Performance of CMOS Circuits

Instructed by Shmuel Wimer Eng. School, Bar-Ilan University

Credits: David Harris Harvey Mudd College

(Some material copied/taken/adapted from Harris' lecture notes)

Dec 2010

Outline

- □ Gate and Diffusion Capacitance
- □ RC Delay Models
- Power and Energy
- Dynamic Power
- Static Power
- □ Low Power Design

MOSFET Capacitance



Dec 2010

- Any two conductors separated by an insulator have capacitance
- Gate to channel capacitor is very important
 - Creates channel charge necessary for operation
- Source and drain have capacitance to body
 - Across reverse-biased diodes
 - Called diffusion capacitance because it is associated with source/drain diffusion

Gate Capacitance

Approximate channel as connected to source

- $\Box C_{gs} = \varepsilon_{ox}WL/t_{ox} = C_{ox}WL = C_{permicron}W$
- $\hfill\square\hfill C_{permicron}$ is typically about 2 fF/µm



Dec 2010



Accumulation occurs when Vg is negative (for P material). Holes are induced under the oxide. Cgate = CoxA where Cox = ε SiO2 ε O/tox



Depletion occurs when Vg is near zero but < Vtn. Here the Cgate is given by CoxA in series with depletion layer capacitance Cdep



Inversion occurs when V_g is positive and > Vtn (for P material). A model for inversion in comprised of Cox A connecting from gateto-channel and Cdep connecting from channel-to-substrate.



Normalized gate capacitance versus Gate voltage Vgs. High freq behavior is due to the distributed resistance of channel

Normalized Experimental MOS Gate Capacitance Measurements vs Vds, Vgs

For $V_{ds} = 0$, the total gate capacitance $C_{ox} A$ splits equally to the drain and source of the transistor.



Dec 2010

Performance of CMOS Circuits

For $V_{ds} > 0$, the gate capacitance tilts more toward the source and becomes roughly 2/3 $C_{ox}A$ to the source and 0 to the drain for high V_{ds} .



Dec 2010

Performance of CMOS Circuits

Higher V_{gs} – Vt forces this tilting to occur later, since the device is linear up to V_{gs} – Vt = Vds.



Dec 2010

MOS Transistor Gate Capacitance Model

Gate capacitance has different components in different modes, but total remains constant.



Dec 2010

Gate capacitance has different components in different modes, but total remains constant.

CAPACITANCE				
Parameter	ю.;	Off	Non-saturated	Saturated
C_{gb}		$\frac{\varepsilon A}{t_{ox}}$	0	0
C_{gs}		0 .	$\frac{\epsilon A}{2t_{ox}}$	$\frac{2\varepsilon A}{3t_{ox}}$
C_{gd}		0	$\frac{\varepsilon A}{2t_{ox}}$	0 (finite for short channel devices)
$C_g = C_{gb} + C_{gs} + C_{gd}$		$\frac{\epsilon A}{t_{ox}}$	$\frac{\varepsilon A}{t_{ox}}$	$\frac{2\varepsilon A}{3t_{ox}} \rightarrow \frac{.9 \varepsilon A}{t_{ox}} \text{ (short channel)}$

Dec 2010

Diffusion Capacitance

- $\square C_{sb}, C_{db}$
- Undesirable, called *parasitic* capacitance
- Capacitance depends on area and perimeter
 - Use small diffusion nodes
 - Comparable to C_g
 - for contacted diff
 - $-\frac{1}{2}C_{g}$ for uncontacted
 - Varies with process

Diffusion Capacitance (Cont'd)



Effective Resistance

- □ Shockley models have limited value
 - Not accurate enough for modern transistors
 - Too complicated for much hand analysis
- Simplification: treat transistor as resistor
 - Replace $I_{ds}(V_{ds}, V_{gs})$ with effective resistance R
 - $I_{ds} = V_{ds}/R$
 - R averaged across switching of digital gate
- □ Too inaccurate to predict current at any given time
 - But good enough to predict RC delay

RC Delay Model

Use equivalent circuits for MOS transistors

- Ideal switch + capacitance and ON resistance
- Unit nMOS has resistance R, capacitance C
- Unit pMOS has resistance 2R, capacitance C
- □ Capacitance proportional to width
- Resistance inversely proportional to width



RC Values

□ Capacitance

 $-C = C_g = C_s = C_d = 2 \text{ fF}/\mu \text{m of gate width}$

– Values similar across many processes

- Resistance
 - R \approx 6 K $\Omega^* \mu m$ in 0.6um process
 - Improves with shorter channel lengths
- Unit transistors
 - May refer to minimum contacted device (4/2 λ)
 - Or maybe 1 μm wide device
 - Doesn't matter as long as you are consistent

Inverter Delay Estimate

□ Estimate the delay of a fanout-of-1 inverter



Transient Response

- \Box DC analysis tells us V_{out} if V_{in} is constant
- **Transient analysis** tells us $V_{out}(t)$ if $V_{in}(t)$ changes
 - Requires solving differential equations
- □ Input is usually considered to be a step or ramp
 - From 0 to $V_{\text{DD}}\,\text{or}\,\text{vice}\,\text{versa}$

Inverter Step Response

□ Ex: find step response of inverter driving load cap



Dec 2010

Inverter Step Response

Ex: find step response of inverter driving load cap



Dec 2010

Delay Definitions



Dec 2010





Dec 2010

Delay Definitions (Cont'd)

- □ **t**_{pdr}: *rising propagation delay*
 - Maximum time from input crossing 50% to rising output crossing 50%
- □ **t**_{pdf}: falling propagation delay
 - Maximum time from input crossing 50% to falling output crossing 50%
- □ **t**_{pd}: average propagation delay
 - $t_{pd} = (t_{pdr} + t_{pdf})/2$
- □ t_r: rise time
 - From output crossing 0.2 V_{DD} to 0.8 V_{DD}

Delay Definitions (Cont'd)

- □ t_f: fall time
 - From output crossing 0.8 V_{DD} to 0.2 V_{DD}
- □ **t**_{cdr}: *rising contamination delay*
 - Minimum time from input crossing 50% to rising output crossing 50%
- □ t_{cdf}: falling contamination delay
 - Minimum time from input crossing 50% to falling output crossing 50%
- □ **t**_{cd}: average contamination delay
 - t_{cd} = (t_{cdr} + t_{cdf})/2

Simulated Inverter Delay

- Solving differential equations by hand is too hard
- □ SPICE simulator solves the equations numerically
 - Uses more accurate I-V models too!
- But simulations take time to write



Dec 2010

Delay Estimation

- □ We would like to be able to easily estimate delay
 - Not as accurate as simulation
 - But easier to ask "What if?"
- The step response usually looks like a 1st order RC response with a decaying exponential.
- □ Use RC delay models to estimate delay
 - C = total capacitance on output node
 - Use *effective resistance* R
 - So that $t_{pd} = RC$
- □ Characterize transistors by finding their effective R
 - Depends on average current as gate switches

RC Delay Model

Use equivalent circuits for MOS transistors

- Ideal switch + capacitance and ON resistance
- Unit nMOS has resistance R, capacitance C
- Unit pMOS has resistance 2R, capacitance C
- □ Capacitance proportional to width
- Resistance inversely proportional to width



Example: 3-input NAND

Sketch a 3-input NAND with transistor widths chosen to achieve effective rise and fall resistances equal to a unit inverter (R).



3-input NAND Caps

Annotate the 3-input NAND gate with gate and diffusion capacitance.



3-input NAND Caps (Cont'd)

Annotate the 3-input NAND gate with gate and diffusion capacitance.



Dec 2010

Elmore Delay



Dec 2010

Elmore Delay

- ON transistors look like resistors
- Pullup or pulldown network modeled as RC ladder
- □ Elmore delay of RC ladder

$$t_{pd} \approx \sum_{\text{nodes } i} R_{i-to-source} C_i$$

$$= R_1 C_1 + (R_1 + R_2) C_2 + \dots + (R_1 + R_2 + \dots + R_N) C_N$$



For a step input Vin, the delay at any node can be estimated with the Elmore delay equation $\mathbf{t}_{Di} = \Sigma \mathbf{C}_i \Sigma \mathbf{R}_k$

For example, the Elmore delay at node 7 is give by:

R1(C1 + C2 + C3 + C4 + C5 + C6 + C7 + C8) +



Dec 2010



The 50 percent delay from A to B is

 $T_{AB} = R_1(C_2 + C_3 + C_4 + C_5 + C_6 + C_7 + C_8 + C_9 + C_{10})$

$$+R_2\left(\frac{C_2}{2} + C_3 + C_4 + C_5 + C_6 + C_7 + C_8 + C_9 + C_{10}\right)$$

$$+R_3\left(\frac{C_3}{2}+C_4+C_5+C_6+C_7+C_8+C_9+C_{10}\right)$$

Dec 2010



The 50 percent delay from B to D is

$$T_{BD} = R_4 \left(\frac{C_4}{2} + C_7 + C_8 + C_9 \right) + R_7 \left(\frac{C_7}{2} + C_9 \right).$$

Dec 2010



The 50 percent delay from B to E is

$$T_{BE} = R_5 \left(\frac{C_5}{2} + C_6 + C_{10} \right) + R_6 \left(\frac{C_6}{2} + C_{10} \right).$$

Dec 2010

Example: 2-input NAND

Estimate rising and falling propagation delays of a 2input NAND driving *h* identical gates.



Example: 2-input NAND

 Estimate rising and falling propagation delays of a 2-input NAND driving *h* identical gates.



$$\begin{array}{c} \mathbf{x} \quad \frac{R/2}{2C} \mathbf{y} \\ \mathbf{k} \\ \mathbf{k}$$

Dec 2010

Delay Components

- Delay has two parts
 - Parasitic delay
 - 6 or 7 RC
 - Independent of load
 - Effort delay
 - 4h RC
 - Proportional to load capacitance

Contamination Delay

- Best-case (contamination) delay can be substantially less than propagation delay.
- Ex: If both inputs fall simultaneously



Dec 2010

Diffusion Capacitance

- we assumed contacted diffusion on every s / d
- Good layout minimizes diffusion area
- □ Ex: NAND3 layout shares one diffusion contact
 - Reduces output capacitance by 2C
 - Merged uncontacted diffusion might help too



Layout Comparison

- □ Layout representation by *stick diagram*. What CKT?
- Which layout is better?



Power and Energy

- □ Power is drawn from a voltage source attached to the V_{DD} pin(s) of a chip.
- □ Instantaneous Power: $P(t) = i_{DD}(t)V_{DD}$

Energy:
$$E = \int_{0}^{T} P(t) dt = \int_{0}^{T} i_{DD}(t) V_{DD} dt$$
Average Power:
$$P_{\text{avg}} = \frac{E}{T} = \frac{1}{T} \int_{0}^{T} i_{DD}(t) V_{DD} dt$$

Dynamic Power

- 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 Tf_{sw} times over an interval of T



Dec 2010

Performance of CMOS Circuits

$$P_{\text{dynamic}} = \frac{E}{T} = \frac{1}{T} \int_{0}^{T} i_{DD}(t) V_{DD} dt =$$

$$\frac{V_{DD}}{T}\int_{0}^{T}i_{DD}(t)dt = \frac{V_{DD}}{T}\left[Tf_{sw}CV_{DD}\right] = CV_{DD}^{2}f_{sw}$$



Dec 2010

Activity Factor

- Suppose the system clock frequency = f
- \Box Let $f_{sw} = \alpha f$, where $\alpha = activity factor$
 - If the signal is a clock, $\alpha = 1$
 - If the signal switches once per cycle, $\alpha = \frac{1}{2}$
 - Static gates:
 - Depends on design, but typically $\alpha = 0.1$
 - Dynamic gates:
 - Switch either 0 or 2 times per cycle, $\alpha = \frac{1}{2}$

Dynamic power:

$$P_{\rm dynamic} = \alpha C V_{DD}^2 f$$

Dec 2010

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

Power Dissipation Sources

- $\square P_{total} = P_{dynamic} + P_{static}$
- **Dynamic power:** $P_{dynamic} = P_{switching} + P_{shortcircuit}$
 - Switching load capacitances
 - Short-circuit current
- **Static power:** $P_{\text{static}} = (I_{\text{sub}} + I_{\text{gate}} + I_{\text{junct}} + I_{\text{contention}})V_{\text{DD}}$
 - Sub-threshold leakage
 - Gate leakage
 - Junction leakage
 - Contention current (ratioed logic)

Dynamic Power Example

- □ 1 billion transistor chip
 - 50M logic transistors
 - Average width: 12 λ
 - Activity factor = 0.1
 - 950M memory transistors
 - Average width: 4 λ
 - Activity factor = 0.02
 - 1.0 V 65 nm process
 - C = 1 fF/ μ m (gate) + 0.8 fF/ μ m (diffusion)
- Estimate dynamic power consumption @ 1 GHz. Neglect wire capacitance and short-circuit current.

Power Estimate Ex (Cont'd)

$$C_{\text{logic}} = \left(50 \times 10^6\right) \left(12\lambda\right) \left(0.025\,\mu m \,/\,\lambda\right) \left(1.8\,f F \,/\,\mu m\right) = 27nF$$

$$C_{\rm mem} = (950 \times 10^6) (4\lambda) (0.025 \,\mu m / \lambda) (1.8 \, fF / \mu m) = 171 nF$$

$$P_{\text{dynamic}} = \left[0.1C_{\text{logic}} + 0.02C_{\text{mem}}\right] \left(1.0\right)^2 \left(1.0GHz\right) = 6.1W$$

Dec 2010

Dynamic Power Reduction

 $\Box P_{\text{switching}} = \alpha C V_{DD}^{2} f$

- **Try to minimize:**
 - Activity factor
 - Capacitance
 - Supply voltage
 - Frequency

Activity Factor Estimation

- $\Box \quad \text{Let } P_i = \text{Prob}(\text{node } i = 1)$
- $\Box \ \alpha_i = P_i^{*}(1 P_i)$
- **Completely random data has P = 0.5 and** α = 0.25
- Data is often not completely random
 - e.g. MSBs of 64-bit words in memory address bus. MSBs of data representing measurements of physical phenomena.
- Data propagating through ANDs and ORs has lower activity factor
 - Depends on design, but typically $\alpha \approx 0.1$

Switching Probability

Gate	P _Y
AND2	$P_A P_B$
AND3	$P_{\mathcal{A}}P_{B}P_{C}$
OR2	$1 - \overline{P}_A \overline{P}_B$
NAND2	$1 - P_A P_B$
NOR2	$\overline{P}_{\mathcal{A}}\overline{P}_B$
XOR2	$P_{\!\mathcal{A}}\overline{P}_B + \overline{P}_{\!\mathcal{A}}P_B$

What is the switching probability?

Dec 2010

Example

- A 4-input AND is built out of two levels of gates
- □ Estimate the activity factor at each node if the inputs have P = 0.5



Dec 2010

Clock Gating

- □ The best way to reduce the activity is to turn off the clock to registers in unused blocks
 - Saves clock activity ($\alpha = 1$)
 - Eliminates all switching activity in the block
 - Requires determining if block will be used



Capacitance

- □ Gate capacitance
 - Fewer stages of logic
 - Small gate sizes
- □ Wire capacitance
 - Good floorplanning to keep communicating blocks close to each other
 - Drive long wires with inverters or buffers rather than complex gates

Voltage / Frequency

- Run each block at the lowest possible voltage and frequency that meets performance requirements
- Voltage Domains
 - Provide separate supplies to different blocks
 - Level converters required when crossing from low to high V_{DD} domains
- Dynamic Voltage Scaling – Adjust V_{DD} and f according to workload



Static Power

- Static power is consumed even when chip is quiescent
 - Ratioed circuits burn power in fight between ON transistors. Occurs when output is low (0).
 - Leakage draws power from nominally OFF devices

Static Power Example

- □ Revisit power estimation for 1 billion transistor chip
- □ Estimate static power consumption
 - Subthreshold leakage
 - Normal V_t: 100 nA/μm
 - High V_t : 10 nA/ μ m
 - High Vt used in all memories and in 95% of logic gates
 - Gate leakage 5 nA/μm
 - Junction leakage negligible

Solution

$$W_{\text{normal-V}_{t}} = (50 \times 10^{6})(12\lambda)(0.025\,\mu\text{m}/\lambda)(0.05) = 0.75 \times 10^{6}\,\mu\text{m}$$
$$W_{\text{high-V}_{t}} = \left[(50 \times 10^{6})(12\lambda)(0.95) + (950 \times 10^{6})(4\lambda) \right] (0.025\,\mu\text{m}/\lambda) = 109.25 \times 10^{6}\,\mu\text{m}$$
$$I_{sub} = \left[W_{\text{normal-V}_{t}} \times 100\,\text{nA}/\mu\text{m} + W_{\text{high-V}_{t}} \times 10\,\text{nA}/\mu\text{m} \right] / 2 = 584\,\text{mA}$$
$$I_{gate} = \left[(W_{\text{normal-V}_{t}} + W_{\text{high-V}_{t}}) \times 5\,\text{nA}/\mu\text{m} \right] / 2 = 275\,\text{mA}$$
$$P_{static} = (584\,\text{mA} + 275\,\text{mA})(1.0\,\text{V}) = 859\,\text{mW}$$

Dec 2010

Leakage Control

$$I_{ds} = I_{ds0} e^{\frac{V_{gs} - V_t}{nv_T}} \left(1 - e^{\frac{-V_{ds}}{v_T}} \right)$$

Leakage and delay trade off

- Aim for low leakage in sleep and low delay in active mode
- □ To reduce leakage:
 - Increase V_t : *multiple* V_t
 - Use low V_t only in critical circuits
 - Increase V_s: *stack effect*
 - Input vector control in sleep

Gate Leakage

- \Box Extremely strong function of t_{ox} and V_{gs}
 - Negligible for older processes
 - Approaches subthreshold leakage at 65 nm and below in some processes
- □ An order of magnitude less for pMOS than nMOS
- □ Control leakage in the process using $t_{ox} > 10.5$ Å
 - High-k gate dielectrics help
 - Some processes provide multiple tox
 - e.g. thicker oxide for 3.3 V I/O transistors
- □ Control leakage in circuits by limiting V_{DD}

Power Gating

- Turn OFF power to blocks when they are idle to save leakage
 - Use virtual V_{DD} (V_{DDV})
 - Gate outputs to prevent invalid logic levels to next block



- Voltage drop across sleep transistor degrades performance during normal operation
 - Size the transistor wide enough to minimize impact
- □ Switching wide sleep transistor costs dynamic power
 - Only justified when circuit sleeps long enough

Low Power Design

- Reduce dynamic power
 - $\square \alpha$: clock gating, sleep mode
 - C: small transistors (esp. on clock), short wires
 - V_{DD} : lowest suitable voltage
 - f: lowest suitable frequency
- □ Reduce static power
 - Selectively use ratioed circuits (minimize)
 - Selectively use low V_t devices (minimize)
 - Leakage reduction:

stacked devices, body bias, low temperature