

Logic Design

# **Prof. Dejan Marković** ee216a@gmail.com

#### Agenda

- Logic design concepts
  - Fast vs. slow inputs
  - Transmission gates
- More delay models
  - Simulation-based
  - Library models
- Optimal gate sizing
  - Min delay
  - Delay > min delay

#### **Faster Inputs** 越东方 ontput A3 友的路,

Faster

input

both input

**in,** –|| M,

の先いり、可以ラクソカ

of schurge to ground

So VX clischarge to Q

回汕2, 马属遗荡度

output copacitor

ふかちいし、他会生

 $in_1 \dashv$ 

21/2 PB B B K dischar 20 V & Frib, WAS

- in I first, then in, • Higher input of a stack is faster (in, is faster than  $in_1$ )
  - V<sub>x</sub> pre-discharged
- 協友pinz 发到, • Also body effect on M2 change YD, ME
  - $V_x > 0$  increases  $V_{T2}$  and  $\sqrt{2}$   $\sqrt{1}$ takes longer to discharge 要更大时间
- Use higher inputs for the LATE arrivals
  - Synthesis tools do this

#### **Mixing Static CMOS & Transmission Gates**

- Static CMOS gates drive a TG switch network
- **Output of network drives a static CMOS gate**
- Example: a 2:1 Mux



#### **Example 3.1:** Transmission Gate Delay

- Delay depends on the gate that drives the switch
- Example: path B selected, Out<sub>1</sub> pulling up

$$t_{p} = t_{p1} + t_{p2}$$



#### **Modeling the Delay**

• Focus on the 1<sup>st</sup> (compound) stage (2<sup>nd</sup> stage is a simple inverter)



Assume: Inverters •  $W_p = 3 \mu m$ •  $W_N = 1 \mu m$ Trans. Gates •  $W_p = 2 \mu m$ •  $W_N = 2 \mu m$ 

#### **Capacitance Components**



#### **Calculate the Capacitances**



#### Calculate RC Time Constant τ<sub>1</sub>



# Simulation-Based Delay Model

#### **Simulation-Based Parameter Extraction**



# **Extracting C**<sub>gate</sub> by Simulation



•  $t_{p1} = t_{p2} \rightarrow C_{gate}$  is equivalent cap of the green gate

• Limitations: the model depends on signal rise times, voltage, temperature, process parameter variations

## **Extracting C**<sub>parasitic</sub> by Simulation

- Delay of the unloaded gate
  - Beware of signal slopes
- 2 Another way: find C<sub>parasitic</sub> from intersection of delay(fanout) line at fanout = 0
- **3** Play with AS, PS, AD, PD model parameters
  - AS = AD = PS = PD = 0 [overlap + junction caps]
  - $C_{gate} \& C_{par} \text{ for } t_{pLH} \text{ and } t_{pHL} \text{ are different [WHY?]}$ 
    - For delay analysis, we use the average value

#### **Extracting R<sub>eff</sub> by Simulation**

- Use the delay formula to find R<sub>eff</sub>
  - R<sub>eff</sub> = t<sub>p,gate</sub> / (0.69 C<sub>gate</sub>)
  - Reminder: set AS, AD, PS, PD to 0!



- **2** Slope of delay(fanout) line  $\propto$  fanout  $R_{eff} \cdot C_{gate}$ 
  - R<sub>eff</sub> for transistor stacks could be similarly found
    - Make sure you factor in the parasitic cap

# The "Flow"

#### • Hand design

- Simple RC models provide intuition about circuits
- Tradeoff analysis
- Dominant effects
- Reasonable starting point in the design process
- Run simulations to refine the model
- **Simulation** (to confirm your intuition)
  - If the simulation results are way off, there is a bug (check schematics, simulation files, your models)
  - Don't forget the corners

#### **Process Variations**

- Not all devices (primarily transistors) are created equal
  - Even if they nominally have the same exact drawing
  - Two on the same die (wafer, lot) can differ
- Variations cause transistors to have different V<sub>T</sub> and μ, speeding up (F) or slowing down (S) the devices
  - Doping conc. (μ)
  - W/L



#### • t<sub>ox</sub>

#### **More Corners (PVT Variation)**



#### Typically, given corner parameters for devices

 Characterize effective parameters across corners



cache

core

120C

**70C** 

#### **Process Corners Combined**



Includes:

- Process
- Voltage
- Temp
- variations

**S:** low 
$$V_{DD}$$
, high T  $(t_{setup})$   
**F:** high  $V_{DD}$ , low T  $(t_{hold})$ 

#### **Delay and Transition Time**

- Gate delay is approximated as RC elements
- The transition time is proportional to the same RC
  - Using ideal RC, the 10-90% is roughly 2.2RC
- For a gate (inverter) driving another gate,
  - Transition time is roughly 2t<sub>DELAY</sub>
  - Valid only for RC network
    - [What would happen in an inductive network?]
- Useful for modeling noise coupling
  - Aggressor transition time determines injected noise

#### **Multi-Stage Gate Delay**

• Static CMOS: total delay = sum of stage delays

$$t_{total} = \sum_{n} t_{n}$$

- The R of the gate
- The C of the <u>subsequent</u> gate(s)



#### Fan-In and Fan-Out

- Fan-In: the number of inputs (logic gale)
  - An indication of the input load that the gate presents to a predecessor gate
    - Because the series stack is roughly the number of inputs
    - Later we will use Logical Effort to embed this concept
- Fan-Out: effective output load (relative to Inv)
  - An indication of the loading (gate type dependent)
  - Useful to normalize the loading to C<sub>gate</sub> of an inverter with equal drive strength as the gate
    - FO =  $C_{LOAD}/C_{INV}$ where  $C_{INV} = C_0'(W_P + W_N)$  and  $R_{INV} = R_{PULL_UP}$

#### **Example 3.2:** Fanout Calculation

# NAND gate driving 5 equal NAND gates • NAND gate: $W_p = 5$ , $W_n = 5$ For Son Total input width = $10 \text{ MeV}_{10}$ • Total load gate width = $5 \cdot 10 = 50$ • Equivalent Inv: $W_p = 5$ , $W_n = 2.5$ • Equivalent (in terms of PU/and/PD strength)

- Total gate width = 7.5
- Fanout = 50/7.5 = 6.6

# **Another Metric: FO4 Inverter Delay**

• Measures quality of design independent of technology



Cadence 90nm technology: FO4 = 33ps Ring Osc Stage = 13ps

Mini assignment:

Simulate FO4 delay in 32nm (slide 2.34)

#### **Model Used in Cell Libraries**

#### Propagate two quantities:

- Delay
- Signal slope

#### Modeled as linear or table lookups

#### **Gate Delay Estimation**

- Depends on transistor sizes, input signal slew, output capacitive load and PVT conditions
- For each PVT, timing library is provided by vendor
  - Delay, output slew, dynamic and static power are characterized for different input slew and output loads
- Characterization values are stored as 2-D look-up tables
  - Output delay = f (input slew, output load capacitance)

#### **Gate Delay Estimation**



Interpolation with 4 points

Courtesy: Synopsys

#### **Design Problem: Interconnect Delay**

- Timing assumption during pre-layout synthesis widely differs from the post-layout reality
- This happens because the interconnect delay dominates the overall delay in scaled technologies
- As a result, timing closure becomes a challenge

#### **CAD Problem: Interconnect Delay**

- New industry trends = new IC design flows
- The major contributing factors:
  - Interconnect delay and area
  - The number of metal layers
  - Planning requirements
- Need to consider interconnect delay early in design

fan-Out

- A few techniques for reducing interconnect delay:
  - Chip-level signal planning
  - Over-the-Cell routing

## Wire Load Models (WLMs)

- In the absence of physical design information DC uses statistically generated WLMs to estimate wire lengths
  - WLMs provide a fanout vs. length relationship
- Interconnect delay estimation:
  - By knowing fanout, estimate the wire lengths
  - Next, given unit-length R and C and the estimated length, estimate R and C to give an estimated delays
- Note: WLMs are area dependent
  - Unit-length R and C increase with area

#### **Example WLM**

wire\_load('30x30') {
 capacitance : 3.0 ; /\* C per unit length \*/
 resistance : 30.0 ; /\* R per unit length \*/
 area : 1.5 ; /\* area per unit length \*/
 slope : 1.5 ; /\* extrapolation slope \*/
 fanout\_length(1,1) ; /\* fanout/length pairs \*/
 fanout\_length(2,2.2);
 fanout\_length(3,3.3);
 fanout\_length(4,4.4); }

# Model from a library compiler (LC) source file:



Courtesy: Synopsys

#### **Hierarchical Wire Load Models**

DC supports **3 modes** for nets that cross hier. boundaries



TOP

30x30

TOP

20x20

**30x30** 

 $10 \times 10$ 

# Speed Optimization via Gate Sizing

#### **Speed Optimization via Gate Sizing**

#### • Gate sizing basics

- P:N ratio
- Complex gates
- Velocity saturation
- Tapering
- Developing intuition
  - Number of stages vs. fanout
  - Popular inverter chain example

#### **Basic Gate Sizing Relationships**

- Rise and fall delays are determined by the pull-up and pull-down "strength"
  - Besides the W/L, strength depends on  $\mu$ ,  $V_T$
  - PMOS is weaker because of lower μ<sub>p</sub>
     Larger P network than N network
- Increasing size of gate can reduce delay
  - $R_{on} \propto 1/W$
  - BUT it can slow down the gate driving it
    - C  $\propto$  W. So be careful!

#### P:N Ratio for "Equal" Rise and Fall Delay

- Good to have  $t_{pHL} \approx t_{pLH}$ 
  - Don't need to worry about a worst-case sequence
  - Size P's to compensate for mobility
    - $C_{OX}$ ,  $V_T$ , L are roughly the same

$$R_{on} \propto \frac{1}{I} \propto \frac{1}{\mu W}$$

n

• Make  $R_p = R_n$ •  $R_n/R_p = 1$ 

$$(\mu_p W_p/\mu_n W_n = \beta/k)$$

$$\frac{1}{I} \propto \frac{1}{\mu W}$$
$$t_{pLH} = t_{pHL}$$
$$\frac{W_p}{W} = \beta = k = \frac{\mu_n}{\mu}$$

kenne = f=tr

# **Complex Gate Sizing**

- N-stack series devices need N times lower resistance
   N×Width
- Make worst-case strength of each path equal
  - Multi-input transition can result in stronger network
- Long series stacking is VERY bad



# **Accounting for Velocity Saturation**

- Series stacking is actually less velocity saturated
  - If we use R<sub>no\_stack</sub> = (4/3)R<sub>stack</sub>
  - Adjust the size of non-stacked devices to account for velocity saturation



#### **P:N Ratio for Minimum Delay?**



#### **P:N Ratio for Equal/Minimum Delay**

• Delay of 2 inverters:



$$t_{pLH} = t_{pHL}$$
$$\frac{W_p}{W_n} = \beta = k = \frac{\mu_n}{\mu_p}$$

NMOS: more drive for a given size, so it is better to use more NMOS

Min 
$$(t_{pLH} + t_{pHL})/2$$
  
 $\frac{W_p}{W_n} = \beta = \sqrt{k} = \sqrt{\frac{\mu_n}{\mu_p}}$ 

#### FO4 Inverter Delay vs. P:N Ratio β

- Optimal β = sqrt(μ) for minimum delay
  - Curve is relatively flat so not a strong delay tradeoff



#### An Idea: Tapering

- C closer to the v-source: less effect on delay
  - N = 2:  $t_p = R_1(C_1) + (R_1 + R_2)(C_2)$ 
    - C<sub>1</sub> has less effect on delay than C<sub>2</sub>
- So taper stacked devices for speedup
  - Make the bottom ones bigger
    - R<sub>1</sub> (many occurrences) has less resistance
    - C<sub>3</sub> (multiplying larger R) has smaller capacitance



In reality, tapering doesn't win as much because layout is less compact when stacking unequal sized transistors (causing more C)

#### **Example 3.3:** Gate Design for Min Delay

- Which scenario (A, B) has faster delay?
  - Drive the same output load



Let's analyze building blocks: NAND, NOR, INV

#### **A1:** Delay of N-input NAND



#### **A1:** Delay of N-input NAND

1

$$t_{p} = \sum_{i=1}^{N-1} iR_{i}C_{i} + R_{tot}(C_{N} + NC_{0}\frac{\beta W_{n}}{f} + C_{gate})$$

$$P_{tot} = \frac{P_{0}}{W_{n}f} \quad P_{i} = \frac{P_{0}}{M_{N}f} \int_{W_{n}}^{N-1} iR_{0}C_{0} + R_{0}(NC_{0} + NC_{0}\beta + C_{gate}\frac{f}{W_{n}})$$

$$n \text{ transistor elements}$$

$$NMOS \quad PMOS \quad Output$$

$$t_{p} = R_{0}C_{0}\left(\frac{N(N-1)}{2} + N\right) + R_{0}C_{0}\left(\beta N + \frac{C_{gate}}{C_{0}}\frac{f}{W_{n}}\right)$$

#### **A2:** Delay of Inverter



#### **Inverter:**

• 
$$R_{INV} = R_0 / W_n$$

- $C_{L,INV} = C_{par} + C_{gate} = C_0(\underbrace{W_n(1+\beta)}_{n} + \underbrace{f W_n(1+\beta)}_{n})$   $t_{INV} = R_0C_0(1+\beta)(1+f)$

$$C_{\text{gate,INV}} = C_0(W_n(1+\beta))$$

#### **B2:** Delay of NOR2



- $R_{NOR} = R_0 / W_n$
- $C_{L,NOR} = C_{par} + C_{gate} = C_0(W_n(2+2\beta) + f W_n(1+\beta))$
- $t_{NOR} = R_0 C_0 (1+\beta)(2+f)$

$$C_{gate,NOR} = C_0(W_n(1+2\beta))$$

-0

#### Delay Comparison ( $\beta = 2$ )

#### A: N-in NAND + INV

$$t_{pA} = R_0 C_0 \left( \frac{N(N-1)}{2} + 3N + 3f \right) + R_0 C_0 (3f+3)$$

#### **B:** N/2-in NAND + NOR

$$t_{pB} = R_0 C_0 \left( \frac{N\left(\frac{N}{2} - 1\right)}{4} + \frac{3N}{2} + 5f \right) + R_0 C_0 (3f + 6)$$

Mini assignment:

Assume f = k N = ? @ crossover

#### **Sweep N:** Delay Comparison ( $\beta = 2$ )

N	t <sub>pA</sub>	t <sub>pB</sub>	<b>f</b> =3	f = 4	f = 5
4	21 + 6f	13 + 8f	39   37	45   45	51   53
6	36 + 6f	21 + 8f	54   42	60   53	66   61
8	55 + 6f	30 + 8f	73   54	79   62	85   70
			A   B	A   B	A   B

## Min t<sub>p</sub>: don't build gates with fan-in > 4

#### **Transmission Gate Sizing**

Try to make a TG with 
$$R_{PD} = R_{PU}$$

#### P:N ratio of k is not good for delay

- NMOS still has some pull-up strength (even if not all the way to V<sub>DD</sub>)
  - No need to have wide PMOS
- PMOS has some pull-down (but very weak)
  - PD slightly faster (NMOS is stronger)

#### **TG Sizing: Some Common Numbers**

**2x penalty, weak transition**  

$$R_{N,DN} = R_0 k\Omega - \mu m | R_{N,UP} = 2R_0 k\Omega - \mu m$$
  
 $R_{P,UP} = 2.5R_0 k\Omega - \mu m | R_{P,DN} = 5R_0 k\Omega - \mu m$ 

• Let's try 
$$W_p = W_n$$
 Slightly

- Parallel Up,  $R_{TGUP} = R_{N,UP} ||R_{P,UP} = 1.1R_0 \neq faster$
- Parallel Down,  $R_{TGDN} = R_{N,DN} ||R_{P,DN} = 0.83R_0$ for transmission gate
- So, using  $W_p/W_n = 1$  is fairly reasonable
- Actual size may depend on the process technology

## **Delay Analysis (So Far): Summary**

- Device R and C determine circuit performance
- Elmore Delay approximation: initial insight into design
  - Step response, does not account for signal slopes
  - Need slope correction [see slide 2.63]
- The sizing of the transistors (a first glimpse)
  - Determines V<sub>M</sub>
  - Determines R<sub>DRV</sub> as well C<sub>Load</sub> it presents to the preceding gate which effects the delay
- Large fan-in gates: large self-loading and loading to the preceding gate
  - Split into 2 gates when fan-in > 4

#### **Speed Optimization via Gate Sizing**

- Gate sizing basics
  - P:N ratio
  - Complex gates
  - Velocity saturation
  - Tapering

#### • Developing intuition

- Number of stages vs. fanout
- Popular inverter chain example

#### **Problem Statement**



- Given:
  - Arbitrary logic function
  - Gate-level implementation
- How do we decide the relative size of each gate?

#### **Simplified Problem: Buffering**



- <u>Assume:</u>  $\beta = k$ 
  - R<sub>0</sub> = Pull down for NMOS with size W<sub>0</sub> or PMOS with size βW<sub>0</sub>
  - C<sub>0</sub> = Gate capacitance of N+PMOS of size W<sub>0</sub>, βW<sub>0</sub>
     Ignore Source/Drain & Wire Capacitance for now
     τ<sub>0</sub> = R<sub>0</sub>C<sub>0</sub>
- <u>Goal</u>: find  $a_1, a_2, ..., a_{N-1}$  to minimize delay

#### **Delay Calculation**



#### **Optimal Fanout for Given N**

- Fanout calculation:
  - Stage 1 =  $a_1$ , Stage 2 =  $a_2/a_1$
- Assuming that the fanout of each stage is equal, a<sub>0</sub>

$$a_{1} = a_{0}, a_{2} = a_{0}^{2}, a_{3} = a_{0}^{3}$$
$$\int_{C_{out}}^{C_{out}} = C_{0} a_{0}^{N}$$

- Total Delay = Sum (Delay of stage 1:N)
  - Delay =  $\tau_0 Na_0$

• Since 
$$C_{in} = C$$
  
•  $C_{out} / C_{in} = a_0^N$ 

$$a_0 = \left(\frac{C_{out}}{C_{in}}\right)^{\frac{1}{N}}$$

#4 stoje-

#### **Optimum Number of Stages**

For an arbitrary N: 
$$N = \frac{ln\left(\frac{C_{out}}{C_{in}}\right)}{ln(a_0)}$$
$$T_{o} a_{o} N$$
  
Min delay: 
$$\frac{\partial Delay}{\partial a_0} = \frac{\partial \left[\tau_0 a_0 \frac{ln\left(\frac{C_{out}}{C_{in}}\right)}{ln(a_0)}\right]}{\partial a_0} = 0$$
$$\frac{\partial Delay}{\partial N} \frac{1}{\tau_0} = \left(\frac{C_{out}}{C_{in}}\right)^{\frac{1}{N}} - \frac{\left(\frac{C_{out}}{C_{in}}\right)^{\frac{1}{N}} ln\left(\frac{C_{out}}{C_{in}}\right)}{N} = 0$$
$$\implies a_0 = e$$

#### **Optimum Number of Stages**



Set number of stages to reach optimal fanout

#### Intuition:

What if we increase the size of a stage by  $(1+\Delta)$ ?

- $R_{DRV} \propto 1/(1+\Delta)$ •  $C_{Load}$  (previous stage)  $\propto (1+\Delta)$
- **Delay** is summed and would increase

#### (Constant FO) Mathematically

- Delay =  $\tau_0(a_1 + a_2/a_1 + a_3/a_2 + a_4/a_3 + a_5/a_4 + ...)$
- $dDelay/da_1 = 0$ 
  - dDelay/da<sub>1</sub> =  $\tau_0(1 a_2/a_1^2)$

• So 
$$a_1^2 = a_2$$

- $dDelay/da_2 = 0$ 
  - dDelay/da<sub>2</sub> =  $\tau_0(1/a_1 a_3/a_2^2)$
  - So  $a_2^2 = a_1 a_3$ , thus  $a_1^3 = a_3$

# ⇒ Min delay: equal fanout

#### **Optimal Buffering with Self-Loading**

#### • Intuition: without self-loading

- Delay decreases proportional to the decrease in N
- But increasing fanout increases delay proportionally
- The two are equal @ the optimal N and fanout
- Intuition: with self-loading
  - Increasing FO no longer increases delay proportionally
    - Delay =  $R_0(FO \cdot C_0 + C_{par})$
  - New N would be less and FO is bigger

(Buffering with Self-Loading) Mathematically

All equations remain the same except Delay

$$\frac{\partial Delay}{\partial a_0} = \frac{\partial \left[ \tau_0 \left( a_0 + \frac{C_{par}}{C_0} \right) \frac{ln \left( \frac{C_{out}}{C_{in}} \right)}{ln(a_0)} \right]}{\partial a_0} = 0$$

$$integrad{bmatrix} a_0 = e^{1 + \frac{p}{a_0}}$$

P

## **Optimal Buffering as** *fn* (Self-Loading)





# Optimal FO for t<sub>p</sub> > t<sub>p,min</sub>

#### **Buffer Optimization for Energy-Delay**

- Minimizing for Energy is trivial: min gate size
- Instead minimize Energy-Delay product



#### **Issue with Optimal Energy-Delay**

- Constant fanout is not a good assumption 又常學最
- Intuition: • Intuition: • Size of the final stage maters the most (use max fanout)
  - Reduce fanout of prior stages to compensate...
- Example: C<sub>in</sub> = 1, C<sub>out</sub> = 1000, 4 gate stages

Design	Gate 1	Gate 2	Gate 3	Gate 4	EDP
C <sub>in</sub>   FO	1   5.6	5.6   5.6	31.6   5.6	177.8   5.6	32200
C <sub>in</sub>   FO	1   4.8	4.8   4.9	23.1   5.4	124.5   <mark>8.0</mark>	31100

## **Tapered fanout** reduces EDP

#### What if Fanout is Low?

- **Example:** large N (each stage drives small fanout)
  - Delay is logic limited, so reduce N
  - Balance Fanout so that they are equal
- Try to adjust N such that each stage has FO ~ 4
  - More complex logic for smaller N
- OK, but not very systematic...
  - Logical effort (next lecture) to formalize sizing

#### **Energy-Delay Optimization (Lecture 14)**

**Variables:** Gate size, Supply Voltage, Threshold Voltage



#### Summary

- Delay of a logic network depends significantly on the relative size of logic gates
  - Not transistors within a gate
- Inverter buffering is a simple example of the analysis
  - The analysis leads to FO ~ 4 as being optimal fanout for driving larger capacitive loads
  - The number of stages is optimized when FO ~ 4
  - For delays longer than minimum, tapered FO works best for minimizing power (more on this in Lec 14)