

Delay Bound Based CMOS Gate Sizing Technique

Alexandre Verle, Xavier Michel, Philippe Maurine, Nadine Azemard, Daniel

Auvergne

▶ To cite this version:

Alexandre Verle, Xavier Michel, Philippe Maurine, Nadine Azemard, Daniel Auvergne. Delay Bound Based CMOS Gate Sizing Technique. ISCAS: International Symposium on Circuits and Systems, May 2004, Vancouver, BC, Canada. pp.189-192, 10.1109/ISCAS.2004.1329494. lirmm-00108856

HAL Id: lirmm-00108856 https://hal-lirmm.ccsd.cnrs.fr/lirmm-00108856v1

Submitted on 11 Sep 2019 $\,$

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

DELAY BOUND BASED CMOS GATE SIZING TECHNIQUE

A. Verle, X. Michel, P. Maurine, N. Azémard, D. Auvergne

LIRMM, UMR CNRS/Université de Montpellier II, (C5506),

161 rue Ada, 34392 Montpellier, France

ABSTRACT

In this paper we address the problem of delay constraint distribution on CMOS combinatorial paths. We first define a way to determine on any path the reasonable bounds of delay characterizing the structure. Then we define two constraint distribution methods that we compare to the equal delay distribution and to an industrial tool based on Newton-Raphson like algorithms. Validation is obtained on a $0.25\mu m$ process by comparing the different constraint distribution techniques on various benchmarks.

1. INTRODUCTION

The goal of gate sizing is to determine optimum sizes for the gates in order that the path delay respects the constraints with the minimum area/power cost. Another parameter to consider is the feasibility of the constraint imposed on a given path. For that, indicators must be found to explore the design space and to select one solution among the available optimization alternatives such as sizing, buffering or technology remapping. The target of this paper is twofold: defining the delay bounds of a given path and determining a way for distributing a delay constraint on this path with the minimum area/power cost.

The problem of transistor sizing has been widely investigated using non linear programming techniques [1] or heuristics based on simple delay models [2]. Recently, in a pedagogical application [3] of the τ model, Sutherland [4], describing the gate delay as the product of electrical and logical efforts, proposed to minimize the delay on a path by imposing an equal effort that is a constant delay on all the elements of the path.

This way to select cell sizes can be proven mathematically exact [3] for a fanout-free path constituted of ideal gates (without parasitic capacitance or divergence branch). However this evenly budget distribution is far to be the optimal one with respect to delay and area for a real path, on which divergence branches and routing capacitance are not negligible. Starting from the definition of the design space in terms of minimum and maximum delay permissible on a given path, we propose in this paper a design space exploration method allowing an area/power efficient distribution of constraint on a combinatorial path.

The delay bound determination and the constraint distribution method are based on a realistic delay model [5] that is input slope dependent and able to distinguish

between falling and rising signals. This model is shortly presented in part 2. In part 3 we give a method for defining delay bounds on a path. Different approaches for distributing a delay constraint are considered in part 4 and compared in part 5 on different benchmarks of increasing complexity. We finally conclude in part 6.

2. GATE DELAY MODELING

As previously mentioned sizing at the physical level imposes to use a realistic delay computation that must consider a finite value of the gate input transition time. As developed in [5] we introduce the input slope effect and the related input-to-output coupling in the model as:

$$t_{\rm HL,LH}(i) = \frac{v_{\rm TN,P}}{2} \tau_{\rm INLH,LH}(i-1) + (1 + \frac{2C_{\rm M}}{C_{\rm M} + C_{\rm L}}) t_{\rm HL,LHstep}(i)$$
(1)

where $v_{TN,P}$ are the reduced value (V_T/V_{DD}) of the threshold voltage of the N,P transistors. $\tau_{INHL,LH}$ is the duration time of the input signal, taken to be twice the value of the step response of the controlling gate. C_M is the coupling capacitance between the input and output nodes. C_L is the output loading capacitance. Indexes (i), (i-1) specify the switching and the controlling gates, respectively.

Following [2], the step response of each edge is defined by the time interval necessary to load (unload) the gate output capacitance under the maximum current, I_{MAX} , available in the structure:

$$t_{\rm HL,LHstep} = \frac{C_{\rm L} \cdot \Delta V}{I_{MAX}}$$
(2)

Following the elegant model of [3] the evaluation of this step response on logic gates supplies a general expression given by:

$$t_{\rm HL,LHstep} = \tau \cdot S_{\rm HL,LH} \cdot \frac{C_{\rm L}}{C_{\rm IN}} .$$
(3)

where τ is a time unit characterizing the process. C_{IN}, the gate input capacitance, is defined in terms of the P/N width ratio k. For simplicity, the S factors (logical effort of [3]) include all the current capability difference between the pull up (pull down) transistor equivalent to the corresponding serial array. These factors are configuration ratio dependent and characterize for each edge, the ratio of current available in an inverter and a gate of identical size.

Then considering an array of gates, the delay path can easily be obtained from (1) and (3) as a technology independent posynomial representation:

$$\frac{t_{\text{HL,LH}}}{\tau} = \theta = S'_{i} \cdot \frac{C_2 + C_{\text{Pl}}}{C_1} + \dots + S'_{i-i} \cdot \frac{C_i + C_{\text{Pl-1}}}{C_{i-1}} + \dots + S'_{n} \cdot \frac{C_{\text{L}}}{C_{n}}$$
(4)

where the S_i ' include the logical effort and the input ramp effect, C_i represents the input capacitance of the gate and C_{pi} the output node total parasitic capacitance, including the interconnect and branching load.

3. DELAY BOUND DEFINITION

We consider realistic combinatorial paths on which two parameters are known and imposed:

- the output load capacitance of the last gate, that is determined by the input capacitance of the output register,

- the input capacitance of the first gate imposed by the loading conditions of the input register.

In that condition the path delay is bounded. These bounds can be determined, considering that the delay of a path (4) is a convex function of the gate input capacitance. This is illustrated in Fig.1 that gives the variation of the path delay with respect to the transistor sizing of a combinatorial path constituted of 13 gates. Note that the slope of the curve corresponds to the sensitivity of the path delay to the transistor sizing.

As shown the delay value decreases from a maximum value down to a minimum value that will be determined below. The maximum delay has been obtained imposing all the transistor sizes at the minimum allowed by the technology. This maximum value is a "reasonable" one but not the absolute maximum value. It is always possible to get a much greater value by loading minimum gates with an infinitely sized driver. This curve illustrates what we define by exploring the design space:

- near the maximum value, Θ_{Max} , of the delay the path sensitivity to the gate sizing is very important, a small variation of the gate input capacitance results in a large change in delay,

- at the contrary near the minimum Θ_{Min} the sensitivity is becoming very low and in that range any delay improvement is highly area/power expensive.

Evaluating the feasibility of a delay constraint Θ_c imposes to compare its value to the preceding bounds. If the Θ_c value is closed to the maximum Θ_{Max} the constraint satisfaction will be obtained at reasonable cost by transistor sizing otherwise it would be more profitable to reconfigure the logic or to insert buffers [6]. Let us define these bounds. As previously mentioned we consider for Θ_{Max} the "reasonable" value obtained when all the gates are implemented with transistors of minimum size. For the minimum bounds we just use the posynomial property [7] of (6). Canceling the derivatives of (4) with respect to the gate input capacitances C_i we obtain a set of linked equations such as:

$$S'_{i-1} \cdot \frac{C_i}{C_{i-1}} - S'_i \cdot \frac{C_{i+1} + C_{P_i}}{C_i} = 0$$

$$S'_i \cdot \frac{C_{i+1}}{C_i} - S'_{i+1} \cdot \frac{C_{i+2} + C_{P_{i+1}}}{C_{i+1}} = 0...$$
(5)

Cell sizes can then be selected to match the minimum delay, by visiting all the gates in a topological order, starting from the output, such as:



Fig.1. Illustration of the variation of the path delay with the gate sizing.

This results in a set of *n* linked equations that can be easily solved by iterations from an initial solution that considers C_{i-1} known and equal to a reference value C_{REF} . This reference can be set equal to the minimum value available (C_{MIN}) in the library or to any other one.



Fig.2. Illustration of the research of minimum delay on an array of ten gates for different values of the initial reference capacitance; the output load of each gate on the array is given in unit of C_{MIN} .

In Fig.2, we illustrate the variation of the calculation convergence with a choice of C_{REF} . As shown, whatever is the value of C_{REF} (C_{MIN} to $100C_{MIN}$), we always obtain a fast convergence to the minimum.

4. CONSTRAINT DISTRIBUTION

Determining the possible bounds of delay for a given path topology, the next step is to evaluate the feasibility of a constraint to be imposed on a path. The theory of constant effort or constant delay [4,8] provides an easy way to select the cell size for each stage but for real configuration it is far from the optimum and often results in oversized structures. For that we propose two techniques for the gate size selection in order to satisfy a constraint that we will compare in the next part to the constant delay method.

To define the first method we consider that imposing equal delay to the gates with an important value of the logical effort (S'_i), results in an important over sizing of these complex gates. The determination of the lowest delay bound directly provides the optimal delay distribution on the path that appears to be the fastest one. So we can use this distribution to define for each gate a weight or gain Θ_i relative to this distribution $\Theta_{Min}=\Sigma\Theta_{Mini}$. In that case we propose to distribute the delay constraint Θ_c on a path using a weight defined with respect to the minimum delay distribution as:

$$\theta_{i} = \frac{\theta_{Mini} \cdot \theta_{c}}{\sum_{i} \theta_{Mini}}$$
(7)

This guarantees the conservation of the path delay distribution, obtained at the optimal solution, for any value of the constraint. Then processing backward from the output of the path, this directly gives, for each gate, the value of Θ_i that determines from (1,3) the size of the corresponding gate.

The second method of equal sensitivity is directly deduced from (5). Instead to search for the minimum we impose the same path delay sensitivity to the sizing, by solving:

$$S_{i-1} \cdot \frac{1}{C_{i-1}} - S_i \cdot \frac{C_{i+1} + C_{P_i}}{C_i^2} = a \dots$$
(8)

where "a" is a constant, representing the slope of the curve of Fig.3, that represents the variation of the delay between the bounds previously defined. Following the procedure used for the first method, the size of the gates is obtained from the iterated solution of (8) using as initial solution the sizing for the maximum delay value (all gates sized at C_{REF}). The different points on the curve of Fig.3 have been obtained from (8), by varying the value of "a", until "a" = 0 to get the minimum.

As expected, for a given value of the sensitivity factor "a", this curve represents the locus of the minimum delay solutions. No inferior solution can be found. Thus, varying the "a" value gives the possibility to explore the design space and to determine the minimum area sizing condition satisfying the delay constraint.

In Fig.3, we compare our approach to an industrial optimization tool (Amps from Synopsys). As shown the two methods give nearly equivalent design range exploration, however the equal sensitivity method results in a minimum area solution.



5. VALIDATION

In order to validate these sizing and constraint distribution techniques we compare on different benchmarks the minimum delay value and the area obtained using the three investigated methods:

- equal distribution of delays (C), $(\Theta_i = \Theta_c/n)$, [4] where *n* is the number of stages,

- weighted distribution (7), (B),

- equal gate sensitivity (8), (A),

- and using an industrial tool based on a Newton-Raphson based algorithm (D), [9] (Amps from Synopsys).

These benchmarks are constituted of array of gates (Nand, NOR, 2 and 3 inputs) with different loading conditions.

The comparison of the minimum delay values obtained with each technique is given in Table 1 for different paths. The targeted process is the STM 0.25 μ m with τ = 7.05ps. As shown the lowest minimum value of delay is obtained with both the weighted and the equal sensitivity techniques.

This ascertains the method used to determine the lowest bound of delay on a logical path. Note that around the minimum value of delay the area penalty is, of course, very large. This value of delay must be more considered as an indicator for the feasibility of the constraint than as a design target.

The next step is to compare for a given delay constraint the area of implementation obtained with the different distribution techniques. For that we impose on the different benchmarks a delay constraint defined between the bounds previously defined. Then we compare in Table 2 the area corresponding to the gate sizing allowing, with the different techniques, to match the constraint. We can observe that if for a weak constraint value the different techniques appear quite equivalent, for a tighter constraint the equal sensitivity distribution technique (A) allows a match with a much smaller area than the others. Note that all the values given in Tables 1 and 2 are obtained from Spice simulations (MM9 model) of the different benchmarks.

The weighted distribution (B) still results in a quite equivalent area but the equal delay distribution (C) and Amps (D) may result for quite complex paths in an

important increase of area. For some constraint values they may fail to get a solution. Note that the equal sensitivity method is mathematically quasi-optimal and always gives slightly better results than the weighted distribution. However this last method, defined from a minimum delay solution obtained for a sensitivity value equal to zero, can be much more easily implemented.

Gate Nb	Θ_{MAX} (ps)	Area µm	Siz. Tech.	Θ_{MIN} (ps)	Area µm
			A	620	987
			В	620	987
9	1874	42	С	676	391
			D	633	632
			Α	698	1448
			В	698	1448
11	2085	46	С	777	440
			D	937	348
			Α	923	4337
		ĺ	В	923	4337
15	3479	3479	С	1023	1083
			D	960	3067
			Α	1192	8419
			В	1192	8419
21	4583	94	С	1484	1039
			D	1693	1152
			A	1503	21578
			В	1503	21578
31	6560	138	С	1881	2226
			D	2426	1826

Table 1

An illustration of these results is given in Fig.4 where we show for the path constituted of 31 gates, the complete exploration of the design space using the preceding constraint distribution methods. As shown for a delay constraint smaller than $\Theta_{max}/2$ the gain in area (power) using the equal sensitivity or the weighted distribution method is quite significant.





6. CONCLUSION

Based on a simple realistic delay model for gates, we have first determined an easy way to characterize the feasibility of a delay constraint imposed on a combinatorial path. We have defined reasonable maximum and real minimum delay bounds. Then we proposed two techniques to match a delay constraint on a path: the equal sensitivity and the weighted method that is a budgeting method. We have applied these methods on different benchmarks with various constraint conditions and compared the resulting implementation area with that obtained from an equal delay distribution and with an industrial tool. If for weak constraints the different methods are quite equivalent, for values near the minimum, the proposed methods always find a solution and result in an important area/power saving. Another point to be clarified further is to define at which distance of the minimum delay value it is area/power efficient to impose a constraint.

7. REFERENCES

 J. M. Shyu &Al, A. Dunlop, "Optimization-based transistor sizing" IEEE J. Solid State Circuits, vol.23, n°2, pp.400-409, 1988.

[2] J. Fishburn, A. Dunlop, "TILOS: a posynomial programming approach to transistor sizing" in Proc. Design Automation Conf. 1985,pp.326-328.

[3] C. Mead, M. Rem, "Minimum propagation delays in VLSI", ", IEEE J. Solid State Circuits, vol.SC17, n°4, pp.773-775, 1982.

[4] I. Sutherland, B. Sproull, D. Harris, "Logical Effort: Designing Fast CMOS Circuits", Morgan Kaufmann Publishers, INC., California, 1999.

[5] K. O. Jeppson, "Modeling the influence of the transistor gain ratio and the input-to-output coupling capacitance on the CMOS inverter delay", IEEE J. Solid State Circuits, vol.29, pp.646-654, 1994.

[6] S. Chakraborty, R. Murgai "Layout driven timing optimization by generalized DeMorgan transform" IWLS 2001, pp.53-59.

[7] M. Ketkar, K. Kasamsetty, S. Sapatnekar "Convex delay models for transistor sizing" Proc. of the 2000 Design Automation Conf. pp.655-660.
[8] J. Grodstein & al"A delay model for logic synthesis of continuously-sized networks", ICCAD 95, Nov 95.

[9] R. K. Brayton, R. Spence "Sensitivity and Optimization" Elsevier 1980

1 able 2									
Gate Nb.	Siz. Tech.	$\frac{\Theta_{MAX}}{\Theta_{MIN}}$	$\frac{\Theta_{C}}{\Theta_{MIN}}$	Area (µm)	$\frac{\Theta_{C}}{\Theta_{MIN}}$	Area µm			
9	A B C D	3	1.4	137 147 161 144	1.02	535 560 Fail 632			
11	A B C D	11	2.15	66 80 94 70	1.1	310 330 440 Fail			
15	A B C D	3.8	1.4	302 310 410 324	1.04	1333 1407 Fail 3067			
21	A B C D	3.8	2.1	196 214 230 198	1.31	553 558 715 1152			
31	A B C D	4.4	3.1	364 400 427 377	1.26	1275 1361 1970 Fail			