



EECS 598 VLSI for Wireless Communication and Machine Learning

Digital Integrated Circuit Overview

Prof. Hun-Seok Kim

hunseok@umich.edu

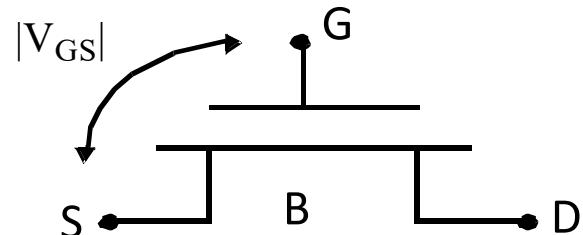
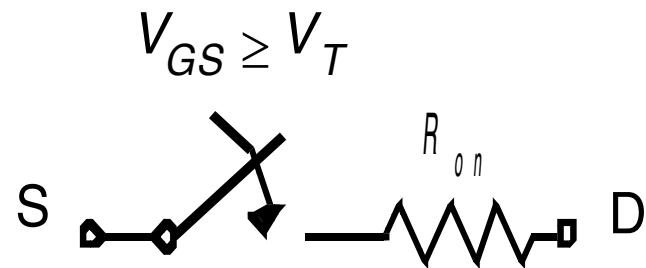


What is a Transistor?

A Switch!

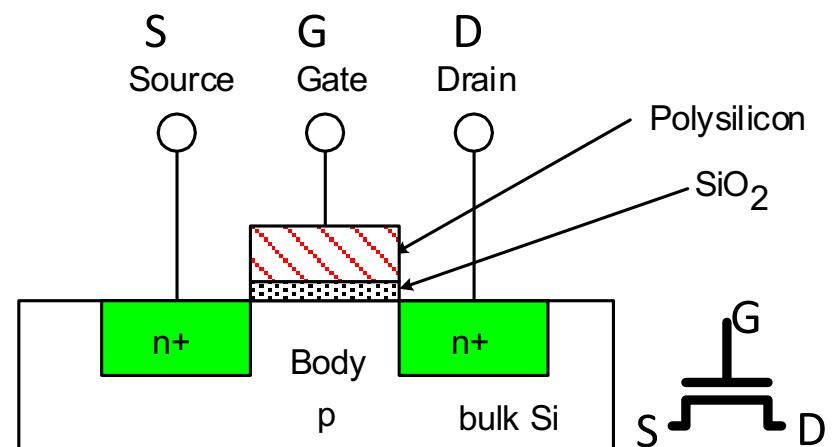
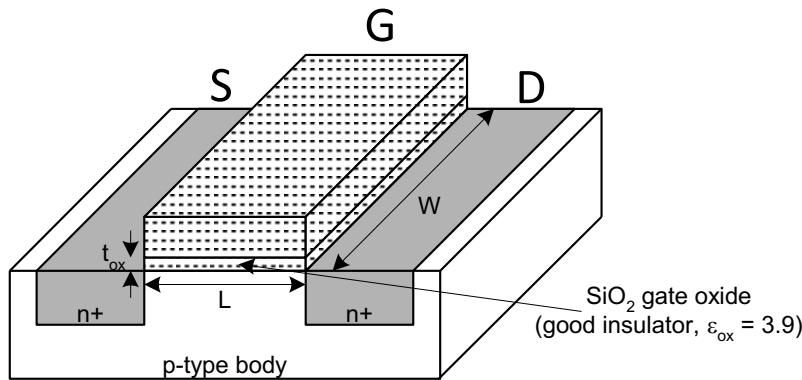


An MOS Transistor



MOS Transistor

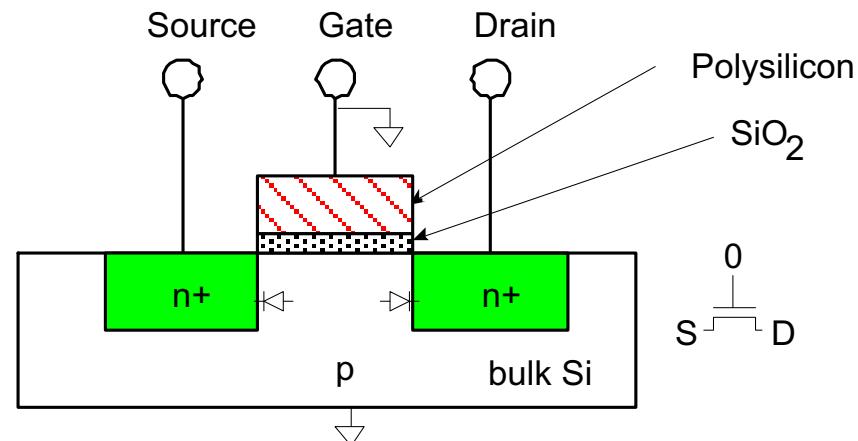
- Four terminals: gate, source, drain, body
- Gate – oxide – body stack looks like a capacitor
 - Gate and body are conductors
 - SiO_2 (oxide) is a very good insulator
 - Called metal – oxide – semiconductor (MOS) capacitor
 - Even though gate is often not made of metal





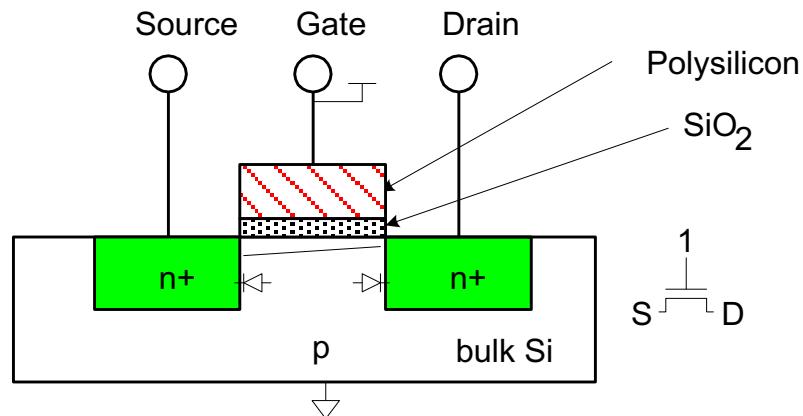
NMOS Operation

- Body is usually tied to ground (0 V)
- When the gate is at a low voltage:
 - P-type body is at low voltage
 - Source-body and drain-body diodes are OFF
 - No current flows, transistor is OFF



NMOS Operation cont.

- When the gate is at a high voltage:
 - Positive charge on gate of MOS capacitor
 - Negative charge attracted to body
 - Inverts a channel under gate to n-type
 - Now current can flow through n-type silicon from source through channel to drain, transistor is ON

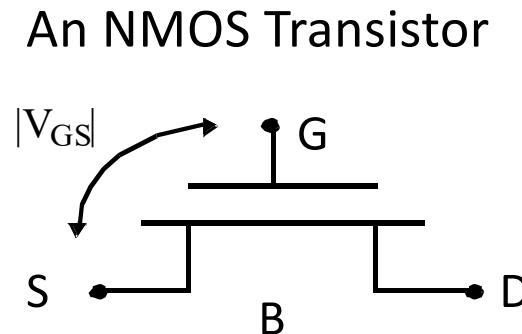




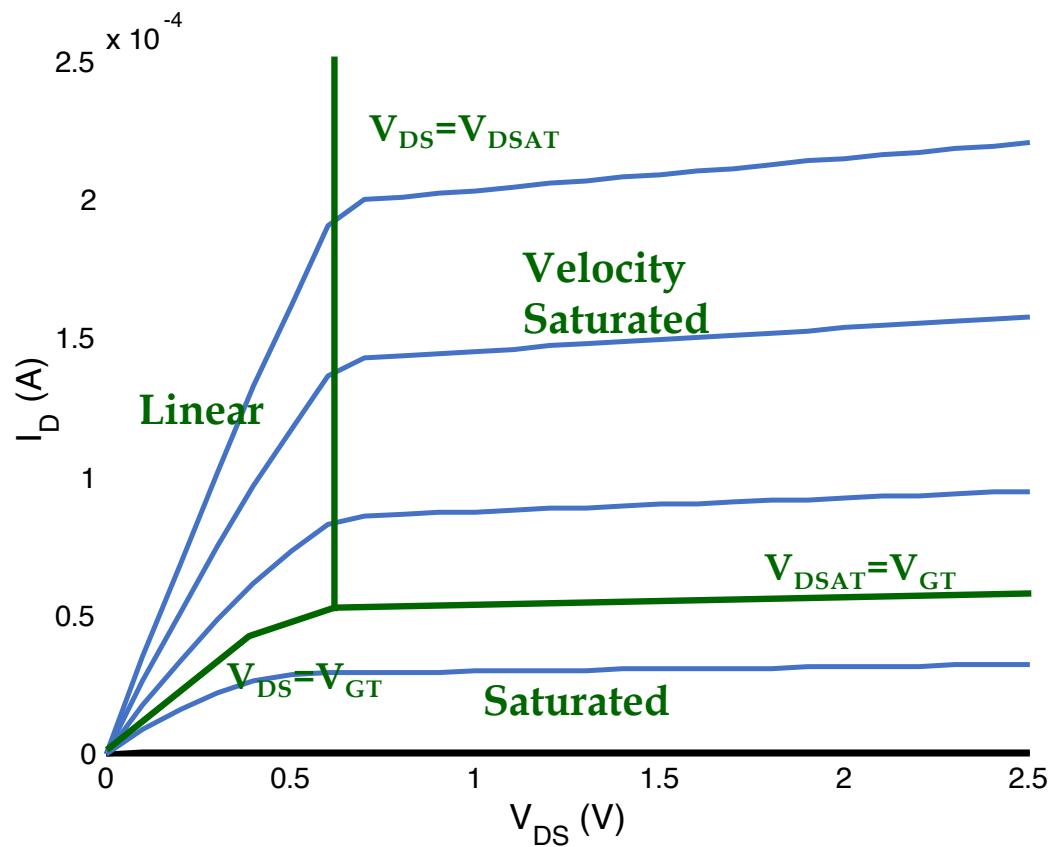
MOS FET IV characteristics (conductivity / resistance)

- Current I_D is a function of V_{GS} , V_{DS} , and V_T
- Transistor conductivity / resistance is controlled by terminal voltages
- Let
 - $V_{GT} = V_{GS} - V_T$
 - $V_{min} = \min(V_{GT}, V_{DS}, V_{DSAT})$
 - $V_T = V_{T0} + \gamma(\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|})$
- Then

$$\bullet I_{DS} = \begin{cases} 0 & V_{GT} \leq 0 \\ k' \left(\frac{W}{L} \right) \left(V_{GT} V_{min} - \frac{V_{min}^2}{2} \right) (1 + \lambda V_{DS}) & V_{GT} > 0 \end{cases}$$



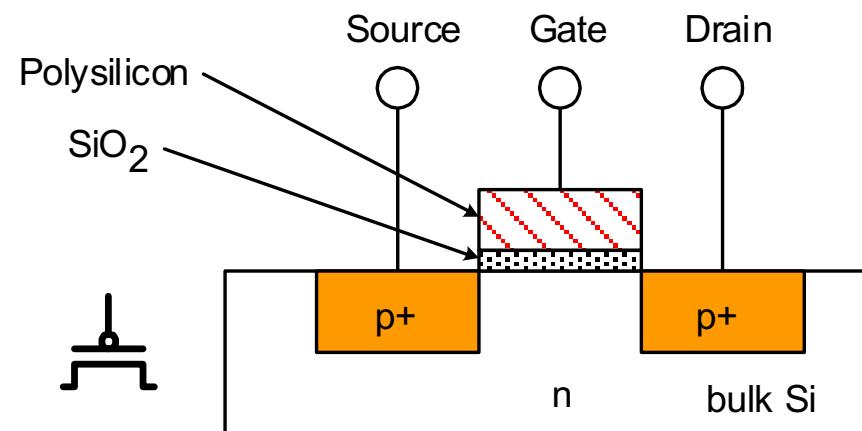
MOS FET IV characteristics (conductivity / resistance)





PMOS Transistor

- Similar, but doping and voltages reversed
 - Body tied to high voltage (V_{DD})
 - Gate low: transistor ON
 - Gate high: transistor OFF
 - Bubble indicates inverted behavior





Transistors as Switches

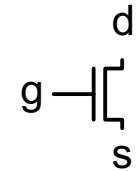
- We can view MOS transistors as electrically controlled switches
- Voltage at gate controls path from source to drain

notation

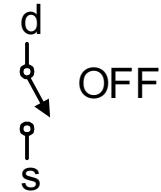
$I \Leftrightarrow V_{DD} \Leftrightarrow \text{high}$

$0 \Leftrightarrow GND \Leftrightarrow \text{low}$

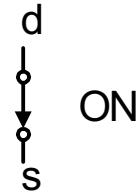
NMOS



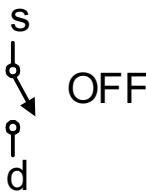
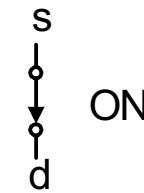
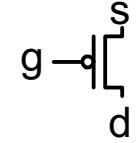
$g = 0$



$g = 1$

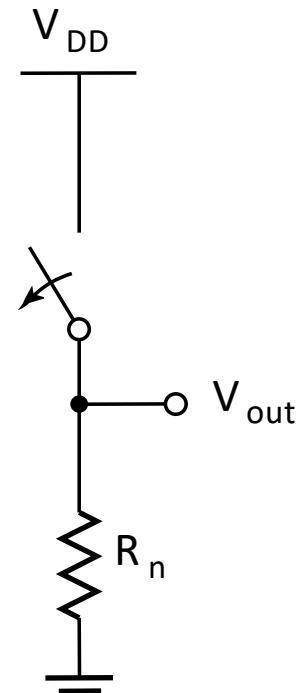
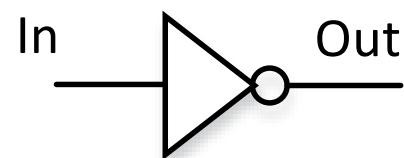
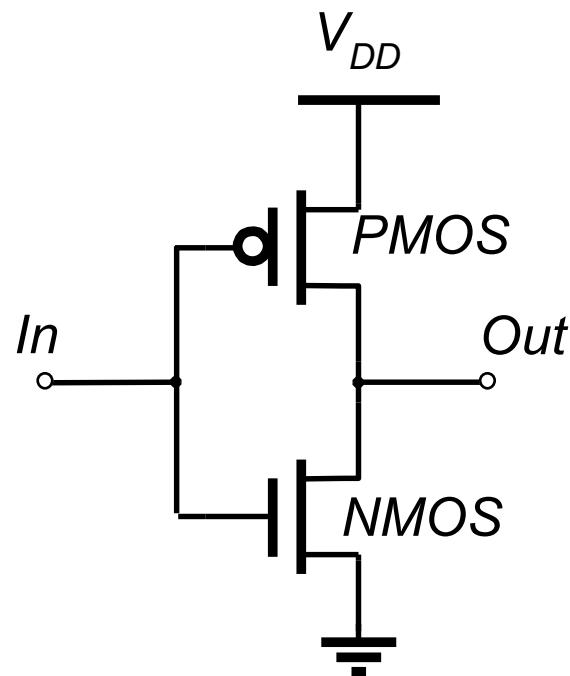


PMOS

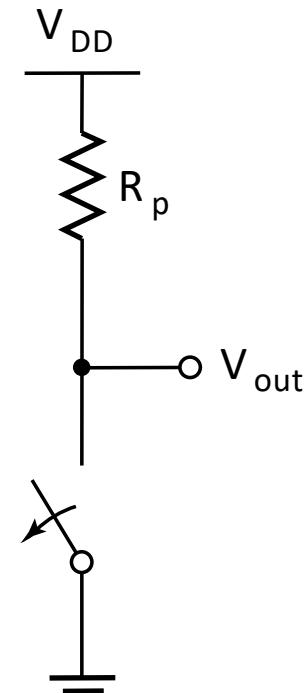




CMOS Inverter



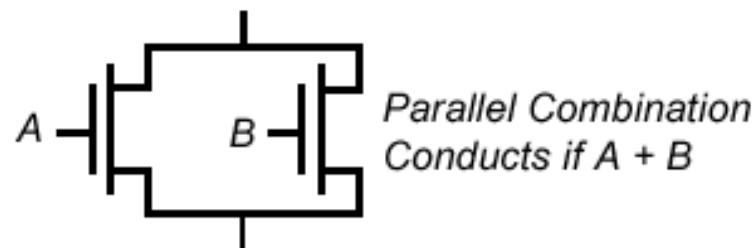
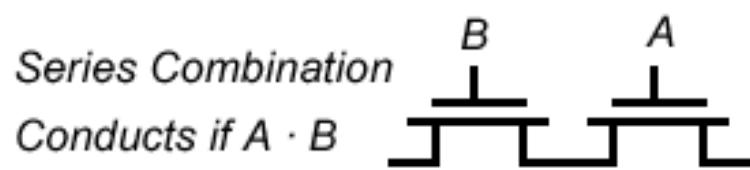
$$V_{in} = V_{DD}$$



$$V_{in} = 0$$



Two major functions: AND, OR



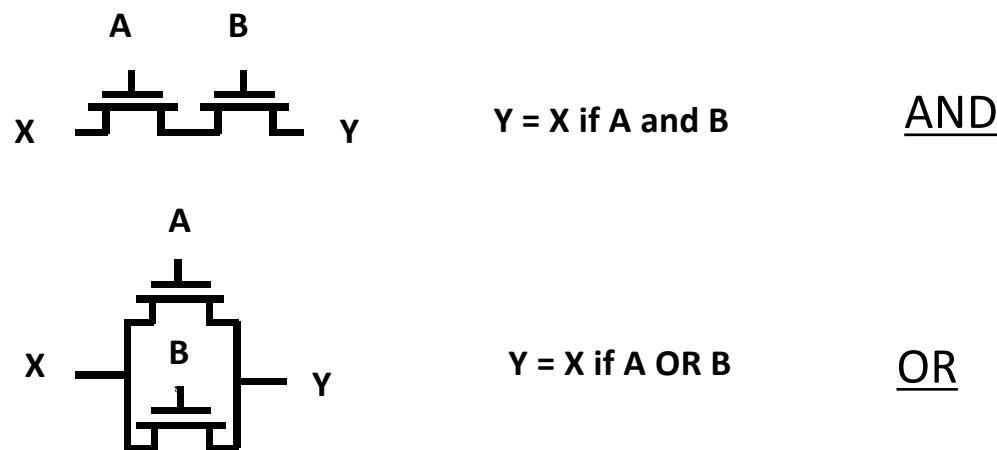
CMOS combinational gates are based on these two combinations of transistors (series & parallel)



NMOS Transistors in Series/Parallel Connection

Transistors can be thought as a switch controlled by its gate signal

NMOS switch closes when switch control input is high





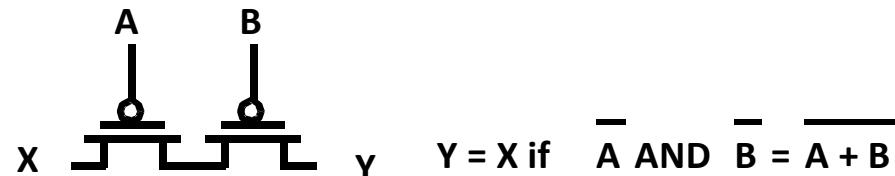
PMOS Transistors in Series/Parallel Connection

PMOS switch closes when switch control input is low

De Morgan

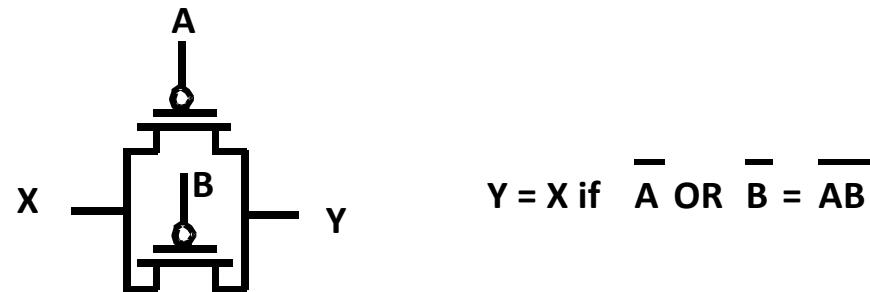
$$\overline{A+B} = \overline{A} \cdot \overline{B}$$

$$\overline{AB} = \overline{A} + \overline{B}$$



$$Y = X \text{ if } \overline{A} \text{ AND } \overline{B} = \overline{A+B}$$

NOR

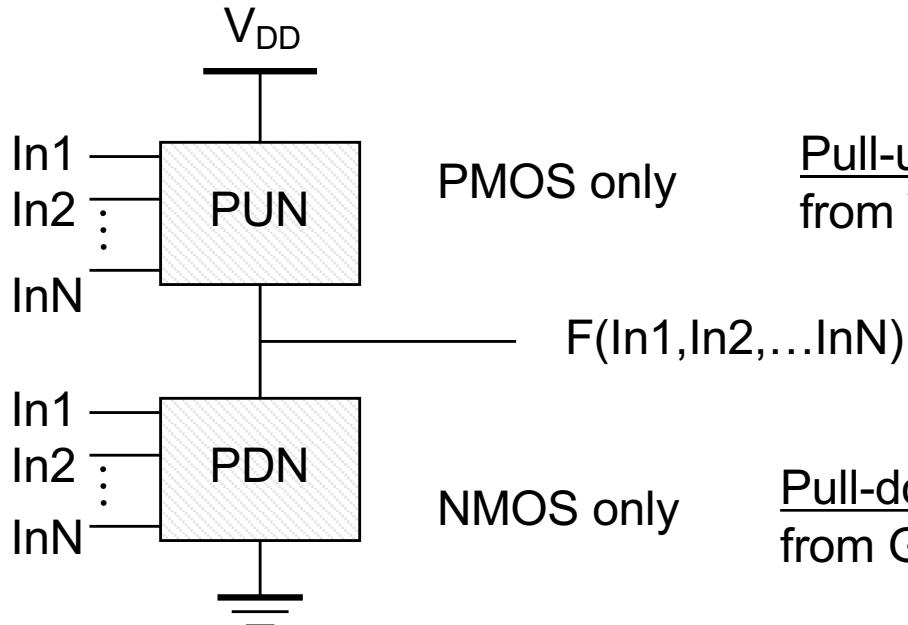


$$Y = X \text{ if } \overline{A} \text{ OR } \overline{B} = \overline{AB}$$

NAND



Static CMOS General Form



PMOS only

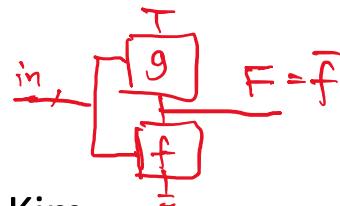
Pull-up network: make a connection from V_{dd} to F when F(In1,In2...)=1

NMOS only

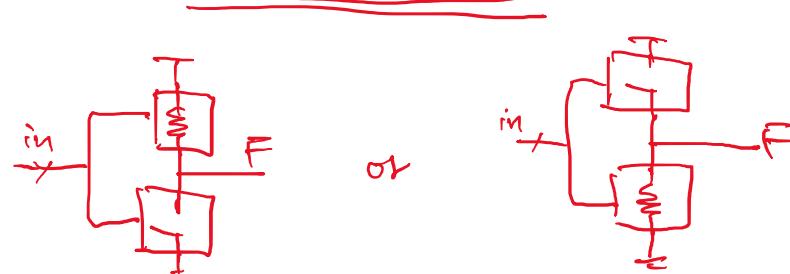
Pull-down network: make a connection from Ground to F when F(In1,In2...)=0

PUN and PDN are dual networks

$$F(A, B, \dots) = \overline{f(\bar{A}, \bar{B}, \dots)} = g(\bar{A}, \bar{B}, \dots)$$



EECS 598 Kim





Complementary CMOS Logic Style Construction

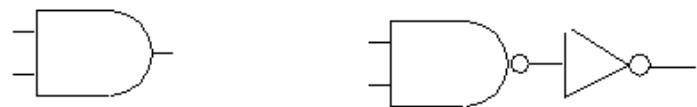
- PUN is the dual of the PDN

Can be shown using DeMorgan's Theorem

$$\overline{A + B} = \bar{A}\bar{B}$$

$$\overline{AB} = \bar{A} + \bar{B}$$

- The complementary gate is inverting

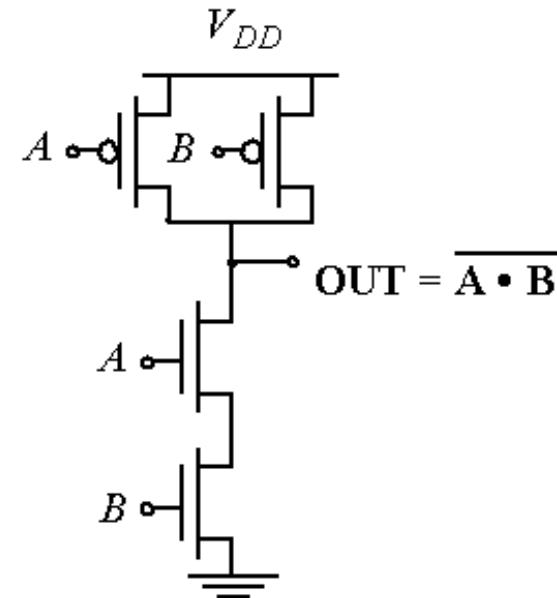


$$\mathbf{AND} = \mathbf{NAND} + \mathbf{INV}$$

Example Gate: NAND (most common gate in digital circuits)

A	B	Out
0	0	1
0	1	1
1	0	1
1	1	0

Truth Table of a 2 input NAND gate



PDN: Connects OUT to ground when $A \bullet B = 1$

PUN: Connects OUT to V_{dd} when $\bar{A} + \bar{B} = 1$

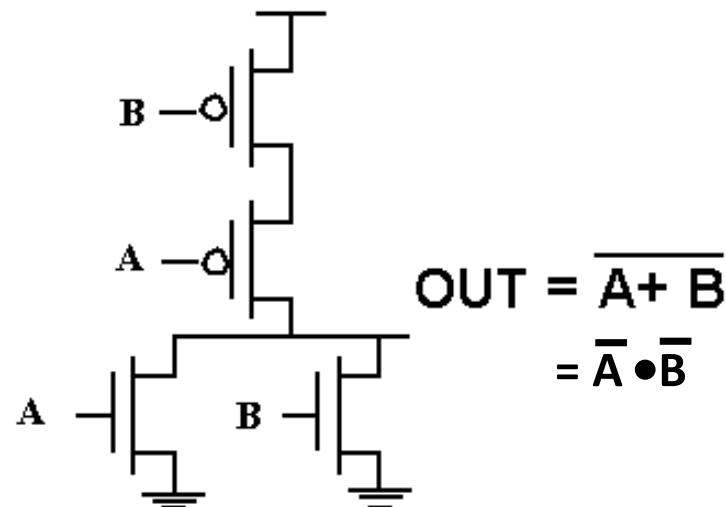
So $OUT = \text{Complement of PDN function}$

Also $OUT = \text{PUN function with each input inverted}$

Example Gate: NOR

A	B	Out
0	0	1
0	1	0
1	0	0
1	1	0

Truth Table of a 2 input NOR gate



N-input gates require $2N$ inputs for static CMOS

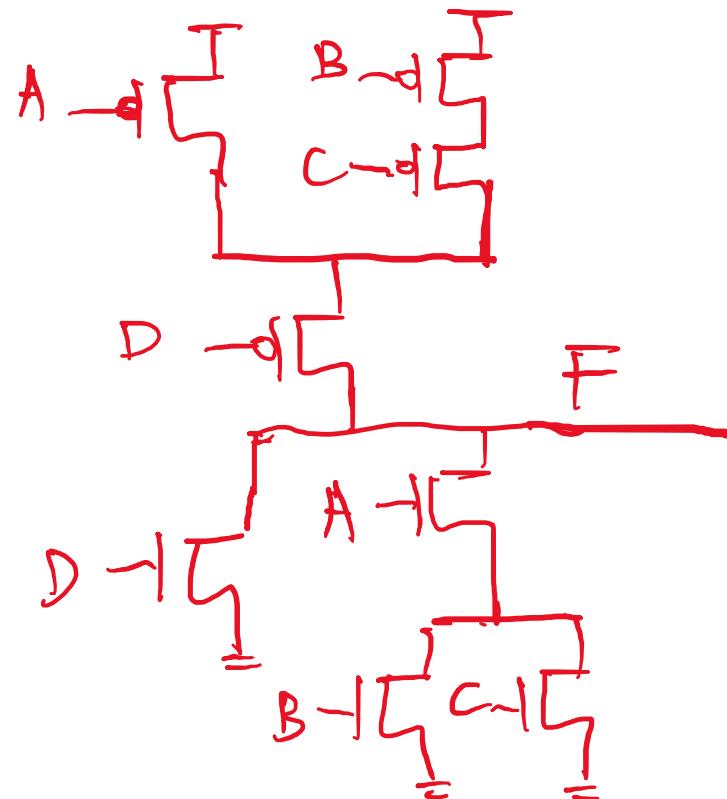


Building Complex Gates

$$F = \overline{D + A(B+C)}$$

$$f = D + A(B+C)$$

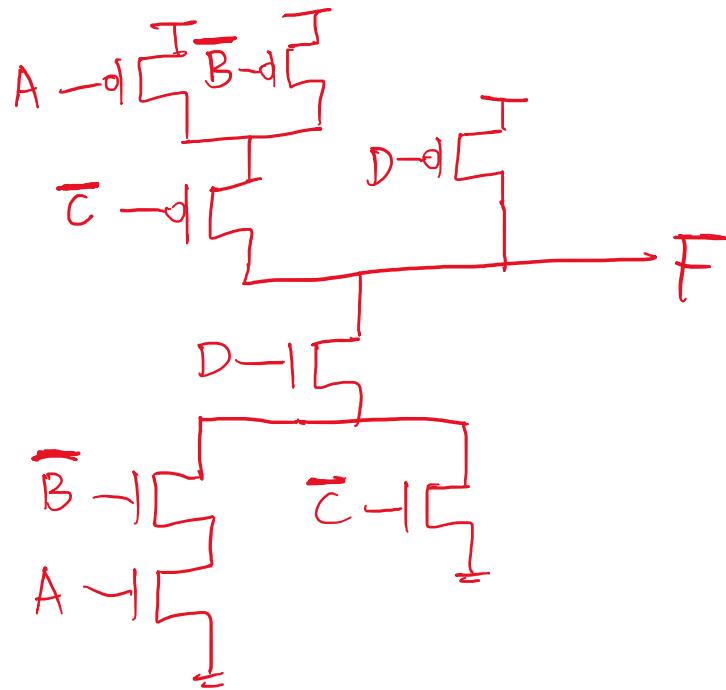
$$g = \overline{D}(\overline{A} + \overline{B}\overline{C})$$



$$F = (\bar{A} + B)C + \bar{D} = \overline{(\bar{A} + B)C + \bar{D}} = \overline{(A \cdot \bar{B} + \bar{C})D} = \overline{(\bar{A} + B)C + \bar{D}}$$

$$f = (A \cdot \bar{B} + \bar{C})D$$

$$g = (\bar{A} + B)C + \bar{D}$$



note: inputs
 are A, \bar{B}, \bar{C}, D
 for both PDN & PUN

inputs must be
 the same for PDN, PUN
 of a CMOS gate

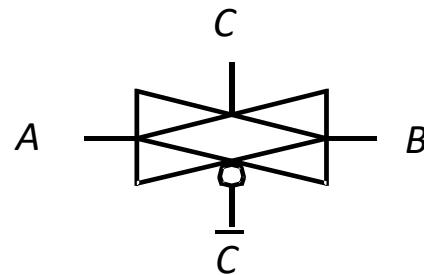
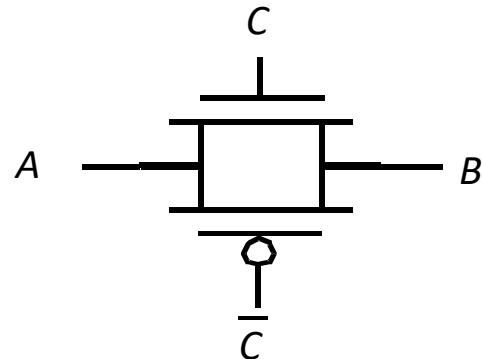


Static Properties of CMOS Gates (same as inverters)

- Full rail-to-rail swing (high noise margins)
- Always a connection to V_{dd} or ground in steady state (low output impedance)
- No direct path between V_{dd} and ground → No static power dissipation (ignoring leakage currents)
- Gates are inverting → if we want an AND gate, we need to add an inverter after a NAND gate



Transmission Gate

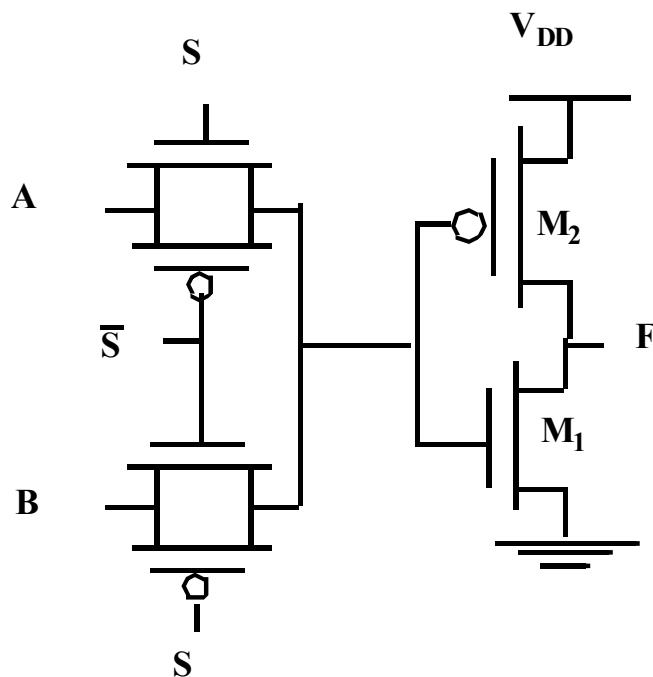


Multiplexers are commonly used in digital design

Function: Select one of N inputs based on a controlling signal (S)

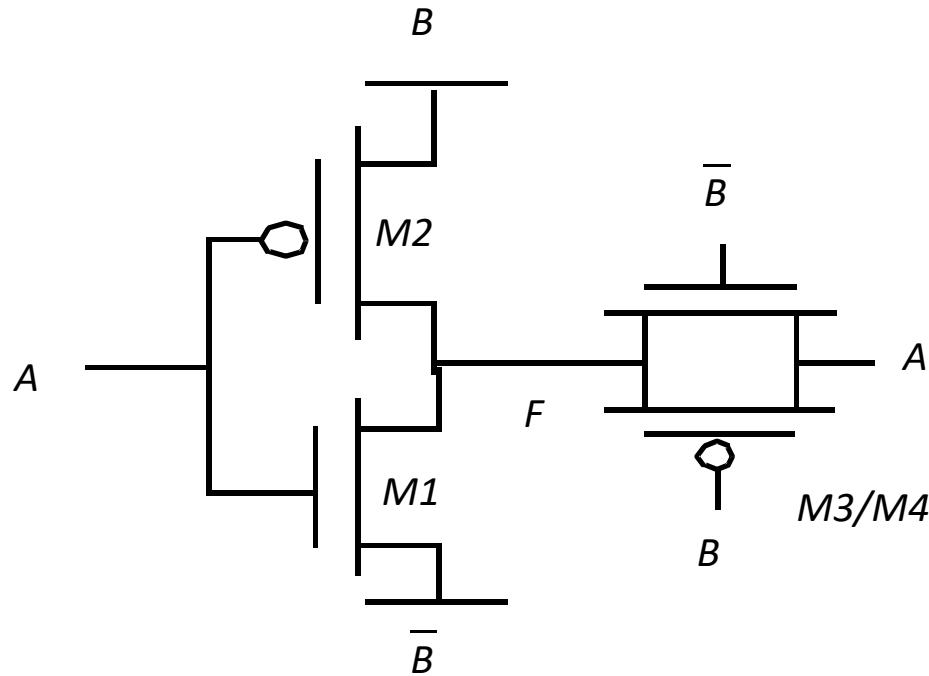
Here $N = 2$

of controlling signals = $\log_2 N$





Transmission Gate XOR



XOR \rightarrow exclusive OR

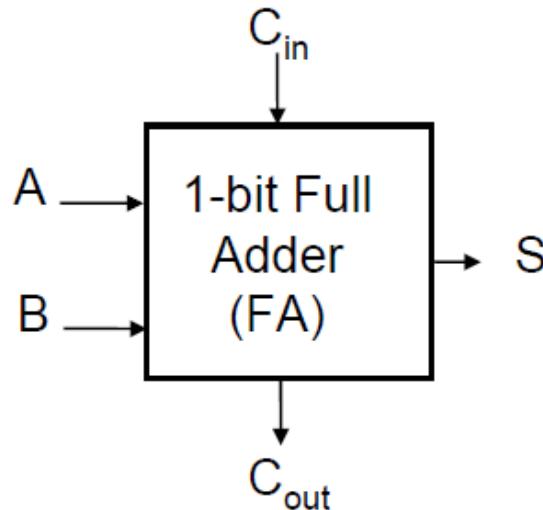
Same inputs: $F = 0$

Different inputs: $F = 1$

Fairly compact compared to
static CMOS XOR gates (8
transistors)

Full Adder

- $A[n-1:0] + B[n-1:0] = S[n-1:0]$

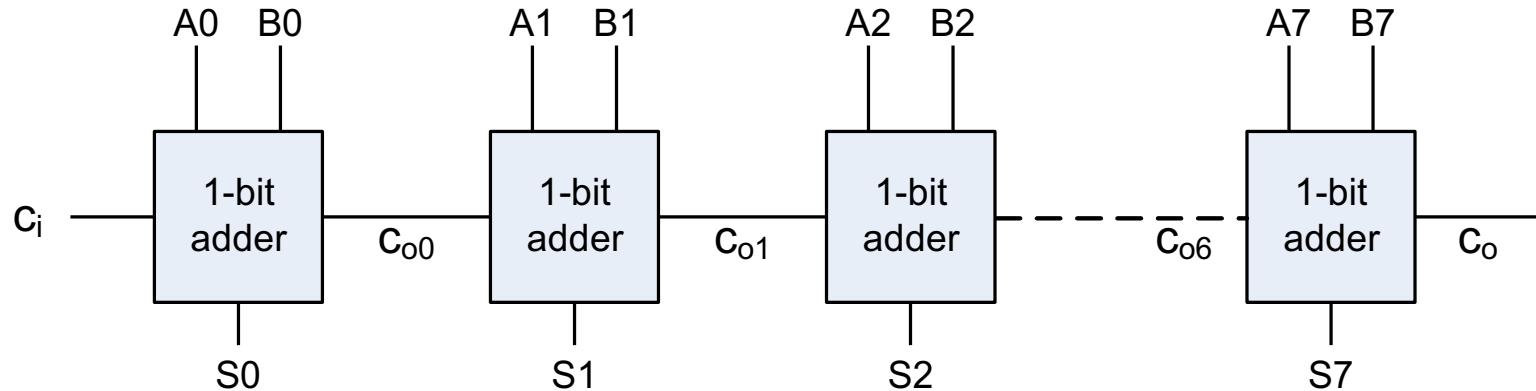


A	B	C _{in}	C _{out}	S	carry status
0	0	0	0	0	kill
0	0	1	0	1	kill
0	1	0	0	1	propagate
0	1	1	1	0	propagate
1	0	0	0	1	propagate
1	0	1	1	0	propagate
1	1	0	1	0	generate
1	1	1	1	1	generate

- $C_o = AB + BC_i + AC_i = AB + (A + B)C_i$
- $S = A \oplus B \oplus C_i = ABC_i + \overline{C_o}(A + B + C_i)$
- $G = AB, K = \overline{A} \overline{B}, P = A \oplus B$



Ripple Carry Adder



Worst case delay linear with the number of bits

$$t_d = O(N)$$

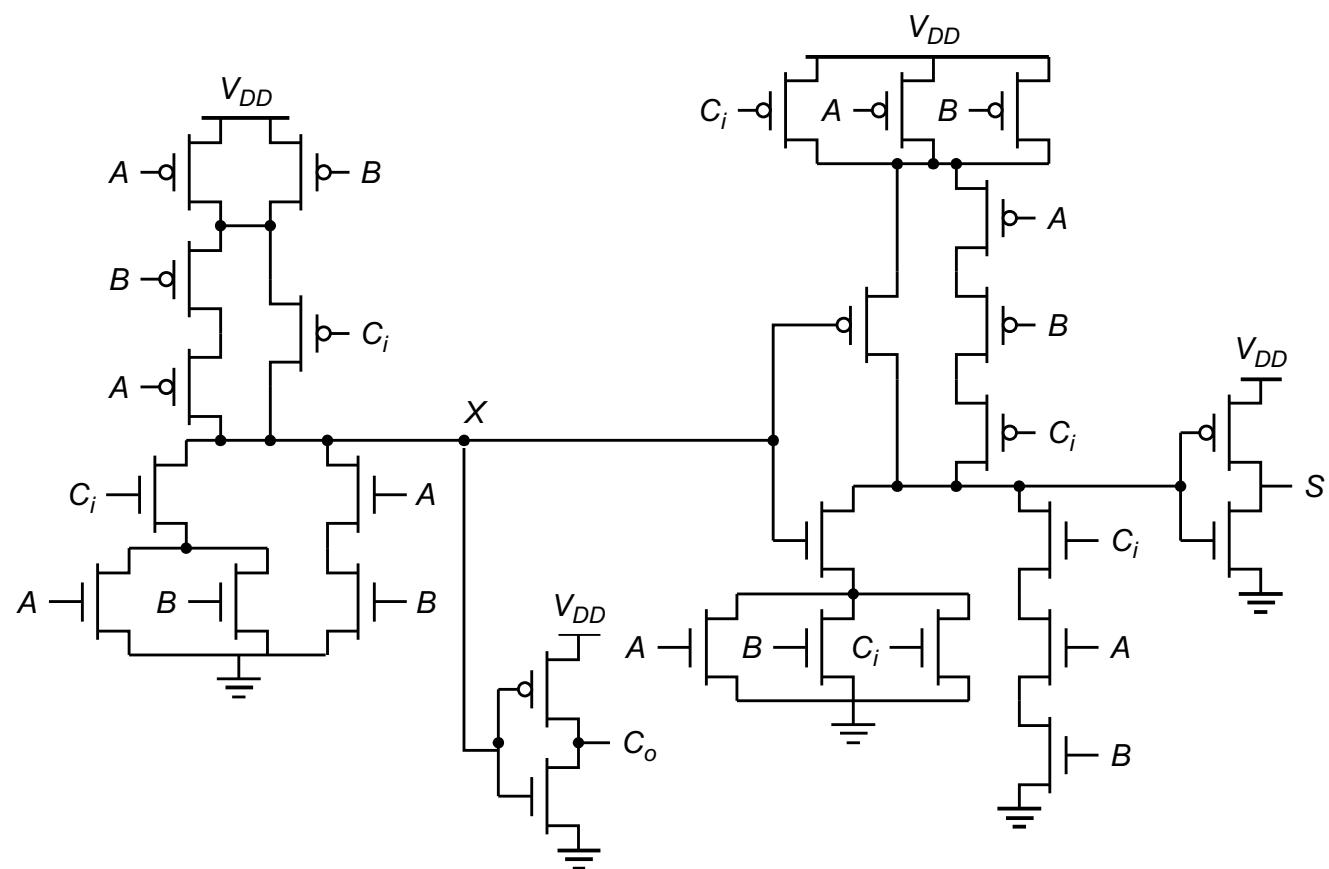
$$t_{adder} = (N-1)t_{carry} + t_{sum}$$

Goal: Make the fastest possible carry path circuit



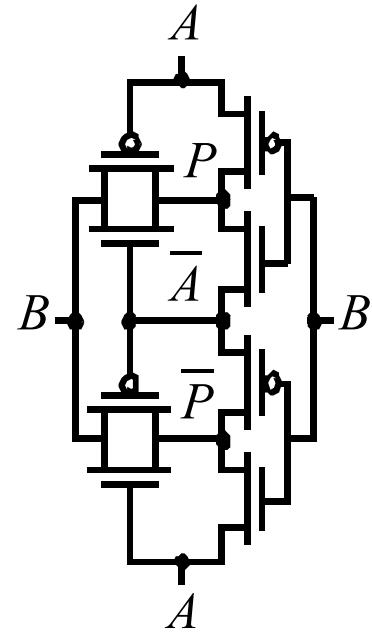
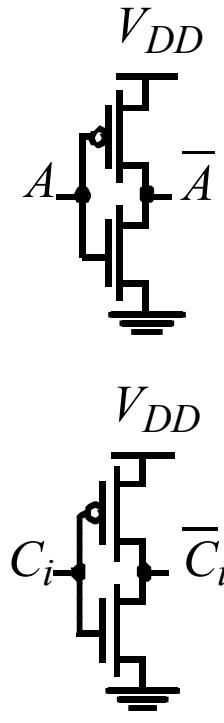
Static CMOS Full Adder

- $C_o = AB + BC_i + AC_i = AB + (A + B)C_i$
- $S = A \oplus B \oplus C_i = ABC_i + \overline{C_o}(A + B + C_i)$
- 28 transistors

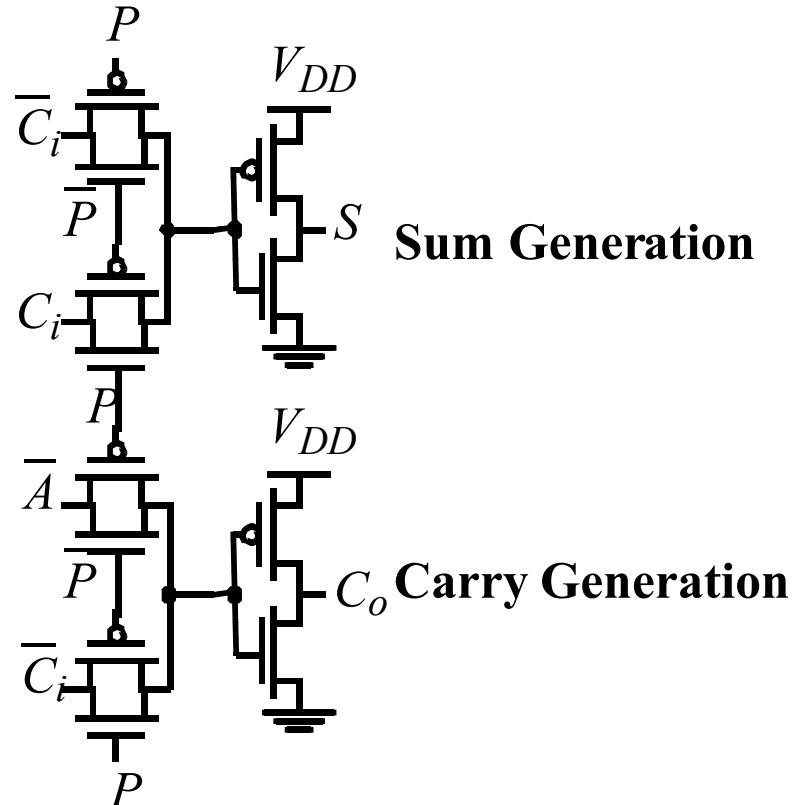




Transmission Gate Full Adder



Setup



Sum Generation

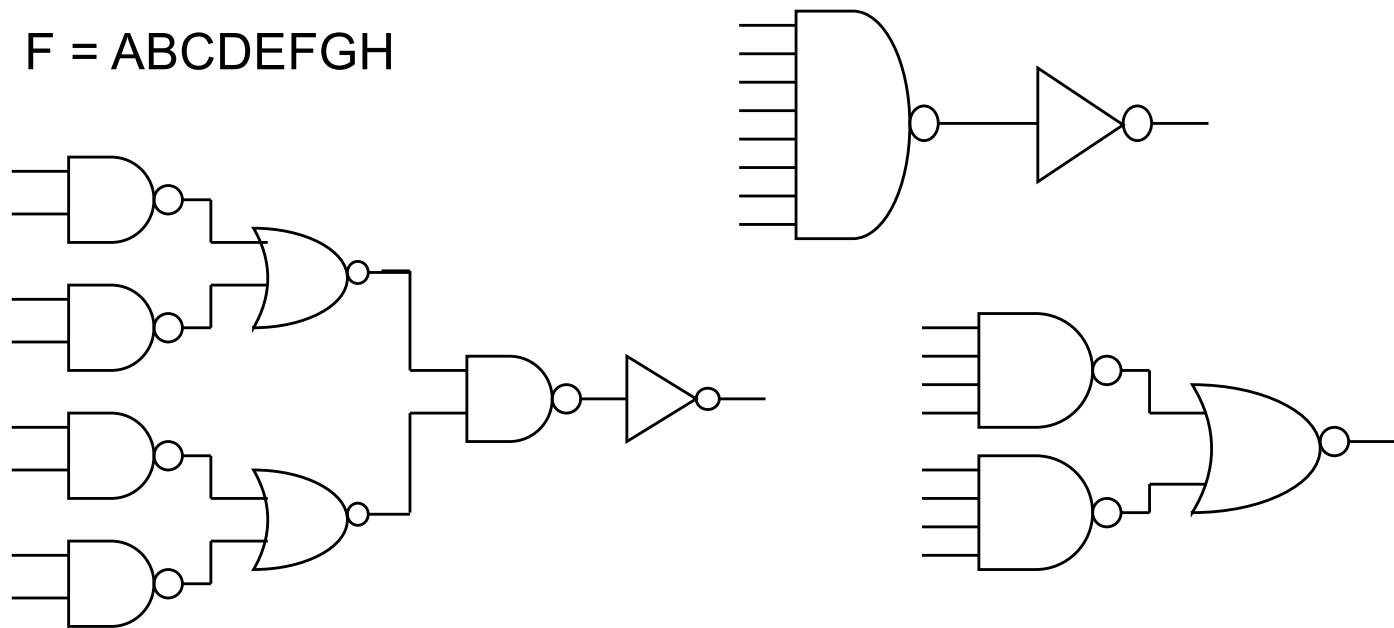
C_o Carry Generation



Synthesis using cell library

- Hierarchical logic structures

$$F = ABCDEFGH$$



Synthesis using standard ‘cell library’



Key VLSI design goals

What are key design metrics for VLSI for wireless communication or machine learning or any other applications?

- Functionality
- Performance (speed, delay)
- Power / Energy
- Area

What factors determine the performance (speed, delay), power(energy), and area of a circuit?

How are these metrics related?

MOS FET IV characteristics (conductivity / resistance)

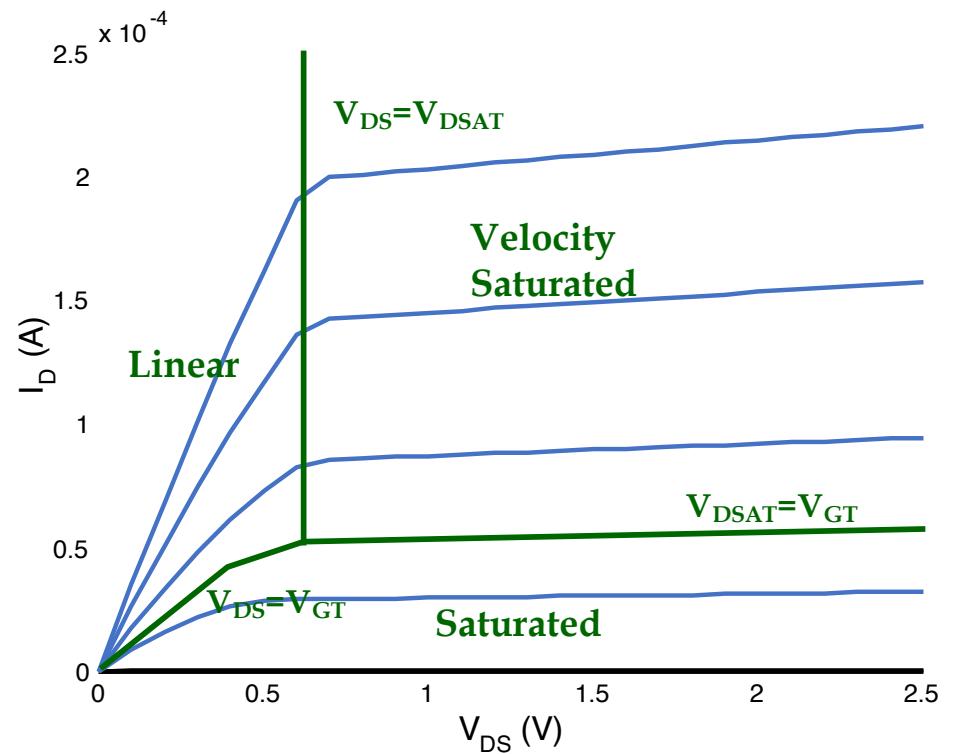
- Current I_D is a function of V_{GS} , V_{DS} , and V_T
- Transistor conductivity / resistance is controlled by terminal voltages and transistor size

Let

- $V_{GT} = V_{GS} - V_T$
- $V_{min} = \min(V_{GT}, V_{DS}, V_{DSAT})$
- $V_T = V_{T0} + \gamma(\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|})$

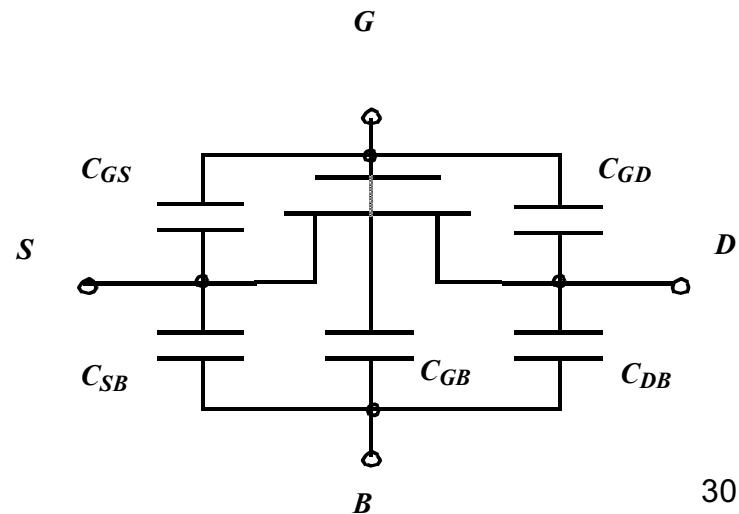
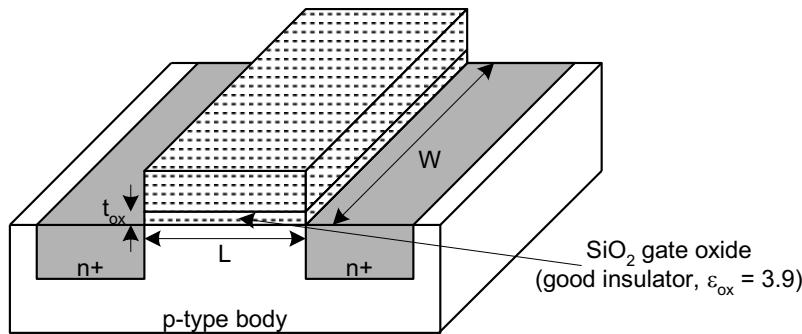
Then

$$I_{DS} = \begin{cases} 0 & V_{GT} \leq 0 \\ k' \left(\frac{W}{L} \right) \left(V_{GT} V_{min} - \frac{V_{min}^2}{2} \right) (1 + \lambda V_{DS}) & V_{GT} > 0 \end{cases}$$

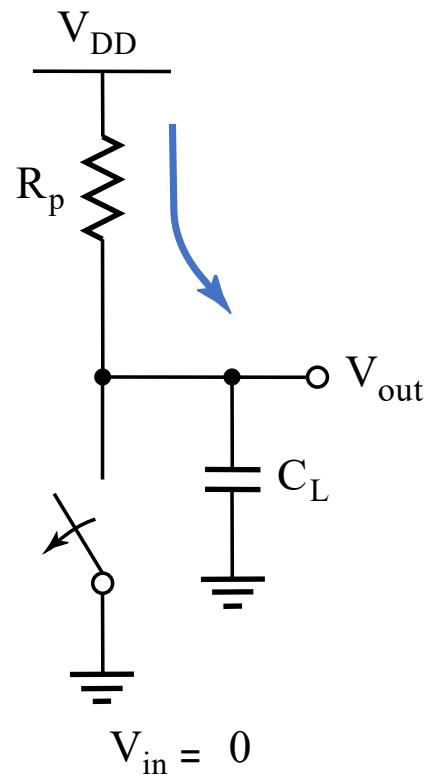
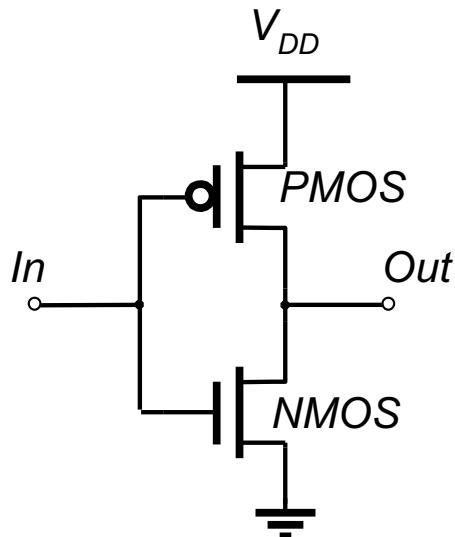


Capacitance

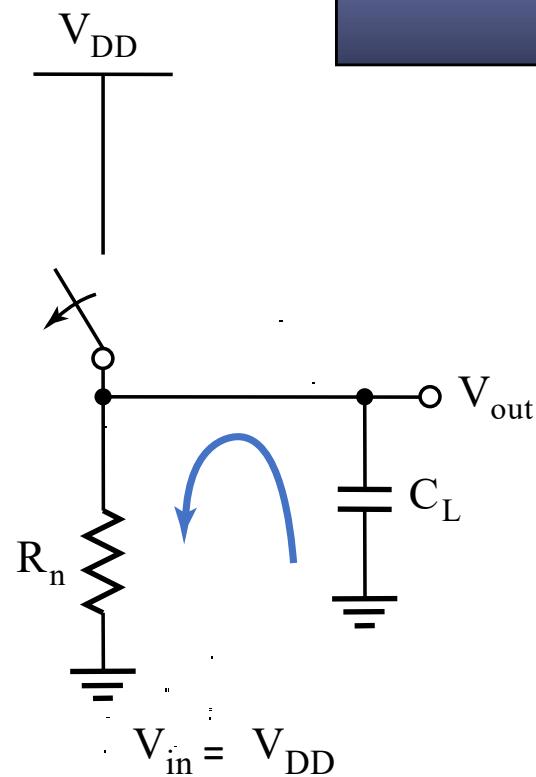
- Any two conductors separated by an insulator have capacitance
- MOS transistor gate to channel capacitor is very important
 - Creates channel charge necessary for operation
- Source and drain have capacitance to body
 - Across reverse-biased diodes
 - Called junction or diffusion capacitance; associated with source/drain diffusion and p-n junction
- MOS transistor cap is proportional to the transistor width (W) given length (L): $C \propto WL$



CMOS Inverter: Transient Response



(a) Low-to-high



(b) High-to-low

$$t_p = f(R_{on} \cdot C_L)$$

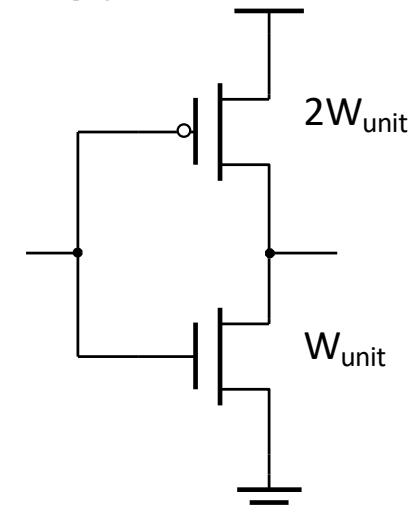
$$= 0.69 R_{on} C_L$$

A ‘unit’ size inverter

- A ‘unit’ size inverter for 250nm technology
 - $L_N = L_P = L_{\text{unit}} = 0.25\mu\text{m}$
 - $W_N = L_{\text{unit}}$, $W_P = 2L_{\text{unit}}$
 - NMOS internal cap $C_{dn} = C_{\text{unit}}$
 - PMOS internal cap $C_{dp} = 2C_{\text{unit}}$
 - Internal Cap $C_{int} = C_{dn} + C_{dp} = 3C_{\text{unit}}$
 - Equivalent R: $R_n = R_p = R_{\text{unit}}$
 - Same pull-up and pull-down currents

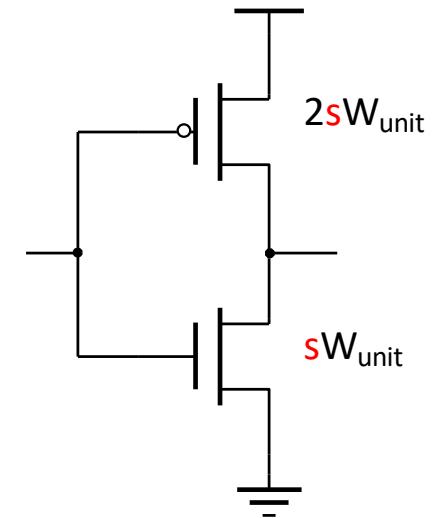
$$t_{pHL} = 0.69R_{\text{unit}}(C_{int} + C_L) = 0.69R_{\text{unit}}(3C_{\text{unit}} + C_L)$$

$$t_{pLH} = 0.69R_{\text{unit}}(C_{int} + C_L) = 0.69R_{\text{unit}}(3C_{\text{unit}} + C_L)$$



Inverter Scaling

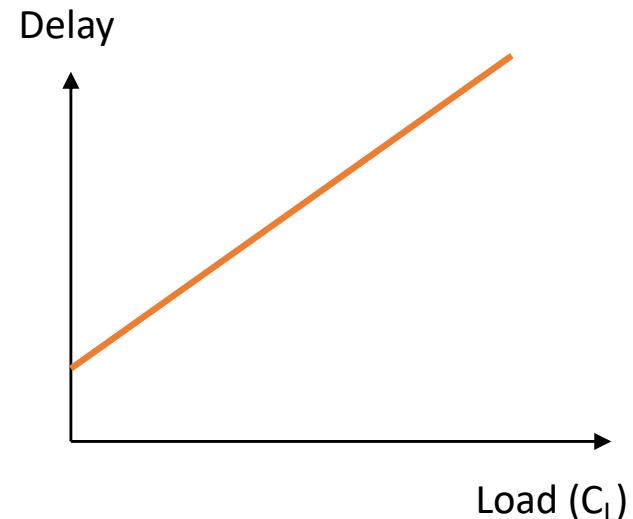
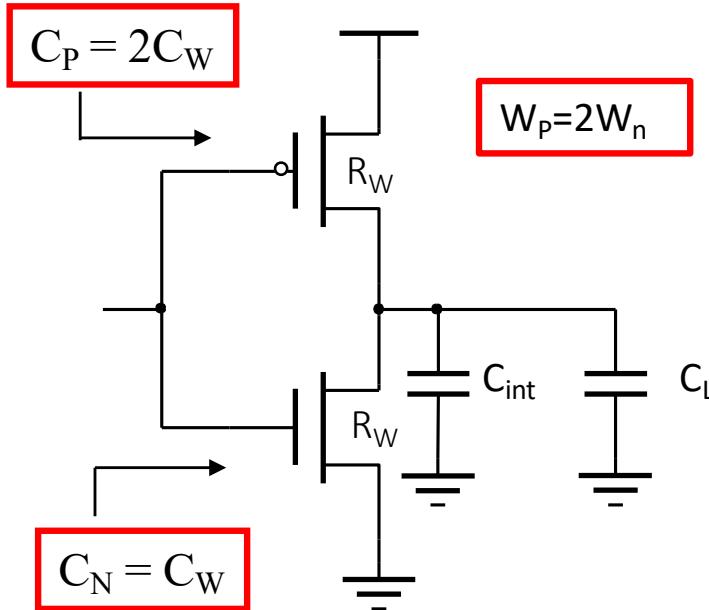
- Scaling factor s
 - $L_N = L_P = L_{\text{unit}} = 0.25\mu\text{m}$
 - $W_N = sL_{\text{unit}}, W_P = 2sL_{\text{unit}}$
 - NMOS internal cap $C_{dn} = sC_{\text{unit}}$
 - PMOS internal cap $C_{dp} = 2sC_{\text{unit}}$
 - Internal cap $C_{\text{int}} = C_{dn} + C_{dp} = 3sC_{\text{unit}}$
 - Equivalent R: $R_n = R_p = R_{\text{unit}}/s$



$$t_{pHL} = 0.69R_{\text{unit}}/s (sC_{\text{int}} + C_L) = 0.69R_{\text{unit}}/s (3sC_{\text{unit}} + C_L)$$

$$t_{pLH} = 0.69R_{\text{unit}}/s (sC_{\text{int}} + C_L) = 0.69R_{\text{unit}}/s (3sC_{\text{unit}} + C_L)$$

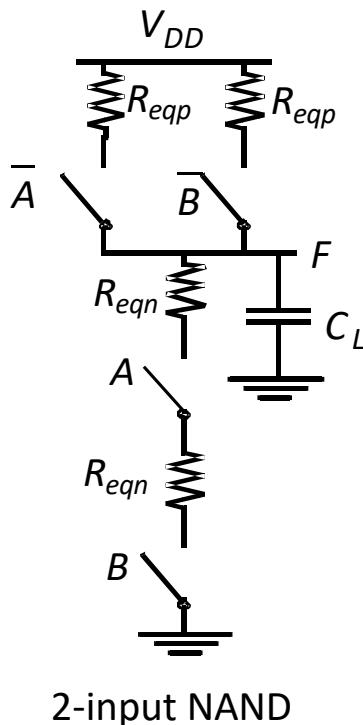
Inverter with Load



$$\begin{aligned} \text{Delay} &= 0.69R_W(C_{\text{int}} + C_L) = 0.69R_WC_{\text{int}} + 0.69R_WC_L = 0.69R_WC_{\text{int}}(1 + C_L/C_{\text{int}}) \\ &= \text{Delay (Internal)} + \text{Delay (Load)} \end{aligned}$$

What's the desired size of this inverter to minimize the delay given C_L ?
 What about the size of the previous stage that drives the input of this inverter?

Analysis of Propagation Delay



**1. Determine “Worst Case Input” transition
 (Delay depends on input values)**

2. Example: t_{pLH} for 2-input NAND

- Worst case when only ONE PMOS Pulls up the output node
- For 2 PMOS devices in parallel, the resistance is lower

$$t_{pLH} = 0.69 * 0.5 * R_{eqp} C_L$$

3. Example: t_{pHL} for 2-input NAND

- Worst case : TWO NMOS in series

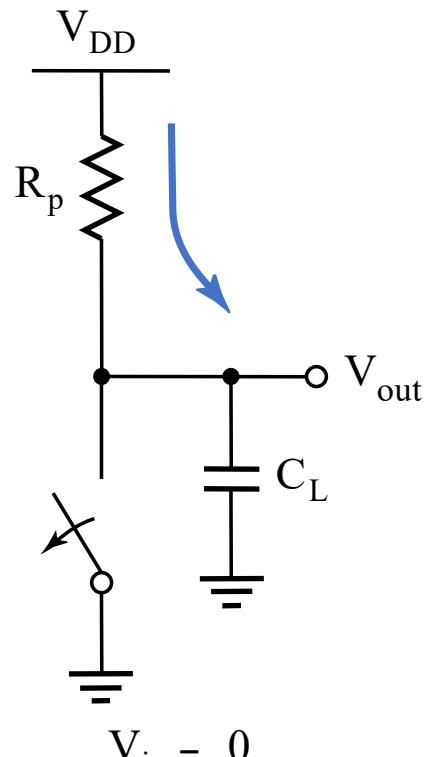
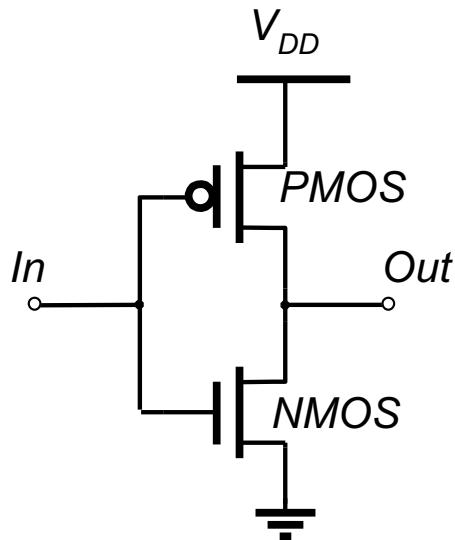
$$t_{pHL} = 0.69(2R_{eqn})C_L$$



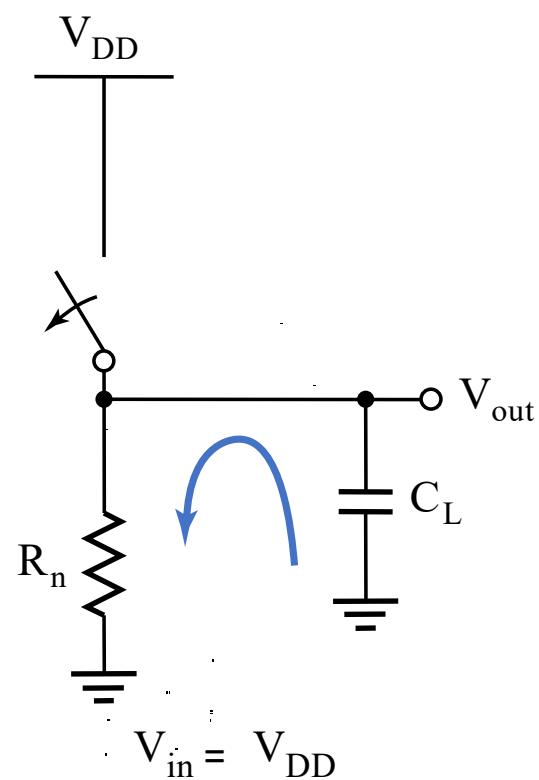
Power Components in CMOS

- Dynamic power consumption (dominant)
 - Charging and discharging capacitances through devices
- Short-circuit power (manageable)
 - Direct path between V_{dd} and ground during switching transients
- Static power (exponentially increasing)
 - Leakage currents exist and consume power when not switching

CMOS Inverter: Energy Consumption



(a) Low-to-high



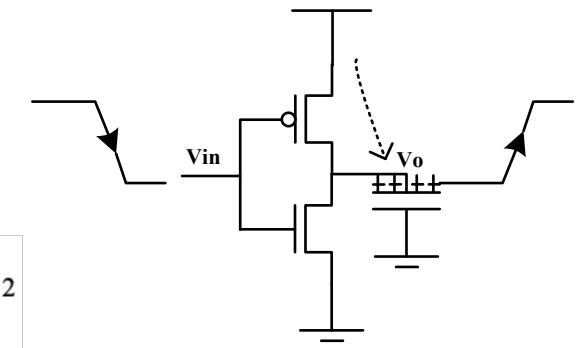
(b) High-to-low

Dynamic Power Consumption (1/2)

- Input 1 → 0

– Energy drawn from power supply:

$$E_{\text{supply}} = \int_0^{\infty} P(t) \cdot dt = \int_0^{\infty} V_{dd} i(t) \cdot dt = \int_0^{V_{dd}} V_{dd} C \cdot \frac{dV_O}{dt} \cdot dt = V_{dd} C \cdot \int_0^{V_{dd}} dV_O = CV_{dd}^2$$



– Energy consumed by PMOS:

$$E_{\text{PMOS}} = \int_0^{\infty} P(t) \cdot dt = \int_0^{\infty} (V_{dd} - V_O) \cdot i_p(t) \cdot dt = C \cdot \int_0^{V_{dd}} (V_{dd} - V_O) \cdot dV_O = \frac{1}{2} CV_{dd}^2$$

– Energy in the capacitor:

$$E_{\text{CAP}} = \int_0^{\infty} V_O \cdot i_C(t) \cdot dt = C \cdot \int_0^{\infty} V_O \cdot \frac{dV_O}{dt} \cdot dt = C \cdot \int_0^{V_{dd}} V_O \cdot dV_O = \frac{1}{2} CV_{dd}^2$$

– Power:

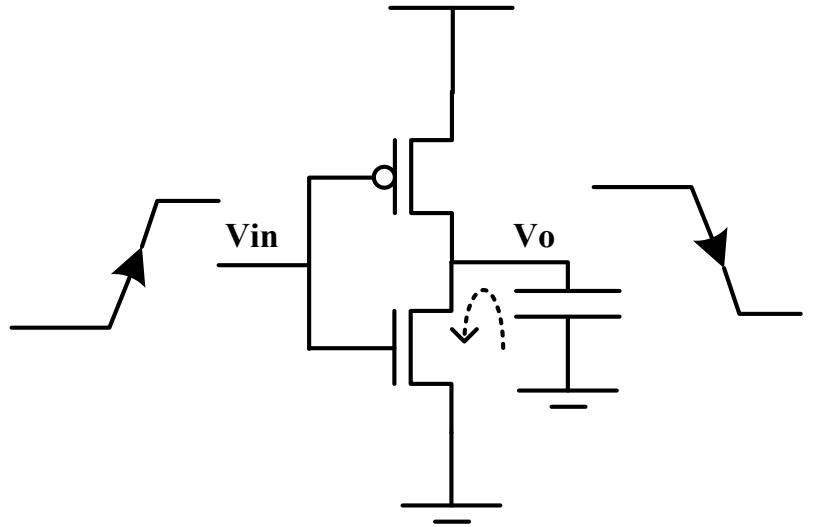
$$P = f \cdot E_{\text{PMOS}} = \frac{1}{2} f CV_{dd}^2$$



Dynamic Power Consumption (2/2)

- Input 0→1
 - Energy drawn from supply: 0
 - Energy consumed by NMOS equals the energy stored on capacitance:

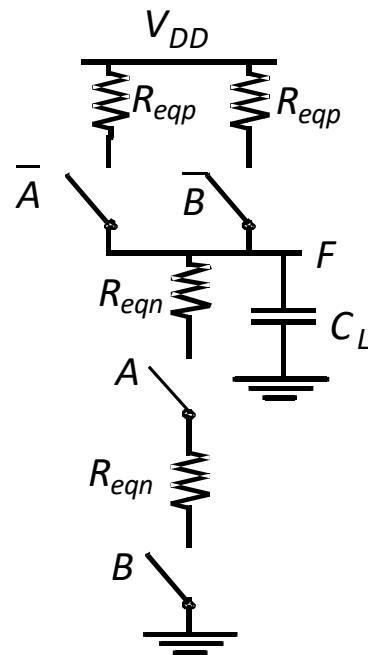
$$E_{NMOS} = \int V_o i(t) \cdot dt = \frac{1}{2} C V_{dd}^2$$



- Power:

$$P = f \cdot E_{NMOS} = \frac{1}{2} C V_{dd}^2 f$$

CMOS Gate Energy / Power

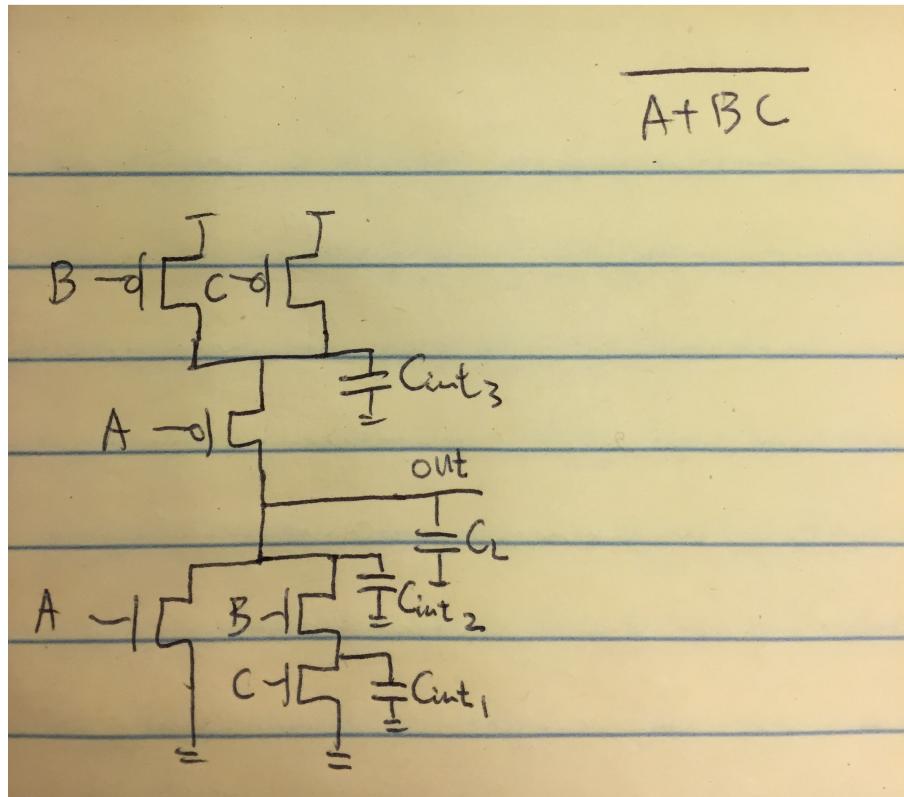


2-input NAND

- Energy consumption for output $0 \rightarrow 1$
- Energy consumption for output $1 \rightarrow 0$



CMOS Gate Energy Consumption





Energy vs. Power

Each rising transition on C_L requires $C_L * V_{dd}^2$ of energy BUT 1/2 of this energy is lost (to heat) while the other half is *stored* on the capacitor

For every transition, $C * V_{dd}^2 / 2$ of energy is **consumed**

For every period (both a $L \rightarrow H$ and $H \rightarrow L$ transition on C_L), $C * V_{dd}^2$ of energy is consumed from the supply

Then, power = rate of energy consumption so:

$$\text{Energy} = P_{dyn} T_{sw} = C_L * V_{dd}^2$$

$$\rightarrow P_{dyn} = C_L * V_{dd}^2 * f_{sw}$$

Where f_{sw} is the frequency with which C_L switches ($L \rightarrow H + H \rightarrow L$)



Switching Activity, α_{sw}

Let $f_{sw} = \alpha_{sw} * f_{clock}$ since we usually know clock frequency of a design (e.g. 3.6 GHz Core i7)

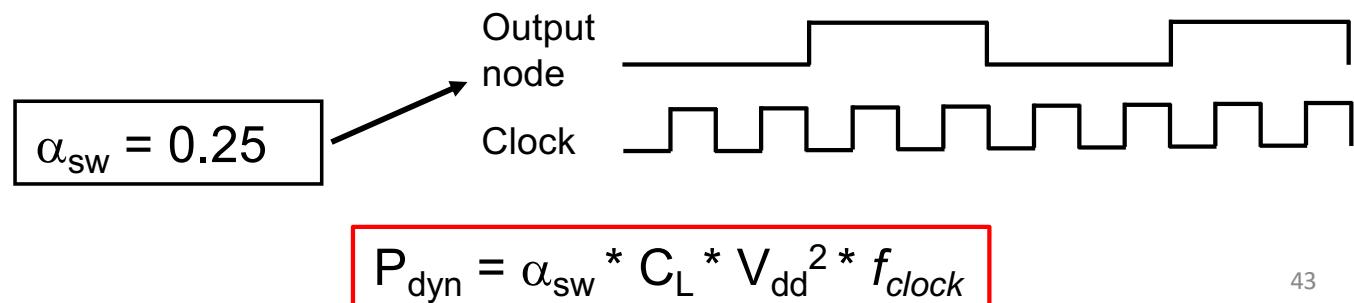
$$0 < \alpha_{sw} < 1$$

For $\alpha_{sw} = 0$, the circuit never switches so no dynamic power is consumed

For $\alpha_{sw} = 1$, the node switches as often as the clock (the circuit cannot switch more often than this) so $f_{sw} = f_{clock}$

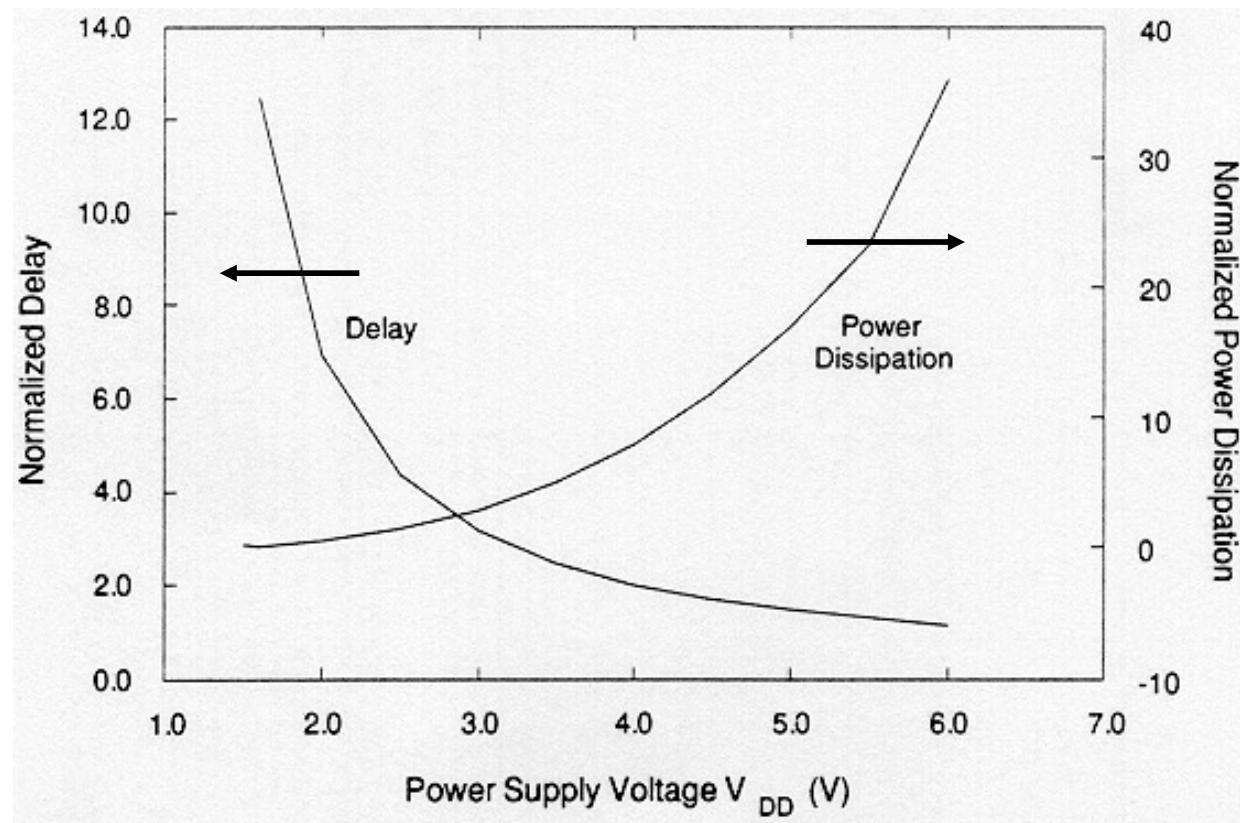
Most cases \rightarrow somewhere in between

Lower α_{sw} = lower power

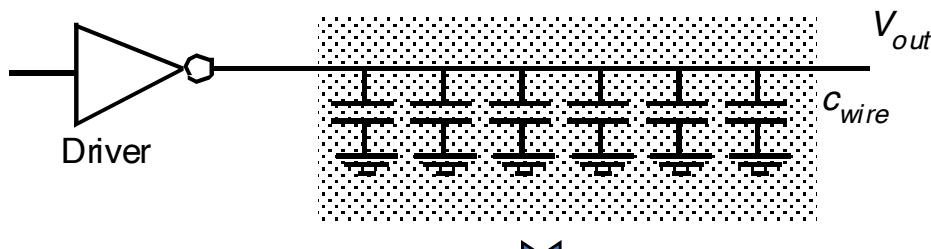
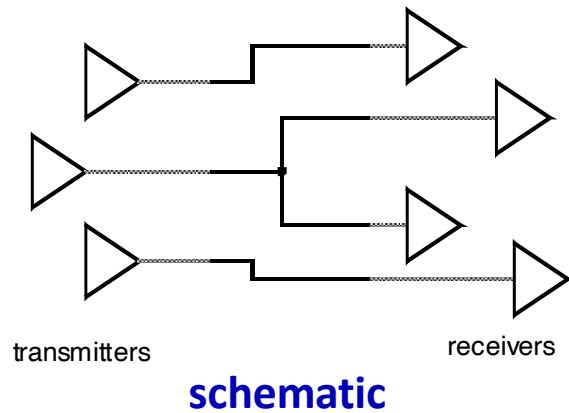


Fundamental Tradeoff: Power vs. Delay

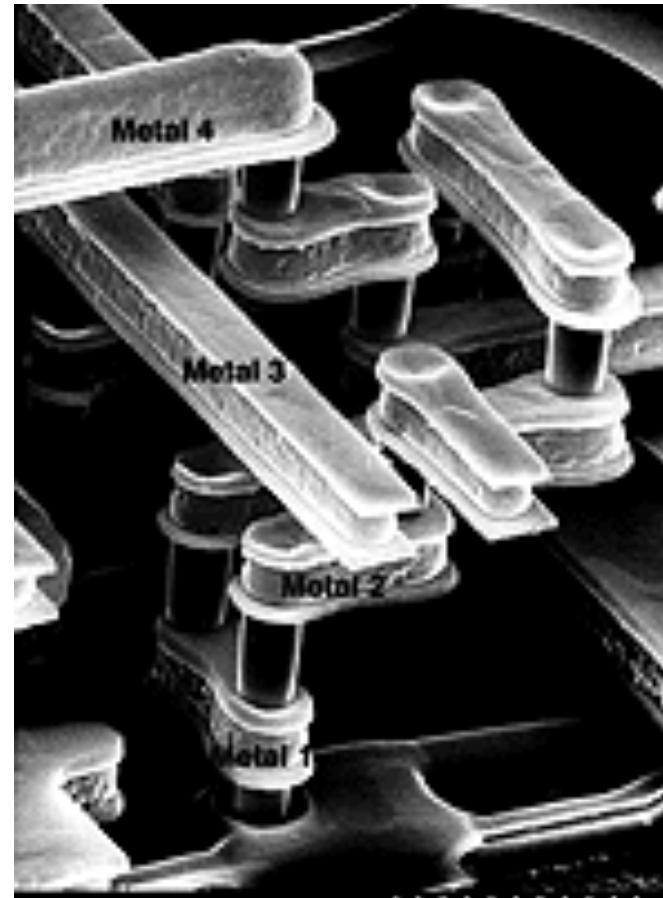
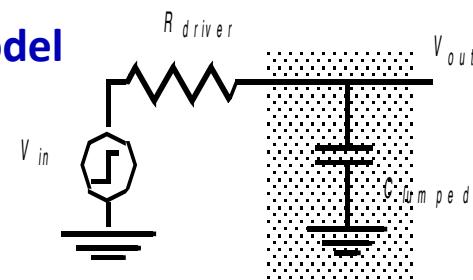
Power and delay tradeoff via voltage scaling



The Wire



Lumped model

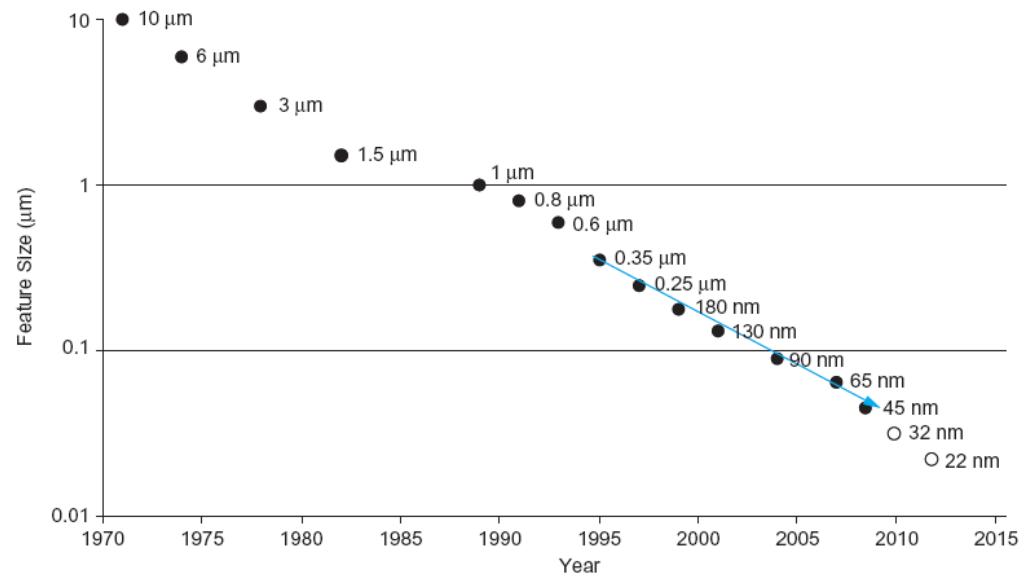


physical

Capacitance and resistance are proportional to wire length

Process Technology Scaling

- The only constant in VLSI is constant change
- Feature size shrinks by 30% every 2-3 years
 - Transistors become cheaper
 - Transistors become faster and lower power
 - Wires do not improve (and may get worse)
- Scale factor S
 - Typically $S = \sqrt{2}$
 - Technology nodes





Process Technology Scaling

newer technology

$$S, U \geq 1$$

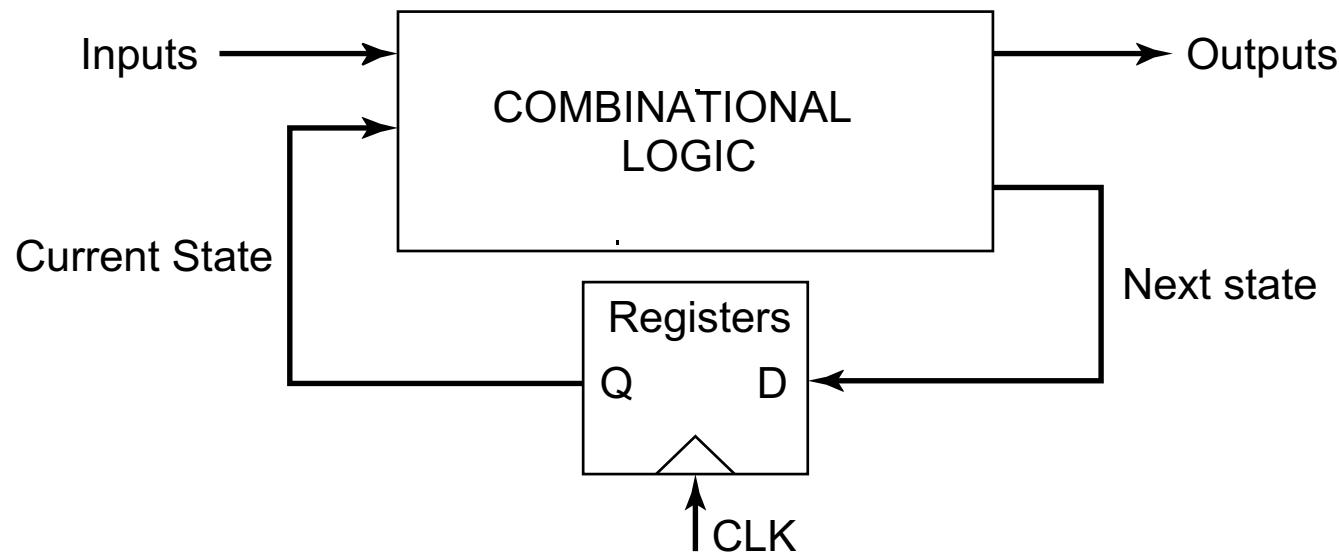
S: feature size scaling, U: voltage scaling

Parameter	Sensitivity	Dennard Scaling (U=S)	General Scaling (U,S)	Fixed Voltage Scaling (U = 1)
L: Length		1/S	1/S	1/S
W: Width		1/S	1/S	1/S
t_{ox} : gate oxide thickness		1/S	1/S	1/S
V_{DD} : supply voltage		1/S	1/U	1
V_t : threshold voltage		1/S	1/U	1
NA: substrate doping		S	S^2/U	S^2
β	$W/(Lt_{ox})$	S	S	S
I_{on} : ON current	$\beta(V_{DD}-V_t)V_{min}$	1/S	1/U	1
R: effective resistance	V_{DD}/I_{on}	1	1	1
C: gate capacitance	WL/t_{ox}	1/S	1/S	1/S
τ : gate delay	RC	1/S	1/S	1/S
f: clock frequency	$1/\tau$	S	S	S
E: switching energy / gate	CV_{DD}^2	$1/S^3$	$1/(SU^2)$	1/S
P: switching power / gate	Ef	$1/S^2$	$1/U^2$	1
A: area per gate	WL	$1/S^2$	$1/S^2$	$1/S^2$
Switching power density	P/A	1	S^2/U^2	S^2
Switching current density	I_{on}/A	S	S^2/U	S^2

47



Sequential Logic



- *Combinational logic*
 - output depends on current inputs
- *Sequential logic*
 - output depends on current and previous inputs
 - Requires separating previous, current, future data



Naming Conventions

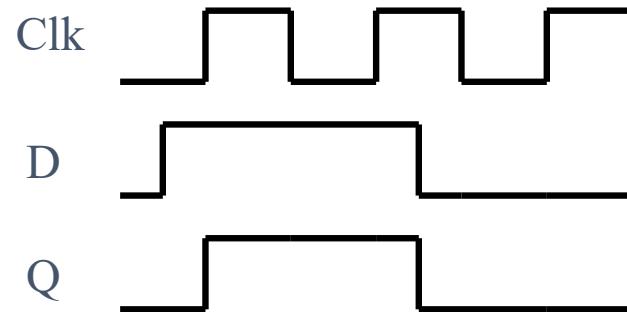
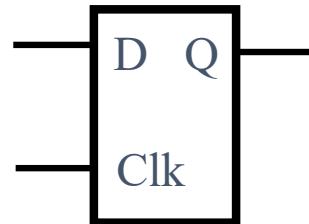
- In our text:
 - a latch is **level sensitive**
 - a register is **edge-triggered**
- There are many different naming conventions
 - For instance, many books call edge-triggered elements **flip-flops**
 - Registers and flip-flop terms are often used interchangeably therefore



Latch versus Register

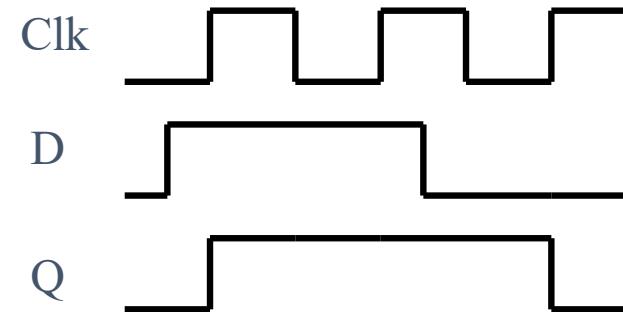
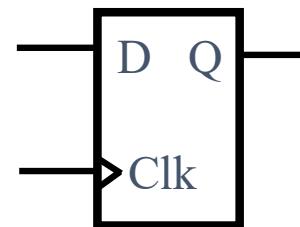
- ❑ Latch

stores data when
clock is low



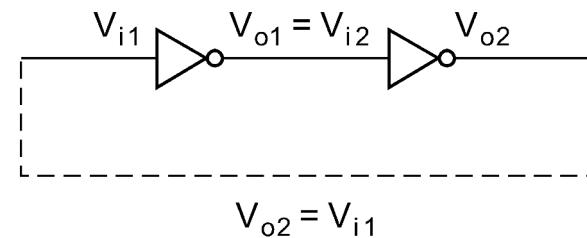
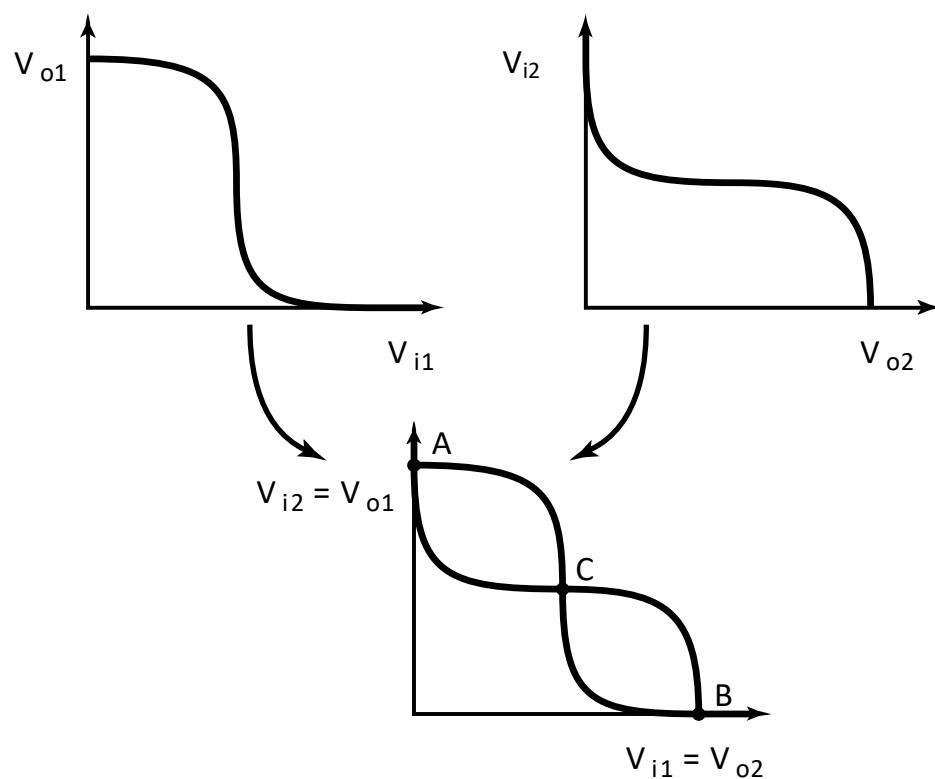
- ❑ Register

stores data when
clock rises



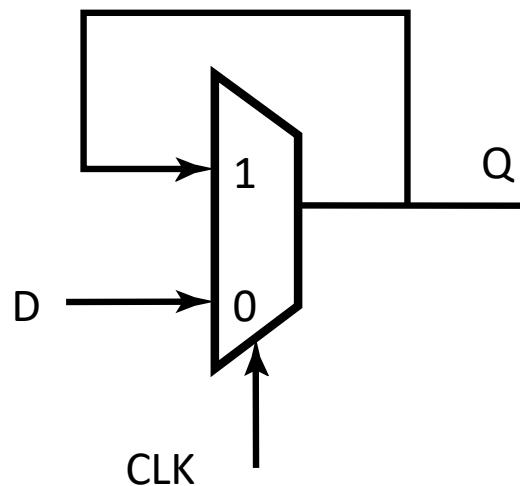


Bit storage mechanism: Bi-Stability



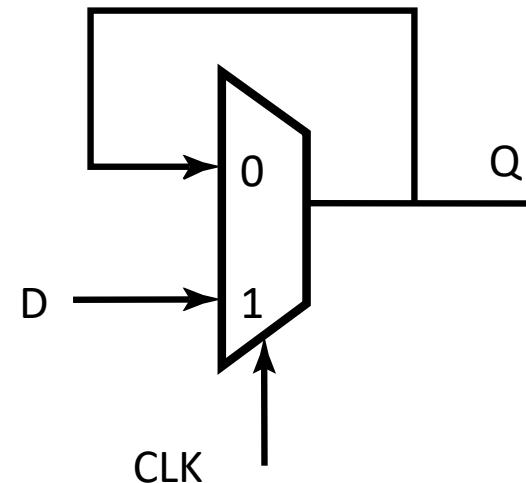
Mux-Based Latches

Negative latch
(transparent when CLK= 0)



$$Q = \overline{Clk} \cdot D + Clk \cdot Q$$

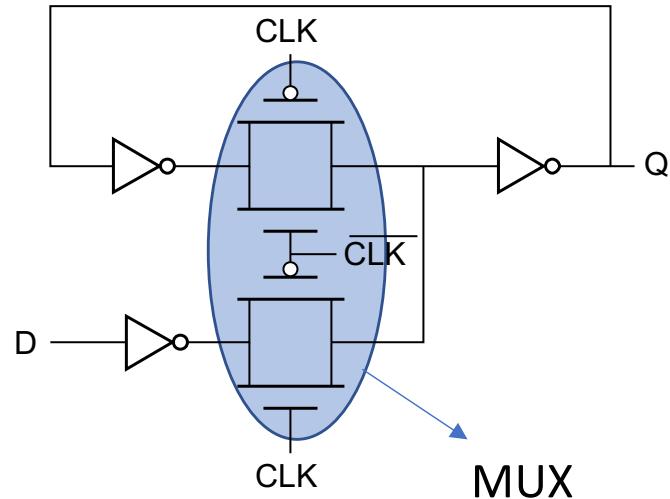
Positive latch
(transparent when CLK= 1)



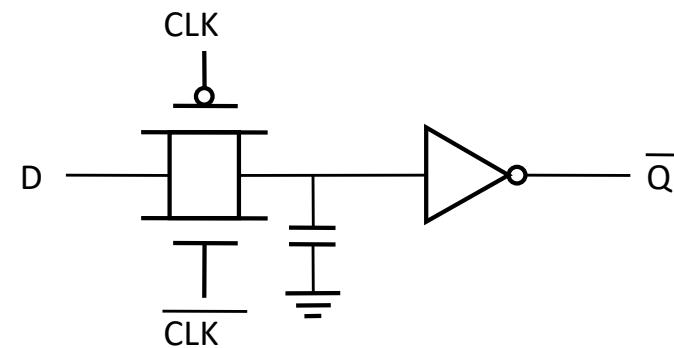
$$Q = Clk \cdot D + \overline{Clk} \cdot Q$$

Storage Mechanisms

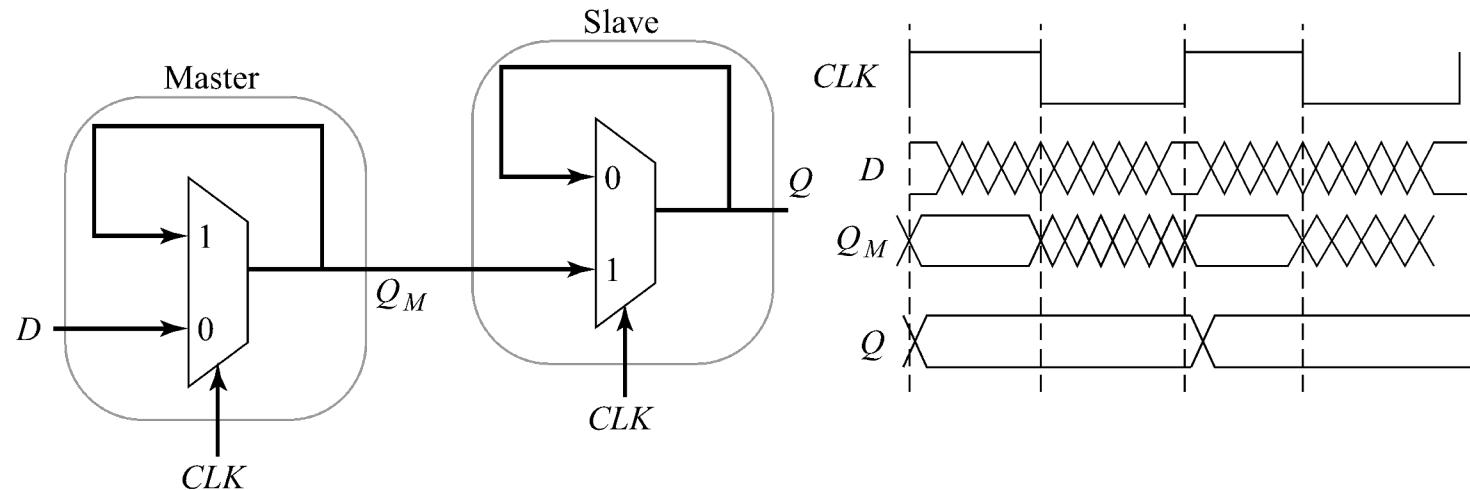
Static



Dynamic (charge-based)



Master-Slave (Edge-Triggered) Register

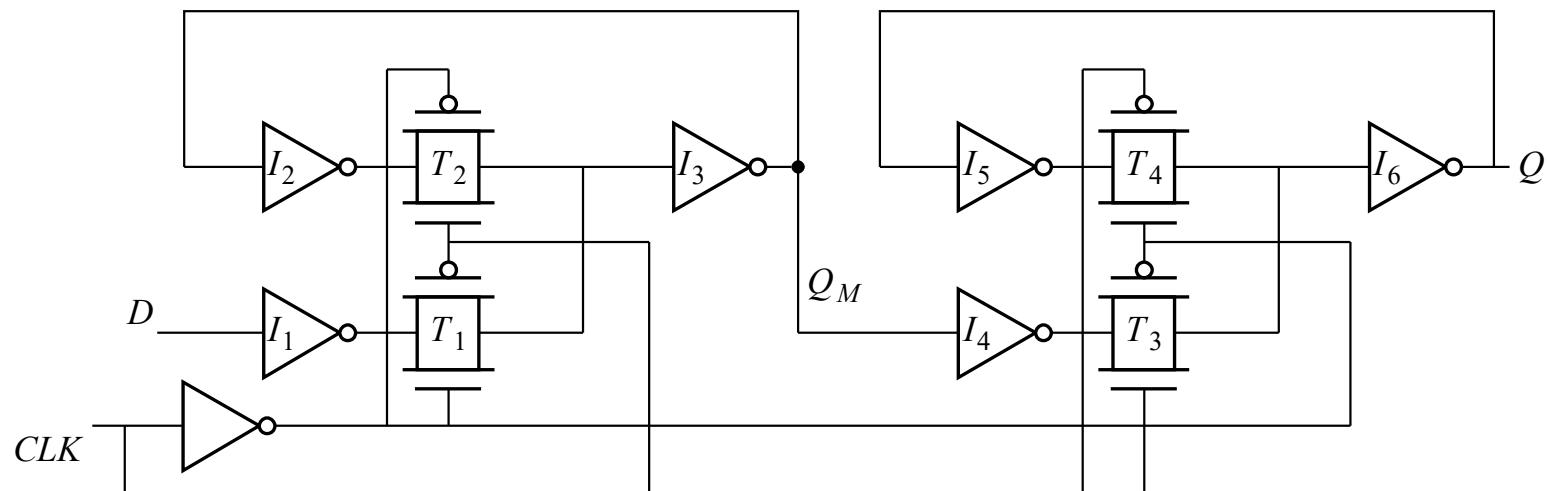


Two opposite latches trigger on edge
 Also called master-slave latch pair



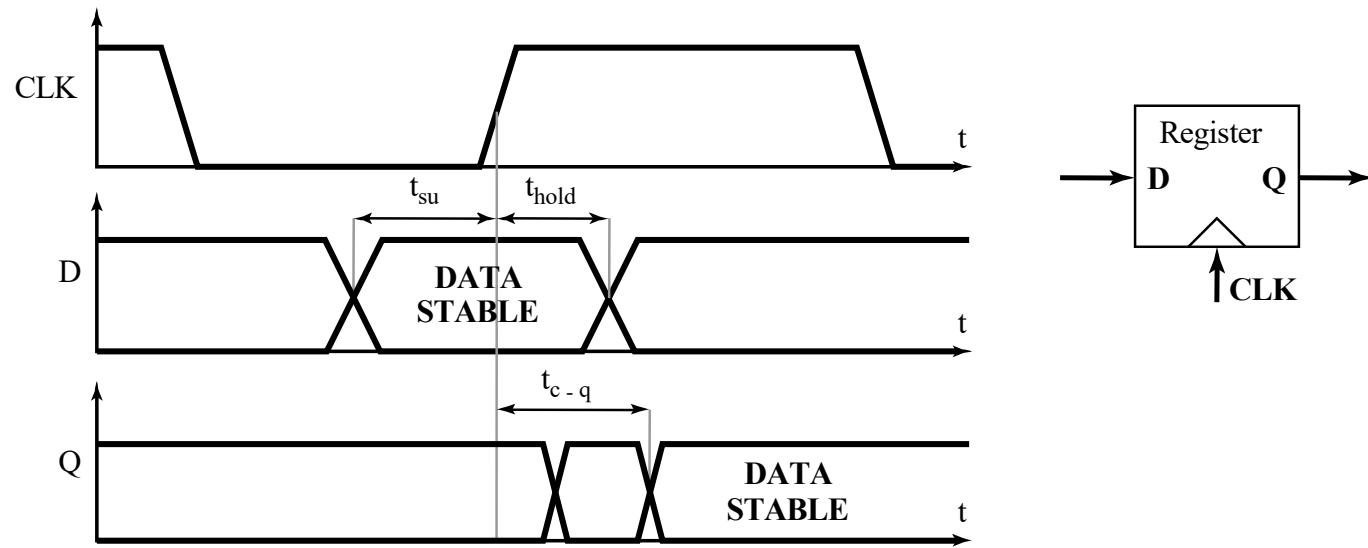
Master-Slave Register

Multiplexer-based latch pair



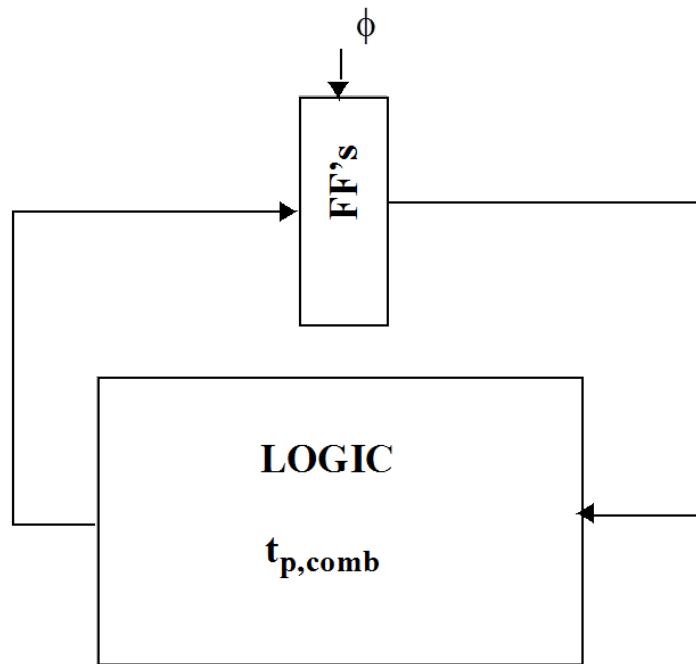


Register Timing Definitions





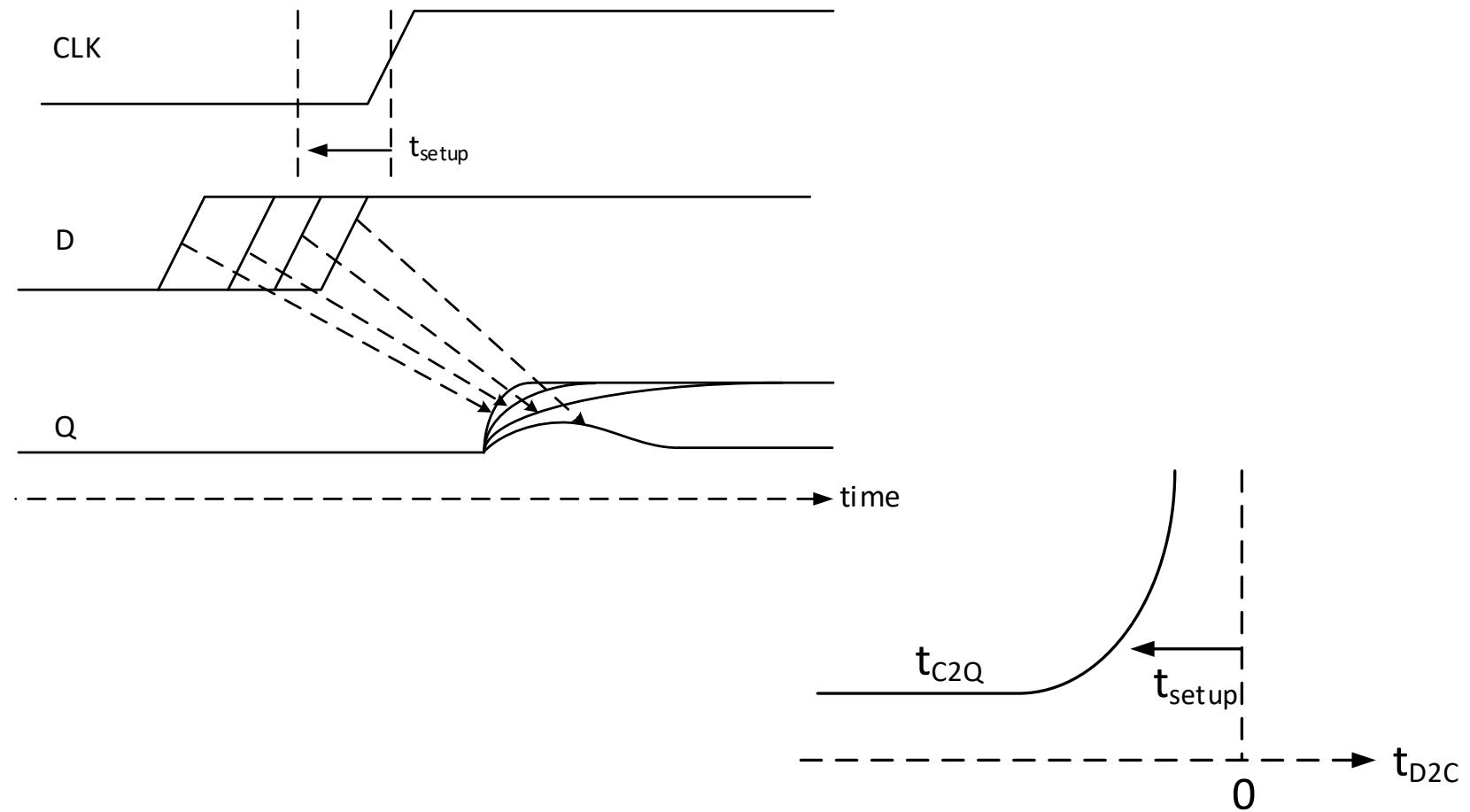
Maximum Clock Frequency



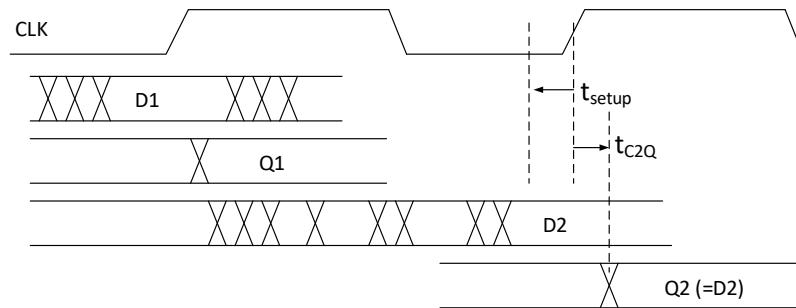
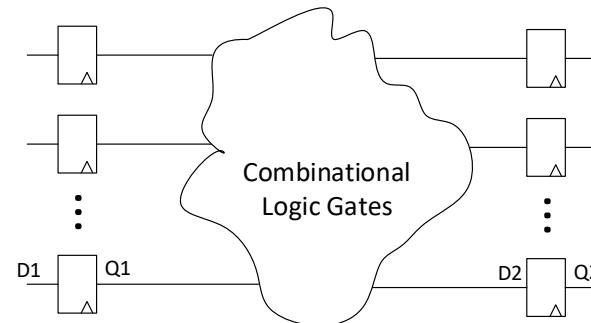
$$t_{\text{clk-Q}} + t_{p,\text{comb}} + t_{\text{setup}} = T$$



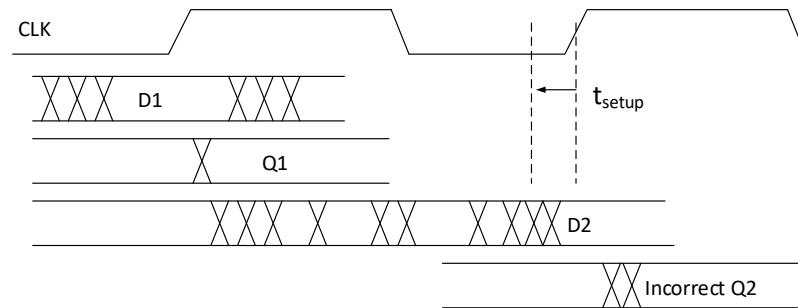
Setup Time



Setup Time Constraint

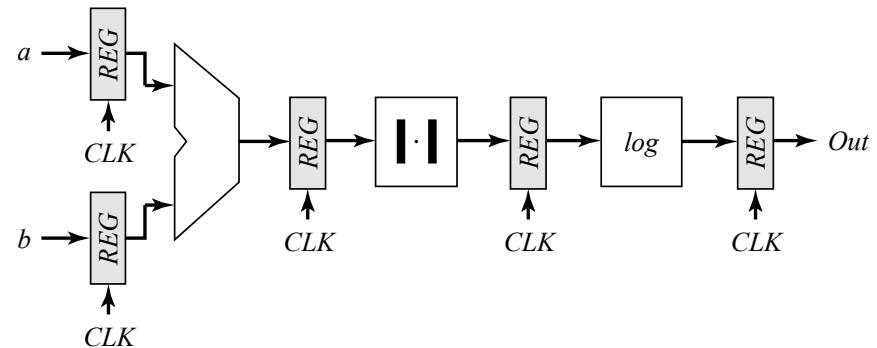
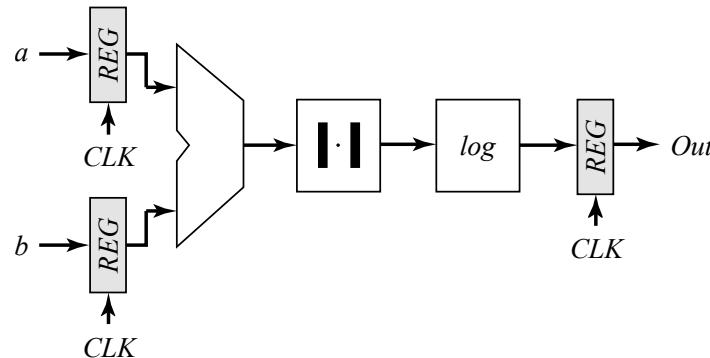


Setup time constraint satisfied



Setup time constraint violated

Pipelining



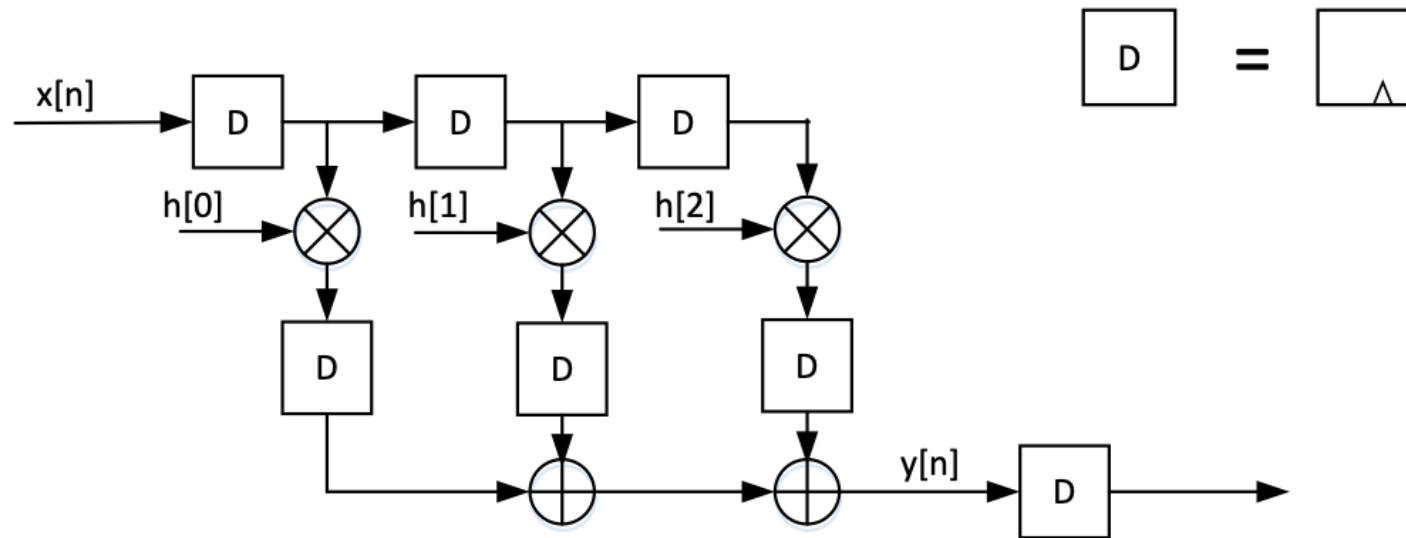
Reference

Clock Period	Adder	Absolute Value	Logarithm
1	$a_1 + b_1$		
2	$a_2 + b_2$	$ a_1 + b_1 $	
3	$a_3 + b_3$	$ a_2 + b_2 $	$\log(a_1 + b_1)$
4	$a_4 + b_4$	$ a_3 + b_3 $	$\log(a_2 + b_2)$
5	$a_5 + b_5$	$ a_4 + b_4 $	$\log(a_3 + b_3)$

Pipelined



Pipeline Example



$$y[n] = h[2]x[n-2] + h[1]x[n-1] + h[0]x[n]$$

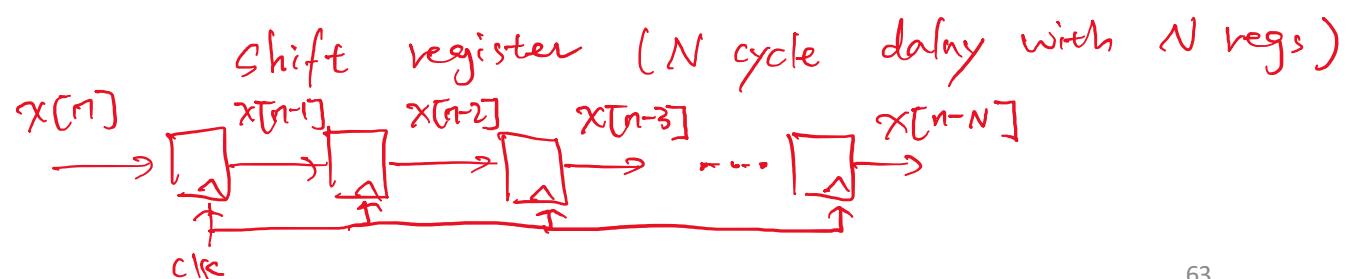
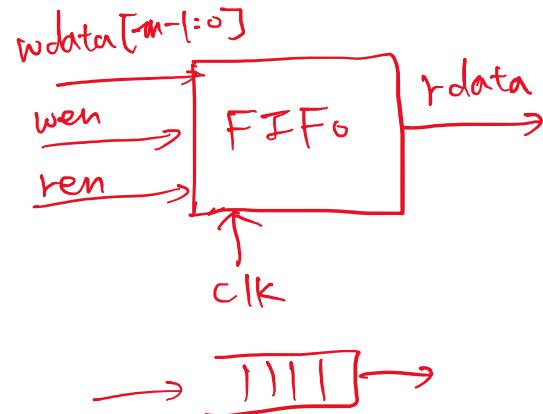
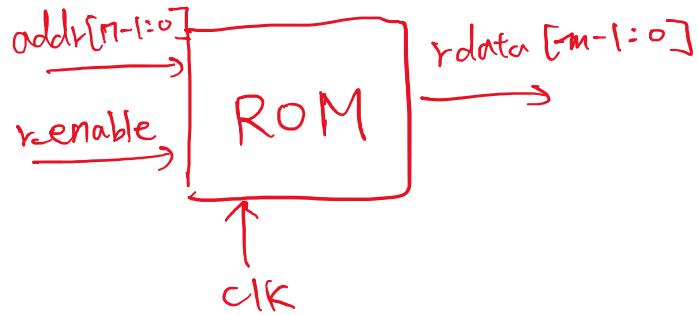
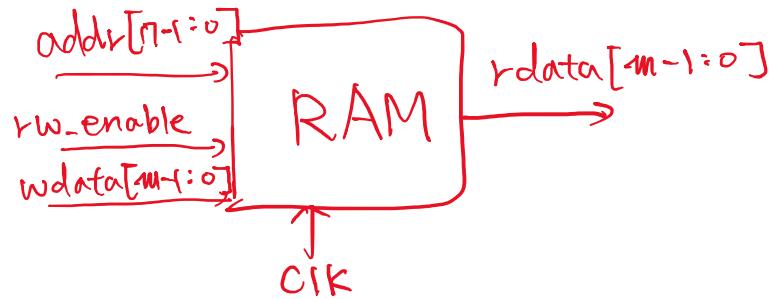


Semiconductor Memory Classification

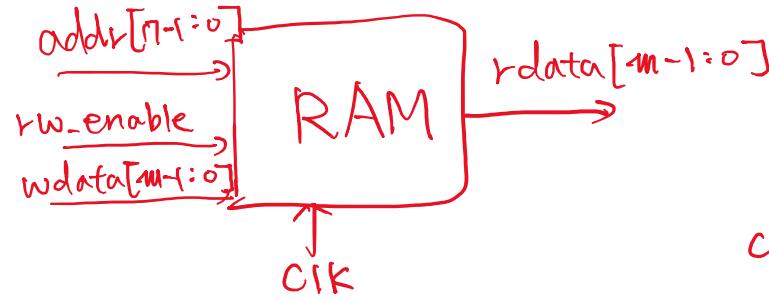
Read-Write Memory		Non-Volatile Read-Write Memory	Read-Only Memory
Random Access	Non-Random Access	EPROM E ² PROM FLASH	Mask-Programmed Read Only Memory (PROM)
SRAM	FIFO LIFO Shift Register CAM		

Key design metric: AREA

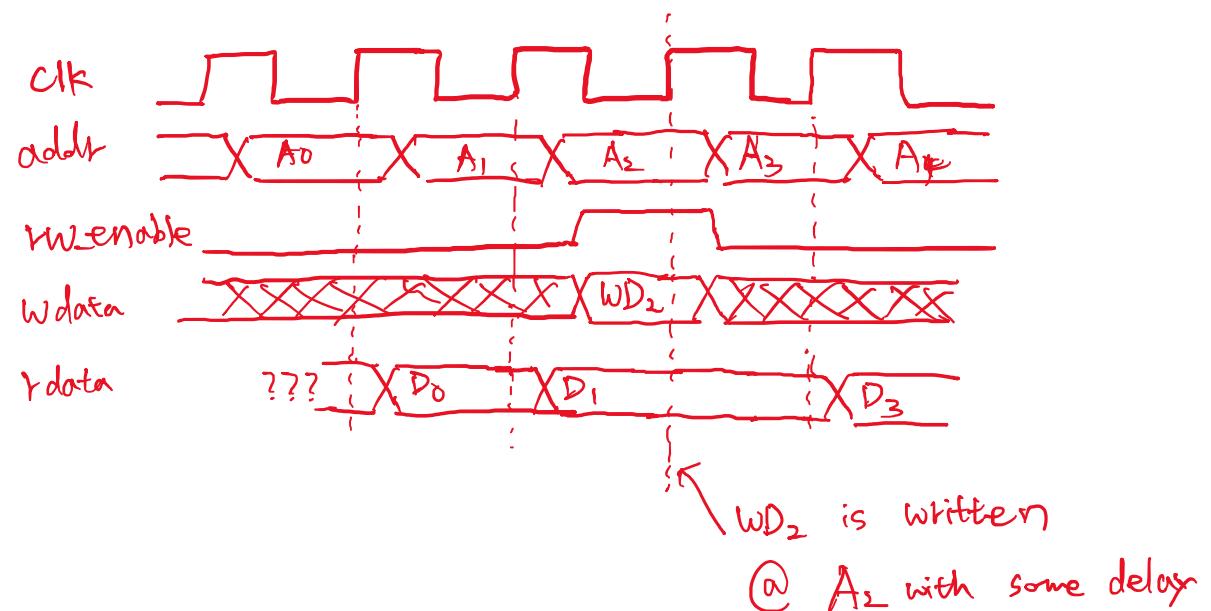
Memory Interface



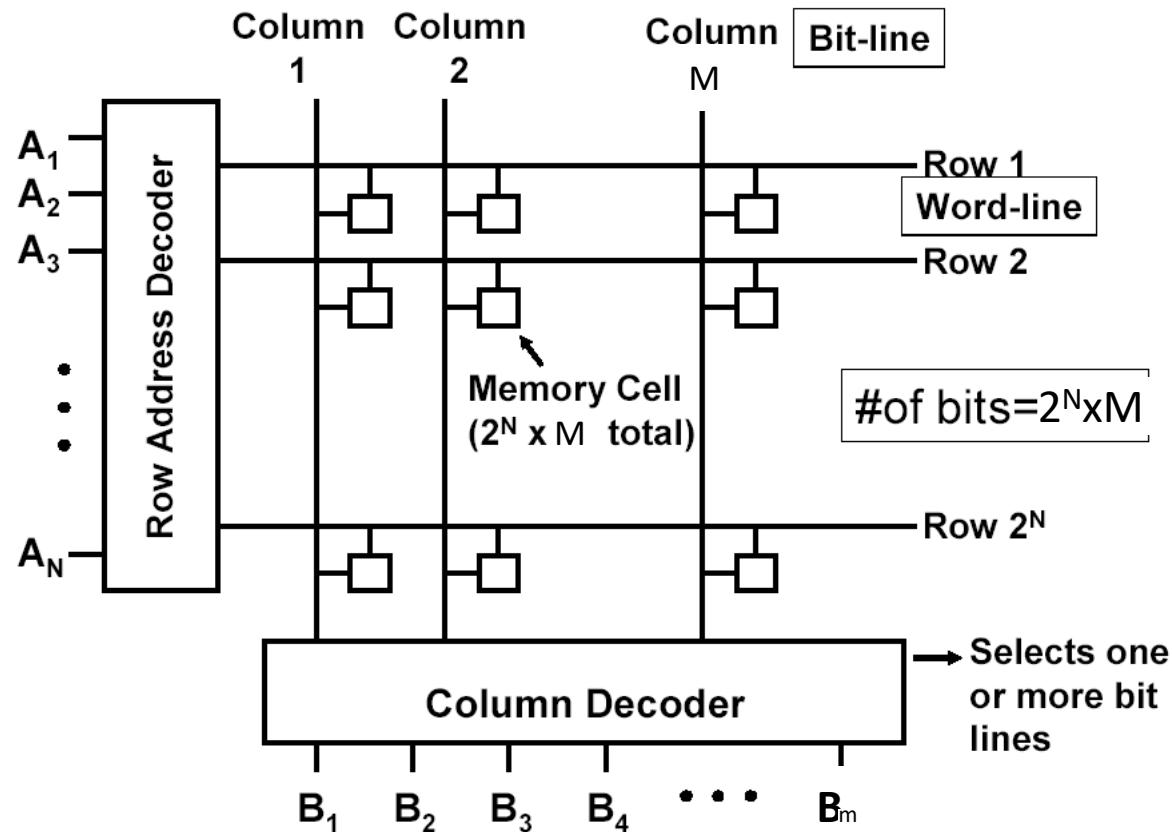
Memory Timing: Definitions



$$rw_enable = \begin{cases} 0 & \text{when read} \\ 1 & \text{when write} \end{cases}$$

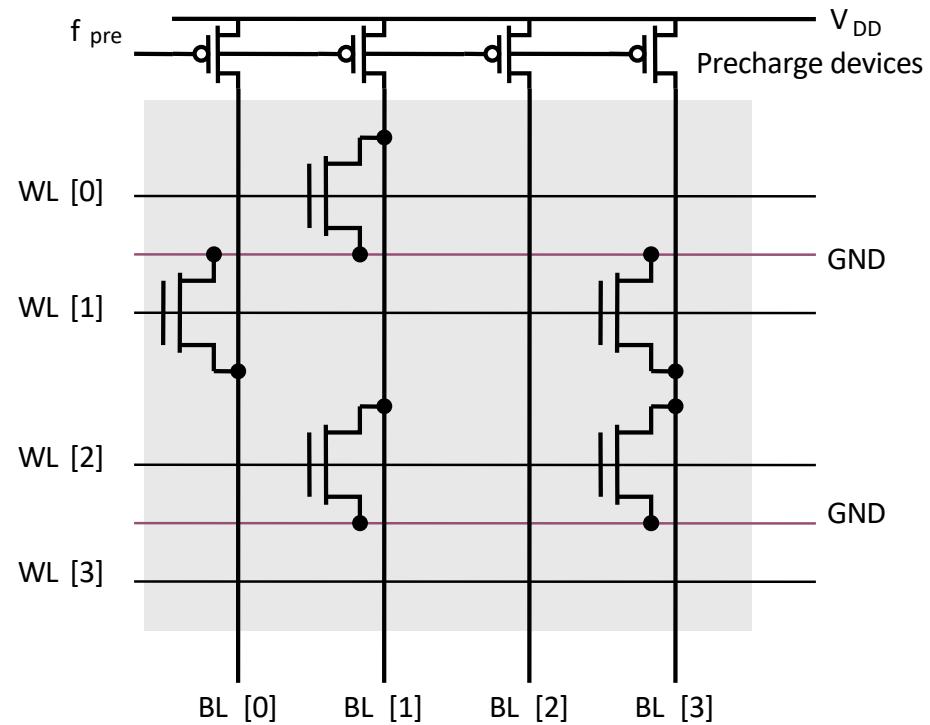


General Memory Structure





Precharged MOS NOR ROM





Address Decoders

M:2^M decoder:

Collection of 2^M M-input complex logic gates
Organized in regular and dense fashion

here
M is # of bits for address
2^M is # of rows

(N)AND Decoder

$$WL_0 = \overline{A_0} \overline{A_1} \overline{A_2} \overline{A_3} \overline{A_4} \overline{A_5} \overline{A_6} \overline{A_7} \overline{A_8} \overline{A_9}$$

W₀ = 1 when Addr = 000...0

$$WL_{511} = \overline{A_0} \overline{A_1} \overline{A_2} \overline{A_3} \overline{A_4} \overline{A_5} \overline{A_6} \overline{A_7} \overline{A_8} \overline{A_9}$$

W₅₁₁ = 1 when Addr = 111...1

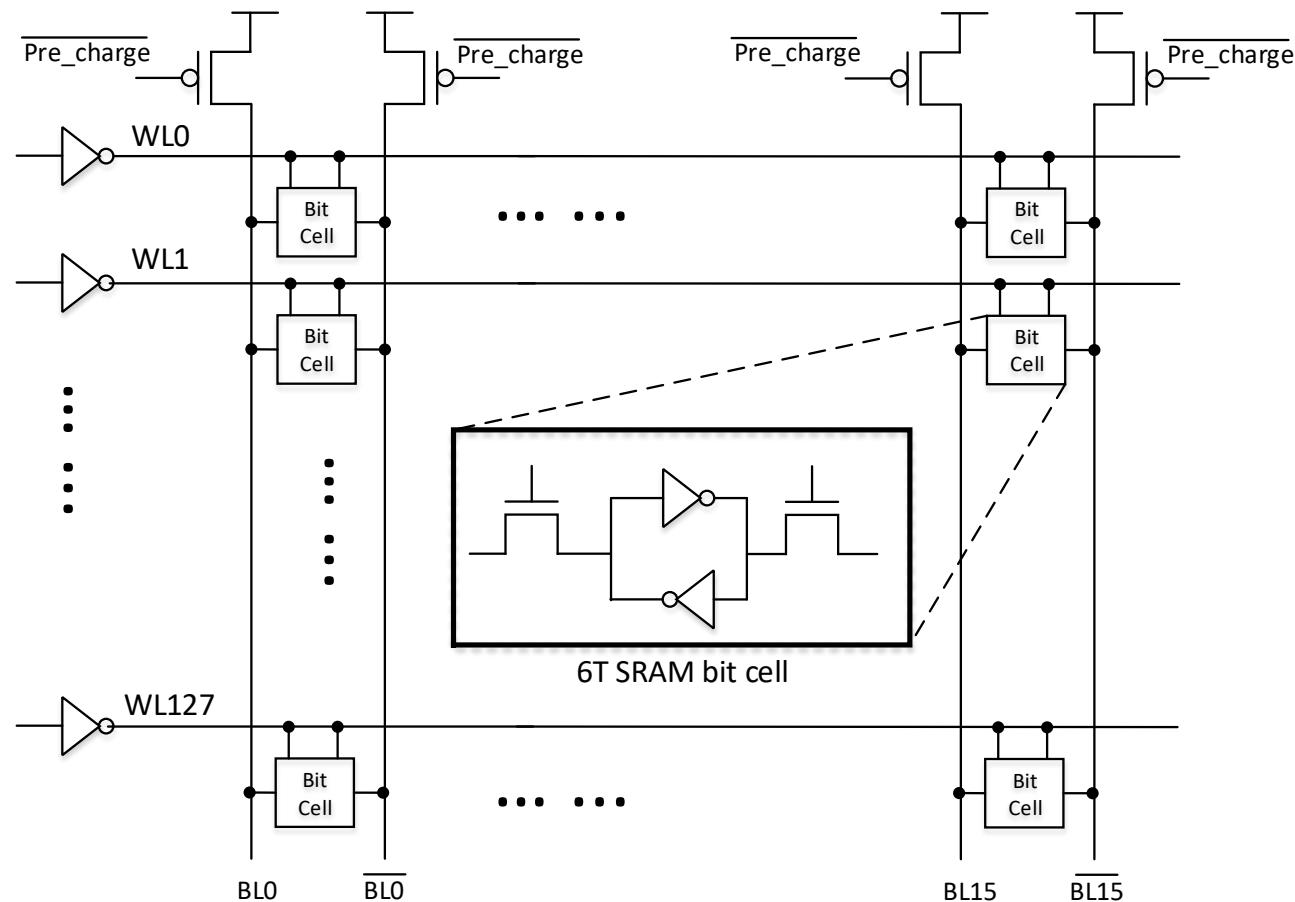
NOR Decoder

$$WL_0 = \overline{A_0 + A_1 + A_2 + A_3 + A_4 + A_5 + A_6 + A_7 + A_8 + A_9}$$

$$WL_{511} = \overline{\overline{A_0} + \overline{A_1} + \overline{A_2} + \overline{A_3} + \overline{A_4} + \overline{A_5} + \overline{A_6} + \overline{A_7} + \overline{A_8} + \overline{A_9}}$$

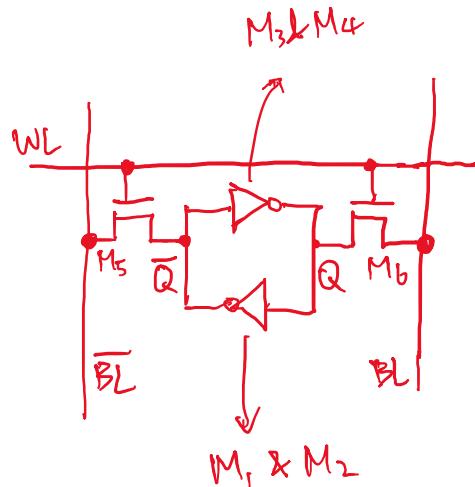


Static RAM (SRAM)



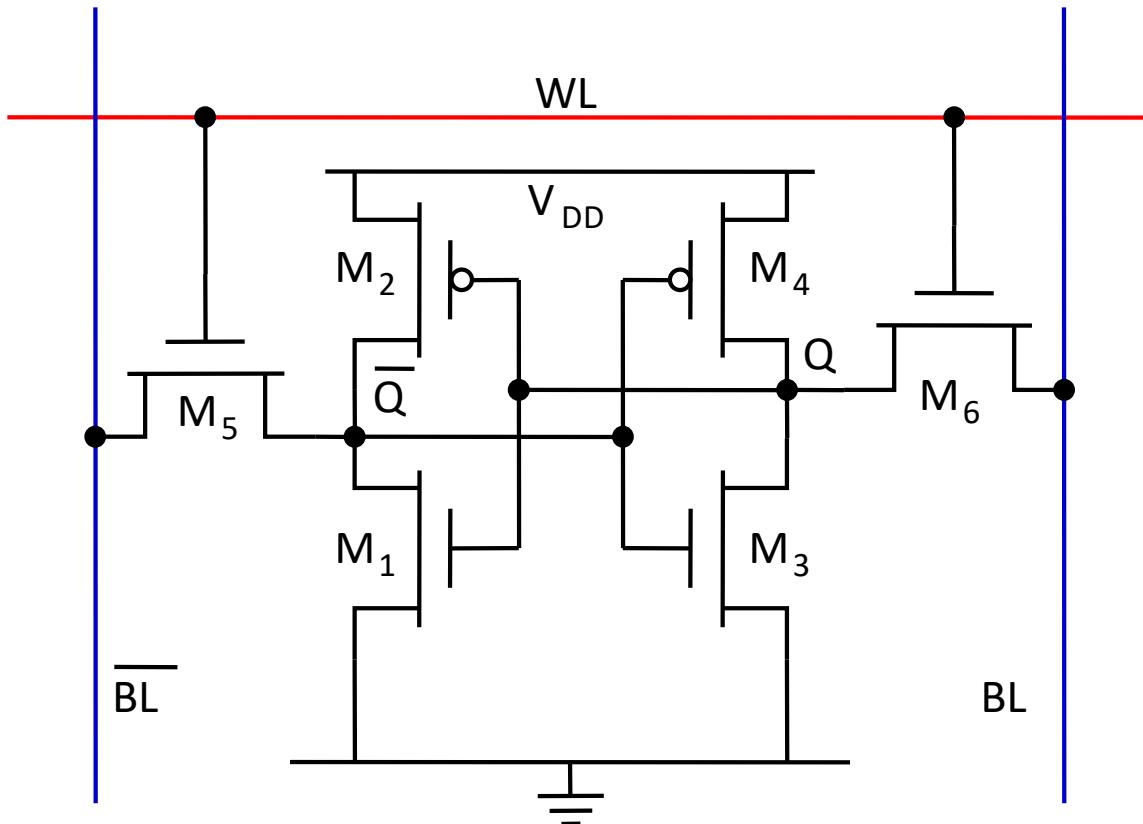


6-transistor CMOS SRAM Cell



read

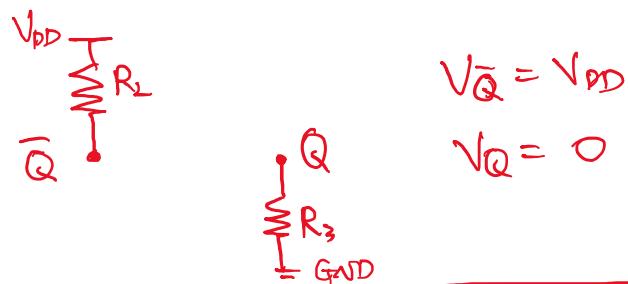
- precharge BL & \overline{BL}
- assert WL high
- sense BL & \overline{BL}
(read)



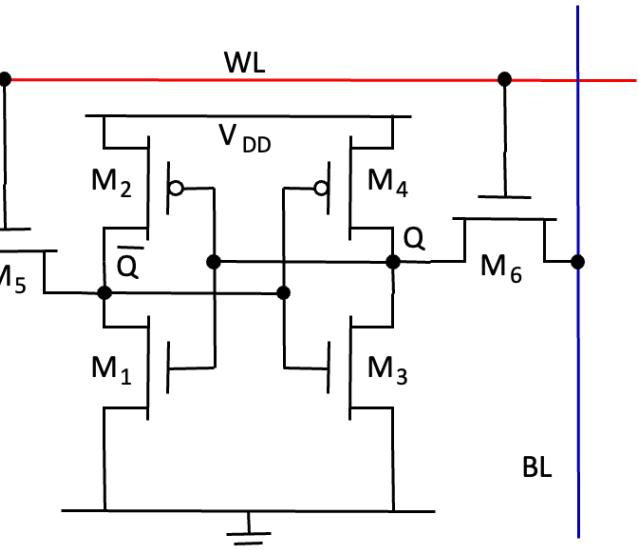
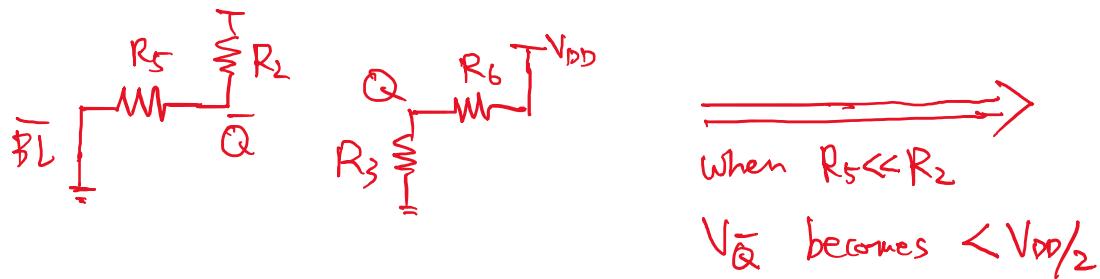
6-transistor CMOS SRAM Cell

Writing $Q = 0 \rightarrow 1$

When $Q = 0$ $M_2 \& M_3$: on
 $\bar{Q} = 1$ $M_1 \& M_4$: OFF

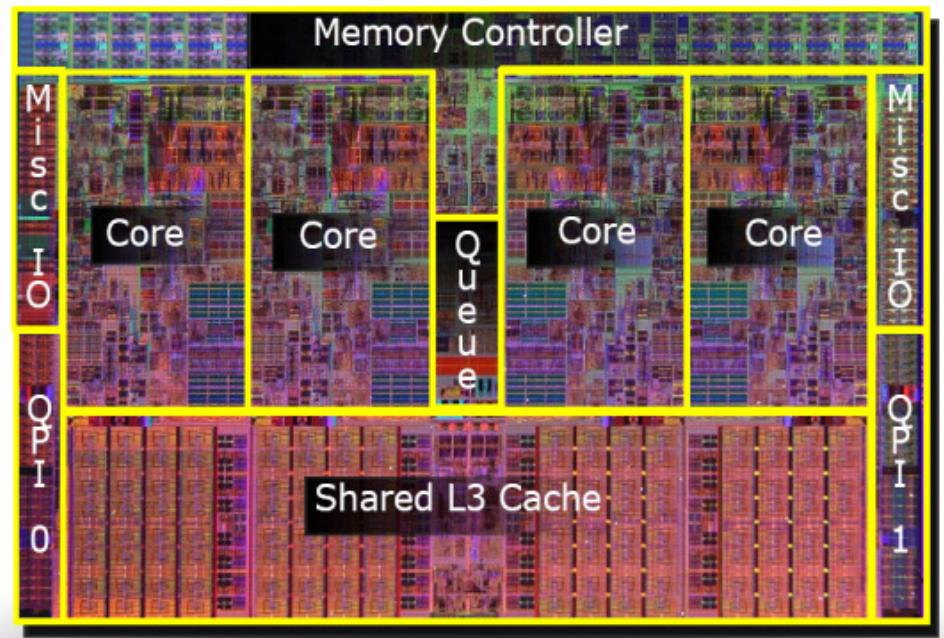


to write $Q = 1$, drive $\begin{cases} WL = V_{DD} \\ BL = V_{DD}, \bar{BL} = 0 \end{cases}$



M_4 turns on $\Rightarrow Q$ becomes 1
 M_3 turns off $\Rightarrow \bar{Q}$ " " 0

On-chip SRAM is very common



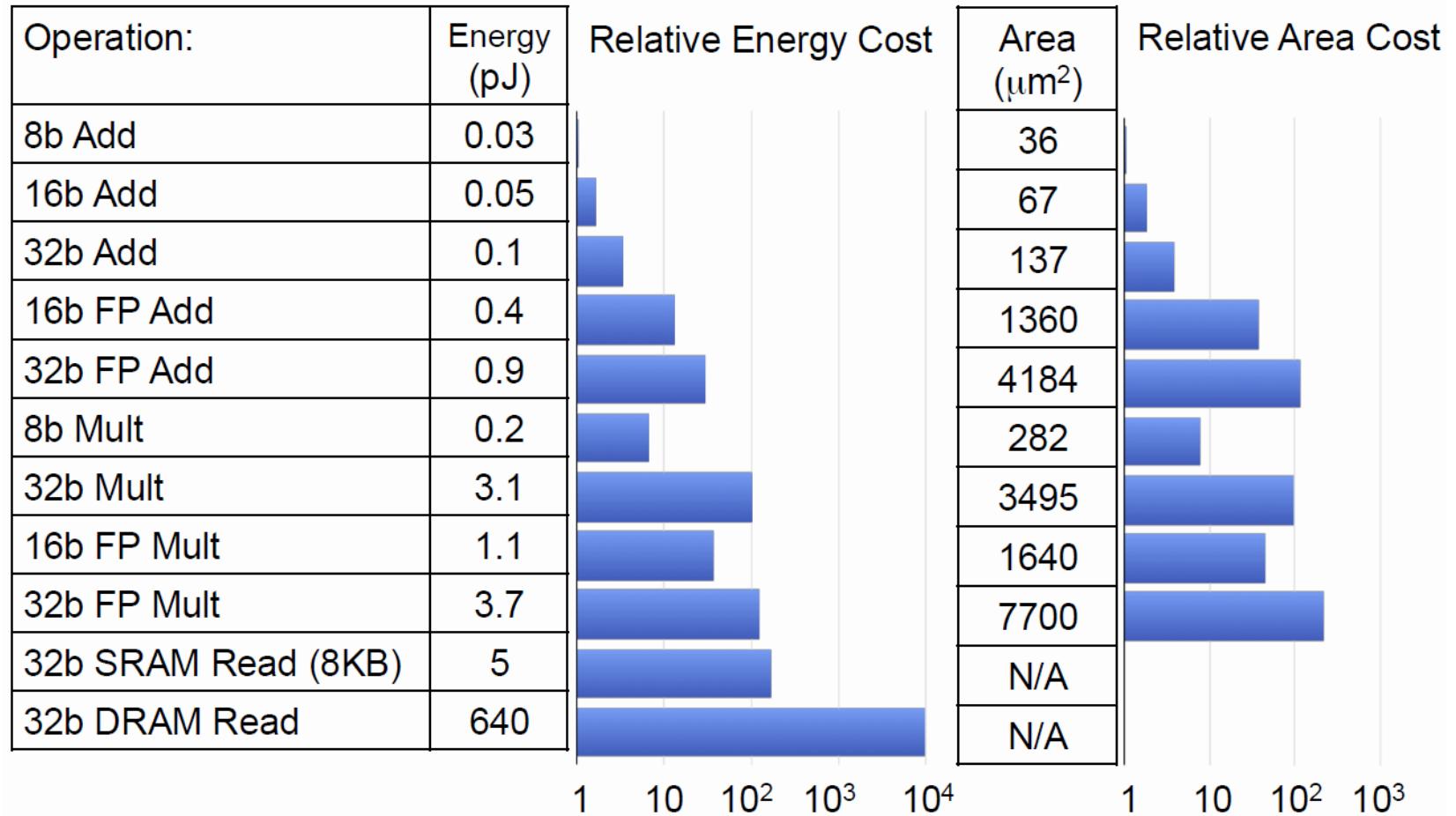
Intel Core i7

Recently used data is stored in on-chip *cache* memory

The larger the cache, the less frequently you need to go off-chip to access data (SLOW)

SRAM is technology compatible with logic

Normalized Energy and Area



[Horowitz, “Computing’s Energy Problem (and what we can do about it)”, ISSCC 2014]