

EECS151 : Introduction to Digital Design and ICs

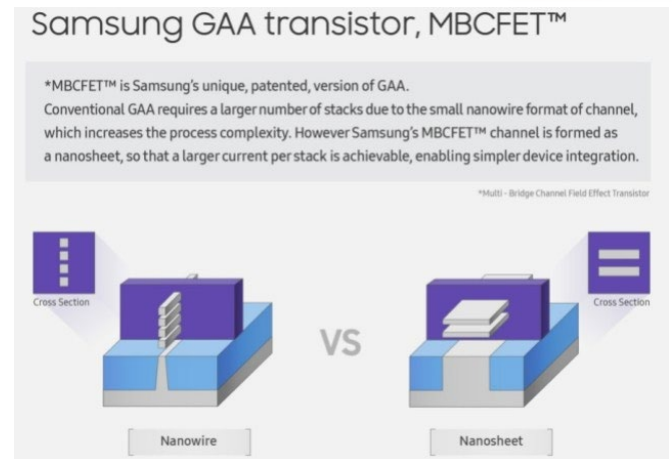
Lecture 15 – Logical Effort

Bora Nikolić



Samsung Foundry Promises Gate All-Around in '22

October 14, 2021, EETimes - Samsung Foundry recently held its Foundry Forum where it revealed some details of its semiconductor process roadmaps and fab expansion. Samsung is being most aggressive pursuing the next generation of transistor technology, with plans to reach mass production ahead of TSMC and Intel. Samsung's 3-nanometer process will use the gate-all-around (GAA) transistor structure, which the foundry calls MBCFET (Multi-bridge channel FET) and will be in production first half of 2022. TSMC will wait another generation until its N2 process to deliver GAA some time in 2023.



EETimes

Review

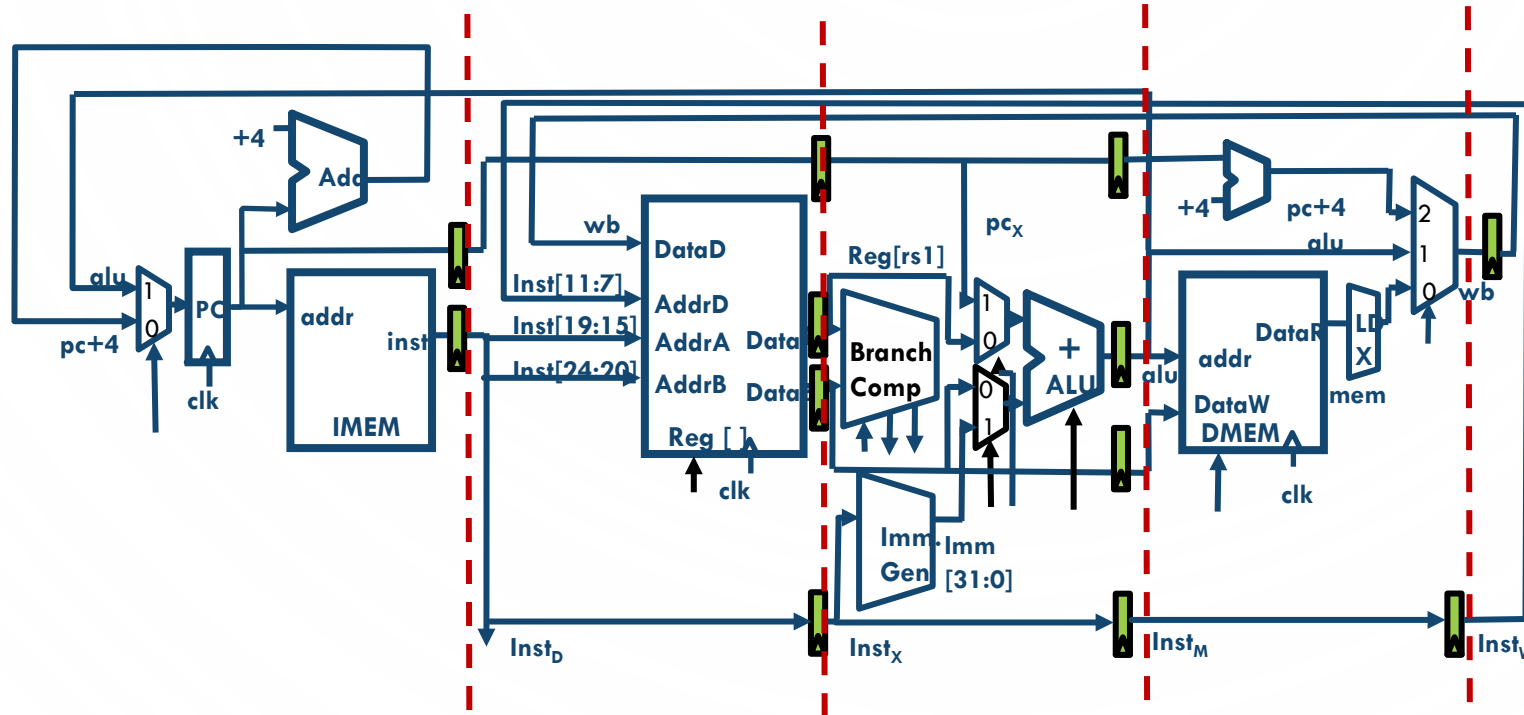
- Delay is a linear function of R and C
- Delay optimization is critical to improve the frequency of the circuit.
- The dimensions of a transistor affect its capacitance and resistance.
- We use RC delay model to describe the delay of a circuit.



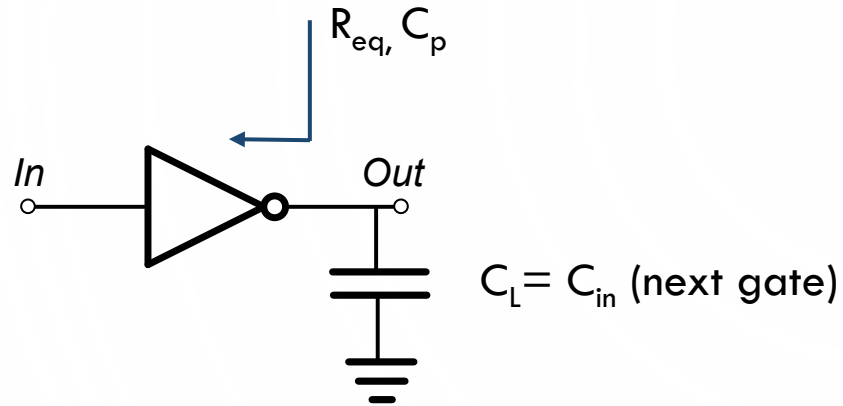
Minimizing Logic Delay

How Do We Optimize the Delay?

- How fast can a pipelined processor run?
- What is the fastest adder?



Inverter RC Delay



- $t_p = R_{eq}(C_p + C_L) = R_{eq}(\gamma C_{in} + C_L)$
 - $\gamma = 1$ (closer to 1.2 in recent processes)
- $t_p = R_{eq} C_{in}(1 + C_L/C_{in}) = \tau_{INV}(1 + f)$
 - Propagation delay is proportional to fanout
- Normalized Delay = $1 + f$

$$\text{Fanout} = f = C_L/C_{in}$$

$$t_p = \tau_{INV}(1 + f)$$

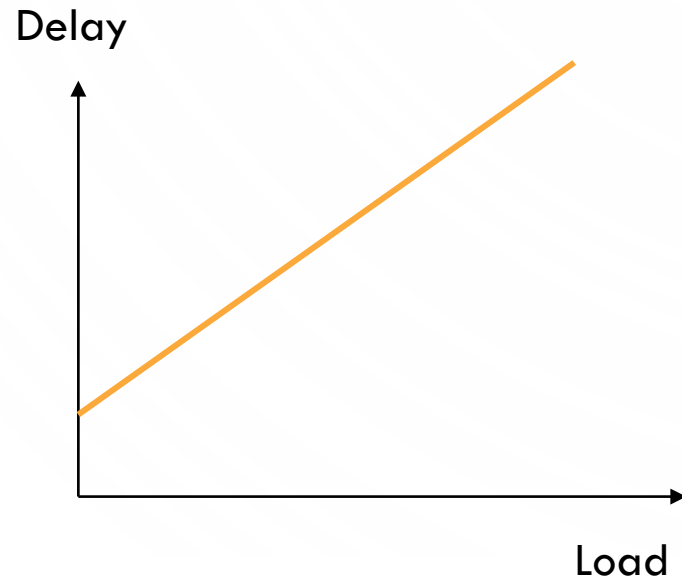
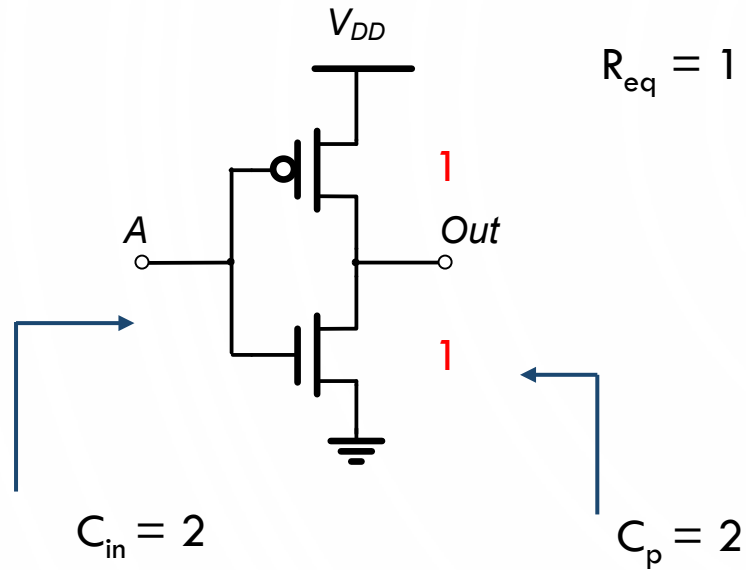
Generalizing to Arbitrary Gates

- Delay has two components: $d = h + p$
- h : effort delay = gf (a.k.a. stage effort)
 - Again has two components
- g : logical effort
 - Measures relative ability of gate to deliver current
 - $g = 1$ for inverter
- f : electrical effort = $C_{\text{out}} / C_{\text{in}}$
 - Ratio of output to input capacitance
 - Sometimes called fanout
- p : parasitic delay
 - Represents delay of gate driving no load
 - Set by internal parasitic capacitance

Note: There are differences in notation between semesters!

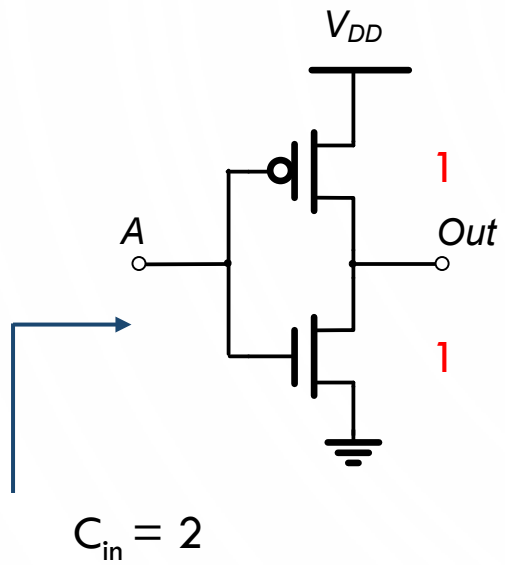
We shouldn't have done this!

Inverter Delay



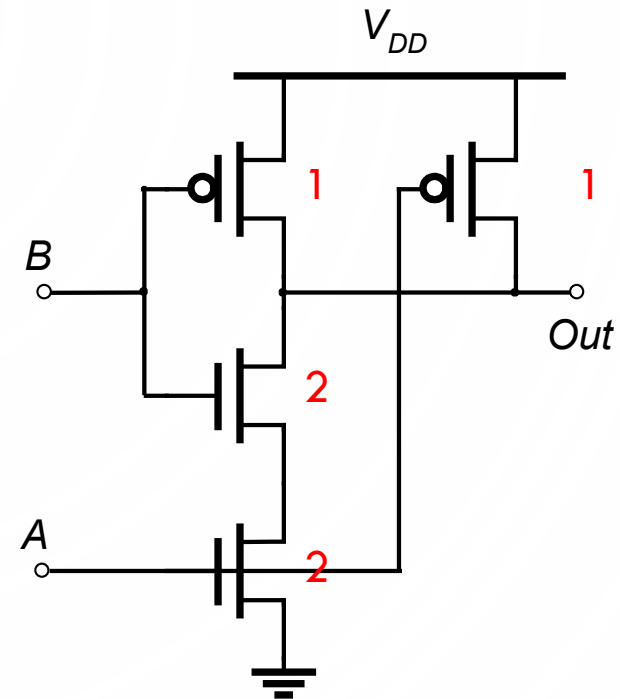
- Parasitic p is the ratio of intrinsic capacitance to an inverter
 - $p(\text{inverter}) =$
- Logical Effort g is the ratio of input capacitance to an inverter
 - $g(\text{inverter}) =$
- Electrical Effort h is the ratio of the load capacitance to the input capacitance
 - $h(\text{inverter}) =$
- $\text{Delay} = p + h = p + g * f = 1 + f$

NAND2 Gate



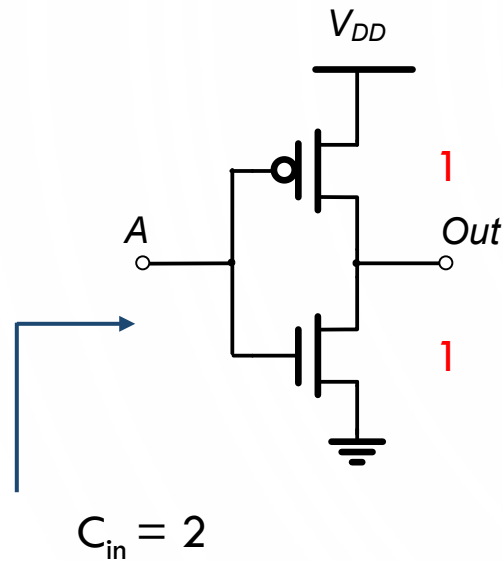
$$R_{eq} = 1$$

$$C_{in} = 3$$



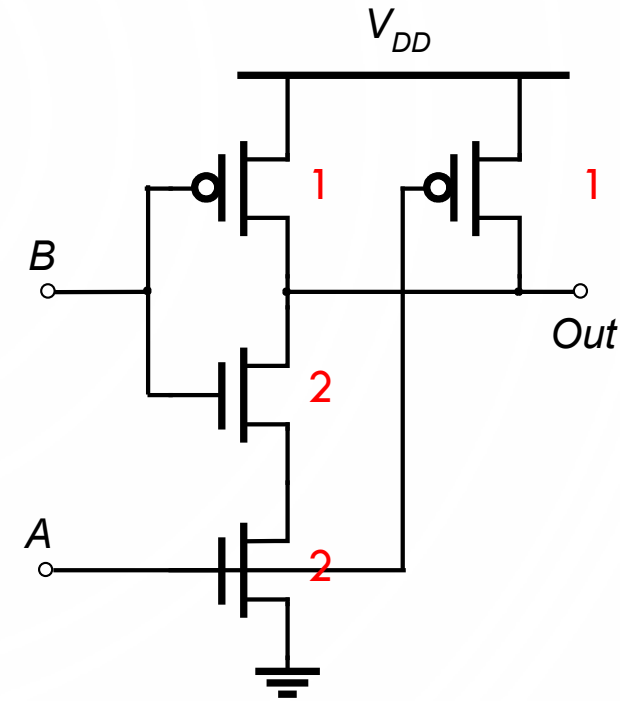
$$R_{eq} = 1$$

Logical Effort of NAND2 Gate



$$R_{eq} = 1$$

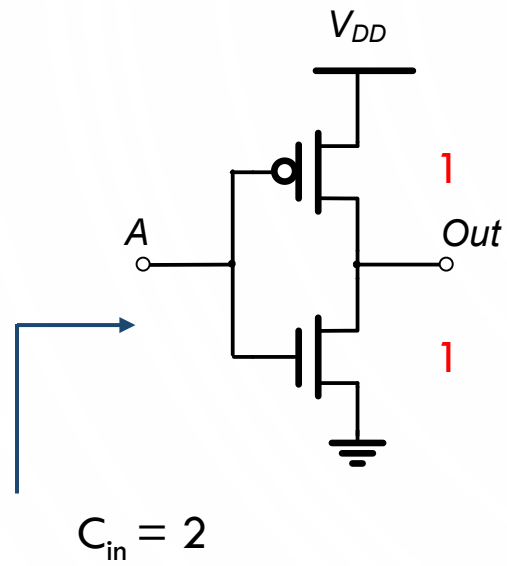
$$C_{in} = 3$$



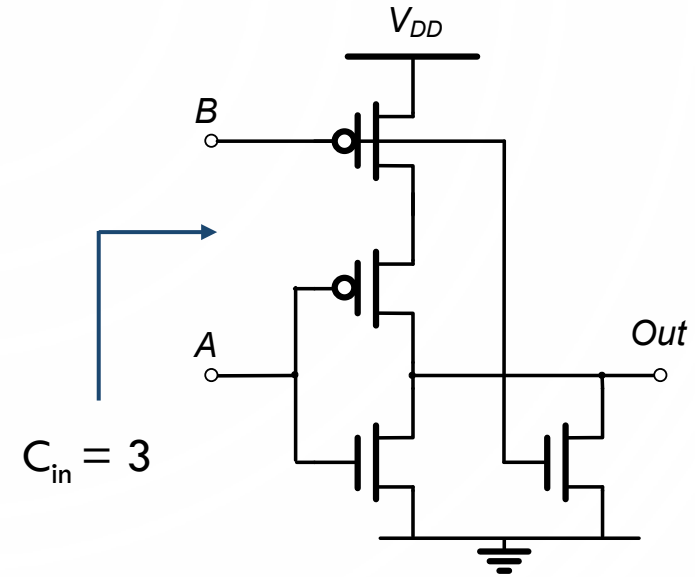
$$R_{eq} = 1$$

- In velocity-saturated devices I_{on} of a stack is $2/3$ (not a half) of two devices
 - So the correct upsizing factor is 1.5 (not 2)
- We will use 2, as it makes calculations easier

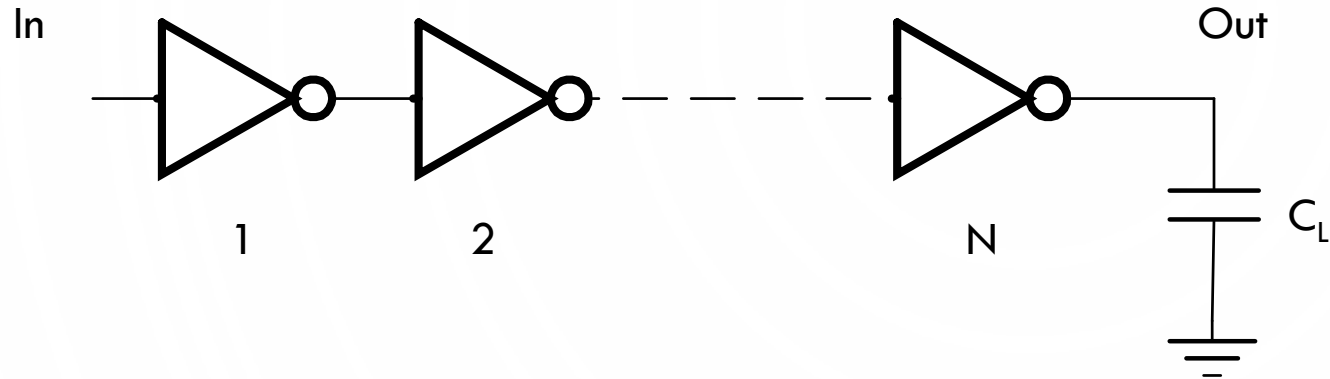
NOR2 Gate



$$R_{eq} = 1$$



Example: Inverter Chain



Logical Effort: $g =$

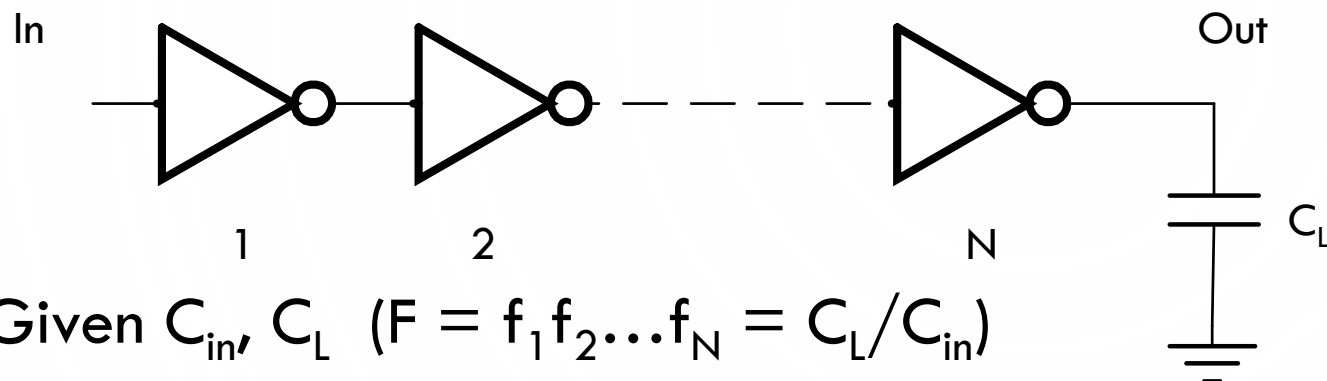
Electrical Effort: $f =$

Parasitic Delay: $p =$

Stage Delay: $d =$

Total Delay: $d_{\text{total}} =$

Optimize Delay of an Inverter Chain



Given C_{in}, C_L ($F = f_1 f_2 \dots f_N = C_L / C_{in}$)

How to optimally size inverter chain to minimize delay?

...There are $N-1$ unknowns: C_2, C_3, \dots, C_N

$$d = (1 + C_2 / C_{in}) + (1 + C_3 / C_2) + \dots + (C_L / C_N)$$

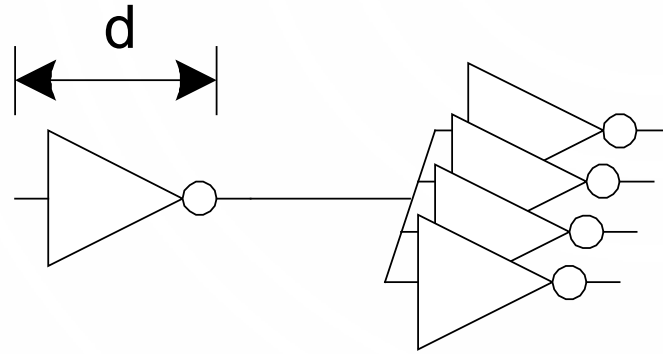
Solution: All delays are equal, $C_2 / C_{in} = C_3 / C_2 = \dots = C_L / C_N$

$$\frac{C_{i+1}}{C_i} = \sqrt[N]{\frac{C_L}{C_{in}}}$$

$$C_i = \sqrt{C_{i+1} C_{i-1}}$$

Example: FO4 Inverter

- Estimate the delay of a fanout-of-4 (FO4) inverter



Logical Effort: $g =$

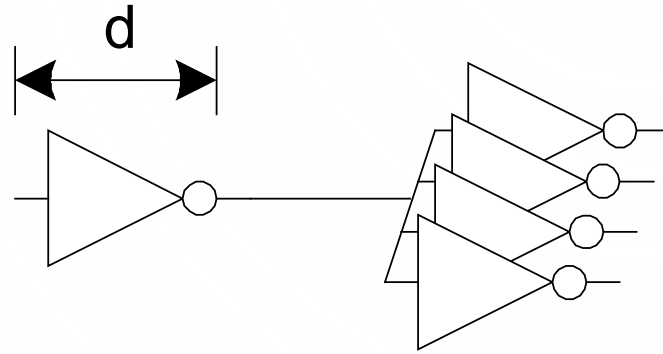
Electrical Effort: $f =$

Parasitic Delay: $p =$

Stage Delay: $d =$

Example: FO4 Inverter

- Estimate the delay of a fanout-of-4 (FO4) inverter



Logical Effort: $g = 1$

Electrical Effort: $f = 4$

Parasitic Delay: $p = 1$

Stage Delay: $d = 5$

Fanout-of-4 is commonly used to normalize the circuit delay across technologies

Multi-stage Logic Networks

- Logical effort generalizes to multistage networks

- *Path Logical Effort*

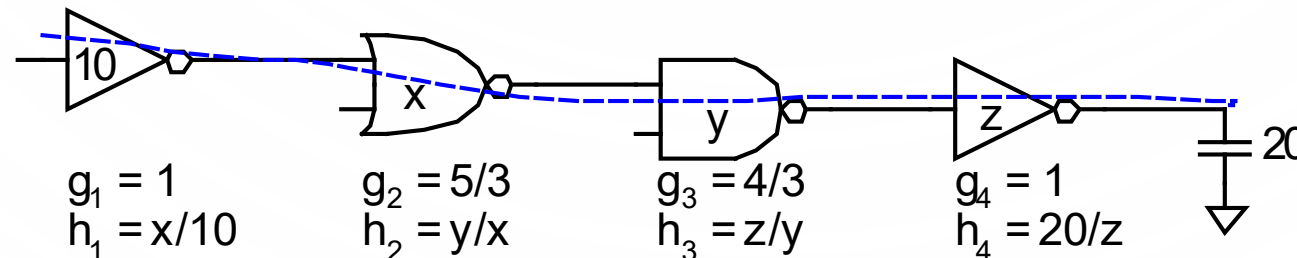
$$G = \prod g_i$$

- *Path Electrical Effort*

$$F = \frac{C_{\text{out-path}}}{C_{\text{in-path}}}$$

- *Path Effort*

$$H = \prod h_i = \prod g_i f_i$$



Branching Effect

$$b = \frac{C_{\text{on path}} + C_{\text{off path}}}{C_{\text{on path}}}$$

$$B = \prod b_i$$

$$G = 1$$

$$F = 90 / 5 = 18$$

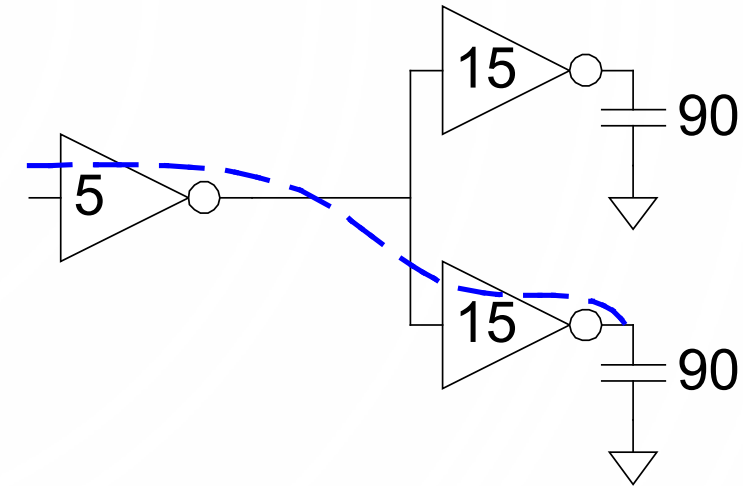
$$GF = 18$$

$$f_1 = (15 + 15) / 5 = 6$$

$$f_2 = 90 / 15 = 6$$

$$B = 2$$

$$F = g_1 g_2 h_1 h_2 = 36 = BGH$$



Designing Fast Circuits

$$D = \sum d_i = H + P$$

- Delay is smallest when each stage bears same effort

$$\hat{h} = g_i f_i = H^{\frac{1}{N}}$$

- Thus minimum delay of N stage path is

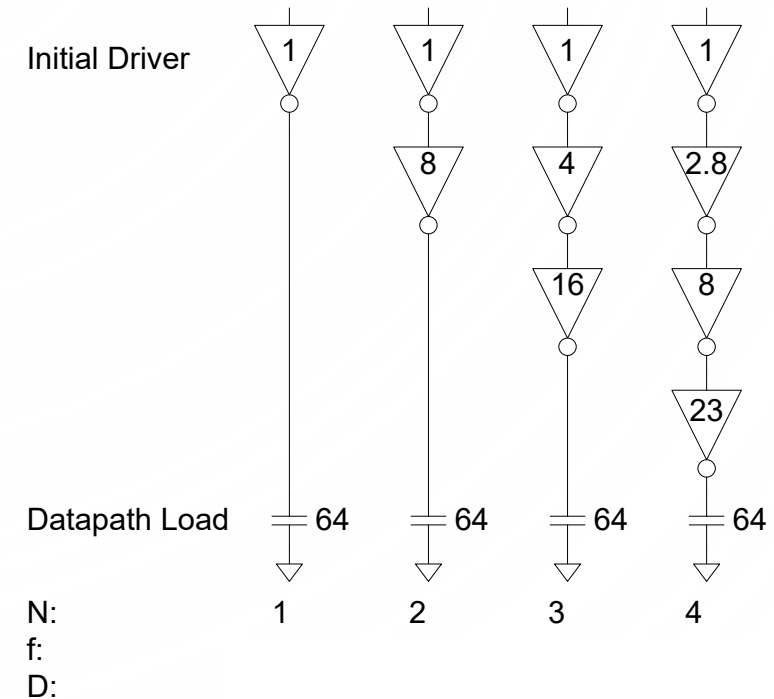
$$D = NH^{\frac{1}{N}} + P$$

- And we can find the gate sizes that result in optimal delay

Example: Best Number of Stages

- How many stages should a path use?
 - Minimizing number of stages is not always fastest
- Example: drive 64-bit datapath with unit inverter

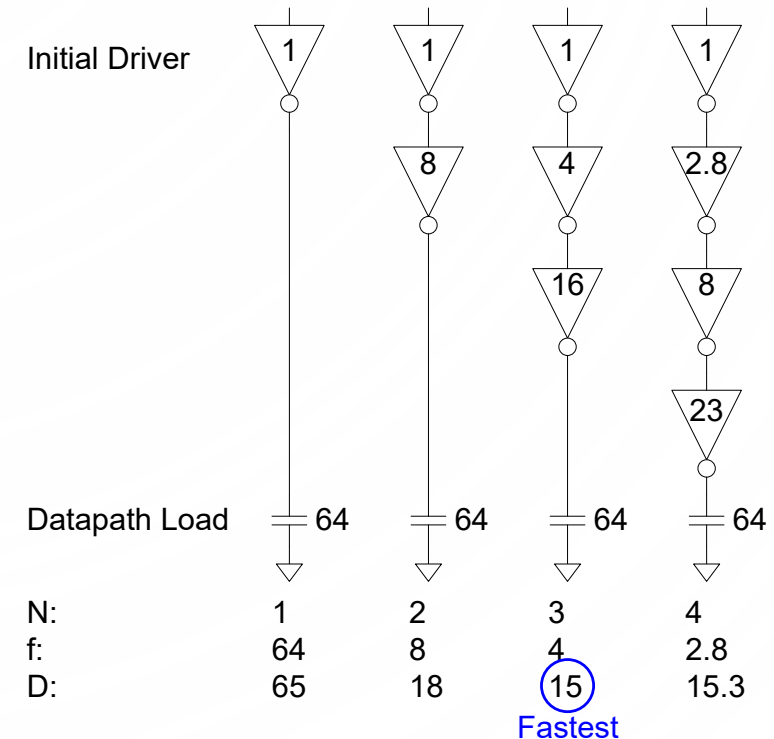
$$D = NF^{1/N} + P$$
$$= N(64)^{1/N} + N$$



Example: Best Number of Stages

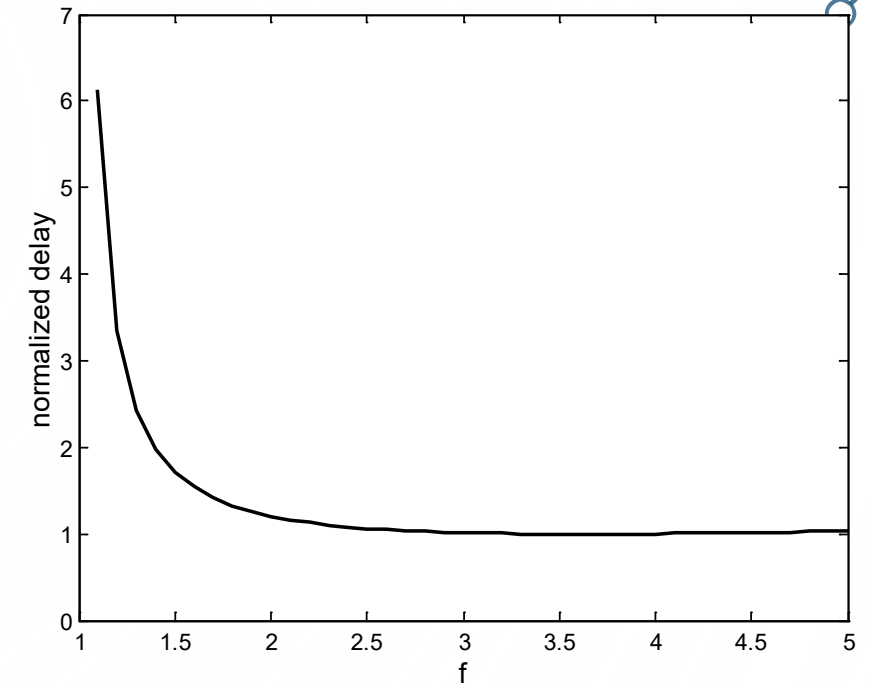
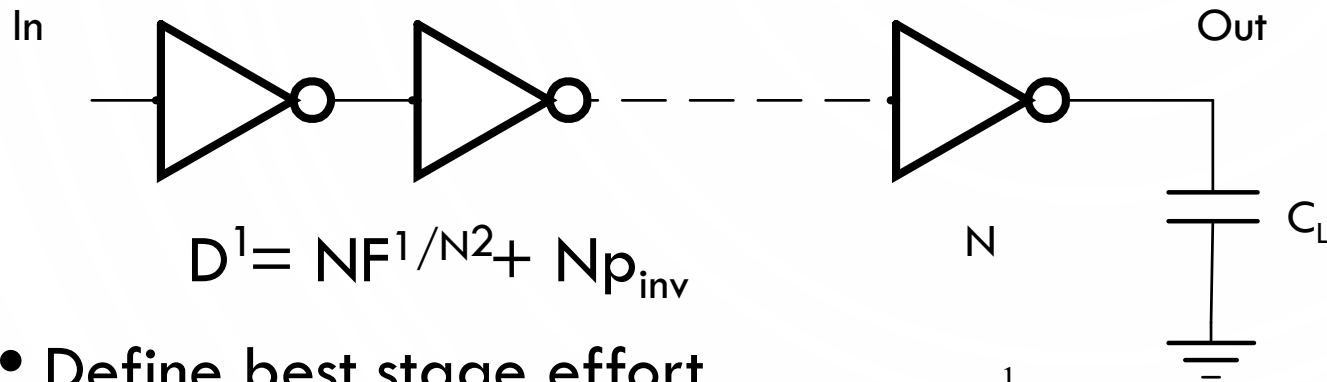
- How many stages should a path use?
 - Minimizing number of stages is not always fastest
- Example: drive 64-bit datapath with unit inverter

$$D = NF^{1/N} + P$$
$$= N(64)^{1/N} + N$$



Best Stage Effort

- How many stages should a path use?
 - To drive given capacitance



- Define best stage effort
- Neglecting parasitics ($p_{inv} = 0$), we find $\rho = e = 2.718$
- For $p_{inv} = 1$, solve numerically for $\rho = 3.59$
- Choose 4 – less stages, less energy
- Extends to any logic path with $h = 4$

Logical Efforts Method

- 1) Compute path effort
- 2) Estimate best number of stages
- 3) Sketch path with N stages
- 4) Estimate least delay
- 5) Determine best stage effort
- 6) Find gate sizes

$$H = GBF$$

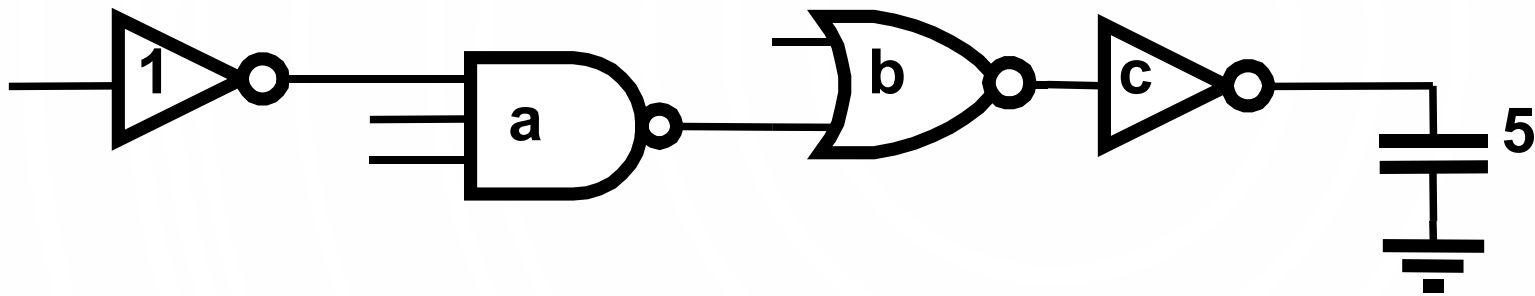
$$N = \log_4 H$$

$$D = NH^{\frac{1}{N}} + P$$

$$\hat{h} = H^{\frac{1}{N}}$$

$$C_{in_i} = \frac{g_i C_{out_i}}{\hat{h}}$$

Example: Optimize Delay



$$g = 1$$
$$f = a$$

$$g = 4/2$$
$$f = b/a$$

$$g = 3/2$$
$$f = c/b$$

$$g = 1$$
$$f = 5/c$$

Effective fanout, $F =$

$G =$

$H =$

$h =$

$a =$

$b =$

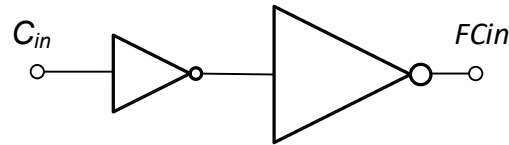
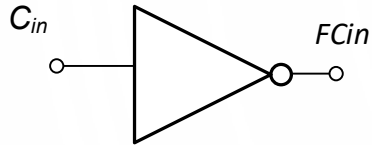
Example: What is the fastest NAND8?

Administrivia

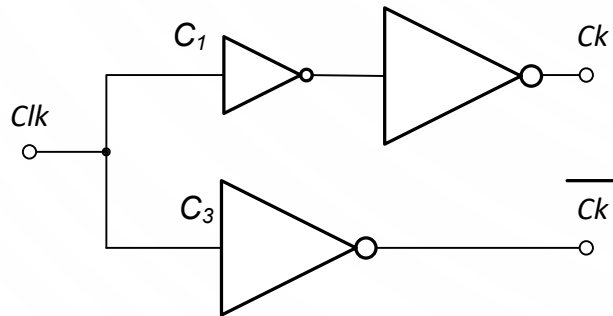
- Homework 6 due this week
- Projects (ASIC and FPGA) start this week

Logical Effort Design Examples

- For which F should we buffer?



- Sizing inverter fork

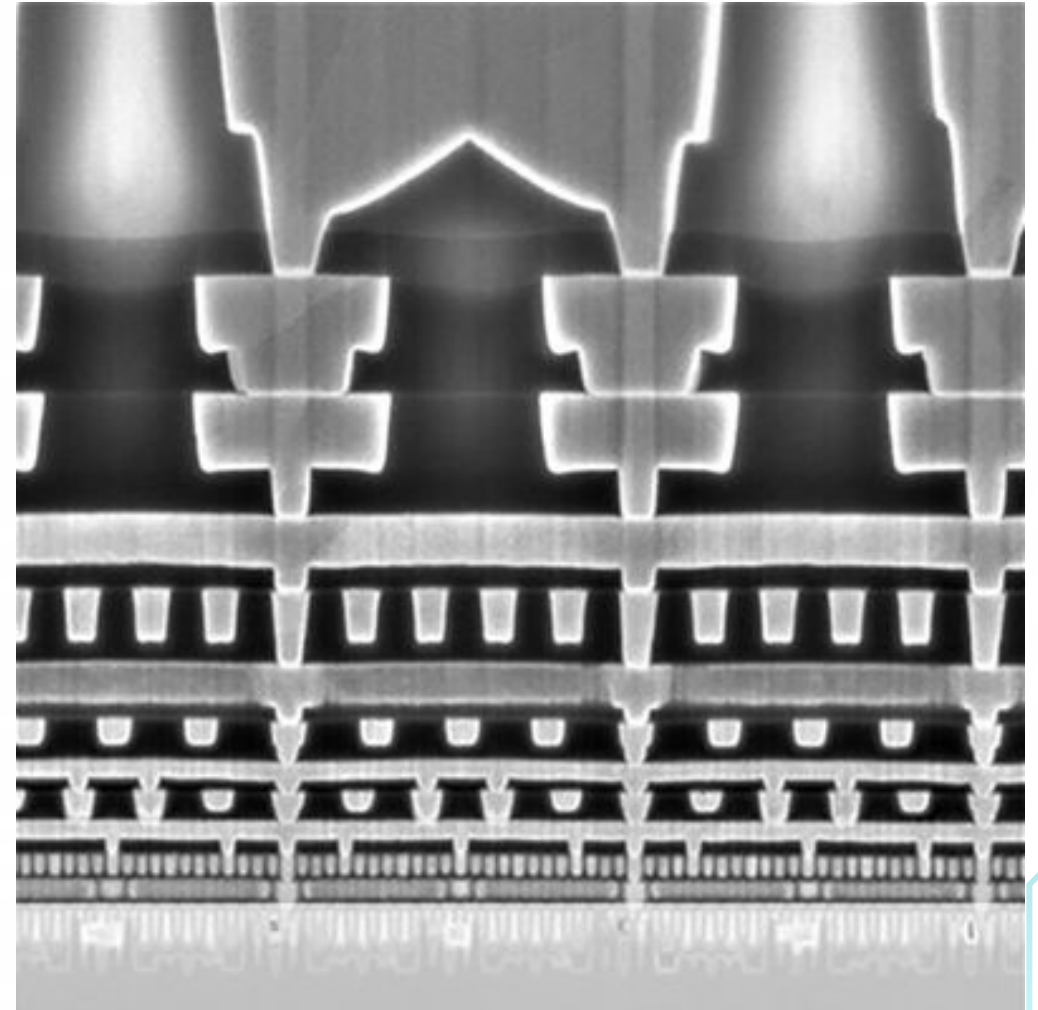




Wires

A modern technology is mostly wires

- Transistors are little things under the wires
- Many layers of wires
- Wires are as important as transistors
 - Speed and power



Wire Resistance

- $\rho = \text{resistivity } (\Omega \cdot \text{m})$

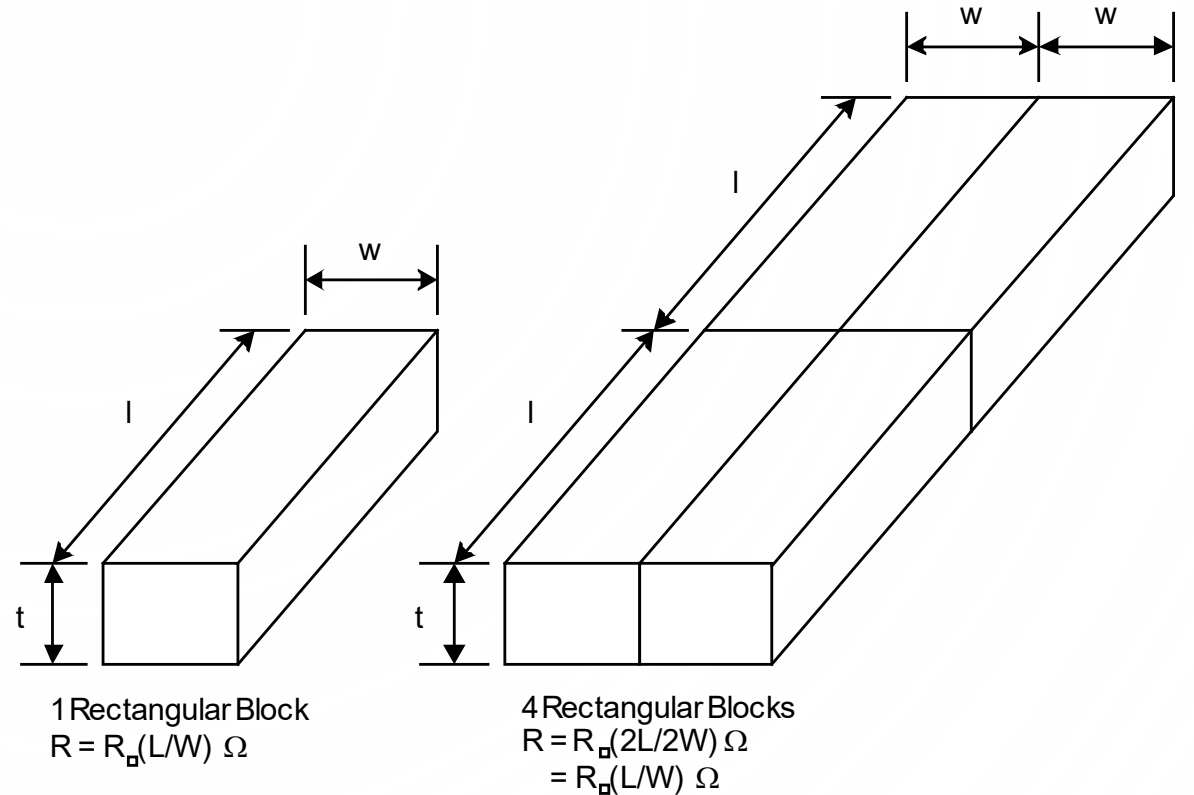
$$R = \frac{\rho}{t} \frac{l}{w} = R_{\square} \frac{l}{w}$$

- $R_{\square} = \text{sheet resistance } (\Omega/\square)$

- \square is a dimensionless unit(!)

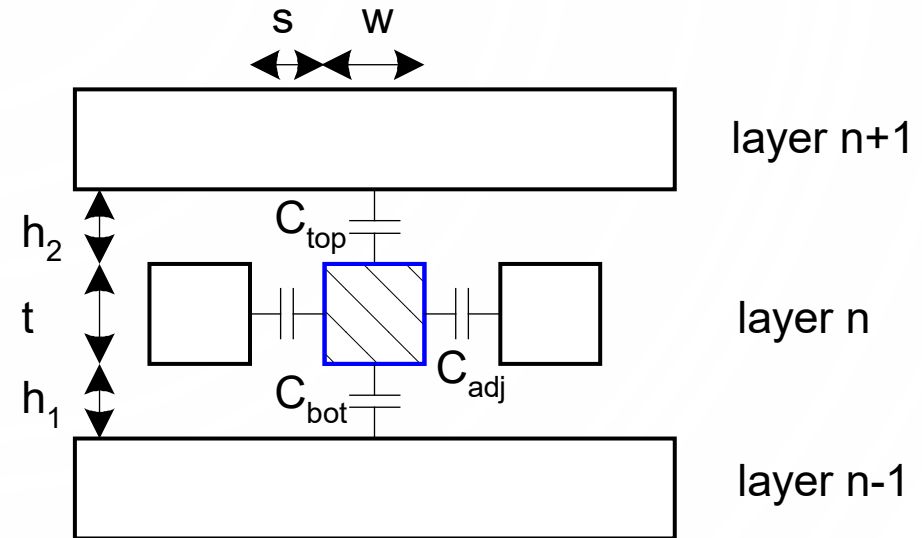
- Count number of squares

- $R = R_{\square} * (\# \text{ of squares})$



Wire Capacitance

- Wire has capacitance per unit length
 - To neighbors
 - To layers above and below
- $C_{\text{total}} = C_{\text{top}} + C_{\text{bot}} + 2C_{\text{adj}}$

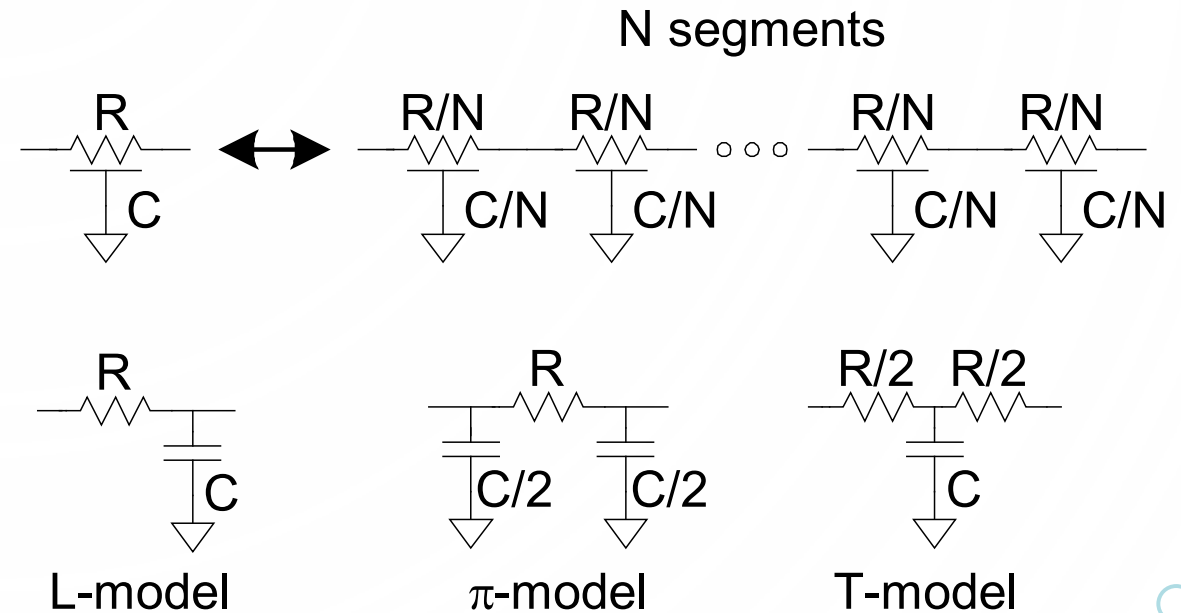




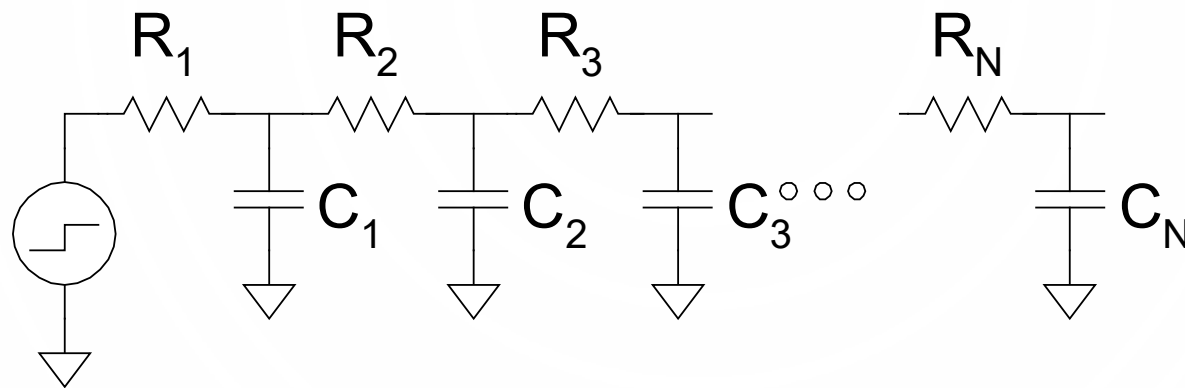
Wire Delay

Wire RC Model

- Wires are a distributed system
 - Approximate with lumped element models
- 3-segment pi-model is accurate to 3% in simulation

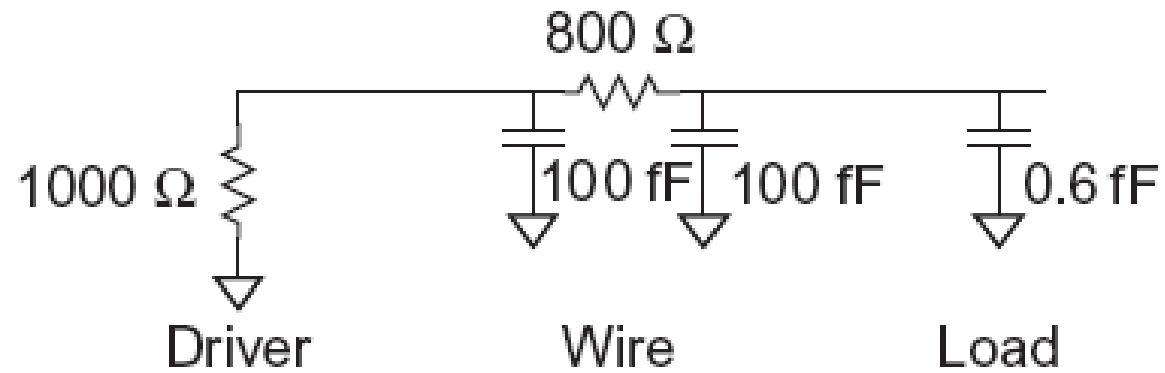


Elmore Delay for RC Tree

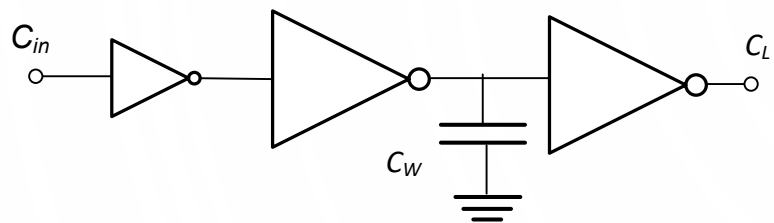


$$t_{pd} \approx \sum_{\text{nodes } i} R_{i\text{-to-source}} C_i$$
$$= R_1 C_1 + (R_1 + R_2) C_2 + \dots + (R_1 + R_2 + \dots + R_N) C_N$$

Example: RC Delay with Wire and Gate



Logical Effort with Wires



Summary

- Two delay components in logical effort:
 - Parasitic delay (p)
 - Effort delay (F)
 - Logical effort (g): intrinsic complexity of the gate
 - Electrical effort (h): load capacitance dependent
- To minimize the delay all stages should have the same effort (h)
- Ideal effort is 4
- Wires are modelled as RC
 - Most commonly just C for hand analysis