

Power – Performance Optimization for Custom Digital Circuits

Radu Zlatanovici and Borivoje Nikolić

University of California, Berkeley, CA 94720
{zradu,bora}@eecs.berkeley.edu

Abstract. This paper presents a modular optimization framework for custom digital circuits in the power – performance space. The method uses a static timer and a nonlinear optimizer to maximize the performance of digital circuits within a limited power budget by tuning various variables such as gate sizes, supply, and threshold voltages. It can employ different models to characterize the components. Analytical models usually lead to convex optimization problems where the optimality of the results is guaranteed. Tabulated models or an arbitrary timing signoff tool can be used if better accuracy is desired and although the optimality of the results cannot be guaranteed, it can be verified against a near-optimality boundary. The optimization examples are presented on 64-bit carry-lookahead adders. By achieving the power optimality of the underlying circuit fabric, this framework can be used by logic designers and system architects to make optimal decisions at the microarchitecture level.

1 Introduction

Integrated circuit design has seamlessly entered the power-limited scaling regime, where the traditional goal of achieving the highest performance has been displaced by optimization for both performance and power. Solving this optimization problem is a challenging task due to a combination of discrete and continuous constraints and the difficulty in incorporating costs for both energy and delay in the objective functions.

System designers typically bypass this problem by forming a hybrid metric, such as MIPS/mW, for evaluating candidate microarchitectures. Similarly, designs at the circuit level are evaluated based on metrics that combine energy and delay, such as the energy-delay product (EDP). A circuit designed to have the minimum EDP, however, may not be achieving the desired performance or could be exceeding the given energy budget. As a consequence, a number of alternate optimization metrics have been used that generally attempt to minimize the $E^m D^n$ product [1]. By choosing parameters n and m a desired tradeoff between energy and delay can be achieved, but the result is difficult to propagate to higher layers of design abstraction.

In contrast, a more systematic and general solution to this problem minimizes the delay for a given energy constraint [2]. Note that a dual problem to this one, minimization of the energy subject to a delay constraint yields the same solution.

Custom datapaths are an example of power-constrained designs where the designers traditionally iterate in sizing between schematics and layouts. The initial design is sized using the wireload estimates and is iterated through the layout phase until a set delay goal is achieved. The sizing is refined manually using the updated wireload estimates. Finally, after minimizing the delay of critical paths, the non-critical paths

are balanced to attempt to save some power, or in case of domino logic to adjust the timing of fast paths. This is a tedious and often lengthy process that relies on the designer's experience and has no proof of achieving optimality. Furthermore, the optimal sizing depends on the chosen supply and transistor thresholds. An optimal design would be able to minimize the delay under power constraints by choosing supply and threshold voltages, gate sizes or individual transistor sizes, logic style (static, domino, pass-gate), block topology, degree of parallelism, pipeline depth, layout style, wire widths, etc.

Custom circuit optimization under constraints has been automated in the past. IBM's EinsTuner [3] uses a static timing formulation and tunes transistor sizes for minimal delay under total transistor width constraints. It uses simulation instead of modeling for best accuracy, but it only guarantees local optimality. TILOS [4] solves a convex optimization problem that results from the use of Elmore's formula for gate delays. While the models are rather inaccurate due to their simplicity, the result is guaranteed to be globally optimal.

This paper builds on similar ideas and presents a modular design optimization framework for custom digital circuits in the power – performance space that:

- Formulates the design as a mathematical optimization problem;
- Uses a static timer to perform all circuit-related computations, thus relieving the designer from the burden of providing input patterns;
- Uses a mathematical optimizer to solve the optimization problem numerically;
- Adjusts various design variables at different levels of abstraction;
- Can employ different models in the timer in order to balance accuracy and convergence speed;
- Handles various logic families (static, dynamic, pass-gate) due to the flexibility of the modeling step;
- Guarantees the global optimality of the solution for certain families of analytical models that result in the optimization problem being convex;
- Verifies a near-optimality condition if optimality cannot be guaranteed.

2 Design Optimization Framework

The framework is built around a versatile optimization core consisting of a static timer in the loop of a mathematical optimizer, as shown in Fig. 1.

The optimizer passes a set of specified design variables to the timer and gets the resulting cycle time (as a measure of performance) and power of the circuit, as well as other quantities of interest such as signal slopes, capacitive loads and, if needed, design variable gradients. The process is repeated until it converges to the optimal values of the design parameters, as specified by the desired optimization goal. The circuit is defined using a SPICE-like netlist and the static timer employs user-specified models in order to compute delays, cycle times, power, signal slopes etc.

Since the static timer is in the main speed-critical optimization loop, it is implemented in C++ to accelerate computation. It is based on the conventional longest path algorithm. The current timer does not account for false paths or simultaneous arrivals, but it can be easily substituted with a more sophisticated one because of the modularity of the optimization framework.

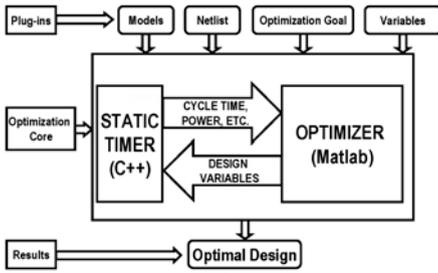


Fig. 1. Design optimization framework

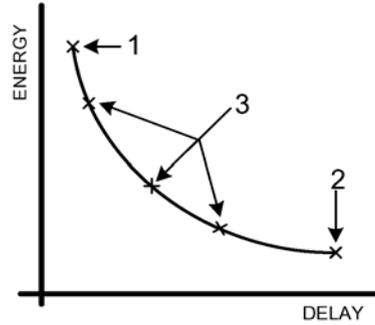


Fig. 2. Typical optimal energy – delay tradeoff curve for a combinational circuit

Adjust GATE SIZES in order to Minimize DELAY subject to:

Maximum ENERGY PER TRANSITION, Maximum internal slopes, Maximum output slopes, Maximum input capacitances, Minimum gate sizes

Additional constraints on signal slopes and minimum gate sizes are inserted in order to ensure manufacturability and correct circuit operation. By solving this optimization problem for different values of the energy constraint, the optimal energy-delay tradeoff curve for that circuit is obtained, as shown in Fig. 2.

The optimal tradeoff curve has two well defined end-points: point 1 represents the fastest circuit that can be designed; point 2 represents the lowest power circuit, mainly limited by minimum gate sizes and signal slope constraints. The points in-between the two extremes (marked “3” on the graph) correspond to minimizing various $E^m D^n$ design goals (such as the well known energy – delay product, EDP).

3 Models

The choice of models in the static timer greatly influences the convergence speed and robustness of the optimizer. Analytical or tabulated models can be used in the optimization framework, depending on the desired accuracy and speed targets. Table 1 shows a comparison between the two main choices of models.

3.1 Analytical Models

In our initial optimizations we use a simple, yet fairly accurate analytical model. This model allows for a convex formulation of the resulting optimization problem, where the gate sizes are the optimization variables. The model has three components: a delay equation (1), a signal slope equation (2), and an energy equation (3):

$$t_D = p + g \frac{C_{load}}{C_{in}} + \eta \cdot t_{slope_in} \tag{1}$$

$$t_{slope_out} = \lambda + \mu \frac{C_{load}}{C_{in}} + \nu \cdot t_{slope_in} \tag{2}$$

$$E = \sum_{all_nodes} \alpha_i C_i V_{DD}^2 + T_{cycle} \sum_{all_gates} W_j P_{leak,j} \tag{3}$$

Equation (1) is an extension of the simple linear model used in the method of logical effort [5], or the level-1 model with limited accuracy in commercial logic synthesis tools[6]. Equations (1) and (2) are a straightforward first order extension to these models that accounts for signal slopes.

The capacitance of a node is computed using (4):

$$C_{node} = \sum_{gate_inputs_at_node} k_i W_i + C_{wire} \quad (4)$$

where W_i are the corresponding gate sizes.

Each input of each gate is characterized for each transition by a set of seven parameters: p, g, η for the delay, λ, μ, ν for the slope and k for the capacitance. Each gate is also characterized by an average leakage power P_{leak} measured when its relative size is $W=1$. Each node of the circuit has an activity factor α , which is computed through logic simulation for a set of representative input patterns.

All the above equations can be written as posynomials in the gate sizes, W_i :

$$t_D = p + g \frac{\sum k_i W_i + C_{wire}}{kW_{current}} + \eta \cdot t_{slope_in} \quad (5)$$

$$t_{slope_out} = \lambda + \mu \frac{\sum k_i W_i + C_{wire}}{kW_{current}} + \nu \cdot t_{slope_in} \quad (6)$$

If t_{slope_in} is a posynomial, then t_D and t_{slope_out} are also posynomials in W_i . By specifying fixed signal slopes at the primary inputs of the circuit, the resulting slopes and arrival times at all the nodes will also be posynomials in W_i . The maximum delay across all paths in the circuit will be the maximum of several posynomials, hence a generalized posynomial. A function f is a generalized posynomial if it can be formed using addition, multiplication, positive power, and maximum selection starting from posynomials [7].

The energy equation is also a generalized posynomial: the first term is just a linear combination of the gate sizes while the second term is another linear combination of the gate sizes multiplied by the cycle time, that in turn is related to the delay through the critical path, hence also a generalized posynomial.

The optimization problem described in Sect. 2 using the above models has generalized posynomial objective and constraint functions:

Adjust W_i in order to *Minimize* $max(t_{arrival, primary_outputs})$ subject to:

$$E \leq E_{max}, t_{slope, primary\ outputs} \leq t_{slope_out, max}, t_{slope, internal\ nodes} \leq t_{slope\ internal, max}, C_{primary\ inputs} \leq C_{in, max}, W_i \geq I.$$

Table 1. Comparison between analytical and tabulated models

ANALYTICAL MODELS	TABULATED MODELS
- limited accuracy	+ very accurate
+ fast parameter extraction	- slow to generate
+ provide circuit operation insight	- no insight in the operation of the circuit
+ can exploit mathematical properties to formulate a convex optimization problem	- can't guarantee convexity; optimization is "blind"

Such an optimization problem with generalized posynomials is called a *generalized geometric program* (GGP) [7]. It can be converted to a convex optimization problem using a simple change of variables:

$$W_i = \exp(z_i) \quad (7)$$

With this change of variables the problem is tractable and can be easily and reliably solved by generic commercial optimizers. Moreover, since in convex optimization any local minimum is also global, the optimality of the result is *guaranteed*.

This delay model applies to any logic family where a gate can be represented through channel-connected components [8], as in the case of complementary CMOS or domino logic. The limitation of this approach is that it uses linear approximations for the delay, signal slopes, and capacitances. Fig. 3 shows a comparison of the actual and predicted delay for the rising transition of a gate for a fixed input slope and variable fanout. Since the actual delay is slightly concave in the fanout, the linear model is pessimistic at low and high fanouts and optimistic in the mid-range.

3.2 Tabulated Models

If the accuracy of such linear, analytical models is not satisfactory tabulated models can be used instead. For instance, (1), (2) and their respective parameters can be replaced with the look-up table shown in Table 2.

The table can have more or less entries, depending on the desired accuracy and density of the grid. Actual delays and slopes used in the optimization procedure are obtained through linear interpolation between the points in the table. The grid is non-uniform, with more points in the mid-range fanouts and slopes, where most designs are likely to operate. Additional columns can be added to the tables for different logic families – for instance relative keeper size for dynamic gates.

The resulting optimization problem, even when using the change of variables from (7), cannot be proven to be convex. However, since the analytical models describe the correct behavior of the circuits (although not absolutely accurate), the resulting optimization problem is *nearly-convex* and can still be solved with very good accuracy and reliability by the same optimizers as before [9]. The result of the nearly-convex problem can be checked against a *near-optimality boundary*. The example in Fig. 4 shows a comparison of the analytical and tabulated models and the corresponding near-optimality boundary.

The figure shows the energy-delay tradeoff curves for an example 64-bit Kogge-Stone carry tree in static CMOS using a 130nm process. The same circuit is optimized using each of the two model choices discussed in this section.

Table 2. Example of a tabulated delay and slope model (NOR2 gate, input A, rising transition)

C_{load}/C_{in}	$t_{slope\ in}$	t_D	$t_{slope\ out}$
1	20 ps	19.3 ps	18.3 ps
...
10	200 ps	229.6 ps	339.8 ps

Both models show that the fastest static 64-bit carry tree can achieve the delay of approx. 560ps, while the lowest achievable energy is 19pJ per transition. The analytical models are slightly optimistic because the optimal designs exhibit mid-range gate

fanouts where the analytical models tend to underestimate the delays (Fig 3.). However, the models indeed exhibit the correct behavior without being grossly inaccurate.

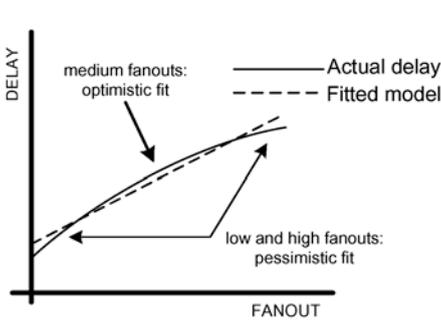


Fig. 3. Accuracy of fitted models

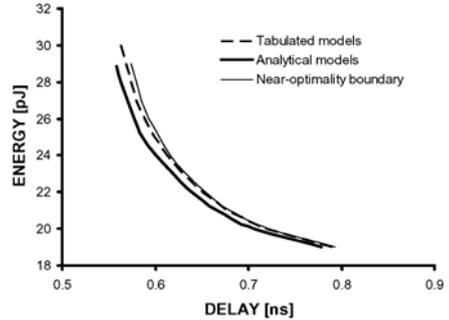


Fig. 4. Analytical vs. tabulated models and near-optimality boundary

The near optimality boundary is obtained by using tabulated models to compute the delay and energy of the designs resulted from the optimization with analytical models. This curve represents a set of designs optimized using analytical models, but evaluated with tabulated models. Since those designs are guaranteed to be optimal for analytical models, the boundary is within those models' error of the actual global optimum. However, if an optimization using the correct models (tabulated) converges to the correct solution, it will always yield a better result than a re-evaluation of the results of a different optimization using the same models. Therefore, if the optimization with tabulated models is to converge correctly the result must be within the near-optimality boundary (e.g. smaller delay for the same energy).

If a solution obtained using tabulated models is within the near-optimality boundary it will be deemed “near-optimal” and hence acceptable.

In a more general interpretation, optimizing using tabulated models is equivalent to optimizing using a trusted timing signoff tool whose main feature is very good accuracy. The result of such an optimization is not guaranteed to be globally optimal. The near-optimality boundary is obtained by running the timing signoff tool on a design obtained from an optimization that can guarantee the global optimality of the solution. The comparison is fair because the power and performance figures on both curves are evaluated using the same (trusted and accurate) timing signoff tool.

3.3 Model Generation and Accuracy

Tabulated models are generated through simulation. The gate to be modeled is placed in a simple test circuit and the fanout and input slope are adjusted using perl scripts. The simulator is invoked iteratively for all the points in the table and the relevant output data (delay, output slope) is stored. This can be lengthy (although parallelizable) if the grid is very fine and the number of points large. This characterization is similar to the one performed for the standard-cell libraries, and yields satisfactory accuracy.

For the analytical models data points are obtained through simulation in the same manner as for tabulated models. Least squares fitting (in Matlab) is used to obtain the

parameters of the models. The number of points required for a good fit (50 – 100, depending on the model) is less than the number of points needed for tabulated models (at least 1000) and thus the characterization time for analytical models is one order of magnitude shorter.

The error of the analytical models depends on their complexity and on the desired data range. The models in (1) and (2) are accurate within 10% of the actual delays and slopes for the range specified in Table 2. The energy equation (3) is accurate within 5% for fast slopes but its accuracy degrades to 12% underestimation at slow input slopes due to the crowbar current (which is not accounted for). The maximum slope constraints for output and internal nodes ensure such worst cases do not occur in usual designs.

4 Results

We use the presented optimization framework to optimize a 64-bit adder, which is a very common component of custom datapaths. The critical path of the adder consists of the carry computation tree and the sum select [10]. Tradeoffs between the performance and power can be performed through the selection of circuit style, logic design of carry equations, selection of a tree that calculates the carries, as well as through sizing and choices of supply voltages and transistor thresholds.

Carry-lookahead adders are frequently used in high-performance microprocessor datapaths. Although adder design is a well-documented research area [11,12,13,14], fundamental understanding of their energy-delay performance at the circuit level is still largely invisible to the microarchitects. The optimization framework presented in this paper provides a means of finding the energy budget breakpoint where the architects should change the underlying circuit design.

Datapath adders are good example for the optimization because their layout is often bit-sliced. Therefore, the critical wire lengths can be estimated pre-design and are a weak function of gate sizing. The optimization is performed on two examples:

1. A 64-bit Kogge-Stone adder carry tree implemented in standard static CMOS, using analytical models to tune gate sizes, supply and threshold voltages;
2. 64-bit carry lookahead adders implemented in domino and static CMOS, using tabulated models.

4.1 Tuning Sizes, Supply and Threshold Using Analytical Models

In order to tune supply and threshold voltages, the models must include their dependencies. A gate equivalent resistance can be computed from analytical saturation current models (a reduced form of the BSIM3v3 [15,16]):

$$R_{EQ} = \frac{1}{V_{DD}/2} \int_{V_{DD}/2}^{V_{DD}} \frac{V_{DS} dV_{DS}}{I_{DSAT}} = \frac{3}{4} \frac{V_{DD}(\beta_1 V_{DD} + \beta_0 + V_{DD} - V_{TH})}{W \cdot K (V_{DD} - V_{TH})^2} \left(1 - \frac{7V_{DD}}{9V_A}\right) \quad (8)$$

Using (8), supply and threshold dependencies can be included in the delay model. For instance (1) becomes (9), with (2) having a very similar expression:

$$t_D = c_2 R_{EQ} + c_1 R_{EQ} \frac{C_{load}}{C_{in}} + (\eta_0 + \eta_1 V_{DD}) \cdot t_{slope_in} \quad (9)$$

The model is accurate within 8% of the actual (simulated) delays and slopes around nominal supply and threshold, over a reasonable yet limited range of fanouts (2.5 – 6). For a +/- 30% range in supply and threshold voltages the accuracy is 15%.

Fig. 5 shows the optimal energy-delay tradeoff curves of a 64-bit Kogge-Stone carry tree implemented in static CMOS in three cases:

1. Only gate sizes are optimized for various fixed supplies and the nominal threshold;
2. Gate sizes and supply are optimized for nominal threshold;
3. Gate sizes, supply and threshold voltage are optimized jointly.

Fig. 6 shows the corresponding optimal supply voltage for case 2 and Fig. 7 shows the corresponding optimal threshold for case 3 normalized to the nominal threshold voltage of the technology.

A few interesting conclusions can be drawn from the above figures:

- The nominal supply voltage is optimal in exactly one point, where the $V_{DD} = 1.2V$ curve is tangent to the optimal V_{DD} curve. In that point, the sensitivities of the design to both supply and sizing are equal [2];
- Power can be reduced by increasing V_{DD} and downsizing if the V_{DD} sensitivity is less than the sizing sensitivity;
- The last picosecond is very expensive to achieve because of the large sizing sensitivity (curves are very steep at low delays);
- The optimal threshold is well below the nominal threshold. For such a high activity circuit, the power lost through increased leakage is recuperated by the downsizing afforded by the faster transistors with lower threshold. Markovic et al, [2], came to a similar conclusion using a slightly different approach.

4.2 Tuning Sizes in 64-Bit CLA Adders Using Tabulated Models

Using tabulated models as described in Sect. 3, various adder topologies implemented in different logic families are optimized in the energy–delay space under the typical loading for a microprocessor datapath. Details about the logic structure of the adders can be found in [17]. Fig. 8 shows the energy – delay tradeoff curves for a few representative adder configurations. Radix-2 (R2) adders merge 2 carries at each node of the carry tree. For 64 bits, the tree has 6 stages of relatively simple gates. Radix-4 (R4) adders merge 4 carries at each stage, and therefore a 64-bit tree has only 3 stages but the gates are more complex. In the notation used in Fig. 8 classical domino adders use only (skewed) inverters after a dynamic gate, whereas compound domino use more complex static gates, performing actual radix-2 carry-merge operations [18].

Based on these tradeoff curves, microarchitects can clearly determine that under these loading conditions radix-4 domino adders are always preferred to radix-2 domino adders. For delays longer than 12.5 FO4 inverter delays, a static adder is the preferred choice because of its lower energy.

Static adders are generally low power but slow, while domino logic is the choice for short cycle time. The fastest adder implements Ling’s pseudo-carry equations in a domino radix-4 tree with a sparseness factor of 2 [17].

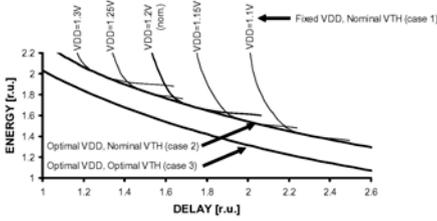


Fig. 5. Energy - delay tradeoff curves for different sets of optimization variables

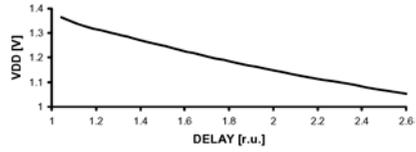


Fig. 6. Optimal supply voltage for designs sized with nominal threshold voltage

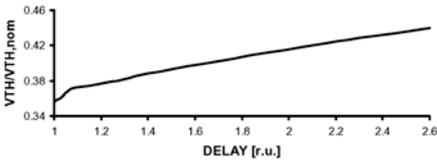


Fig. 7. Optimal threshold voltage when all optimizations are performed simultaneously

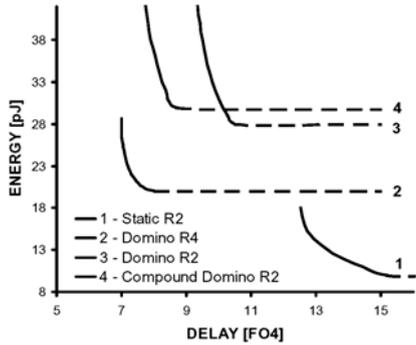


Fig. 8. Energy – delay tradeoff curves for selected 64-bit CLA adders

4.3 Runtime Analysis

The complexity and runtime of the framework depend on the size of the circuit. Small circuits are optimized almost instantaneously. A 64-bit domino adder with 1344 gates (a fairly large combinational block) is optimized on a 900MHz P3 notebook computer with 256MB of RAM in 30 seconds to 1 minute if the constraints are rather lax. When the constraints are particularly tight and the optimizer struggles to keep the optimization problem feasible, the time increases to about 3 minutes. A full power – performance tradeoff curve with 100 points can be obtained in about 90 minutes on such a machine. For grossly infeasible problems the optimizer provides a “certificate of infeasibility” in a matter of seconds.

For large designs the framework allows gate grouping. By keeping the same relative aspect ratio for certain groups of gates, the number of variables can be reduced and the runtime kept reasonable. Gate grouping is a natural solution for circuits with regular structure. All the adders optimized in Sect. 4.1 and 4.2 use gate grouping for identical gates in the same stage.

5 Conclusions

This paper presents a design optimization framework that tunes custom digital circuits based on a static timing formulation. The framework can use a wide variety of models and tune different design variables. The problem solved is generally an en-

ergy-constrained delay minimization. Due to the flexibility in choosing models, the framework can easily handle various logic families.

If analytical models are used the optimization is convex, can be easily and reliably solved, and its results are guaranteed to be optimal. The accuracy of the modeling can be improved by using look-up tables, at the cost of the optimality guarantee as well as increased characterization time and complexity. More generally, the optimization can be run on any trusted and accurate timing signoff tool, with the same tradeoffs and limitations as for tabulated models. Results obtained using tabulated models (or with the said “trusted and accurate timing signoff tool”) can be verified against a near-optimality boundary computed from results guaranteed optimal in their class. If the results fall within that boundary they are considered near-optimal and therefore acceptable.

The framework was demonstrated on 64-bit carry-lookahead adders in 130nm CMOS. A static Kogge-Stone tree was tuned using analytical models by adjusting gate sizes, supply voltage, and threshold voltage. Complete domino and static 64-bit adders were also tuned in a typical high performance microprocessor environment using tabulated models by adjusting gate sizes.

We are extending this framework to optimize the latch positions in pipelined datapaths. By building on the combinational circuit optimization, this tool would allow microarchitects a larger freedom in trading off cycle time for latency.

Acknowledgement

This work was supported in part by NSF grant ECS-0238572.

References

1. Penzes, P. I.; Martin, A. J.: Energy – Delay Efficiency of VLSI Computations, Proc. Great Lakes Symposium on VLSI, 2002,104-111
2. Markovic, D. et al.: Methods for True Energy – Performance Optimization. IEEE Journal of Solid State Circuits vol. 39 Issue 8 (Aug. 2004) 1282 – 1293
3. Conn, A. R.,et al.: Gradient – Based Optimization of Custom Circuits Using a Static Timing Formulation. Proceedings of Design Automation Conference DAC’99, 452 – 459
4. Fishburn, J. P.; Dunlop, A. E.:TILOS: A Posynomial Programming Approach to Transistor Sizing. IEEE International Conference on Computer – Aided Design ICCAD’85, 326-328
5. Sutherland I., Sproul R., Harris D.: Logical Effort, Morgan-Kaufmann, 1999
6. Synopsys® Design Compiler User’s Manual Version 2004.12
7. Boyd, S.; Vandenberghe, L: Convex Optimization, Cambridge University Press, 2003
8. Zlatanovici, R; Master thesis, UC Berkeley, 2002
9. Mathworks, Matlab® Optimization Toolbox User’s Guide Version 3
10. Rabaey, J. M.; Chandrakasn, A.; Nikolic, B: Digital Integrated Circuits: A Design Perspective, 2ndedition, Prentice-Hall 2003
11. Kogge, P. M; Stone, H. S.: A Parallel Algorithm for Efficient Solution of a General Class of Recursive Equations, IEEE Transactions on Computer,s August 1973, 786-793
12. Park, J.; Ngo, H. C.; Silberman, J. A.; Dhong, S. H.: 470ps 64bit Parallel Binary Adder, 2000 Symposium on VLSI Circuits,192-193
13. Han T.; Carlson, D. A.: Fast Area Efficient VLSI Adders, 8th Symposium on Computer Arithmetic 1987, 49-56

14. Naffziger, S.: A Sub-nanosecond 0.5 μm 64b Adder Design, International Solid-State Circuits Conference, 1996, 210-211
15. Toh, K. Y.; Ko, P. K.; Meyer R. G.: An Engineering Model for Short-channel CMOS Devices. IEEE Journal of Solid State Circuits vol. 23 Issue 4 (Aug. 1998) 950 – 958
16. Garrett, J; Master thesis, UC Berkeley, 2004
17. Zlatanovici, R.; Nikolic, B.: Power – Performance Optimal 64-bit Carry-lookahead Adders. European Solid State Circuit Conference ESSCIRC 2003, 321 – 324
18. Dao, H. Q; Zeydel, B. R.; Oklobdzija, V. G.: Energy Minimization Method for Optimal Energy – Delay Extraction. European Solid State Circuit Conference ESSCIRC 2003, 177-180