

Fast and Exact Transistor Sizing Based on Iterative Relaxation

Vijay Sundararajan¹, Sachin S. Sapatnekar², Keshab K. Parhi²

¹Texas Instruments, Wireless Infrastructure Branch
Dallas, TX 75243

²Dept. of ECE, University of Minnesota
Minneapolis, MN 55455

E-mail: vijay@ti.com, sachin@ece.umn.edu, parhi@ece.umn.edu

Abstract— This paper presents MINFLOTRANSIT, a new transistor sizing tool for fast sizing of combinational circuits with minimal cost. MINFLOTRANSIT is an iterative relaxation based tool that has two alternating phases. For a circuit with $|V|$ transistors and $|E|$ wires, the first phase (*D-phase*) is based on minimum cost network flow, which in our application, has a worst-case complexity of $O(|V||E|\log(\log(|V|)))$. The second phase (*W-phase*) has a worst case complexity of $O(|V||E|)$. In practice, during our simulations both the *D-phase* and *W-phase* show a near linear run-time dependence on the size of the circuit, comparable to TILOS. Simulation results show excellent run-time behavior for MINFLOTRANSIT on all the ISCAS85 benchmark circuits. For reasonable delay targets MINFLOTRANSIT shows up to 16.5% area savings (in relatively large circuits) over a circuit sized using a TILOS-like algorithm. In our opinion the primary contribution of this paper is to take advantage of the structure of the transistor sizing problem and devise an iterative relaxation based gradient descent approach (*D-phase*) that has excellent convergence properties.

I. INTRODUCTION

As evidenced by the successive announcement of ever faster computer systems, increasing the speed of VLSI systems is one of the major requirements for VLSI system designers today. Faster integrated circuits are making possible newer applications that were traditionally considered difficult to implement in hardware. In this scenario of increasing circuit complexity, reduction of circuit delay in integrated circuits is an important design objective. Transistor sizing is one such task that has been employed for speeding up circuits for quite some time now [1]. Given the circuit topology, the delay of a combinational circuit can be controlled by varying the sizes of transistors in the circuit. Here, the size of a transistor is measured in terms of its channel width, since the channel lengths of MOS transistors in a digital circuit are generally uniform. In any case, what really matters is the ratio of channel width to channel length, and if channel lengths are not uniform, this ratio can be considered as the size. In coarse terms, the circuit delay can usually be reduced by increasing the sizes of certain transistors in the circuit from the minimum size. Hence, making the circuit faster usually entails the penalty of increased circuit area relative to a minimum sized circuit and the area-delay trade-off involved here is the problem of transistor size optimization. A related problem to transistor sizing is called gate sizing, where a logic gate in a circuit is modeled as an equivalent inverter and the sizing optimization is carried out on this modified circuit with equivalent inverters in place of more complex gates. There is, therefore, a reduction in the number of size parameters corresponding to every gate in the circuit. Needless to say, this is an easier problem to solve than the general transistor sizing problem. Note that gate sizing mentioned here is distinct from library specific gate sizing that is a discrete optimization problem targeted to selecting appropriate gate sizes from an underlying cell library. The gate sizing problem targeted here is one of continuous gate sizing where the gate sizes are allowed to vary in a continuous

manner between a minimum and a maximum size.

There has been a large amount of work done on transistor sizing [1], [2], [3], [4], [5], [6], [7], [8], [9] that underlines the importance of this optimization technique. Starting from a minimum sized circuit, TILOS [1] uses a greedy strategy for transistor sizing by iteratively sizing transistors in the critical path. A sensitivity factor is calculated for every transistor in the critical path to quantify the gain in circuit speed achieved by a unit upsizing of the transistor. The most sensitive transistor is then bumped up in size by a small constant factor to speed up the circuit. This process is repeated iteratively until the timing requirements are met. The technique is extremely simple to implement and has run-time behavior proportional to the size of the circuit. Its chief drawback is that it does not have guaranteed convergence properties and hence is not an exact optimization technique. Among the past approaches, only [2] and [8] are exact optimization techniques and the other techniques do not have proven convergence properties. While [2] was the first ever polynomial time technique reported for addressing the problem exactly, it does not have very good run-time behavior. The technique presented in [8] has shown impressive run-time behavior for sizing large adders. The run-time behavior of this technique on more complex circuits such as multipliers and controllers was not demonstrated. Moreover, the approach appears to be amenable only to tackling sizing problems where the gate delay is expressed using Elmore delays. The approach in [9] is a unique approach in that it uses a very accurate dynamic timing analyzer. While this renders the solutions exact, however there is a huge penalty paid in times of computation speed. What is required is an approach that can handle a loosely restricted class of delay models and at the same time can perform sufficiently fast optimizations. The result of such an optimization can then be used as an input to a more detailed approach like [9] to obtain exact solutions with a reasonable computational delay.

This paper presents a novel way of solving the transistor sizing problem exactly and in an extremely fast manner. The proposed approach has some similarity in form to [6], [7], [10] which will be subsequently explained, but the similarity in content is minimal and the details of implementation are vastly different. In essence, the proposed technique and the techniques in [6], [7], [10] are iterative relaxation approaches that involve a two-step optimization strategy. The first-step involves a delay budgeting step where optimal delays are computed for transistors/gates. The second step involves sizing transistors optimally under this “constant delay” model to achieve these delay budgets. The two steps are iteratively alternated until the solution converges, i.e., until the delay budgets calculated in the first step are exactly satisfied by the transistor sizes determined by the second step. The primary features of the proposed approach are:

- It is computationally fast and is comparable to TILOS in its run-time behavior.
- It can be used for true transistor sizing as well as the relaxed problem of gate-sizing. Additionally, the approach can easily

incorporate wire sizing, as outlined in section II-A.

- It can be adapted for more general delay models than the Elmore delay model, see Appendices A-C for details.

The starting point for the proposed approach is a fast guess solution. This could be obtained, for example, from a circuit that has been optimized using TILOS to meet the given delay requirements. The proposed approach, as outlined earlier, is an iterative relaxation procedure that involves an alternating two-phase relaxed optimization sequence that is repeated iteratively until convergence is achieved. The two-phases in the proposed approach are:

- The *D-phase* where transistor sizes are assumed fixed and transistor delays are regarded as variable parameters. Irrespective of the delay model employed, this phase can be formulated as the dual of a min-cost network flow problem. Using $|V|$ to denote the number of transistors and $|E|$ the number of wires in the circuit, this step in our application has worst-case complexity of $O(|V||E|\log(\log|V|))$ [11].
- The *W-phase* where transistor/gate delays are assumed fixed and their sizes are regarded as variable parameters. As long as the gate delay can be expressed as a separable function of the transistor sizes, this step can be solved as a Simple Monotonic Program (SMP) [12]. The complexity of SMP is similar to an all-pairs shortest path algorithm in a directed graph, [12], [13], i.e., $O(|V||E|)$.

The objective function for the problem is the minimization of circuit area. In the W-phase, this objective is addressed directly, and in the D-phase the objective is chosen to facilitate a move in the solution space in a direction that is known to lead to a reduction in the circuit area.

II. PROPOSED APPROACH

The transistor size optimization problem can be stated as,

$$\begin{aligned} & \underset{\text{all transistors, } i}{\text{minimize}} && \sum x_i, \\ & \text{subject to: } && \text{Delay}(\text{Circuit}) < T, \\ & && \text{minsize} \leq x_i \leq \text{maxsize}, \end{aligned} \quad (1)$$

where x_i refers to the size of transistor i , T is timing requirement specified as an input to the optimization and minsize , maxsize are, respectively, minimum and maximum bounds on the sizes of transistors in the circuit.

A. Equivalent DAG for Transistors/Gates

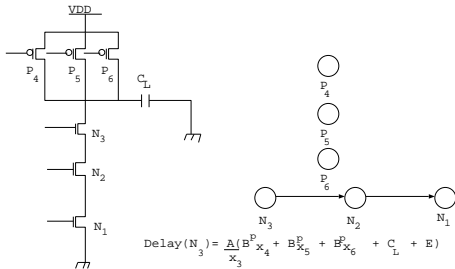


Fig. 1. The DAG corresponding to a 3-input static CMOS nand gate.

The proposed approach requires transistor delay to be expressible as a sum of *simple monotonic functionals* of transistor sizes, which are defined as follows:

Definition II.1: A function $D_i(x_1, \dots, x_n)$ is a simple monotonic functional if it can be rewritten as $D_i =$

$g(x_i)q(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$, where $g(x_i)$ is a monotonic decreasing function of x_i and $q(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ is a monotonic increasing function of each x_j $j \in \{1, \dots, n\}$ $j \neq i$.

Definition II.2: A function $D(x_1, \dots, x_n)$ is termed decomposable into simple monotonic functionals if $D = \sum_{i \in \{1, \dots, n\}} D_i$ where each D_i , as in definition (II.1), is a simple monotonic functional. D_i is termed the simple monotonic projection of D on x_i .

In order to mathematically model the transistor size optimization problem, every static CMOS gate in the circuit is first converted in to an equivalent Directed Acyclic Graph (DAG) model, shown in figure 1, as follows. There is a vertex in the DAG corresponding to every transistor. An edge is drawn between an NMOS (PMOS) transistor and another NMOS (PMOS) transistor provided there is a discharging (charging) path consisting of the two transistors. The edge is always directed from the transistor higher up in the discharging (charging) path to the transistor lower down in the discharging (charging) path. Every vertex of this DAG has a delay attribute associated with it. This delay attribute is given by the *simple monotonic projection* of the worst case discharging (charging) path delay through the transistor corresponding to the vertex on to the size of this transistor. Note that, as will be evident soon, the rise and fall delays are implicitly distinguished due to the fact that the DAG corresponding to every gate has separate components for pullup and pulldown networks.

For ease of exposition, we will henceforth consider the commonly used Elmore delay model that can be decomposed in to *simple monotonic functionals*. Assuming x_i to be the size of transistor $N_i(P_i)$ in the pulldown(pullup) network of the 3-input NAND gate shown in figure 1, it can be shown [1] that the pulldown Elmore delay can be expressed as,

$$\begin{aligned} \text{delay}^{\text{pulldown}} = & \left(\frac{A}{x_1}\right)(Bx_1 + Cx_2) + \left(\frac{A}{x_1} + \frac{A}{x_2}\right)(Bx_2 + Cx_3 + D) + \\ & \left(\frac{A}{x_1} + \frac{A}{x_2} + \frac{A}{x_3}\right)(Bx_3 + B^P x_4 + B^P x_5 + B^P x_6 + C_L + E), \end{aligned} \quad (2)$$

where A , B and C are constant coefficients, the resistance, drain and source capacitance, respectively, of a unit NMOS transistor. D and E are related to the wire capacitances. Similarly, B^P is the drain capacitance of a unit sized PMOS transistor and C_L is the load capacitance. Under this model if wire sizes were considered to be variables also then the form of (2) would remain similar. Rewriting the expression in (2) we get,

$$\begin{aligned} \text{delay}^{\text{pulldown}} = & \frac{A}{x_1}(Bx_2 + Bx_3 + Cx_2 + Cx_3 + D + E + \\ & B^P x_4 + B^P x_5 + B^P x_6 + C_L) + 3AB + \frac{A}{x_2}(Bx_3 + Cx_3 + D + E + \\ & B^P x_4 + B^P x_5 + B^P x_6 + C_L) + \frac{A}{x_3}(B^P x_4 + B^P x_5 + \\ & B^P x_6 + C_L + E). \end{aligned} \quad (3)$$

Since A , B , C , D , E , B^P , C_L are non-negative quantities, the Elmore delay model admits a *simple monotonic decomposition*. The above fact is illustrated in figure 1 where the delay corresponding to NMOS transistor N_3 is explicitly shown. The constant terms used in this expression have the same connotation as in (3). We claim that such a DAG model can always be developed for any complex static CMOS gate consisting of a series/parallel network of transistors as long as the underlying delay model admits a *simple monotonic decomposition*. This is a reasonable requirement since the reduction in the gate delay with an increase in its size can be modeled by the function g in

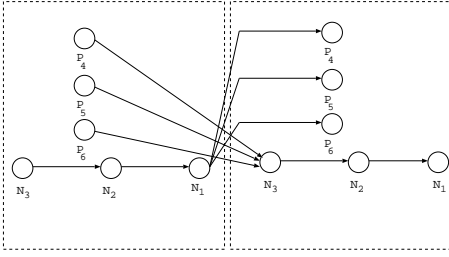


Fig. 2. The DAG corresponding to a circuit consisting of two 3-input static CMOS nand gates in series.

Definition II.1, and the increase in the gate delay with increasing fanout gate sizes can be modeled by the function q .

In addition, if wire sizing were also to be performed together with transistor sizing, then we could model the problem by augmenting the DAG corresponding to a gate by adding vertices corresponding to each wire. The edges emanating from and incident on these wires will be similarly constructed as for transistors. The delay attribute of a vertex corresponding to any wire can also be similarly defined as for that of a vertex corresponding to a transistor. We conclude that modeling the problem of wire sizing along with transistor sizing may use the same framework as transistor sizing alone, and the approach developed in this paper can simultaneously handle both. For ease of exposition, from here onwards, wire sizing will not be considered for the remainder of the paper.

Note that the DAG corresponding to a static CMOS gate has at least two disjoint connected components, as shown in figure 1, corresponding to the pulldown network of NMOS transistors (i.e., vertices N_1, N_2, N_3) and the pullup network of PMOS transistors (i.e., vertices P_1, P_2, P_3) corresponding to the gate. The portion of the DAG representing the NMOS pulldown networks corresponds to falling transitions and the portion of the DAG representing the PMOS pullup network is related to rising transitions at the output of the gate. Note that there are several vertices in the DAG of a gate that only have edges emanating from them and have no edges terminating on them; we refer to these vertices as the **root vertices** of the gate DAG. Also, note that there are several vertices in the DAG of a given gate that have only edges terminating on them and no edges emanating from them; these vertices constitute the **leaf vertices** of the gate DAG.

B. A DAG for the Circuit

The entire circuit consisting of static CMOS gates can be represented with an equivalent DAG, $G = (V, E)$, by connecting the component DAG's of individual gates. The construction of the circuit DAG is as follows, the vertex set V of the circuit DAG is simply the union of the vertex sets of DAG's corresponding to the gates in the circuit. The edge set E of the circuit DAG is constructed as follows. For every wire connecting the output of one gate to the input of another there will be a set of edges in the circuit DAG that go from the NMOS (PMOS) DAG components of the first gate to the PMOS (NMOS) DAG components of the second gate. So corresponding to every wire connecting the output of the first gate to a given NMOS (PMOS) transistor in the second gate, there will be edges emanating from all the leaf vertices of the PMOS (NMOS) DAG of the first gate. These edges terminate on all the root vertices of the NMOS (PMOS) DAG component of the second gate that are connected to the given transistor in the second gate. Figure 2 illustrates the construction of the circuit DAG for a circuit consisting of two

3-input nand gates in series.

C. Two-Phase Optimization

Note that the delay corresponding to a vertex, i , in the circuit DAG, whose corresponding transistor has a size x_i , can always be expressed as,

$$delay(i) = \frac{\sum_{j \in S(V(G))} a_{ij} x_j + b_i}{x_i}, \quad (4)$$

where $S(V(G))$ denotes some subset of $V(G)$ that is located in the neighborhood of the vertex i . In particular, this subset consists of the vertices corresponding to all those transistors whose sizes directly affect the delay of the transistor corresponding to vertex i . Also, note that all the coefficients a_{ij} , b_i in (4) are non-negative. Rearranging (4), we have,

$$delay(i) \cdot x_i - \sum_{j \in S(V(G))} a_{ij} x_j = b_i. \quad (5)$$

Denoting a diagonal matrix whose $(i, i)^{th}$ entry is $delay(i)$ by D , a matrix whose $(i, j)^{th}$ entry is a_{ij} by A , a column vector whose i^{th} component is b_i by \mathbf{B} and a column vector whose i^{th} component is x_i by \mathbf{X} we can rewrite (5) as

$$(D - A)\mathbf{X} = \mathbf{B}. \quad (6)$$

The formulation in (6) can be written as long as the delay model admits a simple monotonic decomposition. Note that, as discussed in Appendices A-C, more general delay models are admissible under the framework developed in this paper.

It can be shown that for strict gate sizing the matrix $(D - A)$ can be written as an upper triangular matrix. This is due to the fact that the adjacency matrix of the circuit DAG can be always written as an upper triangular matrix [14]. On the other hand it can be shown that for transistor sizing the matrix $(D - A)$ can be written as a block upper triangular matrix. This is possible because under the present delay model the delay component associated with a transistor is independent of the sizes of the transistors that topologically precede it after we levelize the gate level circuit netlist. The delay components of the transistor within a gate may, however, depend on the sizes of transistors above and below it and belonging to the same gate as itself. Therefore, if we reorder the matrix $(D - A)$ by having the rows corresponding to transistors closer to the primary inputs PI 's higher than the transistors closer to the primary outputs PO 's then the matrix $(D - A)$ is guaranteed to be block upper triangular. The individual block sizes are no more than the maximum of the sum of the number of transistors in a gate and the number of transistors in its immediate fanout. We will henceforth assume that $(D - A)$ can always be represented in an upper triangular form or block upper triangular form. With this assumption we can state that if D is a constant matrix, then as long as $(D - A)$ is invertible, a system of equations of the form $(D - A)\mathbf{X} = \mathbf{B}$ can be solved for the variables \mathbf{X} by a backward substitution process beginning from the bottom row and proceeding upwards and progressively solving for all x_i . From a circuit point of view, this process proceeds in a backward breadth-first manner beginning with the primary outputs and proceeding backwards in order of decreasing levels of logic of the circuit. This elimination process has $O(|F|N)$ computational complexity, where N is the number of components in the vector \mathbf{X} and $|F|$ is bounded (in gate sizing) by the maximum fanout of any gate in the circuit. Note that as long as all vertices of the circuit DAG have a non-zero delay, which is always the case, the (block) upper triangular matrix $(D - A)$ for (transistor) gate sizing will always be invertible.

Now assume that we start with some initial sizing solution \mathbf{X}_0 and some delay matrix D_0 satisfying (6). We now resize the transistors slightly so that the new delay matrix is $D_0 + \delta D$ and the new size vector is $\mathbf{X}_0 + \delta \mathbf{X}$, so that we then have,

$$\begin{aligned} (D_0 + \delta D - A)(\mathbf{X}_0 + \delta \mathbf{X}) &= \mathbf{B}, \\ (D_0 - A)(\mathbf{X}_0 + \delta \mathbf{X}) + \delta D(\mathbf{X}_0 + \delta \mathbf{X}) &= \mathbf{B}, \\ \mathbf{B} + (D_0 - A)\delta \mathbf{X} + \delta D\mathbf{X}_0 + \delta D\delta \mathbf{X} &= \mathbf{B}, \\ (D_0 - A)\delta \mathbf{X} &\approx -\delta D\mathbf{X}_0, \\ \delta \mathbf{X} &\approx -(D_0 - A)^{-1}\delta D\mathbf{X}_0, \end{aligned} \quad (7)$$

where the term $\delta D\delta \mathbf{X}$ has been ignored, assuming small perturbations in δD and $\delta \mathbf{X}$. In other words, we make the following observations:

- For an infinitesimal resizing of the transistors corresponding to the vertices in the circuit DAG, the vector of infinitesimal changes in transistor sizes can be represented as a linear vector function of the infinitesimal changes in transistor delays.
- As a result of (1), we see that the sum of all the components of $\delta \mathbf{X}$, which represents the sum of the change in sizes of the transistors corresponding to all the vertices in the circuit, can be expressed as a linear function of diagonal entries of the matrix δD .

It can be shown that since all components of \mathbf{X}_0 are positive, all components of $-(D_0 - A)^{-1}\mathbf{X}_0$ will be negative. Hence, the sum of all the components of $\delta \mathbf{X}$, can be expressed as a linear function of diagonal entries of the matrix δD , where the coefficient corresponding to each diagonal element of δD is negative. While we have discussed an Elmore delay model so far, we can extend the above discussion to more general delay models as detailed in Appendix C.

This motivates a two-phase strategy for solving the transistor size optimization problem. In the D-phase, as above, we assume fixed transistor sizes and redistribute the delay budgets in such a manner as to minimize the resultant change in transistor sizes. In the W-phase, on the other hand, we try to find the minimal-sized circuit that satisfies the modified delay budget obtained after the D-phase. The two-phases are alternated till convergence is achieved and the delay budgets output by the D-phase are exactly satisfied by the transistor sizes calculated by the W-phase.

C.1 D-phase

In the D-phase as mentioned earlier, we map the optimization problem on to the delay (D) domain by, a) capturing the delay profile at a given (area, delay) configuration of the circuit, b) computing the slack for all the circuit elements for the given delay configuration, c) re-distribute the slack in the circuit by doing microscopic delay budget reallocation globally over the whole circuit, d) driving the delay reallocation toward a lower area circuit by modeling infinitesimal area changes exactly as a linear function of infinitesimal delay changes.

It is observed that the delay reallocation phase is akin to the dual of a min-cost network flow optimization problem, which aids in solving for the D-phase with extremely fast execution times.

Before proceeding further we first need to develop some terminology. First assume that the circuit has been sized to meet delay requirements using an algorithm such as TILOS. We now define three attributes for every vertex in the circuit DAG G . For a vertex i , these are the arrival time $AT(i)$, the required time $RT(i)$ and the slack, $sl(i)$. Additionally, every edge $e_{ij} \in E$ has the attribute *edge-slack*, $esl(e_{ij})$. The entire circuit graph G has an additional attribute $CP(G)$ that refers to the delay

of the critical path of the corresponding circuit. We will now define all of these attributes formally.

$$\begin{cases} AT(i) &= \text{external time of arrival, } u \in PI, \\ &= \max_{v \in fanin(i)} (AT(j) + delay(j)), \text{ else} \\ CP(G) &= \max_{u \in V} (AT(i) + delay(i)), \\ RT(i) &= CP(G) - delay(i), \text{ } u \in PO, \\ RT(i) &= \min_{v \in fanout(i)} (RT(j) - delay(i)), \text{ else} \\ sl(i) &= RT(i) - AT(i), \\ esl(e_{ij}) &= RT(j) - AT(i) - delay(i). \end{cases} \quad (8)$$

where PI and PO denote respectively the primary inputs and primary outputs of the circuit. We call a circuit safe when all vertices $i \in V$ have $sl(i) \geq 0$ and all edges have $esl(e_{ij}) \geq 0$. The D-phase involves minimally altering the delay budgets of transistors in the circuit to move towards a feasible minimum area solution. For this to be possible, we need to capture the slack (available delay budget) for every transistor and also present a strategy to alter/redistribute these delay budgets. In the next section, we will first present an approach to capture the slack in a circuit in terms of fictitious buffer-like entities called *Fictitious Specific Delay Units (FSDU's)*. Next, an approach called *FSDU-displacement* will be presented, which redistributes the delay budgets for every transistor in such a manner that a lower area solution (from the present solution) is achieved that is also timing feasible.

Delay Balancing

A given circuit DAG G can be transformed to a functionally equivalent circuit DAG G' by introducing dummy units of appropriate delay on to each edge in the circuit DAG in such a manner that for every $e_{ij} \in E$, $esl(e_{ij}) = 0$ and $CP(G') = CP(G)$ [15]. This process is known as delay balancing. For our purposes, we do not explicitly insert physical delays. Instead, we use the *concept* of delay balancing as a tool to capture all the slack in the circuit DAG. This captured slack is then used for the D-phase optimization. The delay units used for delay balancing are, therefore, *fictitious* entities whose only purpose is to model the slack present in the circuit DAG. We refer to these fictitious delay units as FSDUs (Fictitious Specific Delay Units). Figure 3 shows a circuit DAG and figure 4 shows its delay balanced counterpart; the ‘‘square boxes’’ on the edges of the circuit in figure 4 represent the FSDUs on that edge.

Starting with a given circuit DAG there are several possible

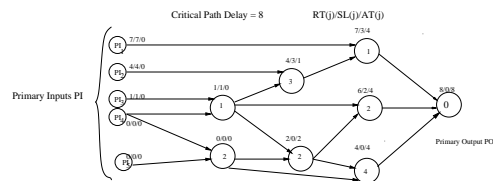


Fig. 3. An example of a circuit DAG the integer numbers within each vertex represent its delay and each vertex i has the triplet $(RT/SL/AT)$ above it.

ways to produce a delay balanced graph. Any such delay balanced graph will from now on be referred to as a delay balanced configuration.

FSDU-Displacement

We define *FSDU-Displacement*, a circuit DAG transformation technique, as a mapping $r:V \rightarrow \mathbf{Z}$, $\{\mathbf{Z}: \text{the set of integers}\}$ such that the delay of the FSDU on the edge e_{ij} after FSDU-Displacement, $FSDU^r(e_{ij})$, is related to the delay of the FSDU

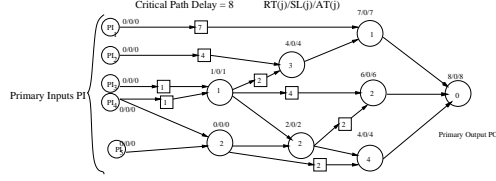


Fig. 4. The circuit DAG in figure 3 after delay balancing. The square boxed integers on edges represent the FSDUs added to the edges for delay balancing.

before FSDU-Displacement, $FSDU(e_{ij})$, by,

$$FSDU^r(e_{ij}) = FSDU(e_{ij}) + r(j) - r(i). \quad (9)$$

We state the following without proof, the proof can be found for similar results for retiming in [16].

Theorem II.1: All legal delay balanced configurations for a given circuit-graph G are *FSDU-Displaced* versions of each other.

Theorem II.2: The net change in the delay of any structural path from a vertex i to another vertex j after *FSDU-Displacement* is always $r(j) - r(i)$.

The above theorem gives rise to the following corollary.

Corollary II.1: If we connect all the leaf vertices corresponding to primary output nodes of a given circuit DAG to a common dummy vertex O through dummy edges and if we restrict $r(O)$ to be exactly 0 and also restrict $r(I)$ for every input vertex $I \in PI$ to be exactly 0, then the critical path of the transformed circuit DAG after FSDU-displacement remains unaltered.

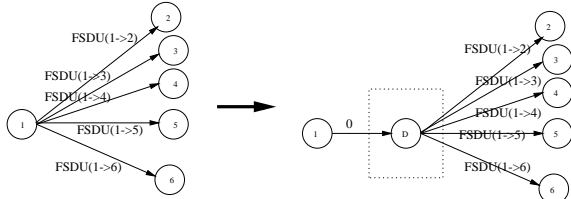


Fig. 5. Illustration of circuit DAG transformation related with the addition of a dummy vertex at the output of every vertex.

Before we develop a formal mathematical model, we first modify the circuit DAG by adding a dummy vertex $Dmy(i)$ of delay 0 units at the output of every vertex i in the circuit DAG. A dummy edge connects vertex i to its corresponding dummy vertex $Dmy(i)$. All fanout edges that initially originated from vertex i now originate from $Dmy(i)$. This is done so as to be able to effectively model $\min_{j \in Fanout(i)} FSDU(i \rightarrow j)$, which represents the maximum delay of the FSDU's in the fanout edges of vertex i in the D-phase optimization. Using (7) and Theorem II.2, it turns out that in the D-phase optimization, $\sum_i \delta \mathbf{X}_i \equiv \sum_i cost_i \delta \mathbf{D}_i \equiv \sum_i cost_i (r(Dmy(i)) - r(i))$. The relation $\Delta \mathbf{D}_i = r(Dmy(i)) - r(i)$ only holds if $cost(i) \geq 0$. This implies that as long as the problem is a *maximization* problem and if $cost(i) \geq 0$, then $r(Dmy(i)) - r(i)$ correctly models $\min_{j \in Fanout(i)} FSDU(i \rightarrow j)$.¹ Figure 5 illustrates this circuit DAG transformation with an example. Now we can summarize the D-phase as follows:

¹ See Appendix A for the proof of the fact that $cost(i) \geq 0$ for every vertex i .

D-phase

(1) Produce any valid delay balanced configuration of the given circuit DAG. We use a depth first FSDU insertion heuristic for this purpose [15].

(2) Now starting from the delay balanced configuration in (1) above, let $\delta \mathbf{X} = -(D_0 - A)^{-1} \delta D \mathbf{X} = -\mathbf{C}^T \text{diag}(\delta D)$ where $\mathbf{C}^T = -(D_0 - A)^{-1} \mathbf{X}$, all other symbols are as defined earlier and $\text{diag}(\delta D)$ is a column vector consisting of the diagonal elements of δD . Note that minimizing $\sum \delta \mathbf{X}_i \equiv$ minimizing $\sum \mathbf{X}_i$. Now, for every vertex i let $\delta \mathbf{D}_i = r(Dmy(i)) - r(i)$, which means that the delay of the FSDU at the output of a vertex is the change in its delay after the D-phase.

(3) To maintain the requirement that δD will be small, introduce the following constraints for every vertex.

$$FSDU(i \rightarrow Dmy(i)) + r(Dmy(i)) - r(i) \geq MIN\Delta D(i),$$

$$FSDU(i \rightarrow Dmy(i)) + r(Dmy(i)) - r(i) \leq MAX\Delta D(i),$$

where $MIN\Delta D(i)$ and $MAX\Delta D(i)$ bound the change in delay of vertex i from both sides, i.e., decrease or increase of vertex delay. The technique used to determine $MIN\Delta D(i)$ and $MAX\Delta D(i)$ is described in Appendix B.

(4) For every edge $e(Dmy(i) \rightarrow j)$, introduce the causality constraint that states that the slack for all edges in the original DAG will be non-negative after the D-phase.

$$FSDU^r(Dmy(i) \rightarrow j) = FSDU(Dmy(i) \rightarrow j) + r(j) - r(Dmy(i)) \geq 0.$$

(5) Now solve the following optimization problem, whose dual is a min-cost network flow problem [17],

$$\begin{aligned} & \text{minimize} && \sum_{\text{over vertices } i} \mathbf{x}_i \\ & \equiv \text{maximize} && \sum_{\text{over vertices } i} C_i \cdot (r(Dmy(i)) - r(i)) \\ & \text{subject to:} && \\ & FSDU(i \rightarrow Dmy(i)) + r(Dmy(i)) - r(i) && \geq MIN\Delta D(i), \\ & FSDU(i \rightarrow Dmy(i)) + r(Dmy(i)) - r(i) && \leq MAX\Delta D(i), \\ & \text{For all edges } Dmy(i) \rightarrow j, && \\ & FSDU^r(Dmy(i) \rightarrow j) = FSDU(Dmy(i) \rightarrow j) && \\ & + r(j) - r(Dmy(i)) && \geq 0. \end{aligned} \quad (10)$$

Note that the D-phase optimization is in the form of the dual of a minimum cost network flow problem, [11]. Also, the constant terms in the RHS of the constraints in the D-phase can be integerized by appropriate scaling, i.e., by multiplying every constant term by some power of 10, and then rounding off the product. By choosing appropriate powers of 10 arbitrary accuracy can be maintained with almost no penalty in computational requirements. In this way, fast methods devised for integerized minimum cost network flow approaches [11] can be fruitfully employed in solving the D-phase optimization problem.

C.2 W-phase

Once the D-phase has computed new delays (delay budgets) for all the vertices in the circuit DAG, we need to find feasible sizes for the transistors corresponding to every vertex in the circuit DAG to satisfy the delay requirements while using up minimal area. In effect we have to solve the following problem,

$$\text{minimize} \quad \sum_{\text{over vertices } i} x_i,$$

$$\begin{aligned}
\text{subject to: } & \frac{\sum_{j \in S(V(G))} a_{ij} x_j + b_i}{x_i} \leq \text{delay}(i), \\
& \equiv \sum_{j \in S(V(G))} a_{ij} x_j + b_i \leq \text{delay}(i) \cdot x_i, \\
& \text{minsize} \leq x_i \leq \text{maxsize}. \tag{11}
\end{aligned}$$

It turns out that due to the non-negativity of a_{ij} , $\text{delay}(i)$ and the coefficients of x_i in the objective function, this optimization problem can be modeled as a Simple Monotonic Program (SMP) [12]. This kind of problem can be solved by a constraint relaxation procedure with worst case complexity of $O(|V||E|)$ where $|E|$ is the number of constraints and $|V|$ is the number of variables. The detail of this relaxation procedure are being omitted for lack of space, but can be found in [12]. In the W-phase, due to the restrictions on the magnitude of the change in delay budgets computed in the D-phase, the magnitude of the change in x_i , i.e., δx_i will be small. To sum up, the W-phase finds a set of sizes for the transistors in the circuit that is a minimum area solution for satisfying the delay requirements calculated by the D-phase.

D. Putting it All Together

Having, defined the D-phase and W-phase of the optimization strategy, we are now in a position to finally describe *MINFLOTRANSIT*, our Min-cost Flow based Transistor sizing Tool.

MINFLOTRANSIT

1. Size the circuit to meet delay requirements using *TILOS*.
2. Iteratively perform alternately the D-phase and W-phase optimizations, solving the problems formulated in (10) and (11) respectively.
3. Stop the iterations when the area improvement after the W-phase is negligible.

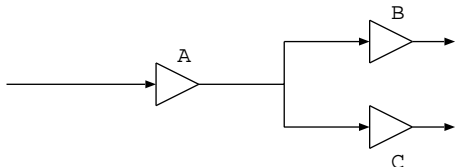


Fig. 6. An example illustrating the global perspective taken by *MINFLOTRANSIT* which *TILOS* tends to overlook.

We now present an example that qualitatively illustrates the improvements provided by *MINFLOTRANSIT* over *TILOS*.

Example II.1: Figure 6 shows a simple three gate circuit to be sized. *TILOS* is a sensitivity based greedy heuristic that proceeds by bumping up in each pass the size of that transistor/gate that leads to maximal benefit in speed for a unit increase in area. Such a transistor/gate is called the transistor/gate with the highest sensitivity. Assume that in figure 6, both *B* and *C* are gates with identical sensitivity and *A* has a lower sensitivity. Therefore the paths $A \rightarrow B$ and $A \rightarrow C$ are both critical. *TILOS*, due to its greedy nature, will bump up the sizes of *B* and *C* in alternate passes, whereas it should be intuitively clear that sizing up *A*, even though it has lower sensitivity may be the better option as it speeds up both paths $A \rightarrow B$ and $A \rightarrow C$ simultaneously. In the D-phase, *MINFLOTRANSIT* explicitly includes constraints to evaluate the benefits of altering the sizes of gates *A*, *B* and *C*. It is there-

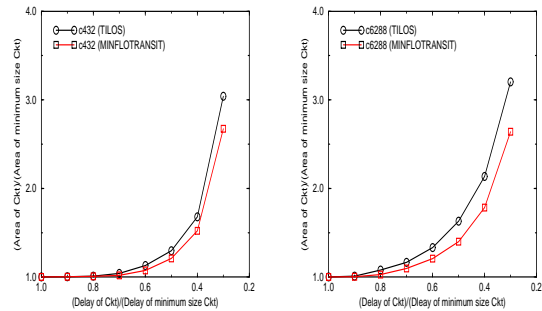


Fig. 7. Comparative area-delay curves for gate sizing of two ISCAS85 benchmark circuits. The total device area of the circuits after transistor sizing with *TILOS* and *MINFLOTRANSIT* is plotted against delay, normalized with respect to the delay of a minimum sized circuit. Even though the curves look close the area benefits are actually significant. For example in the case of *c6288*, for a circuit with 0.5 times the delay of the minimum sized circuit, the area savings of *MINFLOTRANSIT* over *TILOS* is 14.2%.

fore able to identify whether sizing gate *A*, in spite of its lower sensitivity, will be advantageous.

Theorem II.3: *MINFLOTRANSIT* produces minimum transistor sizing for any delay constraints.

Proof: Let us assume that we are in some intermediate iteration at a non-optimal point. We iteratively apply the D-phase, followed by the W-phase, and it is sufficient to show that the application of each of these steps causes the objective function to reduce, while maintaining feasibility. The D-phase uses a Taylor series approximation to the constraint in Equation (6) to represent the objective function entirely in terms of the delay variables. This approximation is valid within a radius of convergence of ϵ around the current point, \mathbf{X} , corresponding to some radius of δ around the vector of delays. Therefore, a solution to the D-phase is a valid solution to the original problem as long as the allowable delay change is bounded by a quantity that lies within a δ -ball of the current delays. This is achieved forcing $\text{MAX}\Delta D(i)$ and $\text{MIN}\Delta D(i)$ to be small for all i .

If the current solution is not optimal, due to the convexity of the problem [1], [2], there must be another feasible point in the neighborhood of the current point that has a smaller objective function value, and this point will be found by the D-phase. In the W-phase that follows, the solution found in the D-phase is a feasible solution. Since the W-phase does not limit the change in the delay or the transistor sizes as greatly as the D-phase, its solution must have an objective function value that is no larger than the solution of the D-phase.

Therefore, since the objective function value decreases in each phase, the procedure is guaranteed to find an optimal solution to the problem.

III. SIMULATION RESULTS

Simulation results were obtained on all the combinational circuits in the ISCAS85 benchmark suite, some combinational circuits from the *mcnc91* benchmark suite and also on ripple carry adders of 32-256 bits. The results shown in this section are for gate sizing which as mentioned before is a special case of true transistor sizing. We implemented the *TILOS* algorithm as described in [18]. Starting with all transistors at minimum size, a bumpsize of 1.1 was used for the initial sizing. Iterative application of D-phase and W-phase optimization was then carried out for optimal transistor sizing. Figure 7 shows the area delay curve for representative benchmark circuits for both *TILOS* and

MINFLOTRANSIT. The technology parameters used for simulation were obtained from [19] for 0.13 μ technology. As can be seen a clear gain in performance is seen when using MINFLOTRANSIT as opposed to using TILOS.

TABLE I

THE AREA SAVINGS IN % OF MINFLOTRANSIT OVER TILOS IS LISTED. THE CPU TIME REQUIRED BY TILOS AND THE EXTRA TIME REQUIRED BY MINFLOTRANSIT OVER AND ABOVE THAT REQUIRED BY TILOS ARE LISTED. THE CRITICAL PATH OF A MINIMUM SIZED CIRCUIT IS DENOTED BY D_{min} .

Circuit	# Gates	Area savings over TILOS	Delay Specs.	CPU (TIME) (TILOS)	CPU (TIME) (OURS)
adder32	480	$\leq 1\%$	$0.5D_{min}$	2.2s	5s
adder256	3840	$\leq 1\%$	$0.5D_{min}$	262s	608s
cm163a	65	2.1%	$0.55D_{min}$	0.13s	0.32s
cm162a	71	10.4%	$0.5D_{min}$	0.23s	0.96s
parity8	89	37%	$0.45D_{min}$	0.68s	2.15s
frg1	177	1.9%	$0.7D_{min}$	0.55s	1.49s
population	518	6.7%	$0.4D_{min}$	57s	179s
pmult8	1431	5%	$0.5D_{min}$	637	1476s
alu2	826	2.6%	$0.6D_{min}$	28s	71s
c432	160	9.4%	$0.4D_{min}$	0.5s	4.8s
c499	202	7.2%	$0.57D_{min}$	1.47s	11.26s
c880	383	4%	$0.4D_{min}$	2.7s	8.2s
c1355	546	9.5%	$0.4D_{min}$	29s	76s
c1908	880	4.6%	$0.4D_{min}$	36s	84s
c2670	1193	9.1%	$0.4D_{min}$	27s	69s
c3540	1669	7.7%	$0.4D_{min\ size}$	226s	651s
c5315	2307	2%	$0.4D_{min\ size}$	90s	201s
c6288	2416	16.5%	$0.4D_{min\ size}$	1677s	4138s
c7552	3512	3.3%	$0.4D_{min\ size}$	320s	683s

Table I lists the area savings of MINFLOTRANSIT over TILOS and the CPU times required in an Ultrasparc 10 Sun workstation for sizing the ISCAS85 benchmark circuits. The tabulated results are for sizing solutions where the area penalty is within 1.5 – 1.75 times that of a minimum sized circuit and all these correspond to points where the area penalty of sizing is reasonable. For the adders the improvement in area over TILOS is marginal thereby suggesting that adders can be easily sized by using greedy heuristics. This is not surprising since ripple carry adders have a single dominant critical path which can, possibly, be sized optimally using a heuristic like TILOS. On the other hand, however, as can be clearly seen for the ISCAS85 benchmark circuits and mnc benchmark circuits, the area savings vary from 1.9% – 37%. The overall time required by MINFLOTRANSIT is almost always (except in the small circuits c432, c499) within 2 – 4 times of TILOS. The circuit c6288 shows an unusually high time requirement for sizing, possibly due to the fact that this circuit (a type of multiplier) has a large number of paths, many of them reconvergent. Therefore, a number of competing paths can become critical at any instance and sizing this circuit is consequently harder. It is also notable that TILOS performs poorly as compared to MINFLOTRANSIT for this particular circuit. In all our simulations only a few tens of iterations were required by MINFLOTRANSIT (except the steepest portions of the area delay curve where no more than 100 iterations were required).

IV. CONCLUSIONS

We presented a new transistor sizing tool MINFLOTRANSIT that has shown impressive run-time behavior over various benchmarks in the ISCAS85 benchmark suite. MINFLOTRANSIT is a two-phase iterative-relaxation based technique. The first, D-phase has a worst case complexity of

$O(|V||E|\log(\log|V|))$, the second, W-phase has a worst case complexity of $O(|V||E|)$. The run-time behavior of this tool is comparable to TILOS but it is guaranteed to produce optimal transistor sizes for meeting the delay constraints. Although Elmore delay models were used in this paper for illustration, MINFLOTRANSIT is valid for a larger class of delay models characterized by the *monotonic decomposition* property in Definition II.2.

We further demonstrated approaches to extend the paper to more general classes of delay models with more parametric constraints, e.g., input rise time constraints.

Our work does not claim convergence to optimal solution for a general delay model but only convergence to a local optimum. Note that this is hardly more remarkable than a gradient descent based approach or an approach based on Lagrangian relaxation for piecewise differentiable functions. In fact we believe that we are only posing gradient descent as the D-phase. And by taking advantage of the structure of the sizing problem the D-phase is modeled as a minimum cost network flow optimization problem. Thus the primary modest contribution of this paper in our opinion is increasing the speed of gradient descent for transistor sizing optimization.

APPENDIX A

In this appendix we will show that the elements of the matrix $(D - A)^{-1}$ are always non-negative. This will in turn prove the critical assumption in the D-phase optimization that infinitesimal change in individual transistor sizes can be expressed as a negative combination of infinitesimal change in individual transistor delays.

It has already been shown that in the D-phase with infinitesimal shifts in the delay-budgets the following relation holds,

$$\delta \mathbf{x} = -(D - A)^{-1} \delta D \mathbf{x}, \quad (12)$$

where the terms have their usual meaning. So as long as the elements of the matrix $(D - A)^{-1}$ are non-negative, we are guaranteed that an infinitesimal change in the size of any transistor, say x_i , in the circuit is a negative combination of the delay components associated with various transistors in the circuit. This is a crucial observation because this facilitates the modeling of the D-phase as a minimum cost network flow problem. Recall that the optimization objective in the D-phase which is *minimize* $\sum_i x_i$ is equivalent to *maximize* $\sum_i cost_i d_i$.

It is essential that the cost coefficient associated with d_i , and hence with $(r(Dummy_i) - r(i))$ (given by $cost(i)$) in the objective function be non-negative. This can be guaranteed if all elements of $(D - A)^{-1}$ are non-negative. Non-negativity of $cost(i)$ is required in order for the artifact shown in Fig. 5 to correctly model the minimum of FSDUs available on all the fanout nets from a given node, i . As pointed out in section II-C.1 the delay component associated with any transistor has a very localized dependence on the sizes of transistors in its neighborhood. This observation implies that the matrix $(D - A)$ can be arranged as a block upper triangular matrix. Therefore, inverting $(D - A)$ amounts to inverting the individual blocks in a bottom-up manner (note that this has similar complexity as inverting a truly upper-triangular matrix). As a result as long as the inverse of the individual blocks can be guaranteed to have non-negative elements we can ensure that $(D - A)^{-1}$ has non-negative elements. So we will now try to establish that arbitrary sized individual blocks of the matrix $(D - A)$ indeed have this property. We will establish this fact with mathematical induction on the dimension of the square matrix $(D - A)$. We will first consider a 2×2 matrix to establish a base case.

Consider an inverter where, the size of the NMOS transistor is x_1 and that of the PMOS transistor is x_2 . In general these sizes will be independent of each other. For such an inverter driving a fixed load C_L , it can be shown that the delay component D_1 associated with the NMOS and the delay component D_2 associated with the PMOS are given by relations of the following form.

$$\begin{aligned} D_1 x_1 &= a_1^1 x_1 + a_2^1 x_2 + c_1, \\ D_2 x_2 &= a_1^2 x_1 + a_2^2 x_2 + c_2, \end{aligned} \quad (13)$$

where the terms $a_1^1, a_2^1, a_1^2, a_2^2$ are non-negative constants determined mainly by the drain capacitance of unit sized NMOS and PMOS transistors, the terms c_1, c_2 are non-negative constants dependent on the wiring capacitance and load capacitance driven by the inverter. Note that due

to non-negativity of all the terms in the relation for delay, $D_1 \geq a_1^1$ and $D_2 \geq a_2^2$. The matrix $(D - A)$ is given by,

$$\begin{bmatrix} (D_1 - a_1^1) & -a_1^2 \\ -a_1^2 & (D_2 - a_2^2) \end{bmatrix} \quad (14)$$

The inverse of this matrix is $(D - A)^{-1}$ given by,

$$\begin{bmatrix} \frac{D_2 - a_2^2}{\det(D-A)} & \frac{a_1^2}{\det(D-A)} \\ \frac{a_1^2}{\det(D-A)} & \frac{D_1 - a_1^1}{\det(D-A)} \end{bmatrix} \quad (15)$$

Now all the terms of $(D - A)^{-1}$ are non-negative provided $\det(D - A) > 0$. It will now be shown that $\det(D - A) = (D_1 - a_1^1)(D_2 - a_2^2) - a_1^2 a_2^1 > 0$. The delay component relations from (13) can be written as

$$\begin{aligned} (D_1 - a_1^1)x_1 &\geq a_2^1 x_2, \\ (D_2 - a_2^2)x_2 &\geq a_1^2 x_1, \end{aligned} \quad (16)$$

Multiplying the last two inequalities together (all terms are non-negative) we have,

$$(D_1 - a_1^1)(D_2 - a_2^2)x_1 x_2 \geq a_2^1 a_1^2 x_1 x_2. \quad (17)$$

Canceling out $x_1 x_2$ from both sides we get $\det(D - A) \geq 0$. We point out that the above discussion is valid as long as we have a relation of the form (16), this fact will be employed in Appendix C to extend the results in this section to general convex delay models.

To illustrate that a 3-dimensional matrix $D - A$ also results in a non-negative $(D - A)^{-1}$ we can use the facts we just established for 2-D matrices. This will be useful for us to establish a proof by mathematical induction that the components of $(D - A)^{-1}$ are always non-negative. Now consider an instance where three delay components are involved in the delay relations.

$$\begin{aligned} D_1 x_1 &= a_1^1 x_1 + a_2^1 x_2 + a_3^1 x_3 + c_1, \\ D_2 x_2 &= a_1^2 x_1 + a_2^2 x_2 + a_3^2 x_3 + c_2, \\ D_3 x_3 &= a_1^3 x_1 + a_2^3 x_2 + a_3^3 x_3 + c_3, \end{aligned} \quad (18)$$

We will show that all the co-factors of the matrix $(D - A)$ are non-negative which will also prove that the inverse of $(D - A)$ has to have all non-negative elements. The form of the matrix $(D - A)$ is,

$$\begin{bmatrix} (D_1 - a_1^1) & -a_1^2 & -a_1^3 \\ -a_1^2 & (D_2 - a_2^2) & -a_2^3 \\ -a_1^3 & -a_2^3 & (D_3 - a_3^3) \end{bmatrix} \quad (19)$$

The diagonal co-factors are of the form,

$$\begin{bmatrix} (D_1 - a_1^1) & -a_1^2 \\ -a_1^2 & (D_2 - a_2^2) \end{bmatrix}, \quad (20)$$

which has already been shown to be non-negative when we considered an inverter. The off-diagonal co-factors can be obtained with row, column permutations of a matrix of the following form $(C_{1,3})$,

$$\begin{bmatrix} -a_1^2 & (D_2 - a_2^2) \\ -a_1^3 & -a_2^3 \end{bmatrix}. \quad (21)$$

The determinant of this co-factor is $(a_1^2 a_2^3 + a_1^3 (D_2 - a_2^2))$ which is non-negative. The remaining co-factors are derived from this form with row/column permutations those requiring an odd number of permutations have an offsetting negative sign that makes them non-negative and those requiring an even number of permutations are evidently non-negative. This implies that all the co-factors are non-negative and hence all the elements of the inverse of $(D - A)$ are of the same sign as its determinant. So either the elements of $(D - A)^{-1}$ are all non-positive or they are all non-negative. The former is not possible because in that case the relation $X = (D - A)^{-1}c$ would not work with non-negative X and c , therefore the elements of $(D - A)^{-1}$ are non-negative.

Now we will use induction to prove that for arbitrary dimensions $(D - A)^{-1}$ with usual definitions has non-negative elements.

Inductive hypothesis: 1) The following $n \times n$ matrix, referred to as form I, has a non-negative determinant,

$$\begin{bmatrix} (D_1 - a_1^1) & \cdots & -a_1^n \\ \vdots & \ddots & \vdots \\ -a_1^n & \cdots & -a_n^n \end{bmatrix}. \quad (22)$$

2) For even(odd) n the following matrix, referred to as form II, has a non-positive(non-negative) determinant,

$$\begin{bmatrix} -a_1^2 & (D_2 - a_2^2) & \cdots & -a_n^2 \\ \vdots & \ddots & \vdots & \vdots \\ -a_1^n & \cdots & -a_{n-1}^n & -a_n^n \end{bmatrix}. \quad (23)$$

Inductive step: Assuming the above, we must prove that,

1) The following $(n+1) \times (n+1)$ matrix has a non-negative determinant,

$$\begin{bmatrix} (D_1 - a_1^1) & \cdots & -a_{n+1}^1 \\ \vdots & \ddots & \vdots \\ -a_1^{n+1} & \cdots & -a_{n+1}^{n+1} \end{bmatrix}. \quad (24)$$

2) For even (odd) $n + 1$, the following matrix has a non-negative(non-positive) determinant,

$$\begin{bmatrix} -a_1^2 & (D_2 - a_2^2) & \cdots & -a_{n+1}^2 \\ \vdots & \ddots & \vdots & \vdots \\ -a_1^{n+1} & \cdots & -a_n^{n+1} & -a_{n+1}^{n+1} \end{bmatrix}. \quad (25)$$

Proof: To prove 1) we compute the determinant by expanding the first row, for even $n + 1$ to obtain an expression of the form, $(D_1 - a_1^1) \times (\text{determinant}(\text{Form I})) - a_2^1 \times (\text{determinant}(\text{Form II})) + \cdots - a_{n+1}^1 \times (\text{determinant}(\text{Form II}))$. Which due to the non-positive value of a $n \times n$ determinant of Form II will have a non-negative value. Similarly, for odd $n + 1$ the determinant looks like, $(D_1 - a_1^1) \times (\text{determinant}(\text{Form I})) + a_2^1 \times (\text{determinant}(\text{Form II})) + \cdots + a_{n+1}^1 \times (\text{determinant}(\text{Form II}))$ which again turns out to be non-negative. Hence 1) is proved.

To prove 2) we again expand along the first row and for even (odd) $n + 1$ we can show that the resulting determinant is non-positive (non-negative) proof is similar to 1).

Hence we have proved that all co-factors of an arbitrary dimensioned matrix of the form $(D - A)$ are non-negative, it can be shown that the determinant of $(D - A)$ is always non-negative also which implies that all the elements of $(D - A)^{-1}$ are non-negative which implies that in the D-phase the objective function cost coefficients are all non-negative.

APPENDIX B: ON FEASIBILITY OF THE D-PHASE TO W-PHASE TRANSITION

As might be apparent, the two phases of the transistor sizing algorithm we describe are inter-dependent on each other in a very subtle manner. They are related precisely through the delay constraint inequalities of the D-phase. The relationship is actually hidden in the lower and upper bounds of these inequalities. The computation of these bounds require the introduction of a new parameter, a width increment, that is used in the simple monotonic functional of the delay formulas. An arbitrarily large value for this width increment may result in the generation of a D-phase solution that can not be feasibly realized in the W-phase. The choice of the width increments has therefore got to be judicious. The trick is to choose small enough values for the allowable delay perturbations around the ambient delay values for transistors in the D-phase such that the constraint region defined by the inequalities in the D-phase maps completely on to a feasible region in the solution space of the original sizing problem. Fortunately a certain aspect of the transistor sizing problem gives us a very practical solution to achieving this objective, which will be established shortly. First let us refer back to (11) and assume that the maximum change in the size of any transistor is δ on either side. If a transistor is already at the maximum (minimum) size it can have then we would assume that only an decrement (increment) in size for the transistor is allowed. Assuming for simplicity that all transistors have sizes which are at least δ under (over) the maximum (minimum) sizes allowed for the transistor, the maximum value of the delay given by the above expression can be obtained as,

$$\begin{aligned} \max \text{delay change}(i) &= \frac{\sum_{j \in S(V(G))} a_{ij}(x_j + \delta) + b_i}{x_i - \delta}, \\ \min \text{delay change}(i) &= \frac{\sum_{j \in S(V(G))} a_{ij}(x_j - \delta) + b_i}{x_i + \delta}, \\ \min \text{size} \leq x_i &\leq \max \text{size}. \end{aligned} \quad (26)$$

APPENDIX C

In this appendix, we will extend the results of this paper to more general delay models that admit the simple monotonic decomposition property in Definition II.2. First, the general transistor sizing problem will be discussed for convex delay models that admit the simple monotonic decomposition property in Definition II.2. Next, more general delay models that do not require convexity of the delay expression will be discussed.

C1: Convex Delay Models

With standard notations, assume any non-Elmore delay model that admits the simple monotonic decomposition property in Definition II.2, giving rise to the following transistor delay(component)-size relations.

$$\begin{aligned} \text{delay}(i) &= \frac{P_i(x_1, \dots, x_j, \dots)}{Q_i(x_i)}, \\ \text{i.e., } \text{delay}(i)Q_i(x_i) &= P_i(x_1, \dots, x_j, \dots), \end{aligned} \quad (27)$$

where $j \in S(V(G))$, P_i , Q_i are monotonic increasing (since $\frac{1}{Q_i}$ is monotonic decreasing). Now if we assume infinitesimal changes in delay and the sizes of the transistors, then we have

$$\begin{aligned} (\text{delay}(i) + \Delta \text{delay}(i))(Q_i(x_i) + \delta x_i Q_i'(x_i)) &= P_i(x_1, \dots, x_j, \dots) + \\ &\sum_{\text{all } j, j \neq i} \delta x_j P_i^j(x_1, \dots, x_j, \dots), \end{aligned} \quad (28)$$

where P_i^j is the partial derivative of P_i with respect to x_j . We point out that if $\text{delay}(i)$ were convex then the actual delay increment on the LHS would always be lower-bounded by the linearized increment on the RHS. This is because the curve corresponding to a convex function always lies above tangents at arbitrary points along the curve. Canceling out common terms in (28) and ignoring second-order terms we get,

$$\text{delay}(i)\delta x_i Q_i'(x_i) \approx -\Delta \text{delay}(i)Q_i(x_i) + \sum_{\text{all } j, j \neq i} \delta x_j P_i^j(x_1, \dots, x_j, \dots), \quad (29)$$

If we assume those cases where $\text{delay}(i)$ is a convex function of the transistor sizes, then due to the non-negativity of all the terms (except the first) in the RHS of (29) and due to the fact that LHS in (29) is always lower bounded by the RHS (due to convexity), the derivations in Appendix A will still be valid. The $LHS \geq RHS$ relation implied by convexity will establish a relation similar to (16). Thereafter, the non-negativity of the terms P_i^j is all that is required to establish the rest of the results in Appendix A. We can therefore use the proof in Appendix A to establish that an infinitesimal change in area of the circuit in the D-phase can be expressed as a non-negative combination of the infinitesimal changes in the delay component of transistors in the circuit.

C2: General Non-Convex Delay Models

In practice sometimes the sizing problem is treated as a single size parameter per gate problem. The transistor sizes within a gate are ganged together by fixed ratios and their collective size is represented by a single parameter. The fixed ratio is determined mainly to equalize rise and fall time delays and to make the task of layout easier and more predictable in order to preserve the gains from sizing using a coarse area model of the type used in (1). In such cases, when the delay model admits the simple monotonic decomposition given by Definition II.2, the gate size-delay relations can be expressed in a form similar to (27) and the relations in (28) and (29) still hold. But as already shown in section II-C the incremental delay and incremental size relations for gate sizing can be arranged in an upper-triangular form for strict gate sizing. We also point out the fact that in (29), Q_i' , P_i^j are all non-negative (due to the fact that Q_i , P_i are monotonic increasing functions). It is easy to see that by beginning from the primary output gates (bottom-most rows of the upper triangular delay/size matrix) and by performing a backward elimination process as in section II-C, the size increments for all the gates can be expressed as a non-positive combination of incremental gate delays. This implies that the incremental change in circuit area can be expressed as a non-positive combination of incremental gate delays. This in turn implies that the rest of the discussion in this paper holds and that a MINFLOTRANSIT like approach will converge for general delay models (not necessarily convex) that admit simple monotonic decomposition as given by Definition II.2.

APPENDIX D: GENERALIZING MINFLOTRANSIT

D1: General Delay Models and D-phase Feasibility

In general the Area-Delay relation for the transistor sizing problem can be succinctly represented in the following manner:

$$\mathbf{F}(X, D) = \mathbf{C}, \quad (30)$$

where \mathbf{F} is a matrix and \mathbf{X} , \mathbf{D} have their usual meanings, and \mathbf{C} is a constant vector. Using variational methods the following relationship can be obtained:

$$\begin{aligned} \Delta_x f(X, D)\Delta X + \Delta_D f(X, D)\Delta D &= 0, \\ \Delta X + H\Delta D &= 0, \end{aligned} \quad (31)$$

where $\mathbf{G} = \Delta_x f(X, D)$ and $\mathbf{H} = \Delta_D$ are matrices. We can therefore express either ΔX as a linear function of ΔD or vice-versa as follows:

$$\begin{aligned} \Delta X &= -(\mathbf{G})^{-1}\mathbf{H}\Delta D, \\ \Delta D &= -(\mathbf{H})^{-1}\mathbf{G}\Delta X. \end{aligned}$$

Now, as derived in Appendix A the class of delay models used in this paper have the property that $(\mathbf{G})^{-1}\mathbf{H}\Delta D$ is a non-negative matrix. Recall that this non-negativity property made it feasible to model the objective function of the D-phase as a maximization (minimization) problem with non-positive (non-negative) linear objective function cost coefficients. This coupled with fact that the constraints of the problem were 2-variable difference constraints, allowed us to use minimum cost network-flow techniques to model the D-phase efficiently. It turns out that as long as this non-negativity property holds, as will be shown subsequently, we can always do the following:

1. Given a feasible point in the (\mathbf{D}, \mathbf{X}) space, using a method based on back-substitution we can map a set of feasibility requirements on the \mathbf{X} and hence $\Delta \mathbf{X}$ variables of the form $Q \leq \Delta X \leq T$ to a set of feasibility constraints of the form $R \leq \Delta D \leq U$.
 2. Hence a set of feasibility constraints on the X variables, which is more lucid in the transistor sizing problem can be readily used to generate a set of constraints to be applied on the \mathbf{D} variables for the D-phase optimization problem.
 3. The result of the D-phase optimization as a result will never be infeasible, and hence MINFLOTRANSIT will never return infeasible results when applied to a sub-optimal point in the (\mathbf{D}, \mathbf{X}) space.
- Now to show that 1) above is true, consider the following set of constraints,

$$\begin{aligned} X_{min} &\leq X \leq X_{max}, \\ X_{min} &\leq X_0 + \Delta X \leq X_{max}, \\ X_{min} - X_0 &\leq \Delta X \leq X_{max} - X_0, \\ X_{min} - X_0 &\leq -(\mathbf{G})^{-1}\mathbf{H}\Delta D \leq X_{max} - X_0, \end{aligned}$$

where, X_{min} (X_{max}) is a vector of minimum (maximum) transistor sizes, X_0 is a vector of ambient transistor sizes, and the rest of the terms have their regular meaning. Now we can rearrange the above set of inequalities as follows,

$$X_0 - X_{max} \leq (\mathbf{G})^{-1}\mathbf{H}\Delta D \leq X_0 - X_{min}. \quad (32)$$

Now since $(\mathbf{G})^{-1}\mathbf{H}$ is a non-negative matrix it is possible to identify a feasible rectangular box in D-space that is completely contained within the feasible space for the above set of constraints. This is done by first solving for $X_0 - X_{max} = (\mathbf{G})^{-1}\mathbf{H}\Delta D$ (ΔD_{min}) and for $X_0 - X_{min} = (\mathbf{G})^{-1}\mathbf{H}\Delta D$ (ΔD_{max}). We can then show, due to non-negativity of $(\mathbf{G})^{-1}\mathbf{H}$, that the rectangular box (D-space) defined by $\Delta D_{min} \leq \Delta D \leq \Delta D_{max}$ is a subset of the original constraint region in the X-space. Note that $(\mathbf{G})^{-1}\mathbf{H}$ will in general be a block-triangular matrix as the infinitesimal size increments of a transistor can only affect the delay of the gate of which the transistor is a part or gates which are in a fanin path to that transistor. This property will render the process of inverting $(\mathbf{G})^{-1}\mathbf{H}$ a relatively easy task. In general as long as $(\mathbf{G})^{-1}\mathbf{H}$ is a non-negative matrix the following iterative relaxation procedure will always converge to a local optimum solution.

GENITERRELAX

1. Start with a initial sized circuit that meets the delay requirements.
2. Use the relation, $\Delta X = -\mathbf{G}^{-1}\Delta D$, to drive an optimization problem identical in form to the D-phase described in relation to MINFLOTRANSIT.
3. Compute the value of X' s which satisfy the result of the D-phase.
4. Iterate Steps 2 and 3 till convergence.

Note that the above technique is extremely general in nature and only requires $\mathbf{G}^{-1}\mathbf{H}$ to be a non-negative matrix, the Area-Delay relation can be static or possibly even dynamic (with some obvious restrictions on smoothness) where an oracle will be queried to compute instances of the \mathbf{G} and \mathbf{H} matrices during every iteration of GENITERRELAX. For the non-negativity restriction on $\mathbf{G}^{-1}\mathbf{H}$, all that is required is the following, if all the transistors were to be infinitesimally resized in such a fashion that none of them have a decrease in their delays (the transistors fed by primary input signals are left with a fixed size), then the delay model should be such that all the transistors necessarily have to either remain of the same size or go down in size. In other words if even one transistor were to increase in size then at least one transistors delay has to go down. This in no way is reflective of the Elmore delay model and is a reasonable restriction on a fairly large class of delay models.

One additional problem that remains is induction of additional parametric constraints (rise-time) etc, in to the sizing framework. This can be done as follows, assume that these additional constraints (on rise-time etc,) can be put in the following form,

$$P(X) \leq T, \quad (33)$$

using a similar approach as for the area-delay dependency function we show the following,

$$P(X_0) + \Delta_X P \cdot \Delta X \leq T, \Delta_X P \cdot \Delta X \leq T - P(X_0). \quad (34)$$

As long as $\Delta_X P$ can be constrained to a non-negative (non-positive) matrix we can express the above constraints as a set of simple constraints on the ΔX vector which can be translated as before on to a set of simple constraints on the ΔD vector. In this manner we will be able to include rise-time and other similar constraints in to the sizing framework.

D2: The Area Constrained Transistor Sizing Problem

It turns out that many a time Area and not Performance is the guiding impetus in the transistor sizing problem. In such cases a hard-limit is put on the area of the circuit, say A_H , and the designer's job is to find the fastest circuit that meets the area constraint.

The techniques developed in this paper can be easily adapted to solve this problem in polynomial time.

1. Use MINFLOTRANSIT to design the fastest possible circuit with no hard limited area constraints, except the constraint on individual transistor sizes. Set the delay of the resulting circuit to be D_u , set $D_l = 0$. Set the area to A_c . If $A_c < A_H$ output the circuit as the solution else perform 2).

2. If $A_c \approx A_H$, then output the circuit as the solution else If $A_c \downarrow A_H$, then set $D_u = (D_u + D_l)/2$ and perform 3) else if $A_c < A_H$ then set $D_l = (D_u + D_l)/2$ and perform 3).

3. Perform MINFLOTRANSIT once more but this time use a critical path requirement of $(D_u + D_l)/2$, this results in a D -phase where most of the transistors will have a non-negative slack. Set the resulting area to A_c and return to step 2).

The above technique will converge after at most $\log_2(D_{max})$ iterations, where D_{max} is the delay of the fastest circuit sized most aggressively.

D3: What if the Initial Sizer Fails ?

Sometimes the delay requirements are such that no greedy initial sizer can even return a feasible circuit. In such instances, as stated presently MINFLOTRANSIT does not have an initial starting point. The failure of the initial sizer can be due to a couple of reasons:

1. The delay requirements are unrealistic or even infeasible.
2. The initial sizer is restrictive.

In such cases we can do the following:

1. Size the circuit for the fastest circuit possible using the initial sizer.
2. Now operate MINFLOTRANSIT in a creep mode, whereby in each iteration the D-phase operates with a overall negative slack requirement. This is done by connecting all the Primary output gates to a Dummy output node and keeping a infinitesimally negative slack requirement on all the nets incident on the Dummy output node.
3. Continue performing MINFLOTRANSIT till either i) the circuit achieves the required delay or MINFLOTRANSIT fails due to infeasibility.

In case MINFLOTRANSIT fails then due to the result in the last section we are assured that the initial delay requirement was infeasible.

REFERENCES

- [1] J. P. Fishburn and A. E. Dunlop, "TILOS: A Polynomial Programming Approach to Transistor Sizing," in *Proceedings of the 1985 International Conference on Computer-Aided Design*, November 1985, pp. 326-328.
- [2] S.S Sapatnekar, V. B. Rao, P. M. Vaidya, and S. M. Kang, "An Exact Solution to the Transistor Sizing Problem for CMOS Circuits Using Convex Optimization," *IEEE Transactions on Computer-Aided Design*, vol. 12, no. 11, pp. 1621-1634, November 1993.
- [3] J.-M. Shyu, A. L. Sangiovanni-Vincentelli, J.P. Fishburn, and A.E. Dunlop, "Optimization-based Transistor Sizing," *IEEE Journal on Solid State Circuits*, vol. 23, no. 2, pp. 400-409, April 1988.
- [4] D. P. Marple, "Performance Optimization of Digital VLSI Circuits," *Technical Report CSL-TR-86-308, Stanford University*, October 1986.
- [5] D. P. Marple, "Transistor Size Optimization in the Tailor Layout System," in *Proceedings of the 26th ACM/IEEE Design Automation Conference*, June 1989, pp. 43-48.
- [6] H. Y. Chen and S. M. Kang, "icoach: A circuit optimization aid for cmos high-performance circuits," *Integration, the VLSI Journal*, vol. 10, no. 2, pp. 185-212, January 1991.
- [7] Z. Dai and K. Asada, "MOSIZ: A Two-Step Transistor Sizing Algorithm based on Optimal Timing Assignment Method for Multi-Stage Complex Gates," in *Proceedings of the 1989 Custom Integrated Circuits Conference*, May 1989, pp. 17.3.1-17.3.4.
- [8] C.-P. Chen, C. N. Chu, and D. F. Wong, "Fast and Exact Simultaneous Gate and Wire Sizing by Lagrangian Relaxation," in *Proceedings of the 1998*

IEEE/ACM International Conference on Computer-Aided Design, November 1998, pp. 617-624.

- [9] A. R. Conn, P. K. Coulman, R. A. Haring, G. L. Morrill, C. Visweshwariah, and C. W. Wu, "JiffyTune: Circuit Optimization Using Time-Domain Sensitivities," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 17, no. 12, pp. 1292-1309, December 1998.
- [10] J. Grodstein, E. Lehman, H. Harkness, B. Grundmann, and Y. Watanabe, "A delay model for logic synthesis of continuously sized networks," in *Proceedings of the 1995 International Conference on Computer-Aided Design*, November 1995, pp. 458-462.
- [11] A. V. Goldberg, M. D. Grigoriadis, and R. E. Tarjan, "Use of Dynamic Trees in a Network Simplex Algorithm for the Maximum Flow Problem," *Mathematical Programming*, vol. 50, no. 3, pp. 277-290, June 1991.
- [12] M. C. Papaefthymiou, "Asymptotically Efficient Retiming under Setup and Hold Constraints," *Proceedings of the IEEE/ACM International Conference on Computer-Aided Design*, pp. 288-295, Nov. 1998.
- [13] T. H. Cormen, C. E. Leiserson, and R. L. Rivest, *Introduction to algorithms*, McGraw-Hill, New York, NY, 1990.
- [14] G. Strang, *Linear Algebra and its Applications*, Harcourt Brace Jovanovich, Publishers, San Diego, CA, 1988.
- [15] V. Sundararajan and K. K. Parhi, "Low Power Gate Resizing Using Buffer-Redistribution," in *Proceedings of the Twentieth Anniversary Conference on Advanced Research in VLSI*, March 1999, pp. 170-184.
- [16] C. E. Leiserson and J. B. Saxe, "Optimizing Synchronous Systems," *Journal of VLSI and Computer Systems*, vol. 1, no. 1, pp. 11-67, 1983.
- [17] V. Chvatal, *Linear Programming*, W. H. Freeman and Company, New York, NY, 1983.
- [18] A. E. Dunlop, J. P. Fishburn, D. D. Hill, and D. D. Shugard, "Experiments Using Automatic Physical Design Techniques for Optimizing Circuit Performance," in *Proceedings of the 32nd Midwest Symposium on Circuits and Systems*, Urbana, IL, August 1989, pp. 216-220.
- [19] J. Cong, "Challenges and Opportunities for Design Innovations in Nanometer Technologies," *Technical Report, Semiconductor Research Corporation*, 1997.