

EECS 151/251A

Spring 2024

Digital Design and Integrated Circuits

Instructor:

John Wawrzynek

Lecture 11: Timing Part 1

Announcements

- ❑ Midterm Date/Time:
 - ❑ **Tue Mar 12 2024 7:00-10:00PM**
 - ❑ **MOFF101, VLSB2040**
 - ❑ **Covers through Week6/HW6**
 - ❑ **One handwritten “cheat sheet” (both sides)**
- ❑ Monday Mar 11, in-class MT review

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

□ Physics effect:

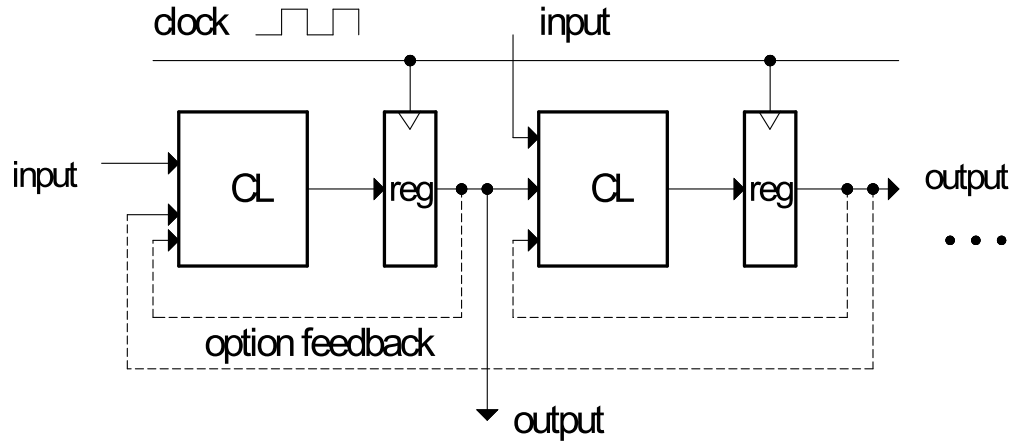
Area \Rightarrow **cost**

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



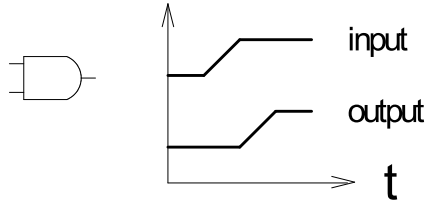
- How do we measure performance?
operations/sec? cycles/sec?
- Performance is directly proportional to clock frequency. Although it may not be the entire story:

Ex: CPU performance

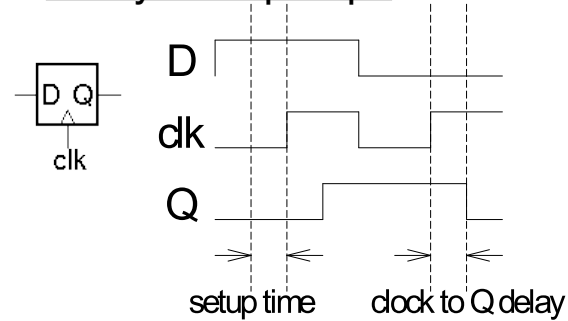
$$= 1 / (\# \text{ instructions} \times \text{CPI} \times \text{clock period})$$

Limitations on Clock Rate

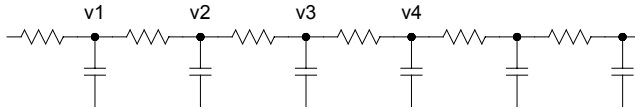
1 Logic Gate Delay



2 Delays in flip-flops



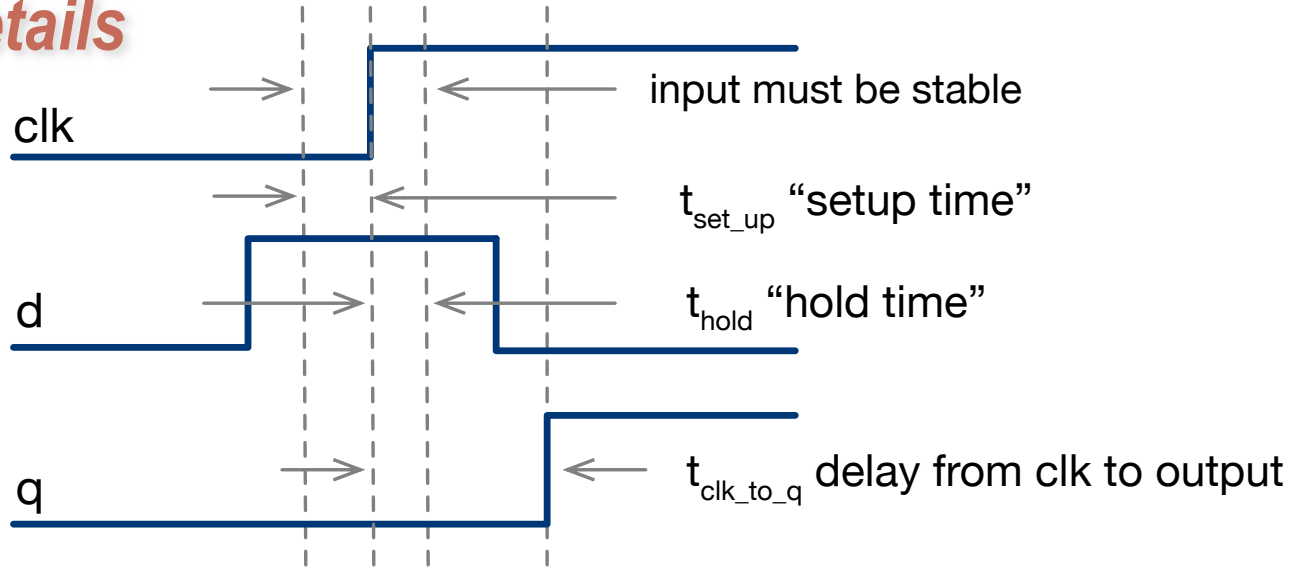
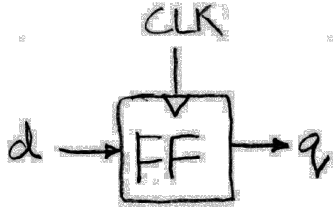
3 Interconnect Delay: wires



1, 2, & 3 all 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).

Register Timing Details



□ Three important times associated with flip-flops:

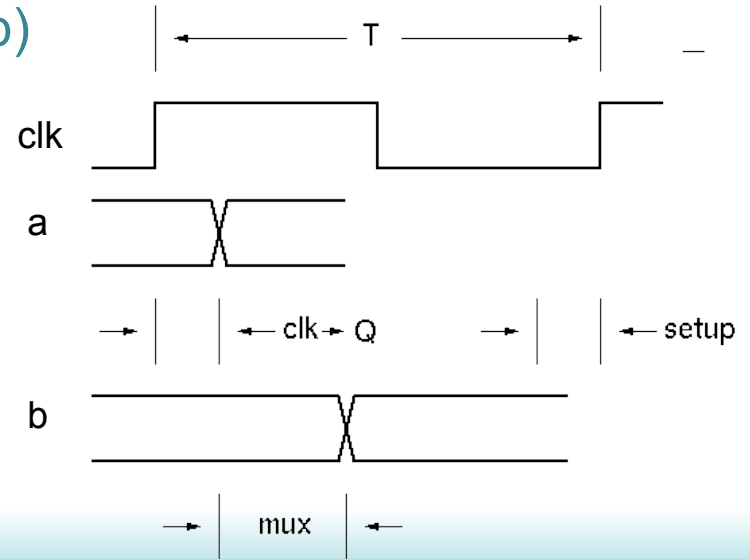
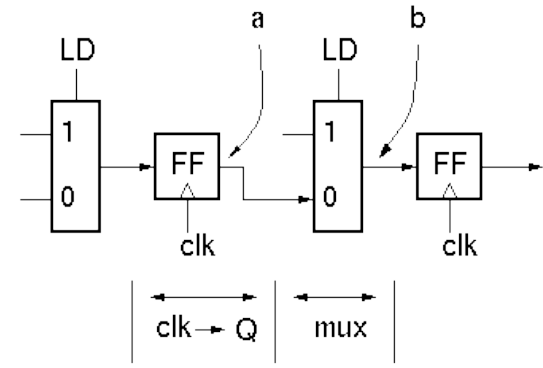
- Setup time - *How long d must be stable before the rising edge of CLK*
- Hold time - *How long D must be stable after the rising edge of CLK*
- Clock-to-q delay – *Propagation delay after rising edge of the CLK*

Example: Timing Analysis

Parallel to serial converter circuit

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

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

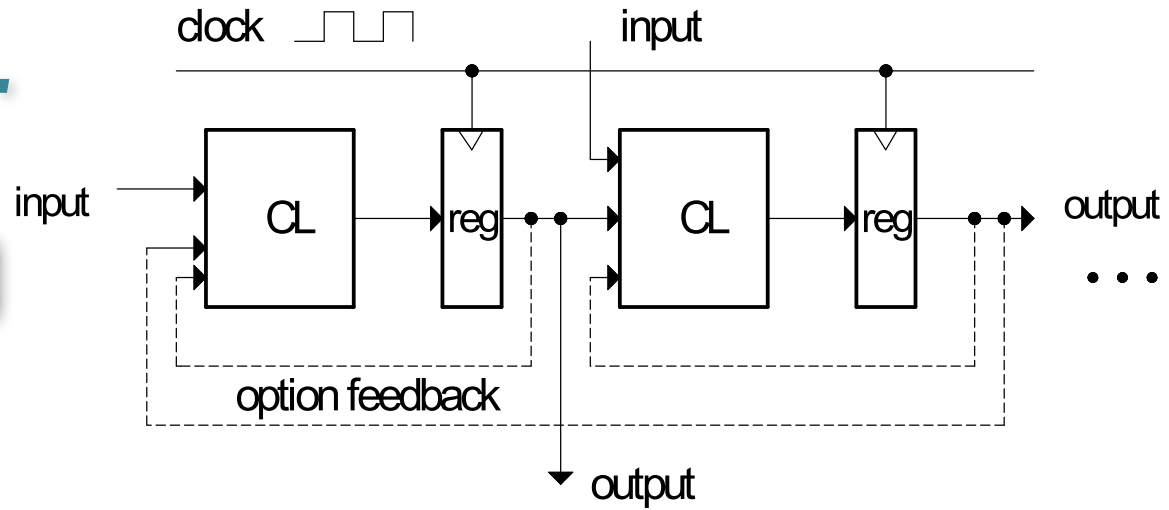


In General ...

For correct operation:

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

for all paths.



□ How do we enumerate **all** paths?

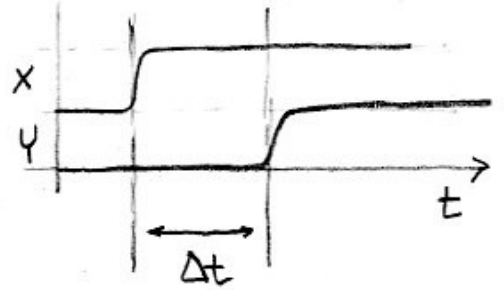
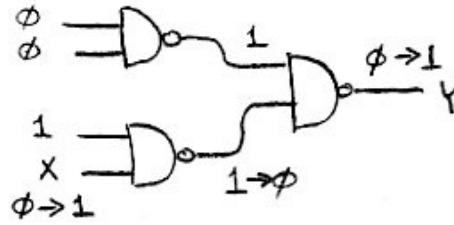
– Any circuit input or register output to any register input or circuit output?

• Note:

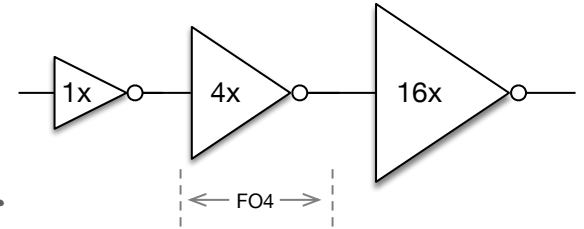
– “setup time” for outputs is a function of what it connects to.

– “clk-to-q” for circuit inputs depends on from where it comes.

“Gate Delay”



- ❑ Modern CMOS gate delays on the order of a few picoseconds. (However, highly dependent on gate design and context.)
- ❑ Often expressed as FO4 delays (fan-out of 4) - as a process dependent delay metric:
 - the delay of an **inverter**, driven by an inverter 4x smaller than itself, and driving an inverter 4x larger than itself.
 - Less than 10ps for a 32nm process. For a 7nm process FO4 is around 2.5ps.



Process Dependent FO4 Delay

Scaling equations for the accurate prediction of CMOS device performance from 180 nm to 7 nm

Aaron Stillmaker^{a,b,*}, Bevan Baas^a

^a Department of Electrical and Computer Engineering, University of California, Davis, One Shields Ave., Davis, CA 95616, USA

^b Department of Electrical and Computer Engineering, California State University, Fresno, 2320 E. San Ramon Ave., Fresno, CA 93740, USA

Characteristics of different technology nodes [23]. The modeled measurements are for a single inverter in an FO4 chain. The energy value is the average energy required for a single inverter transition from low to high, or high to low.

| Production Year | Technology Node (nm) | Technology Type | V _{DD} (V) | Simulated Performance of Inverter | | |
|-----------------|----------------------|-----------------|---------------------|-----------------------------------|-------------|------------|
| | | | | Delay (ps) | Energy (fJ) | Power (μW) |
| 1999 | 180 | Bulk | 1.8 | 77.2 | 27.5 | 105 |
| 2001 | 130 | Bulk | 1.2 | 34.7 | 5.20 | 26.1 |
| 2004 | 90 | Bulk | 1.1 | 26.5 | 2.62 | 13.0 |
| 2007 | 65 | Bulk | 1.1 | 19.8 | 1.72 | 8.58 |
| 2008 | 45 | High-k | 1.1 | 10.9 | 1.05 | 5.19 |
| 2010 | 32 | High-k | 0.97 | 9.8 | 0.51 | 2.47 |
| 2012 | 20 | Multi-Gate | 0.9 | 9.66 | 0.198 | 1.51 |
| 2013 | 16 ^a | Multi-Gate | 0.86 | 6.12 | 0.179 | 1.28 |
| 2013 | 14 ^a | Multi-Gate | 0.86 | 4.02 | 0.144 | 0.995 |
| 2015 | 10 | Multi-Gate | 0.83 | 3.24 | 0.122 | 0.866 |
| 2017 | 7 | Multi-Gate | 0.8 | 2.47 | 0.111 | 0.789 |

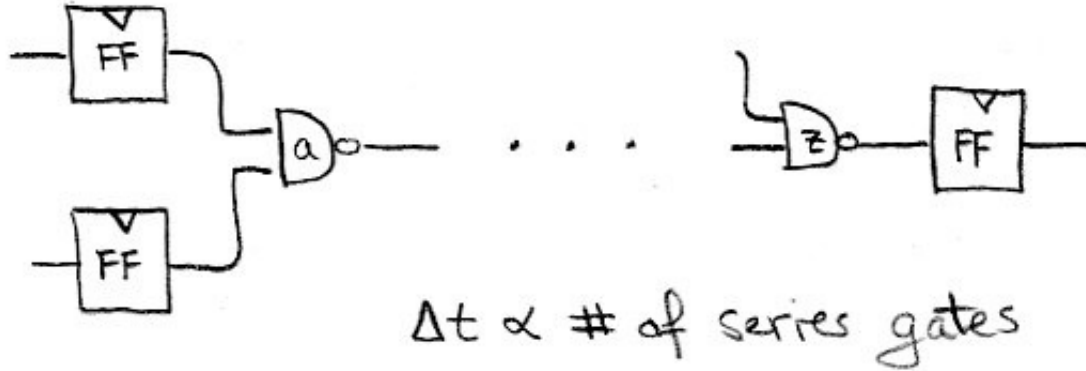
^a The 2013 ITRS report labels a single "16/14" node.

“Path Delay”

- For correct operation:

$$\text{Total Delay} \leq \text{clock_period} - FF_{\text{setup_time}} - FF_{\text{clk_to_q}}$$

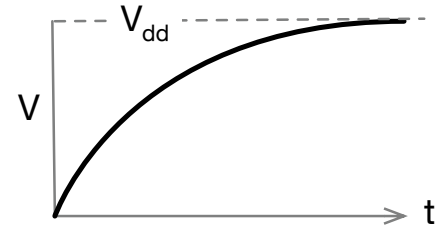
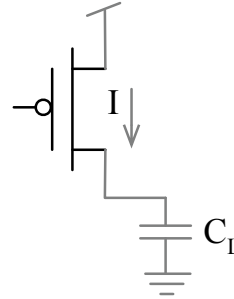
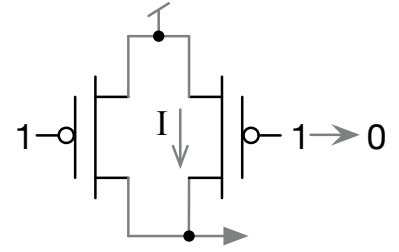
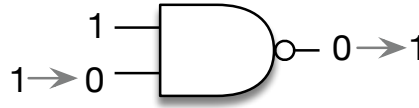
on all paths.



- ▶ High-speed processors critical paths (worst case paths) have around 30 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

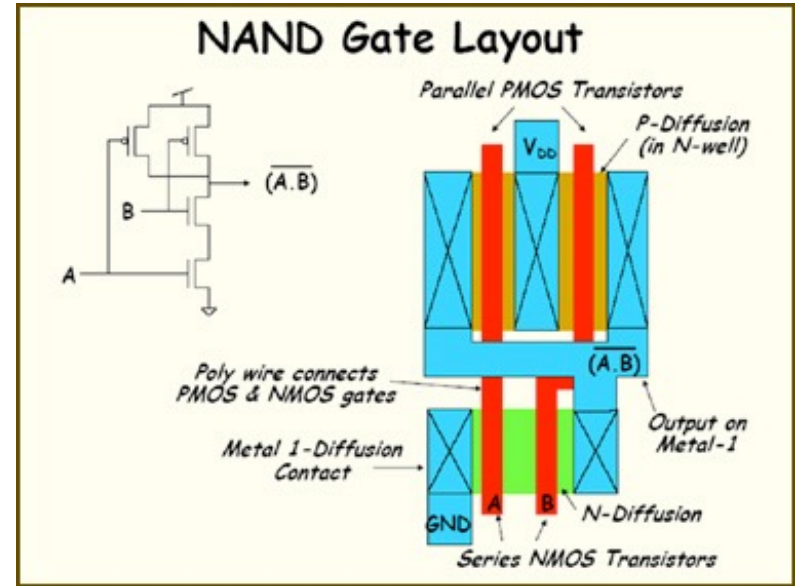


$$\Delta t \propto C_L / I$$

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

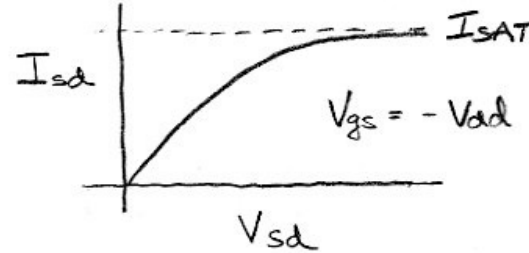
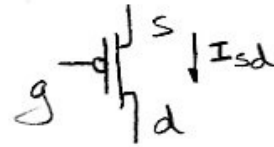
Physical Layout determines FET strength

- “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.



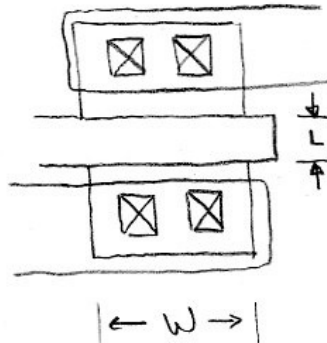
More on transistor Current

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

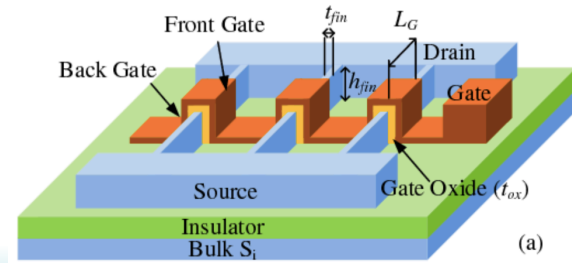


- I_{SAT} depends on process parameters (higher for nFETs than for pFETs) and transistor size (layout):

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



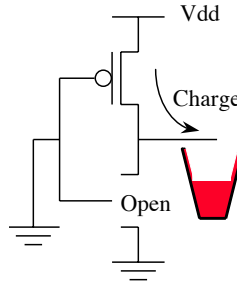
FinFets use multiple “fins” to get wider



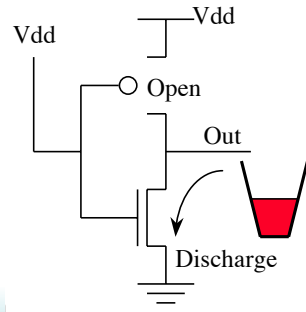
Transistors as water valves.

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

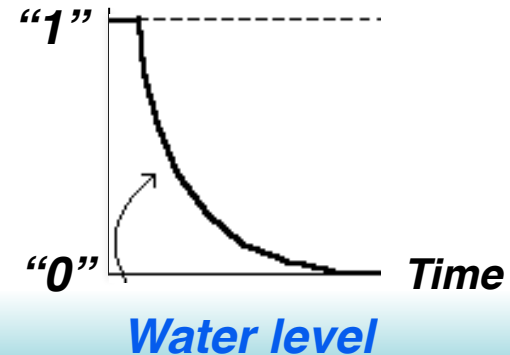
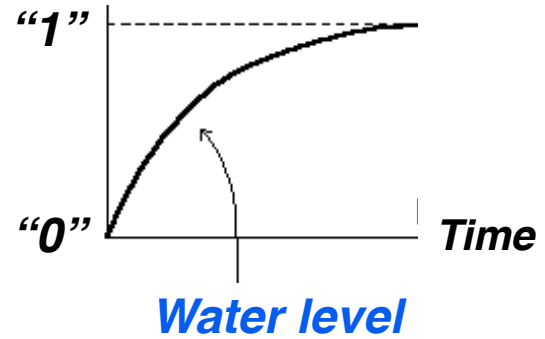
*A “on” p-FET fills
up the capacitor
with charge.*



*A “on” n-FET
empties the bucket.*



(Cartoon physics)



Inverter: Transient Response

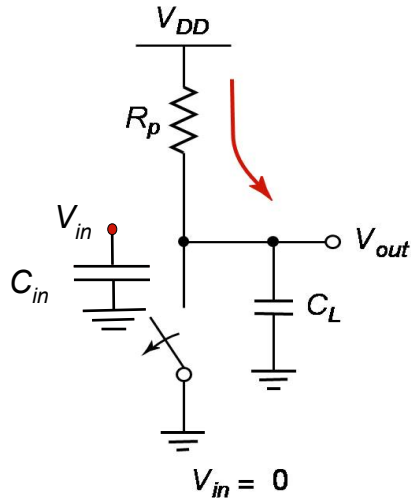
With:

resistive approximation for FETs,
high-to-low (HL)

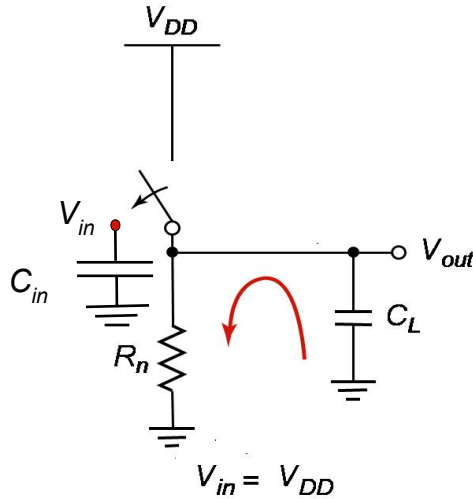


$$V(t) = V_0 e^{-t/RC}$$

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



(a) Low-to-high

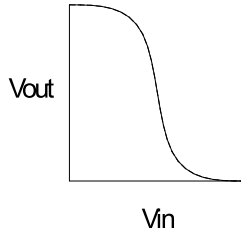
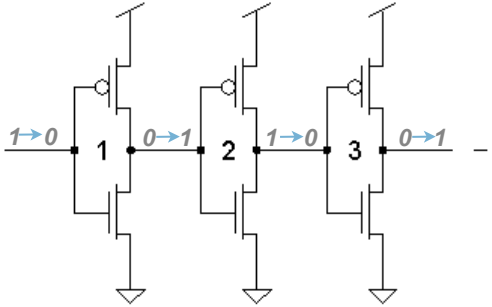
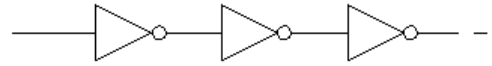


(b) High-to-low

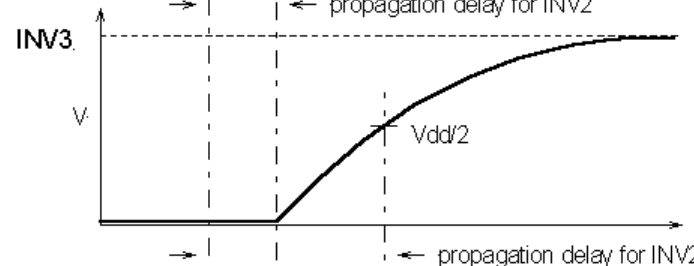
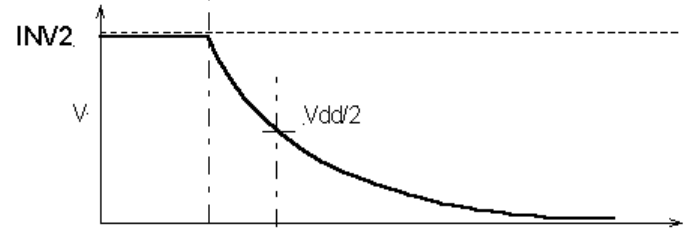
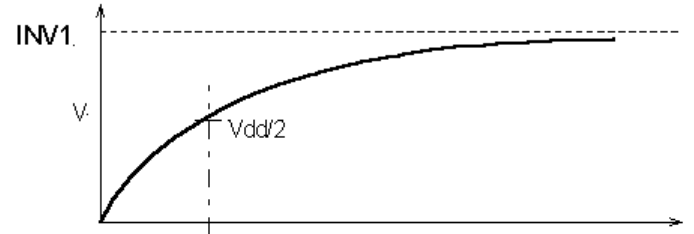
$$t_{pHL} = f(R_{on} C_L) \\ = 0.69 R_n C_L$$

Turning Rise/Fall Delay into Gate Delay

- Cascaded gates:



“transfer curve” for inverter.



← propagation delay for INV2

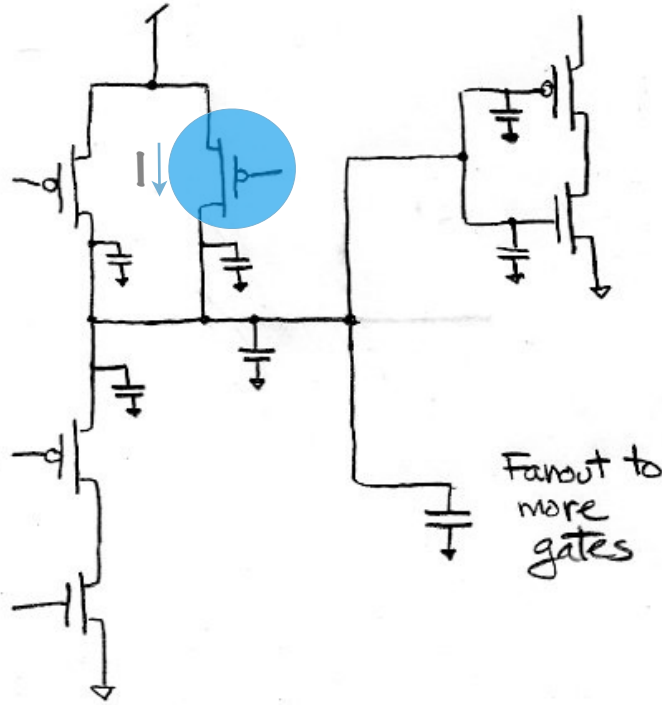
← propagation delay for INV2 & INV3 in series

In general:

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

More on gate delay

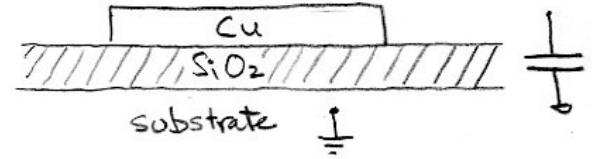
- Everything that connects to the output of a logic gate (or transistor) contributes capacitance:



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

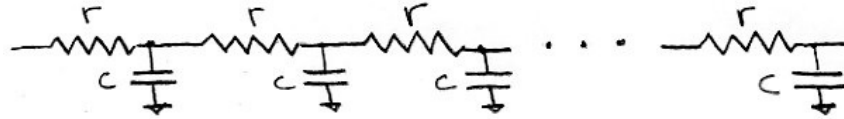
Wires

- As parallel plate capacitors:

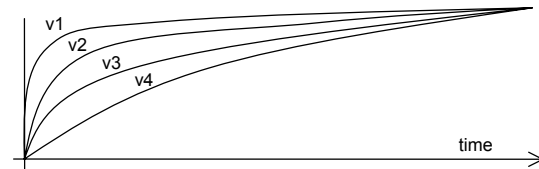
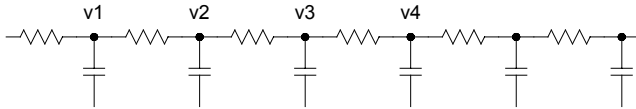


$$C \propto \text{Area} = \text{width} * \text{length}$$

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

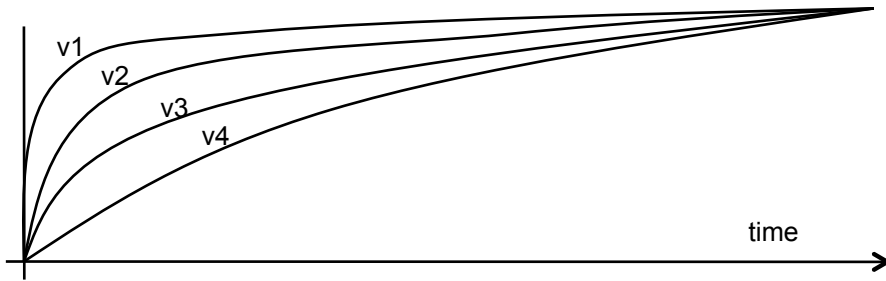
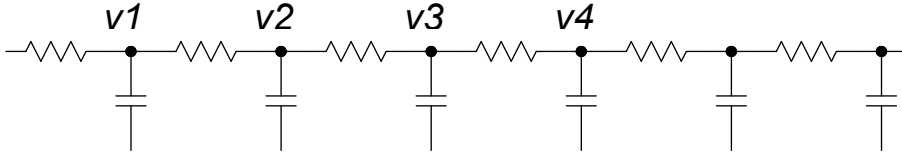


with $r = \text{res}/\text{length}$, $c = \text{cap}/\text{length}$, $\Delta t \propto r c L^2 \cong r c + 2rc + 3rc + \dots$



Wire Delay

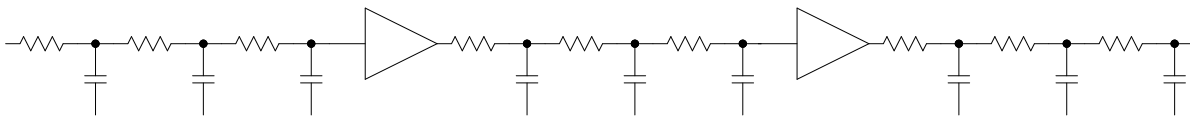
- Wires possess distributed resistance and capacitance
- Time constant associated with distributed RC is proportional to the **square** of the length



- 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
- For **long wires** on ICs with high currents:
 - inductance is also important

Wire Rebuffering

- For **long wires** on ICs:
 - *busses, clock lines, global control signal, etc.*
 - *Resistance is significant, therefore rcL^2 effect dominates.*
 - *signals are typically “rebuffered” to reduce delay:*



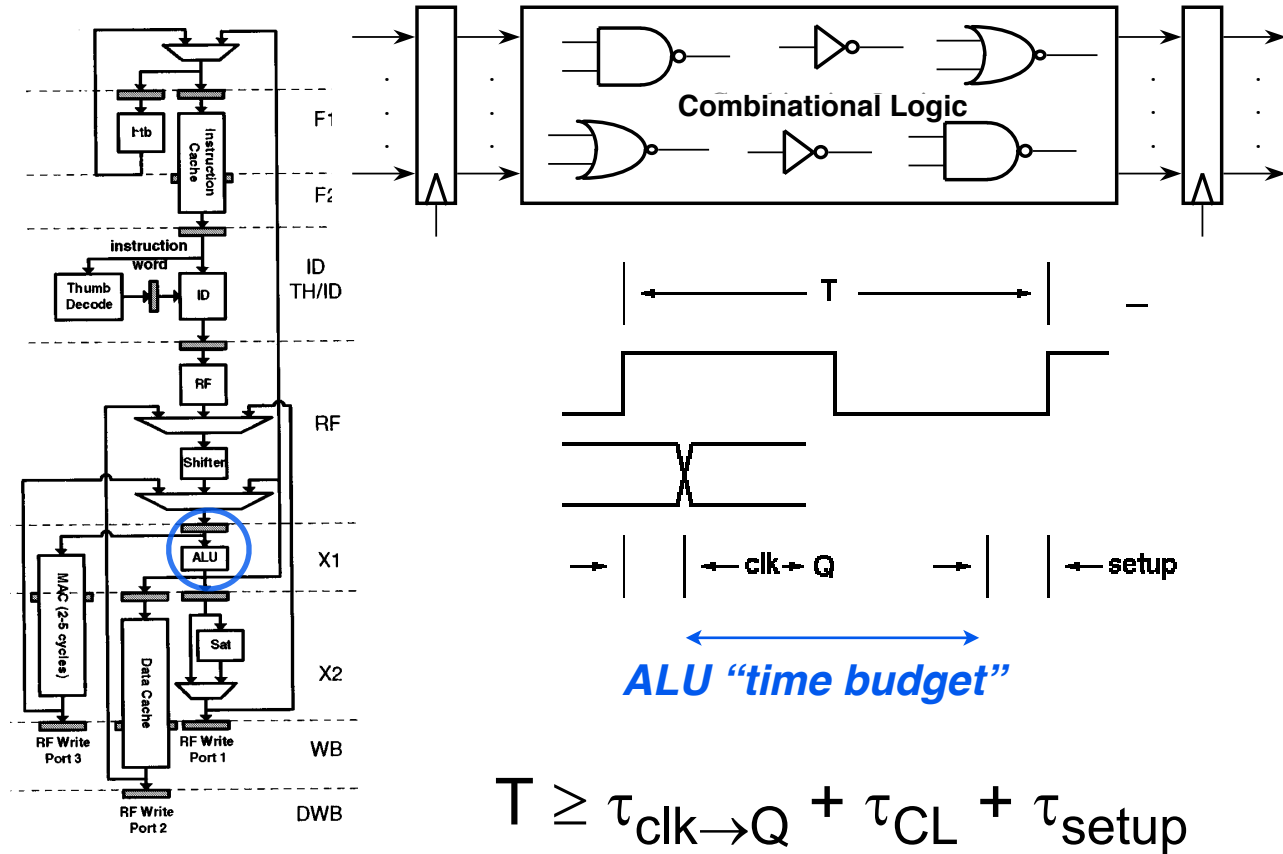
unbuffered wire $\Delta t \propto L^2$

wire buffered into N sections $\Delta t \propto N * (L/N)^2 + (N-1) * t_{buffer}$

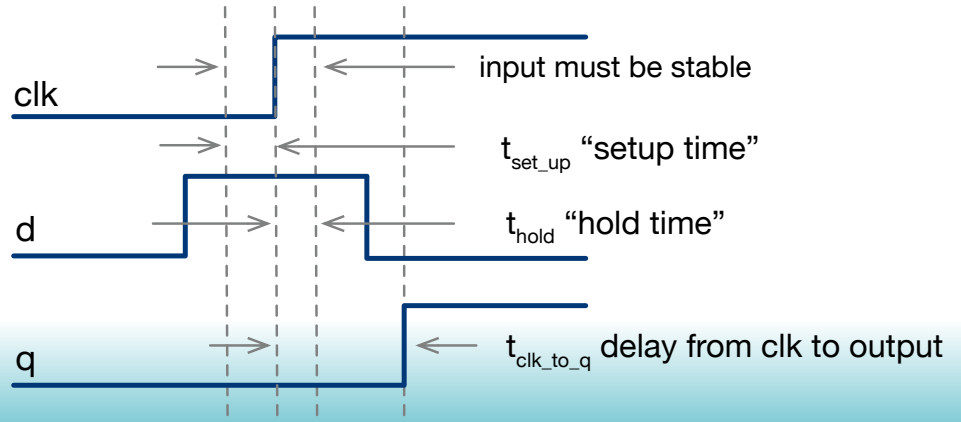
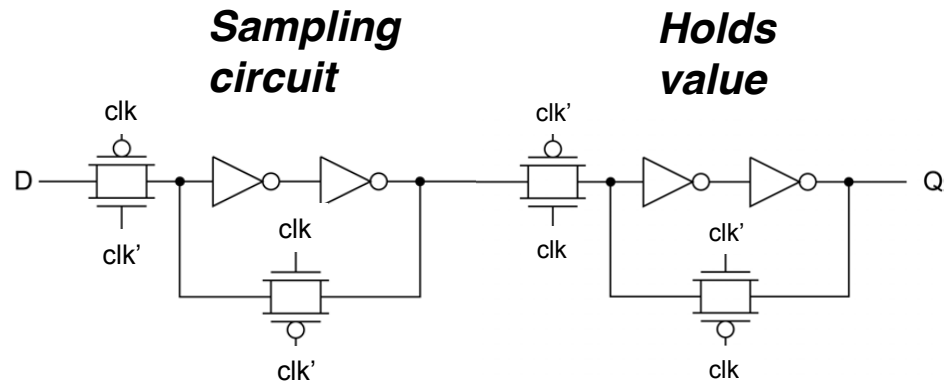
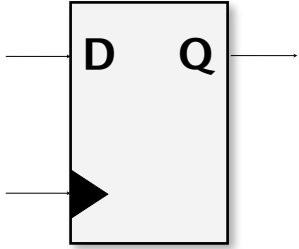
Assuming t_{buffer} is small, $\Delta t \propto L^2/N$

Speedup: $\propto N$

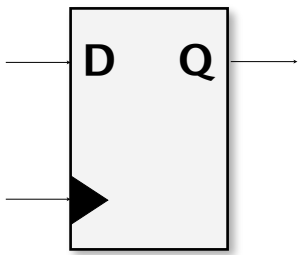
Flip-Flop delays eat into “time budget”



Recall: Positive edge-triggered flip-flop



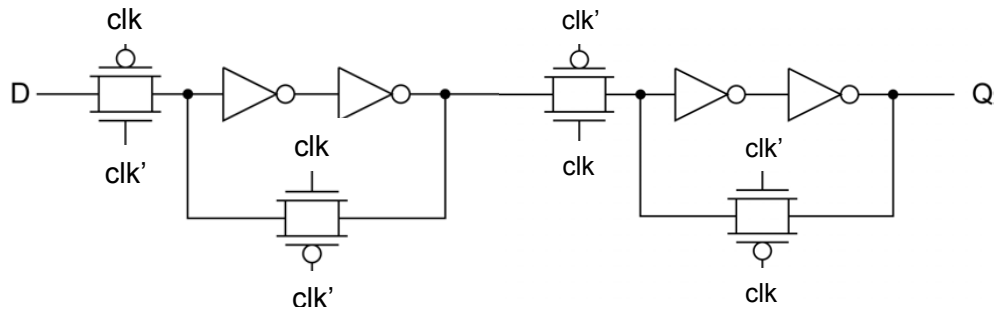
Sensing: When clock is low



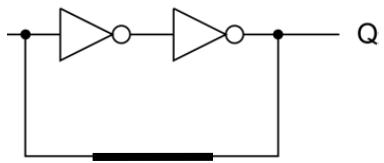
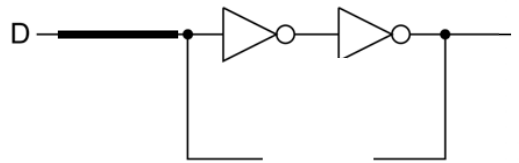
A flip-flop “samples” right before the edge, and then “holds” value.

Sampling circuit

Holds value



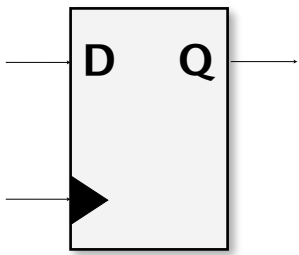
*clk = 0
clk' = 1*



Will capture new value on posedge.

Outputs last value captured.

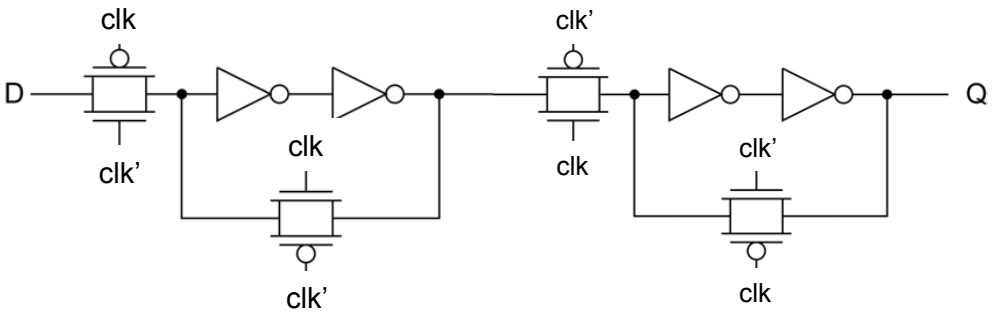
Capture: When clock goes high



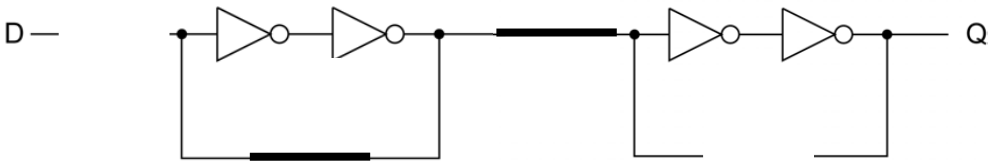
A flip-flop “samples” right before the edge, and then “holds” value.

Sampling circuit

Holds value



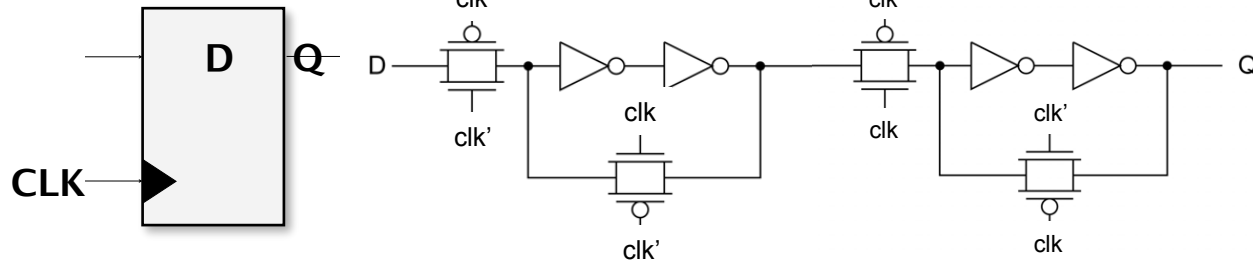
*clk = 1
clk' = 0*



Remembers value just captured.

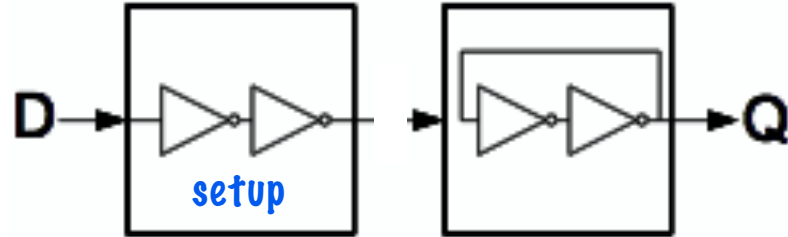
Outputs value just captured.

Flip Flop delays:



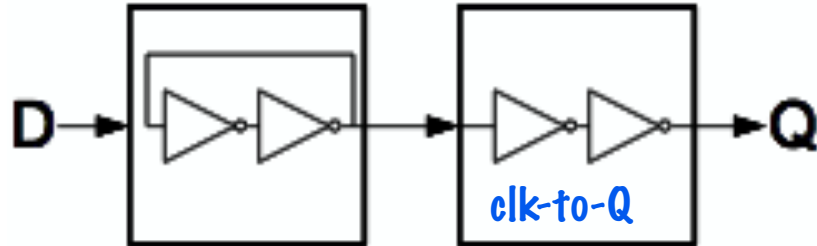
CLK == 0

Sense D, but Q outputs old value.



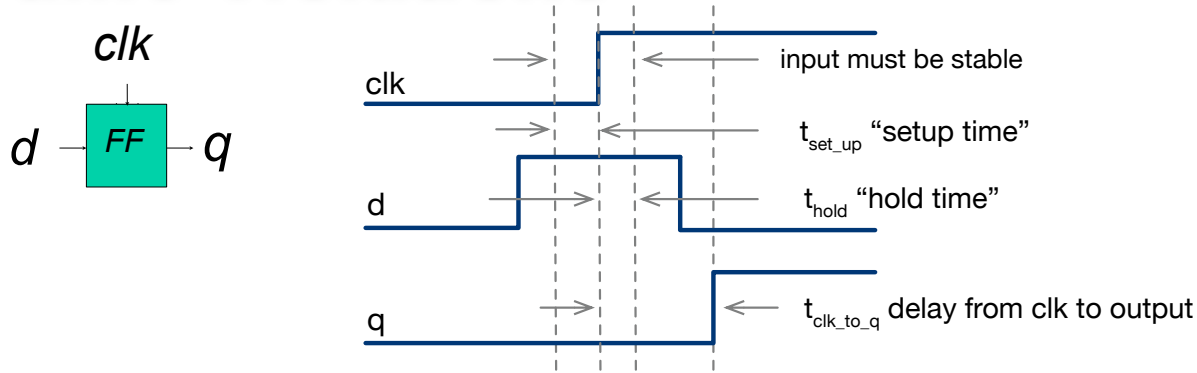
CLK 0->1

Capture D, pass value to Q



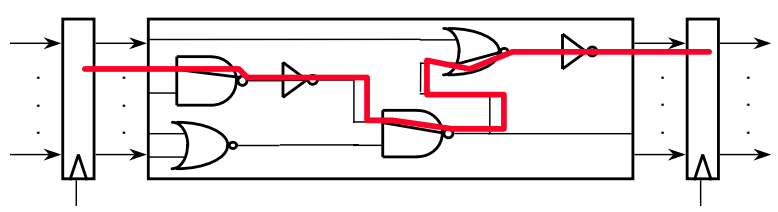
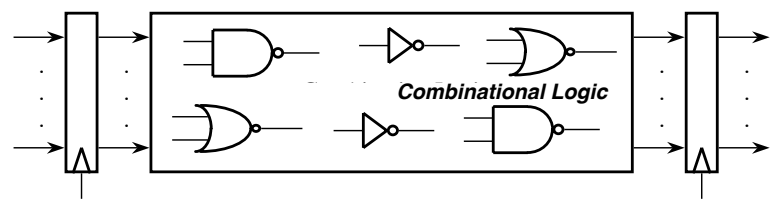
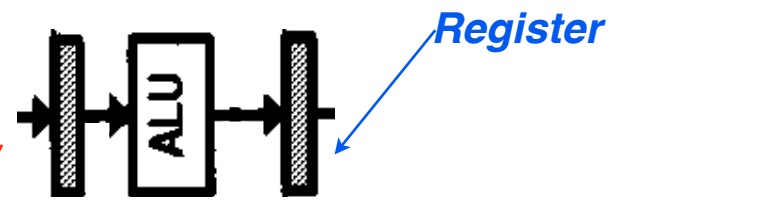
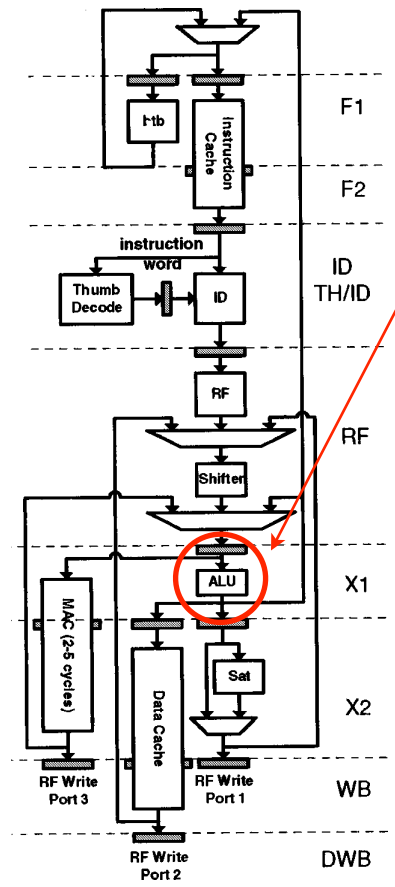
Note: with too much fanout, second stage could fail to capture data properly. Often output is rebuffered.

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.

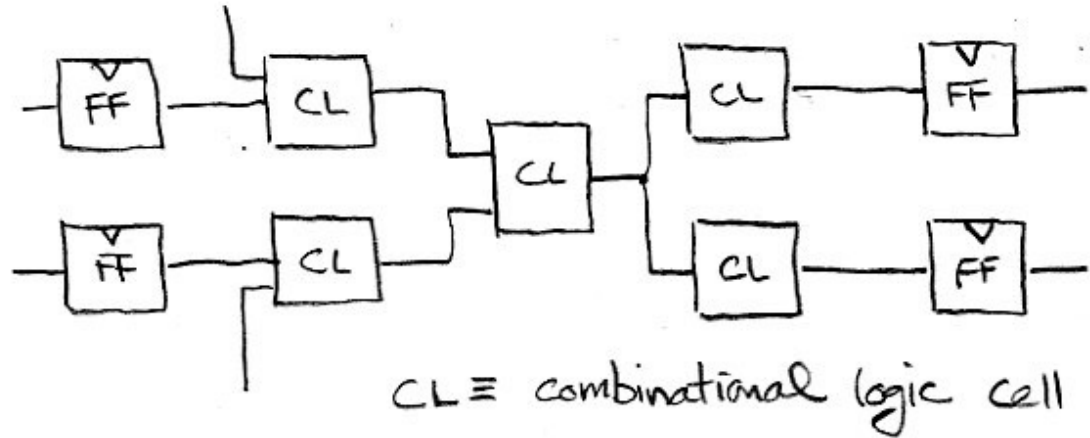
Timing Analysis and Logic Delay



Some path somewhere in the design has the longest delay and is therefore the "critical path".

Components of Combinational Path Delay

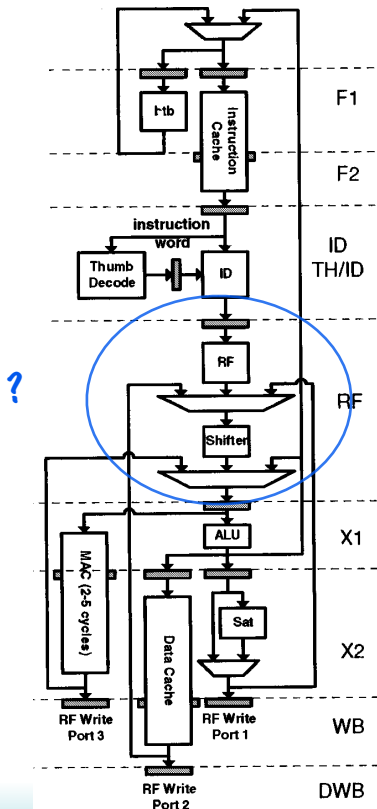
1. # of levels of logic
2. Internal cell delay
3. wire delay
4. cell input capacitance
5. cell fanout
6. cell output drive strength



Who controls the delay in ASIC?

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

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.

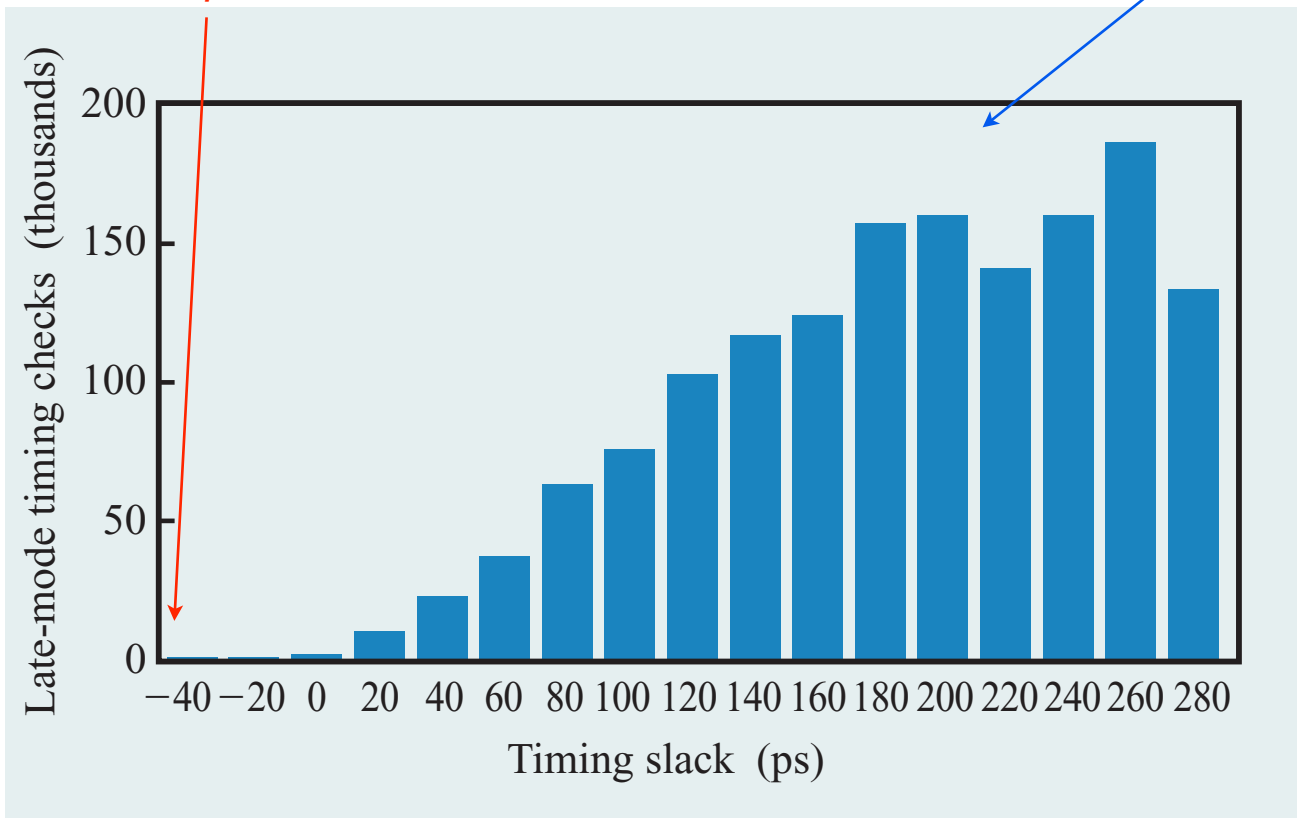
- Design tools help in the search.
- Synthesis tools work to meet clock constraint, report delays on paths,
 - Special static timing analyzers accept a design netlist and report path delays,
 - and, of course, simulators can be used to determine timing performance.

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

Most paths have hundreds of picoseconds to spare.

The critical path



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 HDL 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

Circuit retiming*

Critical path is 5 (ignore FF delay for now).

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.

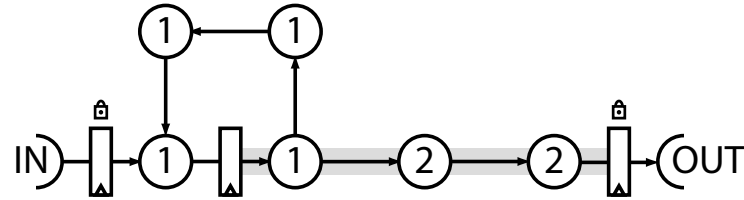


Figure 1: A small graph before retiming. The nodes represent logic delays, with the inputs and outputs passing through mandatory, fixed registers. The critical path is 5.

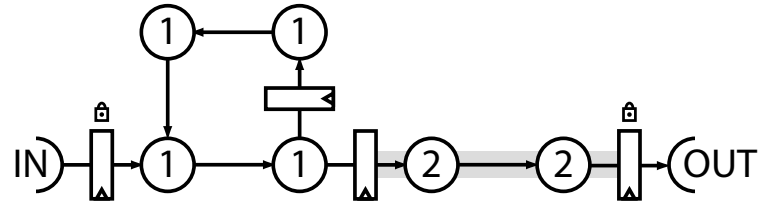
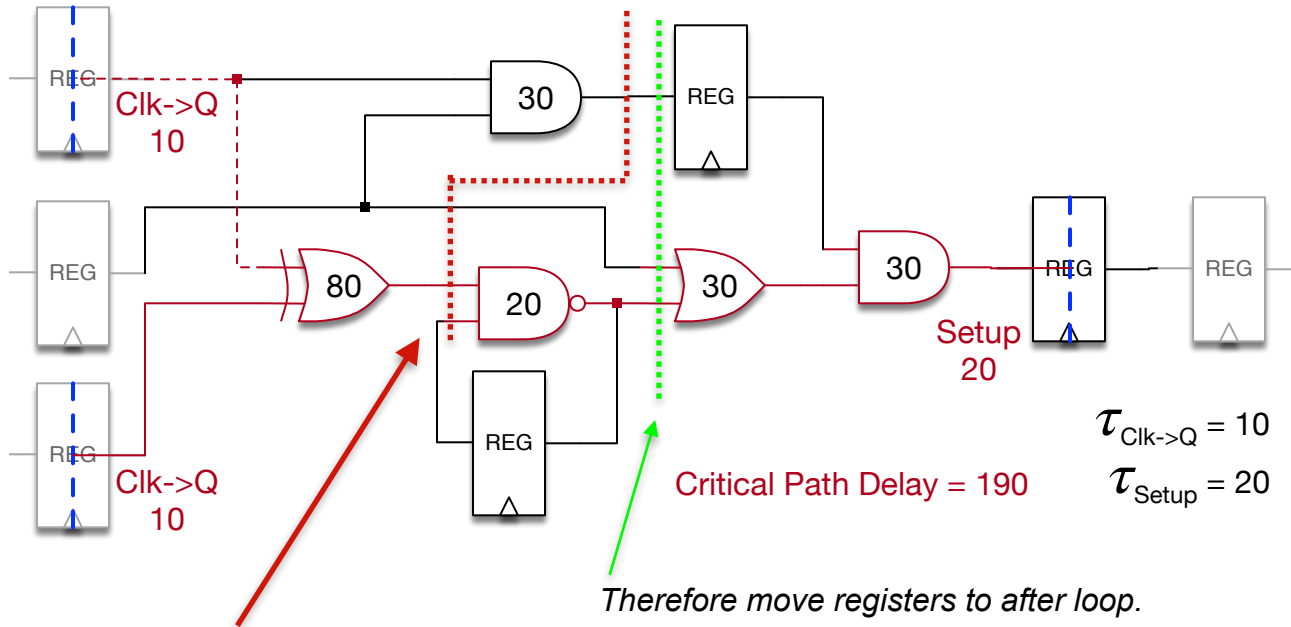


Figure 2: The example in Figure 2 after retiming. The critical path is reduced from 5 to 4.

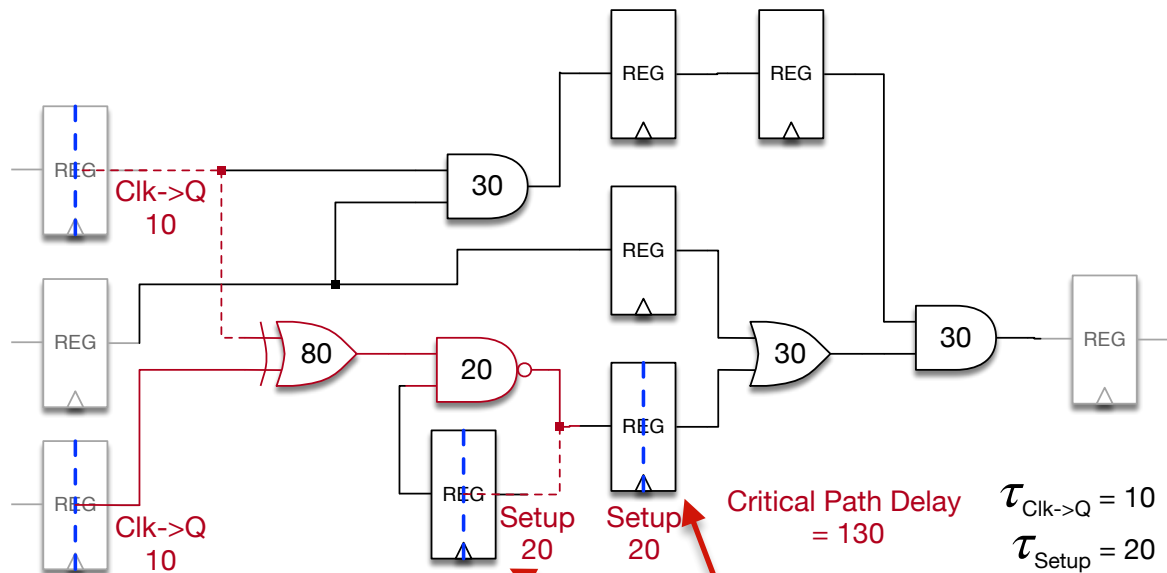
Logic Synthesis tools can do this in simple cases.

Retiming Example



Want to retime to here, however, delay cannot be added to the loop without changing the semantics of the logic. Because of this, many retiming tools stop at loops.

Retiming Example

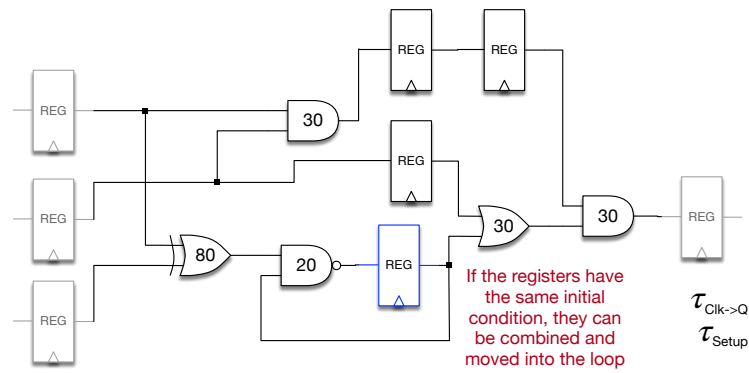


This is the retimed solution that many retiming aware tools will stop at.

This is also the optimal solution when the initial values of the registers are not given.

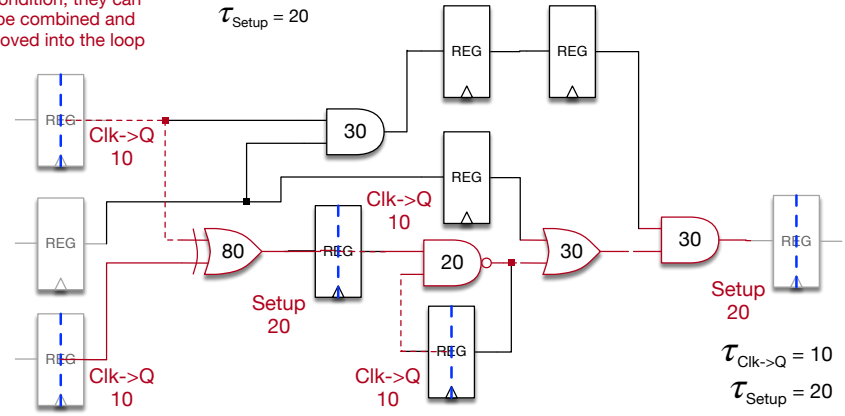
If the registers have the same initial condition, they can be combined and moved into the loop

Retiming Example



If the registers have the same initial condition, they can be combined and moved into the loop

$\tau_{\text{Clk} \rightarrow \text{Q}} = 10$
 $\tau_{\text{Setup}} = 20$

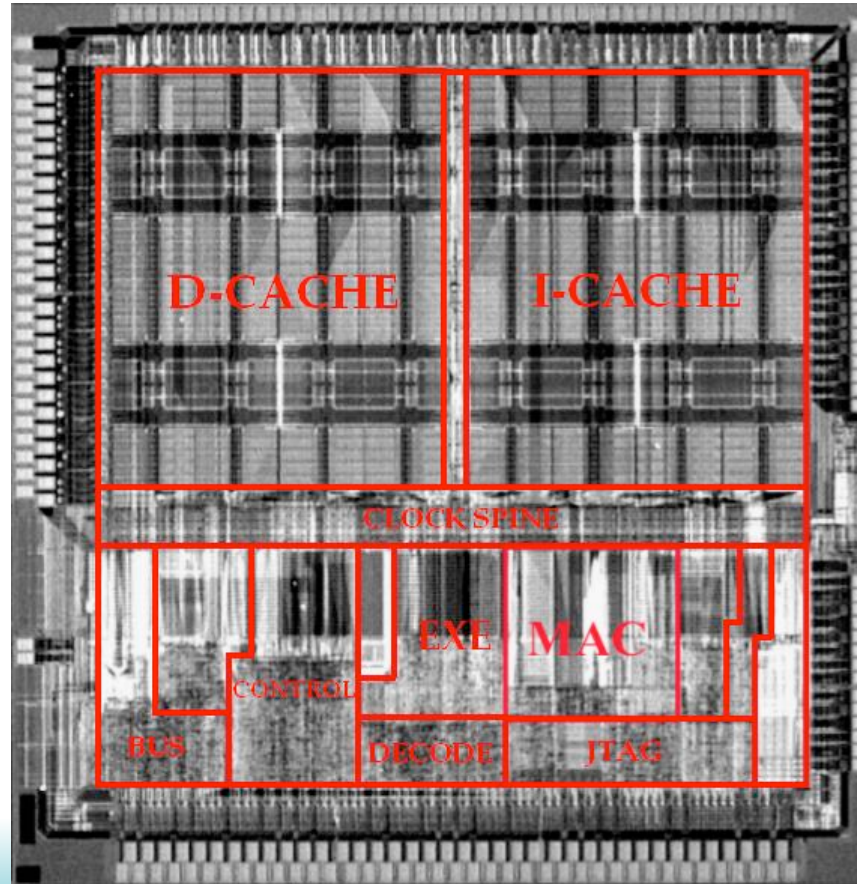
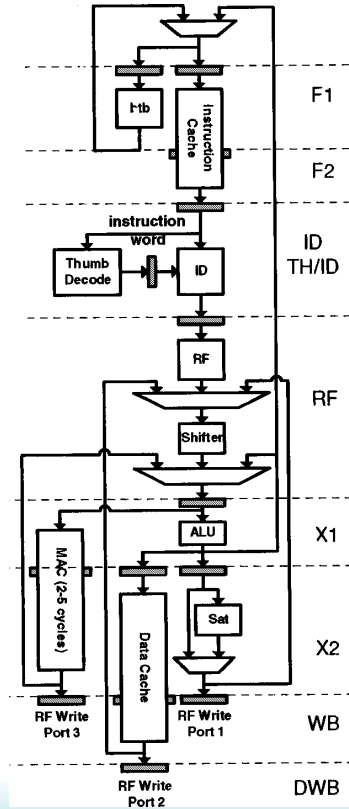


Setup 20
 $\tau_{\text{Clk} \rightarrow \text{Q}} = 10$
 $\tau_{\text{Setup}} = 20$

The register can be moved through the NAND gate, producing, a register at each input. The total delay in the loop is unchanged. Care must be taken to properly set the initial conditions of the registers.

Critical Path Delay = 110

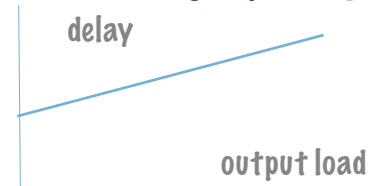
Floorplaning: essential to meet timing.



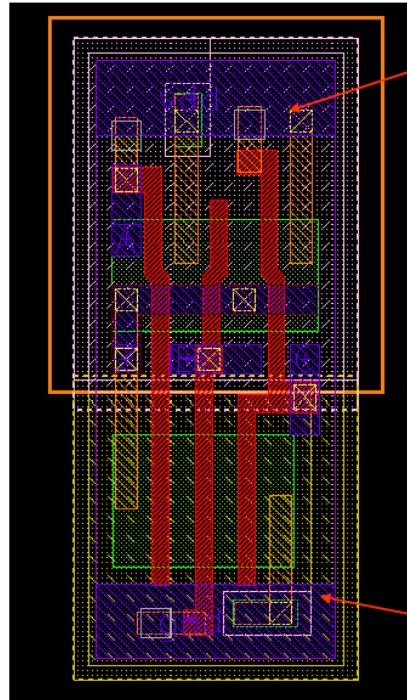
(Intel XScale 80200)

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.**



Standard cell characterization



Power Supply Line (V_{DD}) Delay in (ns)!!

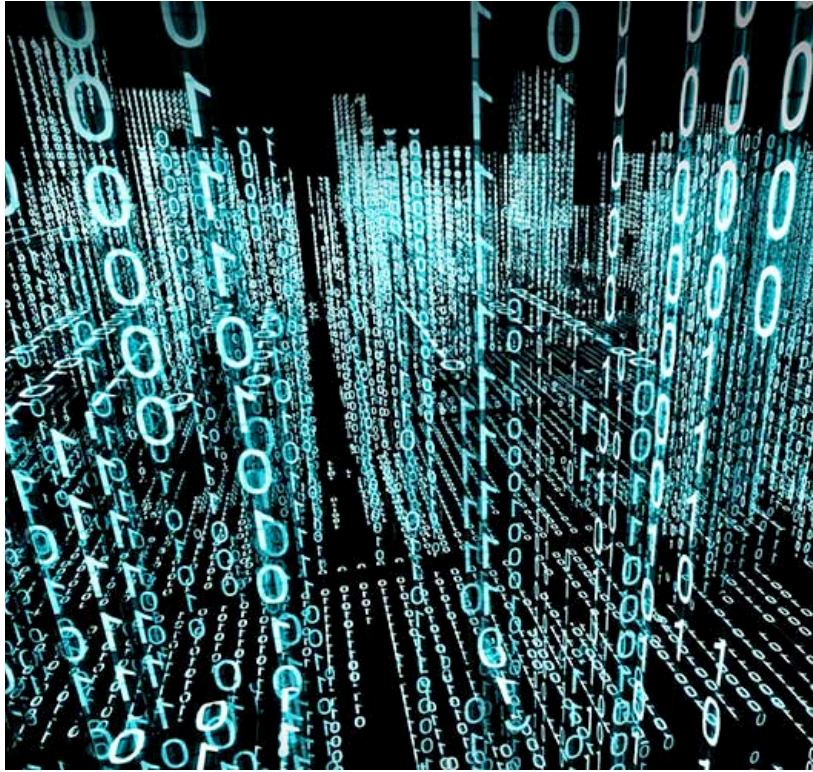
| Path | 1.2V - 125°C | 1.6V - 40°C |
|---------------|----------------------|----------------------|
| $In1-t_{pLH}$ | $0.073+7.98C+0.317T$ | $0.020+2.73C+0.253T$ |
| $In1-t_{pHL}$ | $0.069+8.43C+0.364T$ | $0.018+2.14C+0.292T$ |
| $In2-t_{pLH}$ | $0.101+7.97C+0.318T$ | $0.026+2.38C+0.255T$ |
| $In2-t_{pHL}$ | $0.097+8.42C+0.325T$ | $0.023+2.14C+0.269T$ |
| $In3-t_{pLH}$ | $0.120+8.00C+0.318T$ | $0.031+2.37C+0.258T$ |
| $In3-t_{pHL}$ | $0.110+8.41C+0.280T$ | $0.027+2.15C+0.223T$ |

3-input NAND cell
(from ST Microelectronics):
C = Load capacitance
T = input rise/fall time
Ground Supply Line (GND)

- Each library cell (FF, NAND, NOR, INV, etc.) and the variations on size (strength of the gate) is fully characterized across temperature, loading, etc.

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 may need some help at the floorplan and layout level
 - Hand instantiate cells, layout generators



End of Lecture 11