EE3230 Lecture 5: Circuit Characterization and Performance Estimation II

Ping-Hsuan Hsieh (謝秉璇)

Delta Building R908 EXT 42590 phsieh@ee.nthu.edu.tw

Outline

- Delay estimation
- Logical effort and transistor sizing
- Power dissipation (45
- Interconnect Chb
- Wire engineering
- Design margin
- Reliability
- Scaling



Power and Energy

- V Power is drawn from a voltage source attached to the V_{DD} pin(s) of a chip
 i(t)
- Instantaneous Power: P(t) = v(t) i(t) i(t) vin - voo vin - voo $V \cdot Energy:$ $E = \begin{bmatrix} P(t) at = voo \int_{0}^{1} i(t) at \end{bmatrix}$
 - Average Power: $Pavg = \frac{E}{T} = \frac{1}{T} Voo \int_{0}^{T} ict y dt = Voo T$

Static and Dynamic Power Dissipation



Dynamic Power (I)

- Dynamic power is required to charge and discharge load capacitances when transistors switch
- One cycle involves a rising and falling output.
- On rising output, charge $Q = CV_{DD}$ is required
- On falling output, charge is dumped to GND
 - This repeats f_{sw} times per second





Activity Factor

- Suppose the system clock frequency = f
- Let $f_{sw} = \alpha f$, where α = activity factor
- **e** 26 H² If the signal is a clock, $\alpha = 1$
 - If the signal switches once per cycle, $\alpha = \frac{1}{2}$
 - **Dynamic gates:** Switch either 0 or 1 times per cycle, $a = \frac{1}{2}$
 - φ- <u>Static gates</u>: Depends on design, typically α = 0.1 (M05 random signal, switch with a probability of 0.5 → α=0.25
 - Dynamic power:

Short-Circuit Current

- When transistors switch, both NMOS and PMOS networks may be momentarily ON at once
- Leads to a blip of short-circuit current
- < 10% of dynamic power if rise/fall times are comparable for input and output (well-controlled)



Example (1)

- 200M transistor chip System
- 20M logic transistors
 - Average width: 12λ
 - 180M memory transistors
 - Average width: 4 λ

min channel length

- 1.2-V 100-nm process (
$$\lambda = 0.5^*$$
 feature size = 50nm)

$$- C_g = 2 fF/\mu m$$

p=a. C Voo. fsw

Dynamic Power Consumption (I)

- Static CMOS logic gates: <u>activity factor = 0.1</u>
- Memory arrays: <u>activity factor = 0.05</u> (many banks and partially activated at a time!)
- Estimate dynamic power consumption per MHz.
 (Neglect wire capacitance and short-circuit current)

$$C_{10}q_{12} = (20.10^{6})(12.50nm)(\frac{2-fF}{\mu m}) = 24 nF$$

$$C_{mem} = (180.10^{6})(4.50nm)(\frac{2 fF}{\mu m}) = \eta 2 nF$$

$$\overline{P} = 0.1(24nF)(\frac{2}{10}+0.05)(12nF)(\frac{2}{10}+0.05) = 8.6N$$

$$10q_{12}$$

Static Power Consumption

- Static power is consumed even when chip is quiescent.
- A Ratioed circuits burn power in fight with ON transistors
- Leakage draws power from nominally OFF devices

Ratioed Example

- The chip contains a 32 word x 48 bit ROM
 - Uses 1:32 pseudo-nMOS decoder and bit-line pull-ups
 - In average, one wordline and 24 bitlines are high
- Find static power drawn by the ROM





NAND3

- Subthreshold leakage I_{SUB}
- · 000 Stacking
- 110
- 011
- 111

ABC	ISUB	IGATE	Vx	Vz
000	0.4		Stack	stock
001				
010	0	1.3		
011	3.8		Voy-Vth	Voo-Vth
100				
101				
110	5.6		0	0
111	28		٥	0

NAND3

- Gate leakage I_{GATE}
- 000
- 001
- 010
- 101

ABC	ISUB	IGATE	Vx	Vz
000		0		
001		0		Voo-Vth
010	0	1,3	intermediate	intermediate
011				
100				
101		6.3	D	Voo-Vth
110				
111				

Leakage Example (I)

- The process has two threshold voltages and two oxide thicknesses.
- Subthreshold leakage:
 - 20 nA/ μ m for low V_{th} devices high-performence
 - 0.02 nA/ μ m for high V_{th} devices
- Gate leakage:
 - 3 nA/ μ m for thin oxide
 - 0.002 nA/ μ m for thick oxide
- Memories use low-leakage transistors everywhere

low-power

• Gates use low-leakage transistors on 80% of logic

Leakage Example (II)

- Estimate static power: - High leakage: (20.10) (20%) (12 50mm)=2.4×10 mm - Low leakage: $\binom{6}{20.10} (\frac{80.10}{6}) (\frac{12.50}{12.50}) + \binom{6}{180.10} (\frac{80.10}{6}) (\frac{4.50}{12.50}) = 45.6.10 \mu m$ Isranz= 2.4.10 (2011A+31A)/2 6 + 45.6.10 (0.020A+0.0020A)/2 = 27.6+0.5016=28.1016MA
- Withow leakage devices, P_{static} = 552 mW!!

Low Power Design



stacked devices, body bias, low temperature

Reduce Static Power

• Leakage stack effect



• MTCMOS: multiple threshold CMOS



• Body bias



Outline

- Delay estimation
- Logical effort and transistor sizing
- Power dissipation
- Interconnect
- Wire engineering
- Design margin
- Reliability
- Scaling

Interconnect

- Chips are mostly made of wires called *interconnect*
 - In stick diagrams, wires determine size
 - Transistors are little things under the wires
 - Many layers of wires
- Wires are as important as transistors
 - f Speed FC
 - + Power C
 - Noise (compling)
- Alternating layers run orthogonally

Wire Geometry

- Pitch = w + s
- Aspect ratio: AR = t/w
 - Old processes had AR << 1
 - Modern processes have $AR \approx \frac{2}{2}$
 - Pack in many skinny wires



Layer Stack

- AMI 0.6-µm process has 3 metal layers
- Modern processes use 6-10+ metal layers CMP
- Example: Intel 180 nm process



Wire Resistance



Choice of Metals

- Until 180 nm, most wires were aluminum
- Modern processes often use copper
 - Cu atoms diffuse into silicon and damage FETs
 - Must be surrounded by a diffusion barrier

Metal	Bulk resistivity (μΩ*cm)	
Silver (Ag)	1.6	
Copper (Cu)	1.7	
Gold (Au)	2.2	
Aluminum (Al)	2.8	
Tungsten (W)	5.3	
Molybdenum (Mo)	5.3	

Sheet Resistance

• Typical sheet resistances in 180-nm process

	Layer	Sheet Resistance (Ω/\Box)	
×	Diffusion (silicided)	3-10	Ţ
x	Diffusion (no silicide)	50-200	
v	Polysilicon (silicided)	<u>3</u> -10	46
Ī	Polysilicon (no silicide)	50-400 🧲	
Ī	Metal1	0.08	1629日1天中
Ī	Metal2	0.05	
Ī	Metal3	0.05	
Ī	Metal4	0.03	
	Metal5	0.02	
Ī	Metal6	0.02	

Contact Resistance

- Contacts and vias also have 2-20 Ω
- Use many contacts for lower R
 - Many small contacts for current crowding around periphery



Wire Capacitance

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



Capacitance Trend

- Parallel plate equation: $C = \epsilon A/d$
 - Wires are not parallel plates, but obey trends
 - Increasing area (W, t) increases capacitance
 - Increasing distance (s, h) decreases capacitance
- Dielectric constant
 - $-\epsilon = k\epsilon_0$
 - $-\epsilon_0 = 8.85 \text{ x } 10^{-14} \text{ F/cm}$
 - k = 3.9 for SiO₂
- Processes are starting to use low-k dielectrics
 - $k \approx 3$ (or less) as dielectrics use air pockets

M2 Capacitance Data

- Typical wires have ~ 0.2 $fF/\mu m$
 - Compare to 2 $fF/\mu m$ for gate capacitance



Diffusion and Polysilicon

- Diffusion capacitance is very high (about 2 $fF/\mu m$)
 - Comparable to gate capacitance
 - Diffusion also has high resistance
 - Avoid using diffusion *runners* for wires!
- Polysilicon has lower C but high R
 - Use for transistor gates
 - Occasionally for very short wires between gates

Lumped Element Models

- Wires are a distributed system
 - Approximate with lumped element models



- 3-segment π -model is accurate to 3% in simulation
- L-model needs 100 segments for same accuracy!
- Use single segment π -model for Elmore delay

Example

- M2 wire in 180-nm process
 - <u>5-mm</u> long
- Construct a 3-segment π-model

$$\begin{array}{c} + R_{\Box} = 0.05 \ \Omega/\Box & \rightarrow R = \ \ensuremath{\neg} \$ \ \ensuremath{\neg} \$ \ \ensuremath{\neg} \$ \ \ensuremath{\neg} = \ \ensuremath{\bigcirc} 2 \ fF/\mu m & \rightarrow C = \ \ensuremath{1 \ P \ F} \end{array}$$





Wire RC Delay

 Estimate the delay of a 10x inverter driving a 2x inverter at the end of the 5-mm wire from previous example

F Effective R = 2.5 kΩ \neq μm for gates, C = 2 *f*F/μm

- Unit inverter: $4\lambda = 0.36 \ \mu m \ nMOS, 8\lambda = 0.72 \ \mu m \ pMOS$ $\int_{a}^{b} \frac{1}{2} \frac{1}{2} \frac{1}{2} \int_{a}^{b} \frac{1}{2} \frac{1}{2} \int_{a}^{b} \frac{1}{2} \frac{1}{2} \frac{1}{2} \int_{a}^{b} \frac{1}{2} \frac{1}{2} \frac{1}{2} \int_{a}^{b} \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} \int_{a}^{b} \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} \int_{a}^{b} \frac{1}{2} \frac{1$

Crosstalk compling noise

- Capacitor do not change voltage instantaneously
- A wire has high capacitance to its neighbor
 - When the neighbor (aggressor) switches from 1→0 or 0→1, the wire (victim) tends to switch as well
 - Called *capacitive coupling* or *crosstalk*
- Impacts
 - Cause noise on non-switching wires
 - Increase delay on switching wires

OVA = OVB

o VB

Crosstalk Delay

- Assume layers above and below in average are quiet
 - Second terminal of capacitor can be ignored
 - Modeled as $C_{gnd} = C_{top} + C_{bot}$
- Effective C_{adi} depends on behavior of neighbors
 - Miller Coupling Factor (MCF)



Crosstalk Noise (Floating Victims)

- Crosstalk causes noise on non-switching wires
- If victim is floating
 - modeled as capacitive voltage divider



Crosstalk Noise (Driven Victims)

- Usually victim is driven by a gate that fights noise
 - Noise depends on relative resistances
 - Assume victim driver in linear region and aggressor driver in saturation (considering inverter operation)
 - With equal sizes, $R_{aggressor} = 2.4 \times R_{victim}$ $\delta V_c = \delta V_a \frac{C_{aaj}}{C_{gaa} + C_{aaj}}$ $K = \frac{C_{aggressor}}{T_{victim}} V_{aggressor}$



Coupling Waveforms



Noise Implications

- So what if we have noise?
- If the noise is less than the noise margin, nothing happens
- Static CMOS logic will eventually settle to correct output even if disturbed by large noise spikes
 - But glitches cause extra delay
 - Also cause extra power from false transitions
 - Dynamic logic never recovers from glitches
 - Memories and other sensitive circuits also can produce wrong outputs

Outline

- Delay estimation
- Logical effort and transistor sizing
- Power dissipation
- Interconnect
- Wire engineering <---
- Design margin
- Reliability
- Scaling

Wire Engineering

- Goal: to achieve delay, area, and power goals with acceptable noise
- Degrees of freedom:

Wire Engineering

- Goal: to achieve delay, area, and power goals with acceptable noise
- Degrees of freedom:



Wire Engineering

- Goal: to achieve delay, area, and power goals with acceptable noise
- Degrees of freedom:



Staggered repeater



Repeaters

- R and C are proportional to / (length)
- RC delay is proportional to P
 - Unacceptably long delays for long wires
- Break long wires into N shorter segments
 - Drive each one with an inverter or buffer



Repeater Design

res/length



- How many repeaters should we use?
- How large should each one be? (z)
 - Equivalent Circuit
 - Wire length /
- cop/length ./ • Wire Capacitance $\xi_{w} \cdot I$, Resistance $k_{w} \cdot I$
 - Inverter size W (nMOS = W, pMOS = 2W)
 - <u>Gate</u> capacitance C·W, parasitic capacitance $C \cdot p_{inv} \cdot W$ (drain)
 - Resistance R/W

Repeater Design

Vin

- How many repeaters should we use?
- How large should each one be?
- Equivalent Circuit
 - Wire length I/N
 - Wire Capacitance $C_w \cdot I/N$, Resistance $R_w \cdot I/N$
 - Inverter size W (nMOS = W, pMOS = 2W)
 - Gate capacitance C·W, parasitic capacitance C· p_{inv} ·W
 - Resistance R/W

$$\frac{F_{W}}{N} = \frac{F_{W}}{1} \frac{V}{1} \frac{V}{2N} = \frac{V}{1} \frac{V}{1}$$

2W

Repeater Results

* Elmore Delay

$$t_{pd} = N \left(\frac{R}{W} \left(C_{w} \frac{l}{N} + CW(1+P_{inv}) \right) + R_{w} \frac{g}{N} \left(\frac{C_{w}l}{zN} + CW \right) \right)$$

 $= f(N, W)$
 $- Differentiate With W and N and find out optimum
 $\frac{l}{N} = \int \frac{2 RC(1+P_{inv})}{R_{w}C_{w}}$ with FOQ = SRC (implied Prov=1)
with Folded Transistor Priv=0.5
 $0.11 \int \frac{Foe}{R_{w}C_{w}}$
 $t_{pd} = (2 + J_{2}(1+P_{inv})) \int RCR_{w}C_{w} \leq 1.69 \int FOQ R_{w}(w)$$

$$W = \int \frac{R \cdot C_{w}}{R_{w}C}$$

$$\frac{E}{g} = C_{w} \left(1 + \int \frac{|T|P_{ihv}|}{2}\right) \cdot V_{ev}^{2}$$

$$\stackrel{\simeq}{=} 1.89 \cdot C_{w} \cdot V_{ev}^{2}$$

$$\stackrel{?}{=} 389\% \text{ power compared to no repeaters}$$

$$\Rightarrow ih addition to min (tpa), the T x min (EDP) energy - delay product$$

$$\rightarrow 0.8 \text{m} \cdot \pi De - 1 \text{E} = W = 11 \text{ \mum E3 inverten}$$

$$\frac{tpd/e}{2} = 49 \text{ ps/mm} \quad \overline{f}_{e} = 0.26 \text{ ps/mm}$$

Example

- 65-nm technology, middle routing layer with 2x width, spacing, and height.
 - Rw=200 ohm/mm, Cw= 0.2 pF/mm
 - FO4=15 ps

$$\frac{tpd}{l} = 1.69 \int 15p \cdot 200 \ 0.2p \cong 41 \ ps/mm$$

$$\frac{E}{l} = 0.4 \ pJ/mm$$

$$(pmos = 36 \mu m)$$

$$\frac{l}{N} = 0.45 \ mm \ \mathcal{R} - 1 \ mm \ W = 18 \ \mu m$$

$$\int transform = 18 \ mm \ fr = 18 \ mm$$