



EECS 151/251A Spring 2018 Digital Design and Integrated Circuits

Instructors:
Nick Weaver & John Wawrzynek

Lecture 10

What do ASIC/FPGA Designers need to know about physics?

► Physics effect:

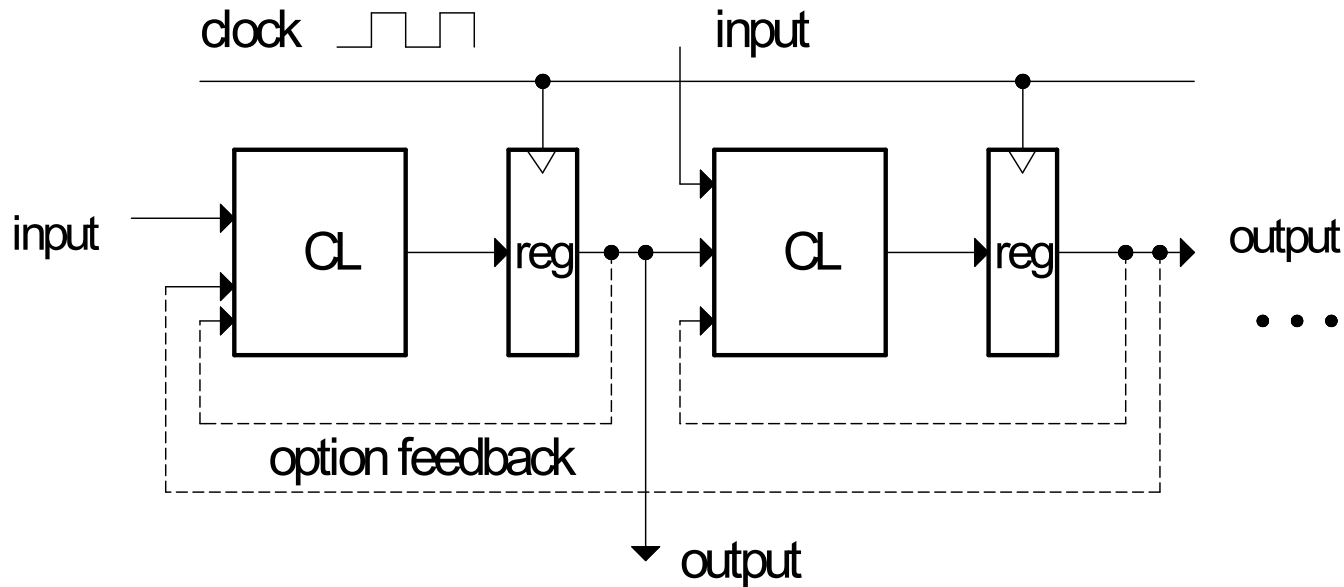
Area \Rightarrow cost

Delay \Rightarrow performance

Energy \Rightarrow performance & cost

- Ideally, zero delay, area, and energy. However, the physical devices occupy area, take time, and consume energy.
- CMOS process lets us build transistors, wires, connections, and we get capacitors, inductors, and resistors whether or not we want them.

Performance, Cost, Power



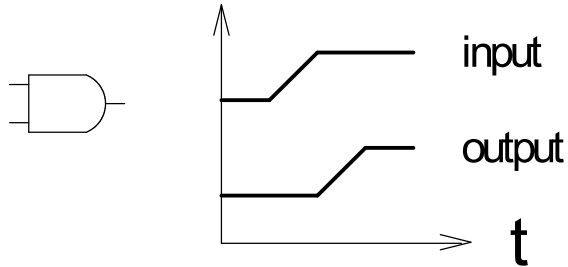
- How do we measure performance?
operations/sec? cycles/sec?
- Performance is directly proportional to clock frequency.
Although it may not be the entire story:

Ex: CPU performance

$$= \# \text{ instructions} \times \text{CPI} \times \text{clock period}$$

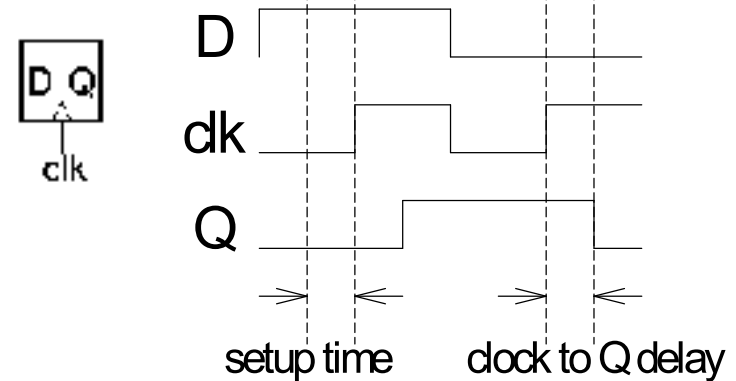
Limitations on Clock Rate

1 Logic Gate Delay



What are typical delay values?

2 Delays in flip-flops

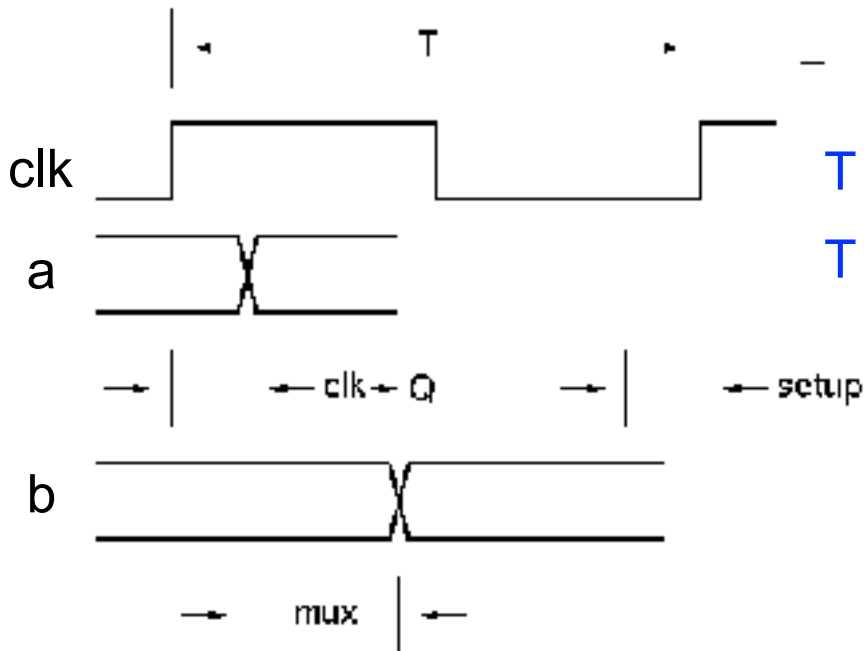
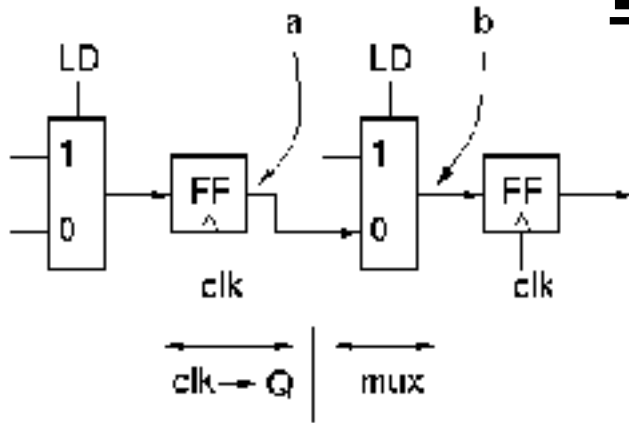


Both times contribute to limiting the clock period.

- What must happen in one clock cycle for correct operation?
 - All signals connected to FF (or memory) inputs must be ready and "setup" before rising edge of clock.
 - For now we assume perfect clock distribution (all flip-flops see the clock at the same time).

Example

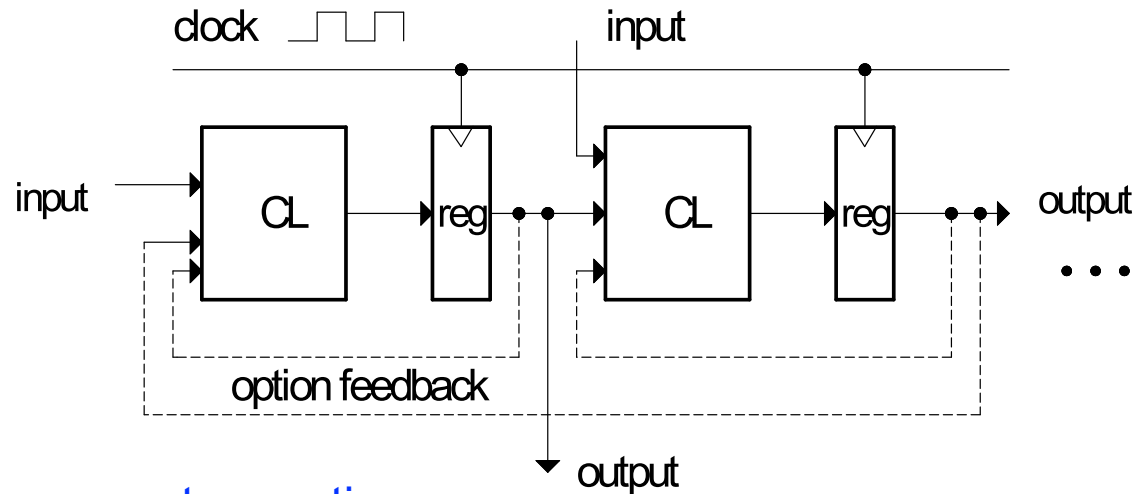
Parallel to serial converter circuit



$$T \geq \text{time}(\text{clk} \rightarrow Q) + \text{time}(\text{mux}) + \text{time}(\text{setup})$$

$$T \geq \tau_{\text{clk} \rightarrow Q} + \tau_{\text{mux}} + \tau_{\text{setup}}$$

In General ...

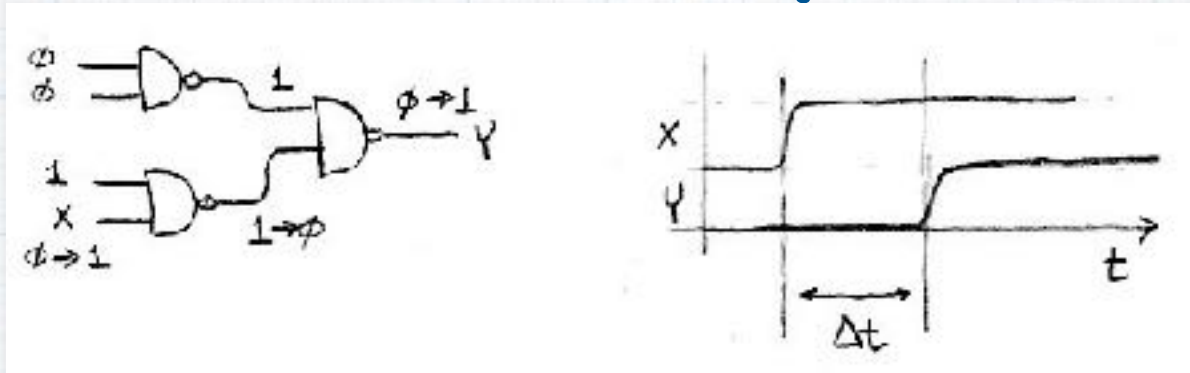


For correct operation:

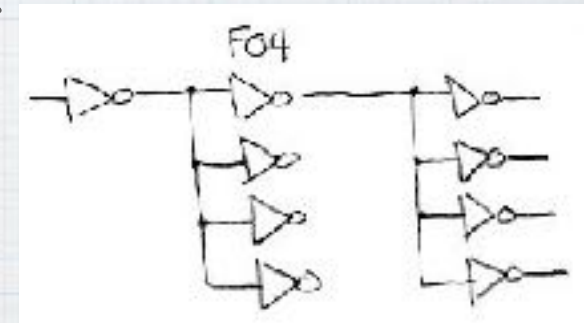
$$T \geq \tau_{\text{clk} \rightarrow \text{Q}} + \tau_{\text{CL}} + \tau_{\text{setup}} \quad \text{for all paths.}$$

- How do we enumerate **all** paths?
 - Any circuit input or register output to any register input or circuit output?
- Note:
 - "setup time" for outputs is a function of what it connects to.
 - "clk-to-q" for circuit inputs depends on where it comes from.

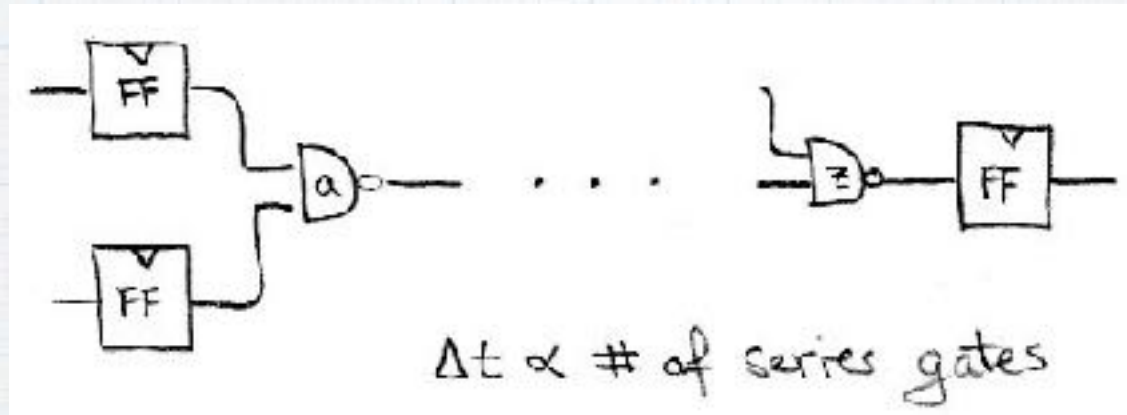
"Gate Delay"



- ▶ Modern CMOS gate delays on the order of a few picoseconds. (However, highly dependent on gate context.)
- ▶ Often expressed as F04 delays (fan-out of 4) - as a process independent delay metric:
 - ▶ the delay of an inverter, driven by an inverter 4x smaller than itself, and driving an inverter 4x larger than itself.
 - ▶ For a 90nm process F04 is around 20ps. Less than 10ps for a 32nm process.



“Path Delay”



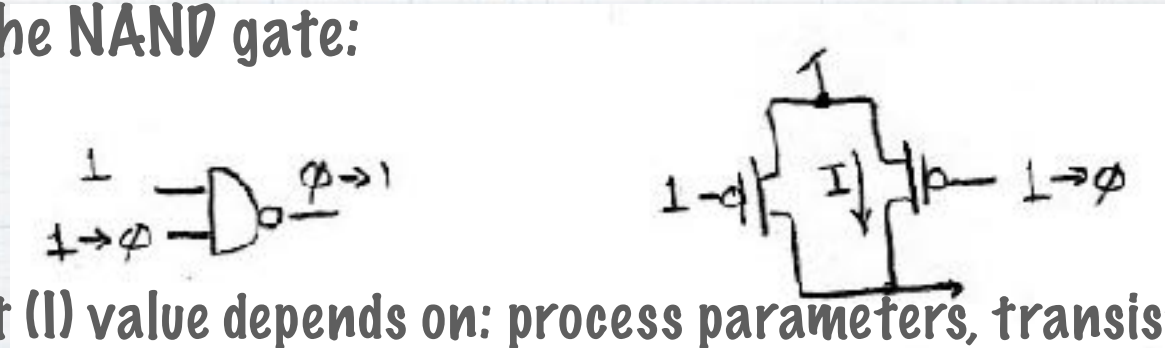
- ▶ For correct operation:

Total Delay \leq **clock_period** - **FF_{setup_time}** - **FF_{clk_to_q}**
on all paths.

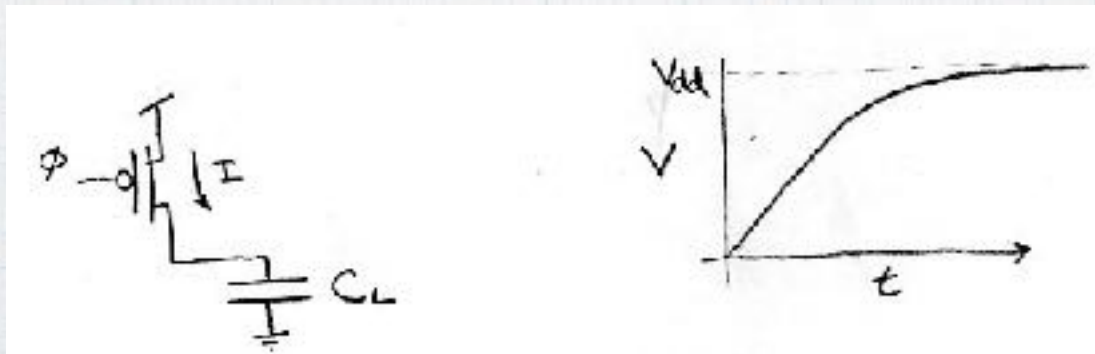
- ▶ High-speed processors critical paths have around 20 F04 delays.

"Gate Delay"

- ▶ What determines the actual delay of a logic gate?
- ▶ Transistors are not perfect switches - cannot change terminal voltages instantaneously.
- ▶ Consider the NAND gate:



- ▶ Current (I) value depends on: process parameters, transistor size

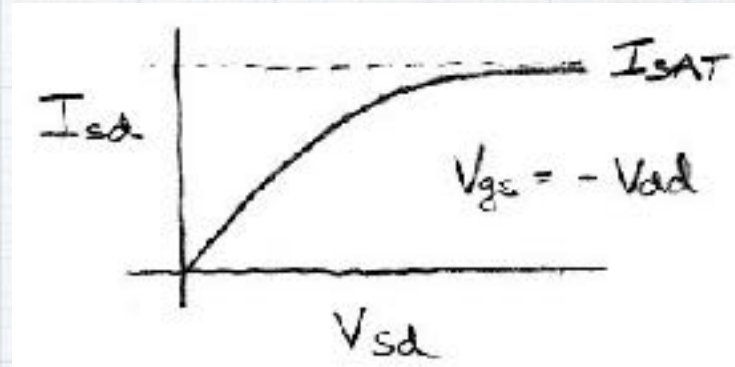
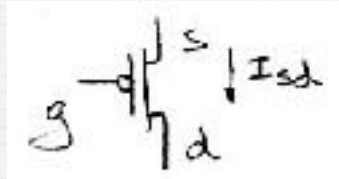


$$\Delta \propto C_L / I$$

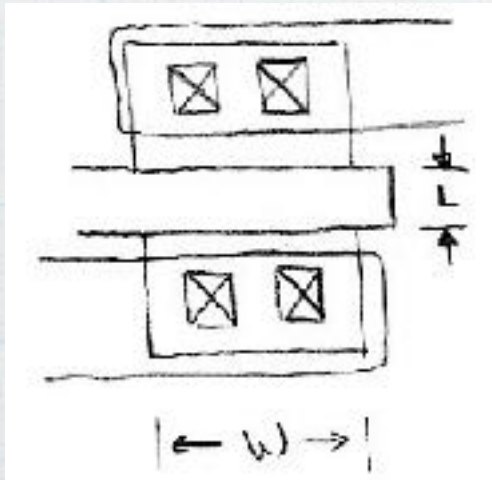
- ▶ C_L models gate output, wire, inputs to next stage (Cap. of Load)
- ▶ C "integrates" I creating a voltage change at output

More on transistor Current

- ▶ Transistors act like a cross between a resistor and “current source”

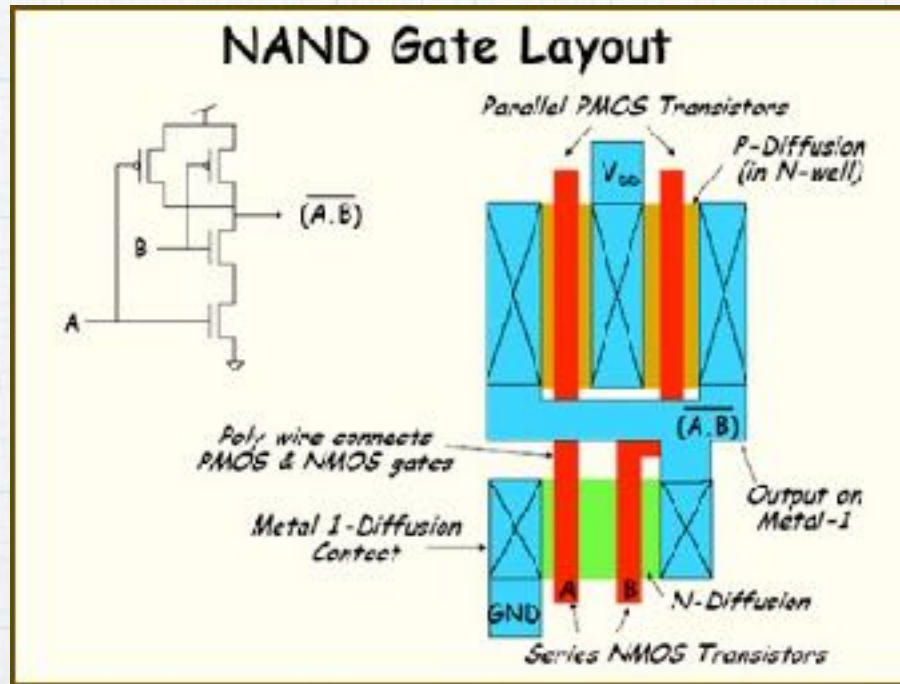


- ▶ I_{SAT} depends on process parameters (higher for nFETs than for pFETs) and transistor size (layout):



$$I_{SAT} \propto W/L$$

Physical Layout determines FET strength

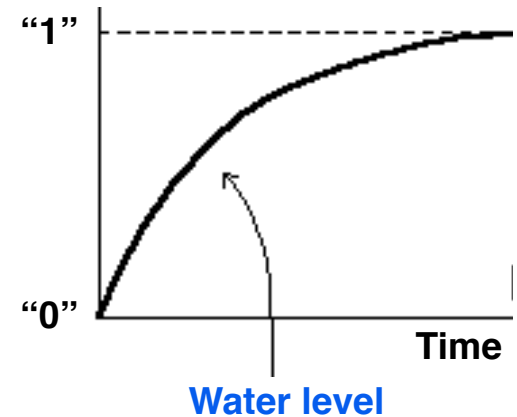
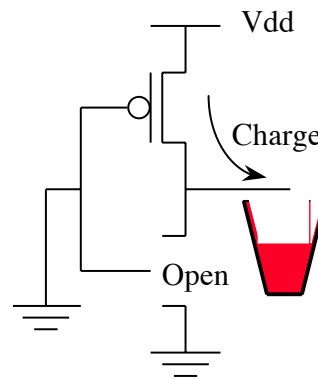


- ▶ “Switch-level” abstraction gives a good way to understand the function of a circuit.
 - ▶ nFET (g=1 ? short circuit : open)
 - ▶ pFET (g=0 ? short circuit : open)
- ▶ Understanding delay means going below the switch-level abstraction to transistor physics and layout details.

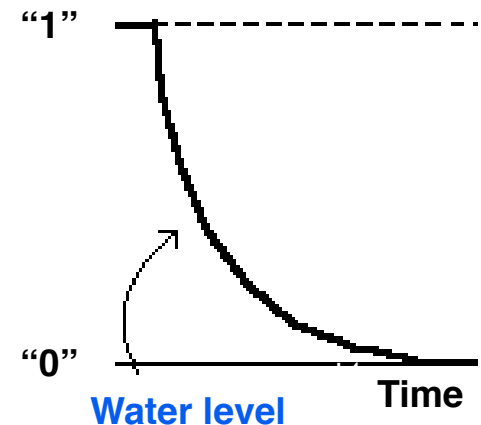
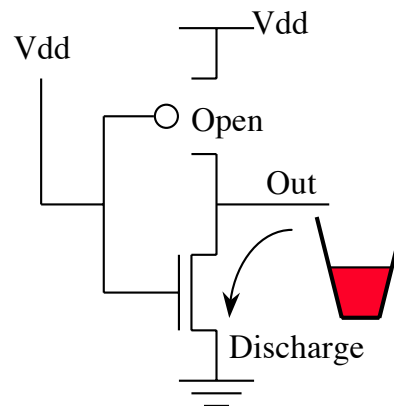
Transistors as water valves. (Cartoon physics)

If **electrons** are water molecules,
transistor strengths (W/L) are pipe diameters,
and **capacitors** are buckets ...

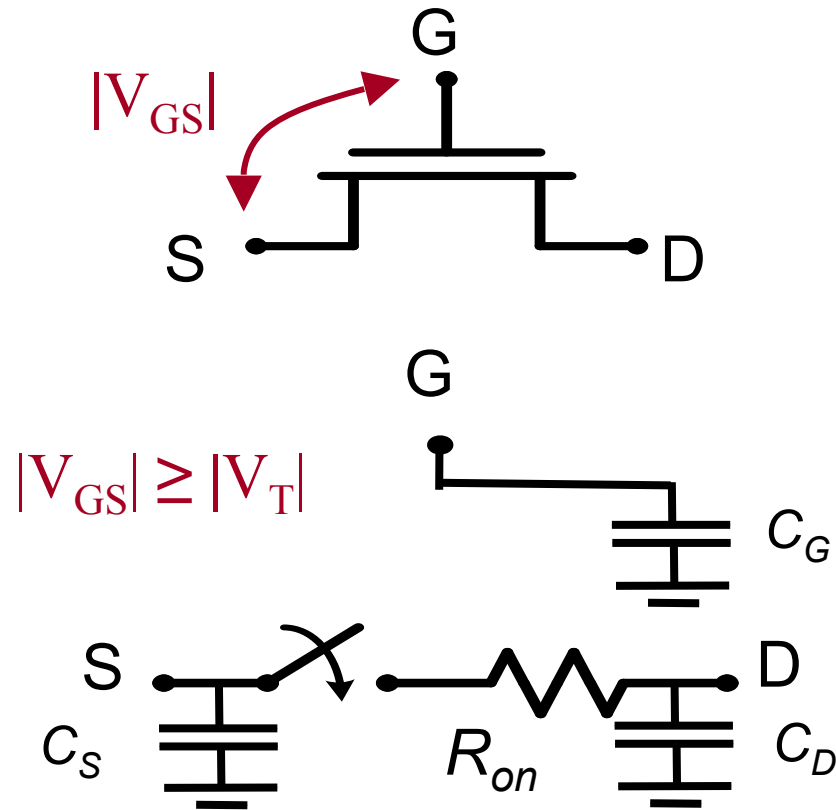
A “on” p-FET fills
up the capacitor
with charge.



A “on” n-FET
empties the bucket.

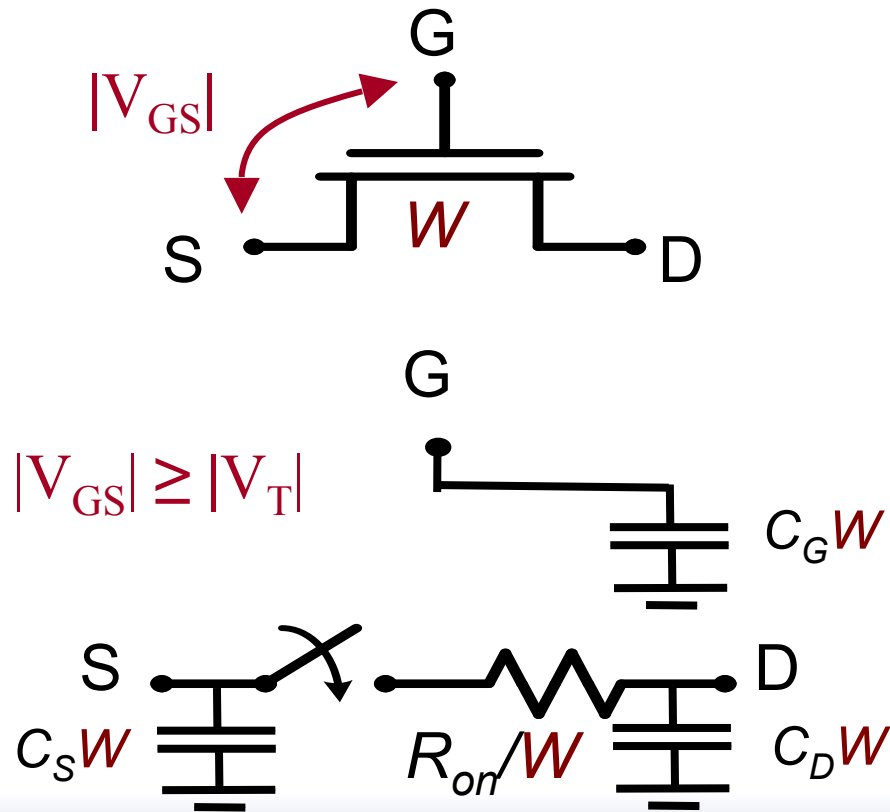


The Switch – Dynamic Model (Simplified)



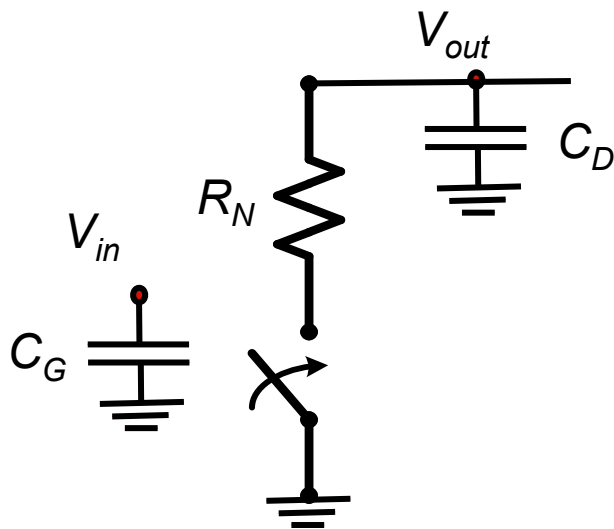
Switch Sizing

What happens if we make a switch W times larger (wider)

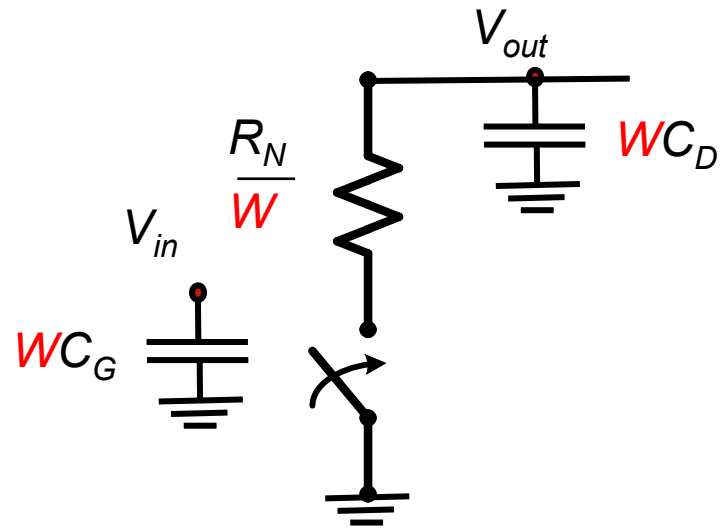


Switch Parasitic Model

The pull-down switch (NMOS)



Minimum-size switch

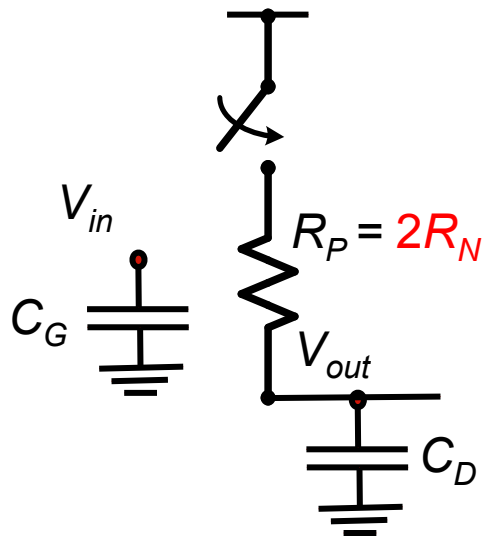


Sizing the transistor (factor W)

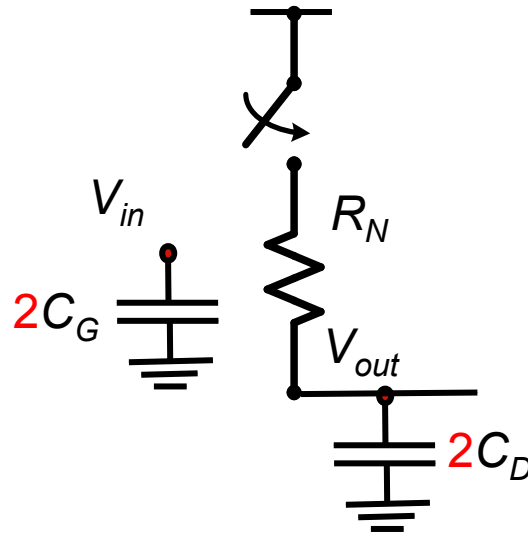
We assume transistors of minimal length (or at least constant length). R 's and C 's in units of per unit width.

Switch Parasitic Model

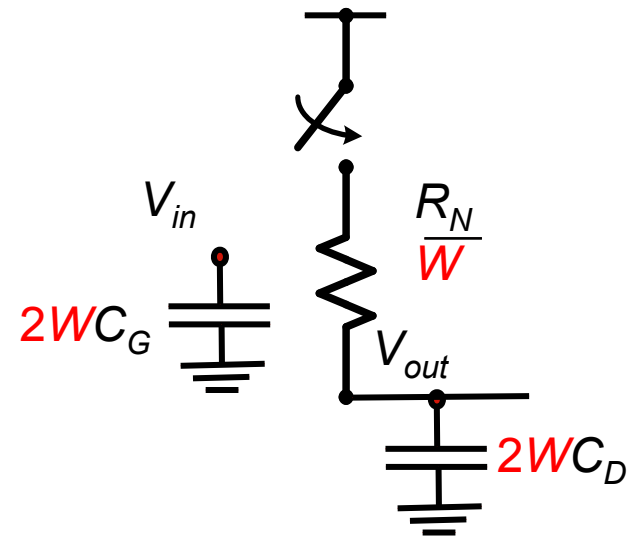
The pull-up switch (PMOS)



Minimum-size switch

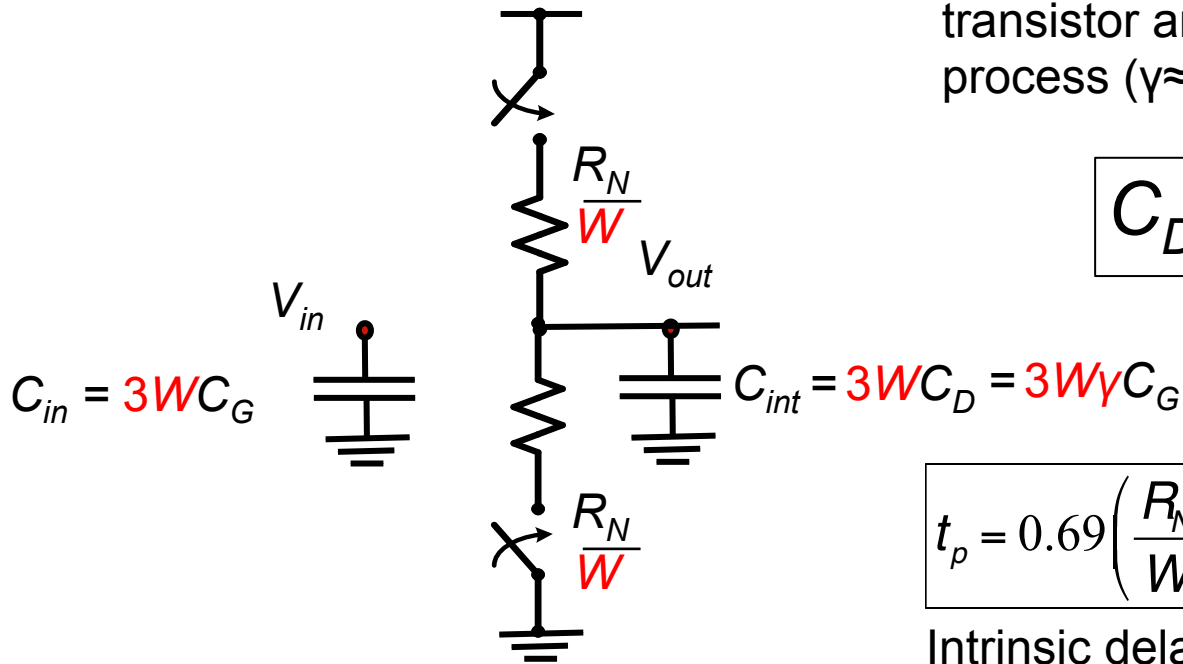


Sized for symmetry



General sizing

Inverter Parasitic Model



Drain and gate capacitance of transistor are **directly** related by process ($\gamma \approx 1$)

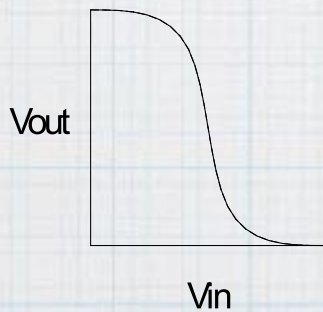
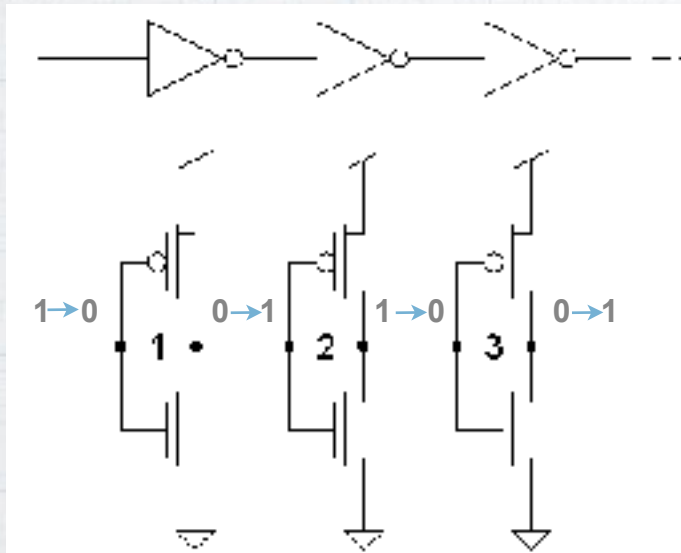
$$C_D = \gamma C_G$$

$$t_p = 0.69 \left(\frac{R_N}{W} \right) (3W\gamma C_G) = 0.69(3\gamma) R_N C_G$$

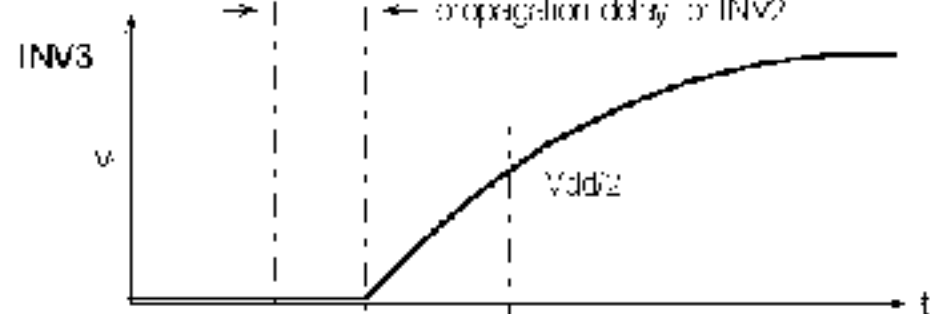
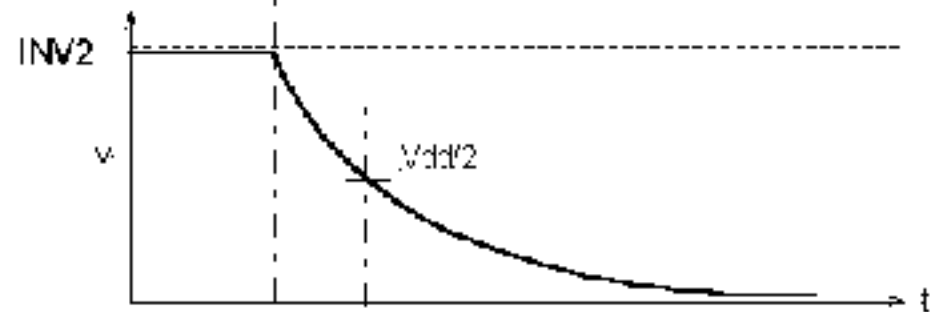
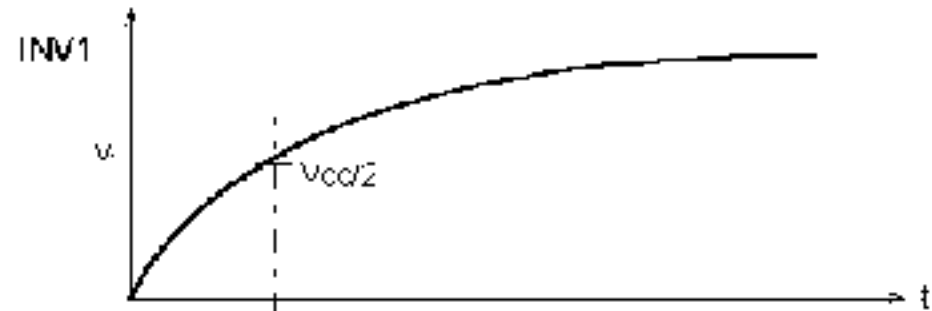
Intrinsic delay of inverter
independent of size

Turning Rise/Fall Delay into Gate Delay

- Cascaded gates:



“transfer curve” for inverter.

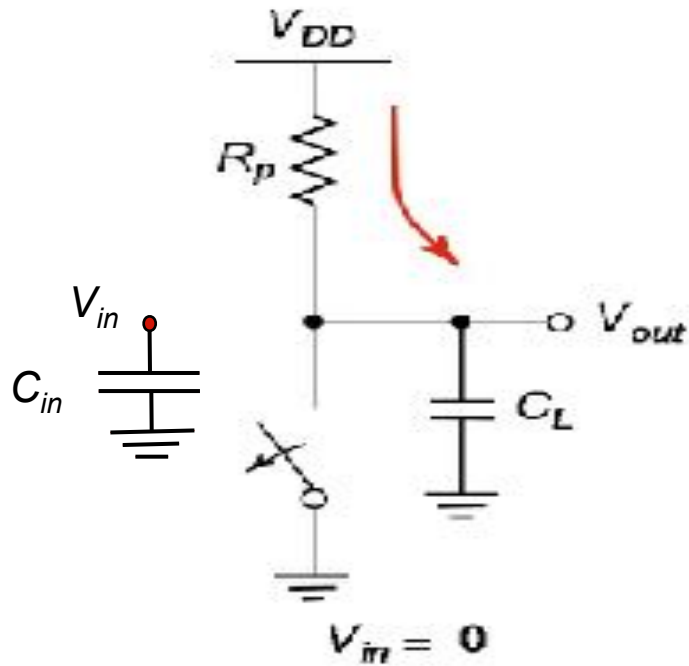


In general
prop. delay = sum of individual prop. delays of gates in series.

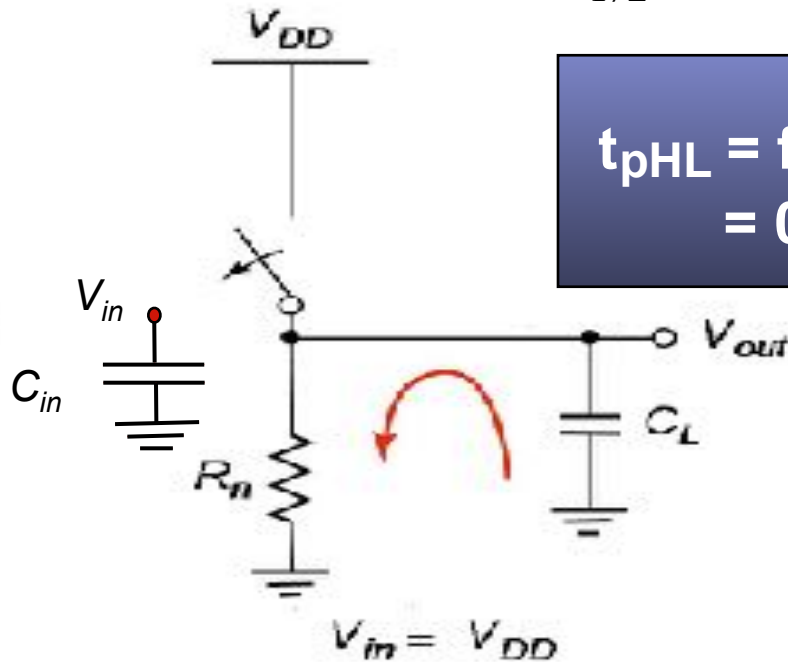
The Switch Inverter: Transient Response

$$V(t) = V_0 e^{-t/RC}$$

$$t_{1/2} = \ln(2) \times RC$$



(a) Low-to-high

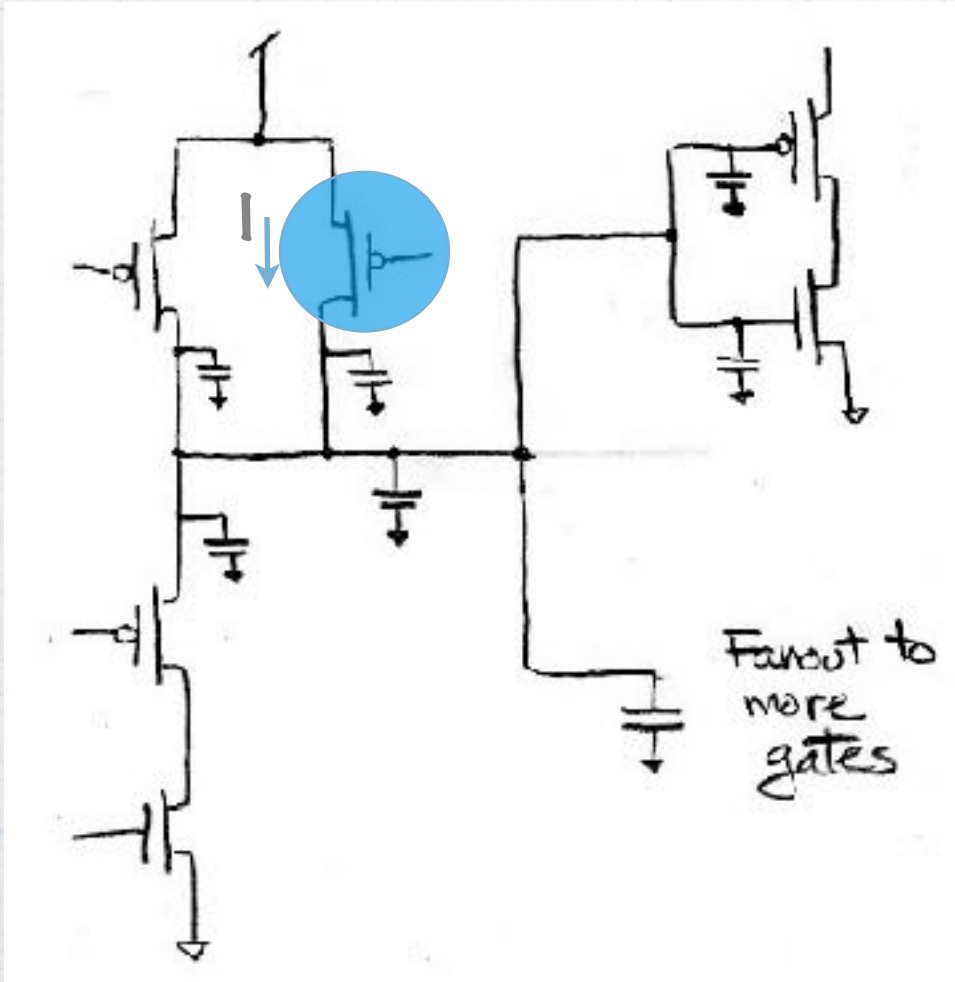


(b) High-to-low

$$t_{pHL} = f(R_{on} C_L) \\ = 0.69 R_n C_L$$

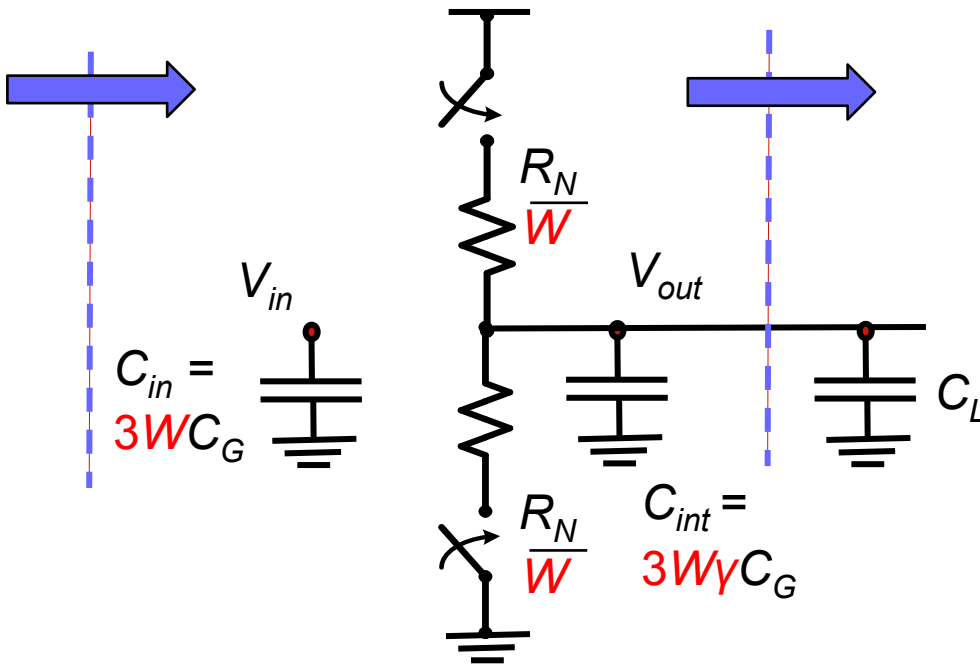
More on C_L

- ▶ Everything that connects to the output of a logic gate (or transistor) contributes capacitance:



- ▶ Transistor drains
- ▶ Interconnection (wires/contacts/vias)
- ▶ Transistor Gates

Inverter with Load Capacitance



$$\begin{aligned}
 t_p &= 0.69 \left(\frac{R_N}{W} \right) (C_{int} + C_L) \\
 &= 0.69 \left(\frac{R_N}{W} \right) (3W\gamma C_G + C_L) \\
 &= 0.69 (3C_G R_N) \left(\gamma + \frac{C_L}{C_{in}} \right) \\
 &= t_{inv} \left(\gamma + \frac{C_L}{C_{in}} \right) = t_0 (\gamma + f)
 \end{aligned}$$

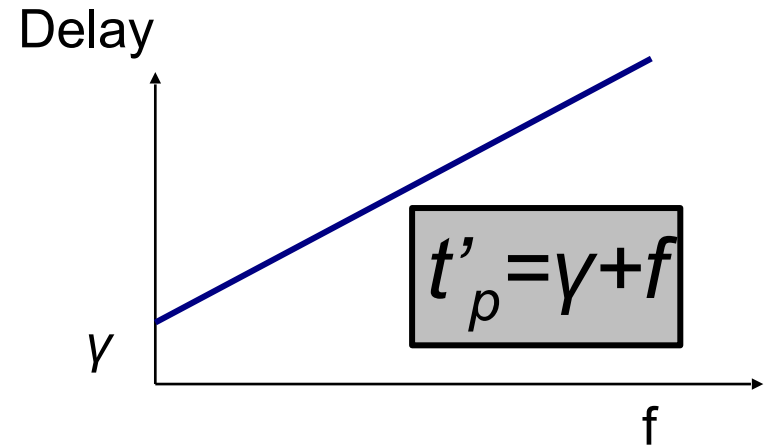
f = fanout = ratio between load and input capacitance of gate

Inverter Delay Model

$$t_p = t_{inv}(\gamma + f)$$

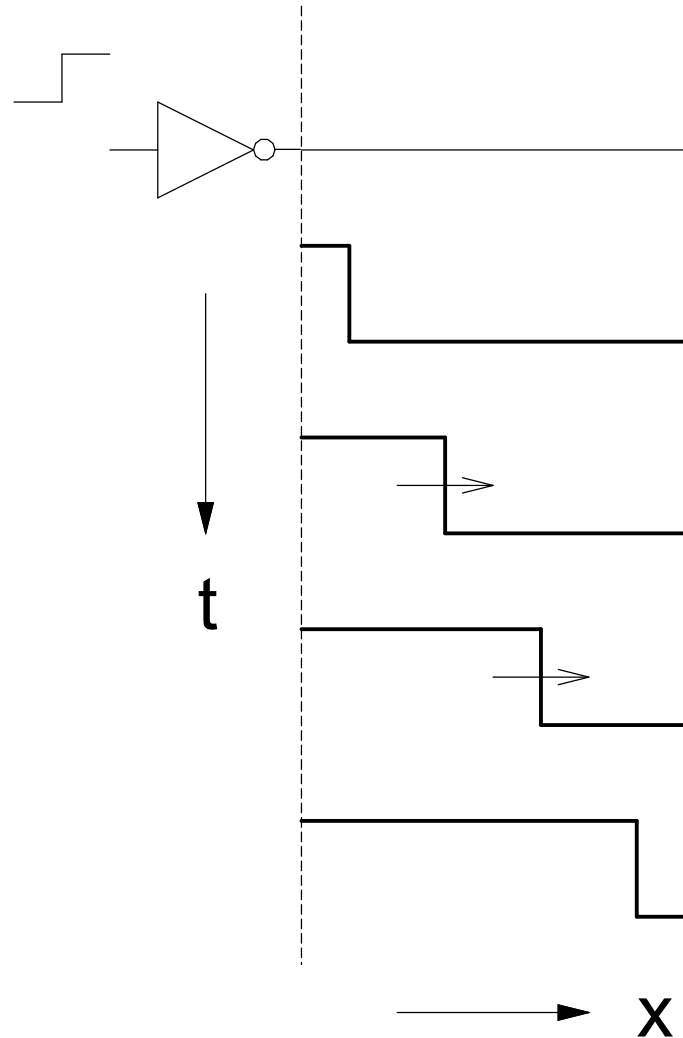
t_{inv} technology constant

- Can be dropped from expression
- Delay unit-less variable (expressed in unit delays)



Question: how does transistor sizing (W) impact delay?

Wire Delay

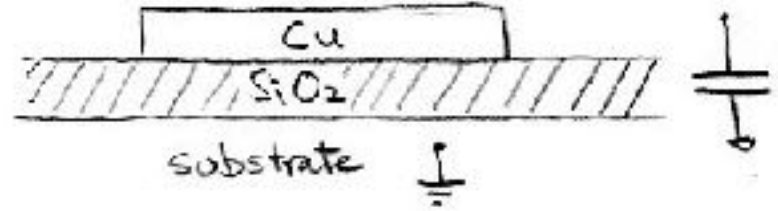


- Ideally, wires behave as “transmission lines”:
 - signal wave-front moves close to the speed of light
 - $\sim 1\text{ft/ns}$
 - Time from source to destination is called the “transit time”.
 - In ICs most wires are short, and the transit times are relatively short compared to the clock period and can be ignored.
 - Not so on PC boards.

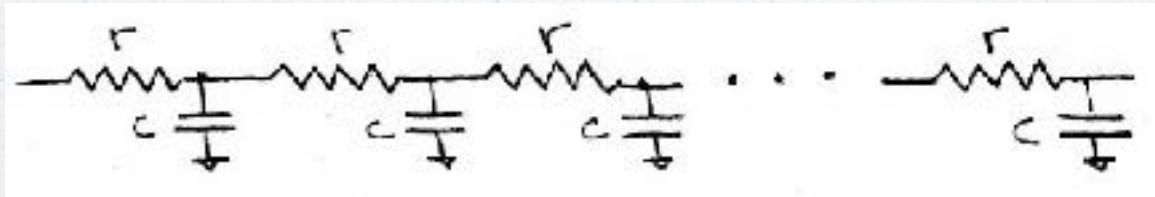
Wires

- ▶ As parallel plate capacitors:

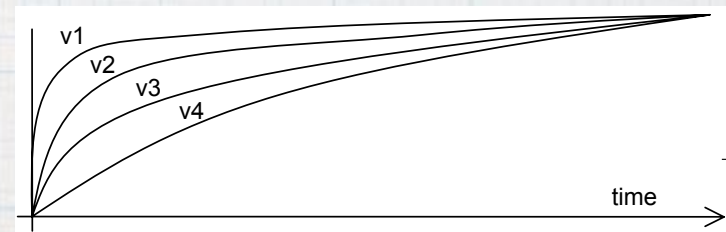
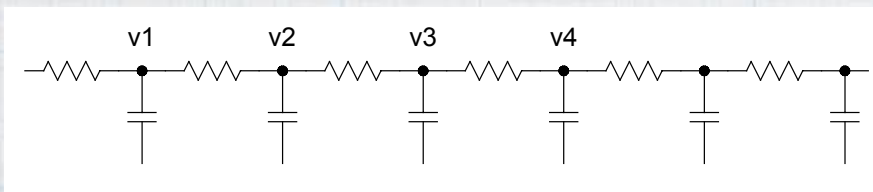
$$C \propto \text{Area} = \text{width} * \text{length}$$



- ▶ Wires have finite resistance, so have distributed R and C:



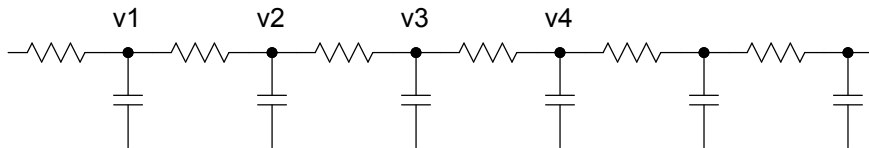
with $r = \text{res}/\text{length}$, $c = \text{cap}/\text{length}$, $\Delta \propto rL^2 \cong rc + 2rc + 3rc + \dots$



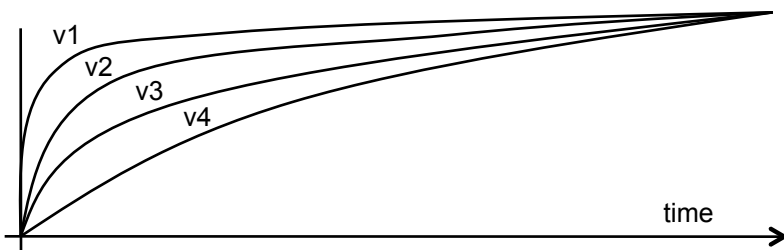
Wire Delay

- Even in those cases where the transmission line effect is negligible:

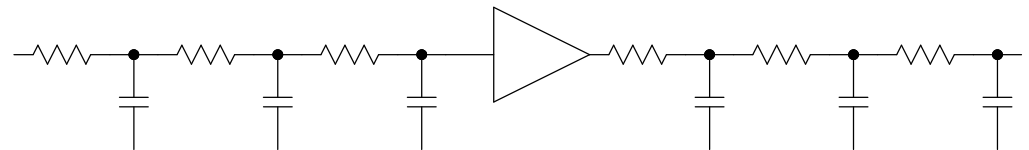
- Wires possess distributed resistance and capacitance



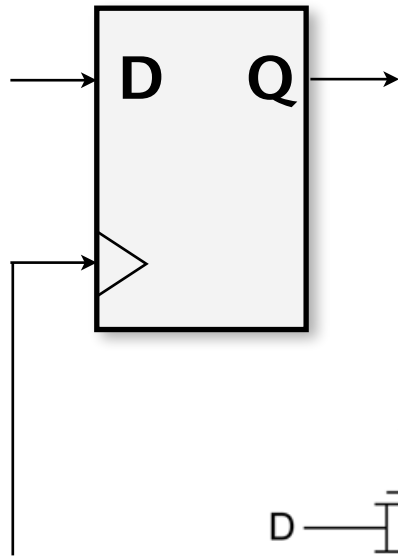
- Time constant associated with distributed RC is proportional to the *square* of the length



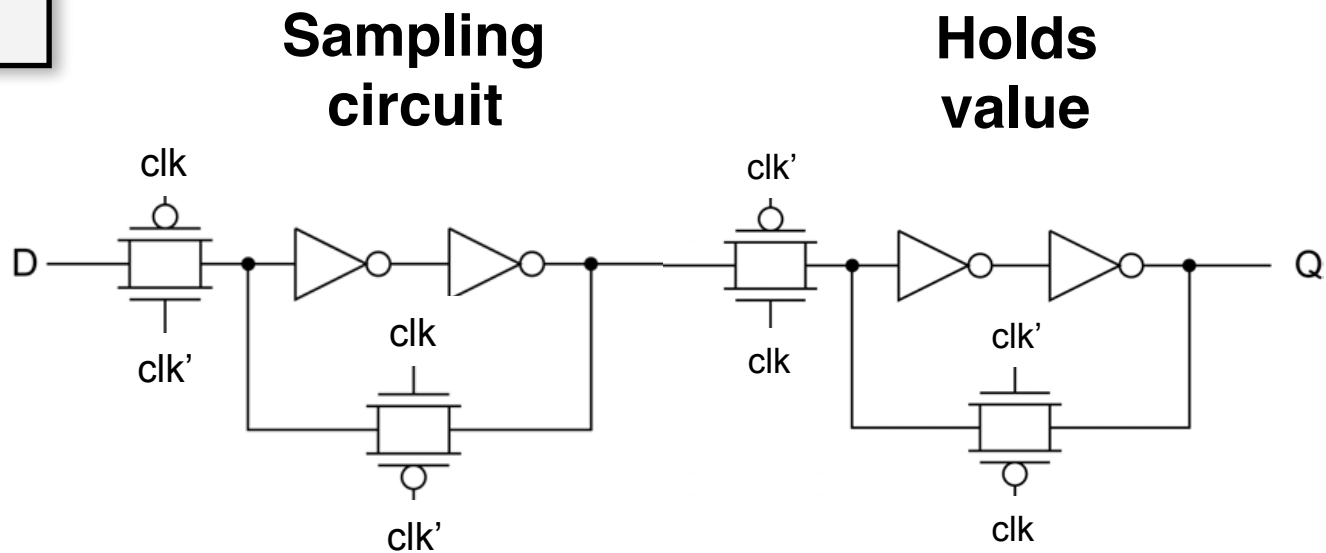
- For **short wires** on ICs, resistance is insignificant (relative to effective R of transistors), but C is important.
 - Typically around half of C of gate load is in the wires.
- For **long wires** on ICs:
 - busses, clock lines, global control signal, etc.
 - Resistance is significant, therefore distributed RC effect dominates.
 - signals are typically “rebuffered” to reduce delay:



Recall: Positive edge-triggered flip-flop

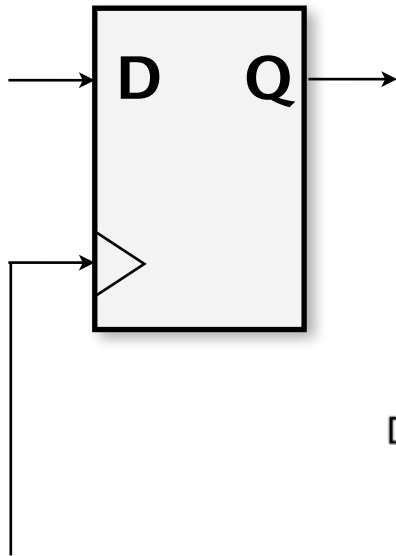


A flip-flop “samples” right before the edge, and then “holds” value.

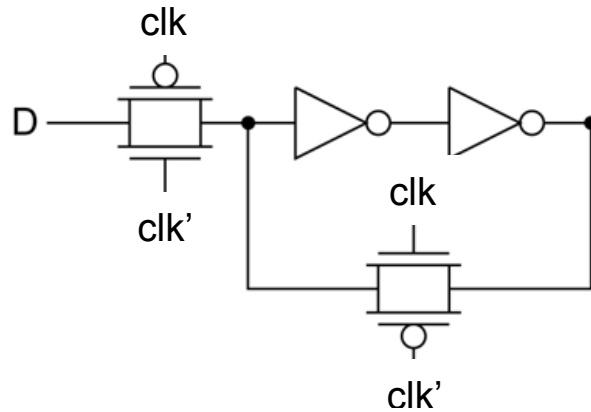


Sensing: When clock is low

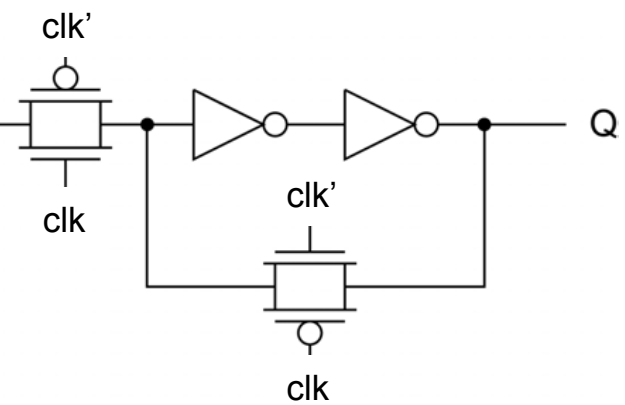
A flip-flop “samples” right before the edge, and then “holds” value.



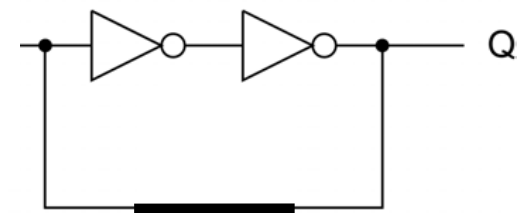
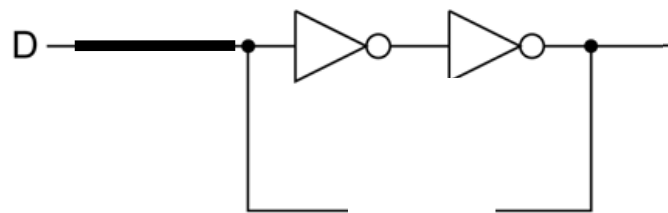
Sampling circuit



Holds value



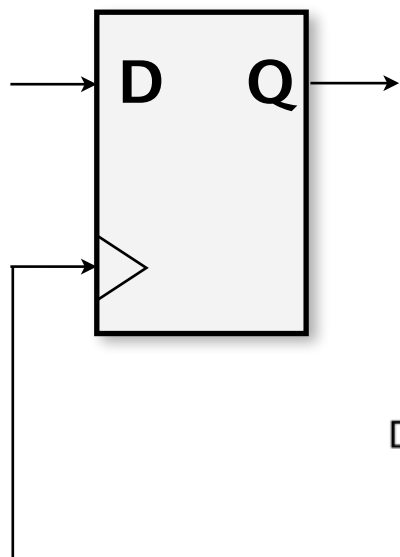
$clk = 0$
 $clk' = 1$



Will capture new value on posedge.

Outputs last value captured.

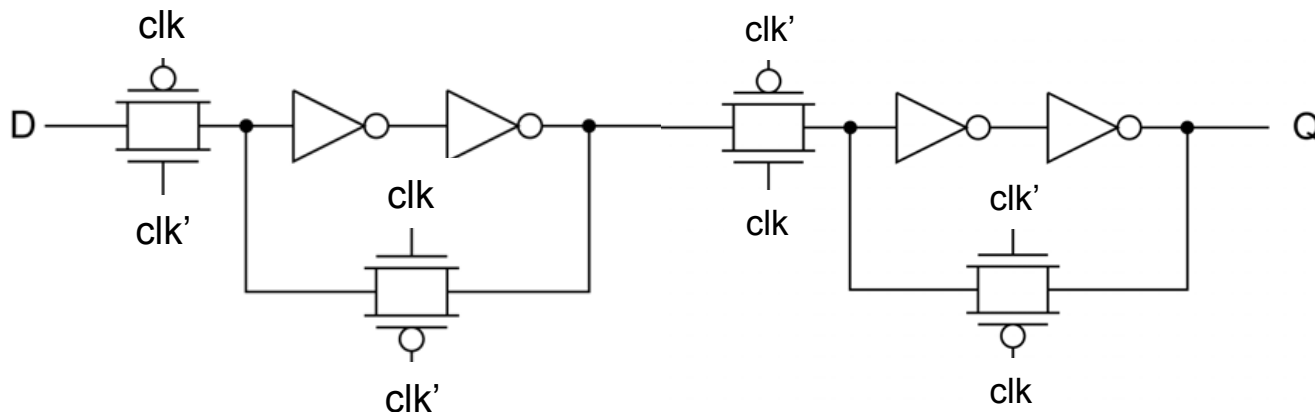
Capture: When clock goes high



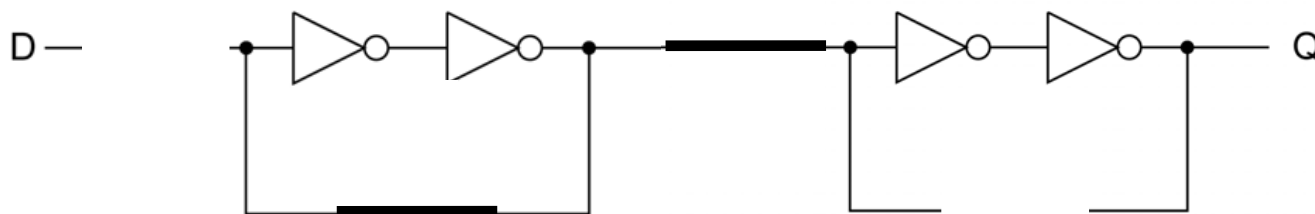
A flip-flop “samples” right before the edge, and then “holds” value.

Sampling circuit

Holds value



$clk = 1$
 $clk' = 0$



Remembers value just captured.

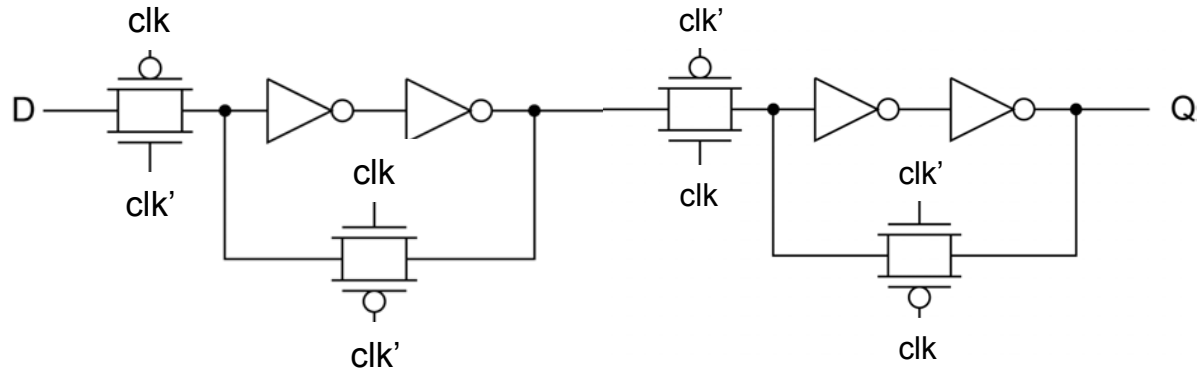
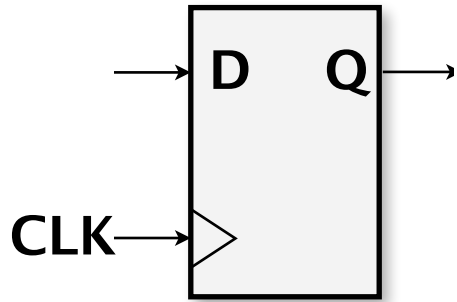
Outputs value just captured.

Flip Flop delays:

clk-to-Q?

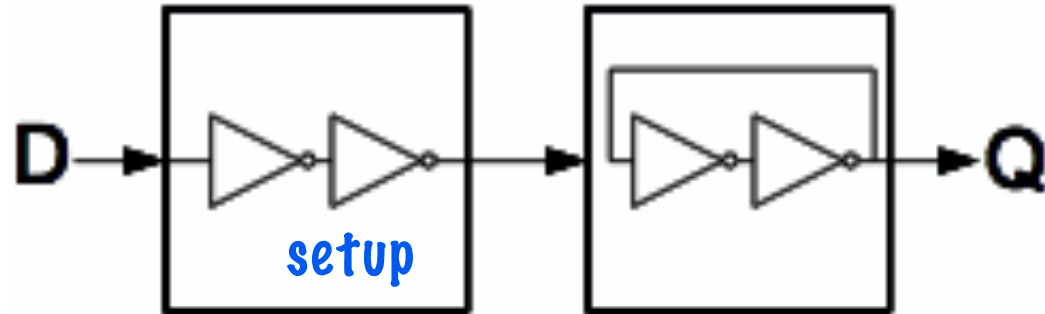
setup?

hold?



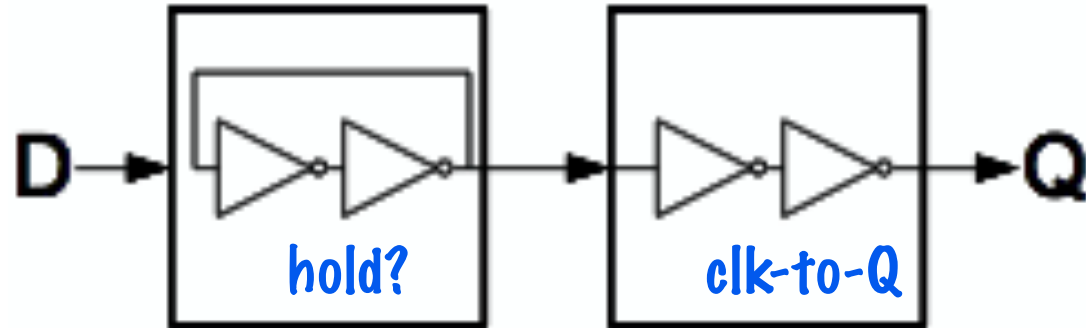
CLK == 0

Sense D, but Q outputs old value.

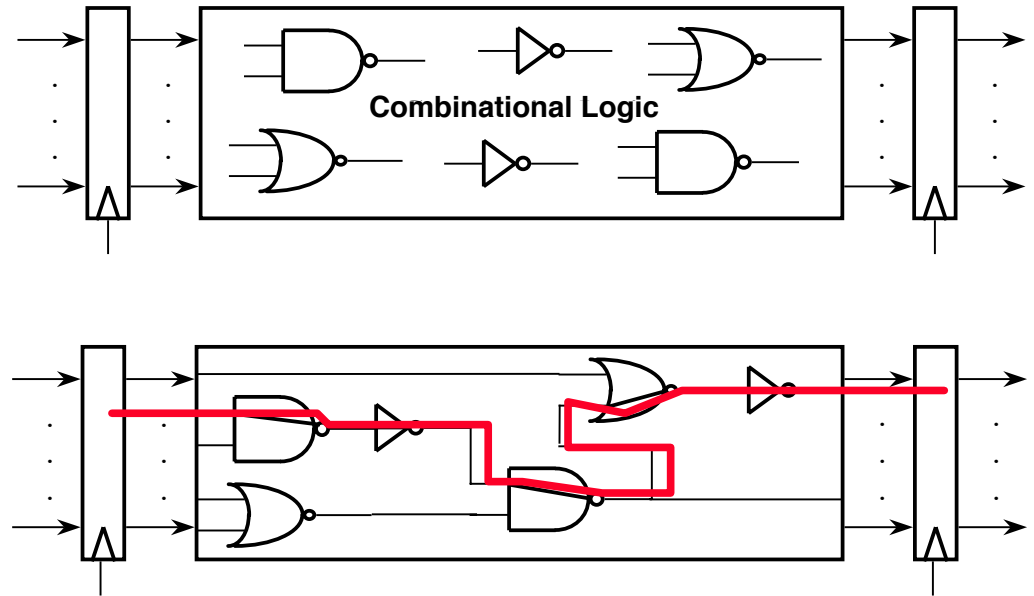
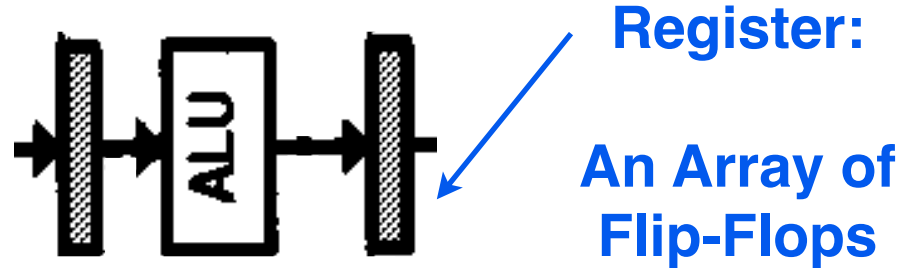
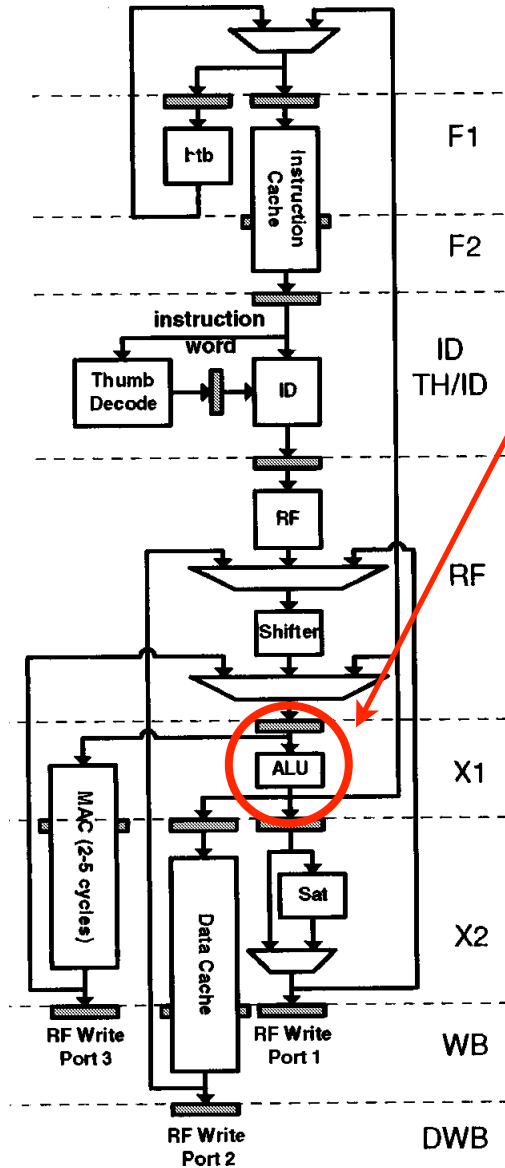


CLK 0->1

Capture D, pass value to Q



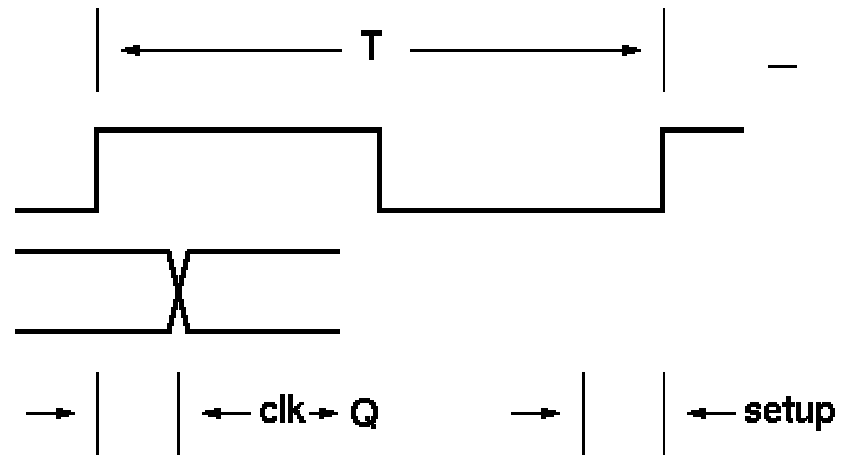
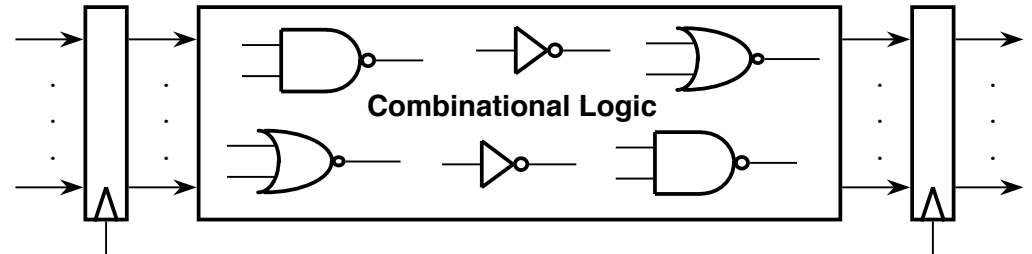
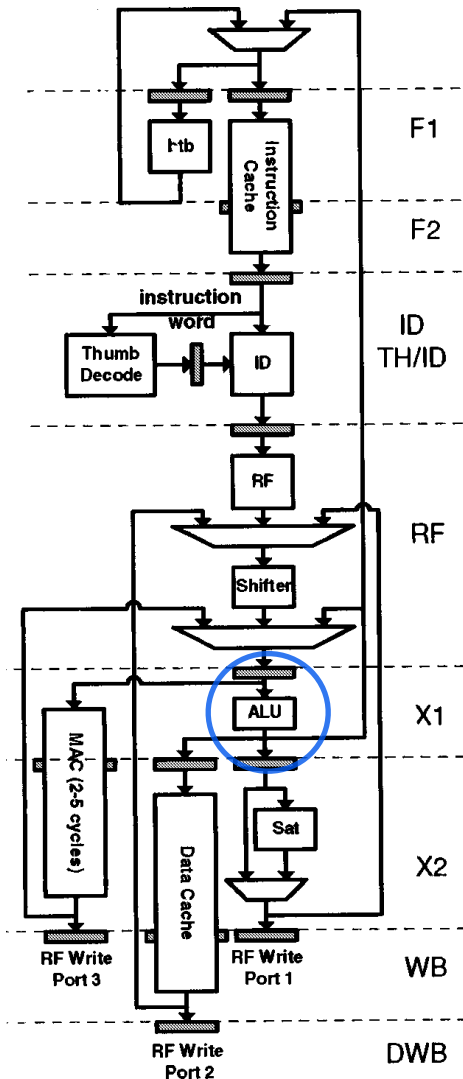
Timing Analysis and Logic Delay



If our clock period $T >$ worst-case delay through CL, does this ensure correct operation?



Flip-Flop delays eat into “time budget”

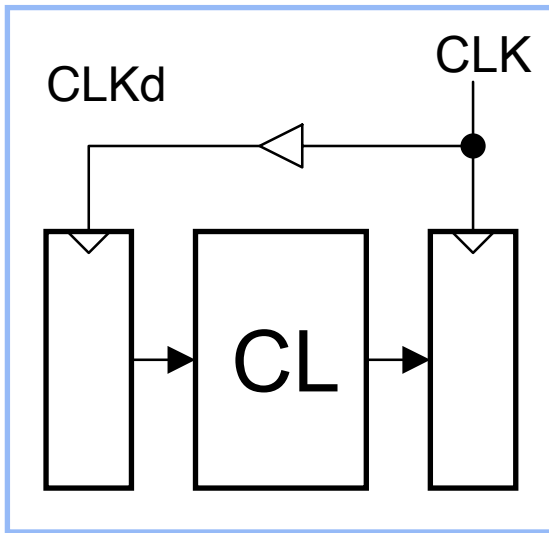


ALU “time budget”

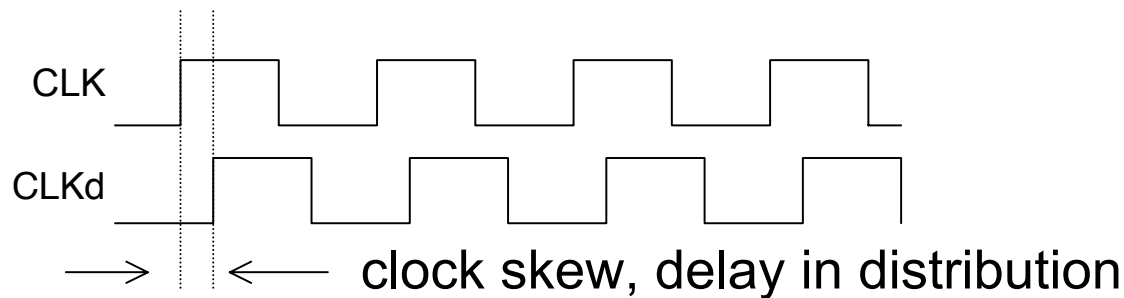
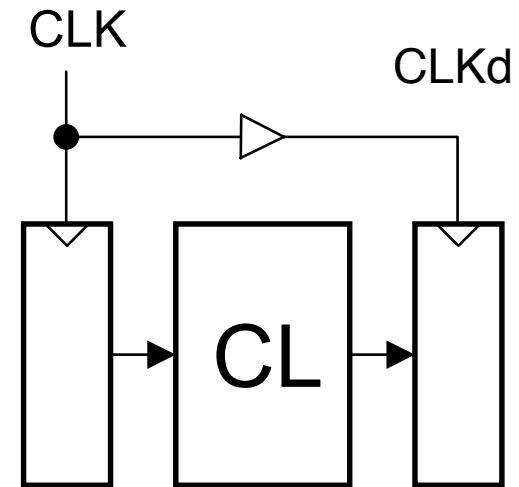
$$T \geq \tau_{\text{clk} \rightarrow Q} + \tau_{\text{CL}} + \tau_{\text{setup}}$$



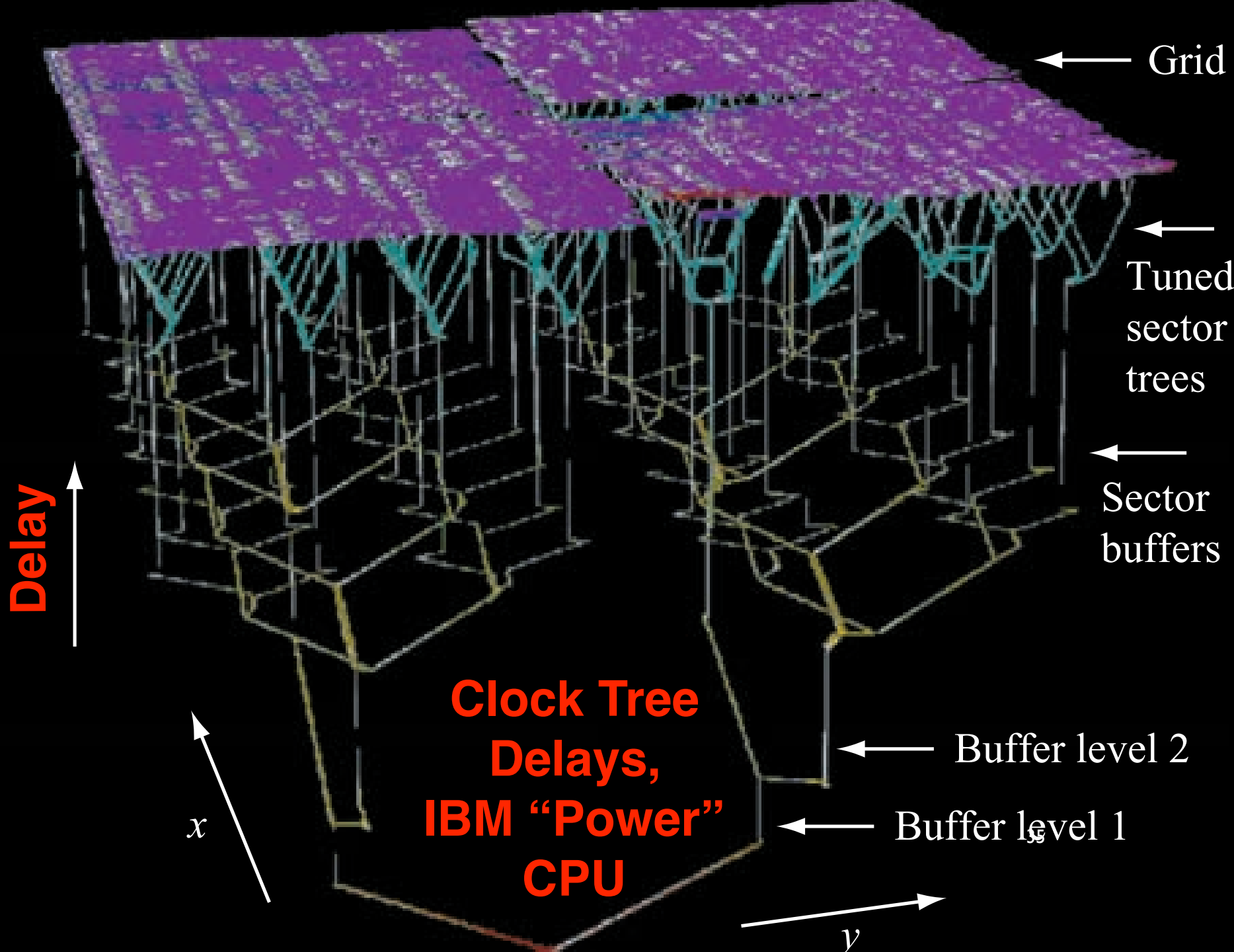
Clock skew also eats into “time budget”

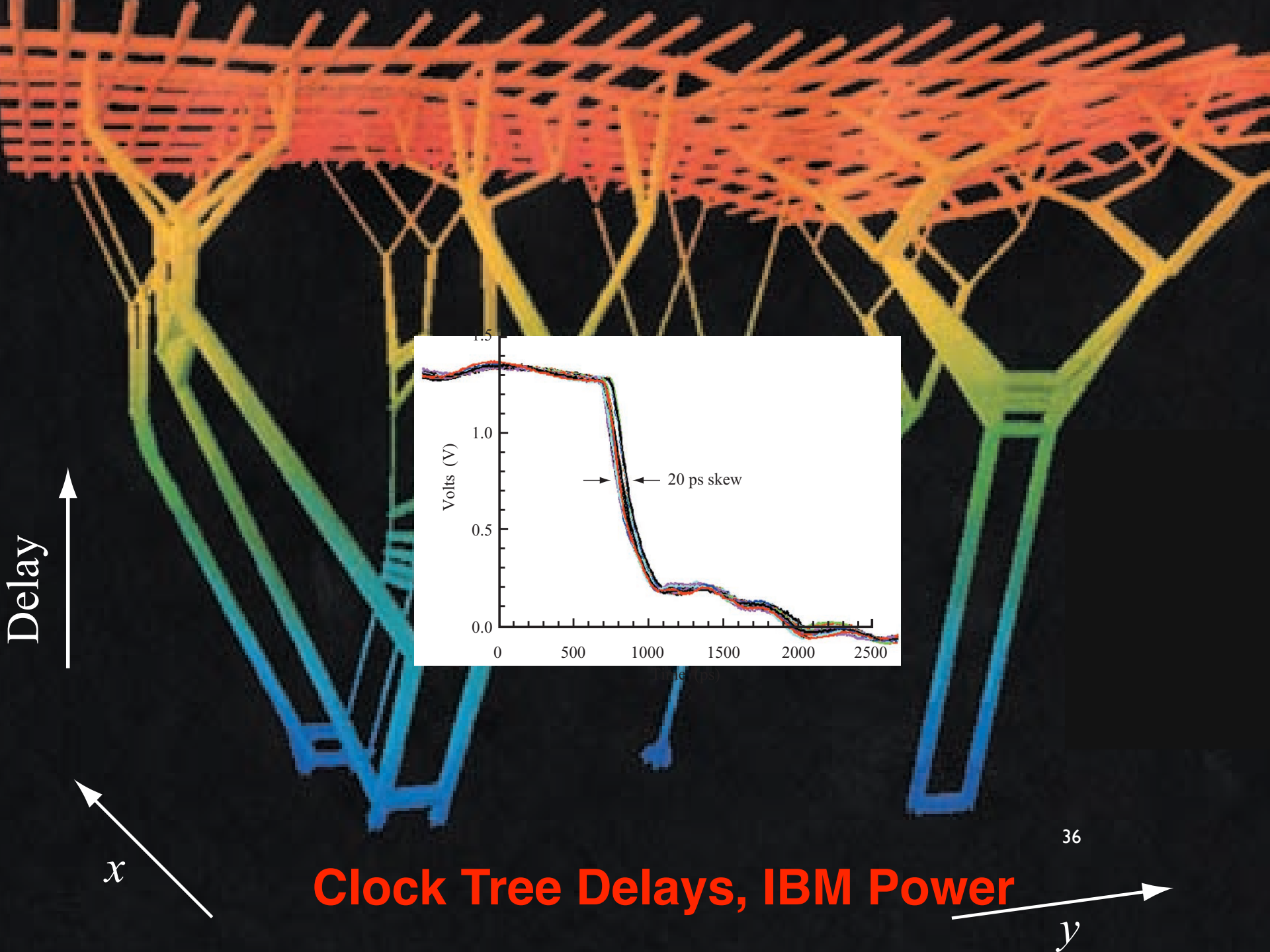


**As $T \rightarrow 0$,
which circuit
fails first?**



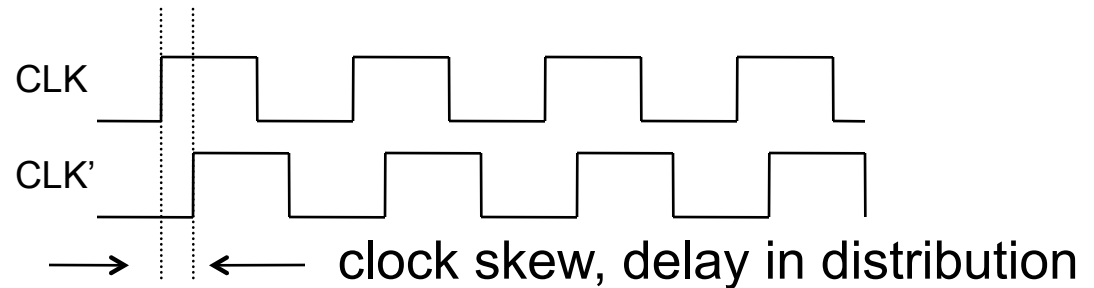
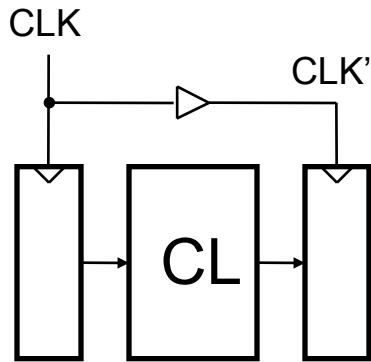
$$T \geq T_{CL} + T_{\text{setup}} + T_{\text{clk} \rightarrow Q} + \text{worst case skew.}$$





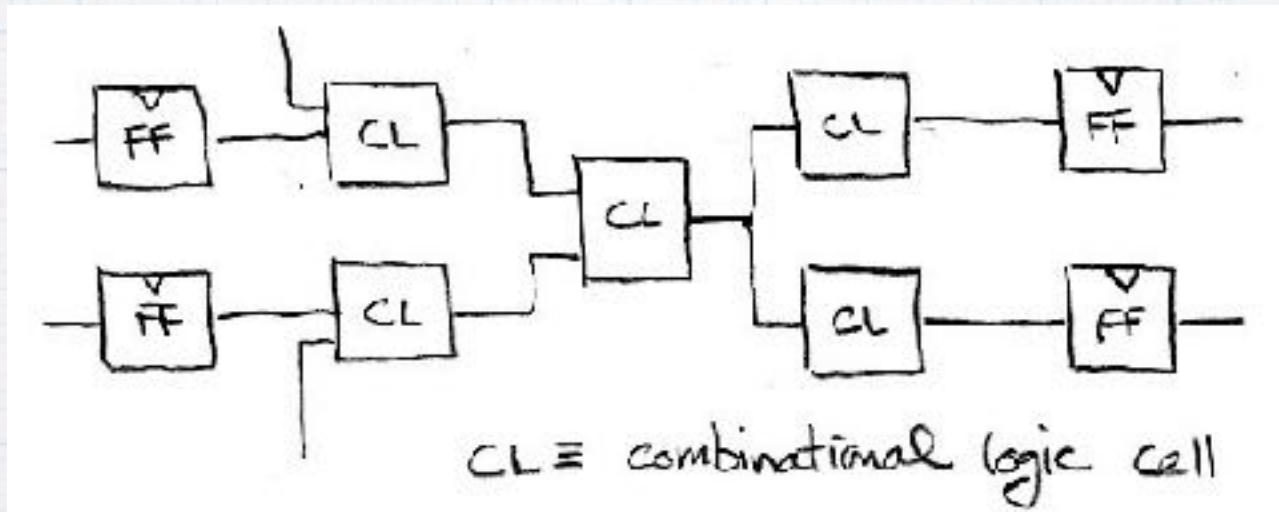
Clock Tree Delays, IBM Power

Clock Skew (cont.)



- Note reversed buffer.
- In this case, clock skew actually provides *extra time* (adds to the effective clock period).
- This effect has been used to help run circuits at higher clock rates. Risky business!

Components of Path Delay

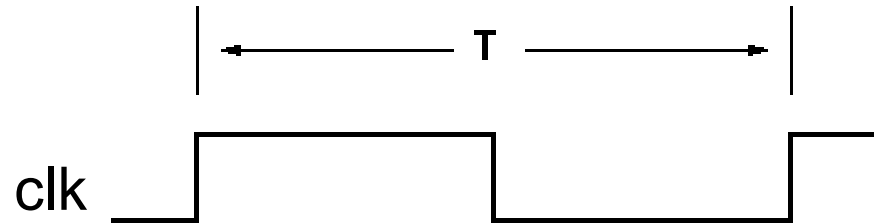
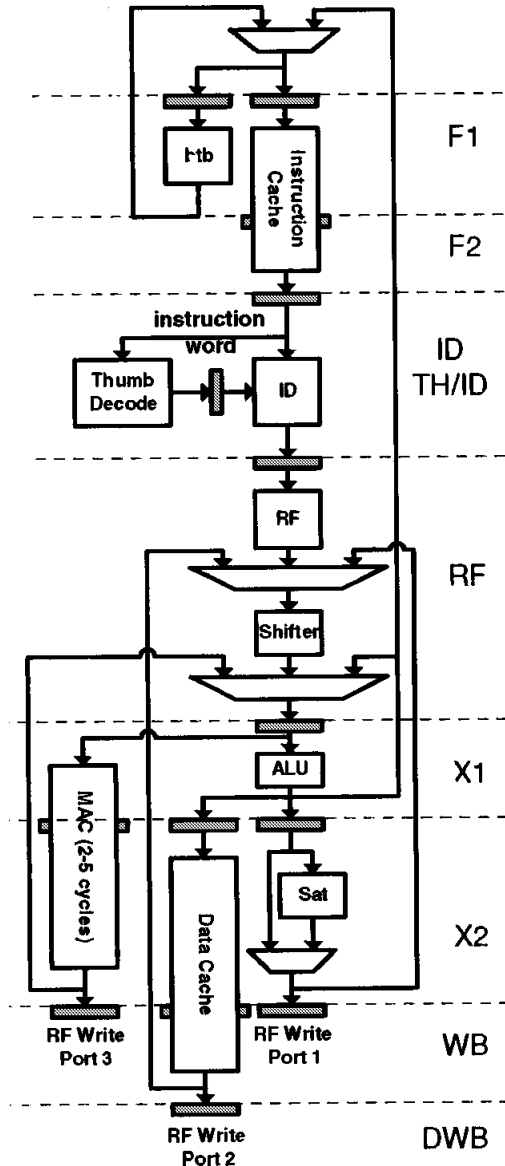


1. # of levels of logic
2. Internal cell delay
3. wire delay
4. cell input capacitance
5. cell fanout
6. cell output drive strength

Who controls the delay?

| | foundary engineer (TSMC) | Library Developer (Aritsan) | CAD Tools (DC, IC Compiler) | Designer (you!) |
|----------------------------------|--------------------------|-----------------------------|-----------------------------|------------------|
| 1. # of levels | | | synthesis | RTL |
| 2. Internal cell delay | physical parameters | cell topology, trans sizing | cell selection | |
| 3. Wire delay | physical parameters | | place & route | layout generator |
| 4. Cell input capacitance | physical parameters | cell topology, trans sizing | cell selection | instantiation |
| 5. Cell fanout | | | synthesis | RTL |
| 6. Cell drive strength | physical parameters | transistor sizing | cell selection | instantiation |

From Delay Models to Timing Analysis



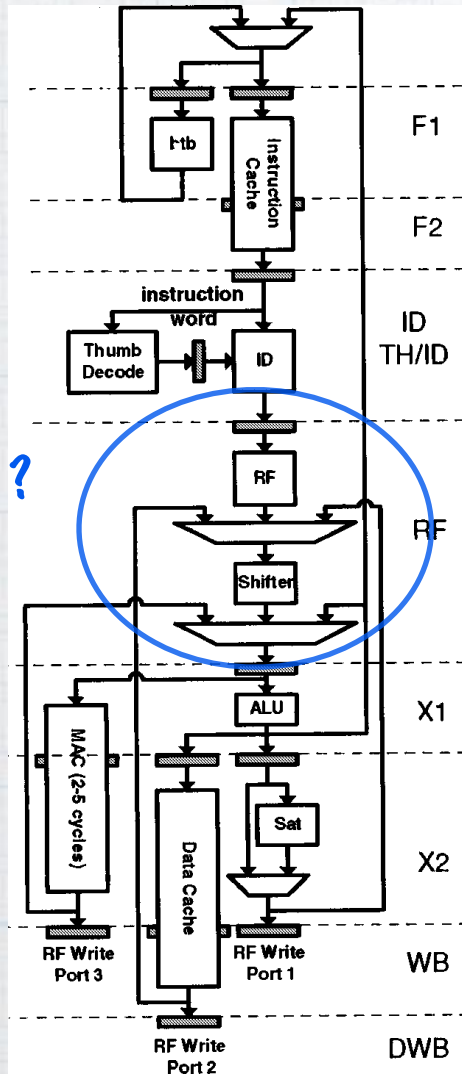
Timing Analysis

What is the smallest T that produces correct operation?
Or, can we meet a target T ?

| f | T |
|---------|-----------|
| 1 MHz | 1 μ s |
| 10 MHz | 100 ns |
| 100 MHz | 10 ns |
| 1 GHz | 1 ns |



Timing Closure: Searching for and beating down the critical path



Must consider all connected register pairs, paths, plus from input to register, plus register to output.

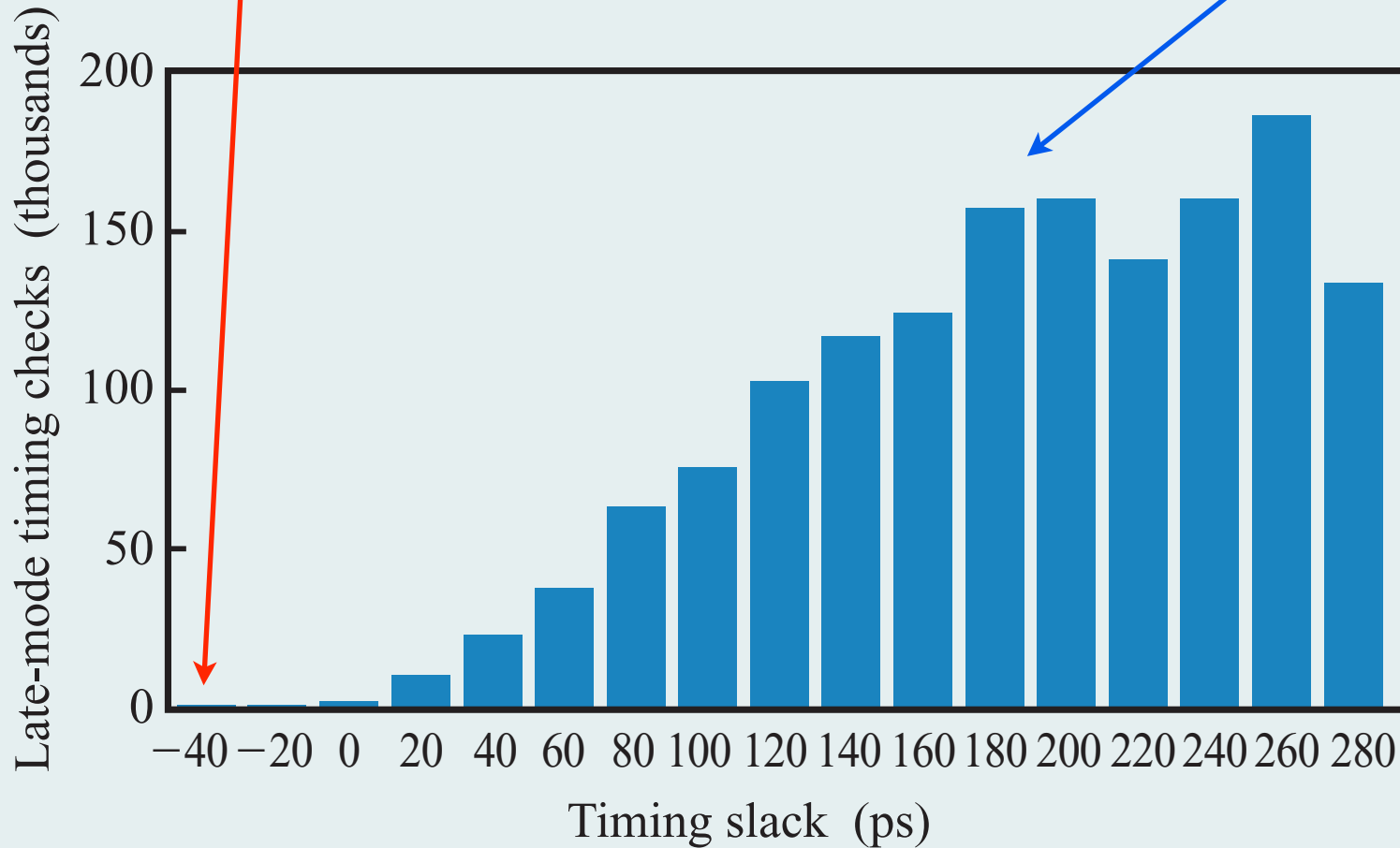
- Design tools help in the search.
- Synthesis tools work to meet clock constraint, report delays on paths,
- Special static timing analyzers accept a design netlist and report path delays,
- and, of course, simulators can be used to determine timing performance.

Tools that are expected to do something about the timing behavior (such as synthesizers), also include provisions for specifying input arrival times (relative to the clock), and output requirements (set-up times of next stage).

Timing Analysis, real example

The critical path

Most paths have hundreds of picoseconds to spare.



From "The circuit and physical design of the POWER4 microprocessor", IBM J Res and Dev, 46:1, Jan 2002, J.D. Warnock et al.

Timing Optimization

As an ASIC/FPGA designer you get to choose:

- ▶ The algorithm
- ▶ The Microarchitecture (block diagram)
- ▶ The RTL description of the CL blocks (number of levels of logic)
- ▶ Where to place registers and memory (the pipelining)
- ▶ Overall floorplan and relative placement of blocks

How to retime logic

Critical path is 5.
We want to improve it without changing circuit semantics.

Add a register, move one circle.
Performance improves by 20%.

Circles are combinational logic, labelled with delays.

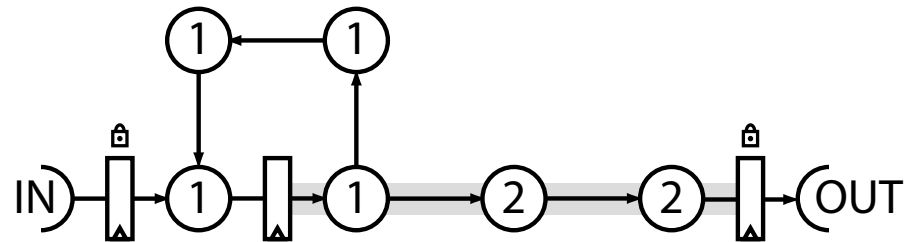


Figure 1: A small graph before retiming. The nodes represent logic delays, with the inputs and outputs passing through mandatory, fixed registers. The critical path is 5.

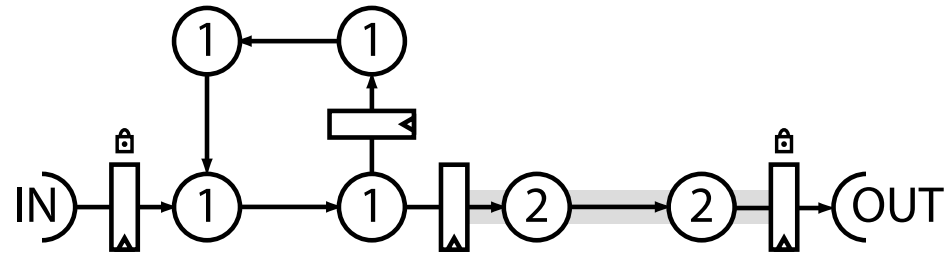
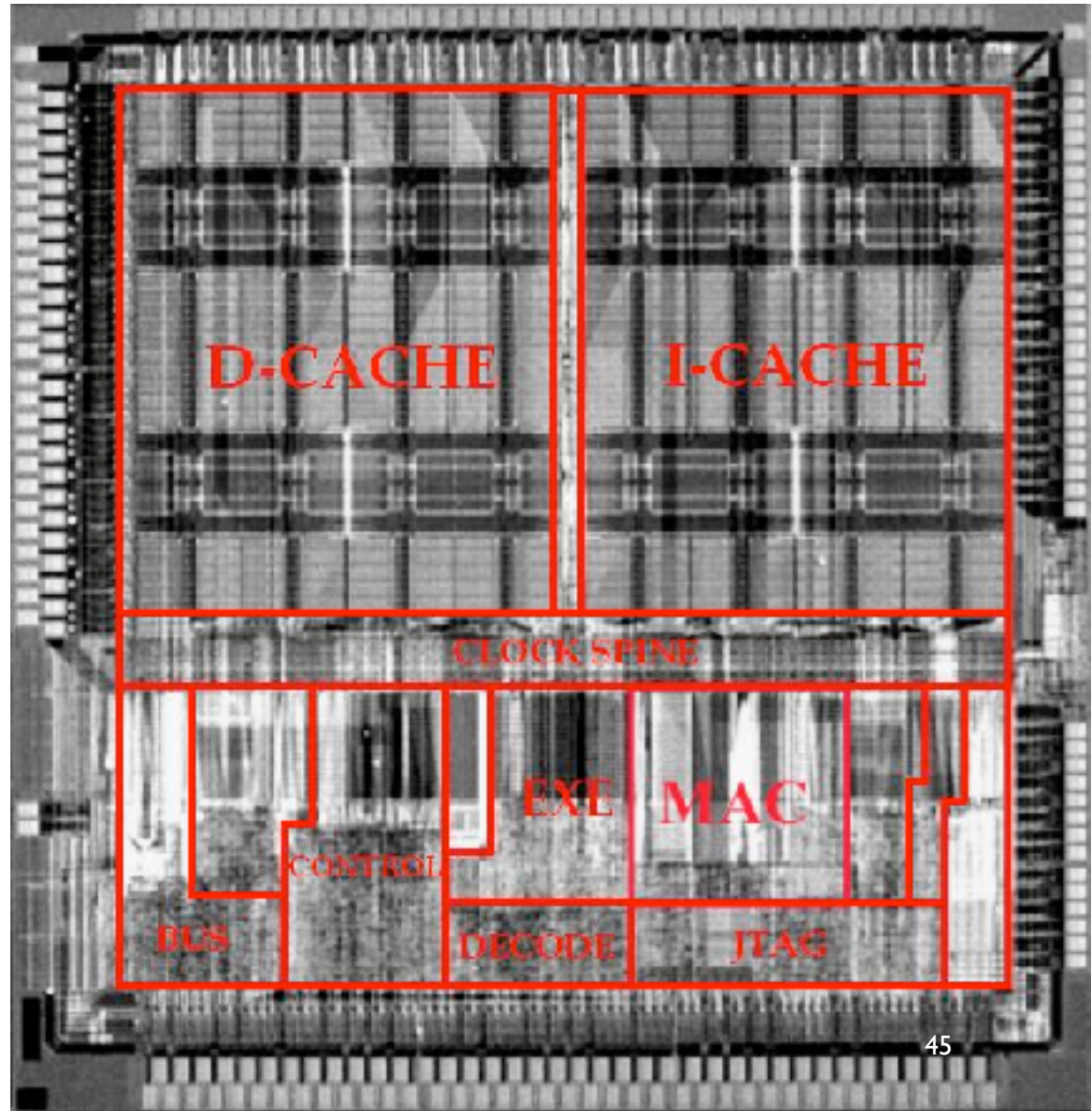
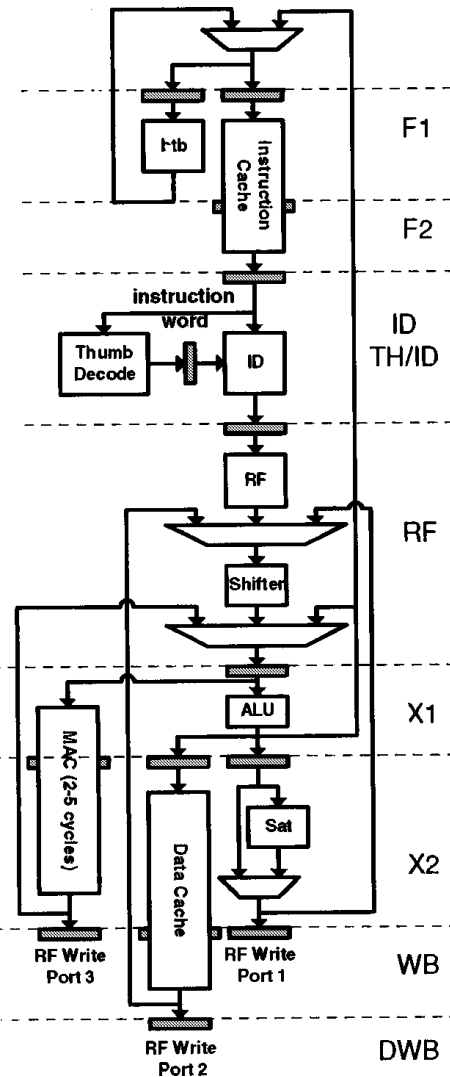


Figure 2: The example in Figure 2 after retiming. The critical path is reduced from 5 to 4.

Logic Synthesis tools can do this in simple cases.

Floorplaning: essential to meet timing.



(Intel XScale 80200)

Timing Analysis Tools

- ▶ **Static Timing Analysis:** Tools use delay models for gates and interconnect. Traces through circuit paths.

- ▶ **Cell delay model capture**

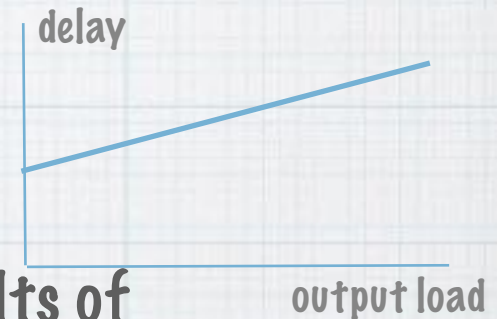
- ▶ For each input/output pair, internal delay (output load independent)
- ▶ output dependent delay

- ▶ **Standalone tools (PrimeTime) and part of logic synthesis.**

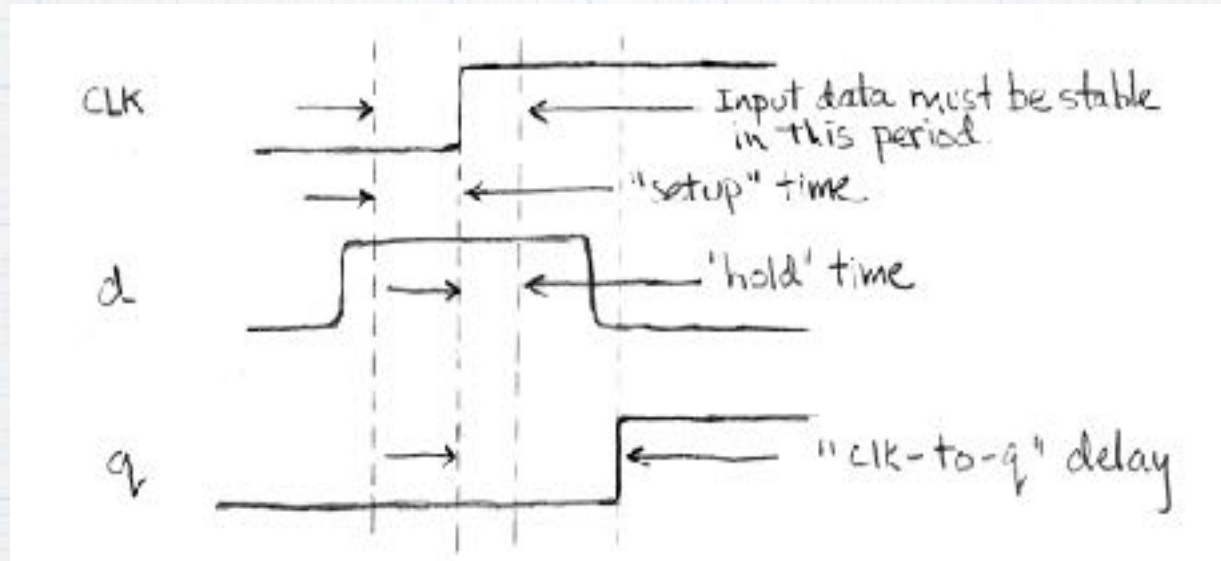
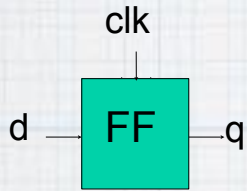
- ▶ **Back-annotation takes information from results of place and route to improve accuracy of timing analysis.**

- ▶ **DC in “topographical mode” uses preliminary layout information to model interconnect parasitics.**

- ▶ Prior versions used a simple fan-out model of gate loading.



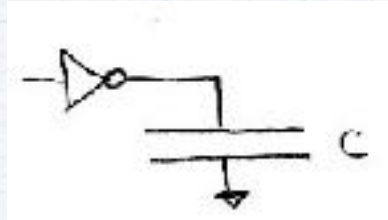
Hold-time Violations



- ▶ Some state elements have positive hold time requirements.
 - ▶ How can this be?
- ▶ Fast paths from one state element to the next can create a violation. (Think about shift registers!)
- ▶ CAD tools do their best to fix violations by inserting delay (buffers).
 - ▶ Of course, if the path is delayed too much, then cycle time suffers.
 - ▶ Difficult because buffer insertion changes layout, which changes path delay.

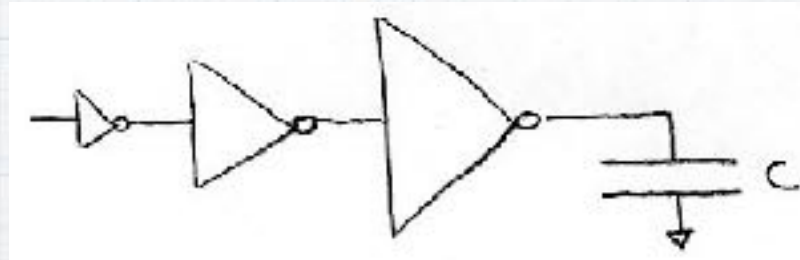
Driving Large Loads

- ▶ Large fanout nets: clocks, resets, memory bit lines, off-chip
- ▶ Relatively small driver results in long rise time (and thus large gate delay)



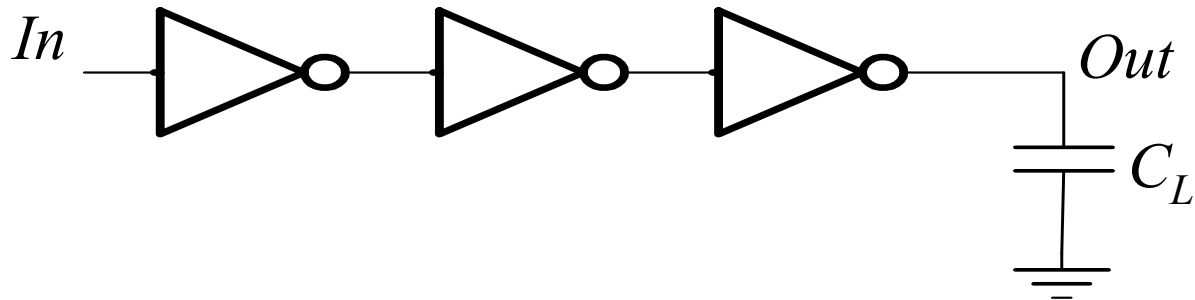
- ▶ Strategy:

Staged Buffers



- ▶ How to optimally scale drivers?
- ▶ Optimal trade-off between delay per stage and total number of stages?

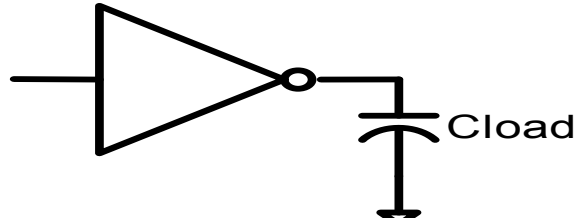
Inverter Chain



- For some given C_L :
 - How many stages are needed to minimize delay?
 - How to size the inverters?
- Anyone want to guess the solution?

Careful about Optimization Problems

- Get fastest delay if build one **very** big inverter
 - So big that delay is set only by self-loading

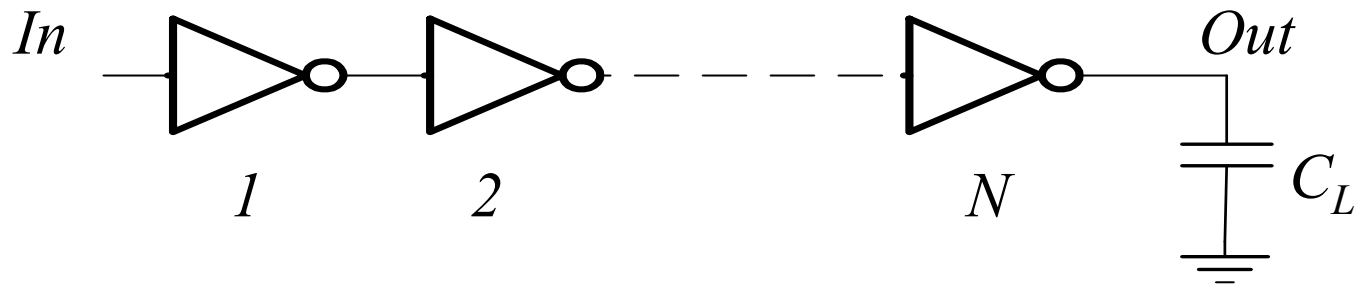


- Likely not the problem you're interested in
 - Someone has to drive this inverter...

Delay Optimization Problem #1

- You are given:
 - A fixed number of inverters
 - The size of the first inverter
 - The size of the load that needs to be driven
- Your goal:
 - Minimize the delay of the inverter chain
- Need model for inverter delay vs. size

Apply to Inverter Chain



$$t_p = t_{p1} + t_{p2} + \dots + t_{pN}$$

$$t_{pj} = t_{inv} \left(\gamma + \frac{C_{in,j+1}}{C_{in,j}} \right)$$

$$t_p = \sum_{j=1}^N t_{p,j} = t_{inv} \sum_{i=1}^N \left(\gamma + \frac{C_{in,j+1}}{C_{in,j}} \right), \quad C_{in,N+1} = C_L$$

Optimal Sizing for Given N

- Delay equation has $N-1$ unknowns, $C_{in,2} \dots C_{in,N}$
- **To minimize the delay**, find $N-1$ partial derivatives:

$$t_p = \dots + t_{inv} \frac{C_{in,j}}{C_{in,j-1}} + t_{inv} \frac{C_{in,j+1}}{C_{in,j}} + \dots$$
$$\frac{dt_p}{dC_{in,j}} = t_{inv} \frac{1}{C_{in,j-1}} - t_{inv} \frac{C_{in,j+1}}{C_{in,j}^2} = 0$$

Optimal Sizing for Given N (cont'd)

- Result: every stage has equal fanout (f):

$$\frac{C_{in,j}}{C_{in,j-1}} = \frac{C_{in,j+1}}{C_{in,j}}$$

- Size of each stage is geometric mean of two neighbors:

$$C_{in,j} = \sqrt{C_{in,j-1} C_{in,j+1}}$$

- Equal fanout \rightarrow every stage will have same delay

Optimum Delay and Number of Stages

- When each stage has same fanout f :

$$f^N = F = C_L / C_{in,1}$$

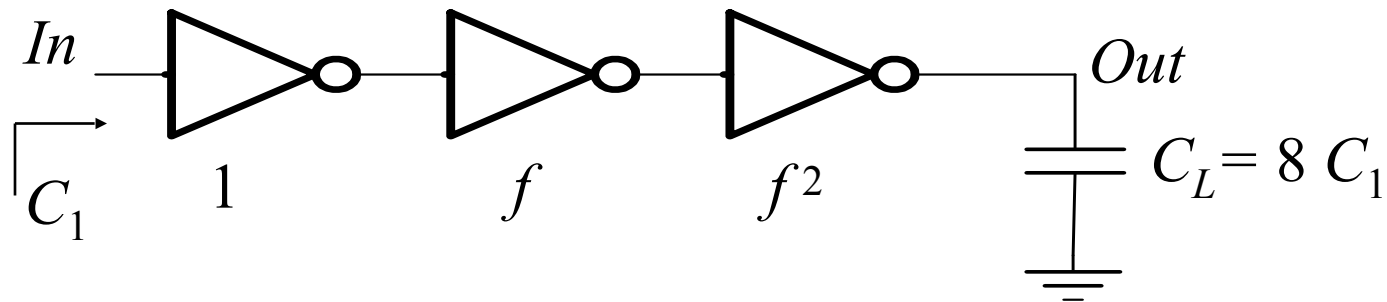
- Fanout of each stage:

$$f = \sqrt[N]{F}$$

- Minimum path delay:

$$t_p = Nt_{inv} \left(\gamma + \sqrt[N]{F} \right)$$

Example



C_L/C_1 has to be evenly distributed across $N = 3$ stages:

Delay Optimization Problem #2

- You are given:
 - The size of the first inverter
 - The size of the load that needs to be driven
- Your goal:
 - Minimize delay by finding optimal number and sizes of gates
- So, need to find N that minimizes:

$$t_p = Nt_{inv} \left(\gamma + \sqrt[N]{C_L / C_{in}} \right)$$

Untangling the Optimization Problem

- Rewrite N in terms of fanout/stage f :

$$f^N = C_L / C_{in} \rightarrow N = \frac{\ln(C_L / C_{in})}{\ln f}$$

$$t_p = N t_{inv} \left((C_L / C_{in})^{1/N} + \gamma \right) = t_{inv} \ln(C_L / C_{in}) \left(\frac{f + \gamma}{\ln f} \right)$$

$$\frac{\partial t_p}{\partial f} = t_{inv} \ln(C_L / C_{in}) \cdot \frac{\ln f - 1 - \gamma / f}{\ln^2 f} = 0$$

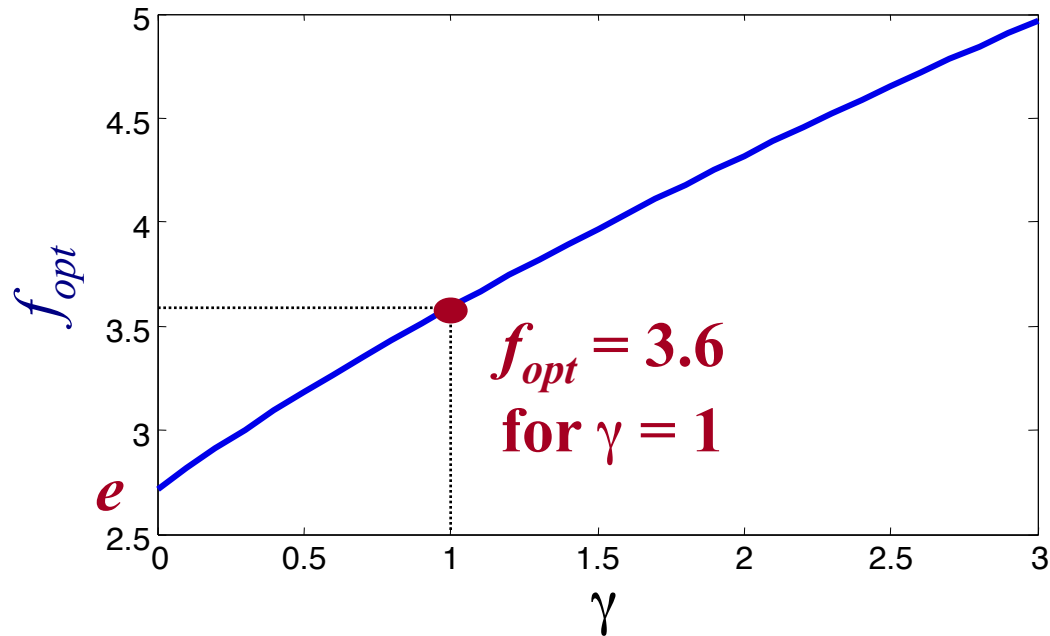
$$f = \exp(1 + \gamma / f) \quad (\text{no explicit solution})$$

For $\gamma = 0$, $f = e$, $N = \ln(C_L / C_{in})$

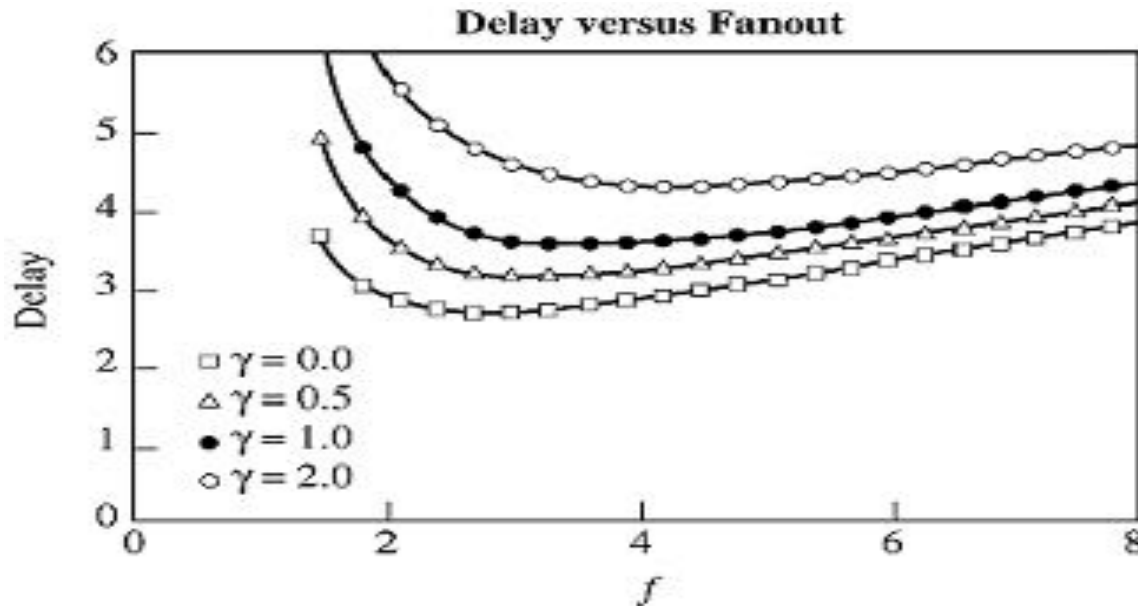
Optimum Effective Fanout f

- Optimum f for given process defined by γ

$$f = \exp(1 + \gamma / f)$$



In Practice: Plot of Total Delay



[Hodges, p.281]

- Why the shape?
- Curves very flat for $f > 2$
 - Simplest/most common choice: $f = 4$

Normalized Delay As a Function of F

$$t_p = Nt_{inv} \left(\gamma + \sqrt[N]{F} \right), F = C_L / C_{in}$$

| <i>F</i> | Unbuffered | Two Stage | Inverter Chain |
|----------|------------|-----------|----------------|
| 10 | 11 | 8.3 | 8.3 |
| 100 | 101 | 22 | 16.5 |
| 1000 | 1001 | 65 | 24.8 |
| 10,000 | 10,001 | 202 | 33.1 |

($\gamma = 1$)

[Rabaey: page 210]

Conclusion

- ▶ **Timing Optimization:** You start with a target on clock period. What control do you have?
- ▶ **Biggest effect is RTL manipulation.**
 - ▶ i.e., how much logic to put in each pipeline stage.
 - ▶ We will be talking later about how to manipulate RTL for better timing results.
- ▶ **In most cases, the tools will do a good job at logic/circuit level:**
 - ▶ Logic level manipulation
 - ▶ Transistor sizing
 - ▶ Buffer insertion
 - ▶ But some cases may be difficult and you may need to help
- ▶ **The tools will need some help at the floorplan and layout**