

## On Estimating the Performance of VLSI Circuits

Clark D. Thompson  
Prabhakar Raghavan

Computer Science Division  
573 Evans Hall  
University of California  
Berkeley, CA 94720

### ABSTRACT

Metrics of energy, area, and time are defined for a graph-theoretic model of VLSI computation. A number of "technological constant factors" are introduced in order to account for the effects of using different technologies for implementing logic circuits. Different constant factors are seen to be appropriate for different logic families. We examine seven such families: NMOS, CMOS, CMOS-SOS,  $\text{I}^2\text{L}$ , GaAs HEMT, JJ-CIL, and JJ-CS.

### 1. Introduction

The area of VLSI modeling has been the scene of a disparity between circuit theorists and complexity theorists. The former have developed methods for the detailed study of circuit behavior, while the latter have used graph-theoretic models to derive asymptotic results.

The circuit-theoretic approach, while accurate to within the stability and robustness of the numerical methods used, is time-consuming; circuits involving more than a few transistors take extremely long to simulate. Further, in the process of producing the required information they usually generate a good deal of data that is not directly usable. Finally, these circuit techniques do not give any idea of the area of the chips they attempt to model; for this, a detailed layout must be prepared. The factors governing these layouts and simulations vary widely from technology to technology, so that the techniques used are not general in any sense.

Complexity theorists have sought to avoid these limitations, and have succeeded in a way; however, their techniques have resulted in some difficulties of a different nature. This approach<sup>11,14,15</sup> utilizes the concept of embedding a circuit as a graph on a planar grid, making statements about the asymptotic

growth of area and time in these embedded circuits. Elegant proof techniques involving graph theory and information flow can be applied. Unfortunately, the very simplicity of the graph abstraction proves to be its major shortcoming. Gates are abstracted as grid points; in essence all the area is assumed to be taken up by wiring.

More important, the physical phenomena underlying the devices used in VLSI make many of the asymptotic results difficult to interpret; we might, for instance be faced with a result such as "The area of a chip solving problem  $P$  is  $N^2$  for input size  $N$ " whereas in fact we actually mean "The area of a chip solving problem  $P$  is  $N^2$  for input size  $N$ , provided the area does not exceed  $100 \text{ mm}^2$ ." While there is more truth to the latter statement, the relation of asymptotic limits to physical constraints is hard to perceive. Further, designers and circuit modelers are often interested in evaluating different approaches to solve a problem that might only differ by a constant factor; such distinctions are lost in asymptotic studies.

We propose a new model here that attempts to tackle some of these problems. While it would be tempting to report that our model captures the best of both worlds by combining the simplicity of the complexity theoretic approach with the realism of the circuit simulators, there are a number of problems that prevent us from attaining this happy state of affairs. We outline these problems later in this paper:

We have, however, developed a graph-theoretic formulation that includes the notions of gate size, wire size, fanout, gate and wire delay, wiring layers and signal propagation. More importantly, we have classified the modes of energy consumption in integrated circuits and supplied the means for evaluating this quantity. Finally, we have equipped the model with a number of "technological constants", thereby furnishing it with the versatility to deal with various fabrication technologies and logic families.

---

This work was supported by the Semiconductor Research Corporation Grant SRC 1-44247-52055.

Section 2 is a statement of the model; this is followed by a list of values for technological constants for present and future technologies in section 3. Section 4 discusses model parameters for seven current technologies. The final section presents an example and concludes by outlining some of the problems faced in applying this model.

**2. Model of Computation**

In the following assumptions, greek letters are used for the technological constant factors. (There are two exceptions:  $\delta$  and  $\Delta$  bear their standard meaning of vertex in- and out-degree.) Sets and their elements are defined by capitalized and lower case roman letters, respectively.

1. *Sources  $O_h$ , sinks  $I_h$ .* A computation graph is a directed hypergraph  $G = (V, H)$ . A hyperedge  $h$  is denoted by an ordered pair  $(O_h, I_h)$  of vertex sets  $O_h \subseteq V, I_h \subseteq V$ . The vertices in  $O_h$  are the sources of  $h$ ; the vertices in  $I_h$  are its sinks.

2. *Edge fanout restrictions  $o_{max}, i_{max}$ .* Each edge  $h$  has at least one, and at most  $O(1)$ , sources and sinks:

$$1 \leq |O_h| \leq o_{max}, \quad 1 \leq |I_h| \leq i_{max}$$

Limits on vertex indegree  $\delta$  and outdegree  $\Delta$  are discussed in Assumption 10.

3. *Vertex widths  $\lambda_{gate}, \lambda_{I/O}$ .* Each vertex in a computation graph is embedded as a square region in the Euclidean plane. No two vertices overlap. The size of an embedded vertex depends upon its functionality: gates  $v \in V_{gate}$  occupy  $\lambda_{gate}^2$  area, while I/O ports  $v \in V_{I/O}$  occupy  $\lambda_{I/O}^2$  area.

4. *Edge width  $\lambda_{wire}$ , number of wiring layers  $\mu$ .* An edge is embedded as a connected set of wire segments. Each wire segment is a rectangle of width  $\lambda_{wire}$  and arbitrary length, placed on one of  $\mu$  planar wiring layers stacked above the plane of the vertices. A wire segment on the bottom wiring layer connects to the vertices it passes over. Two wire segments are connected to each other if they pass over the same point and if they are on either the same layer or an adjacent layer. (Note that  $\lfloor \mu/2 \rfloor$  disconnected wire segments may pass over the same point in the vertex plane. Also note that any hyperedge  $h$  can be embedded as a tree of wire segments passing over the vertices in  $O_h \cup I_h$ .)

5. *Total area  $A$ , maximum total area  $\alpha_{max}$ .* The total area  $A$  of an embedded computation graph is the area of the smallest square that encloses all its vertices and wire segments.

The area of this square is bounded by a technological constant:  $A \leq \alpha_{max}$ .

6. *Maximum edge length  $\lambda_{max}$ .* The total length  $\|h\|$  of an (embedded) edge  $h$  is the sum of the lengths of its wire segments. Edge lengths are bounded by a constant:  $\forall h, \|h\| \leq \lambda_{max}$ .

7. *Votes  $v(t)$ , signals  $h(t)$ .* The state of the computation graph at any time  $t$  is defined by a vector  $(V(t), H(t))$  of votes  $v(t)$  and signals  $h(t)$  associated with each vertex  $v$  and hyperedge  $h$ . The value of a vote or signal is taken from the ternary set  $\{0, 1, u\}$ : logic-0, logic-1, and undetermined. (An alternative formulation, found in<sup>11</sup> and in state-of-the-art circuit simulators, takes signal values from a two-dimensional set of voltages  $V$  and impedances  $R$ .)

8. *Maximum size of voting equivalence class  $\xi_{max}$ , edge delay  $d_h$ , time constants  $\tau_{gate}, \tau_{wire}$  and  $\tau_{fanout}$ , transmission line indicator  $\zeta$ , signal rise time  $r_h$ .*

a. In many technologies, the delay associated with a wire can be decreased by driving that wire with a larger transistor. Such high-power drivers can be represented by several (unit-power) sources with identical voting behavior. We are thus led to the following definition of equivalence classes  $C_{h,i}$ , on the voting behavior of the sources for each edge  $h$ :

$$v_1, v_2 \in C_{h,i} \iff (v_1, v_2 \in O_h) \wedge (\forall t, v_1(t) = v_2(t))$$

A technological limit on driving power translates into a restriction on the size of (i.e., number of vertices in) any voting equivalence class:

$$\forall h, i, |C_{h,i}| \leq \xi_{max}$$

b. At the time of circuit construction ( $t = 0$ ), a fixed but indeterminate delay  $d_h$  is assigned to each edge  $h$ . An edge's delay (in a worst-case analysis) is proportional to its length  $\|h\|$  and number of sinks  $|I_h|$ , and inversely proportional to the size of its smallest equivalence class  $C_h = \min_i |C_{h,i}| \leq \xi_{max}$ :

$$d_h = \tau_{gate} + \frac{\|h\| \tau_{wire} + |I_h| \tau_{fanout}}{C_h} (\pm 50\%)$$

(Indeterminacy is introduced into the definition of edge delays to force "realistic" design practices, e.g., self-timed or clocked logic.)

- c. We define  $r_h$  to be the rise time of a signal on edge  $h$ . For technologies in which wires are transmission lines,  $r_h$  is approximately equal to the gate delay  $\tau_{gate}$ . We indicate this by assigning the value 1 to the 0-1 variable  $\zeta$  (a mnemonic is the common symbol  $Z$  for the impedance of a line). The other technological possibility is that the wires are essentially capacitive in nature (as long as their length does not exceed  $\lambda_{max}$ , as defined in Assumption 6). Thus

$$r_h = \begin{cases} d_h, & \text{if } \zeta = 0 \\ \tau_{gate}, & \text{if } \zeta = 1 \end{cases}$$

- d. The value of a signal  $h(t)$  is determined by the votes of its sources  $O_h$ , with delay  $d_h$ . We prevent the propagation of unreasonably-short signal pulsewidths by requiring the "election results" to be stable for at least  $r_h$  time units.

$$h(t) = \begin{cases} 1, & \text{if } \exists v \in O_h \forall s \in [t-d_h, t-d_h+r_h] v(s) = 1, \text{ else} \\ 0, & \text{if } \exists v \in O_h \forall s \in [t-d_h, t-d_h+r_h] v(s) = 0, \text{ else} \\ u & \end{cases}$$

Note that this formulation allows "wire-oring": the signal on an edge becomes 1 if any of its source votes is 1 for at least  $r_h$  time.

9. *Symmetry indicator*  $\sigma$ .

- a. Not all patterns of voting behavior are allowed in all technologies. One restriction is observed in the so-called "symmetric" technologies ( $\sigma = 1$ ). In these technologies, the effects of logic-1 votes and logic-0 vote are symmetric, making "wire-oring" infeasible. (A system of "majority-rule" is conceivable but not observed in any present-day logic family, possibly because it would reduce noise margins.) To outlaw wire-oring, we permit just one equivalence class per edge:

$$(\sigma = 1) \Rightarrow (\forall t \forall h \forall v_1, v_2 \in O_h v_1(t) = v_2(t))$$

- b. A second type of restriction on allowable voting behavior arises in the asymmetric ( $\sigma = 0$ ) technologies. We must restrict the number of high-power logic-1 votes that appear at one time on an edge, to avoid exceeding the current density limit mentioned in Assumption 8a:

$$(\sigma = 0) \Rightarrow (\forall t \forall h |\{v \in O_h : v(t) = 1\}| \leq \xi_{max})$$

10. *Logic family*  $\phi$ , power supply period  $\tau_{supply}$ , I/O schedule  $S$ , external clock period  $\tau_{I/O}$ .

- a. A logic family  $\phi$  is a technologically-constrained set of triples  $(\delta, f, \Delta)$ . The first and third parts of a triple denote the indegree and outdegree of one type of gate. The second part of a triple defines a functionality, or voting behavior. A gate with the 'and' functionality, for example, is modeled by a vertex whose vote is the logical 'and' of the signals on its in-edges. As another example, the 'latch' function depends on a delayed feedback signal. Finally, the voting behavior of gates in the JJ-CIL technology depends upon the phase of their AC power supply. Thus, in the general case, the functionality  $f_v$  of a gate  $v$  has  $2 + \delta_v$  parameters, and defines the gate's vote as follows:

$$v(t) = f_v(v(t - \tau_{wire}), c(t), h_1(t), h_2(t), \dots, h_{\delta_v}(t))$$

where the phase of the power supply (assumed to have a 90% duty cycle) is

$$c(t) = \begin{cases} 0, & \text{if } t \leq (.1 + \lfloor t/\tau_{supply} \rfloor) \tau_{supply} \\ 1, & \text{otherwise} \end{cases}$$

Note that voting is a zero-delay process, since gate delays were included in the definition of edge delay  $d_h$ .

- b. An I/O port  $v_i \in V_{I/O}$  has  $\delta_v = 1$ ,  $\Delta_v = 1$ . Its voting is determined by an externally-imposed I/O schedule  $S_i \in \{r_0, r_1, r_u, w_0, w_1, w_u\}^*$ . Each literal in  $S_i$  indicates whether the I/O port is to read ( $r_0, r_1, r_u$ ) or write ( $w_0, w_1, w_u$ ) a '0', a '1', or a 'u'. The  $k$ -th literal in  $S_i$  refers to the  $k$ -th external clock period defined by  $t \in ((k-1)\tau_{I/O}, k\tau_{I/O}]$ , where  $\tau_{I/O}$  is a technological constant. Thus, if the  $k$ -th literal is  $r_x$ , the port votes  $v_i(t) = x$  during the  $k$ -th clock period. Alternatively, if the  $k$ -th literal is  $w_y$ , we say the schedule  $S_i$  is "satisfied" only if the port's in-edge  $h$  has signal  $h(t) = y$ , for all times  $t$  in the  $k$ -th clock period. (If the output bit for some time period is  $u$ , i.e. undetermined, we allow  $h(t)$  to be any value.)

11. *Energy consumption*  $E_{standby}$ ,  $E_{1-0}$ ,  $E_{wire}$ ,  $E_{sink}$ ,  $E$ . Four modes of energy consumption are observed in physical realizations of computation graphs.

- a. A constant power dissipation of  $\epsilon_{standby}/\tau_{gate}$  is associ-

ated with every gate. The worst case (over all I/O schedules  $S$ ) total "standby" energy dissipation over the period  $[t_1, t_2]$  is thus defined as

$$E_{standby} = \max_S \sum_{v \in V_{gate}} \frac{(t_2 - t_1)}{\tau_{gate}} \epsilon_{standby}$$

- b. In asymmetric (wire-or) technologies, a gate voting 1 consumes more power than a gate voting 0. We define energy  $\epsilon_{1-0}$  so that the difference between these two levels of power consumption is  $\epsilon_{1-0}/\tau_{gate}$ . Total energy consumption in this mode is thus

$$E_{1-0} = \max_S \sum_{v \in V_{gate}} \int_{\substack{t_1 \leq t \leq t_2 \\ v(t)=1}} \frac{\epsilon_{1-0}}{\tau_{gate}} dt$$

By Assumption 8d, a gate's vote can change a signal only if it persists for at least  $r_h \geq \tau_{gate}$  time. We thus employ the following (approximate) expression for  $E_{1-0}$ :

$$E_{1-0} = \max_S \sum_{v \in V_{gate}} \sum_{\substack{t_1/\tau_{gate} \leq k \leq t_2/\tau_{gate} \\ v(k\tau_{gate})=1}} \epsilon_{1-0}$$

- c. Each change in an edge's signal consumes energy proportional to the length of that edge. Assuming such signal changes occur at a frequency less than  $1/\tau_{gate}$ , we write

$$E_{wire} = \max_S \sum_h \sum_{\substack{t_1/\tau_{gate} \leq k \leq t_2/\tau_{gate} \\ h(k\tau_{gate}) \neq h((k+1)\tau_{gate})}} \|h\| \epsilon_{wire}$$

- d. Energy  $E_{sink}$ , like  $E_{wire}$ , is a form of "switching energy." In this case, the energy consumption is proportional to the number of sinks:

$$E_{sink} = \max_S \sum_h \sum_{\substack{t_1/\tau_{gate} \leq k \leq t_2/\tau_{gate} \\ h(k\tau_{gate}) \neq h((k+1)\tau_{gate})}} |J_h| \epsilon_{sink}$$

- e. The total energy consumed by a computation is  $E = E_{standby} + E_{1-0} + E_{wire} + E_{sink}$ .

### 3. Technological Parameters

In this section we give a list of technological constant values for seven technologies. Present-day values as well as projected values for circuits fabricated in the late eighties are listed.

Table 1 gives approximate values for the "constant factors" of seven VLSI technologies with current fabrication and circuit-design techniques.

An important feature of Table 1 is the diagonal structure of the entries for circuit energies  $\epsilon_{standby}$ ,  $\epsilon_{1-0}$ ,  $\epsilon_{wire}$ , and  $\epsilon_{sink}$ . When calculating total energy, contributions from entries below the diagonal can be ignored. For example, the technologies with  $\epsilon_{standby} > 0$  have a nearly constant power dissipation per gate which does not increase by more than 10% when the gates change their state at Maximum frequency.

Table 1 is not quite a complete list of the parameters in the model. The following are nearly constant over all technologies:

$$\alpha_{max} = 10^8 \text{ um}^2, \quad \tau_{I/O} = 20 \text{ ns}, \\ \mu = 4 \text{ to } 6 \text{ layers}, \quad \lambda_{I/O} = 10^2 \text{ um}.$$

Note that, by Assumption 4,  $\mu = 6$  corresponds to a three-level metal process. The other  $\mu - 3$  layers are made of an insulating material, through which small square holes or "vias" are cut.

Table 2 presents technological constants for the future, (hopefully) valid for late-1980s fabrication. The following are nearly constant over all technologies:

$$\alpha_{max} = 10^9 \text{ um}^2, \quad \tau_{I/O} = 10 \text{ ns}, \\ \mu = 4 \text{ to } 6 \text{ layers}, \quad \lambda_{I/O} = 50 \text{ um}.$$

Table 3 indicates the availability of gates in each of the seven technologies; we use these gates as the basis of our (conservative) estimates for various parameters in section 4. The use of pass-transistors on the input of a 2-input NAND in NMOS gives us a 4-input NAND in this family.

## 4. Model Parameters

This section deals with the study of model parameters for current technologies.

### 4.1. Derivation of Model Parameters

The model we have just delineated includes a number of "technological constants" to cater to the wide variety of technologies and logic families available. We begin this section by providing some intuition into the process of obtaining values for these constants. Following this we deal with some aspects of circuit behavior common to all technologies, notably the properties of interconnect lines and their futuristic trends. A detailed discussion of individual technologies and parameter values concludes the section.

The most fundamental specification for any technology is the *linewidth*, or the size of the smallest feature that can be fabricated. The parameter  $\lambda_{wire}$  is assumed to be twice the linewidth for most technologies. An exception is the Josephson

current-steering technology (or, indeed, any technology that transmits information via a current loop rather than a voltage referenced to a universal ground plane); in this case we have to consider the fact that every physical wire is accompanied by a return path for closing the loop. We then take  $\lambda_{wire}$  to be four times the linewidth.

From the device characteristics, we could determine the types of gates that can be made. In the literature, this specification is usually given along with the process description. We directly have the feasible family of boolean functions ( $\delta, f, \Delta$ ) available in the logic family  $\phi$ . From the dimensions of the devices made in the process, we can then form an estimate of  $\lambda_{gate}$

the limit  $\lambda_{max}$  is to avoid the problem of operating lines in the diffusion regime.

Most descriptions of new logic give delay figures for basic inverters (either by themselves or in a ring-oscillator configuration). Since the figure given for inverter/gate delay usually assumes that the gate is driving an identical gate, we can form a conservative estimate by setting  $\tau_{fanout}$  also equal to this figure. Note that  $\tau_{gate}$  includes two components - one for the intrinsic switching time of the gate (in some sense, a "no-load" switching time), and the other for driving the load (in most technologies, this appears to dominate). Only the latter component is augmented by the addition of fanout; hence the conservative

	I <sup>2</sup> L	JJ-CIL	NMOS	HEMT	JJ-CS	CMOS	CMOS-SOS	units
$\lambda_{gate}$	10	60	70	70	100	100	100	$\mu m$
$\lambda_{wire}$	4	10	4	4	10	4	4	$\mu m$
$\lambda_{max}$	$10^4$	$10^5$	$10^4$	$10^4$	$10^5$	$10^4$	$10^4$	$\mu m$
$\tau_{gate}, \tau_{fanout}$	2000	20	500	50	500	2000	2000	ps
$\tau_{wire}$	1	0.02	1	1	1	1	0.5	ps/ $\lambda_{wire}$
$\tau_{supply}$		1000						ps
$\epsilon_{standby}$	10	0.01	$\sim 0$	$\sim 0$	$\sim 0$	$\sim 0$	$\sim 0$	fJ
$\epsilon_{1-0}$			1	0.1	$\sim 0$	$\sim 0$	$\sim 0$	pJ
$\epsilon_{wire}$					0.02	1	0.01	fJ/ $\lambda_{wire}$
$\epsilon_{sink}$							1	fJ
$\xi_{max}$	1	1	10	10	1	10	10	--
$\sigma$	0	1	0	0	1	1	1	--
$\zeta$	0	1	0	0	0	0	0	--
$\sigma_{max}$	$10^2$	1	$10^2$	$10^2$	1	(= $\xi_{max}$ )	(= $\xi_{max}$ )	--
$\iota_{max}$	1	1	$10^4$	$10^4$	$10^4$	$10^4$	$10^4$	--

Table 1. Current constant factors for seven VLSI technologies.

(Multiplicative factors for units are  $f = 10^{-16}$ ,  $p = 10^{-12}$ ,  $n = 10^{-9}$ , and  $u = 10^{-6}$ .)

by a simple layout process (if gate size/density is not already specified by the manufacturer). The figures reported in Tables 1 and 2 reflect the size of the largest of the gates in  $\phi$  for a given technology. Timing estimates for large networks are generally the hardest to form from device-level data. The problem is compounded by the fact that the behavior of interconnect lines falls into several regimes.<sup>2</sup> The primary objective of our imposition of

nature of this estimate for  $\tau_{fanout}$ .

$\tau_{wire}$  can only be determined if we have the reactive properties of the interconnect lines (capacitance/inductance per unit length). We then compare this figure with the capacitance/inductance (depending on whether signals are propagated as voltages/currents) of a gate input, and determine what fraction of a gate load a unit length of line constitutes.

	I <sup>2</sup> L	JJ-CIL	NMOS	HEMT	JJ-CS	CMOS	CMOS-SOS	units
$\lambda_{gate}$	4	15	7	7	40	10	10	$\mu m$
$\lambda_{wire}$	0.5	2	0.5	0.5	4	0.5	0.5	$\mu m$
$\lambda_{max}$	$5 \cdot 10^4$	$10^6$	$5 \cdot 10^4$	$5 \cdot 10^4$	$10^6$	$5 \cdot 10^4$	$5 \cdot 10^4$	$\mu m$
$\tau_{gate}, \tau_{fanout}$	1000	5	100	10	200	100	100	ps
$\tau_{wire}$	0.1	0.005	0.1	0.1	1	0.1	0.1	ps / $\lambda_{wire}$
$\tau_{supply}$		1000						ps
$\epsilon_{standby}$	1	0.002	$\sim 0$	$\sim 0$	$\sim 0$	$\sim 0$	$\sim 0$	fJ
$\epsilon_{1-0}$			0.05	0.1	$\sim 0$	$\sim 0$	$\sim 0$	pJ
$\epsilon_{wire}$					0.02	0.05	0.001	fJ / $\lambda_{wire}$
$\epsilon_{sink}$							0.05	fJ
$\xi_{max}$	1	1	10	10	1	10	10	--
$\sigma$	0	1	0	0	1	1	1	--
$\zeta$	0	1	0	0	0	0	0	--
$o_{max}$	$10^2$	1	10	$10^2$	1	( $= \xi_{max}$ )	( $= \xi_{max}$ )	--
$l_{max}$	1	1	$10^4$	$10^4$	$10^4$	$10^4$	$10^4$	--

Table 2. Futuristic constant factors for seven VLSI technologies.

$\tau_{wire}$  is then given by the product of this fraction and  $\tau_{fanout}$ .

4.2. Interconnect Lines

We now turn to the characteristics of interconnection lines, as their behavior exerts considerable influence on timing and energy dissipation. Recent work on the properties of intercon-

nection lines on various substrates is reported in.<sup>16</sup> The model of Saraswat and Mohammadi<sup>13</sup> appears somewhat simplistic in that they assume a parallel plate model by which they are able to predict a linear reduction in capacitance.

The resistivity of interconnections depends on the material

$(\delta, \Delta).f$	I <sup>2</sup> L	JJ-CIL	NMOS	HEMT	JJ-CS	CMOS	CMOS-SOS
(1,1).INVERTER					*	*	*
(1,4).INVERTER			*				
(1,1).CLOCKED INVERTER			*				
(2,1).NAND, NOR					*	*	*
(2,1).AND, OR				*		*	*
(2,1).ARBITRARY						*	*
(4,1).NAND					*		*
(4,1).NOR						*	*

Table 3. Gate availabilities in seven VLSI technologies.

(Availability indicated by \*).

used for the lines; while aluminum and polysilicon have conventionally been used, metals like tungsten and titanium (and their silicides) are under investigation. For current linewidths these have a line resistivity of about  $500 \Omega/cm$ . This resistivity increases as lines are made smaller; while the value increases as the square of the scaling factor (a figure  $>1$ , by which all dimensions are divided), it is important to recognize the fact that the absolute resistance for some fixed line-length (say  $\lambda_{gate}$ ) is a more meaningful metric. This figure goes up linearly as dimensions are scaled.

Capacitances are harder to deal with, since they depend on the material used for the substrate and the insulation rather than the line itself. Further, interactions between adjacent lines introduce additional capacitances that are not easy to account for. Finally, a simple parallel plate model<sup>13</sup> is not always appropriate since "fringing effects" become particularly dominant at small linewidths. The parallel plate model predicts that capacitance per unit line-length remains essentially constant, while the capacitance per  $\lambda_{wire}$  scales down linearly. By taking a more rigorous approach, Yuan *et al.*<sup>16</sup> show that even at present (1-5  $\mu m$ ) linewidths, the latter figure does not diminish quite as rapidly for silicon substrates. Gallium arsenide and sapphire substrates are already into this zone of sublinear capacitance scaling. Current values for capacitance per  $\lambda_{wire}$  are in the region of 0.05 - 0.1 femtofarads per  $\lambda_{wire}$ , for various substrates. The work of Yuan *et al.*<sup>16</sup> also indicates that these figures rise by as much as a 100% for gallium arsenide and sapphire substrated when multiple lines are considered (due to inter-line capacitances). We therefore take the 0.1  $fF$  value for current technology.

The results mentioned above show that the RC product (per  $\lambda_{wire}$ ) will not remain constant, due to the fact that the increase in resistance will not quite be counterbalanced by the decrease in capacitance. We are not considering such figures as the delay on the longest line of the largest chip that can be fabricated in a technology (as do Saraswat and Mohammadi), the significance of which is not entirely obvious (large die sizes are not necessarily a part of the scaling process; further, it is not clear that such long lines are inevitable). It appears that with metal lines in current technology, we have the ability to transmit signals at delays below one nanosecond per cm. It also appears that the most important concern for designers of the future is to reduce the series resistance of lines; "tall" lines (small width and large height) could alleviate this problem, but are hard to etch.

### 4.3. On the Modeling of Seven Technologies

In this section we will look into the properties of computational devices manufactured by seven technologies in some detail.

#### 4.3.1. NMOS

The salient features of NMOS logic are noted below, following which we give a listing of model parameters and their future trends. There are two basic ingredients in NMOS logic: *active gates*, and *pass-transistor networks*. In our model, we depict gates by vertices in the hypergraph. Modeling pass-transistors explicitly has proven very complicated, and we choose instead to absorb them into the gates following them. Since good circuit design practice precludes long chains of pass-transistors (to maintain logic levels) we impose the constraint that no more than two pass-transistors can be connected serially without level restoration by buffers/gates.

We anticipate reductions of  $\lambda_{gate}$  and  $\lambda_{wire}$  to about a tenth of their current values by the beginning of the next decade. We also note that modeling the temporal behavior of transistors driving large loads (wire/fan-out) is a complex problem; we have therefore adopted the conservative technique of separating the gate- and wire-delay components. The figures reported for  $\lambda_{max}$  in<sup>2</sup> are 10  $mm$  for present technology and 50  $mm$  in the future. The work of Hoeneisen and Mead<sup>8</sup> and of Hart *et al.*<sup>7</sup> suggests that the ultimate limits on MOS speed will be reached at around a tenth of a nanosecond; this limit stems from the fundamental physical properties of MOS devices.

In NMOS it is possible to connect the outputs of several gates (sharing a common pullup) so that the output of the composite structure (at the lower end of the common pullup) is pulled down if any one of the gates (voters) conducts. In keeping with popular convention, we speak of this as a wire-or configuration (the exact logic function realized is, of course, dependent on which logic convention - positive or negative - is assumed). We restrict the number of voters in a wire-or configuration to a hundred; this is because a pull-down has a reverse-junction leakage current even when "OFF". In scaled transistors, the ratio of OFF-resistance to ON-resistance is even lower, and we expect the maximum possible number of wire-or voters to decrease to about ten. For similar reasons we restrict  $\sigma_{max}$  to be 10;  $i_{max}$  is restricted to  $10^3$  because of the limitation on  $\lambda_{max}$ .

The major phenomenon of interest to us is that an NMOS gate dissipates power when "ON" i.e. when the pulldown net-

work is conducting. In this condition, the power dissipation is governed by the size of the pull-up alone, since the pull-down resistance is small in comparison. Over the next decade, we expect  $\epsilon_{1-0}$  to drop from its present value in the picojoule range to the order of a hundredth of a picojoule. At this point we again run into fundamental limits, and intricate circuit techniques must be resorted to for lowering the energy dissipations of NMOS circuits; these may include the use of coding techniques, intermittently operational circuits (between bouts of activity, the circuit is "rested" to cool down) and other tricks.

#### 4.3.2. CMOS and CMOS-SOS

Logic in CMOS is realized through a combination of complementary gates and pass transistors. We shall follow the same convention for pass-transistors here as we did in the case of NMOS. The dominant effects in CMOS are circuit-dependent rather than device-dependent. It is believed,<sup>8</sup> for instance, that speed of operation would be constrained by power dissipation rather than by device physics. Almost all the dissipation in CMOS is during state transitions; this imposes a maximum rate of occurrence of transitions within a circuit.

We give a 2 ns figure for gate delay with present technology; this could drop by a factor of ten over the next decade. Clearly, wire-oring is not permitted, at least in vanilla CMOS. The parameters  $\sigma_{max}$  and  $i_{max}$  are restricted to ten and one thousand, as with NMOS. Depending on whether we consider ordinary CMOS or CMOS fabricated by the SOS technique, energy dissipation is dominated by  $\epsilon_{wire}$  or  $\epsilon_{sink}$ . This is because CMOS-SOS has very low interconnect capacitance making  $\epsilon_{wire}$  small. The dominant form of energy dissipation in each style is currently of the order of a femtojoule; this could scale down to about 0.05 fJ if the dissipation density is not excessive. The fact that CMOS dissipation occurs mainly during transitions could conceivably be exploited using circuit design and system timing schemes which exploit the speed of individual gates to the fullest extent without causing them to change state too often.

#### 4.3.3. Integrated Injection Logic

Integrated Injection Logic is perhaps the most promising bipolar technology for VLSI, mainly by virtue of its relatively low power dissipation. In addition, this logic family is dense and is particularly suitable for gate arrays. A typical gate consists of a four-output inverter whose outputs can be wire-ored with the outputs of other inverters. While a number of studies give results on the intrinsic limitations of the bipolar components

used for I<sup>2</sup>L, little has been said about the behavior of this logic family in circuit configurations. We rely primarily on the work of Evans<sup>4</sup> and of Hart *et al.*<sup>7</sup> for the constant factor estimates given here.

The parameter  $\lambda_{gate}$  is currently of the order of ten microns. Devices with  $\tau_{gate}$  as low as 2 ns are within the realm of currently available technology, but this is already close to physical limits imposed by collector-base capacitances. It is thus reasonable to conclude that in the case of I<sup>2</sup>L, the benefits of scaling will consist mainly of lower injection current requirements for minimum-delay devices.  $\epsilon_{standby}$ , currently of the order of a hundredth of a picojoule, can be expected to fall to about one femtojoule over the next decade. Dissipation density is likely to be a problem, since I<sup>2</sup>L is a very dense family. The problem of reverse leakage through wire-ored transistors is not as serious as for NMOS, and we permit up to a hundred sources for a hyperedge.

#### 4.3.4. Josephson Junction - Current Injection Logic

IBM has fabricated circuits out of current injecting Josephson junction devices.<sup>5</sup> Projections of circuit performance under device and circuit scaling are given by Ko and Van Duzer.<sup>9</sup>

There are two basic gates, a two-input AND and a two-input OR; inversion is accomplished by means of a "clocked inverter".<sup>5</sup> A feature unique to this technology is that gates are latching; in order to perform a new computation with a gate, it is necessary to turn off the power supply briefly. Circuits using this technology thus use power supplies that are cycled periodically (at approximately one nanosecond intervals in current practice).

Linewidths in this technology are presently on the order of two microns; but the presence of a number of inductances and resistances in the gates raises gate size to about 4000  $\mu m^2$ , yielding a figure of about 60  $\mu m$  for  $\lambda_{gate}$ . Signals propagate along interconnect lines that are essentially transmission lines. However, we impose the restriction that a signal transmitted during one power supply cycle should be received during the same cycle at the "far end". This is a somewhat conservative approach to synchronization, and limits the maximum length of interconnect transmission lines to about a tenth of a meter or 10<sup>6</sup>  $\mu m$ .

Gheewala has split delays in this technology into three components: (i) gate delay, (ii) crossing delay, and (iii) propagation delay. While gate and propagation delay have clear-cut analogues in other technologies, the term crossing delay requires

some explanation. Since information transfer is in the form of a current rather than a voltage, the output of a gate drives the input of another if the latter is a part of the output current loop of the former. A current pulse traveling on this loop suffers a delay in "crossing" the gate being driven; this corresponds to fanout delay in technologies like NMOS and CMOS.

Current values for these parameters are in the range of 10 ps for  $\tau_{gate}$  and for  $\tau_{fanout}$ , and .01 ps/ $\mu m$  for  $\tau_{wire}$ . Power dissipation is not a characteristic of any single state assumed by the gates, but takes on a steady-state form instead. Gheewala reports values in the range of 1-10  $\mu W$  for power dissipation, which yields figures in the range of 0.01-0.1 fJ for  $\epsilon_{standby}$ ; we adopt the upper limit for our estimate.

Ko and Van Duzer<sup>9</sup> suggest that in this technology, area and delay cannot be minimized simultaneously. Their work indicates, however, that gate delay can be brought down to about 5 ps by means of various optimizations, while propagation delay can be cut by a factor of four. Crossing delay seems harder to reduce, and it is likely that a heavy price will be paid for fanout in current injection logic of the future.

#### 4.3.5. Josephson Junction - Current Steering Logic

Current-steering logic is an alternate form of Josephson junction circuitry, developed for use in conjunction with single-flux-quantum memory devices.<sup>6</sup> Current-steered superconducting loops form the basis for logic implementation; this has a useful property we will discuss below. Guéret *et al.* have demonstrated the feasibility of a complete family of logic gates consisting of a two-input AND, a two-input OR and an inverter. It is worth noting that all gates provide both true and complemented outputs, so that in effect we have NANDs and NORs as well; this is because current is switched between two loops each of which could drive other gates (somewhat like ECL).

From the figures reported by Guéret *et al.*<sup>6</sup> we can arrive at an estimate of 500 ps for  $\tau_{gate}$  and  $\tau_{fanout}$ . Values for  $\lambda_{gate}$  are of the order of 100  $\mu m$ , with 10  $\mu m$  lines. Guéret *et al.* state that energy dissipation occurs only during switching events; the magnitude depends on the loop inductance and hence the wire length, so that this is a case of  $\epsilon_{wire}$  dissipation. In addition, however, we expect an  $\epsilon_{standby}$  component in the current source(s) driving the circuit.

Large fanouts call for long output loops, the inductance of which determines the  $\epsilon_{wire}$  dissipation. However, the inductance can be prevented from growing linearly with the fanout by

increasing the width of the line. We thus give a somewhat conservative estimate of  $0.02 fJ / \lambda_{wire}$  for present technology. The loop structure restricts  $o_{max}$  to one.

#### 4.3.6. Gallium Arsenide Logic Circuitry

Recent advances in gallium arsenide technology have paved the way for high performance combinational logic. The two classes of GaAs circuit technology that exhibit the most promise are self-aligned MESFET circuits and High Electron Mobility Transistor (HEMT) circuitry.<sup>1</sup> A more comprehensive review of the various kinds of GaAs circuitry can be found elsewhere.<sup>3</sup> Current linewidths are of the order of a micron, as with other technologies. The most complex GaAs chip known is the 8 x 8 multiplier of Lee *et al.*,<sup>10</sup> with over a thousand gates and measuring 2.7mm by 2.25 mm. A gate density of 33000 gates/ $cm^2$  was reported then, while Abe *et al.* report a somewhat higher density.<sup>1</sup> From these figures it is reasonable to assume a figure of 10  $\mu m$  for  $\lambda_{gate}$  with present technology. Gate delays approaching 10 ps have been reported in recent GaAs literature.<sup>10</sup>

#### 5. Conclusion

Before summing up, we present an example of the application of our model. The example we have chosen is a 1-of-N decoder. The decoder is made up of  $\log N$  modules stacked up. Address bit  $a_i$ ,  $1 \leq i \leq \log N$  fans out to  $2^{i-1}$  leaves; outputs  $D[i, 2k-1]$  and  $D[i, 2k]$  are produced by ANDing  $a_i$  and its complement with the input  $D[i-1, k]$  from the previous level, for  $1 \leq k \leq 2^{i-1}$ . The equations governing the operation of the circuit are:

$$D[0,1] = 1$$

$$D[i, 2k-1] = D[i-1, k] \wedge \bar{a}_i,$$

$$D[i, 2k] = D[i-1, k] \wedge a_i$$

The module for  $i=4$  is shown in Fig 1. It is clear that the output  $D[\log N, k]$  is high if and only if the binary encoding of  $k$  is  $\{a_1, a_2, \dots, a_{\log N}\}$ .

We use the obvious layout for the binary trees; more sophisticated variants<sup>12</sup> could be used to obtain smaller, faster and/or more energy-efficient decoders. It is easy to see that the width of the structure can be approximated by  $N \lambda_{gate}$  and the height by  $\approx \frac{1}{2} \log^2 N [\lambda_{gate} + \lambda_{wire}]$ . The time of operation of the circuit is  $(N + o(\log^2 N)) \frac{\lambda_{gate}}{\lambda_{wire}} \tau_{wire} + \log N \tau_{gate}$ . The

gates in each module are of two kinds - those in the fan-out tree and those at the leaves generating the outputs  $D[i,j]$ . The number of gates in fan-out tree in module  $i$  is  $2^i - 1$ ; summing over  $i$  from 1 to  $\log N$ , we find that there are  $\approx 2N$  gates in the fan-out trees. A similar argument suffices to show that there are  $2N$  gates at the outputs of each module. The energy consumed by this decoder is (in the worst case)

$$\begin{aligned} &\approx 4N \epsilon_{standby} + (2N + \log N) \epsilon_{1-0} \\ &+ (2N + 4\log N) \epsilon_{sink} + \frac{N}{4} \log^2 N \frac{\lambda_{gate}}{\lambda_{wire}} \epsilon_{wire} \end{aligned}$$

The latter terms need some explanation. The  $\epsilon_{1-0}$  component was obtained by assuming that all the address bits  $a_i$  were 1s, so that the fan-out gates each dissipated  $\epsilon_{1-0}$ ; in addition, one of the  $D[i,j]$  gates at the output of each module dissipates  $\epsilon_{1-0}$ . The last two terms are associated with the change of state of the decoder gates (this would typically occur between one computation and the next, so that we are considering these terms on a "per computation" basis). If every address bit were to change between two computations, all the fan-out gates

$E_{wire}$  is thus  $i \frac{N}{2} \frac{\lambda_{gate}}{\lambda_{wire}} \epsilon_{wire}$ . Summing over  $i$  from 1 to  $\log N$ , we obtain the above value for the fourth term.

These figures indicate that (a) this design is quite compact; (b) a technology where  $\tau_{wire}$  is small in relation to  $\tau_{gate}$  would be preferable; even so, for large  $N$  the wire delay component dominates; (c) a technology with small  $\epsilon_{wire}$  is necessary to keep the energy consumption down. An easy improvement would be to compact each of the fan-out trees such that all their gates lie in a straight line; the height then diminishes to

$$2\log N \lambda_{gate} + \frac{1}{2} \log^2 N \lambda_{wire}$$

Our discussion of the model would be incomplete without a list of the problems encountered in using it. So far, we have identified two types of circuits for which our model is inadequate.

Practical memory devices are not made of static latch circuits unless speed is of the essence. In fact, information storage is often realized by means of very much simpler circuitry - such as the single-transistor cell in dynamic MOS memories. The inability of our model to deal with such phenomena forces us to use expensive static latch circuits in our constructions. The second

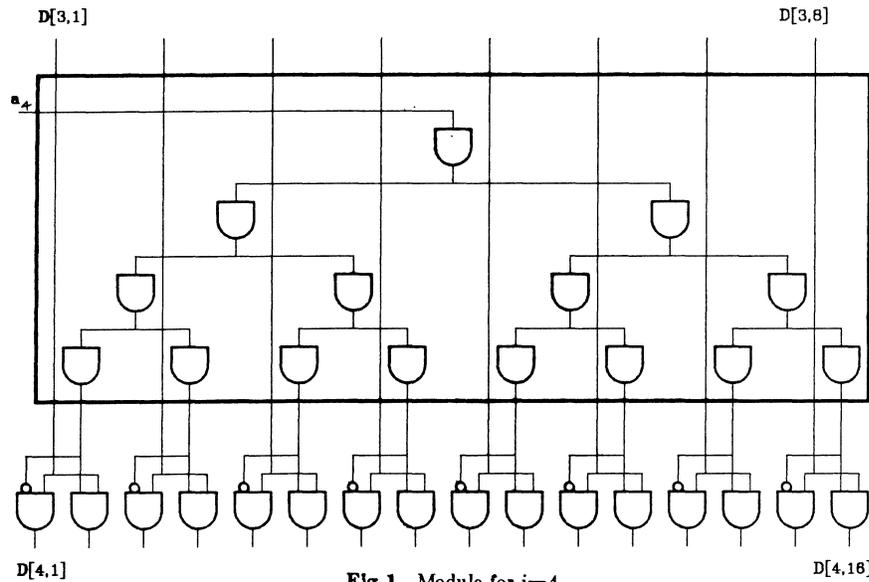


Fig 1. Module for  $i=4$

change state. In addition, two  $D[i,j]$  gates change state in each module; since each of these has two sinks, we have the  $4\log N$  coefficient for  $\epsilon_{sink}$  as well. The total wire-length in the fan-out tree in module  $i$  is  $i \frac{N}{2} \frac{\lambda_{gate}}{\lambda_{wire}}$ . This is because the gates in the bottom row of the last module are spaced at  $\lambda_{gate}$ ; in keeping with the structure of our construction, we require that the width of all the modules be  $N\lambda_{gate}$ . The worst case estimate for

major shortcoming of the model is in the treatment of switch devices, like pass-transistors in NMOS. In practice, the availability of these devices often reduces circuit area significantly.

It should be stressed that ours is a strictly "upper bound" model, and generally overestimates the area/time/energy metric. Our goal in developing this model was to be able to form estimates that are within an order of magnitude of the exact value.

Some of these conservative assumptions were made necessary by our objective of catering to several technologies, given the diversity of the underlying physical phenomena.

A major application of this model would thus consist in the evaluation of different circuit solutions for a given problem. Asymptotic predictions of circuit complexity and performance (under different metrics) can be made, subject to technological limits; for instance, it may be possible to predict that the area of the chip solving a certain problem grows as the square of the size of the input, as long as the area does not exceed  $\alpha_{max}$  and the longest wire in the circuit does not exceed  $\lambda_{max}$ .

Within the constraints of these limitations, our model permits some useful comparisons and estimates; in particular, our provision for the evaluation of energy consumption is perhaps the most general classification of energy dissipation modes in VLSI to date. Our experience in using the model has been that it avoids the minute details that complicate circuit simulation models, while retaining a more realistic picture of reality than existing graph-theoretic models.

#### References

1. Masayuki Abe, Takashi Mimura, Naoki Yokoyama, and Hajime Ishikawa, "New Technology Towards GaAs LSI/VLSI for Computer Applications," *IEEE Trans. on Electron Devices*, vol. ED-29, no. 7, pp. 1088-1093, July 1982.
2. G. Bilardi, M. Pracchi, and F. P. Preparata, "A Critique of Network Speed in VLSI Models of Computation," *IEEE Journal of Solid-State Circuits*, vol. SC-17, no. 4, pp. 696-702, August 1982.
3. Richard C. Eden, Bryant M. Welch, Ricardo Zucca, and Stephen I. Long, "The Prospects for Ultrahigh-Speed VLSI GaAs Digital Logic," *IEEE Journal of Solid-State Circuits*, vol. SC-14, no. 2, pp. 221-239, April 1982.
4. S. A. Evans, "Scaling  $I^2L$  for VLSI," *IEEE Journal of Solid-State Circuits*, vol. SC-14, no. 2, pp. 313-326, April 1979.
5. T. R. Gheewala, "Design of 2.5-Micrometer Josephson Current Injection Logic (CIL)," *IBM J. Res. Develop.*, vol. 24, no. 2, pp. 130-142, March 1980.
6. P. Guéret, A. Moser, and P. Wolf, "Investigations for a Josephson Computer Main Memory with Single-Flux-Quantum Cells," *IBM J. Res. Develop.*, vol. 24, no. 2, pp. 155-166, March 1980.
7. Paul A. H. Hart, Toon Van 'T Hof, and Francois M. Klassen, "Device Down Scaling and Expected Circuit Performance," *IEEE Journal of Solid-State Circuits*, vol. SC-14, no. 2, pp. 343-351, April 1979.
8. B. Hoeneisen and C. A. Mead, "Fundamental Limitations in Microelectronics - I. MOS Technology, II. Bipolar Technology," *Solid-State Electronics*, vol. 15, pp. 819-829, 891-897, 1972.
9. H. Ko and T. Van Duzer, "Miniaturization of Josephson Current Injection (CIL) Logic Circuits," in *International Conference on Computer Design, VLSI in Computers*, IEEE Computer Society, Nov 1983.
10. F. S. Lee, R. C. Eden, S. I. Long, B. M. Welch, and R. Zucca, "High speed LSI GaAs integrated circuits," in *Proc IEEE International Conference on Circuits and Computers*, pp. 697-700, Oct 1980.
11. Thomas Lengauer and Kurt Mehlhorn, "On the Complexity of VLSI Computations," in *VLSI Systems and Computations*, ed. H. T. Kung, Bob Sproull, Guy Steele, pp. 89-99, Computer Science Press, October 1981.
12. M. S. Paterson, W. L. Ruzzo, and L. Snyder, "Bounds on Minimax Edge Length for Complete Binary Trees," in *Proc. 19th Annual ACM Symp. on Theory of Computing*, pp. 293-299, May 1981.
13. Krishna C. Saraswat and Farrokh Mohammadi, "Effect of Scaling of Interconnections on the Time Delay of VLSI Circuits," *IEEE Journal of Solid-State Circuits*, vol. SC-17, no. 2, pp. 275-280, April 1982.
14. J. Savage, "Planar Circuit Complexity and the Performance of VLSI Algorithms," in *VLSI Systems and Computations*, ed. H. T. Kung, Bob Sproull, Guy Steele, pp. 61-68, Computer Science Press, October 1981.
15. C. D. Thompson, "A Complexity Theory for VLSI," Ph. D. Dissertation, CMU-CS-80-140, Computer Science Dept., Carnegie-Mellon University, August 1980.
16. Han-Tzong Yuan, Yung-Tao Lin, and Shang-Yi Chiang, "Properties of Interconnection on Silicon, Sapphire, and Semi-Insulating Gallium Arsenide Substrates," *IEEE Journal of Solid-State Circuits*, vol. SC-17, no. 2, pp. 269-274, April 1982.



DEPARTMENT OF ELECTRICAL ENGINEERING AND COMPUTER SCIENCE

MASSACHUSETTS INSTITUTE OF TECHNOLOGY  
CAMBRIDGE, MASSACHUSETTS 02139

Room 36-575  
October 4, 1983

Dear Author:

Congratulations on the acceptance of your paper for the 1984 conference on Advanced Research in VLSI, to be held at M.I.T. in Cambridge, January 23-25, 1984. You should already have received a letter from me telling of the decision of the Program Committee.

As you know, all papers we received were sent out for review to a panel of referees. In many cases, the referees made comments or suggestions which they felt might help the authors. I am enclosing copies of such comments about your paper with this letter, in the hope that you will find them useful when you revise your paper in its final, camera-ready form. You may use your own judgement whether or not to follow these suggestions. The camera-ready version must be on the model paper which you received and about ten pages in length.

We are looking forward to receiving your final version by November 1.

Sincerely,

A handwritten signature in cursive script that reads "Paul Penfield, Jr.".

Paul Penfield, Jr.  
Professor of Electrical  
Engineering