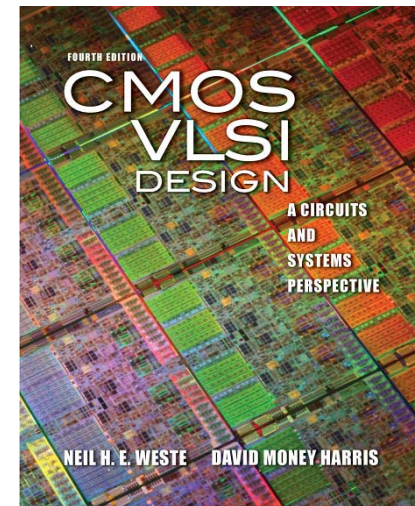


Power dissipation

October 10, 2017

Kjell Jeppson, Lena Peterson

Chapter 5 Power in W & H



Aim of the lecture

- How is power dissipated in an integrated CMOS circuit
 - **Dynamic** power dissipation
 - **Static** power dissipation
- Guidelines for reducing **dynamic** power dissipation
 - Activity factor
 - Capacitance – clock gating
 - Supply voltage: multi-VDD
 - Frequency: dynamic voltage scaling
- Guidelines for reducing **static** power dissipation
 - Reducing subthreshold leakage: multi-VT, stacking effect
 - Reducing gate leakage: high-K materials, oxide thickness
 - Junction leakage
 - Power gating

Why care?

Power in vs power out on chip

- Power grid
- Clock grid
- Cooling
- Hot spots

Clock net

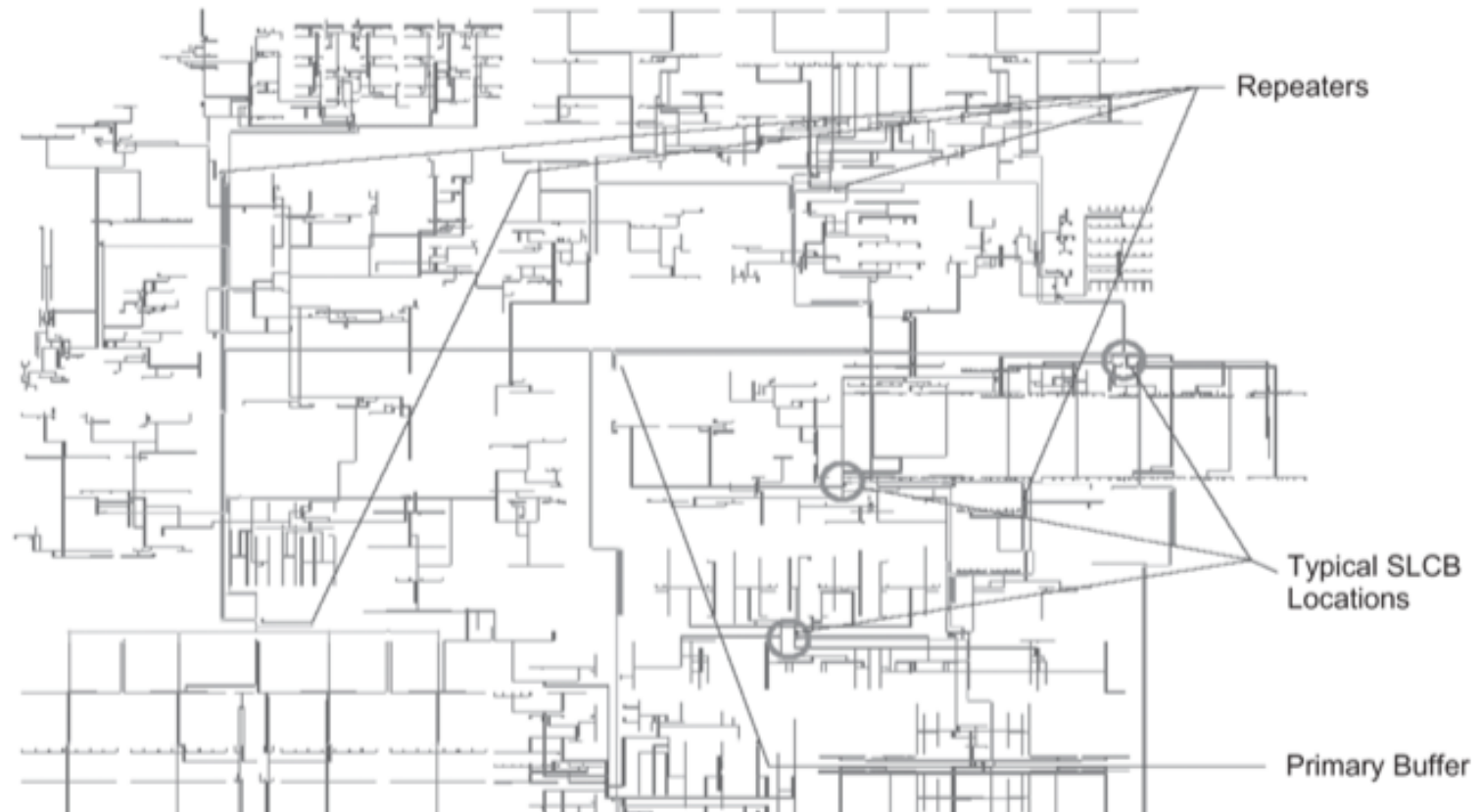


FIGURE 13.25 Itanium 2 modified H-tree

Power and energy

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

Power and energy units

- Power (work done per time unit)
 - Watts (W) = Joule/second (J/s) = Nm/s
 - $1\text{W} = 1\text{ VA} = 1\text{ VC/s}$
- Energy (\approx the ability to do work)
 - $1\text{ Joule (J)} = 1\text{ Nm} = 1\text{ Ws} = 1\text{ CV}$
 - Energy is often given in Whr = 3600 J
 - Note that work and energy has the same unit

Power in circuit elements

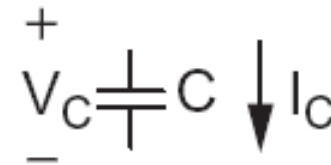
$$P_{V_{DD}}(t) = I_{DD}(t) V_{DD}$$



$$P_R(t) = \frac{V_R^2(t)}{R} = I_R^2(t) R$$



$$E_C = \int_0^{\infty} I(t) V(t) dt = \int_0^{\infty} C \frac{dV}{dt} V(t) dt$$



$$= C \int_0^{V_C} V(t) dV = C \frac{V_C^2}{2} = \boxed{\frac{1}{2} C V_C^2}$$

E_C is energy stored in Capacitor

Note that it is $\frac{1}{2} VQ$

Charging a capacitor

- When the gate output **rises**
 - Energy stored in capacitor:

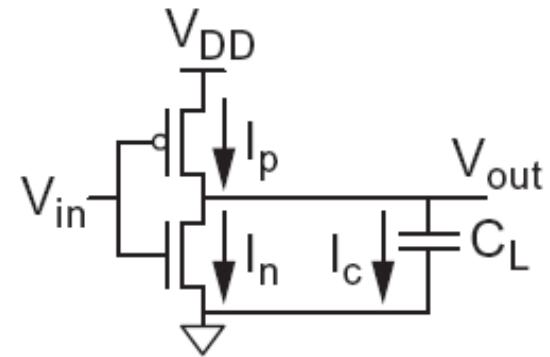
$$E_C = \frac{1}{2} C_L V_{DD}^2$$

- Energy drawn from supply:

$$E_{V_{DD}} = \int_0^{\infty} I(t) V_{DD} dt = \int_0^{\infty} C_L \frac{dV}{dt} V_{DD} dt = C_L V_{DD} \int_0^{V_{DD}} dV = C_L V_{DD}^2$$

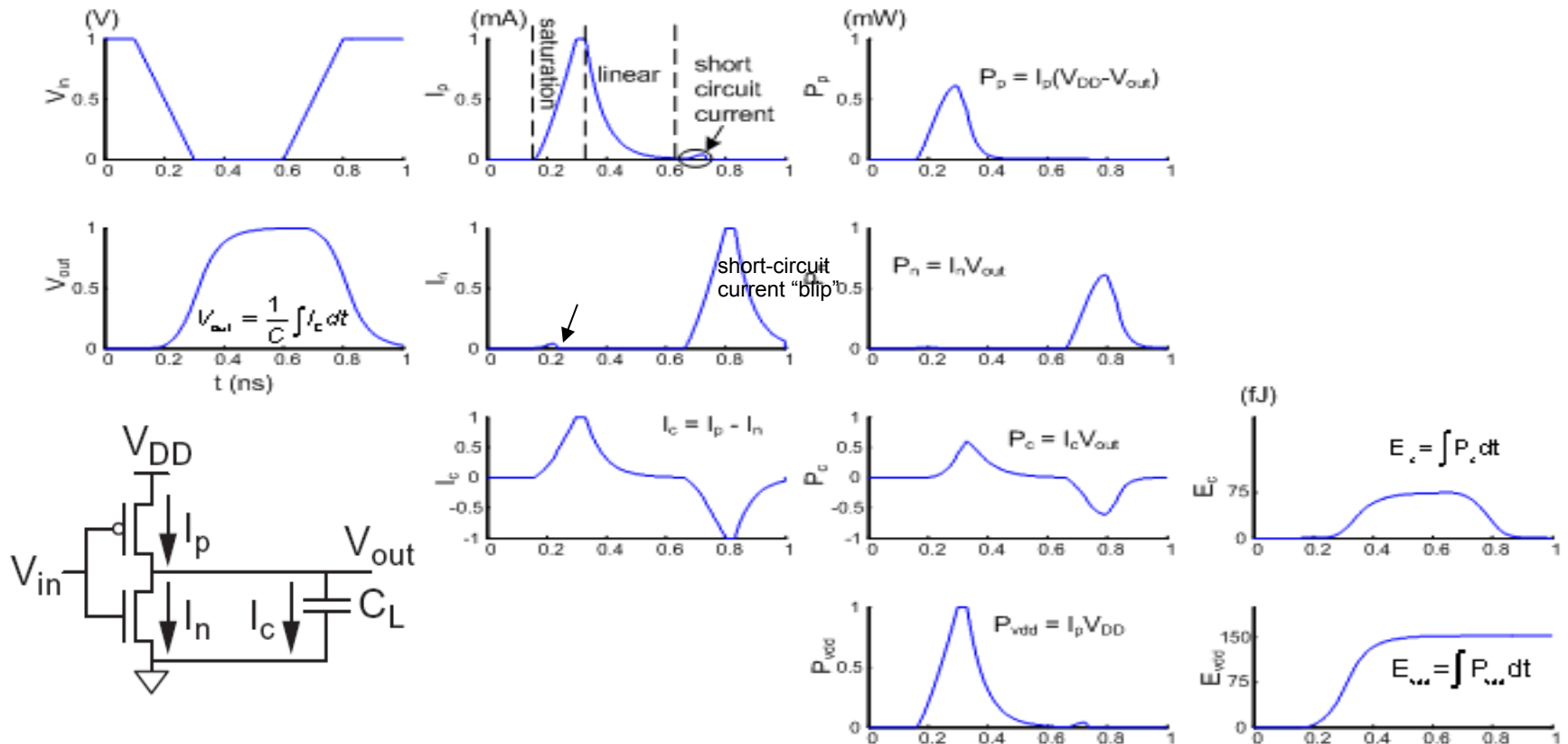
- Half $E_{V_{DD}}$ is dissipated in the pMOS transistor as heat, other half stored in capacitor

- When the gate output **falls**
 - Energy in capacitor is dumped to GND
 - Dissipated as heat in the nMOS transistor

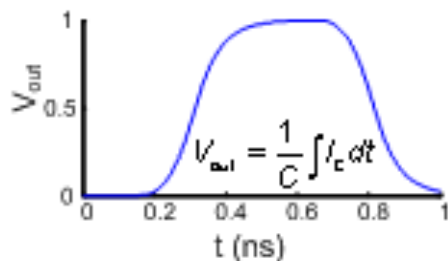
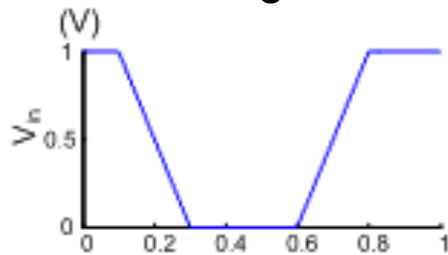


Switching waveforms

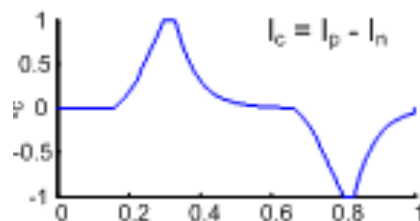
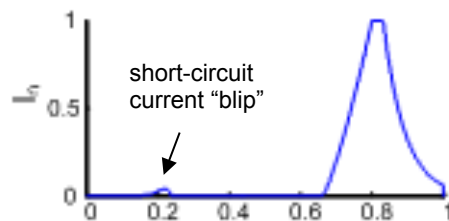
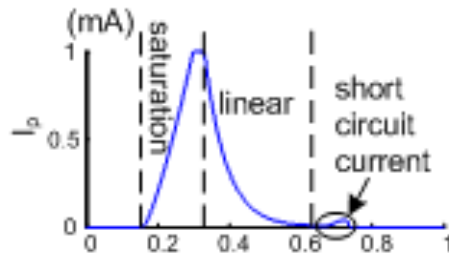
- Example: $V_{DD} = 1.0$ V, $C_L = 150$ fF, $f = 1$ GHz



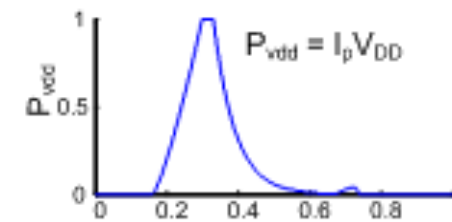
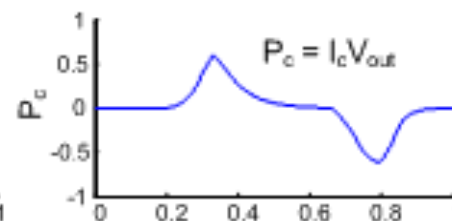
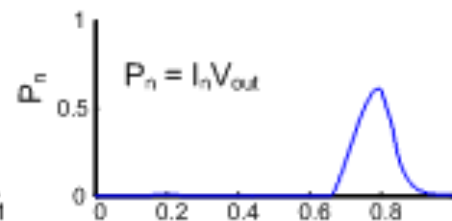
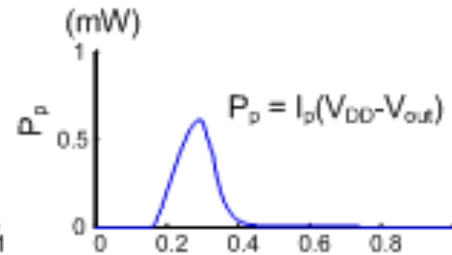
Voltages



Currents

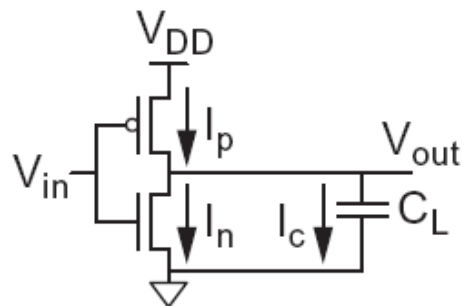
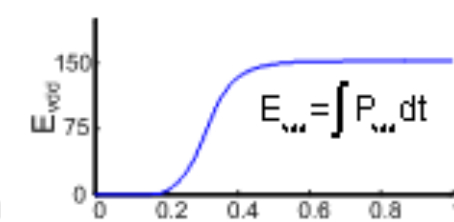
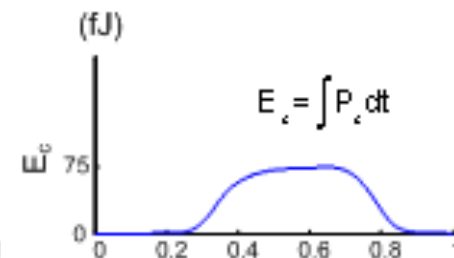


Instantaneous power



Energy

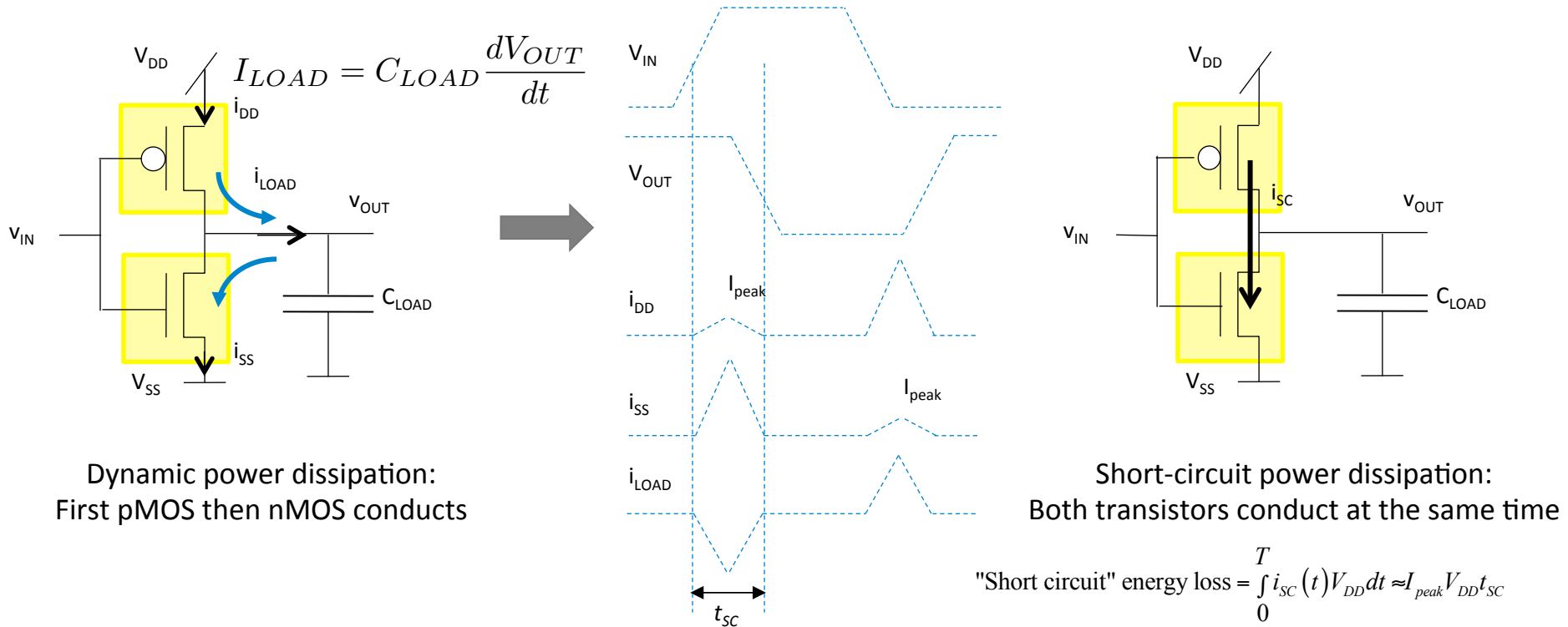
W & H: Fig 5.5



Example of switching waveforms:
 $V_{DD} = 1.0 \text{ V}$, $C_L = 150 \text{ fF}$, $f = 1 \text{ GHz}$

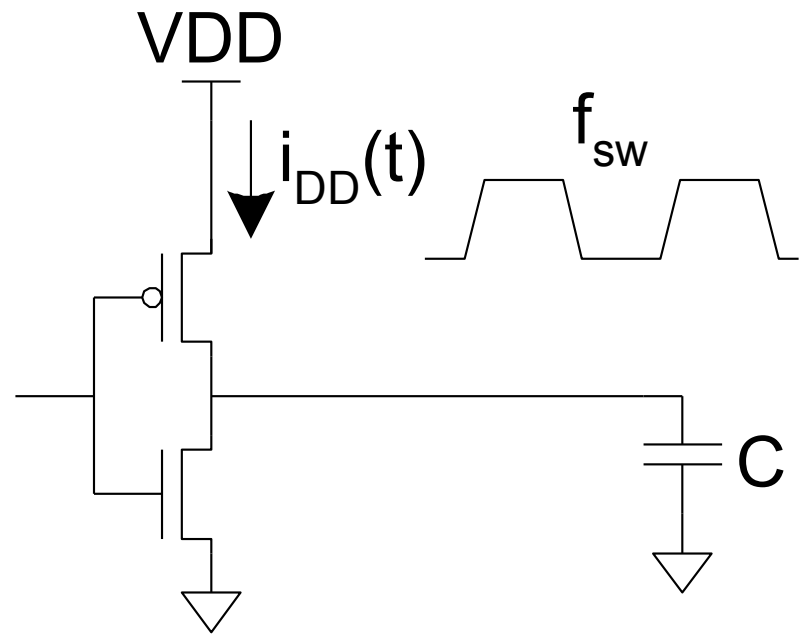
Switching waveforms

- Example: $V_{DD} = 1.0 \text{ V}$, $C_L = 150 \text{ fF}$, $f = 1 \text{ GHz}$



Switching power

$$\begin{aligned} P_{\text{switching}} &= \frac{1}{T} \int_0^T i_{DD}(t) V_{DD} dt \\ &= \frac{V_{DD}}{T} \int_0^T i_{DD}(t) dt \\ &= \frac{V_{DD}}{T} [T f_{\text{sw}} C V_{DD}] \\ &= C V_{DD}^2 f_{\text{sw}} \end{aligned}$$



Switching power due to charging C_L

Energy per transition:

$$E_{V_{DD}} = C_L V_{DD}^2$$

Unit: Joules (J)

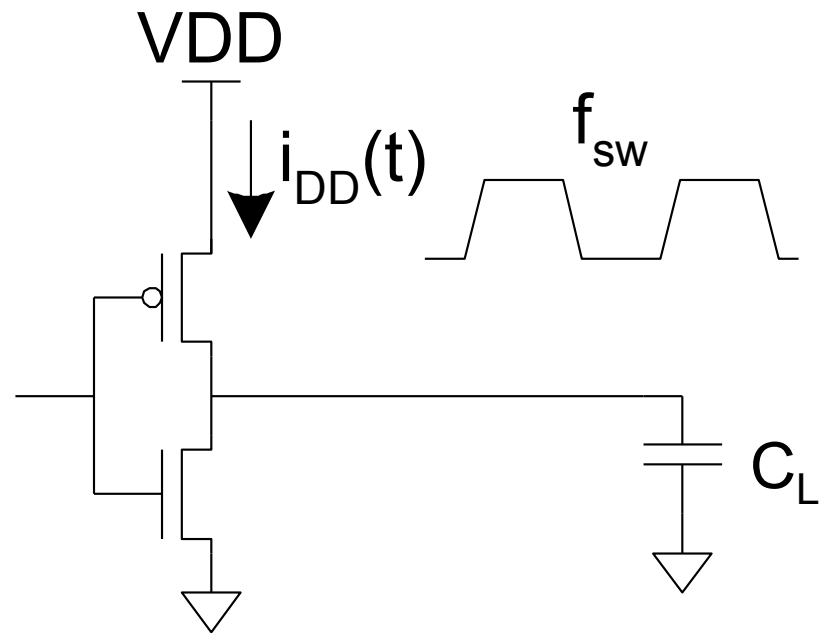
Number of low-to-high transitions per second:

$$f_{sw} \quad \text{Unit: 1/second}$$

Power:

$$P_{V_{DD}} = f_{sw} C_L V_{DD}^2$$

Unit: Joules/second = Watts



Activity factor

- Suppose the system clock frequency = f
- Let $f_{sw} = \alpha f$, where α = activity factor
 - If the signal is a clock, $\alpha = 1$
 - If the signal switches once per cycle, $\alpha = 1/2$
- Dynamic power:

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

Activity factor - definition

α = Probability that output switches from 0 to 1

$\overline{P_i}$ = probability that node i is 1

$P_i = 1 - \overline{P_i}$ = probability that node i is 0

If the probabilities are uncorrelated:

$$\alpha = \overline{P_i} P_i = (1 - P_i) P_i$$

For random data: $P_i = 0.5$ so $\alpha = 0.25$

See W&H section 5.2.1

Power-dissipation sources

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

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
- We will generally ignore this component
- But if rise or fall times are long it may dominate!

Dynamic power: an example

- 1 billion transistor chip
 - 50M logic transistors
 - Average width: $12 \lambda = 12 \times 25 \text{ nm} = 300 \text{ nm}$
 - Activity factor = 0.1
 - 950M memory transistors
 - Average width: $4 \lambda = 100 \text{ nm}$
 - Activity factor = 0.02
 - 1.0 V 65 nm process
 - $C = 1 \text{ fF}/\mu\text{m}$ (gate) + $0.8 \text{ fF}/\mu\text{m}$ (diffusion)
- Estimate dynamic power consumption @ 1 GHz.
Neglect wire capacitance and short-circuit current.

Dynamic power: an example

1.0 V 65 nm process

$C = 1 \text{ fF}/\mu\text{m}$ (gate) + $0.8 \text{ fF}/\mu\text{m}$ (diffusion/parasitic)

Estimate power dissipation at 1 GHz!

1 billion transistor chip

50M logic transistors

Average width: 300 nm Activity factor = 0.1

$$C_{\text{logic}} = (50 \times 10^6)(0.3 \mu\text{m})(1.8 \text{ fF} / \mu\text{m}) = 27 \text{ nF}$$

950M memory transistors

Average width: 100 nm Activity factor = 0.02

$$C_{\text{mem}} = (950 \times 10^6)(0.1 \mu\text{m})(1.8 \text{ fF} / \mu\text{m}) = 170 \text{ nF}$$

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

W & H: Example 5.1

Dynamic power reduction

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

- Try to minimize:
 1. Activity factor
 2. Capacitance
 3. Supply voltage
 4. Frequency

1. Reduce activity factor

- Let $P_i = \text{Probability}(\text{node } i = 1)$
 - $\bar{P}_i = 1 - P_i$ (probability(node $i = 0$))
- $\alpha_i = \bar{P}_i \times P_i$
- Completely random data has $P = 0.5$ and thus $\alpha = 0.25$
- Data is often not completely random
 - e.g. upper bits of 64-bit words representing bank account balances are usually 0
- Data propagating through ANDs and ORs has lower activity factor
 - Depends on design, but typically $\alpha \approx 0.1$

1. Reduce activity factor

α = Probability that output switches from 0 to 1

$\overline{P_i}$ = probability that node i is 1

$P_i = 1 - \overline{P_i}$ = probability that node i is 0

If probabilities are uncorrelated:

$$\alpha = \overline{P_i} P_i = (1 - P_i) P_i$$

For random data: $P_i = 0.5$ so $\alpha = 0.25$

- Data is often not completely random
- After AND or OR gates α is lower than at inputs
- Design dependent but typically $\alpha \approx 0.1$

See W&H section 5.2.1

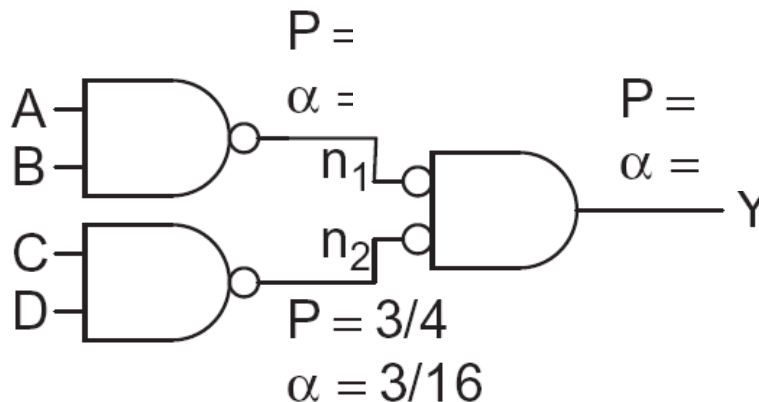
Switching probabilities

Gate	P_Y
AND2	$P_A P_B$
AND3	$P_A P_B P_C$
OR2	$1 - \bar{P}_A \bar{P}_B$
NAND2	$1 - P_A P_B$
NOR2	$\bar{P}_A \bar{P}_B$
XOR2	$P_A \bar{P}_B + \bar{P}_A P_B$

W&H: Table 5.1

Textbook example 5.2a

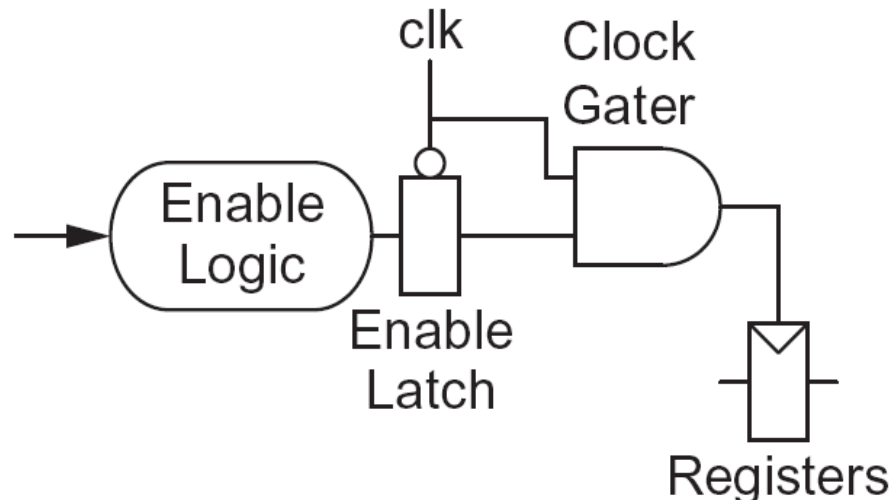
- A 4-input AND is built out of two levels of gates
- Estimate the activity factor at each node if the inputs have $P = 0.5$



$$\text{NAND2: } P_Y = 1 - P_A P_B \quad \text{NOR2: } P_Y = \overline{P_A} \overline{P_B}$$

Clock gating

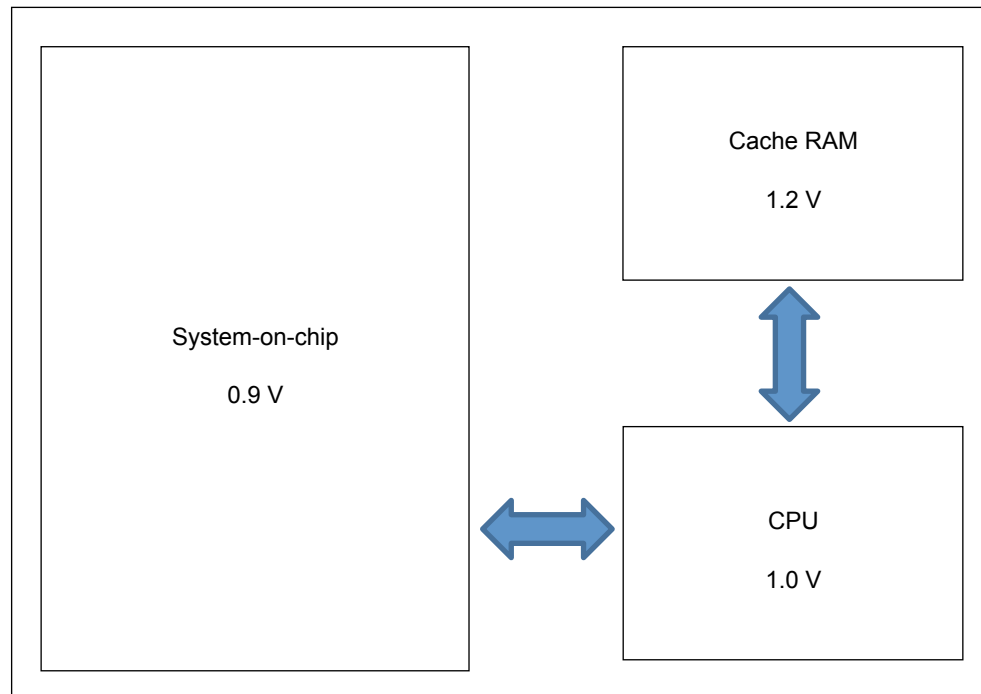
- The best way to reduce activity is to turn off the clock to registers in unused blocks = sleep mode
 - Saves clock activity (clock has $\alpha = 1$)
 - Eliminates all switching activity in the block
 - Requires determining if block will be used



2. Reduce capacitance

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

3. Reduce supply voltage – Multi- V_{DD}

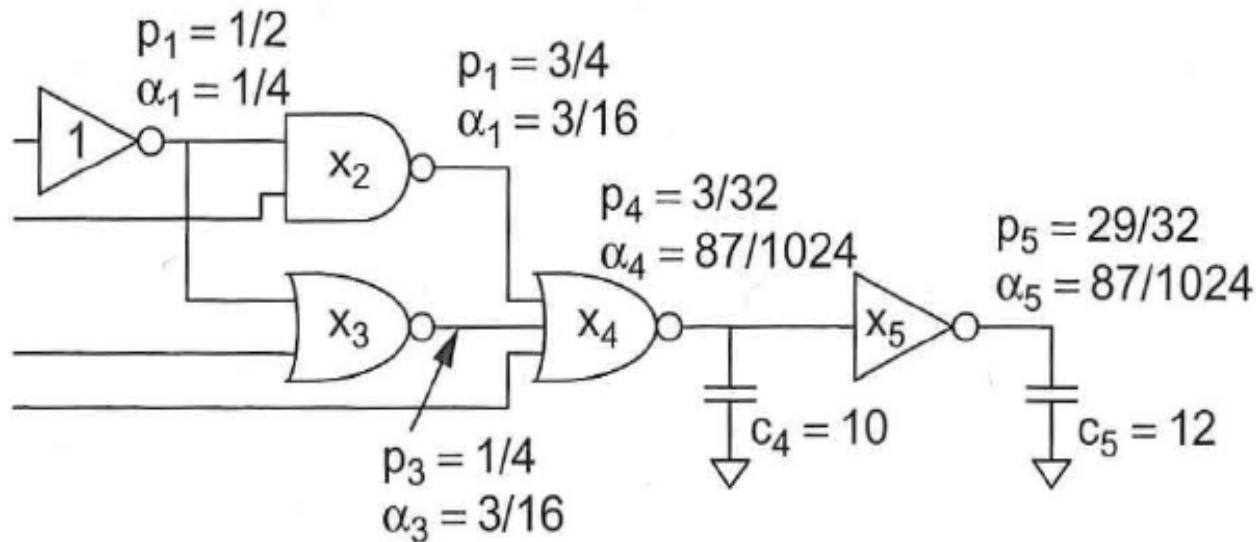


4. Reduce frequency

- Do not run clock faster than necessary!
- Multiple frequency domains
 - Often integer factors & synchronized clocks among domains
- Lower clock frequency can be combined with smaller transistors and lower V_{DD} which saves even more power.

Tradeoff: Textbook example 5.3

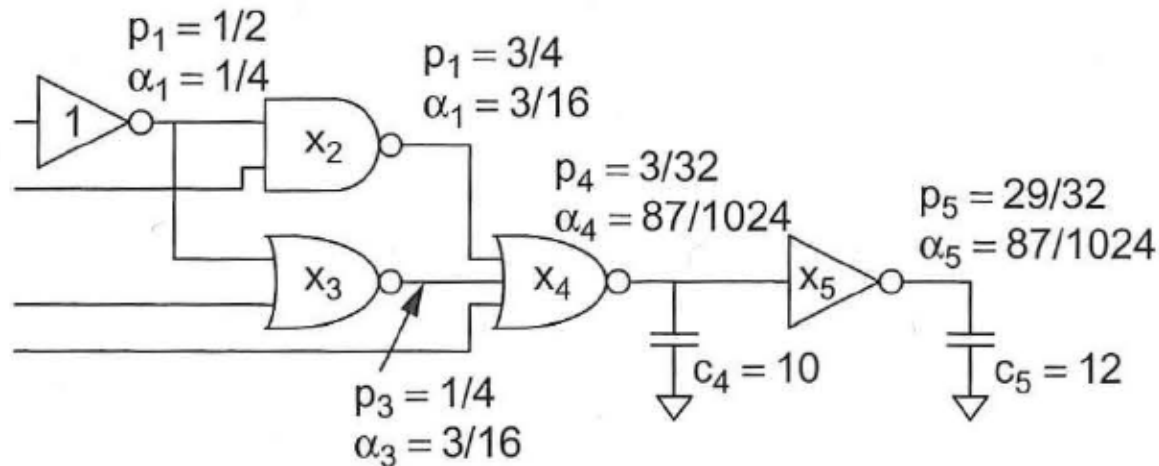
- Generate an energy-delay trade-off curve for the circuit below as delay varies from the minimum possible ($D_{min}=23.44 \tau$) to 50τ . Assume all input probabilities are 0.5.



Textbook example 4.3

$$E = \frac{1}{4} \left(1 + \frac{4}{3}x_2 + \frac{5}{3}x_3 \right) + \frac{3}{16} \left(2x_2 + \frac{7}{3}x_4 \right) + \frac{3}{16} \left(2x_3 + \frac{7}{3}x_4 \right) + \frac{87}{1024} \left(10 + 3x_4 + x_5 \right) + \frac{87}{1024} \left(12 + x_5 \right)$$

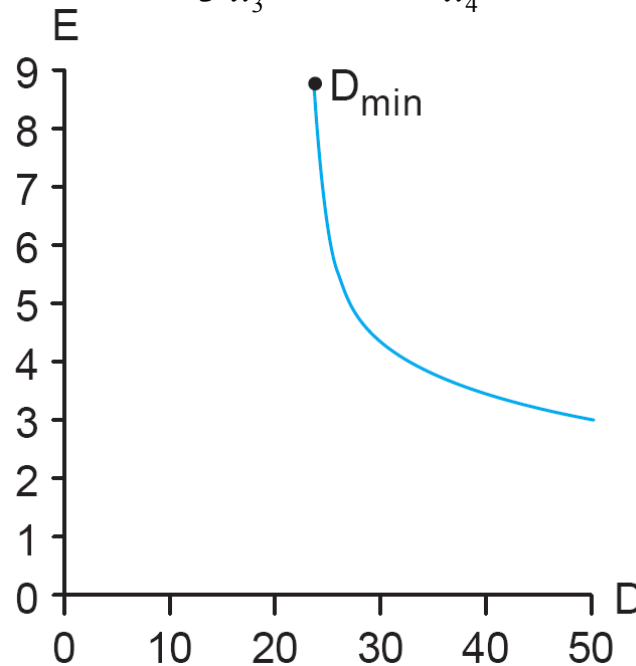
$$d = 1 + \frac{4}{3}x_2 + \frac{5}{3}x_3 + 2 + \frac{7}{3}x_4 + 3 + \frac{10 + x_5}{x_4} + 1 + \frac{12}{x_5}$$



Textbook example 5.3

$$E = \frac{1}{4} \left(1 + \frac{4}{3}x_2 + \frac{5}{3}x_3 \right) + \frac{3}{16} \left(2x_2 + \frac{7}{3}x_4 \right) + \frac{3}{16} \left(2x_3 + \frac{7}{3}x_4 \right) + \frac{87}{1024} (10 + 3x_4 + x_5) + \frac{87}{1024} (12 + x_5)$$

$$d = 1 + x_2 + x_3 + 2 + \frac{5}{3} \frac{x_4}{x_3} + 3 + \frac{3 + 10 + x_5}{x_4} + 1 + \frac{12}{x_5}$$



Textbook example 5.3

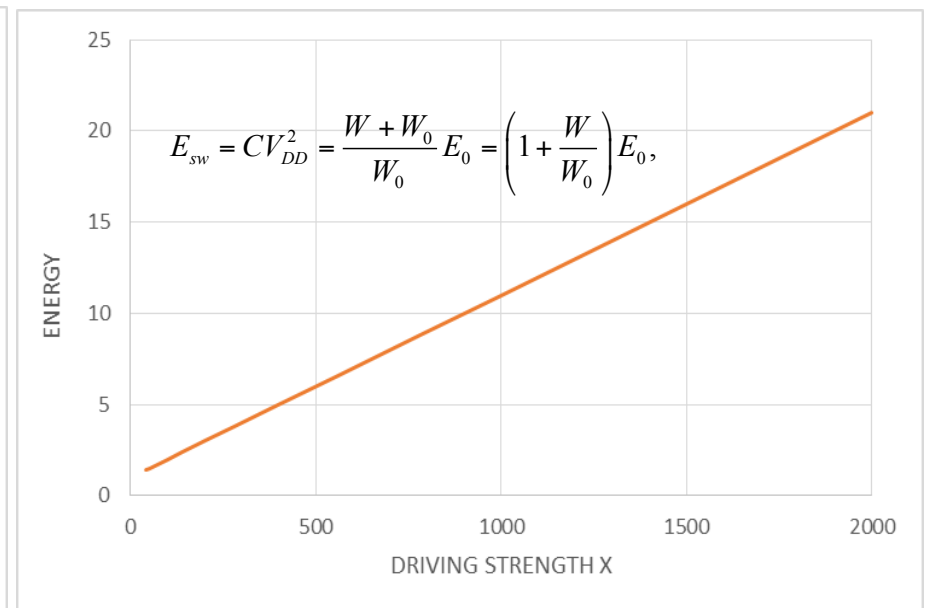
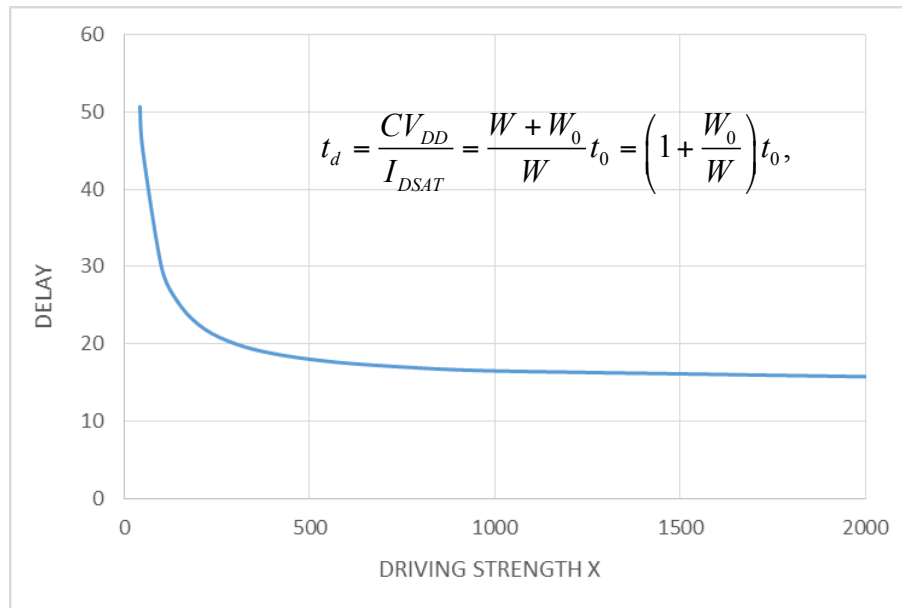
- Can we understand this?
- Model delay and energy assuming capacitance is not fully scalable

$$t_d = \frac{CV_{DD}}{I_{DSAT}} = \frac{W + W_0}{W} t_0, \quad E_{sw} = CV_{DD}^2 = \frac{W + W_0}{W_0} E_0,$$

$$t_d = \frac{CV_{DD}}{I_{DSAT}} = \left(1 + \frac{W_0}{W}\right) t_0, \quad E_{sw} = CV_{DD}^2 = \left(1 + \frac{W_0}{W}\right) E_0,$$

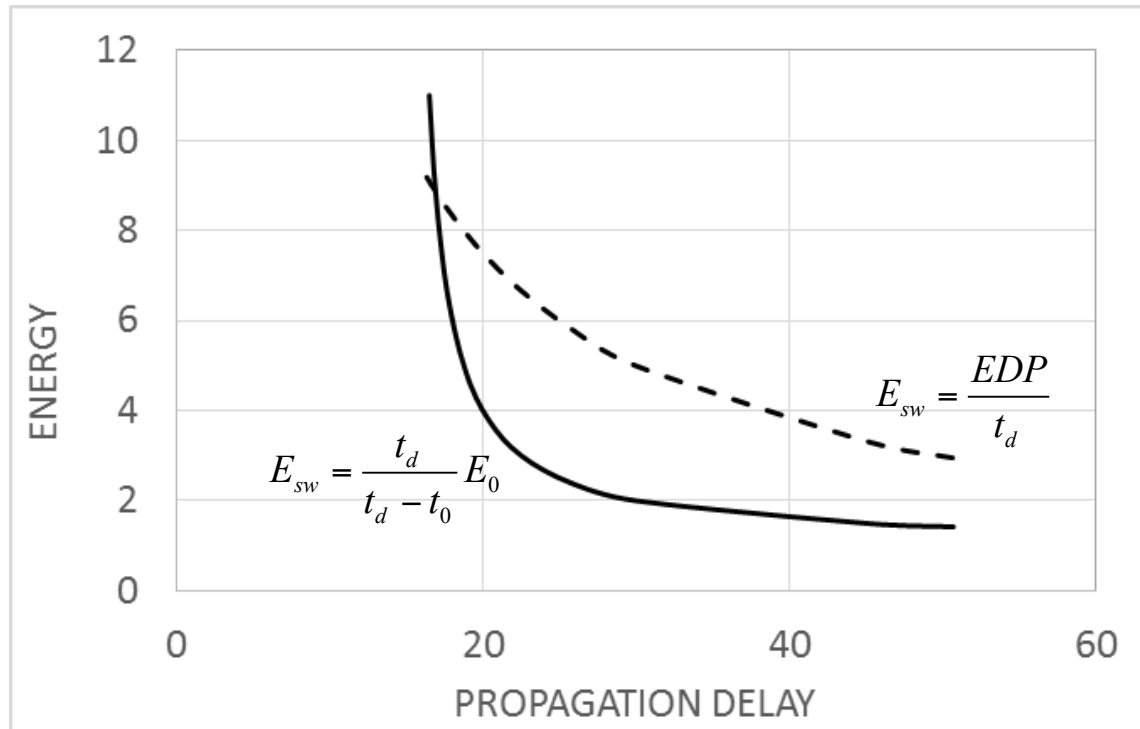
Textbook example 4.3

- Can we understand this?
- Model delay and energy assuming capacitance is not fully scalable



Textbook example 4.3

- Can we understand this?



Adder switching activity

ADD/SUBTRACT		ADD=0		CIN=?		A=		-65		<<<<<< ENTER TWO NUMBERS							
CONTROL SIGNAL:		1		SUB=1		CIN=?		B=		32		<<<<<< -128<NUMBER<128					
								SUM=		-97							
a7	b7	a6	b6	a5	b5	a4	b4	a3	b3	a2	b2	a1	b1	a0	b0		
1	0	0	0	1	1	1	0	1	0	1	0	1	0	1	0		
1	1	0	1	1	0	1	1	1	1	1	1	1	1	1	1		
G7	P7	G6	P6	G5	P5	G4	P4	G3	P3	G2	P2	G1	P1	G0	P0		
1	0	0	1	0	1	1	0	1	0	1	0	1	0	1	0	CIN	P
1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	1
G7:0 P7:0		G6:0 P6:0		G5:0 P5:0		G4:0 P4:0		G3:0 P3:0		G2:0 P2:0		G1:0 P1:0		G0:0 P0:0			
1		0		0		1		1		1		1		1			
SUM7		SUM6		SUM5		SUM4		SUM3		SUM2		SUM1		SUM0			
SUM converted back to decimal:				-97				Both sums are equal?		YES							
								OVERFLOW?		NO							

Adder switching activity

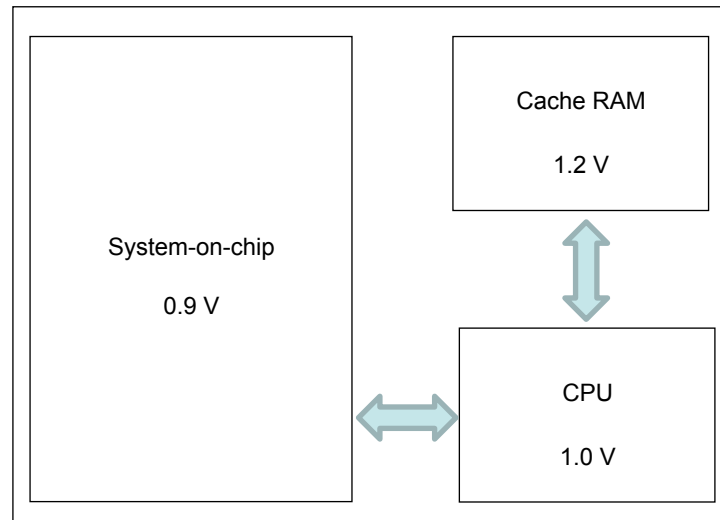
a7	b7	a6	b6	a5	b5	a4	b4	a3	b3	a2	b2	a1	b1	a0	b0		
0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5	0,5		
	0,5		0,5		0,5		0,5		0,5		0,5		0,5		0,5		
G7	P7	G6	P6	G5	P5	G4	P4	G3	P3	G2	P2	G1	P1	G0	P0		
0,25	0,5	0,25	0,5	0,25	0,5	0,25	0,5	0,25	0,5	0,25	0,5	0,25	0,5	0,25	0,5	CIN	
0,4	0,002	0,4	0,004	0,4	0,008	0,401	0,016	0,402	0,031	0,405	0,063	0,414	0,125	0,438	0,25	0,5	0,5
G7:0	P7:0	G6:0	P6:0	G5:0	P5:0	G4:0	P4:0	G3:0	P3:0	G2:0	P2:0	G1:0	P1:0	G0:0	P0:0		
0,5		0,5		0,5		0,5		0,5		0,5		0,5		0,5			
SUM7		SUM6		SUM5		SUM4		SUM3		SUM2		SUM1		SUM0			

Adder switching activity

a7	b7	a6	b6	a5	b5	a4	b4	a3	b3	a2	b2	a1	b1	a0	b0			alpha	cap
0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25	0,25			4	*C1
	0,25		0,25		0,25		0,25		0,25		0,25		0,25		0,25			2	*C2
G7	P7	G6	P6	G5	P5	G4	P4	G3	P3	G2	P2	G1	P1	G0	P0			1,5	*C3
0,188	0,25	0,188	0,25	0,188	0,25	0,188	0,25	0,188	0,25	0,188	0,25	0,188	0,25	0,188	0,25	CIN		2	*C4
0,24	0,002	0,24	0,004	0,24	0,008	0,24	0,015	0,24	0,03	0,241	0,059	0,243	0,109	0,246	0,188	0,5	1	1,9	*C5
G7:0	P7:0	G6:0	P6:0	G5:0	P5:0	G4:0	P4:0	G3:0	P3:0	G2:0	P2:0	G1:0	P1:0	G0:0	P0:0			0,4	*C6
0,25		0,25		0,25		0,25		0,25		0,25		0,25		0,25				2	*C7
SUM7		SUM6		SUM5		SUM4		SUM3		SUM2		SUM1		SUM0				SUM(ai*Ci)	

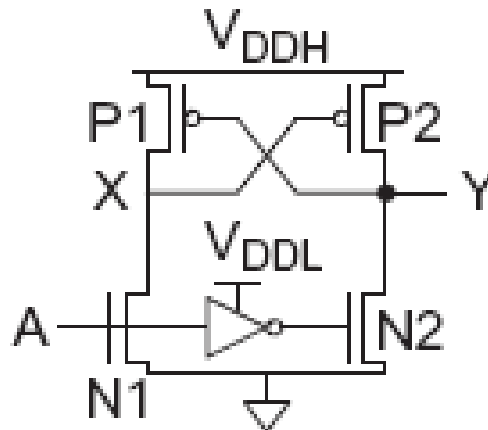
Voltage domains

- Run each block at the lowest possible voltage and frequency that meets performance requirements
- Voltage domains
 - Provide separate supplies to different blocks



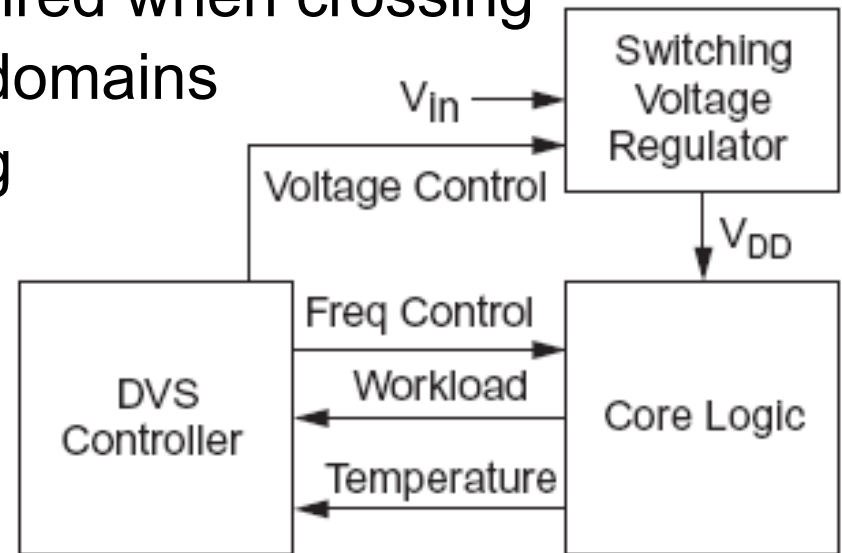
Voltage domains

- 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 (DVS)

- 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



Conclusion

dynamic power reduction

1. Activity factor

- Clock gating
- Reduce logic depth → reduced glitching

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

2. Capacitance

- Minimum transistor sizes
- Short wires

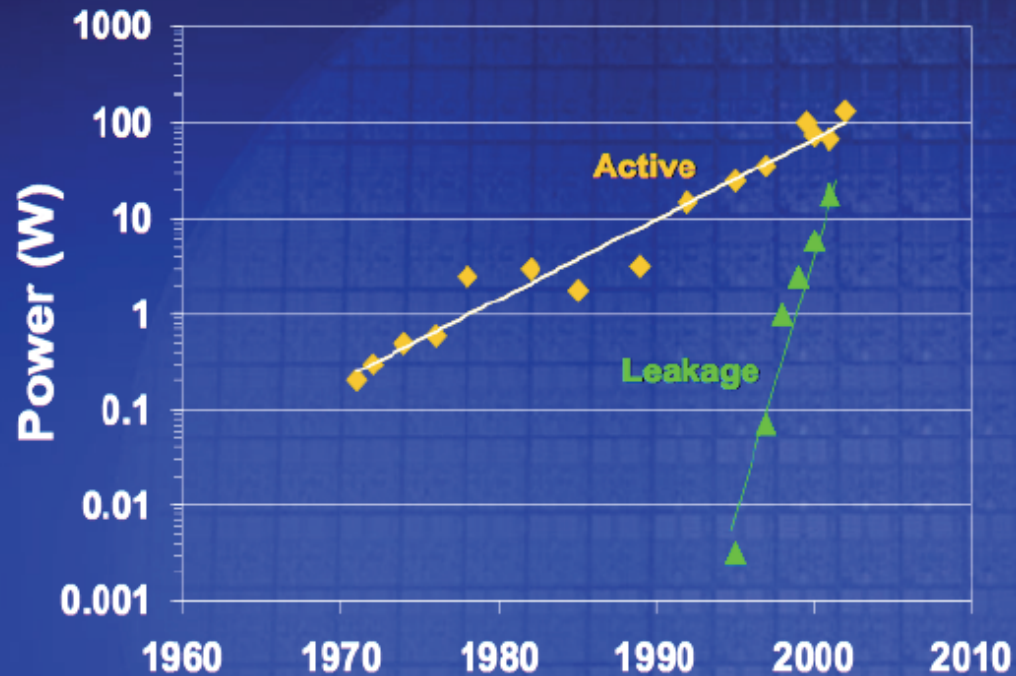
3. Supply voltage

- Multiple V_{DD} domains
- Dynamic voltage/frequency scaling

4. Frequency

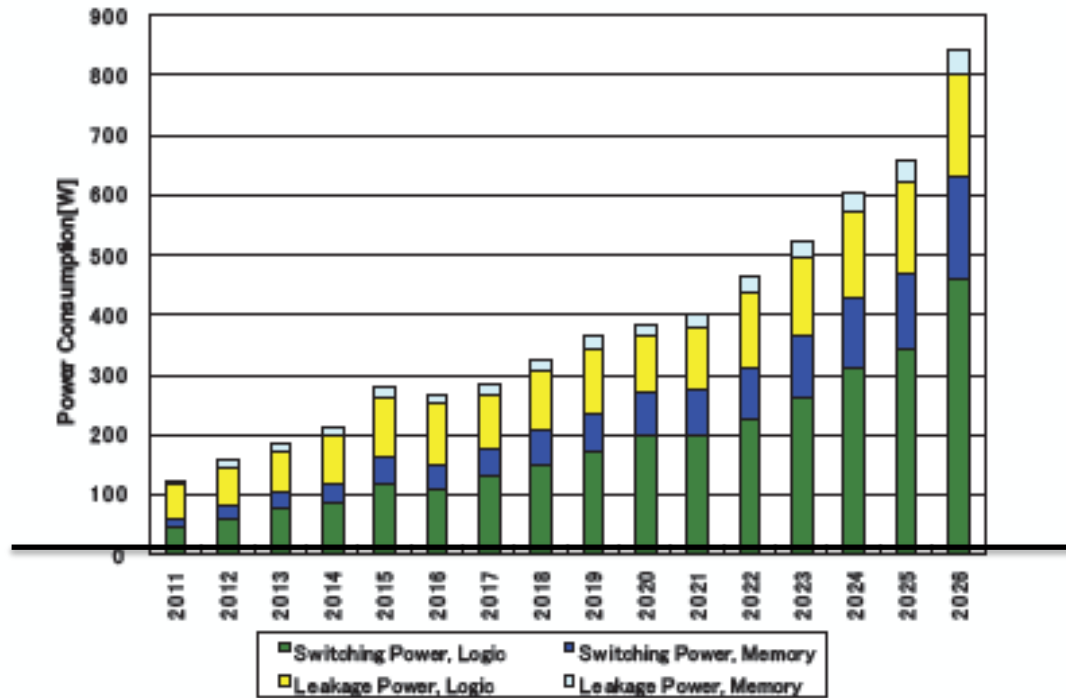
- Multiple clock domains
- Dynamic voltage/frequency scaling

Processor Power (Watts) - Active & Leakage



From DAT093

Processor power predictions

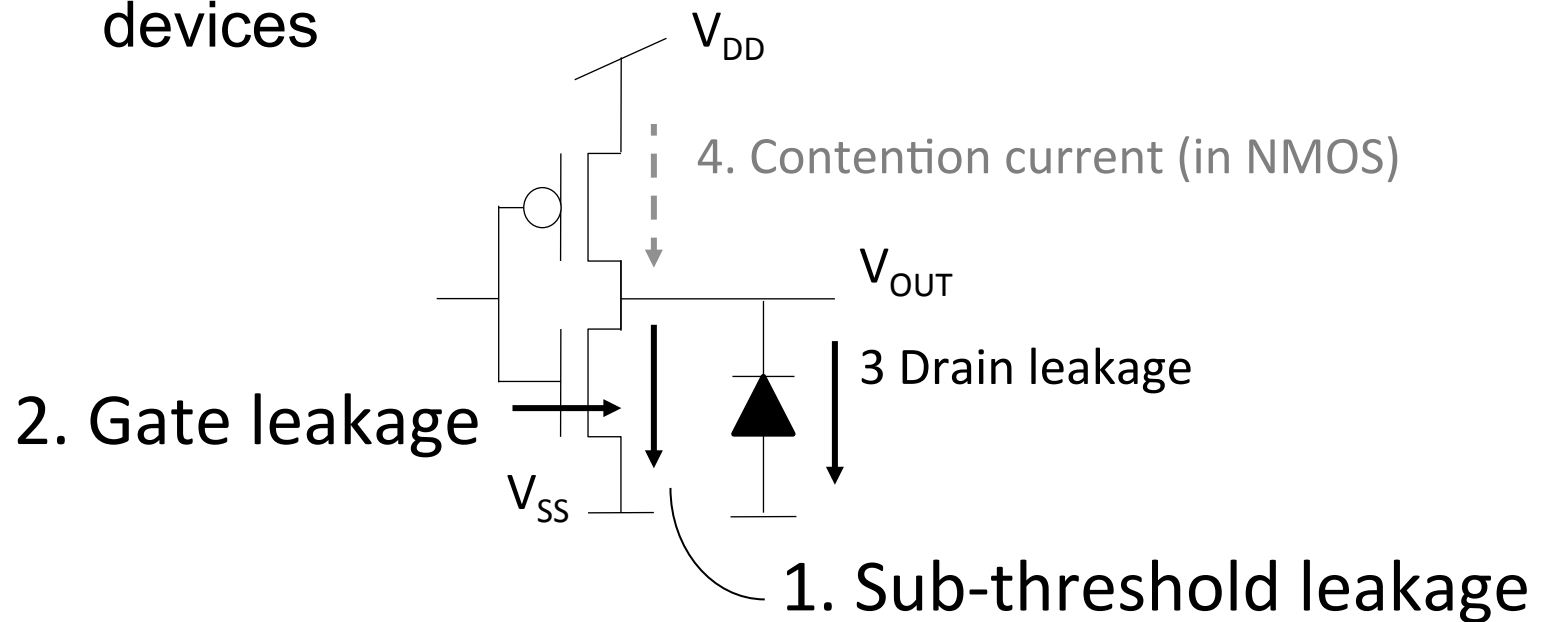


From DAT093

- Drawing out trends of transistor, interconnect, etc
- Note: know how to cool ~150W from processor 🤔
- Cooling clearly limits what can be done!

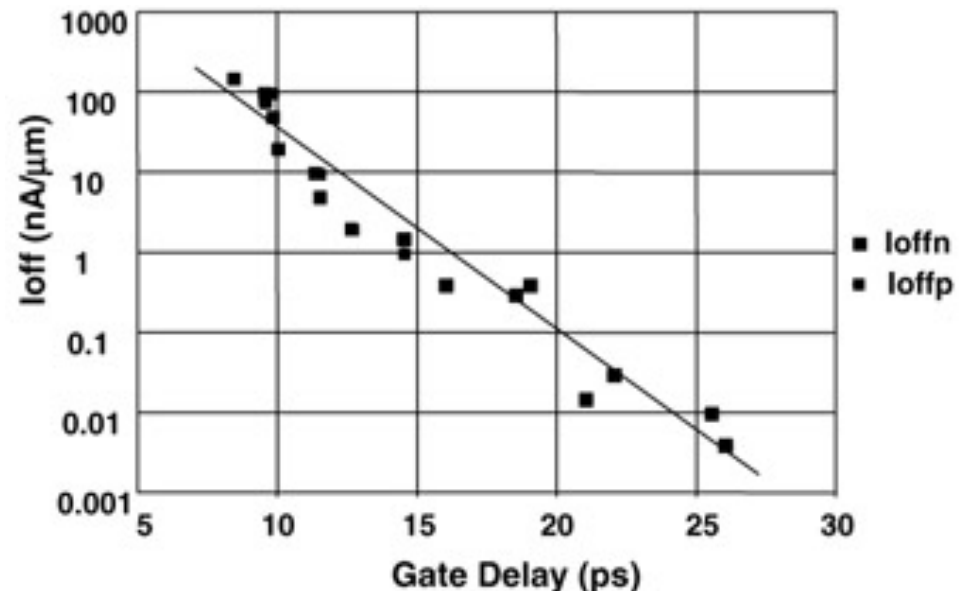
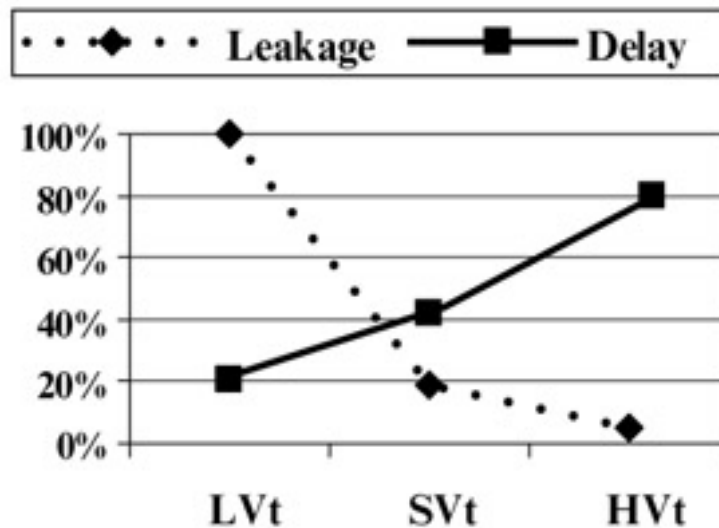
Static power

- Static power is consumed even when chip is quiescent.
 - Leakage draws power from nominally OFF devices



Reducing sub- V_T leakage: Multi- V_T

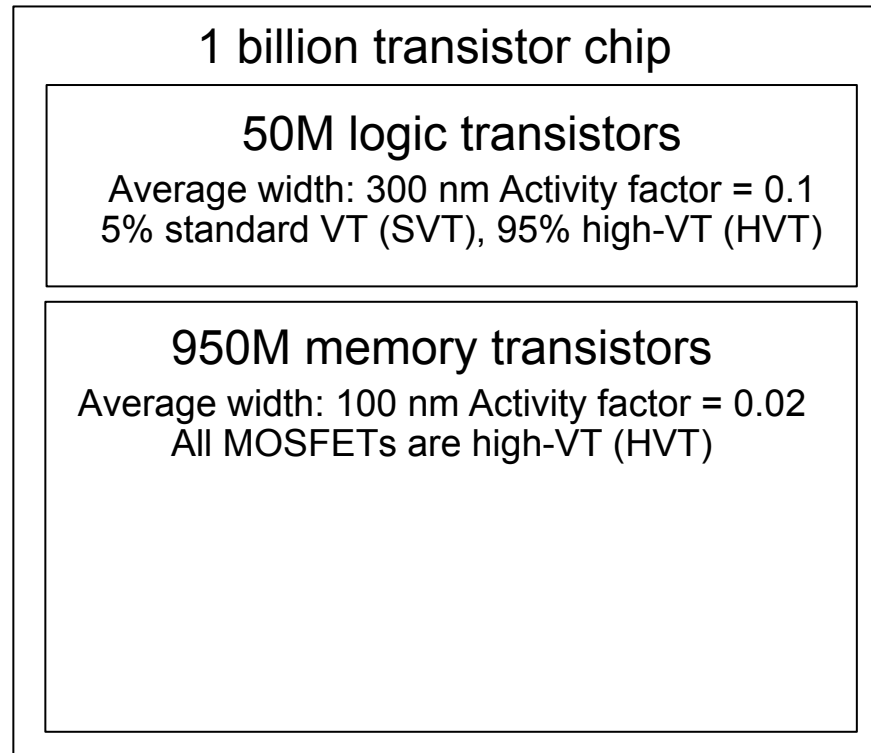
- In 65 nm processes and below, libraries with multiple (typically 3) V_T has become a common way of reducing leakage currents.



Static power: textbook example 5.4

Subthreshold leakage: SVT: 100 nA/ μm ; HVT: 10 nA/ μm
Gate leakage: 5 nA/ μm ; Junction leakage: negligible

Estimate static power consumption!



Static power: textbook example 5.4

Solution: Assumption half of transistors are ON and give gate leakage) half are OFF give subthreshold leakage)

$$W_{\text{SVT}} = 5\% * (50\text{M})(0.3\mu\text{m}) = 0.75 \times 10^6 \mu\text{m} = 75 \text{ cm}$$

$$W_{\text{HVT}} = [95\% * (50\text{M})(0.3\mu\text{m}) + (950\text{M})(0.1\mu\text{m})] = 110 \times 10^6 \mu\text{m} = 110 \text{ m}$$

$$I_{\text{sub}} = [W_{\text{SVT}} \times 100 \text{ nA}/\mu\text{m} + W_{\text{HVT}} \times 10 \text{ nA}/\mu\text{m}] / 2 = 584 \text{ mA}$$

$$I_{\text{gate}} = [(W_{\text{SVT}} + W_{\text{HVT}}) \times 5 \text{ nA}/\mu\text{m}] / 2 = 275 \text{ mA}$$

$$P_{\text{static}} = (584 \text{ mA} + 275 \text{ mA})(1.0 \text{ V}) = 860 \text{ mW}$$

- 860 mW is 14% of the 6.1 W dynamic power calculated earlier!
- Will deplete battery of handheld device rapidly

Subthreshold leakage

- For $V_{DS} > 50 \text{ mV}$

$$I_{sub} \approx I_{off} 10^{\frac{V_{gs} + \sigma(V_{DS} - V_{DD})}{S}}$$

- I_{off} = leakage at $V_{gs} = 0$, $V_{DS} = V_{DD}$

Typical values in 65 nm

$I_{off} = 100 \text{ nA}/\mu\text{m}$ @ $V_T = 0.3 \text{ V}$

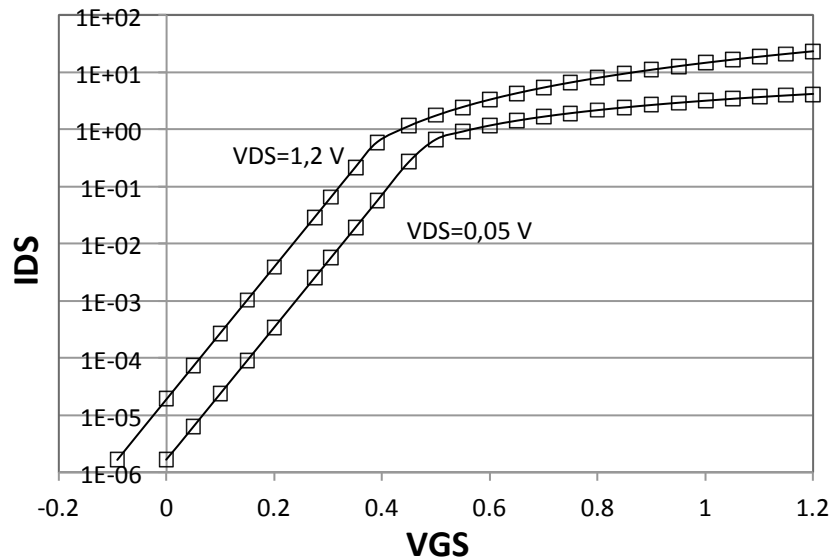
$I_{off} = 10 \text{ nA}/\mu\text{m}$ @ $V_T = 0.4 \text{ V}$

$I_{off} = 1 \text{ nA}/\mu\text{m}$ @ $V_T = 0.5 \text{ V}$

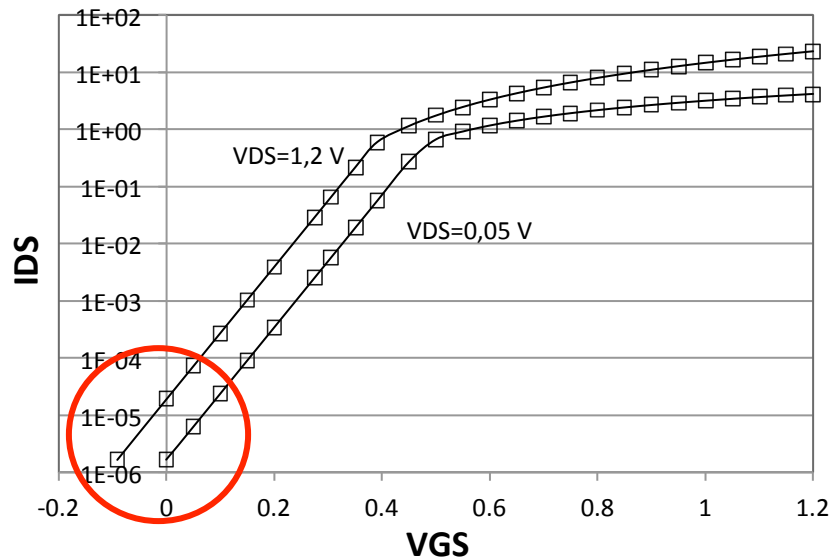
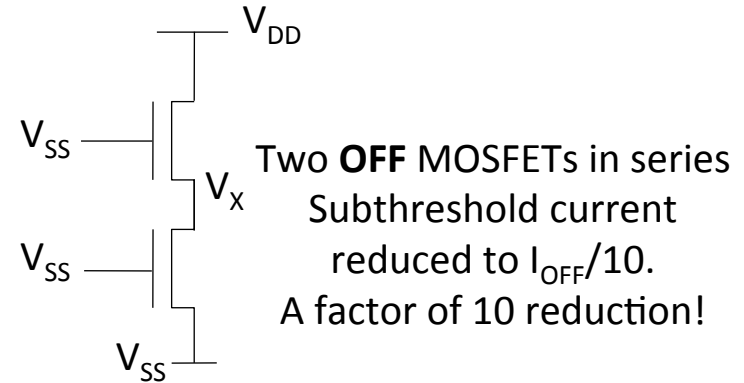
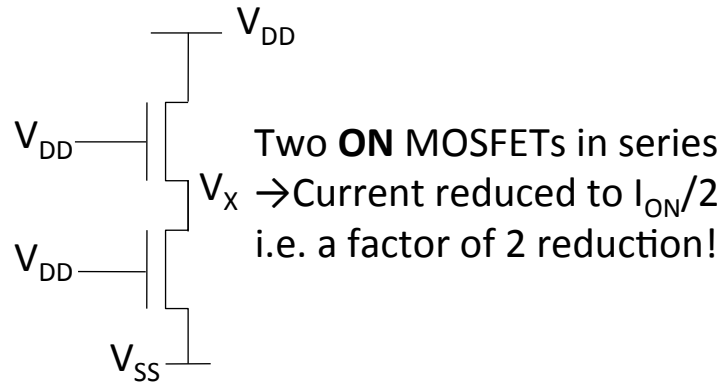
$\sigma = 0.1 \text{ V}^{-1}$

Subthreshold slope:

$S = 100 \text{ mV/decade}$



Subthreshold leakage – stacking effect



How can that be?

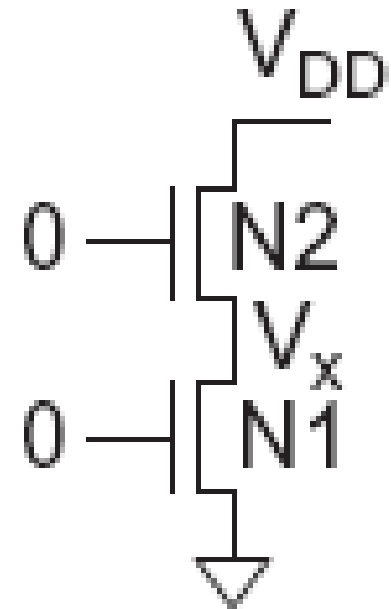
Subthreshold leakage

- Series OFF transistors have less leakage
 - $V_x > 0$, so N2 has negative V_{gs}

$$I_{sub} = I_{OFF} e^{\frac{V_{GS} + \sigma(V_x - V_{DD})}{nV_t}} = I_{OFF} 10^{\frac{V_{GS} + \sigma(V_x - V_{DD})}{S}}$$

$$I_{sub} = \underbrace{I_{off} 10^{\frac{\sigma(V_x - V_{DD})}{S}}}_{N1} = \underbrace{I_{off} 10^{\frac{-V_x + \sigma((V_{DD} - V_x) - V_{DD})}{S}}}_{N2}$$

$$I_{sub} \approx I_{off} / 10$$



- Leakage through 2-stack reduces $\sim 10x$
- Leakage through 3-stack reduces further

2. Gate leakage

- 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 $t_{\text{ox}} > 10.5 \text{ \AA}$
 - High-k gate dielectrics help
 - Some processes provide multiple t_{ox}
 - e.g. thicker oxide for 3.3 V I/O transistors
- Control leakage in circuits by limiting V_{DD}

NAND3 leakage variation due to input pattern

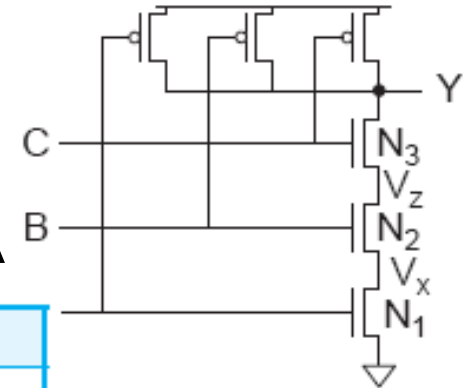
- 100 nm process

Gate leakage: $I_{gn} = 6.3 \text{ nA}$, $I_{gp} = 0$

Subthreshold leakage: $I_{offn} = 5.63 \text{ nA}$, $I_{offp} = 9.3 \text{ nA}$

Input State (ABC)	I_{sub}	I_{gate}	I_{total}	V_x	V_z
000	0.4	0	0.4	stack effect	stack effect
001	0.7	0	0.7	stack effect	$V_{DD} - V_t$
010	0.7	1.3	2.0	intermediate	intermediate
011	3.8	0	3.8	$V_{DD} - V_t$	$V_{DD} - V_t$
100	0.7	6.3	7.0	0	stack effect
101	3.8	6.3	10.1	0	$V_{DD} - V_t$
110	5.6	12.6	18.2	0	0
111	28	18.9	46.9	0	0

All current values are in nA



More than two orders of magnitude difference in leakage current!

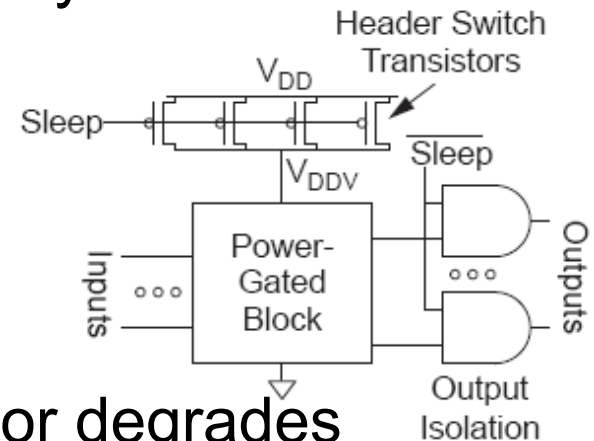
Table 5.2
In W&H

3. Junction leakage

- From reverse-biased p-n junctions
 - Between diffusion and substrate or well
- Ordinary diode leakage is negligible
- Band-to-band tunneling (BTBT) can be significant
 - Especially in high- V_T transistors where other leakage is small
 - Worst at $V_{db} = V_{DD}$
- Gate-induced drain leakage (GIDL) exacerbates
 - Worst for $V_{gd} = -V_{DD}$ (or more negative)

How to reduce: Supply 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



Leakage control

- 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 delay critical circuits
 - Increase V_s : *stack effect*
 - *Input vector control* in sleep
 - Decrease V_b
 - *Reverse body bias* in sleep
 - Or forward body bias in active mode

Summary

Dynamic (active) is mainly due to switching:

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

Static is mainly due to subthreshold leakage :
Exponential dependence

$$I_{sub} \approx I_{off} 10^{\frac{V_{gs} + \sigma(V_{DS} - V_{DD})}{S}}$$

Modern power reduction techniques require knowledge about high-level aspects of the application