be computed after the generation of the congestion map, which can be derived either using probabilistic techniques or by performing routing. Employing these definitions, the problem of delay oriented technology mapping targeting congestion can be defined as follows.

**PROBLEM DEFINITION 1.** *Given a subject graph of a network and a library of gates, generate a mapped netlist that minimizes the total track overflow under specified delay constraints.*

Traditional mapping procedures use a dynamic programming framework that involves two phases, referred to as *matching* and *covering*: in the former, non-inferior mapping solutions are stored during a topological traversal of a circuit, while, in the latter, a mapped network is built by selecting from these solutions during a reverse topological traversal. Usually, technology mapping employs one of the following two classes of delay models: load- or gain-based. In this paper, we consider only load-based delay models, as our algorithm can be easily extended to the one based on the latter. The load-based delay model is shown in Figure 1(a) for a typical standard cell: it shows a straight line with the internal delay of the gate,  $D_{internal}$ , as an intercept on the delay axis, while the slope of the line indicates the effective driver resistance<sup>1</sup>. Technology mapping targeting delay involves storing piece-wise linear load-delay curves,  $\{(l_1, D_1), (l_2, D_2), \dots\}$ , during the matching phase, where  $l_i$  and  $D_i$ , respectively, denote the load and delay coordinates of an end-point of a piece-wise linear segment. At each node, a set of choices that are delay-optimal for the load ranges corresponding to piece-wise linear segments is stored on these curves. These choices are referred to as non-inferior matches. One such curve is shown in Figure 1(b) with three non-inferior matches  $M_1$ ,  $M_2$ , and  $M_3$ , where  $M_1$  is optimal for the load range  $[0, l_1]$ ,  $M_2$  for the range  $(l_1, l_2]$ , and  $M_3$  for larger load values. During the covering phase, when loads are known, delay-optimal matches are chosen from the curves. SIS [16] contains an implementation of a delay oriented mapper based on this scheme, and we employ the same framework for our technology mapping targeting routing congestion.

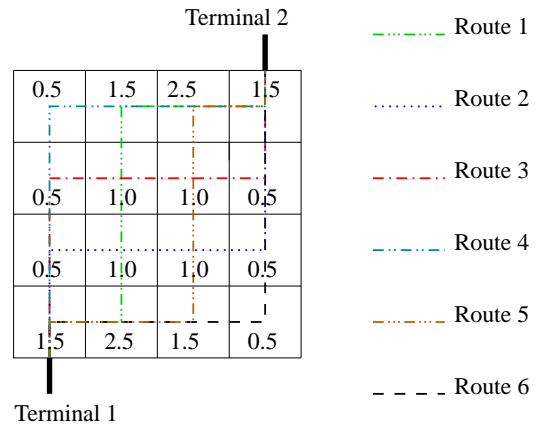


**Figure 1: (a) A load-based delay model for a typical standard cell, such as an inverter. (b) A typical load-delay curve stored during the matching.**

The concept of employing a companion placement for the subject graph to estimate wirelengths or congestion maps is not new. It has been used by previous technology mapping or physical synthesis methods [9–11, 15, 18], and we employ the same concept here. Based on such a companion placement, the congestion maps are generated during the matching phase employing a well known probabilistic method, which is shown to have a good fidelity with post-routing congestion in [17, 19]. The probabilistic estimation technique assumes all routes to be equally possible for a given net and then computes the demands for tracks as a ratio of the number of paths passing through the bin to the total number of paths. Figure 2 shows the congestion computation for a net. Six possible paths, assuming only L- and Z-shaped routes<sup>2</sup>, through different bins in the bounding box of a net are shown in the figure. The

<sup>1</sup>The delay of a cell also depends on the slopes of the input signal transitions, which are considered during precise timing analysis, but are often ignored in the delay models at the technology mapping stage.

<sup>2</sup>L- and Z-shaped routes are shown for illustration purposes only. In practice, we allow the nets to have unlimited bends, as in [19].



**Figure 2: Probabilistic congestion estimation for a net assuming only L- and Z-shaped routes [17]. Only the demands for horizontal tracks in each bin are shown.**

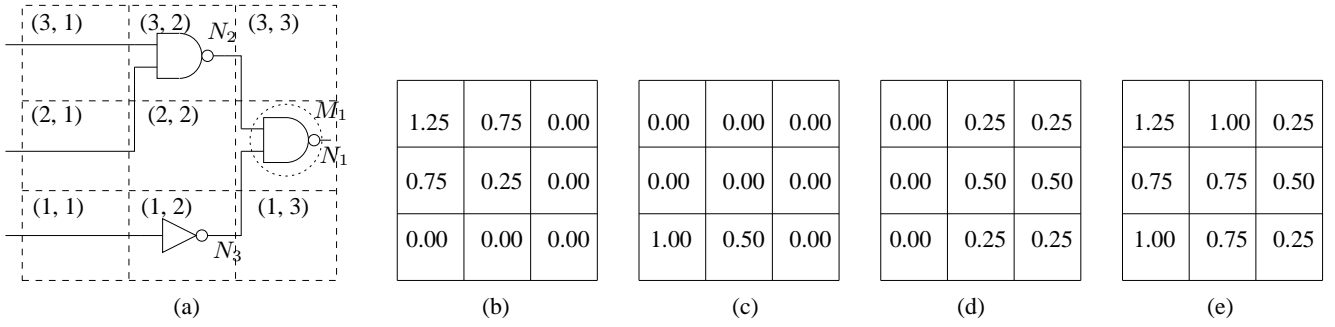
numbers associated with each bin show the demands for the horizontal tracks. For instance, the leftmost bin in the bottom row has a demand of 1.5, since there are three routes, route 1, 5, and 6, which require half track each in that bin. Therefore, the congestion for the bin due to the net is  $\frac{1.5}{6 \times h_{bin}}$ , where,  $h_{bin}$  is the number of horizontal tracks available for the bin.

### 3. CONGESTION MAP GENERATION DURING THE MATCHING PHASE

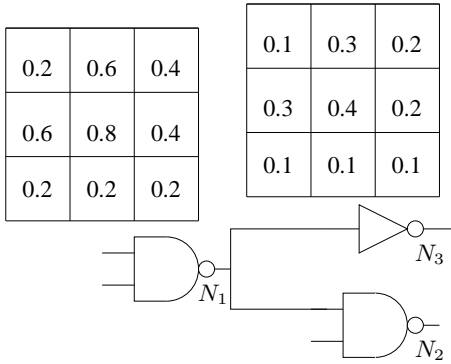
The matching phase of dynamic programming based delay oriented technology mapping typically involves storing a load vs. delay curve at each node. We employ the same method and preserve the curve containing non-inferior matches that minimize the delay for different load values. During the construction of the curve, wire-loads and wire-delays are accounted for based on the companion placement of the underlying subject graph. To evaluate different mapping solutions based on their contribution to the total track overflow, we associate a probabilistic congestion map with each match. This congestion map represents wires due to the mapping solution corresponding to a given match, as explained in the following subsections.

#### 3.1 Bottom-up Congestion Map Construction

To generate a global as well as a partial view of congestion, we propose a bottom-up congestion map construction. A match at a given node is assumed to be placed at the center of gravity of its fanins and fanouts, and multi-terminal nets are modeled employing the well known star topology. Figure 3 shows the creation of a congestion map for a match  $M_1$  at the node  $N_1$ . The match  $M_1$  receives its inputs from two nodes:  $N_2$  and  $N_3$ . During the topological traversal, these nodes are processed before node  $N_1$  and hence, non-inferior matches and the corresponding congestion maps are already stored at these nodes. The maps of horizontal congestion for matches at  $N_2$  and  $N_3$  that minimize the delay for the solution due to the match  $M_1$  are shown in Figure 3(b) and (c), respectively, while the horizontal congestion contribution due to the match  $M_1$  is shown in 3(d). For the purposes of illustration, only the track demands are shown as congestion without loss of generality. In Figure 3(d), the right-most bin in the third row has a horizontal track demand of 0.25, as there are 2 paths from the output of the match at  $N_2$  to an input of  $M_1$ . This results in a probability of  $\frac{1}{2}$  of



**Figure 3: Generating congestion maps during the matching: (a) A choice  $M_1$  at node  $N_1$ . Maps for horizontal congestion for matches at  $N_2$  and  $N_3$  are shown in (b) and (c), respectively, while (d) shows congestion contribution due to  $M_1$  and (e) shows the congestion map at  $N_1$  due to  $M_1$ .**



**Figure 4: The congestion in each bin is divided by the number of fanouts for forward propagation.**

the path through the bin being selected. Furthermore, for that path, only half of one horizontal track is occupied under the assumptions of probabilistic congestion estimation [19], resulting in a demand of 0.25. Figure 3(e) shows the map for horizontal congestion for the solution involving the match  $M_1$ : it is obtained by the simple bin-wise sum of the congestion maps at  $N_1$ ,  $N_2$ , and  $N_3$ . Thus, for instance, the demands in bin (1, 2) in the congestion maps in figures 3(b), (c) and (d) are added to create the demand in bin (1, 2) in the congestion map shown in Figure 3(e). This congestion map represents horizontal track demands due to the subset of wires from the transitive fanin cone of  $N_1$ , as these wires appear in the mapping solution corresponding to the match  $M_1$ . For a different match at node  $N_1$ , this subset of wires will be different leading to a different congestion map. Note that these congestion maps account for the relevant subsets of wires from only the transitive fanin cone of a given node, ignoring wires in the rest of the network. Thus, they do not represent a global picture of congestion; rather, they represent a partial picture that accounts only for the matches and corresponding wires in the transitive fanin cone of a given node. It is clear that a complete picture of congestion that represents all the wires that will appear in the mapped network cannot be obtained until matching and bottom-up congestion map generation is performed for all the nodes. This is the reason why we postpone the total track overflow computation until the covering stage.

### 3.2 Handling Multiple Fanouts

For the multiple fanout points, as shown in Figure 4, the congestion in each bin is divided by the number of fanouts. It allows

the construction of the congestion map for the solutions at primary outputs by simply carrying out the bin-wise sum of the corresponding congestion maps. This heuristic is similar to the one employed in [5] for the area minimization under delay constraints, where in spite of such a division of gate-areas at multiple fanout points, addition of the gate-areas due to points on area-delay curve at primary outputs generates the gate-area for an entire solution.

Thus, employing simple algebraic operations such as addition and multiplication, two-dimensional congestion maps can be created during the matching phase. These maps can be used to optimize any cost function defined on them, for instance, total track overflow or maximum congestion. Moreover, the entire technique to generate congestion maps is quite general and can be applied to create two-dimensional maps for even other physical properties, such as power density.

### 3.3 Analogy with the Classical Matching

The bottom-up congestion map generation is analogous to the work in [5], where area minimization under delay constraints is sought. However, in contrast, our work targets routing congestion under the same constraints, rather than the area as in [5]. The partial congestion maps during the matching in our work correspond to gate-areas for points on the area-delay curve in [5]. Just as the gate-area for a point on the curve models the corresponding gate-areas due to the match at a given node and its transitive fanin cone in [5], the congestion map for a match at a given node accounts for wires due to matches in the transitive fanin cone and represents the corresponding probabilistic congestion of these wires in our work. Just as the gate-area for entire network cannot be predicted until the end of the matching phase in [5], so is the inability of the dynamic programming framework in our work to predict the congestion map of entire network until the end of the matching. Of course, compared to a metric like gate-area, the routing congestion optimization objective is far more complicated, as noted in several recent works [9, 11, 13, 14]. We also employ a better delay computation approach as compared to [5]. In their work, the construction of area-delay curves introduces sub-optimality due to unknown loads in the matching phase. In contrast, we store non-inferior solutions using piece-wise linear load-delay curves, which serves the following purposes: (a) it does not introduce any sub-optimality in delays due to the load computation and therefore, allows one to generate a delay-optimal solution; and (b) because of the ability to generate a delay-optimal solution, it ensures that if the delay-optimal solution cannot meet constraints, no other solution can. Moreover, we utilize the companion placement to account for wire-loads and wire-delays, which are ignored in their work.

### 3.4 Comparison with Competitive Congestion-aware Matching Techniques

The matching procedures employed in previous approaches on congestion-aware mapping either compute the wire-length, as in [8–10], or the congestion cost based on probabilistic estimates, as in [11]. The following are the limitations of these approaches.

1. All of the approaches heuristically modify the area or delay cost function by adding wirelength or congestion penalty terms. This introduces sub-optimality in the estimation of these objective functions, and therefore, none of these approaches can ensure meeting area or delay constraints (under the usual assumptions about delay models and placement fidelity) while minimizing routing congestion.
2. As pointed out earlier, metrics such as wirelength and predicted congestion may not capture an accurate picture of the congestion.

These limitations are overcome easily by our matching procedure, which stores all non-inferior matches along with their probabilistic congestion maps that are constructed using simple algebraic operations such as addition and multiplication in a bottom-up manner. Moreover, although the metrics proposed in previous approaches have been applied during the matching phase, none of them extends these metrics to guide the covering, which actually selects the mapping solution, and thus, these techniques ignore the search space available during the covering. In contrast, we utilize the generated congestion maps during the covering process to compute the total track overflow due to different mapping choices, and to select the one that minimizes this track overflow while satisfying the delay constraints. This is explained in the following section.

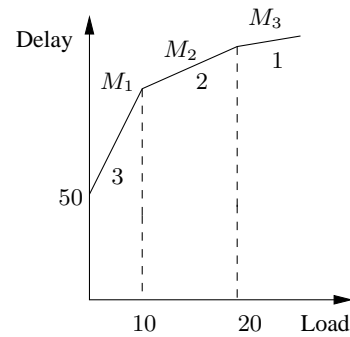
## 4. CONGESTION MINIMIZATION UNDER DELAY CONSTRAINTS DURING THE COVERING PHASE

To preserve the delay-optimality of the solution while improving the congestion, we associate the notion of a slack with all nodes.

**DEFINITION 4.1.** *The slack at a given node is the difference between the required arrival time at that node and the actual arrival time.*

A positive value of slack means that the signal arrives earlier than the required arrival time, while a negative value implies that it arrives later than the required time. During traditional covering, delay-optimal choices that minimize the delay for a given load are chosen at each node. The load-delay curves that are built during the matching phase also assume the same. We observe that the covering need not choose delay-optimal choices at all nodes to respect the delay constraints. At nodes with positive slack, matches that are not necessarily delay-optimal can still be chosen, as long as they meet the delay constraints. Our covering algorithm employs this idea to minimize the total track overflow. We explain the same using the following example.

Consider a load-delay curve, shown in Figure 5, stored during the matching for a node. During the covering, when the node is processed, let us assume that it has a slack of 10 and it has to drive a load of 15. The delays due to matches  $M_1$ ,  $M_2$ , and  $M_3$  for this load are, respectively, 95, 90, and 95. In this case, regular covering will choose match  $M_2$ , since it minimizes the delay, while our congestion-aware covering will choose a match that minimizes the track overflow. Note that choosing  $M_1$  and  $M_3$  does not affect the delay-optimality of the overall solution in this case, since there



**Figure 5:** A piece-wise linear load-delay curve with three matches  $M_1$ ,  $M_2$ , and  $M_3$  optimal for load ranges  $[0, 10)$ ,  $[10, 20)$ , and  $[20, \infty)$ , respectively. The curve has a slope of 3 for the load values  $[0, 10)$ , a slope of 2 for remaining load values  $[10, 20)$ , and the slope of 1 for remaining load values.

is a slack of 10 at the node, and the arrival times at the node due to both,  $M_1$  and  $M_3$ , satisfy this slack constraint.

### 4.1 Algorithm for the Covering

The pseudo-code for the covering that targets track overflow under delay constraints is shown in Algorithm 1. It begins with the computation of the delay-optimal matches at the primary outputs. The slacks ( $S_{o_i}$ ) are then computed for all the outputs. The congestion map  $CM$  for this solution is built by the bin-wise addition of the congestion maps due to delay-optimal matches for all the primary outputs. This congestion map represents the contributions of all the wires that will appear in the mapped network due to the conventional delay-optimal solution. The total track overflow  $OF$  corresponding to this solution is estimated from the congestion map using Equation (1). After this initialization, all the nodes ( $v \in V$ ) are processed in the reverse topological order. First the delay-optimal match  $m_v^{D-optimal}$  is determined for a node, followed by a computation of the slack-constrained congestion-optimal match  $m_v^{C-optimal}$  using a procedure that is explained later. If the overflow  $OF_v^{C-optimal}$  due to the slack-constrained congestion-optimal match  $m_v^{C-optimal}$  is smaller than that due to the delay-optimal match  $m_v^{D-optimal}$ , then  $m_v^{C-optimal}$  is chosen as the optimal match  $m_v^{Optimal}$ . In this case, the congestion map is updated to reflect the change due to this match, the new overflow is stored, and the slacks at the inputs of the match are also updated. Instead, if the delay-optimal match is chosen as the optimal match  $m_v^{Optimal}$ , then the slack is simply propagated to the inputs of the match. Finally, the loads at the inputs of the selected optimal match (be it delay-optimal or slack-constrained congestion-optimal) are incremented to reflect the selection of that match.

### 4.2 Procedure for Finding Slack-constrained Congestion-optimal Match

The pseudo-code for selecting the slack-constrained congestion-optimal match is shown in Algorithm 2. It considers all the matches except the delay-optimal match such that they satisfy the slack constraint (as enforced in line 6 in the pseudo-code). Among these matches, one that results in the lowest total track overflow is stored as the slack-constrained congestion-optimal match. Note that we can store a match that results in the smallest maximum congestion as a congestion-optimal one, leading to the optimization of the maximum congestion. In general, any cost function defined over the congestion maps can be optimized by storing a match that op-

---

**Algorithm 1** Perform the covering targeting congestion minimization under delay constraints

---

**Input:** A Boolean network  $G(V, E)$ , a set of primary outputs  $O \subseteq V$ , sets of non-inferior matches  $M_v$  and their congestion maps  $CM_v$  for  $v \in V$

**Output:** Assignment of congestion-optimal match  $m \in M_v$  to  $v \in V$ , which satisfies delay constraint

- 1: **for**  $\forall o_i \in O$  **do**
- 2:  $m_{o_i}^{D-optimal} \leftarrow \text{DelayOptimalMatch}(M_{o_i}, \text{load}_{o_i})$
- 3:  $S_{o_i} \leftarrow D_{o_i}^{Required} - D_{m_{o_i}^{D-optimal}}$
- 4: **end for**
- 5:  $CM \leftarrow \sum_{i=1}^{|O|} CM_{m_{o_i}^{D-optimal}}$
- 6:  $OF \leftarrow \text{ComputeOverflow}(CM)$
- 7: **for**  $\forall v \in V$ , in reverse topological order **do**
- 8:  $m_v^{D-optimal} \leftarrow \text{DelayOptimalMatch}(M_v, \text{load}_v)$
- 9:  $m_v^{C-optimal} \leftarrow \text{CongestionOptimalMatch}(M_v, s_v)$
- 10: **if**  $OF_v^{C-optimal} < OF$  **then**
- 11:  $m_v^{Optimal} \leftarrow m_v^{C-optimal}$
- 12:  $CM \leftarrow CM - CM_{m_v^{D-optimal}} + CM_{m_v^{C-optimal}}$
- 13:  $OF \leftarrow OF_v^{C-optimal}$
- 14:  $\text{UpdateSlacks}(m_v^{C-optimal}, s_v^{C-optimal})$
- 15: **else**
- 16:  $m_v^{Optimal} \leftarrow m_v^{D-optimal}$
- 17:  $\text{UpdateSlacks}(m_v^{D-optimal}, s_v)$
- 18: **end if**
- 19:  $\text{IncrementLoads}(m_v^{Optimal})$
- 20: **end for**

---

timizes the given objective as a congestion-optimal match. The corresponding slack updates are maintained with the congestion-optimal match. If the track overflow due to this match is lower than that due to the delay-optimal match, then the slack-constrained congestion-optimal match is stored as the optimal match, as described earlier.

### 4.3 Time Complexity of the Algorithm

The time-complexity of the entire mapping algorithm is almost same as that of a conventional mapping. Both conventional as well as our congestion-aware mapping techniques employ the same matching procedure, except that the matching phase in our case computes and stores congestion maps for non-inferior matches. If there are  $N_{Bins}$  number of bins in the layout, then the matching phase would require  $O(N_{Bins}N_{Matches})$  time, where  $N_{Matches}$  is the number of non-inferior matches over entire network, since the computation of a congestion map requires  $O(N_{Bins})$  time. Since  $N_{Bins}$  is a constant (although possibly large as compared to other constants subsumed by  $O()$ ), the matching phase for our approach requires  $O(N_{Matches})$  time for a given layout. During the covering,  $\text{CongestionOptimalMatch}()$  function is called for all the nodes. The function requires  $O(|M_v|N_{Bins})$  time, since, in the worst case, it has to consider all the matches at the node to find the slack-constrained congestion-optimal one. Over all the nodes, therefore, the covering requires  $O(N_{Matches})$  time, which is same as that of the matching phase.

## 5. EXTENSIONS TO THE ALGORITHM

Our algorithm, which is based on the load-based delay model, can be easily extended to the one based on gain-based delay model. Note that during technology mapping based on the gain-based delay model, the sizes of the cells are adjusted depending on the loads that they drive, and this does not affect the wires and hence, the

---

**Algorithm 2** Find congestion-optimal match that satisfies the slack constraint

---

**Input:** A set of matches  $M_v$  for  $v$  and available slack  $s_v$  at node  $v$

**Output:** Congestion-optimal match satisfying the slack and the corresponding total track overflow  $OF_v^{C-optimal}$

- 1: Procedure  $\text{CongestionOptimalMatch}(M_v, s_v)$  {
- 2:  $OF_v^{C-optimal} \leftarrow \infty$
- 3: **if**  $s_v > 0$  **then**
- 4: **for all**  $m \in M_v - m_v^{D-optimal}$  **do**
- 5:  $D_{m_v} \leftarrow \text{DelayDueToMatch}(\text{load}_v)$
- 6: **if**  $D_{m_v} - D_v^{D-optimal} < s_v$  **then**
- 7:  $CM_{new} \leftarrow CM - CM_{m_v^{D-optimal}} + CM_m$
- 8:  $OF_{new} \leftarrow \text{ComputeOverflow}(CM_{new})$
- 9: **if**  $OF_{new} < OF_v^{C-optimal}$  **then**
- 10:  $OF_v^{C-optimal} \leftarrow OF_{new}$
- 11:  $m_v^{C-optimal} \leftarrow m$
- 12:  $s_v^{C-optimal} \leftarrow s_v - (D_{m_v} - D_{m_v^{D-optimal}})$
- 13: **end if**
- 14: **end for**
- 15: **end if**
- 16: **end if**
- 17: }

---

computed routing congestion in the design under the assumptions about the placement of cells. The extension of our tree-mapping algorithm to the DAG-mapping one, such as [6], is, however, not obvious because of the duplication of wires due to the corresponding replication of the underlying logic gates.

In our current implementation, we do not store matches that are potentially good from the congestion standpoint, but have possibly inferior delay characteristics. Consequently, we do not minimize the congestion to the fullest extent possible. This can be remedied by storing a few inferior matches apart from those on the regular load-delay curve, at the cost of extra computation and memory. The memory requirement of our mapper is larger than that for a conventional mapper due to the storage of congestion maps for all non-inferior choices during the matching phase. This requirement may be reduced by storing just the bins corresponding to the bounding box that is affected by the mapping solution. Moreover, the memory efficient variant of the algorithm, which has the same memory capacity as that of the conventional delay oriented mapping, is possible by restricting the selection of the congestion-optimal matches to the primary outputs. In such a variant, the congestion maps for all non-inferior matches at all nodes except for the primary outputs are not required after the maps for the matches for the subsequent nodes in the topological order are computed during the matching phase, and hence, the memory, which is occupied by the corresponding congestion maps that have served the purpose of the forward propagation, can be freed and re-used. We expect our current mapper to be used in an engineering change order (ECO) mode, where at most a few thousand cells corresponding to the congestion hot-spots are (re-)synthesized. It is an ideal candidate for such an application, since it ensures that delay constraints will be met while minimizing total track overflow. This is demonstrated by the experimental results in the following section, which shows that circuits with more than 3500 cells can be handled well by our algorithm.

## 6. EXPERIMENTAL RESULTS

### 6.1 Experimental Setup

The congestion-aware mapping algorithm is implemented in C and incorporated in SIS [16]. We performed experiments on the

Example	Area	Conventional [16] / Ours						
		Overflow	(% gain)	Delay	MC	RU	# of cells	Run-time
	$\mu^2$			ps		%		s
C1355	3439	227 / 134	(40)	789 / 786	1.70 / 1.30	80 / 81	621 / 592	11 / 12
C1908	3616	323 / 225	(30)	1059 / 1042	1.70 / 1.40	80 / 80	578 / 571	12 / 13
C2670	11707	417 / 167	(59)	1258 / 1240	1.65 / 1.20	75 / 77	1482 / 1426	24 / 51
C3540	25994	1078 / 294	(72)	1655 / 1632	2.25 / 1.40	75 / 80	3254 / 3105	90 / 279
C432	1962	66 / 49	(25)	854 / 842	1.40 / 1.20	80 / 82	264 / 311	7 / 9
C499	3550	262 / 135	(48)	823 / 821	1.60 / 1.20	80 / 79	595 / 563	11 / 13
C5315	17265	1100 / 289	(73)	1120 / 1114	2.20 / 1.40	75 / 77	2122 / 2131	38 / 121
C6288	21379	515 / 452	(12)	4771 / 4731	1.70 / 1.40	80 / 80	3737 / 3596	88 / 135
C7552	28223	1343 / 547	(59)	1341 / 1309	1.60 / 1.30	75 / 73	3198 / 3080	132 / 213
C880	3944	378 / 260	(31)	890 / 884	1.70 / 1.20	80 / 76	584 / 575	12 / 13
Average		554 / 255	(44)	1455 / 1439	1.74 / 1.29	78 / 78	164 / 159	42 / 85

**Table 1: Comparison of the conventional mapping with our algorithm. ‘RU’ and ‘MC’ denote the average row utilization and maximum congestion, respectively.**

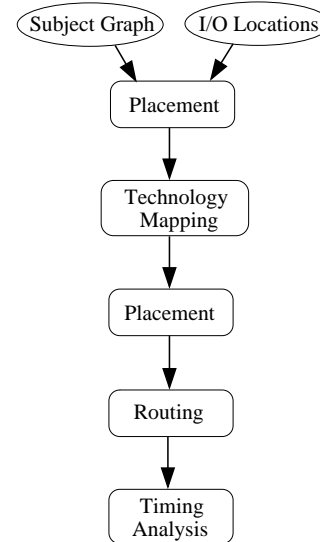
entire<sup>3</sup> set of ISCAS’85 benchmarks employing the design flow shown in Figure 6. For a given benchmark, the subject graph containing primitive gates is placed to create the companion placement. This is followed by either the conventional or our congestion-aware mapping. For the fair comparison, the conventional delay oriented technology mapping algorithm in SIS [16] is modified to use the companion placement information to compute wire-loads and wire-delays, as our congestion-aware technology mapping utilizes the same information not only for the bottom-up congestion map generation but also for the wire-load and wire-delay computation. The comparison with the competitive congestion-aware delay oriented mapping approaches [9, 10] could not be performed because of the unavailability of the access to proprietary benchmarks and tools. After the technology mapping, the resulting netlists are placed and routed. For all the experiments, technology mapping is performed employing lib2.genlib library in SIS, which is characterized for 100 nm technology [20] and contains up to 4 strengths for each cell. To generate delay constraints for a given benchmark, conventional technology mapping is run first and the delay of the most critical primary output is assigned as the required arrival time for all primary outputs. For the placement, we employ the publicly available recursive bisectioning based placer Capo [21], while, for routing, we use a router [22] that is based on non-Hanan routing. The post-routing delays are measured employing a static timing analyzer.

## 6.2 Analysis of Experimental Results

Table 1 shows the comparison of post global routing results due to conventional and our mapping. It shows block area in Column 2 for the benchmarks in Column 1, while Columns 3, 4, 5, 6, 7, and 8 show the total track overflow (with the improvement percentage), the delay, the maximum congestion, the average row utilization, the number of cells, and the run-time, respectively. All the experiments are run on Sun Ultra Sparc 60 machine with 400 MHz clock speed. From the table, we can observe the following:

1. Our mapping algorithm has been consistently able to reduce the track overflow, as shown in Column 3. On an average, the reduction is 44%. This improvement is encouraging, especially since the mapped netlists are re-placed without using the companion placement information employed during the technology mapping. The impressive gains in overflow also underline the ability of the technology mapping to reduce the congestion. The largest improvement, 72%, is in the case of

<sup>3</sup>The results on the smallest benchmark from the ISCAS’85 suite, C17, are not shown, as its implementation requires only a few gates, becoming an un-interesting case to show meaningful comparison.



**Figure 6: Design flow for conventional and congestion-aware technology mapping**

C3540, while the smallest one is in case of C6288. The small improvement in case of C6288 can be attributed to the relatively lower congestion in the design as compared to other benchmarks of a similar size.

2. From Column 4, it is clear that the delays due to congestion-aware mapping have improved slightly, as the algorithm maintains the delay-optimality, resulting in the same delays as the conventional mapping. The subsequent stages, especially the routing, show the effect of reduced congestion, causing smaller detours and hence, smaller delays as compared to the conventionally mapped netlists.
3. The maximum congestion shown in Column 5 has improved due to our mapping in all the cases with an average improvement of 25%, even though the algorithm targets only the total track overflow. This can be ascribed to the correlation between the maximum congestion and the total track overflow: generally, high overflow implies the large maximum congestion and vice-versa.
4. It appears from Column 6 that the average row utilization has increased, as in cases of C2670 and C3540, and has also de-

creased, as in cases of C7552 and C880, with the track overflow. This can be attributed to the corresponding increase or decrease in the cell area. Depending on the context, an area increase or decrease can result in a small track overflow: an example in [11] shows that even an area optimal match, which results in a relatively smaller area, can worsen the congestion, while it is also possible that an increase in area due to more gates, and hence, more wires, proves detrimental to the objective of reducing the congestion. This is why the congestion is relatively insensitive to the average row utilization, when the utilization is not too high.

5. As shown in Column 7, the number of cells in the congestion-aware mapped netlists have decreased slightly in all cases but that of C432. The decrease in the number of cells can be associated with the corresponding reduction in the number of nets, which may be indirectly correlated with metrics such as structural pin-density, and hence, congestion, as pointed out in [14]. In case of C432, however, employing more logic cells has resulted in the congestion alleviation, which can be explained using the discussion in the previous remark on the average row utilization.
6. Finally, as can be seen from Column 8, the run-times for congestion-aware mapping are still comparable to those of the conventional one. On an average, the run-time is 2.02 times worse, but still practical, being within a few minutes on a 400 MHz Sun Ultra Sparc machine, than that for the conventional mapping. It shows that the constants subsumed by  $O()$  in the time complexity expression in Section 4.3 are not too dominating. This bodes well for the applicability of the algorithm.

## 7. CONCLUSION

In this paper, we have presented a technology mapping algorithm targeting routing congestion. We have shown how to overcome the “chicken-and-egg” problem during the mapping and placement stages by generating and propagating congestion maps associated with each choice during the matching phase. We have proposed a covering procedure, which exploits slacks to select the congestion-optimal choices that preserve the delay-optimality of the solution. The experimental results on ISCAS’85 benchmarks prove the efficiency of the algorithm, as they show, on an average, 44% improvement in the track overflow as compared to the conventional mapping. We believe that these results may be further improved by applying the placement legalization instead of the full-scale replacement for the mapped netlists, since it may preserve the gains obtained during the technology mapping. Nonetheless, the impressive gains in the track overflow emphasize the ability of the technology mapping to reduce the routing congestion. The proposed matching and covering techniques are sufficiently general and can be applied to optimize different cost functions defined on congestion maps as well as different physical properties, which can be captured employing two-dimensional maps, such as power-density maps.

## 8. REFERENCES

- [1] G. E. Moore. Cramming more components onto integrated circuits. In *Electronics*, pages 114–117, April 1965.
- [2] P. Saxena, N. Menezes, P. Cocchini, and D. A. Kirkpatrick. Repeater scaling and its impact on CAD. *IEEE Trans. CAD*, 23(4):451–463, April 2004.
- [3] P. Kapur, G. Chandra, and K. C. Saraswat. Power estimation in global interconnects and its reduction using a novel repeater insertion methodology. In *Proc. DAC*, pages 461–466, June 2002.
- [4] K. Keutzer. DAGON: Technology binding and local optimization by DAG matching. In *Proc. DAC*, pages 341–347, June 1987.
- [5] K. Chaudhary and M. Pedram. A near optimal algorithm for technology mapping minimizing area under delay constraints. In *Proc. DAC*, pages 492–498, June 1992.
- [6] Y. Kukimoto, R. K. Brayton, and P. Sawkar. Delay-optimal technology mapping by DAG covering. In *Proc. DAC*, pages 348–351, June 1998.
- [7] J. Y. Lin, A. Jagannathan, and J. Cong. Placement-driven technology mapping for LUT-Based FPGA’s. In *Proc. ISFPGA*, pages 121–126, February 2003.
- [8] D. Pandini, L. T. Pileggi, and A. J. Strojwas. Congestion aware logic synthesis. In *Proc. DATE*, pages 664–671, March 2002.
- [9] D. Pandini, L. T. Pileggi, and A. J. Strojwas. Global and local congestion optimization in technology mapping. *IEEE Trans. CAD*, 22(4):498–505, April 2003.
- [10] L. Stok and T. Kutzschebauch. Congestion aware layout driven logic synthesis. In *Proc. ICCAD*, pages 216–223, November 2001.
- [11] R. Shelar, S. Sapatnekar, P. Saxena, and X. Wang. A predictive distributed congestion metric and its application to technology mapping. In *Proc. ISPD*, pages 210–217, April 2004.
- [12] Q. Liu and M. Marek-Sadowska. Technology mapping with pre-layout wirelength prediction. In *Proc. IWLS*, June 2004.
- [13] P. Kudva, A. Sullivan, and W. Dougherty. Metrics for structural logic synthesis. In *Proc. ICCAD*, pages 551–556, November 2002.
- [14] Q. Liu and M. Marek-Sadowska. Pre-layout wirelength and congestion estimation. In *Proc. DAC*, pages 582–587, June 2004.
- [15] M. Pedram and N. Bhat. Layout driven technology mapping. In *Proc. DAC*, pages 99–105, June 1991.
- [16] E. M. Sentovich *et al.* SIS: A system for sequential circuit synthesis. Memorandum No. UCB/ERL M92/41, May 1992.
- [17] J. Westra, C. Bartels, and P. Groeneveld. Probabilistic congestion prediction. In *Proc. ISPD*, pages 204–209, April 2004.
- [18] W-J. Dai. Hierarchical physical design methodology for multi-million gate chips. In *Proc. ISPD*, pages 179–181, April 2001.
- [19] J. Lou, S. Thakur, S. Krishnamoorthy, and H. S. Sheng. Estimating routing congestion using probabilistic analysis. *IEEE Trans. CAD*, 21(1):32–41, January 2002.
- [20] Berkeley predictive technology model. <http://www-device.eecs.berkeley.edu/~ptm/download.html>.
- [21] A. E. Caldwell, A. B. Kahng, and I. L. Markov. Can recursive bisection alone produce routable placements? In *Proc. DAC*, pages 477–482, June 2000.
- [22] J. Hu and S. Sapatnekar. A timing-constrained simultaneous global routing algorithm. *IEEE Trans. CAD*, 21(9):1025–1036, September 2002.