

# **1.0 Introduction**

Once upon a time, the delay of a circuit was well approximated as the sum of the delays of the gates. Wires could be viewed as equipotential nodes that instantly carried the voltage to all ends. As transistors have shrunk, the relative delay of gates has decreased. At the same time, the delay of wires per unit length has actually increased. Thus, wires now are a very important part of total delay.

Wire delay comes from two sources. One is the intrinsic, speed-of-light delay. The second is the lossy nature of on-chip wires; because the resistance of such wires is very high, the wires form RC circuits. Speed of light delay is proportional to the length of the wire, while RC delay increases with the square of the wire length. Thus, for long wires, RC delay dominates. For shorter wires, speed of light delay still dominates. However, such delay over a short wire is still relatively small compared to gate delay.

In addition to increasing circuit delay, wires are subject to coupling noise. According to Maxwell's equations, changes in voltages and currents create electric and magnetic fields which in turn cause changes in voltages and currents on nearby conductors. These effects are called capacitive and inductive crosstalk, and get worse as conductors move closer together and edge rates get sharper. They now can greatly exceed the noise margins of many circuits unless interconnect is carefully designed.

We will first examine the RC wire effects which are most important in current circuits. We'll see how to compute wire resistance and estimate wire capacitance, then use these two properties in wire delay models. We'll also look at capacitive crosstalk and circuit failures caused by IR drops, electromigration, and self-heating. Since delay and crosstalk are so severe, we will explore how to optimize wire geometries, add repeaters, and drive busses to improve performance.

In the next chapter, we'll consider how RC effects scale with decreasing feature size. Then we'll address inductance, an effect which becomes more significant for wire pitches in the one micron range. We'll look at inductive delay effects, di/dt noise, and inductive crosstalk and mutter some banal warnings since little concrete is known about the subject.

### 2.0 Wire Resistance

The resistance of a wire can easily be computed from its resistivity  $\rho$ , length L, width W, and thickness t:

$$\mathbf{R} = \rho \mathbf{L} / \mathbf{W} \mathbf{t} \tag{EQ 1}$$

Since the thickness is not under designer control, resistance is frequently rewritten as  $R_{\Box}$  L/W, where  $R_{\Box}$  is the sheet resistance having units of ohms / square. Typical values for various interconnect levels are:

Material	R
diffusion	10
polysilicon	10
thin metal (m1, m2,)	0.07
thick metal (top level)	0.03

**TABLE 1. Typical Interconnect Sheet Resistances** 

Without special treatment, diffusion and polysilicon would have much higher resistances that might degrade performance. Therefore, they are usually covered by a thin layer of a refractory metal such as tungsten which alloys with the silicon to form a "silicide" and reduce resistance to approximately 10 ohms/square. Even so, polysilicon should only be used for short interconnect within cells because it is still rather resistive. Diffusion should never be used for routing in high-performance circuits because it has tremendous parasitic capacitance as well as resistance.

Lower metal layers are usually thin, allowing very tight spacing but higher sheet resistance. Metal 1 can be especially thin and resistive. Higher metal layers intended for power distribution and critical signals are thicker and built on a wider pitch, reducing sheet resistance. Most wires are aluminum, but extensive research on other materials with slightly lower resistivity such as copper and gold are beginning to pay off in commercial processes. Such materials can have nasty effects if they contaminate transistors, so exotic schemes are used to prevent their infiltration of the substrate.

Every time a signal crosses metal layers, it must pass through a contact cut, adding further resistance. Such cuts generally have 1-10 ohms of resistance. Low-resistance connections therefore use an array of many cuts.

# 3.0 Wire Capacitance

Wires also have capacitance to the substrate and to other nearby wires. A handy rule of thumb is that most wires have about 0.2 fF/ $\mu$ m of total capacitance, about 1/10 that of a transistor gate. Therefore, when estimating sizes, one can treat 10  $\lambda$  of interconnect as about 1  $\lambda$  of gate loading.

In the good old days, wires were much wider than thick. They could therefore be approximated as having parallel plate capacitance to the grounded substrate. Now, wires often have comparable thickness and width and are closer to adjacent wires than to layers above or below. Moreover, the layers above and below are much closer than the substrate. Thus, a negligible part of a wire's capacitance is to ground; instead it is to nearby wires which may transition and couple noise into the hapless wire.

These capacitances are shown in Figure 1:



If all the capacitance is from a conductor to a solid ground plane, it can be represented as the sum of a parallel plate capacitance on the bottom and fringing capacitance on the sides. There are many physical and empirical models for such computing such capacitances relatively accurately<sup>1</sup>. Unfortunately, this situation is rare in actual designs; instead, coupling occurs to neighbors on both sides and on the top and bottom. No closed-form equations are possible in general; instead, capacitances must be extracted with a 2 or 3-dimensional field solver using Poisson's equation. Many CAD tools avoid this time-consuming process by pre-extracting tables of capacitances for various geometries and interpolating between the tabulated results closest to the actual geometries of the circuit. Such a table for a 0.8 micron 4-level process can be found at:

```
http://www-leland.stanford.edu/class/ee371/lib/wirecap
```

Even in such a process, over 75% of the capacitance of a metal3 line over substrate can be to adjacent lines because the adjacent lines are much closer than the substrate. In more advanced processes, the coupling to adjacent lines gets even larger as the lines get closer. In a later section will analyze the amount of crosstalk such coupling can cause.

### 4.0 Wire Delay Modeling

Wires contribute delay from both the speed-of-light limit and the RC delay. For long wires, the RC term dominates because it increases quadratically with length while speed-of-light delay increases linearly. For short wires, speed-of-light would matter, except that

<sup>1.</sup> See E. Barke, "Line-to-Ground Capacitance Calculation for VLSI: A Comparison," IEEE Transactions on Computer-Aided Design, Vol. 7, No. 2, Feb. 1998, pp. 295-298.

the wires are short enough that the delay is usually assumed to be negligible. This assumption will change in the future as picoseconds of delay cease to be negligible.

### 4.1 Models

The distributed RC delay can be modeled by breaking up a wire into one or more segments and using a lumped model for each segment. If each segment has resistance R and capacitance C, popular lumped models include:

#### FIGURE 2. Lumped models of wires



Note that R and C depend on the length of the segment. The Pi model and T model, named for their shapes, are both good models. The Pi model can be slightly more convenient because it has fewer circuit nodes. The L model is much less accurate, requiring many more segments for the same accuracy of delay estimation. For example, the delay of a line can be modeled to within 5% accuracy with 4 Pi segments or 64 L segments! Thus, it is usually reasonable to model long wires with 4 Pi segments in simulation and with a single Pi segment for hand calculation.

The delay of the driver must also be included when estimating the delay of a wire. For example, we can compute the delay of an inverter driving an L mm wire with resistance R/ mm and capacitance C/mm, inverter drive resistance  $R_{inv}$ , drain capacitance  $C_{diff}$ , and total load on the wire  $C_{load}$ :

#### FIGURE 3. Delay Estimation Example



The Elmore delay of the model is:

$$R_{inv}^{*}(C_{diff} + CL/2) + (R_{inv} + RL)^{*}(CL/2 + C_{load})$$
 (EQ 2)

This can be rewritten as:

$$R_{inv}^{*}(C_{diff} + C_{load} + CL) + RCL^{2/2} + RLC_{load}$$
(EQ 3)

This expression has three terms. The first represents the inverter driving it's own diffusion, the load capacitance, and the lumped capacitance of the wire. The second is the quadratic delay of the wire's self-loading. It is proportional to RC/2, rather than RC, because the capacitance is distributed along the wire rather than lumped at the end. The third term is the extra delay contributed by the wire resistance discharging the load capacitance.

Remember that the units in this equation are somewhat curious.  $R_{inv}$  is  $\Omega$ . R is  $\Omega^*\mu m$ . C<sub>diff</sub> and C<sub>load</sub> are fF, while C is fF/ $\mu m$ .

### 5.0 Capacitive Crosstalk

The RC models presented above assume that all capacitance is to ground, or some other node which is not transitioning. When the capacitance couples to another wire which might be switching, the crosstalk between wires causes noise and change in delay.

When two adjacent wires switch opposite directions, their delay is longer than if one switched while the other remained steady. This is because the capacitance between the wires sees a total of twice as much voltage change and thus requires twice as much current to be charged up. We can model this case for delay with a capacitor to ground twice as big as the capacitance between lines. For example, if a wire had 0.1 pF/mm capacitance to ground and 0.15 pF/mm capacitance/mm to an adjacent wire which might switch in the opposite direction, it can be modeled as a total of 0.4 pF/mm of capacitance to ground for delay purposes. Conversely, if the adjacent wire is switching in the same direction, the coupling capacitor will see zero voltage change and have no effect on delay. In the previous example, it could be modeled as 0.1 pF/mm of total capacitance to ground. At first, this seems good because the delay will be less. However, the delay of the wire becomes much less predictable because its effective capacitance changes by a factor of 4 depending on what the neighboring line does. If delays of two signals must be matched for some reason, the variation in delay caused by such coupling is very important.





Worse yet, the capacitive crosstalk introduces noise that could corrupt a signal. We can conservatively model crosstalk noise with a capacitive voltage divider. Suppose victim

signal A is supposed to be 0. Adjacent aggressor signal B rises. The coupling between B and A tends to cause A to rise as well. The amount A rises depends on the ratio of coupling capacitance  $C_{AB}$  to total capacitance  $C_{A-GND} + C_{AB}$ . It is also influenced by the driver of A, which will fight the crosstalk, trying to maintain A at a 0. If the crosstalk happens fast enough that the capacitive impedances are much lower than the driver resistance, the driver will not help much and the voltage noise on A is:

$$C_{AB} / (C_{AB} + C_{A-GND}) * V_{DD}$$
(EQ 4)

When  $C_{AB} = C_{A-GND}$ , the noise will be 50% and could cause a false transition on a gate at the end of the wire. Most static circuits can tolerate such noise glitches. Dynamic circuits, however, will be corrupted irrecoverably. Therefore, dynamic nodes are especially sensitive to coupling. Moreover, in a bus, a wire can receive coupling from both sides. In fact, it is even possible to receive coupling from orthogonal wires on layers above and below. Usually, this coupling is negligible. However, if a wide bus above or below switches entirely in one direction, the coupling can be significant. Most of the time, this coupling is ignored by verification tools. One day, chips may start to fail in a data-dependent way from such coupling on other layers!

When a node is actively driven, the small driver resistance does reduce the amount of coupling noise. Simulation may be the only way to determine exactly how much the resistance helps.

## 6.0 Other Wire Effects

In addition to delay and coupling noise, wires are susceptible to IR drops, electromigration, and self-heating.

### 6.1 IR Drops

Wires such as power busses which carry large DC currents droop in voltage due to the resistance of the wire multiplied by the current carried on the wire. This is bad because, for instance, it reduces the value of  $V_{DD}$  seen at the center of a chip when the power supply is driven by pads around the perimeter.

IR drops get to be a bigger problem as lower power supplies are used. For a fixed wattage, a lower  $V_{DD}$  means a higher  $I_{DD}$  flows through the power distribution network. For the same power network resistance, a higher IR drop is caused. Moreover, this drop is even a larger fraction of  $V_{DD}$ !

For example, a 100 watt chip operating at 2V requires 50 amps of supply current! If IR drop should be no more than 5% of VDD, this requires a power network with only 2 m $\Omega$  of resistance. To meet the same spec, a 1V chip requires 100 amps of current and can only tolerate 0.5 m $\Omega$  of resistance in the power distribution network! Remembering that a typical top-level metal layer has a resistance of about 30 m $\Omega$  / square, the power network can only have a small fraction of a square of size. Even a solid power plane may be insuffi-

cient; instead power may be distributed in a thick plane in the package and dropped all over the chip with C4 bumps to reduce the resistance.

### 6.2 Electromigration

When high currents flow through a wire, they create an "electron wind" which can apply force on metal atoms in the wire. The atoms may actually move slightly, in an effect known as electromigration. If electromigration continues for a long time, it may actually cause a break in a wire, resulting in an open circuit and chip failure. This is especially unpleasant because it may happen after years of operation.

Such failure is avoided by limiting the current flowing through a wire. DC current is worst because it pushes atoms in a single direction. AC currents are less of a problem for electromigration because they push the atoms back and forth without a net average movement. Thus, power and ground lines are most prone to electromigration failure. A rule of thumb for many processes is to allow no more than 1 mA of current per micron of wire width. Wider wires may tolerate somewhat higher current densities than this. Processes with very thin wires tolerate less current. Check with your vendor or be conservative with wire widths.

### 6.3 Self-heating

IR drops in wires cause the wires to heat up. Oxide is a good insulator, so the heat may build up in the wires, resulting in the wires being hotter than the chip. This can increase wire resistance and degrade performance. In extreme cases, it is possible to experience thermal runaway, in which wires get hotter, thus become more resistive, causing larger IR drops and more heat until the wire melts down!

Such self-heating sets minimum wire widths even for lines carrying only AC signals.

### 7.0 Optimizing Wire Geometry

We have seen that wire resistance and capacitance cause significant amounts of delay and noise. Thus, optimizing wires can be as important as optimizing gates.

The best way to optimize wires is to keep them as short as possible. This requires careful floorplanning and good exploitation of datapath structure. Nevertheless, it is inevitable that a chip will have some long wires. What can we do to minimize delay and noise of these wires?

Wire design is fairly simple because there are only a few degrees of freedom. The designer can control the width of each wire, the spacing between wires, and the arrangement of which wires are adjacent to which other wires.

To reduce delay, the designer may adjust each of these factors. If the delay is dominated by the wire RC product, increasing wire width is often useful. This decreases wire resistance

proportionally, while only slightly increasing total wire capacitance if the capacitance from the top and bottom of the wire is a small part of the whole. Increasing spacing between wires is also good because it decreases the total wire capacitance and the amount of coupling to adjacent wires. This is especially important when the adjacent wires might be switching and Miller-multiplying the coupling capacitance. Finally, interleaving busses such that the adjacent wires are not switching at the same time prevents Miller-multiplication of coupling capacitance. For example, a bus driven in phase 1 could be interleaved with a bus driven in phase 2 so no two adjacent wires switch at once.

The designer also may use these degrees of freedom to reduce crosstalk. The first tactic is to increase the spacing between lines, thus reducing coupling. At some point, this yields diminishing returns. The next option is to increase the width of the wire or to shield the wire with power/ground lines. Increasing the width of a wire is sometimes helpful because it increases the fraction of capacitance to the layers above and below rather than to an adjacent coupling wire. It also reduces the RC delay. In contrast, shielding the wire with a nearby minimum pitch power or ground line is extremely effective because it eliminates almost all capacitive coupling. It also increases the number of wires in the power and ground grids which may improve the power supply network. The drawback is that it may increase the total amount of capacitance on a wire and thus increase RC delay. Shielding every bit of a wide dynamic but also has horrible effects on area.





As mentioned earlier, the capacitance to adjacent wires and to the upper and lower layers is a complex function of wire width and spacing. It may be necessary to use a table of capacitances to understand how width and spacing influence noise and delay. The designer faces a constant tradeoff between delay and area; there is no single pitch optimum for all problems.

# 8.0 Repeaters

Even when a wire is well optimized to minimize its resistance and capacitance per unit length, wire delay still increases quadratically with distance. At some point, it becomes better to break a long wire in half and insert a buffer, called a repeater, to drive the second half. Each half of the wire contributes 1/4 of the RC delay, thus reducing total RC delay by a factor of 2. The buffer, however, introduces extra delay. The best wire length for introducing a repeater therefore is when the repeater delay is less than the amount of wire delay saved. For very long lines, breaking the line into more than 2 pieces and adding multiple repeaters becomes beneficial.

By splitting a wire every time it gets too long and adding a repeater, delay can be made to be proportional to the length, instead of the square of length. This results in a "effective speed of light" in wires driven with optimal repeater placement. The effective speed scales with transistor performance because faster transistors mean faster repeaters and shorter optimal wire lengths between repeaters. In Appendix A, we prove that the effective speed of light is:

$$\sqrt{C_W R_W C_T R_T} (2 + \sqrt{2}) \tag{EQ 5}$$

where  $C_W$ ,  $C_T$ ,  $R_W$ , and  $R_T$  are the capacitances and resistances per unit length of wires and transistors.

For a 0.4 micron technology, this delay is around 50 ps/mm using optimal repeater placement of inverters with 23 micron NMOS devices every 2.5 mm of wire. The repeater placement gets more frequent as transistors get faster and wires get slower. This compares with about 6 ps/mm for the true speed of light on a chip. Increasing wire pitch or building faster transistors improves the RC delay.

When the final load is comparable to the capacitance of a repeater, it is interesting to see that each repeater is sized identically. The tapered buffer chain concept of inverters increasing exponentially in size along the chain is not appropriate because the goal is not to power up to drive a huge load, but rather to break the wire into lengths that can be driven quickly. In fact, raising the size of the repeaters will actually slow the circuit down as well as hurt area and power. Knowing the appropriate repeater sizing and placement for your process is important.

There are several tricky issues with actually using repeaters. Unless an even number of inversions are always used, the logic at the receiving end may get the wrong polarity of signal. Repeaters must also be placed in the layout; this may require extensive communication between designers so that the repeaters for cross-chip interconnect are placed in suitable locations. When long interconnect forks to multiple destinations, such as in critical stall signals, repeater placement to minimize delay while providing the right polarity becomes even trickier.

# 9.0 Advanced Signaling Techniques

There are many tricks to improve interconnect performance further, but each introduces risks. Precharging busses can be very fast and integrates well with dynamic circuits, but is susceptible to coupling noise. Small swing busses can be even faster and more noise-prone. Bi-directional busses may be beneficial for area, but make using repeaters difficult.

#### 9.1 Precharged Busses

Precharged busses are fast for the same reasons precharged gates are fast. They can be driven by fast gates, improving repeater performance. The repeaters can also have switching thresholds adjusted to favor the critical edge, responding after a smaller signal swing.

Unfortunately, precharge busses are not actively driven while in the precharge state. Therefore, crosstalk noise is especially severe. Even if they have a small keeper, the keeper is too weak to cancel much of the crosstalk.

Hence, many designers limit use of precharged busses to within a unit so that the designer has complete control over the coupling sources. Good routing and verification tools are necessary to support global precharged busses.

#### 9.2 Small Signal Swings

The delay of a heavily capacitive system is set by the time charging the capacitance:

Z

$$\Delta t = C/I \,\Delta v \tag{EQ 6}$$

Hence, reducing the  $\Delta v$  required to indicate a logic level change also reduces the time required to propagate the information. This trick is especially popular in SRAM arrays and on I/O pads. In RAM arrays, a differential sense amplifier detects a small change in voltage between the bitline and its complement. If the gain of the amplifier is sufficiently high, only enough voltage difference must develop to overcome offset voltages in the amplifier. Similarly, many low-voltage I/O techniques use swings of one volt or less.

Coupling is especially serious when small swing signals are located near full swing signals, because even a small amount of coupling from a full swing line onto a small signal line could upset the value on the victim line. I/O pads avoid this problem if all external I/O signals are small swing. SRAM bitlines are often twisted in such a way as to guarantee the same coupling onto the true and complementary bitlines so that the net differential noise is zero.

#### 9.3 Bidirectional Busses

Sometimes area can be reduced by sharing a single bus among many seldom-used drivers. Each driver uses a tri-state circuit to only drive the bus when it is granted control.

Such a design has many drawbacks. Each driver must be large enough to drive the heavily loaded bus. The drivers which are not on therefore present large parasitic drain capacitances, slowing the active driver. Moreover, building repeaters on the bus is difficult because the driver might be located at either end of the wire. Therefore, the repeaters cannot be inverters which only drive one direction. Instead, they must be bidirectional circuits which sense a transition starting and pass it on. Such circuits are slower and more prone to noise than regular repeaters. In general, bidirectional busses should be avoided in high performance systems. Point to point links should be used instead. When many seldom-used drivers must be connected to many receivers, the drivers can each drive point-to-point links to a central location where a multiplexor selects the appropriate source and drives it to all the receivers.

## **Appendix A Derivation of Repeater Sizing**

Delay of a single long wire in an integrated circuit is proportional to the square of the length because both the resistance and capacitance increase linearly with length. By breaking the wire into multiple stages and placing repeaters at the start of each stage, the delay can be made linear with length. This appendix explores the optimal sizing of buffers and number of repeaters to minimize the delay. It computes delay ( $ps/\mu m$ ) as a function of process parameters at the optimal sizing.

We begin with several simplifying assumptions to make the analysis cleaner.

- Source / drain diffusion capacitance is negligible
- Inverters are sized for equal rise and fall times (faster results may be possible with unequal times)
- Wire pitch and spacing are preset
- Propagation speeds are low enough that transmission line effects may be neglected
- Load capacitance equals repeater capacitance

The interconnect of total length L is divided into S segments. Each inverter has a pulldown N microns wide and a pullup  $\beta$ N microns wide to achieve equal rise and fall times. The figure below illustrates a wire divided into three segments:

#### FIGURE 6. Interconnect divided into 3 segments



The path can be broken into S identical segments for analysis. A model of one segment is shown below. Resistances are in ohms / micron; capacitances are in pF / micron. Notice that  $C_T$  is the combined capacitance of the PMOS and NMOS transistors in an inverter sized at N=1.

FIGURE 7. Model of single segment



Observing that a distributed RC line has the same delay as a lumped RC element with half the resistance, we can derive the following expression for the delay through a segment. The delay through the entire interconnect is S times the delay through each segment:

$$t_{interconnect} = S \left[ C_T N \left( \frac{R_T}{N} + R_W \left( \frac{L}{S} \right) \right) + C_W \left( \frac{L}{S} \right) \left( \frac{R_T}{N} + \frac{R_W}{2} \left( \frac{L}{S} \right) \right) \right]$$
(EQ 7)

simplifying, 
$$t_{interconnect} = SC_T R_T + L \left( C_T R_W N + \frac{C_W R_T}{N} \right) + L^2 \left( \frac{C_W R_W}{2S} \right)$$
 (EQ 8)

Take partial derivatives with respect to N and S to minimize total delay:

$$\frac{\partial t_{interconnect}}{\partial N} = C_T R_W L - \frac{C_W R_T L}{N^2} = 0 \Longrightarrow N = \sqrt{\frac{C_W R_T}{C_T R_W}}$$
(EQ 9)

$$\frac{\partial t_{interconnect}}{\partial S} = C_T R_T - \frac{C_W R_W L^2}{2S^2} = 0 \Rightarrow S = L \sqrt{\frac{C_W R_W}{2C_T R_T}}$$
(EQ 10)

Hence, the total delay is:

$$L_{interconnect} = L_{\sqrt{\frac{C_{W}R_{W}C_{T}R_{T}}{2}} + L(\sqrt{C_{W}R_{W}C_{T}R_{T}} + \sqrt{C_{W}R_{W}C_{T}R_{T}}) + L_{\sqrt{\frac{C_{W}R_{W}C_{T}R_{T}}{2}}(EQ+1)$$

$$= L_{\sqrt{C_{W}R_{W}C_{T}R_{T}}}(2 + \sqrt{2})$$
(EQ 12)

and the optimal length of each segment is:

$$\sqrt{\frac{2C_T R_T}{C_W R_W}} \tag{EQ 13}$$