

Power Dissipation and Low Power design

Kjell Jeppson

2018 MCC092 lecture

Chalmers University of Technology

First some definitions!

Energy and power in circuit elements

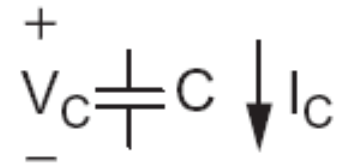
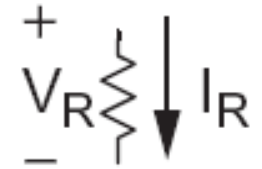
- Power is drawn from voltage sources attached to the V_{DD} pin(s) of a chip in form of a current I_{DD} .
 - This power is dissipated in resistive elements like MOSFETs, or
 - temporarily stored in on-chip capacitors from MOSFET gates and drains, and from interconnect wiring

- Instantaneous Power [Watt (W)]: $P(t) = I(t)V(t)$

- Energy [Joules (J)]:
$$E = \int_0^T P(t) dt$$

- Average Power [Joules/second]:
$$P_{avg} = \frac{E}{T}$$

- Energy stored in a capacitor [J]: E_C



Textbook illustrations

Charging a capacitor

- When the gate output **rises** (V_{OUT} going high)

- Energy drawn from supply :

$$E_{V_{DD}} = \int_0^T P(t) dt = \int_0^T \underbrace{C_{LOAD} \frac{dV_C(t)}{dt}}_{I_C(t)} \times V_{DD} dt = C_{LOAD} V_{DD}^2$$

- Energy stored in capacitor :

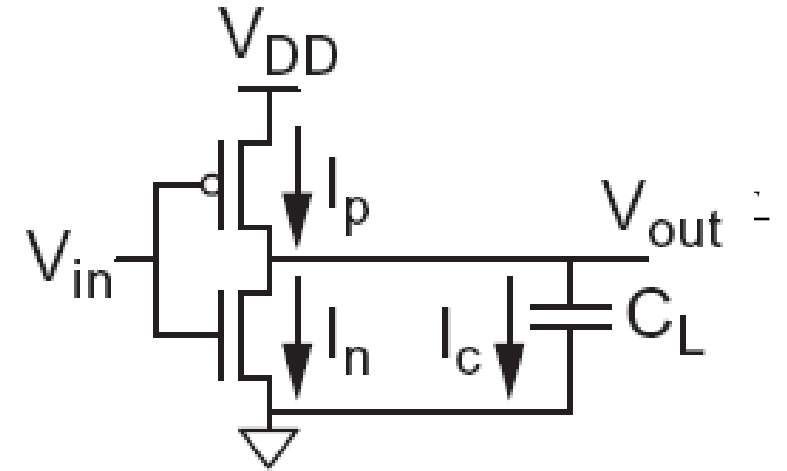
$$E_C = \int_0^T P(t) dt = \int_0^T I_C(t) V_C(t) dt = \int_0^T C_{LOAD} \frac{dV_C(t)}{dt} \times V_C(t) dt = \frac{1}{2} C_{LOAD} V_{DD}^2$$

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

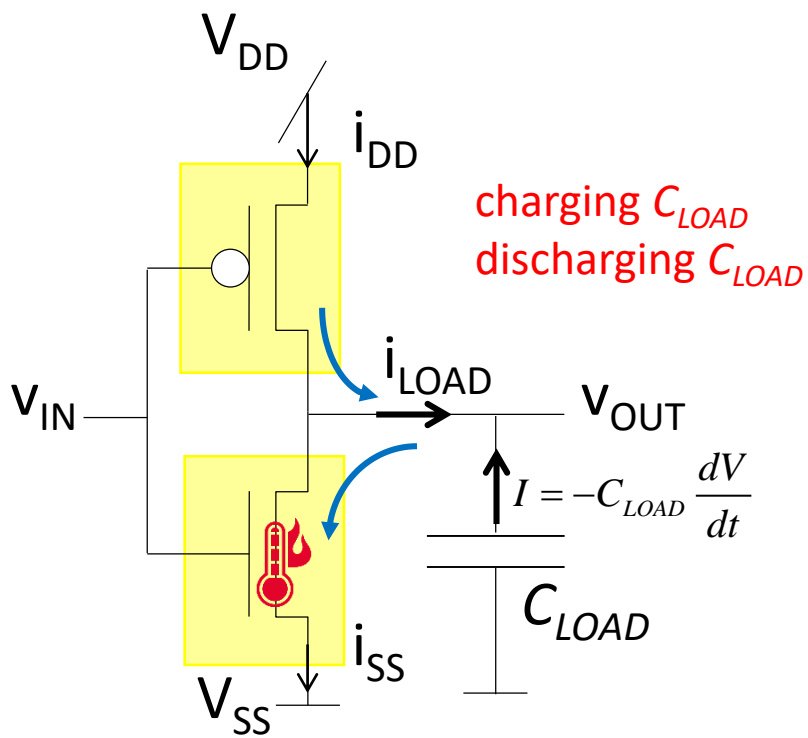
- Half of $E_{V_{DD}}$ is dissipated in the pMOS transistor as heat, other half stored in capacitor

- When the gate output **falls**

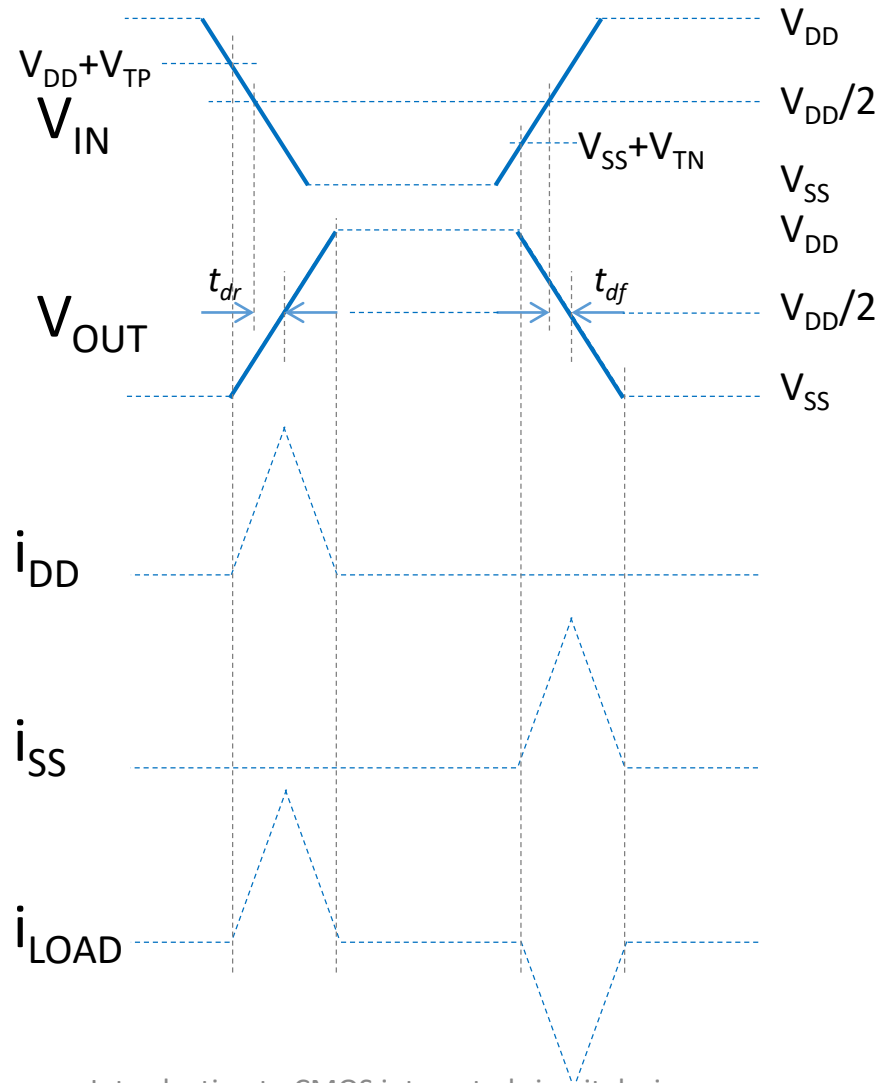
- Energy in capacitor is dumped to GND
- Dissipated as heat in the nMOS transistor



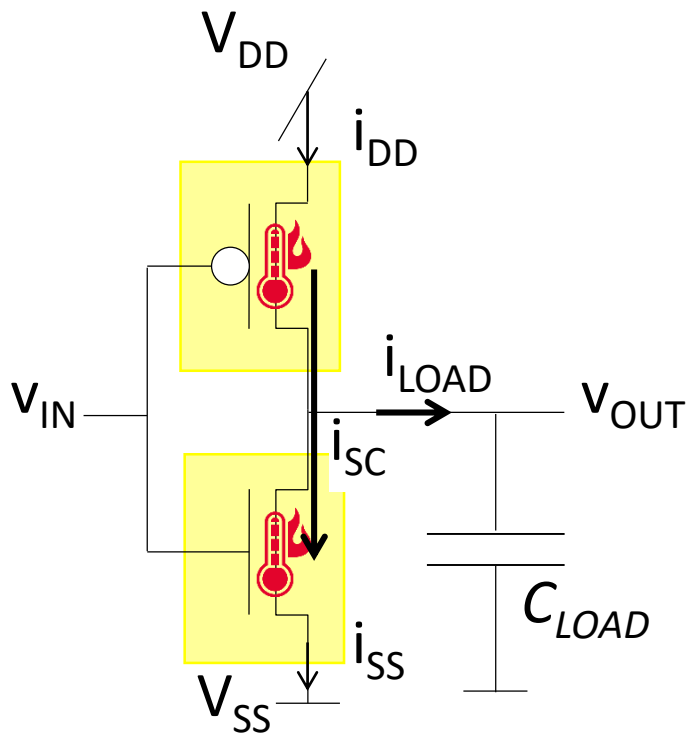
Dynamic Power Dissipation



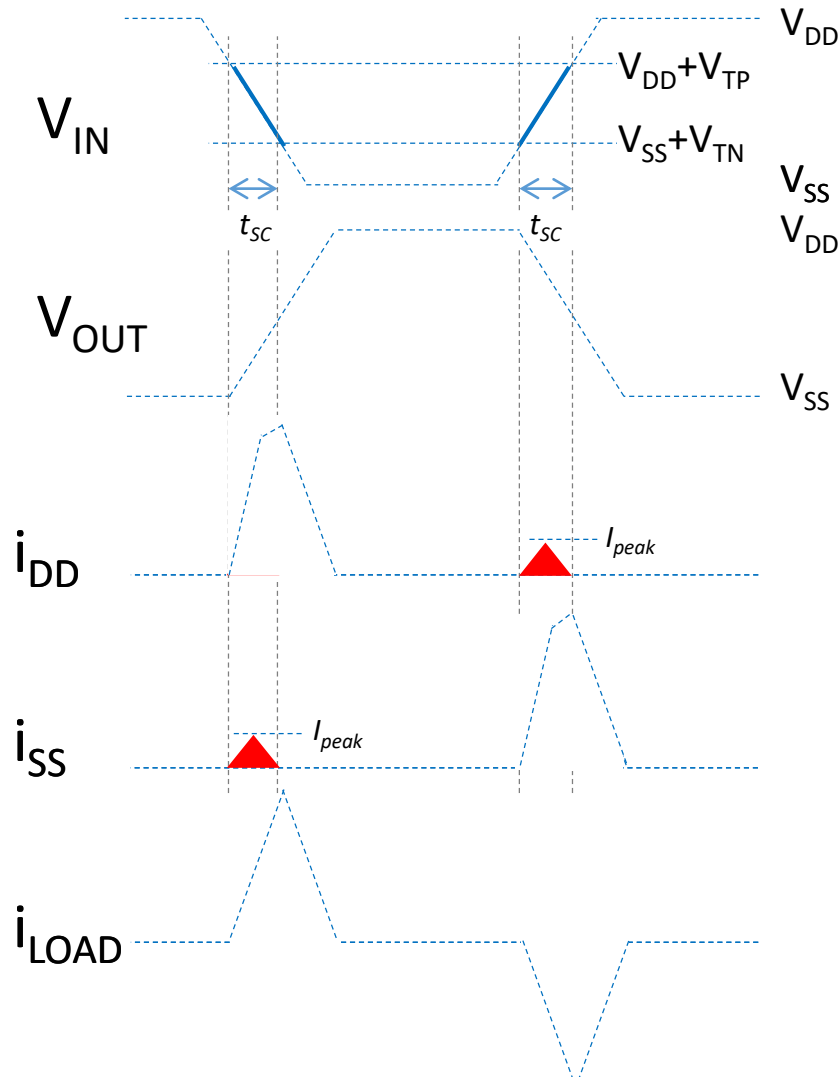
Dynamic power dissipation



Dynamic Power Dissipation – “short-circuit” current



Dynamic power dissipation



When transistors switch, both n-MOSFET and p-MOSFET may be momentarily ON simultaneously

Leads to a blip of “short-circuit” current.

Energy lost in “short-circuit”

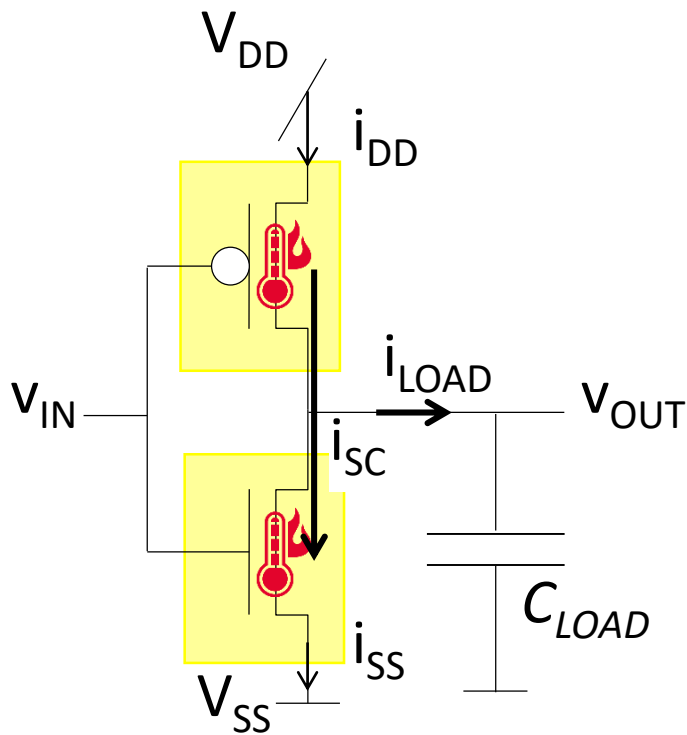
$$E_{SC} = \int_0^T i_{SC}(t) V_{DD} dt$$

$$\approx I_{peak} V_{DD} t_{SC} \ll C_{LOAD} V_{DD}^2$$

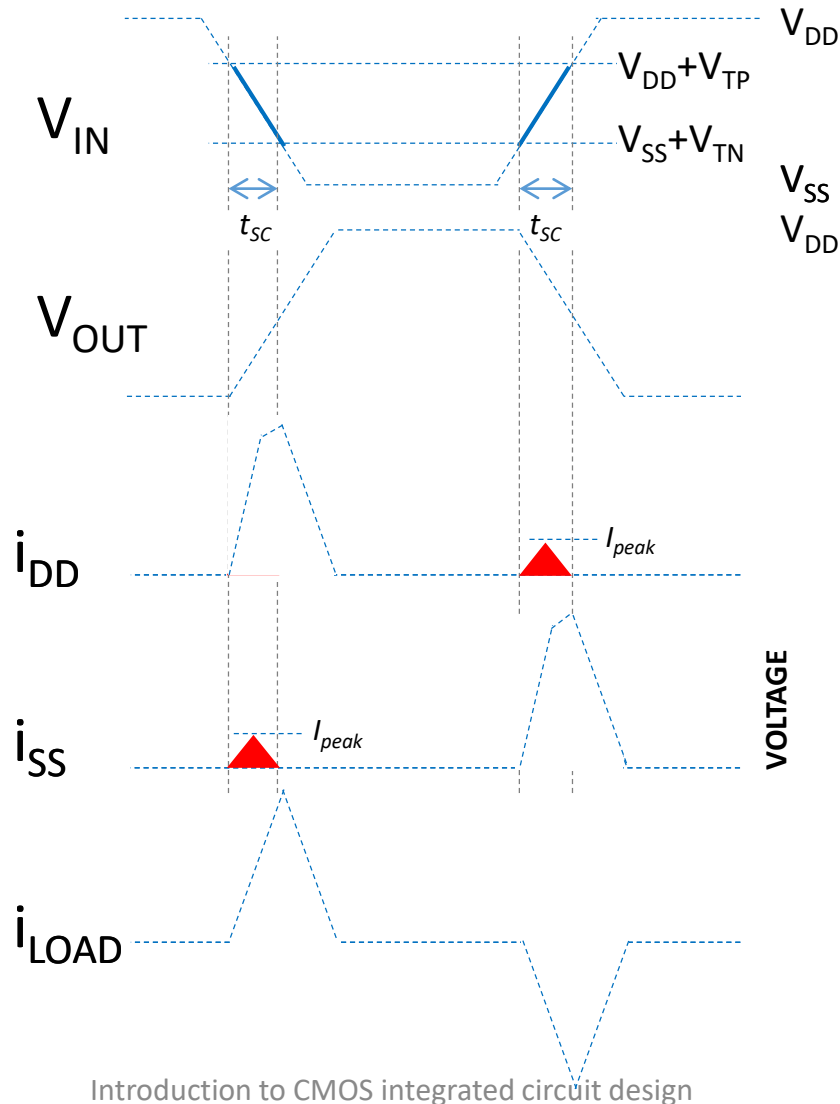
$E_{SC} < 10\%$ of the total dynamic energy if input/output rise & fall times are comparable

We will generally ignore this component of power dissipation

Dynamic Power Dissipation – “short-circuit” current

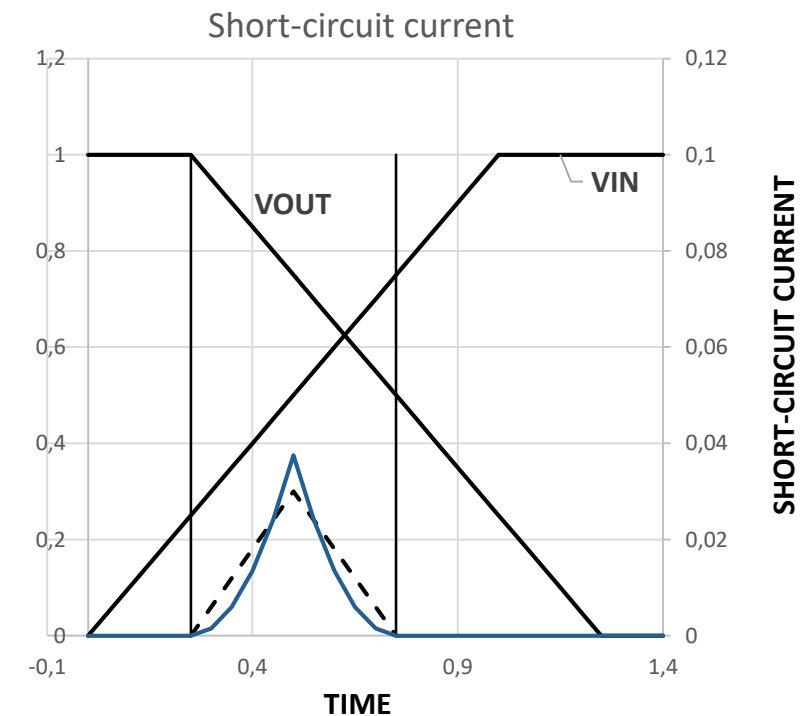


Dynamic power dissipation



When transistors switch, both n-MOSFET and p-MOSFET may be momentarily ON simultaneously

Leads to a blip of “short-circuit” current.



Dynamic power dissipation – in summary

- Power is drawn from the voltage source that is attached to the V_{DD} pin(s) of a chip.
 - This power is dissipated in resistive elements like MOSFETs, or temporarily stored in circuit node capacitors, C_k
- Average power dissipation in node k : $P_k = \frac{E_k}{T} = C_k V_{DD}^2 \times f_{sw} = \alpha_k C_k V_{DD}^2 f_{CLK}$
Energy
- Effective circuit capacitance
 - $C_{eff} = \alpha C_{CHIP}$; total chip cap. $C_{CHIP} = \sum_{k=1}^{\infty} C_k$; activity factor $\alpha = \left(\sum_{k=1}^{\infty} \alpha_k C_k \right) / C_{CHIP}$
- Activity factor - alpha
 - For the clock distribution network that is charged **and** discharged every clock cycle, the activity factor is $\alpha_{clock}=1$
 - For the input of a logic gate going high **or** low during a clock cycle, then $\alpha_{input}=1/2$.

Dynamic power – an example

What if we have a 1 billion transistor chip w/

- 50M logic transistors: $W_{avg}=300$ nm, $\alpha=0.1$
- 950M memory transistors: $W_{avg}=100$ nm, $\alpha=0.02$
- Assume a 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 capacitances and short-circuit currents.

1 billion transistor chip

50M logic transistors

Average width: 300 nm, activity factor: 0.1

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

950M memory transistors

Average width: 100 nm, activity factor: 0.02

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

$$P_{\text{dynamic}} = (\alpha_{\text{logic}} C_{\text{logic}} + \alpha_{\text{mem}} C_{\text{mem}}) \times f \times V_{DD}^2 =$$

$$= \left(\underbrace{0.1 C_{\text{logic}}}_{2.7} + \underbrace{0.02 C_{\text{mem}}}_{3.4} \right) \times (1.0 \text{ GHz}) \times (1.0 \text{ V})^2 = 6.1 \text{ W}$$

global activity factor: $\alpha = (\alpha_{\text{logic}} C_{\text{logic}} + \alpha_{\text{mem}} C_{\text{mem}}) / C_{\text{tot}} \approx 6.1 / 200 = 3.05\%$,

where $C_{\text{tot}} = C_{\text{logic}} + C_{\text{mem}} \approx 200$ nF.

Activity factor

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

$$P_{sw} = \alpha C_{CHIP} V_{DD}^2 f_{CLK}$$

Activity factor -Definition

- alpha (α) = Probability that output switches from 0 to 1
- P_i = probability that node i is 1
- $\overline{P_i} = 1 - P_i$ = 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 $\alpha = 0.25$!
- 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 $\alpha \approx 0.1$

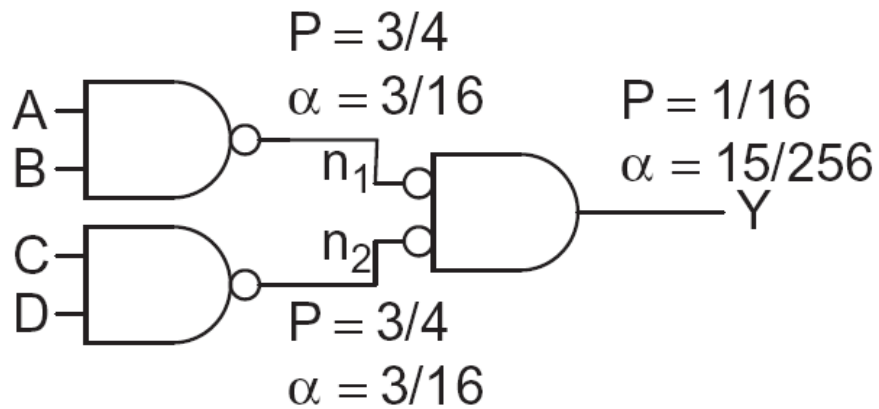
Switching probability

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 $P = 0.5$ i.e. it is just as likely that the input is a one as a zero



2-input NAND

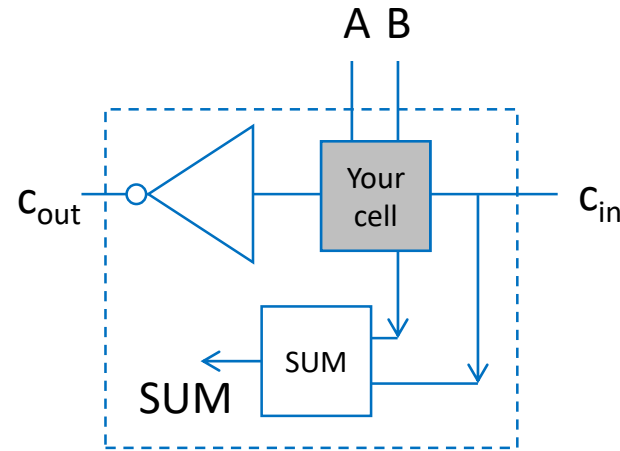
$$P_Y = P_A * P_B$$
$$\alpha = P_Y * (1 - P_Y)$$

2-input NOR

$$P_Y = (1 - P_A) * (1 - P_B)$$
$$\alpha = P_Y * (1 - P_Y)$$

de Morgan's theorem $Y = \overline{\overline{AB} + \overline{CD}} = ABCD$

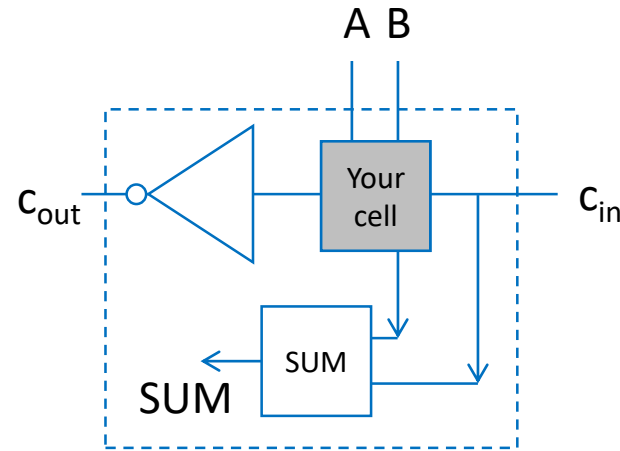
Lab session carry chain



A	B	CIN	COU	SUM
0	0	0	0	0
0	0	1	0	1
1	0	0	0	1
1	0	1	1	0
0	1	0	0	1
0	1	1	1	0
1	1	0	1	0
1	1	1	1	1

	a7	b7	a6	b6	a5	b5	a4	b4	a3	b3	a2	b2	a1	b1	a0	b0
P input=1	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
input alpha	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
A B input caps	28	28	28	28	28	28	28	28	28	28	28	28	28	28	28	28
output node cap	10	14	10	14	10	14	10	14	10	14	10	14	10	14	10	14
output node parasitic	10	28	10	28	10	28	10	28	10	28	10	28	10	28	10	28
sum output node cap	20	42	20	42	20	42	20	42	20	42	20	42	20	42	20	42
input alpha*cap	7,00	7,00	7,00	7,00	7,00	7,00	7,00	7,00	7,00	7,00	7,00	7,00	7	7	7	7
P output=1	0,50	0,50	0,50	0,50	0,50	0,50	0,50	0,50	0,50	0,50	0,50	0,50	0,5	0,5	0,5	0,5
output alpha	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
output alpha*cap	5,00	10,50	5,00	10,50	5,00	10,50	5,00	10,50	5,00	10,50	5,00	10,50	5	10,5	5	10,5
inverter driving capability	x=	10			f=	1	GHz		alpha*C=	85	fF					
logic gate driving capability	y=	7			VDD=	1,2			total cap=	340	fF					
	pinv=	1							alpha=	0,25						
	CIN=	0,36	fF/X						Power=	alpha*C*f*VDD*VDD=					122	uW

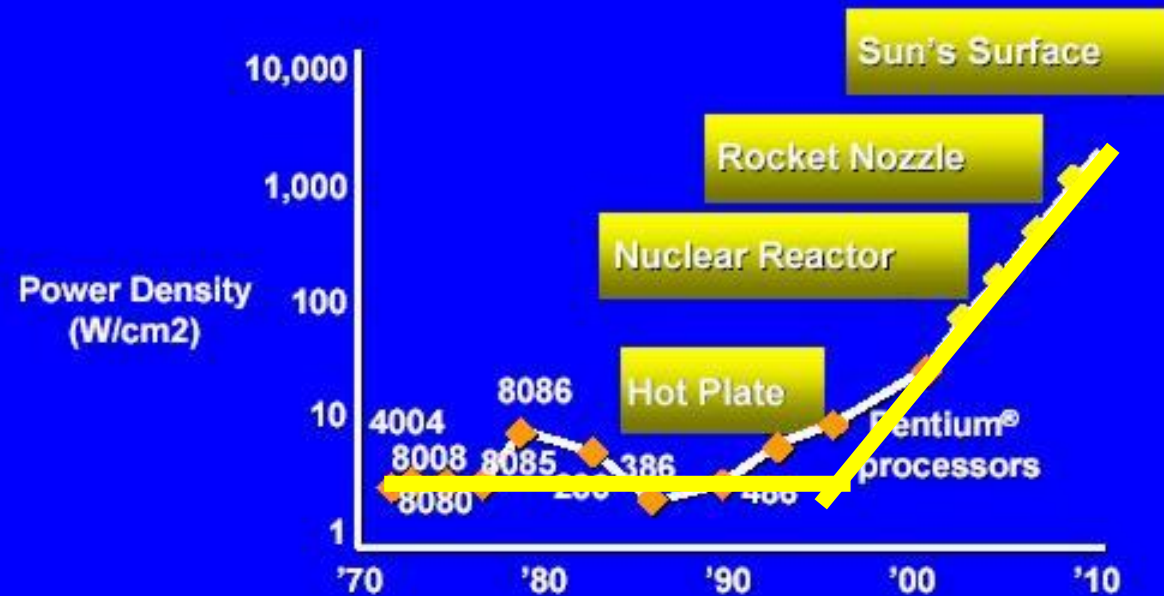
Lab session carry chain



A	B	CIN	COU	SUM
0	0	0	0	0
0	0	1	0	1
1	0	0	0	1
1	0	1	1	0
0	1	0	0	1
0	1	1	1	0
1	1	0	1	0
1	1	1	1	1

	a7	b7	a6	b6	a5	b5	a4	b4	a3	b3	a2	b2	a1	b1	a0	b0	CIN
P input=1	0,105	0,105	0,131	0,131	0,164	0,164	0,205	0,205	0,256	0,256	0,32	0,32	0,4	0,4	0,5	0,5	0,5
input alpha	0,094	0,094	0,114	0,114	0,137	0,137	0,163	0,163	0,19	0,19	0,218	0,218	0,24	0,24	0,25	0,25	0,25
A B input caps	28	28	28	28	28	28	28	28	28	28	28	28	28	28	28	28	
output node cap	10	14	10	14	10	14	10	14	10	14	10	14	10	14	10	14	
output node parasitic	10	28	10	28	10	28	10	28	10	28	10	28	10	28	10	28	
sum output node cap	20	42	20	42	20	42	20	42	20	42	20	42	20	42	20	42	
input alpha*cap	2,63	2,63	3,19	3,19	3,84	3,84	4,56	4,56	5,33	5,33	6,09	6,09	6,72	6,72	7	7	
P output=1	0,02	0,98	0,03	0,97	0,05	0,95	0,10	0,90	0,17	0,83	0,28	0,72	0,4	0,6	0,5	0,5	
output alpha	0,02	0,02	0,03	0,03	0,05	0,05	0,09	0,09	0,14	0,14	0,20	0,20	0,24	0,24	0,25	0,25	
output alpha*cap	0,32	0,68	0,57	1,20	1,01	2,13	1,76	3,70	2,83	5,95	4,00	8,40	4,8	10,08	5	10,5	
inverter driving capability	x=	10			f=	1	GHz		alpha*C=	51	fF						
logic gate driving capability	y=	7			VDD=	1,2			total cap=	340	fF						
	pinv=	1							alpha=	0,15							
	CIN=	0,36	fF/X						Power=	alpha*C*f*VDD*VDD=						73 uW	

Power Density Extrapolation



Gelsinger's Slide from ISSCC 2001

12

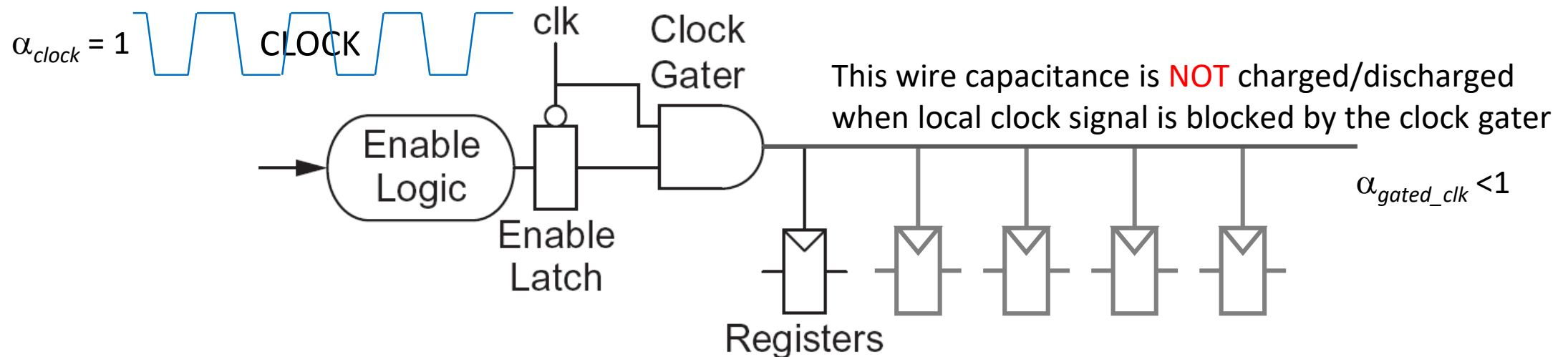


Power dissipation – how can it be reduced?

- Lower the power supply voltage ($P \sim V_{DD}^2$) – but power gives speed
- Minimize C_{LOAD}
 - Minimize gate capacitance
 - Fewer stages of logic and small gate sizes
 - Minimize wire capacitance
 - Good floor planning to keep communicating blocks close to each other
 - Careful drive of long wires (repeaters?)
- Reduce clock frequency – but again we want speed
- Reduce activity factor – no unnecessary switching!

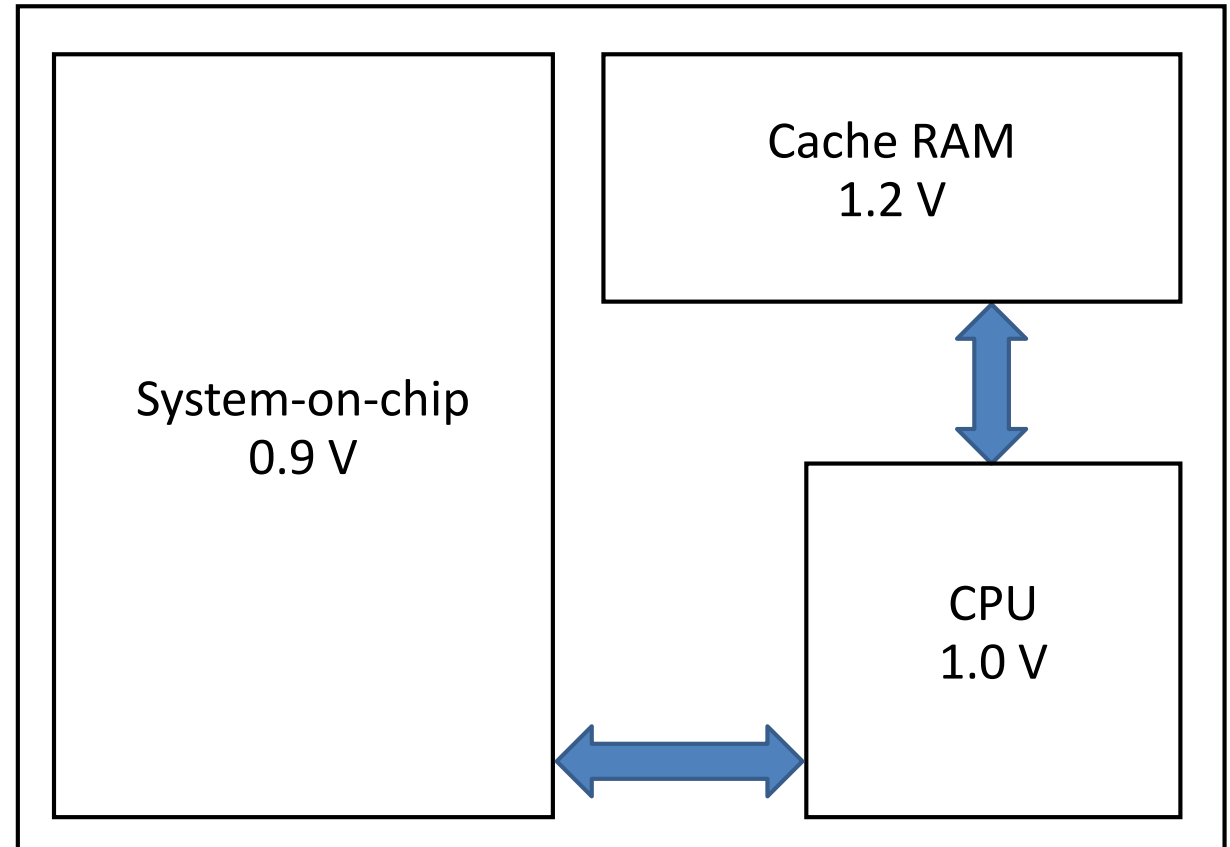
Clock gating

- The most efficient way to reduce the switching activity is to turn off the clock to registers in unused blocks – sleep mode
 - (+) Saves clock activity ($\alpha_{clock} = 1$)
 - (+) Eliminates all switching activity in clock-gated blocks
 - (-) Requires determining if block will be used



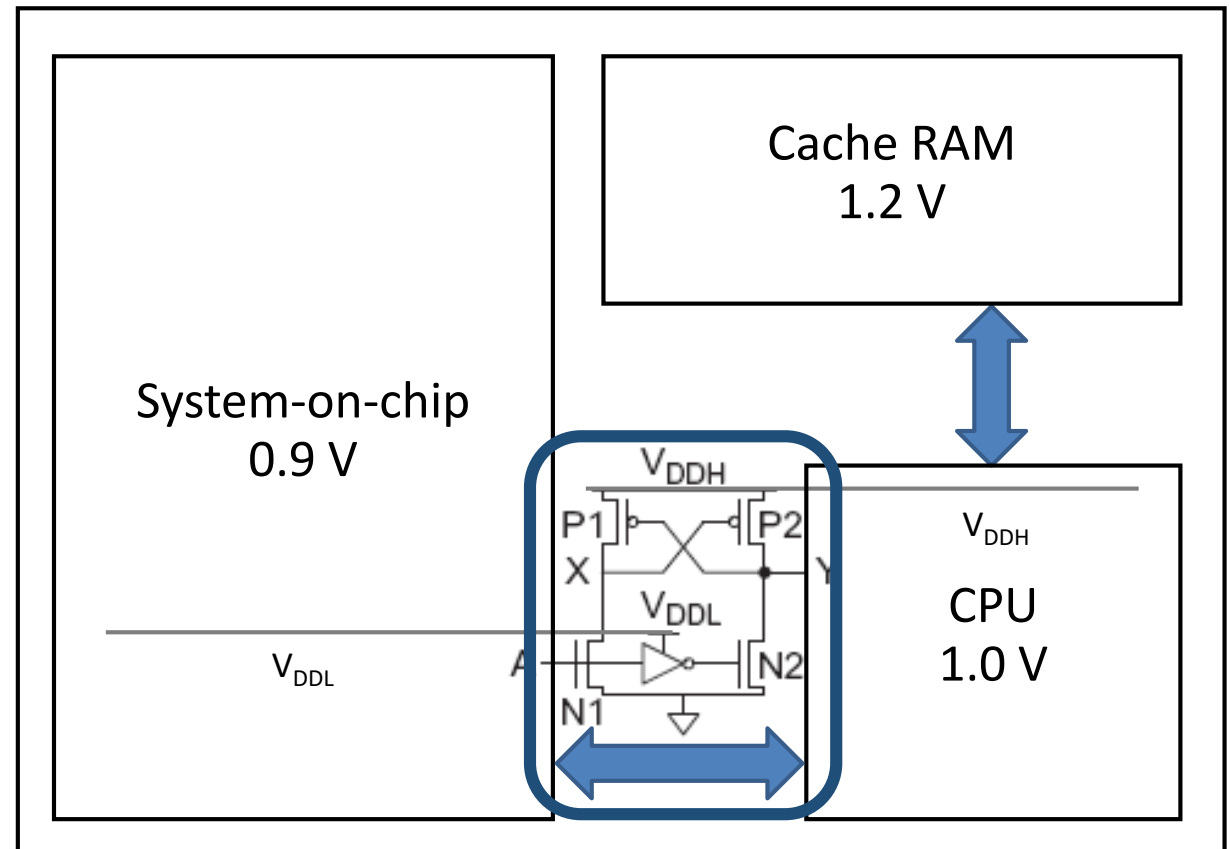
Reduce supply voltage – Multi-VDD

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



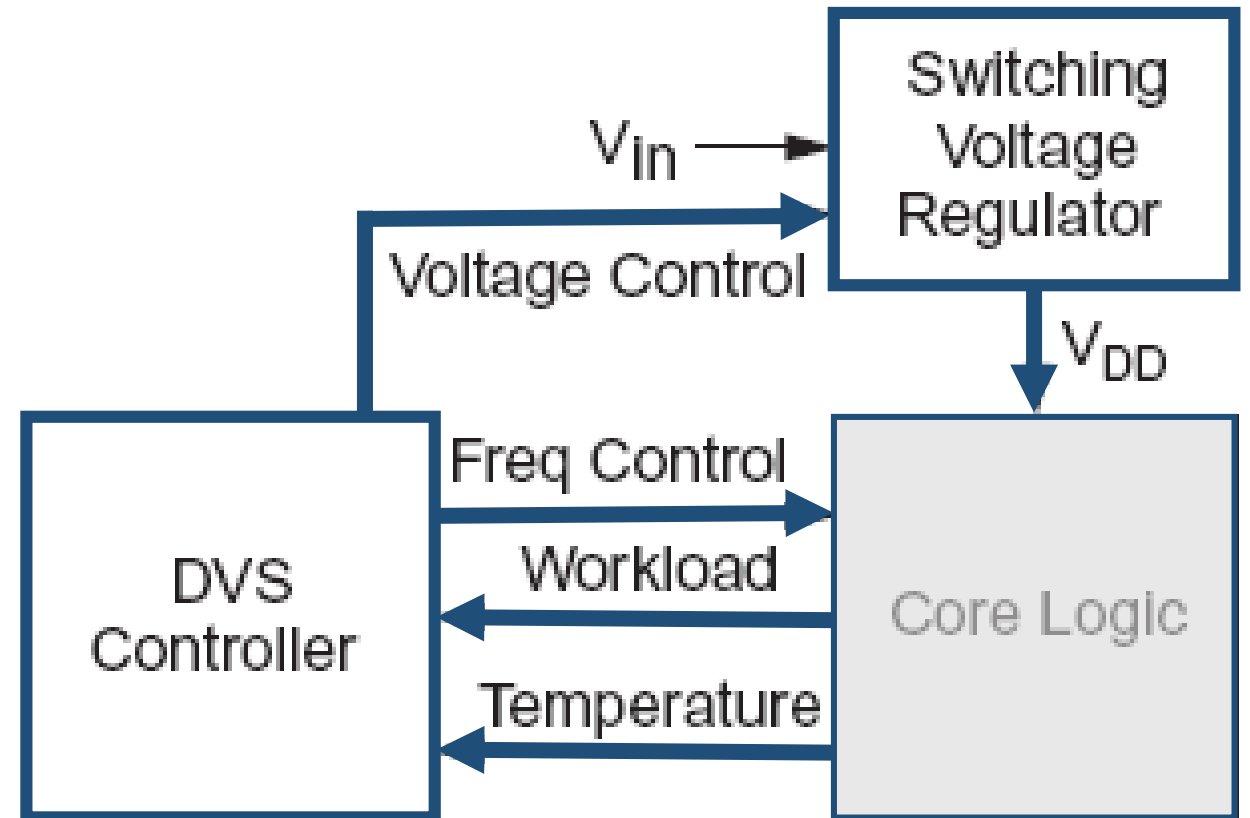
Reduce supply voltage – Multi-VDD

- 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



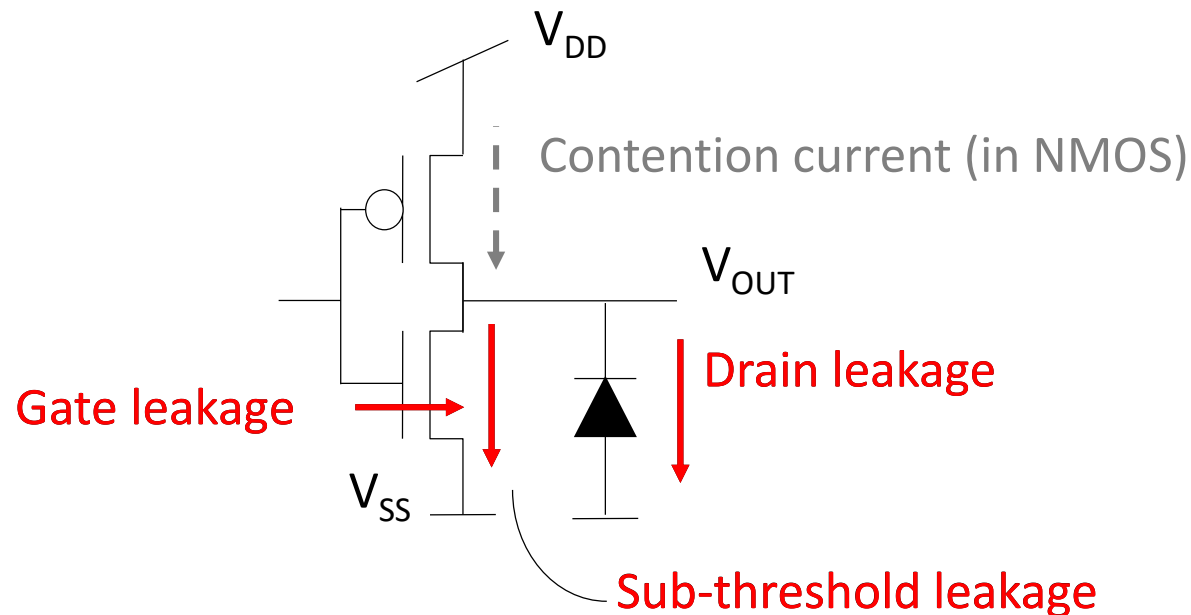
Dynamic voltage scaling - DVS

- 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



Static power dissipation

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

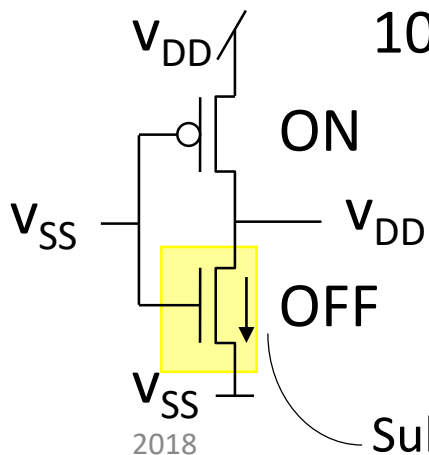


Gate leakage

- 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
- Keep thickness of gate dielectric $>10.5 \text{ \AA}$ to prevent tunneling
 - High-k gate dielectrics help
 - Some processes provide multiple oxide thicknesses t_{ox}
 - e.g. thicker oxide for 3.3 V I/O FETs
- Control leakage in circuits by limiting V_{DD}

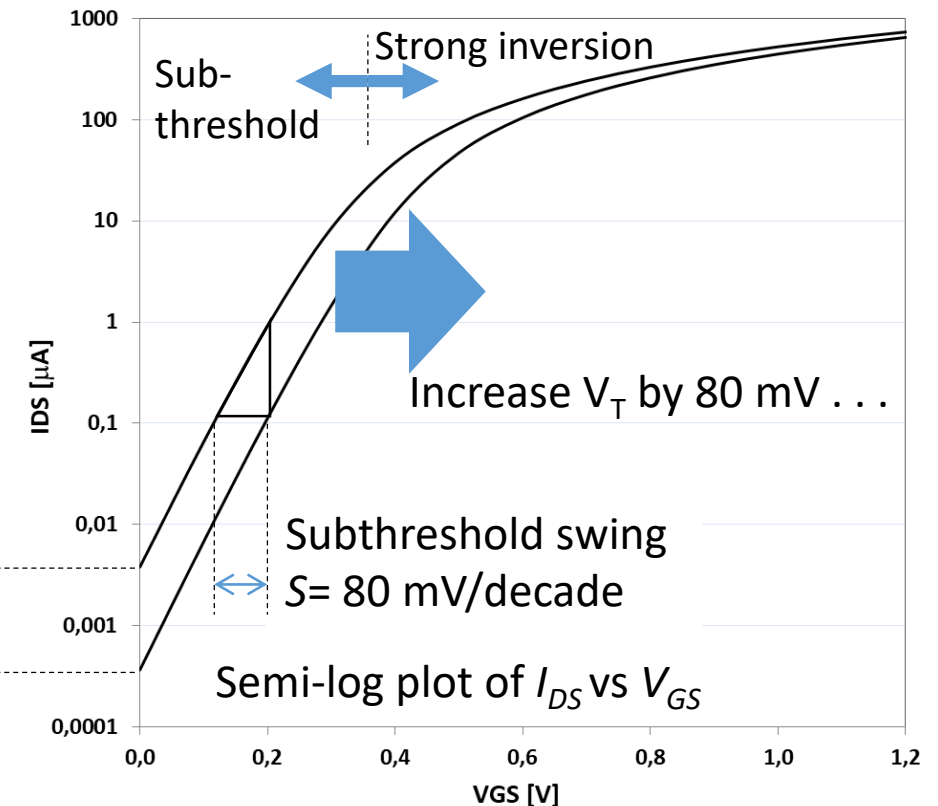
Reducing sub-threshold leakage: Multi-VT

- As geometries have shrunk to 65 nm and below, libraries with multiple VT has become a common way of reducing leakage currents
- This semi-log plot of I_{DS} vs V_{GS} illustrates the sub-VT leakage
- If a device technology has a subthreshold swing of 80 mV per decade
 - an 80 mV increase in V_T results in a 10X reduction of sub-VT leakage!



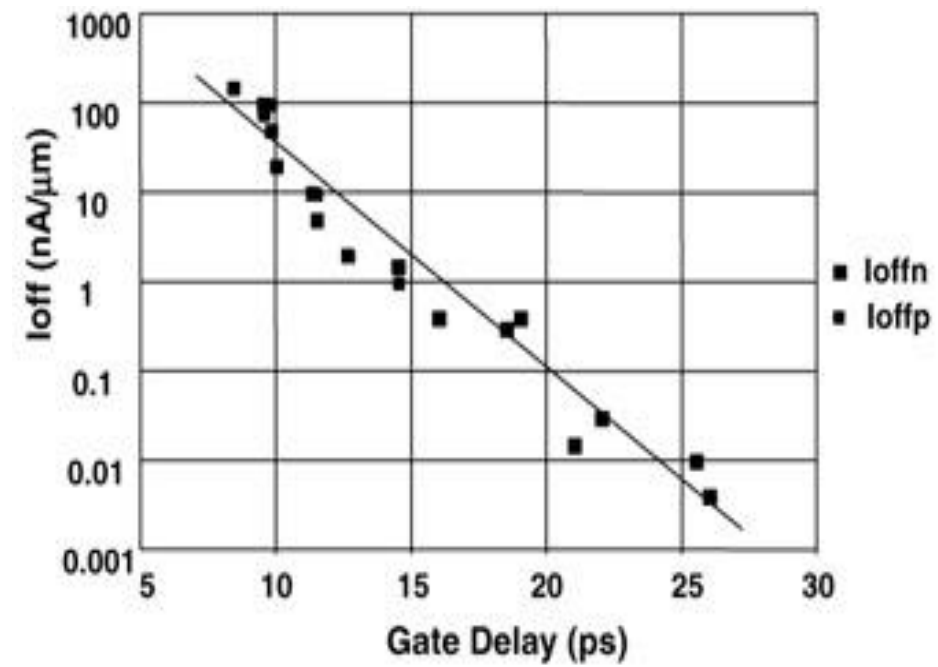
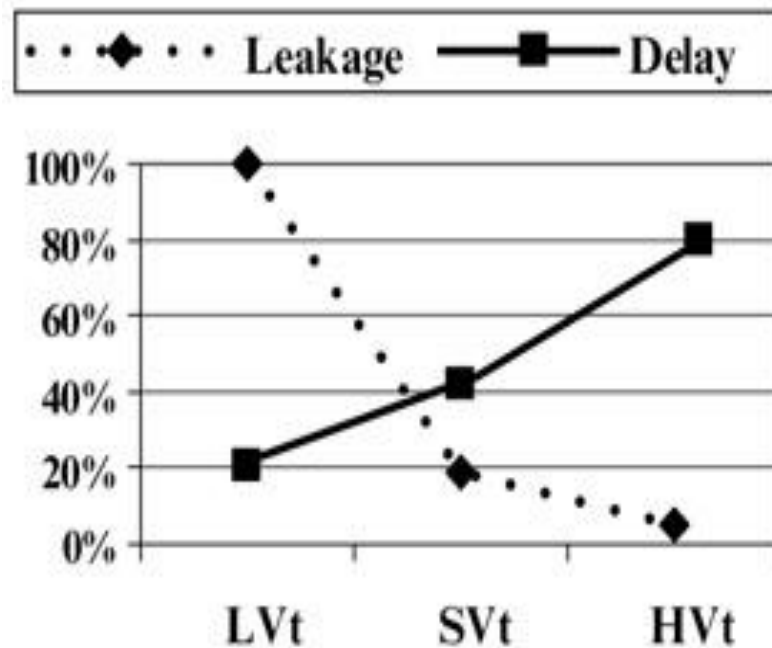
... to get a factor of 10 in leakage reduction

Subthreshold leakage: I_{OFF}



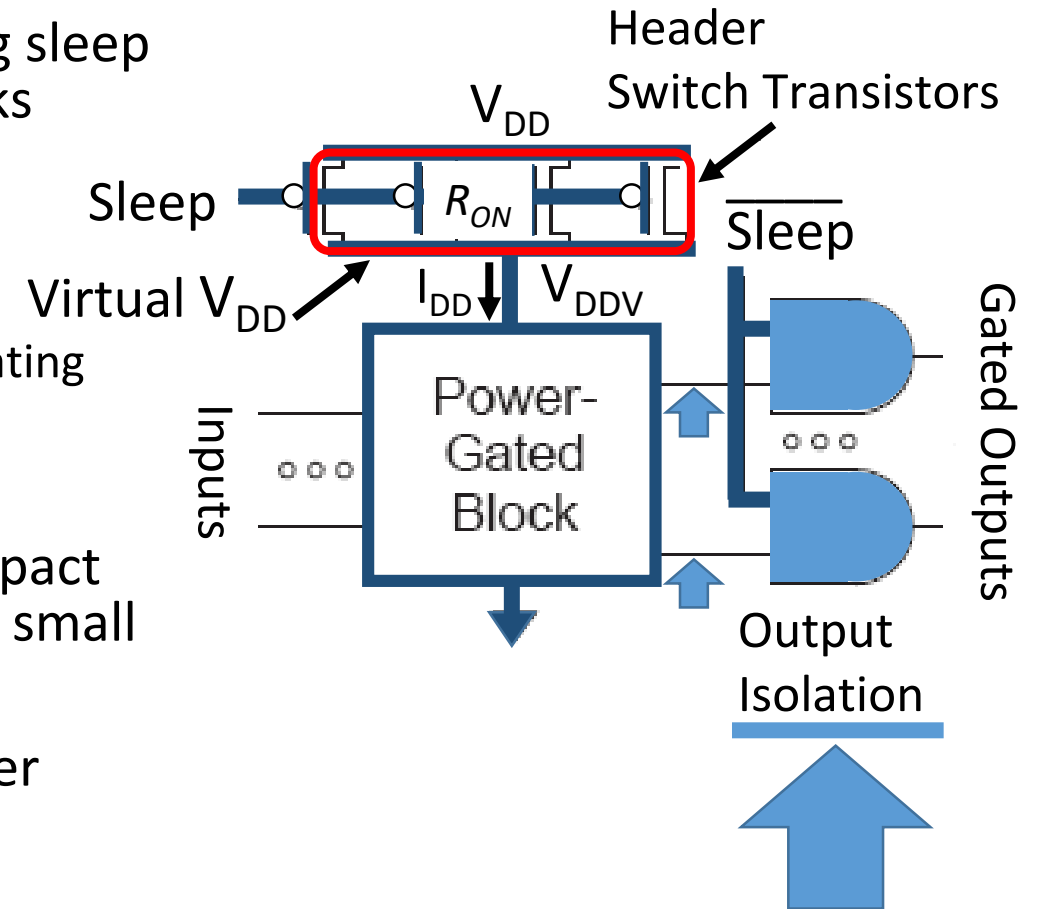
Reducing sub-threshold leakage: Multi-VT

- As geometries have shrunk to 65 nm and below, libraries with multiple VT has become a common way of reducing leakage currents



Power gating

- Turn OFF power to idle blocks to save leakage using sleep transistors between V_{DD} and the power-gated blocks
 - Use virtual V_{DD} (V_{DDV}) for power gated block
 - Gate block outputs to prevent invalid logic levels to next block
 - Cannot leave outputs from unpowered blocks floating
- Voltage drop across sleep transistors degrades performance during normal operation
 - Size the transistor wide enough to minimize impact by keeping difference between supply voltages small
 - $V_{DD} - V_{DDV} = R_{ON} \times I_{DD}$
- Switching wide sleep transistor costs dynamic power
 - Only justified when circuit sleeps long enough



Static power: textbook example 4.4 (5.4)

Let's have a look again at our 1 billion transistor chip w/

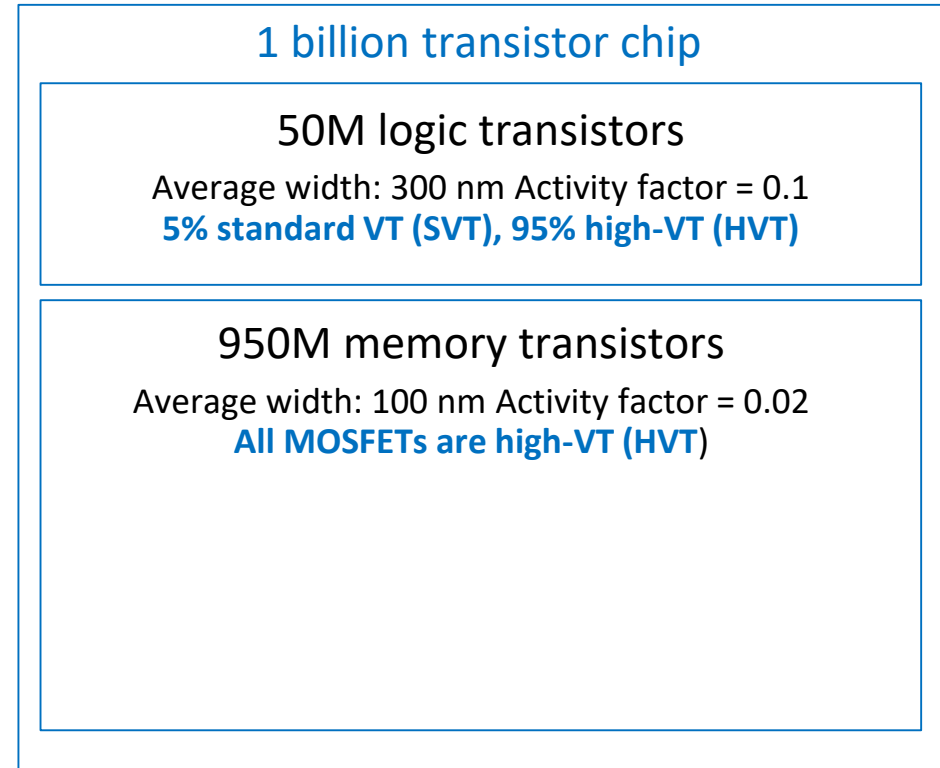
- 50M logic transistors: $W_{avg}=300$ nm, $\alpha=0.1$
- 950M memory transistors: $W_{avg}=100$ nm, $\alpha=0.02$
- Assume a 1.0 V, 65 nm process
- $C = 1$ fF/ μ m (gate) + 0.8 fF/ μ m (diffusion)

What would be its static power consumption?

All MOSFETs are High-VT, except 5% of logic MOSFETs on critical timing paths that are Standard-VT

Leakage information

- Subthreshold leakage:
 - SVT: 100 nA/ μ m; HVT: 10 nA/ μ m
- Gate leakage: 5 nA/ μ m;
- Junction leakage: negligible



Static power: textbook example 4.4 (5.4)

SOLUTION:

- Total gate width of 1 billion MOSFETs is $15+95=110$ m
 - Total channel width of 2.5M Standard-VT logic MOSFETs is $W_{SVT}=2.5M \times 0.3\mu m=0.75$ m
 - Total channel width of high-VT MOSFETs is $W_{HVT}=109.25$ m
- Subthreshold leakage (assuming 50% are turned OFF):
 - $I_{sub} = \underbrace{\frac{1}{2}W_{SVT} \times 100 \text{ nA}/\mu m}_{\text{SVT leakage } 37.5 \text{ mA}} + \underbrace{\frac{1}{2}W_{HVT} \times 10 \text{ nA}/\mu m}_{\text{HVT leakage } 546 \text{ mA}} \approx 585 \text{ mA}$
- Gate leakage (again assuming 50% are turned OFF):
 - $I_{gate} = \frac{1}{2} \underbrace{(W_{SVT} + W_{HVT})}_{\text{total width } 110 \text{ m}} \times \underbrace{5 \text{ nA}/\mu m}_{\text{gate leakage}} = 275 \text{ mA}$
- Total static power dissipation:
 - $P_{static} = \underbrace{585 \text{ mA} + 275 \text{ mA}}_{I_{sub} + I_{gate}} \times \underbrace{1.0 \text{ V}}_{V_{DD}} = 860 \text{ mW}$
- These 860 mW of static power dissipation is about 14% of the 6.1 W dynamic power we calculated previously and will deplete batteries of handheld devices rapidly!

1 billion transistor chip

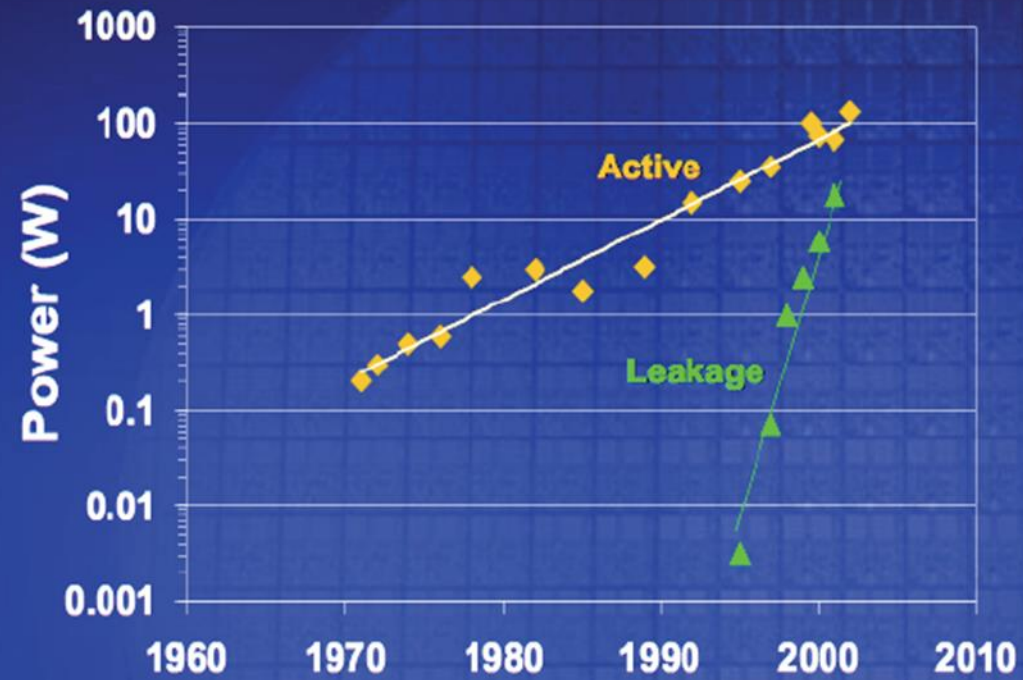
50M logic transistors

Average width: 300 nm Activity factor = 0.1
5% standard VT (SVT), 95% high-VT (HVT)
Total logic static power: ~150 mW

950M memory transistors

Average width: 100 nm Activity factor = 0.02
All MOSFETs are high-VT (HVT)
Total width: 95 m
Subthreshold leakage: 475 mW
Gate leakage: half of sub-VT leakage ~235 mW
Total memory static power: ~710 mW

Processor Power (Watts) - Active & Leakage



From DAT093

Textbook example 4.3 (5.3)

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

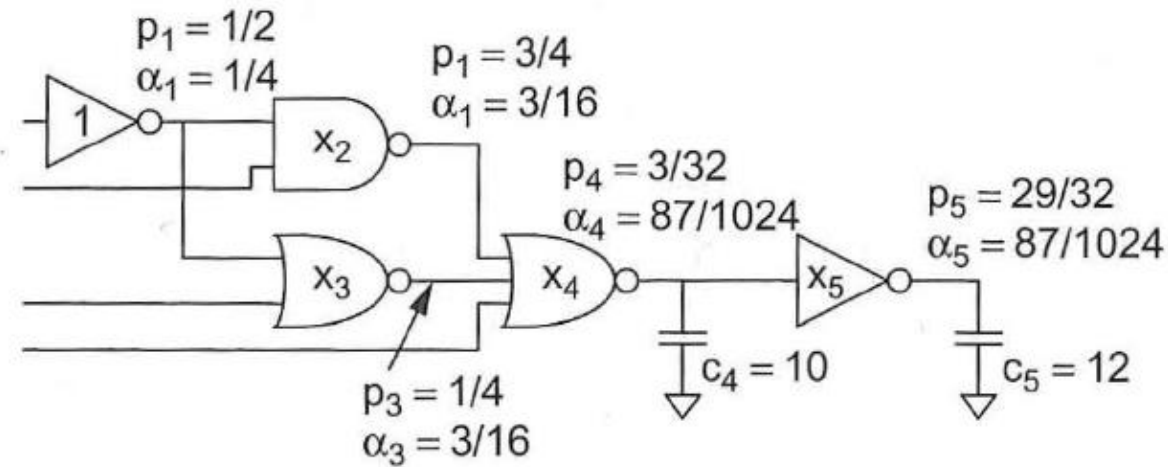


FIGURE 4.12 Activity factors

Textbook example 4.3 (5.3)

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

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

x=drive

$$d = p + g \cdot \frac{g'x'}{gx} = p + g' \cdot \frac{x'}{x}$$

$$C_{IN} = gx; C_{OUT} = px$$

$$g_2 = 4/3$$

$$p_2 = 2$$

$$g_3 = 5/3$$

$$p_3 = 2$$

$$g_4 = 7/3$$

$$p_4 = 3$$

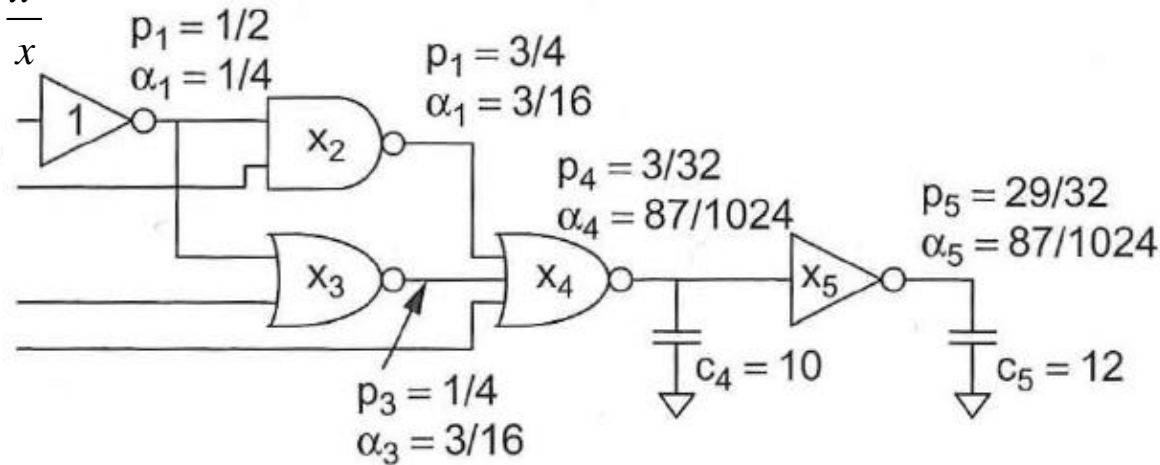
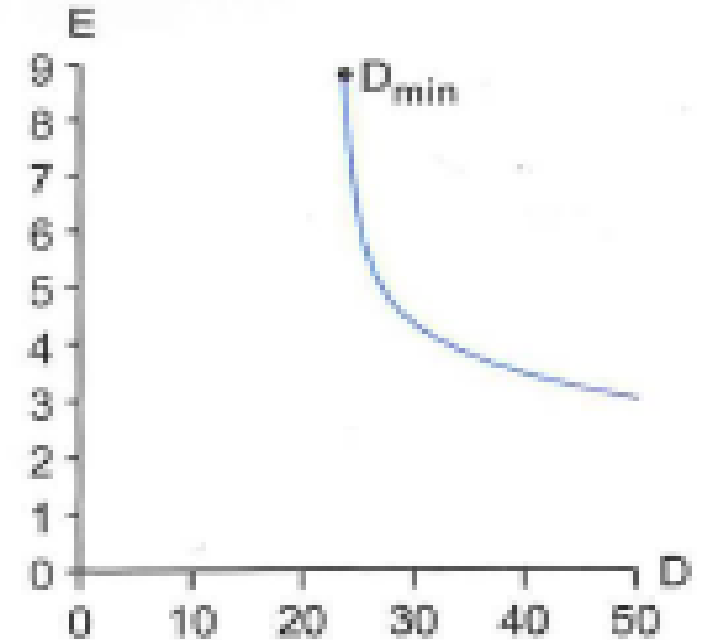


FIGURE 4.12 Activity factors



Textbook example 4.3 (5.3)

$$E = \frac{1}{4}(1 + x_2 + x_3) + \frac{3}{16}\left(\frac{2}{4/3}x_2 + x_4\right) + \frac{3}{16}\left(\frac{2}{5/3}x_3 + x_4\right) + \frac{87}{1024}\left(10 + \frac{9}{7}x_4 + x_5\right) + \frac{87}{1024}(12 + x_5)$$

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

x = input size

$$d = p + g \frac{x'}{x}$$

$$C_{IN} = x; C_{OUT} = \frac{p}{g}x$$

$$g_2 = 4/3$$

$$p_2 = 2$$

$$g_3 = 5/3$$

$$p_3 = 2$$

$$g_4 = 7/3$$

$$p_4 = 3$$

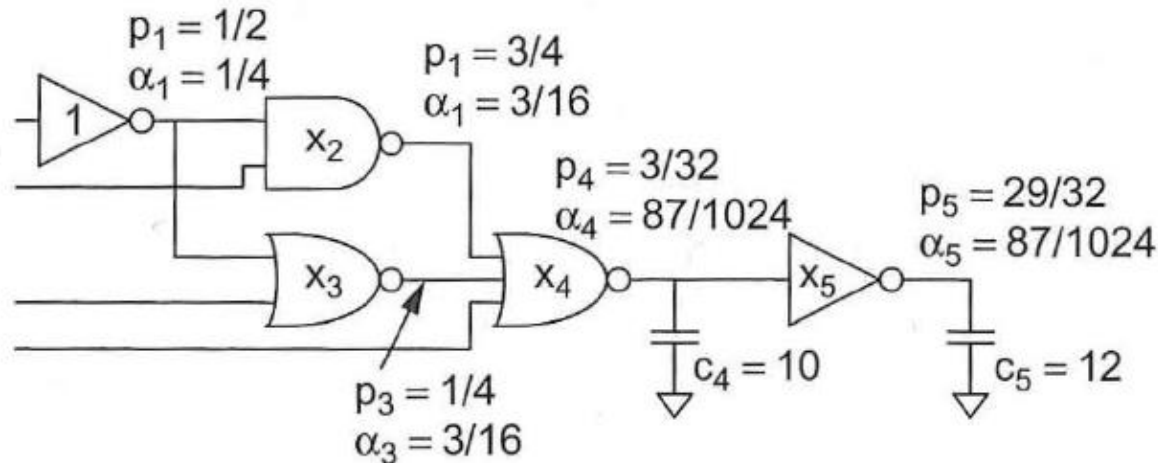
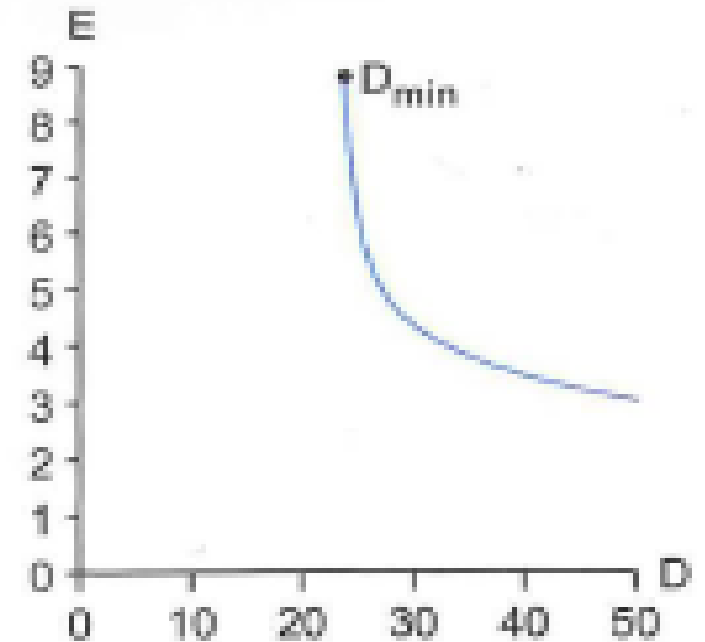


FIGURE 4.12 Activity factors



In summary

- Dynamic power dissipation
 - Where the main source is repeatedly charging the capacitive circuit nodes
 - blips of "short-circuit" current
- Static power dissipation
 - Drain junction leakage (often negligible)
 - Gate leakage (up to 1/3 of static power)
 - Subthreshold leakage (main source)
- Low power design – how can power dissipation be reduced?
 - Clock gating for reducing switching activity
 - Multi-VDD - divide chips into voltage domains using reduced supply voltages
 - DVS -Dynamic voltage scaling – monitor work load and temperature
 - Multi-VT - Reducing sub-threshold leakage
 - Power gating - Turning power OFF to idle blocks sending them into sleep mode