

ESE3700: Circuit-Level Modeling, Design, and Optimization for Digital Systems

Lec 7: February 12, 2025

Layout and Area, MOS Scaling





Today

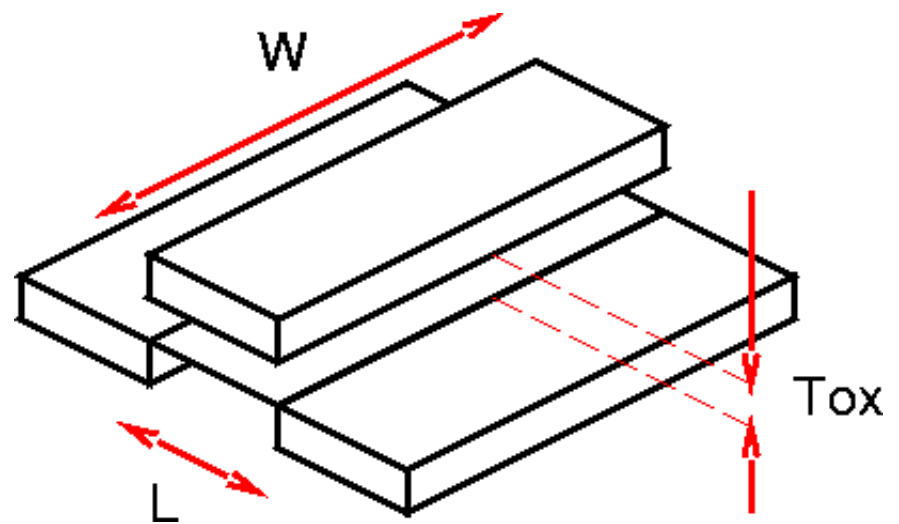
- Layout
 - Transistors
 - Gates
- Design rules
- Standard cells
- VLSI Scaling Trends/Disciplines



Transistor



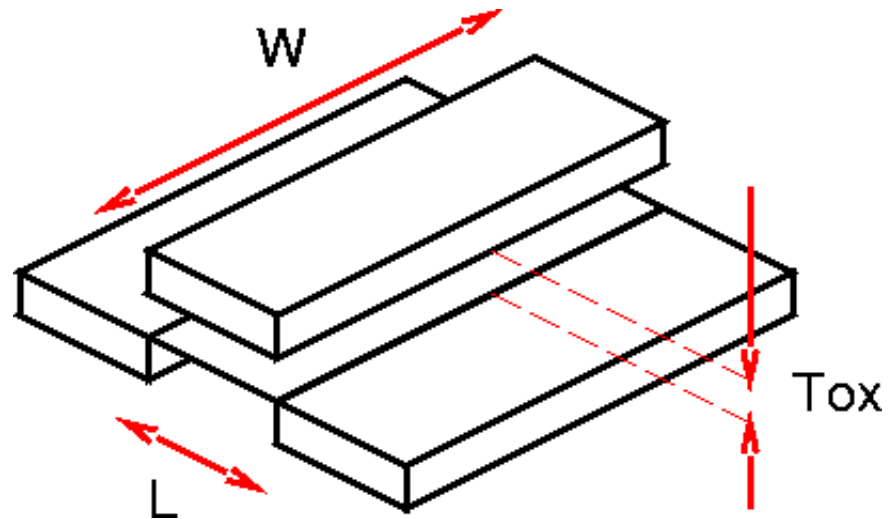
Side view



Perspective view

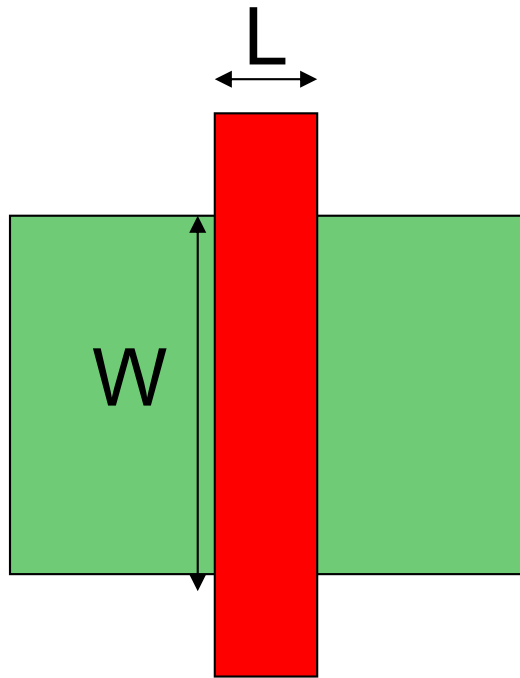
Layout

- ❑ Sizing & positioning of transistors
- ❑ Designer controls W , L
- ❑ t_{ox} fixed for process
 - Sometimes thick/thin oxide “flavors”

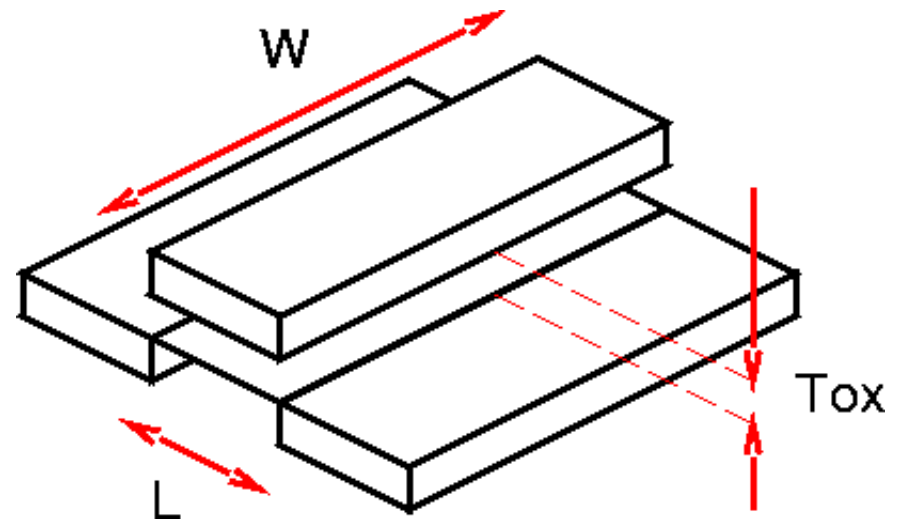




NMOS Geometry



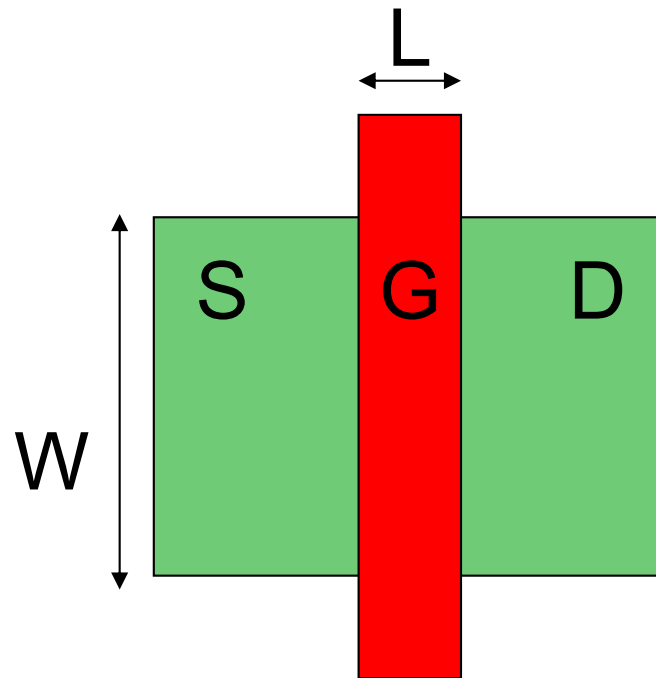
Top view



Perspective view

NMOS Geometry

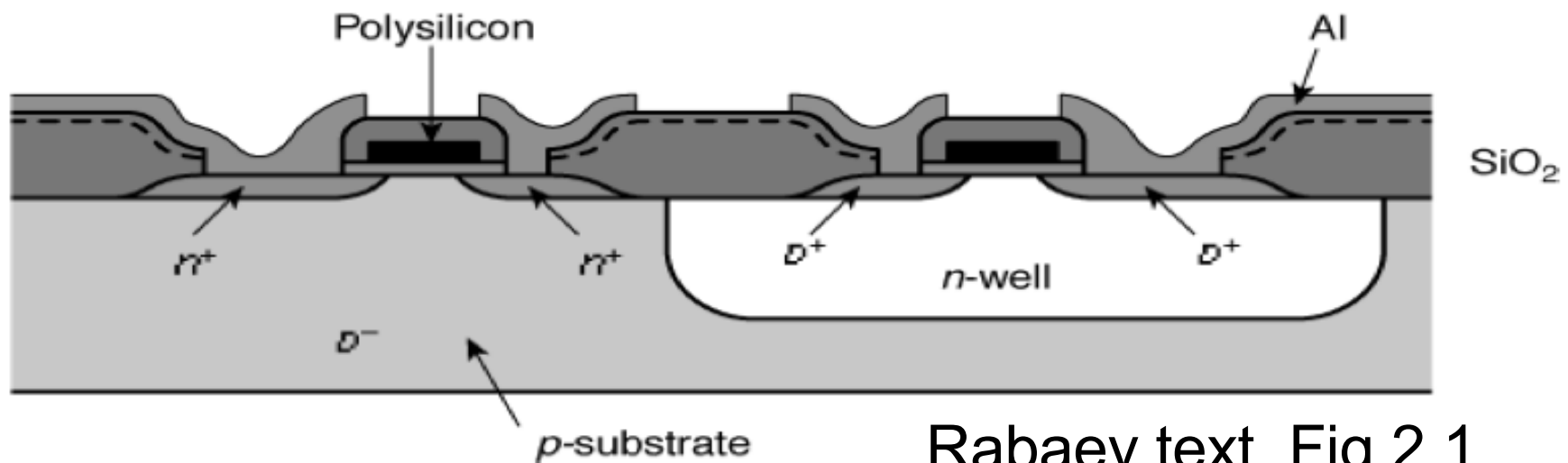
- Color scheme
 - Red: gate (polysilicon material)
 - Green: source and drain areas (n type diffusion)



Top view

NMOS vs PMOS

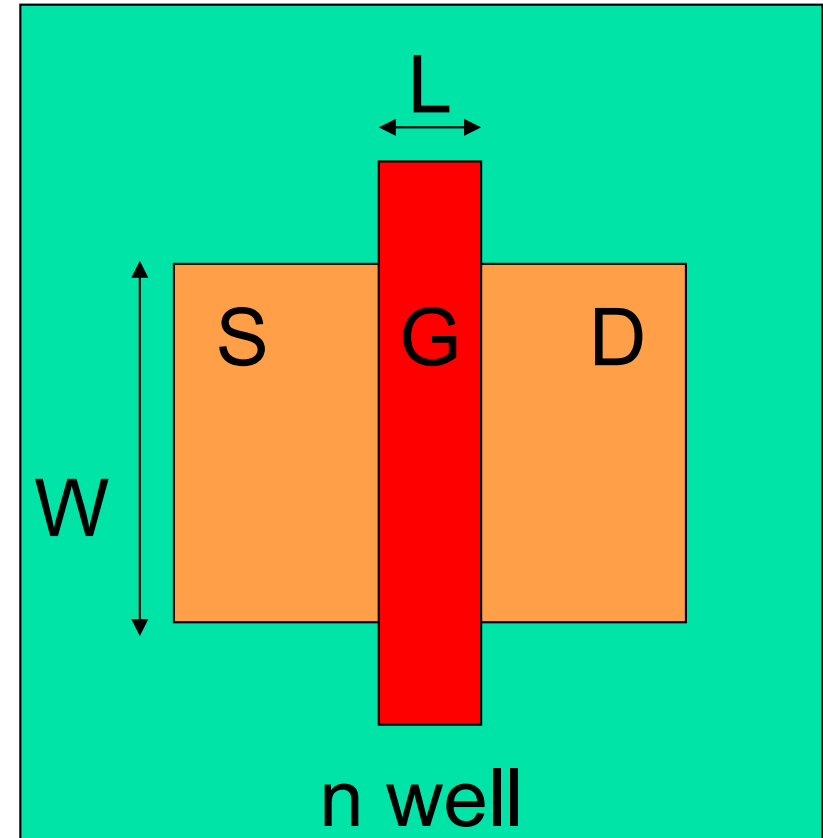
- ❑ NMOS built on p substrate
- ❑ PMOS built on n substrate
 - Needs an N-well



Rabaey text, Fig 2.1

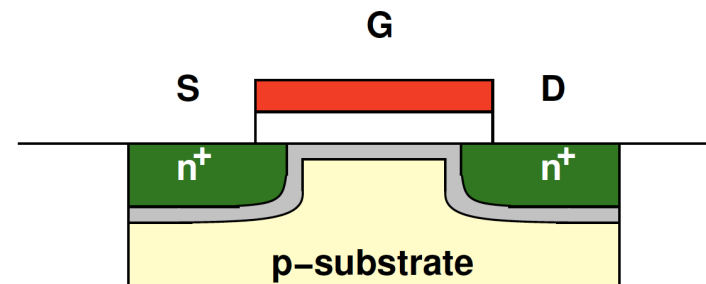
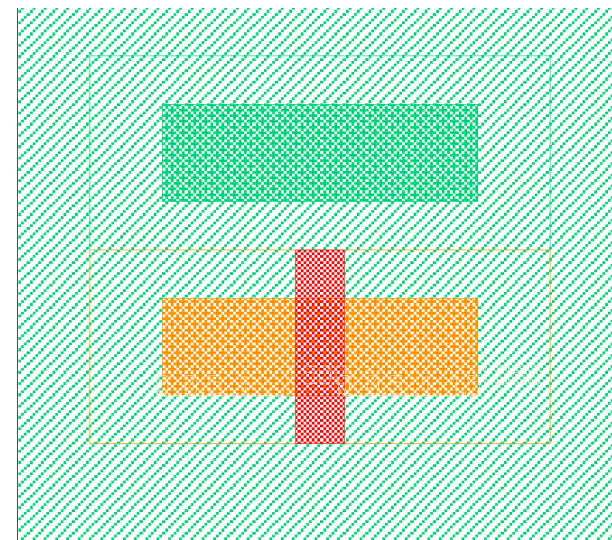
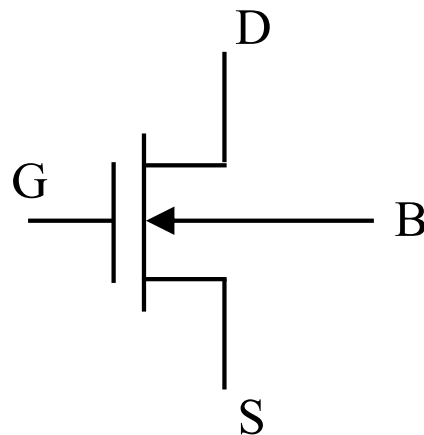
PMOS Geometry

- Color scheme
 - Red: gate
 - Orange: source and drain areas (p type)
 - Green: n well
- NMOS built on p wafer
 - Must add n well material to build PMOS



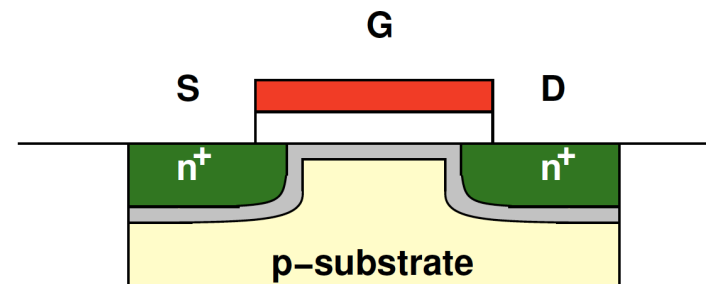
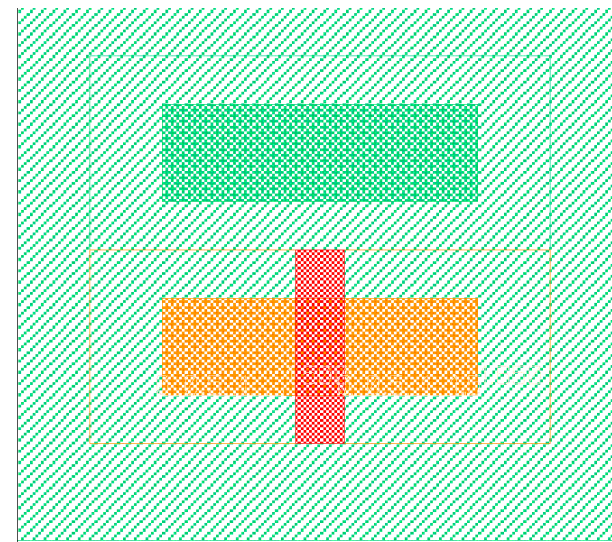
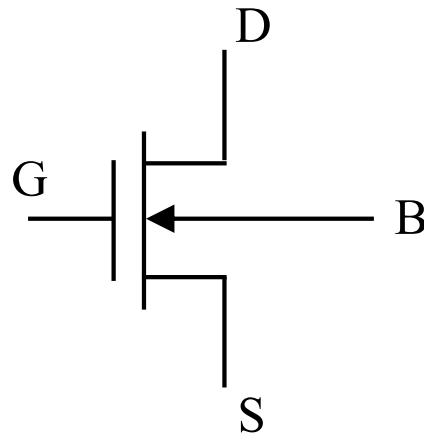
Body Contact

- ❑ “Fourth terminal”
- ❑ Needed to set voltage around device
 - PMOS: $V_b = V_{dd}$
 - NMOS: $V_b = GND$
- ❑ At right: PMOS (orange) with bulk contact (dark green)



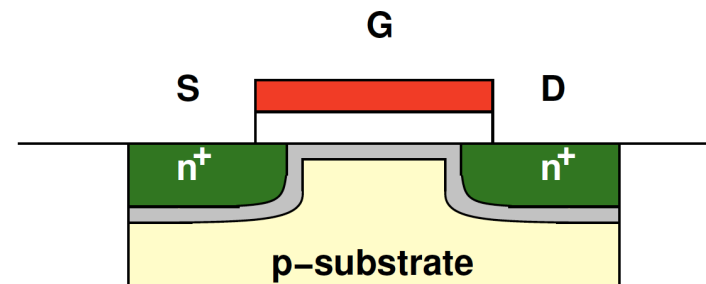
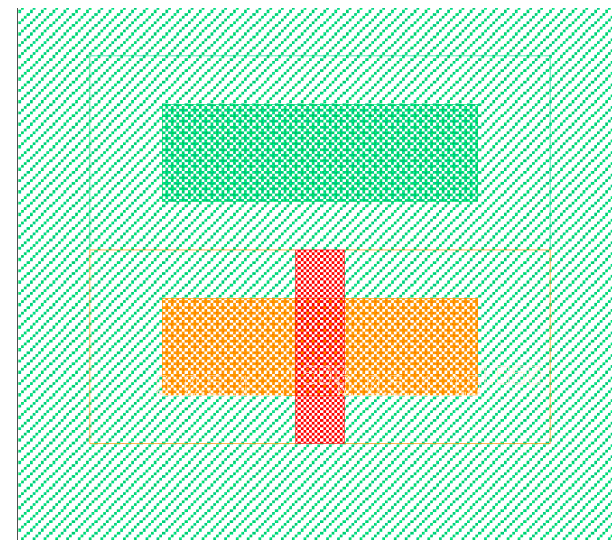
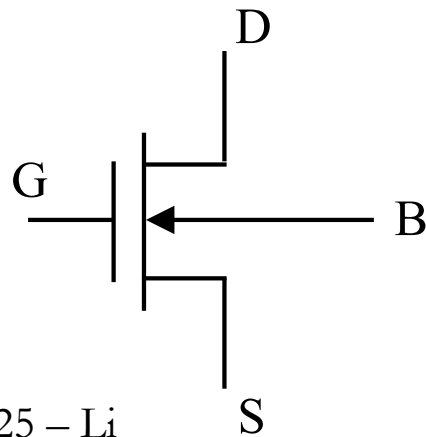
Body Contact

- ❑ Needed to set voltage around device
 - PMOS: $V_b = V_{dd}$
 - NMOS: $V_b = \text{GND}$
- ❑ What happens if NMOS body contact is V_{dd} ?



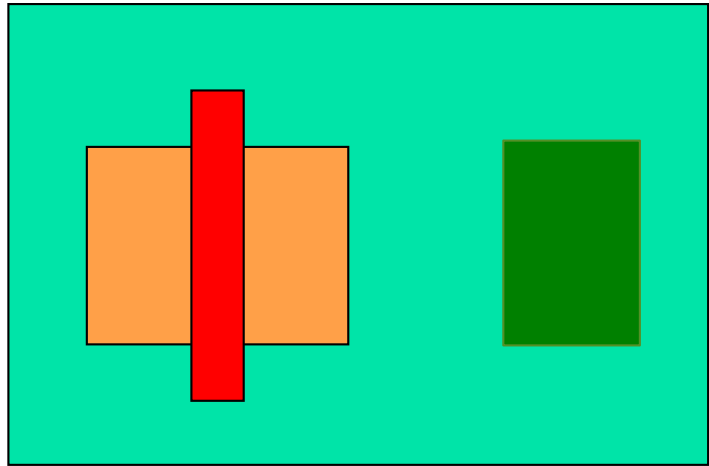
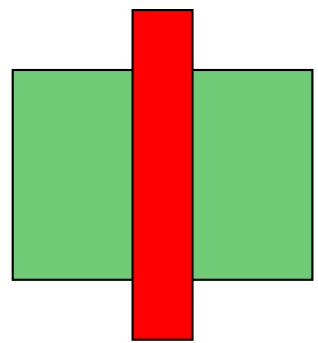
Body Contact

- ❑ Needed to set voltage around device
 - PMOS: $V_b = V_{dd}$
 - NMOS: $V_b = \text{GND}$
- ❑ What happens if NMOS body contact is V_{dd} ?
 - Polarity of field wrong
 - Increase V_{th} (need higher voltage to invert the channel)





Transistor Geometry

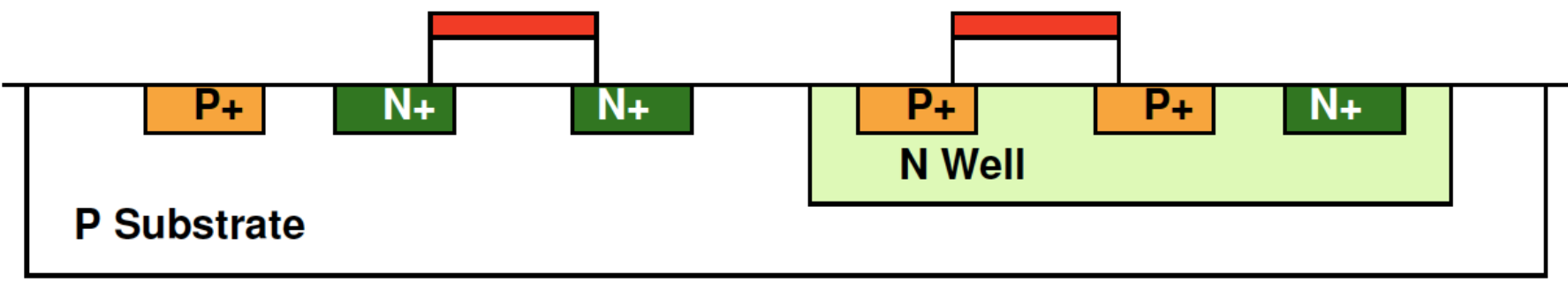


NMOS

PMOS

B S G D

D G S B

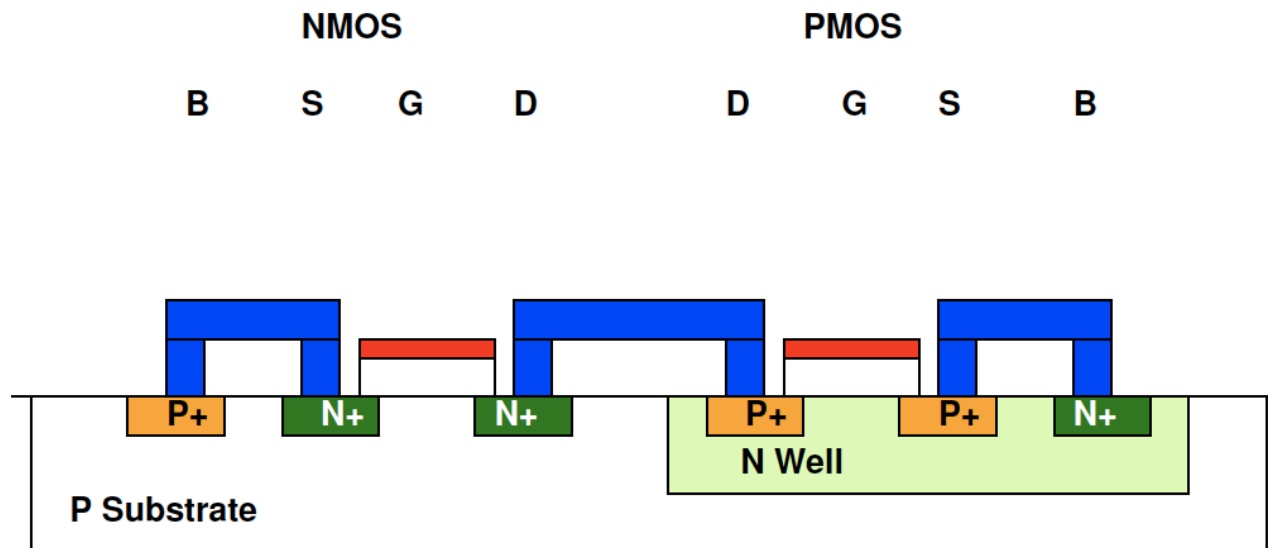


P Substrate

N Well

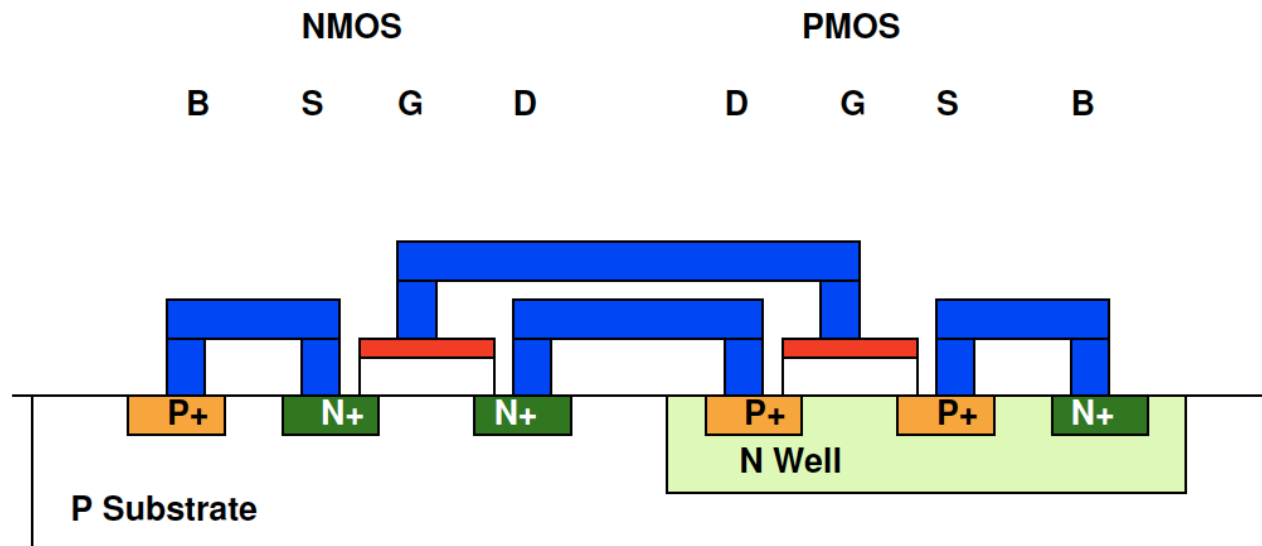
Interconnect

- Connect transistors
 - Different layers of metal
 - “Contact” - metal to transistor
 - “Via” - metal to metal

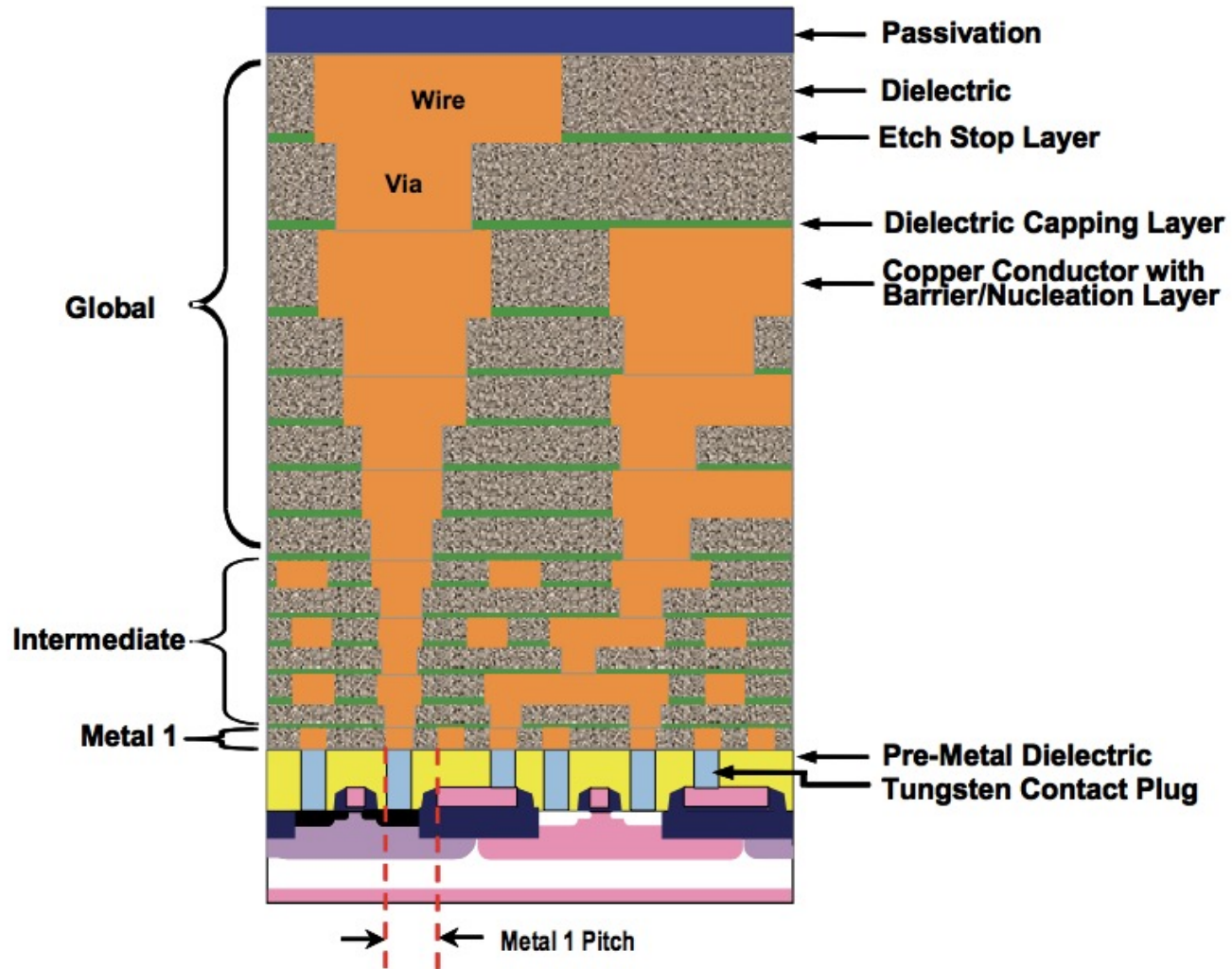


Interconnect

- Connect transistors
 - Different layers of metal
 - “Contact” - metal to transistor
 - “Via” - metal to metal

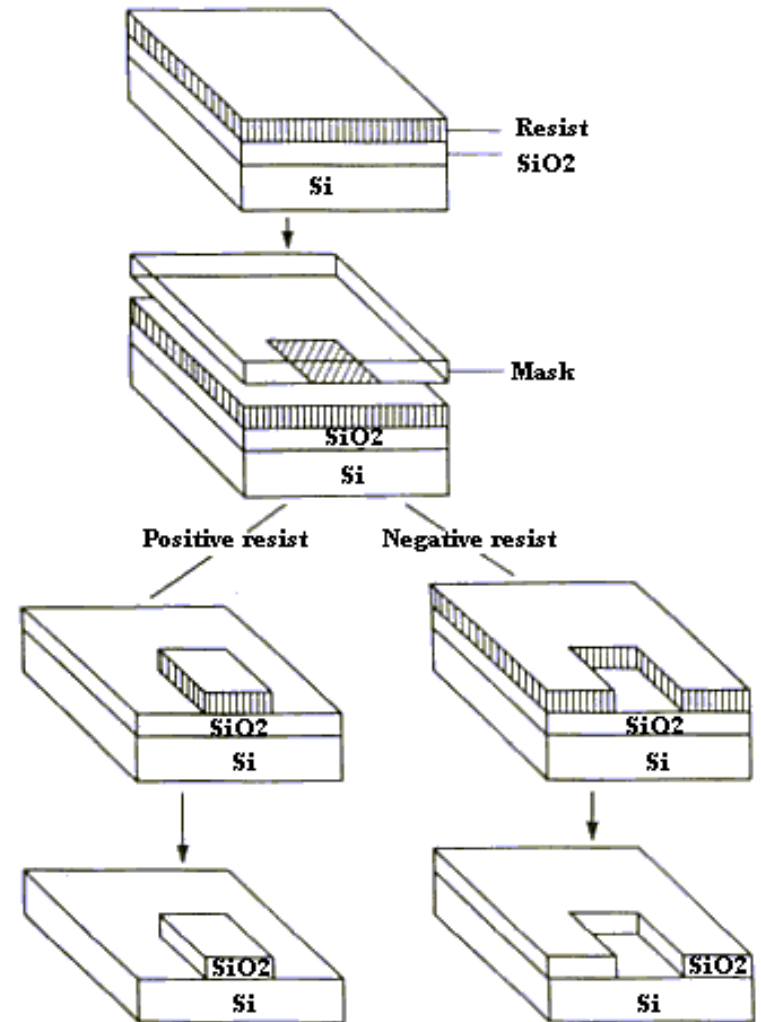


Interconnect Cross Section

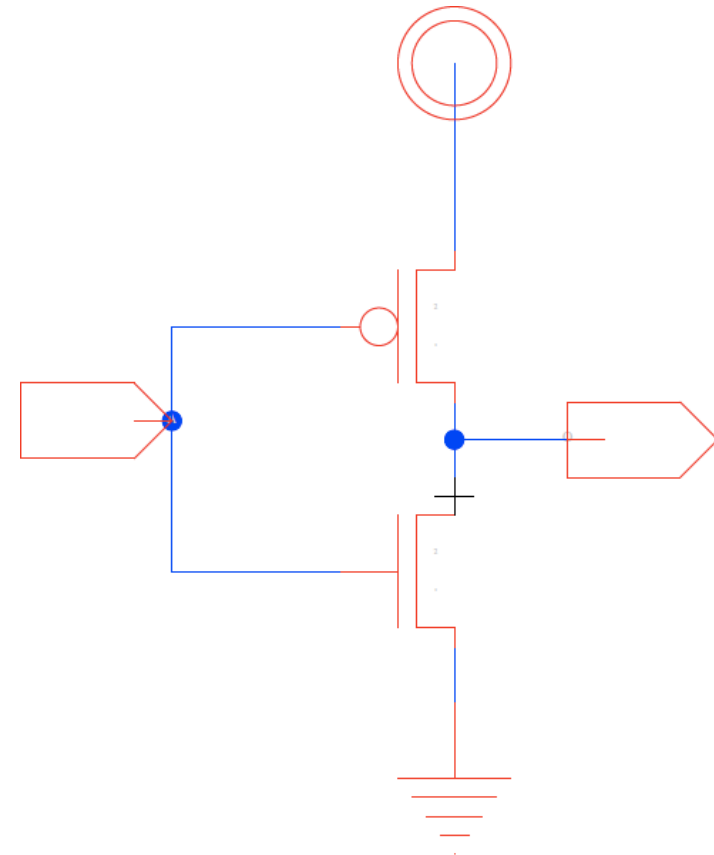
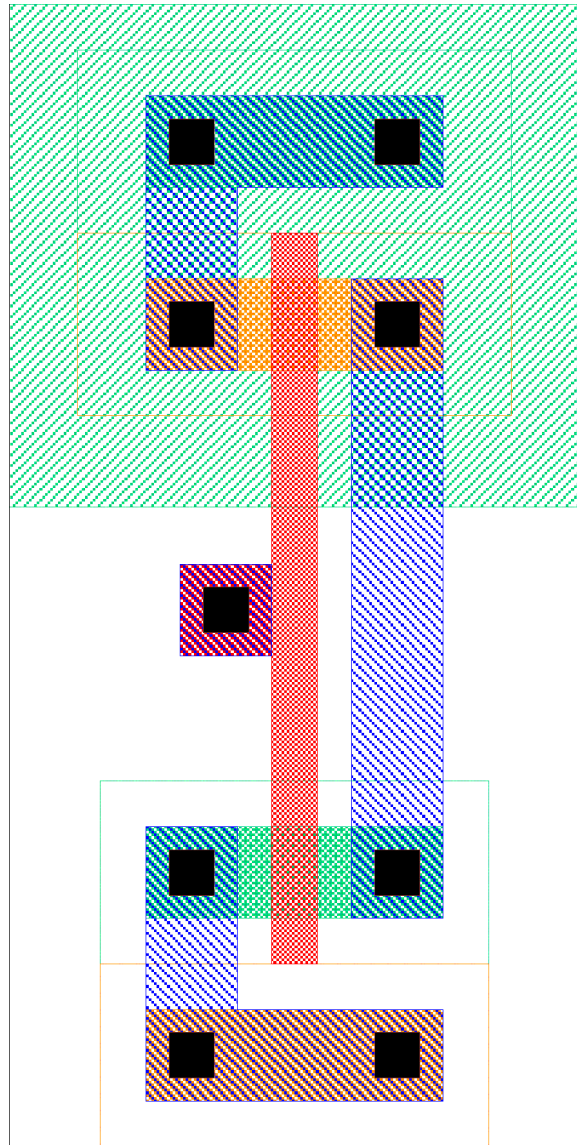


Masks

- ❑ Define areas want to see in layer
 - Think of “stencil” for material deposition
- ❑ Use photoresist (PR) to form the “stencil”
 - Grow PR over entire wafer
 - Expose PR through mask
 - PR dissolves in exposed areas
 - Material is deposited/etched
 - Only “sticks” in area w/ dissolved PR



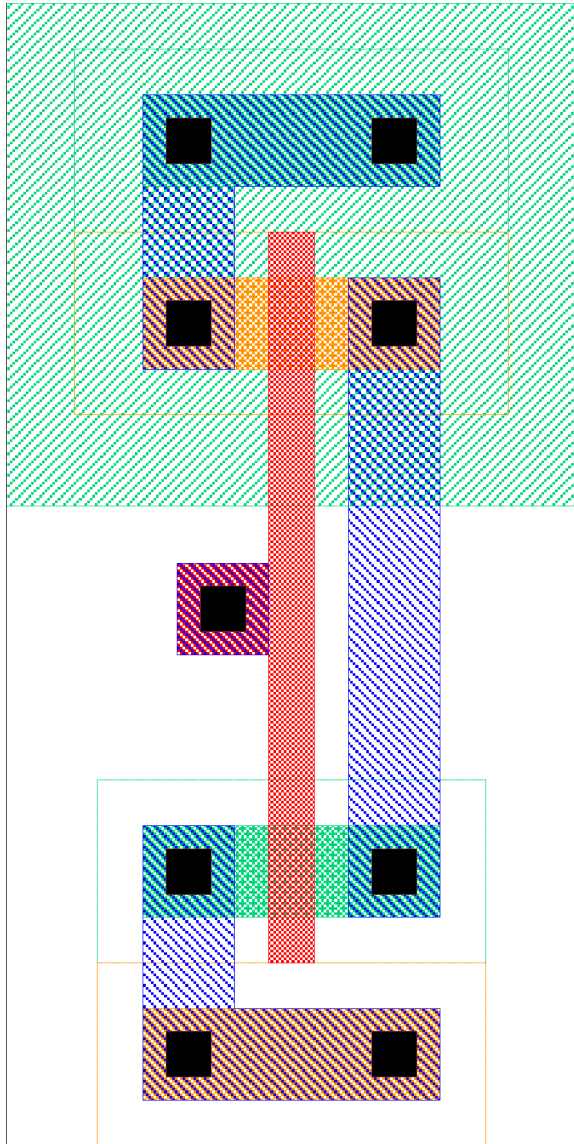
Reverse Engineer Inverter Layout (Preclass 1)



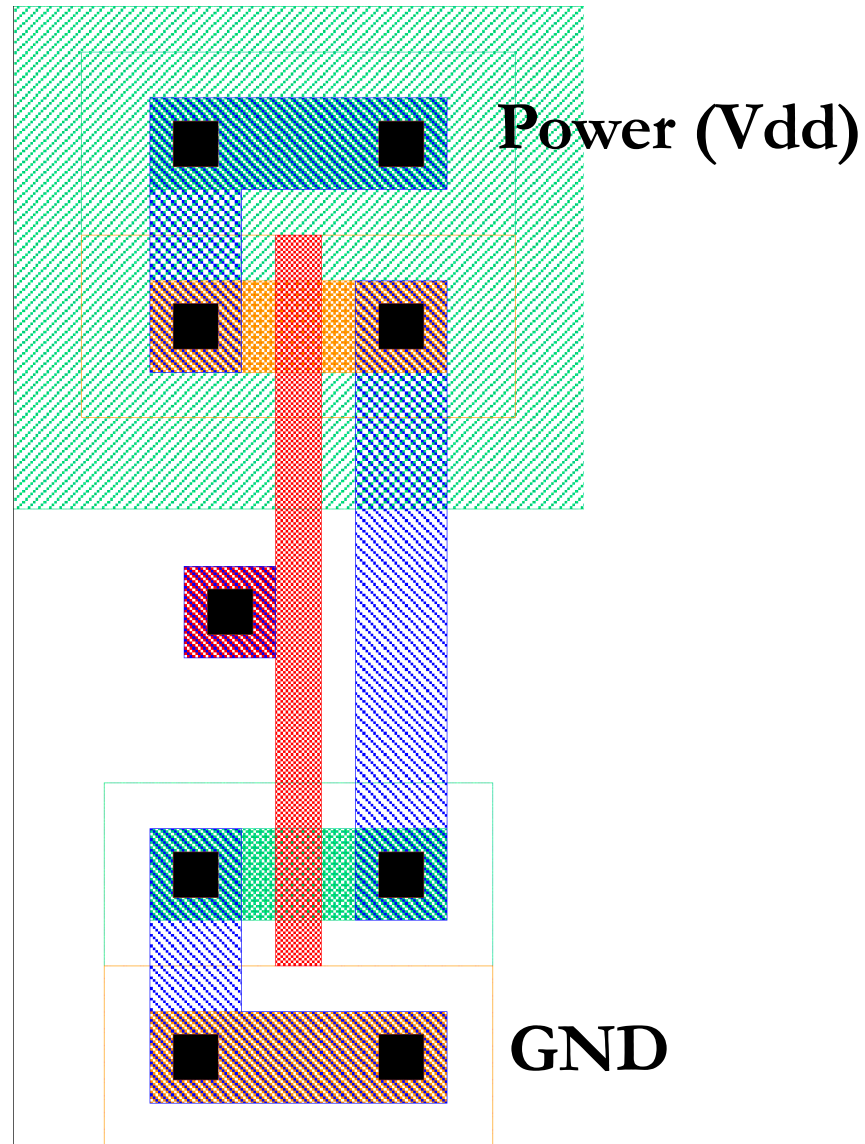


Layout Revisited

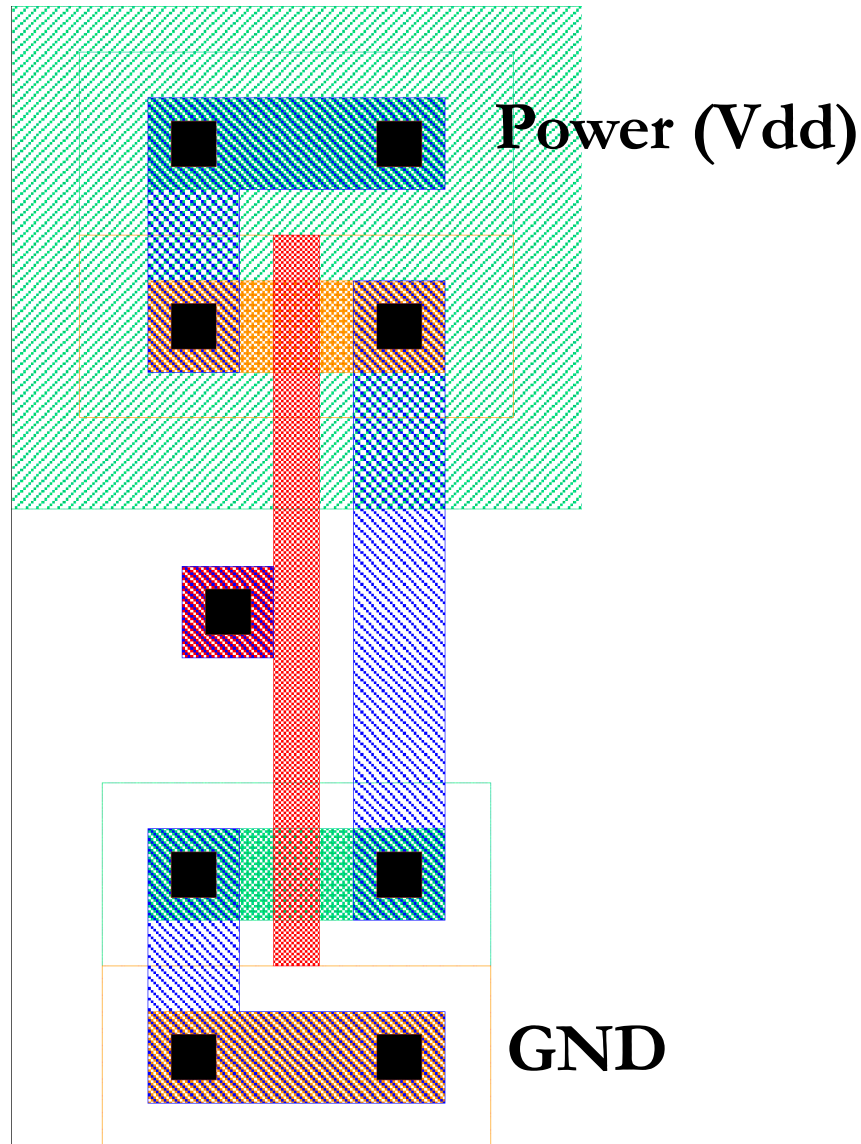
- How to “decode” circuit from layout?



Reverse Engineer Inverter Layout



Reverse Engineer Inverter Layout

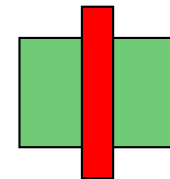
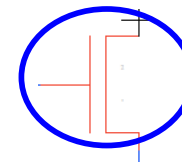
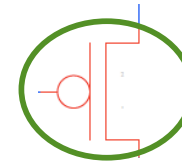
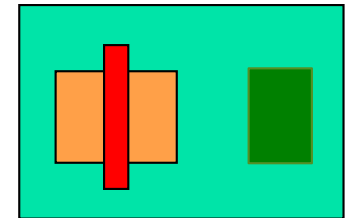
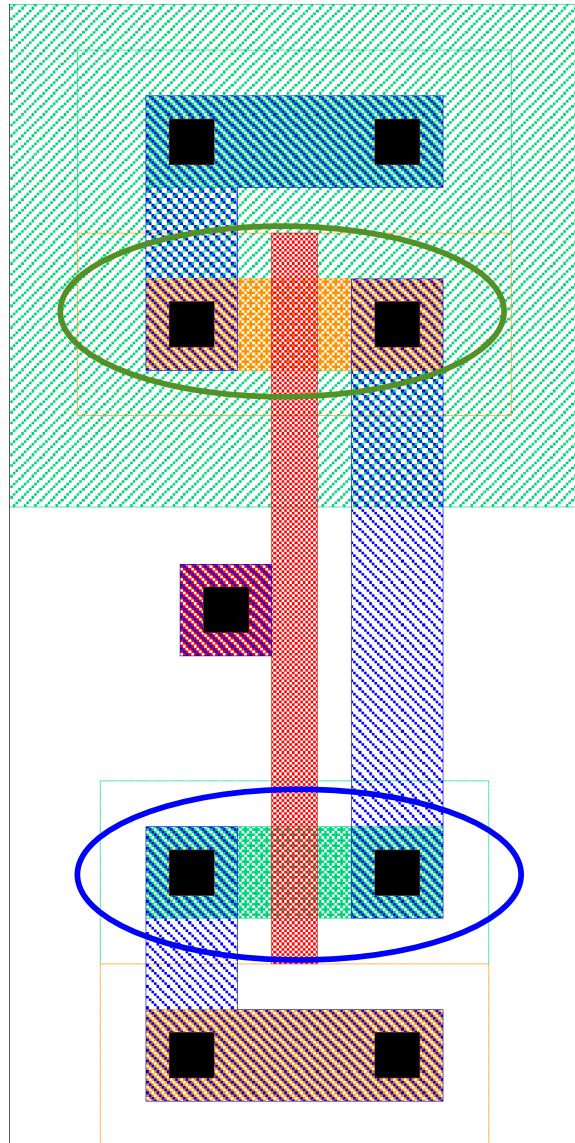


- Where is PMOS transistor?
- NMOS?



Layout to Circuit

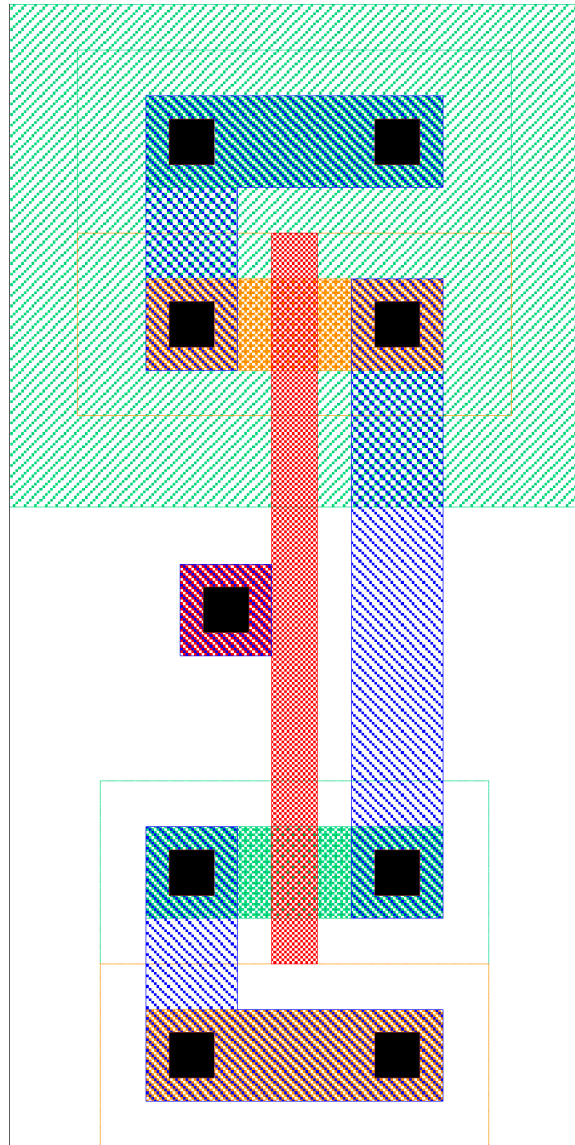
- 1. Identify transistors





Inverter Layout

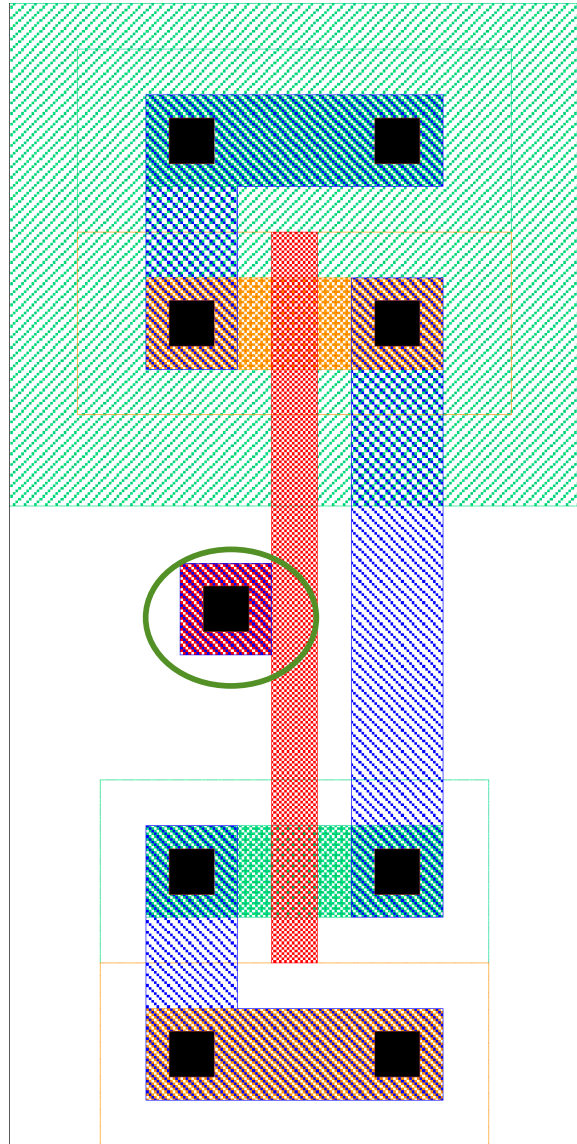
□ Where is Input?





Inverter Layout

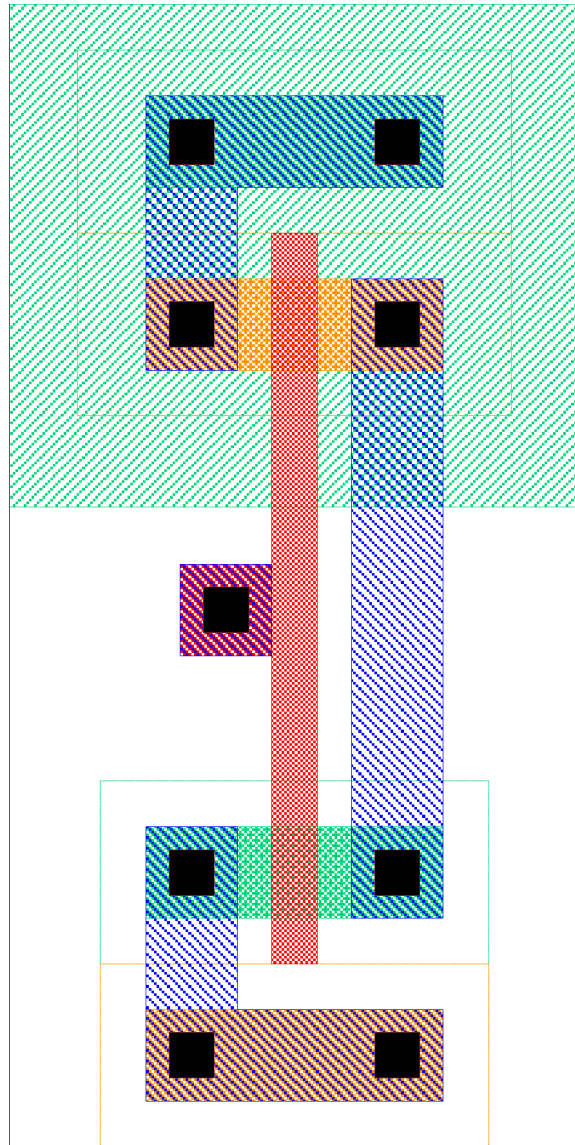
□ Where is Input?





Inverter Layout

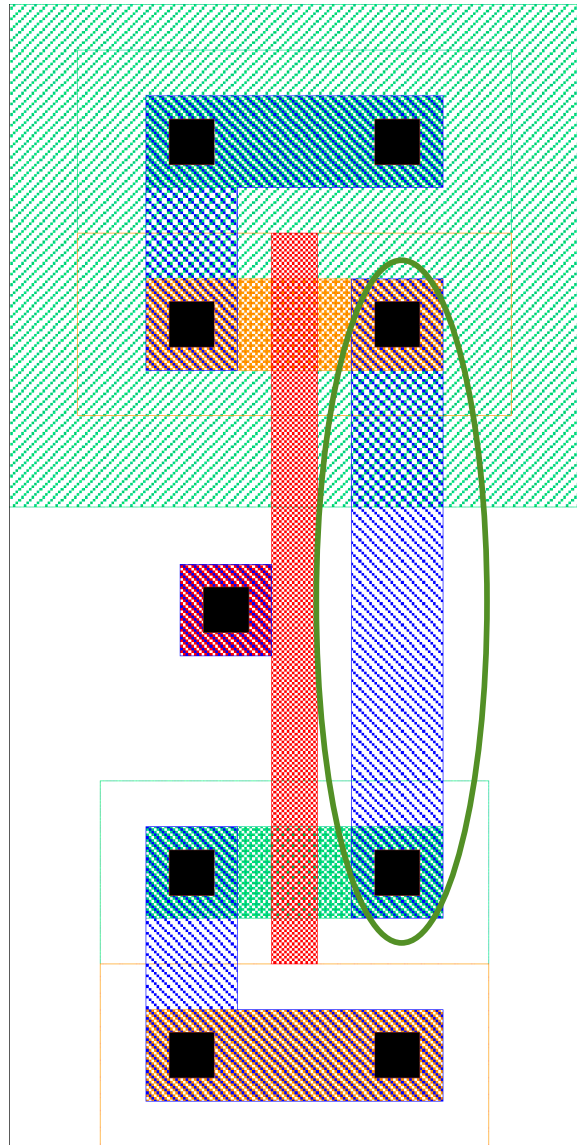
□ Where is Output?





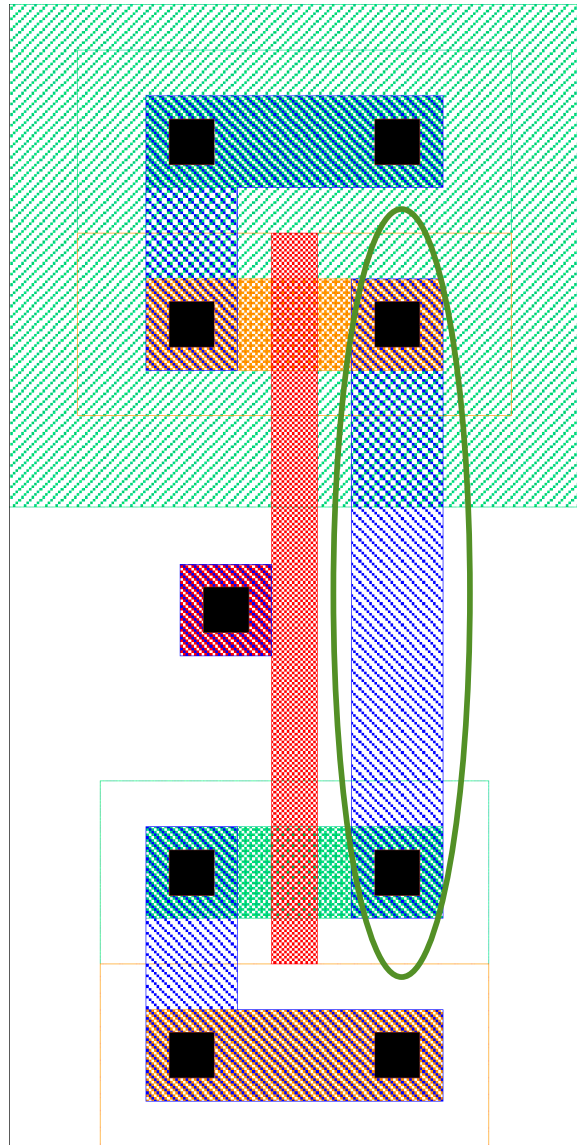
Inverter Layout

□ Where is Output?

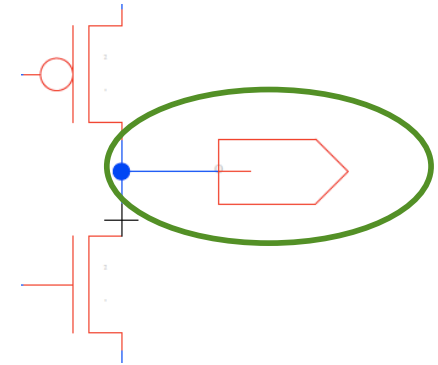




Layout to Circuit

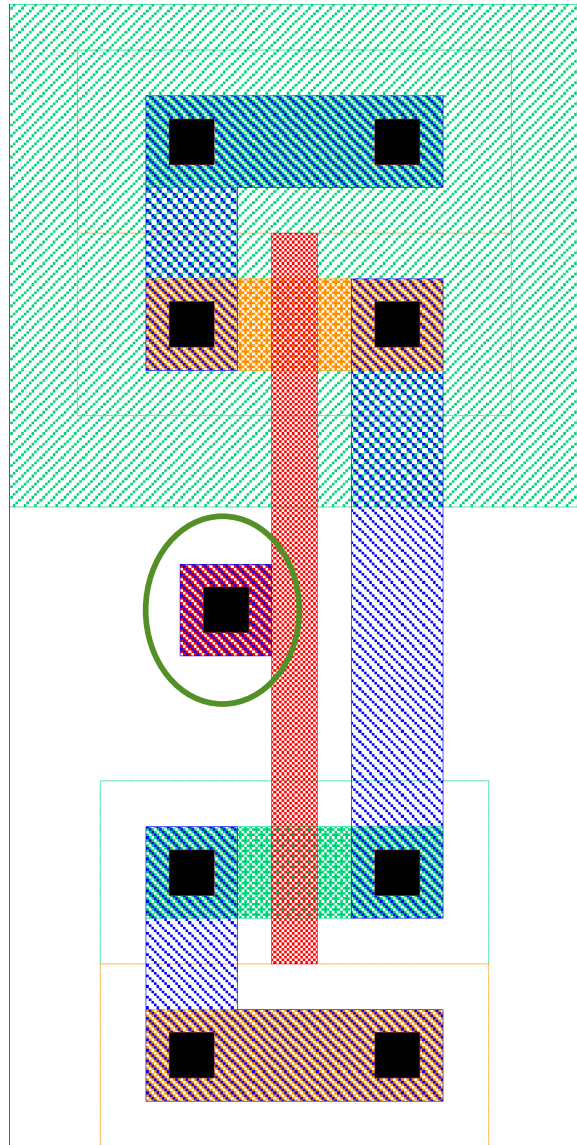


- 2. Add connections
 - Drain connection



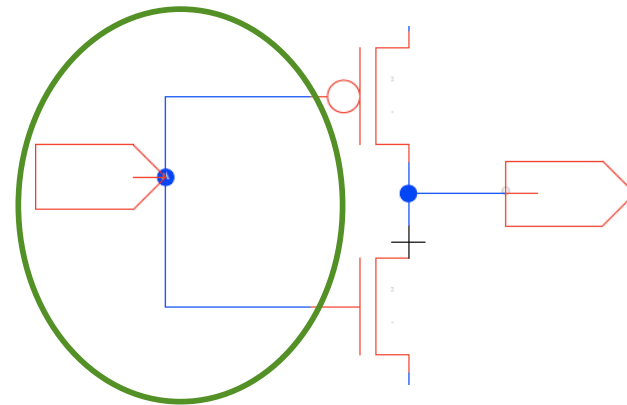
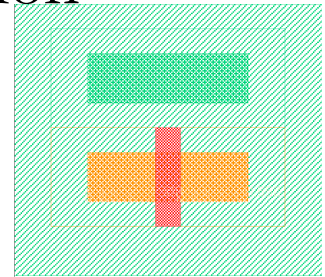


Layout to Circuit



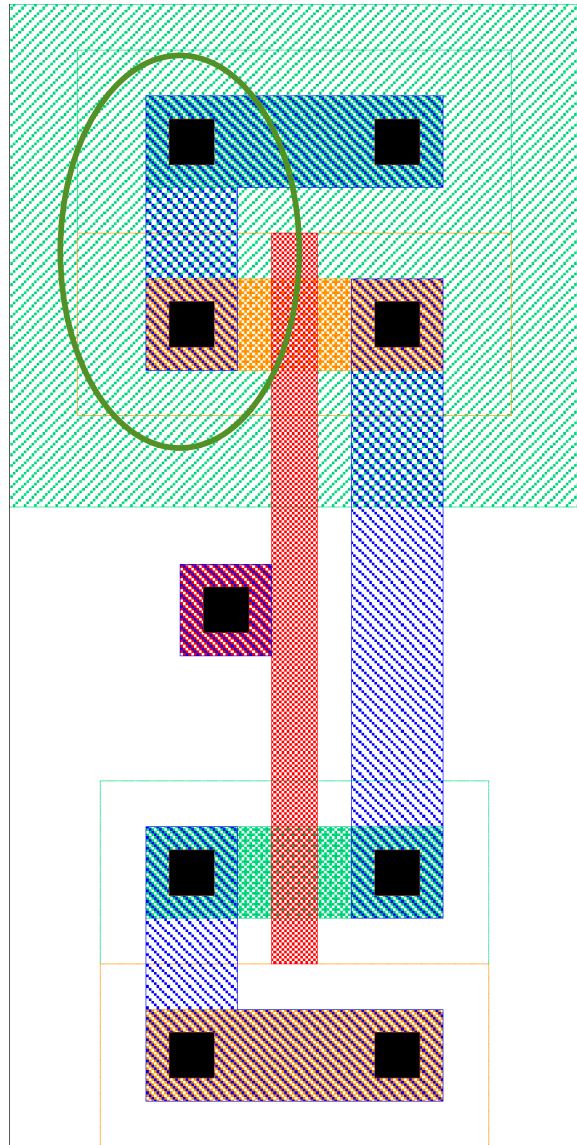
□ 2. Add connections

■ Gate connection

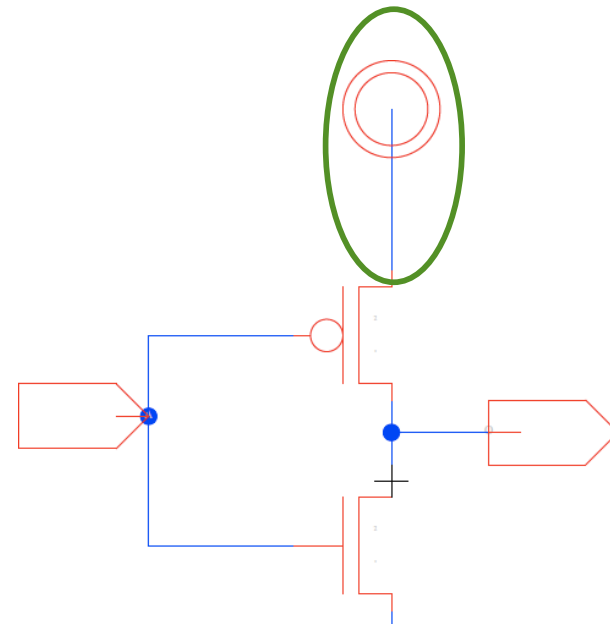




Layout to Circuit

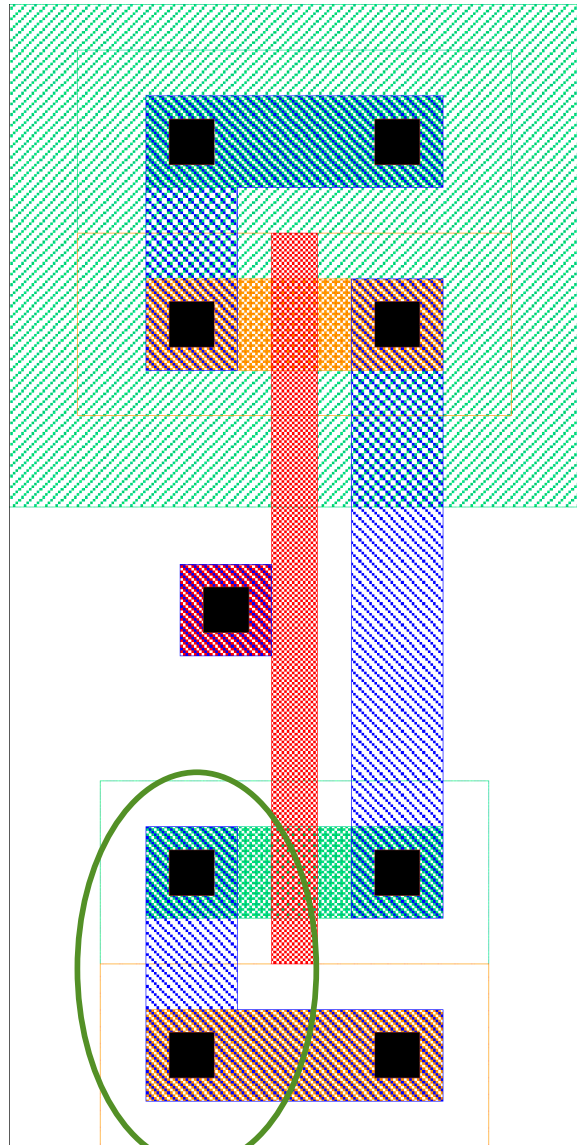


- 2. Add connections
 - pMOS-source to VDD

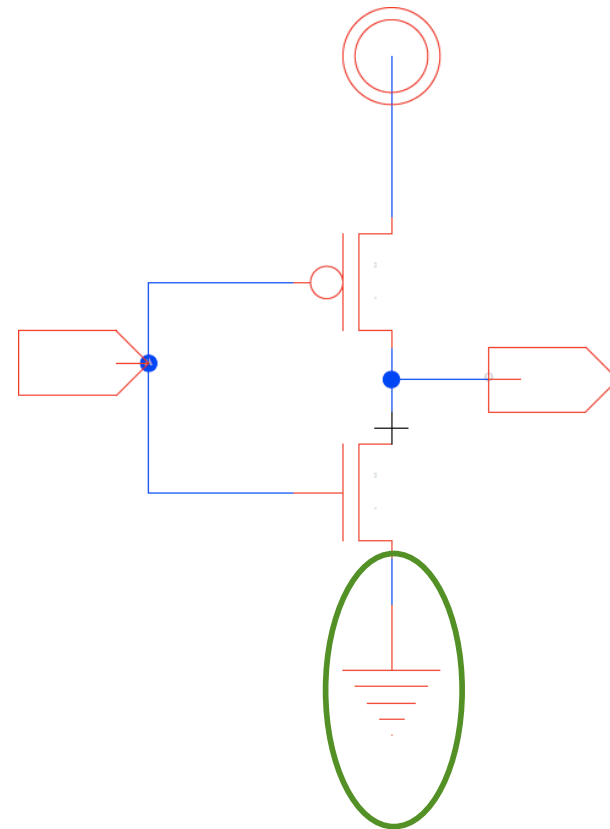




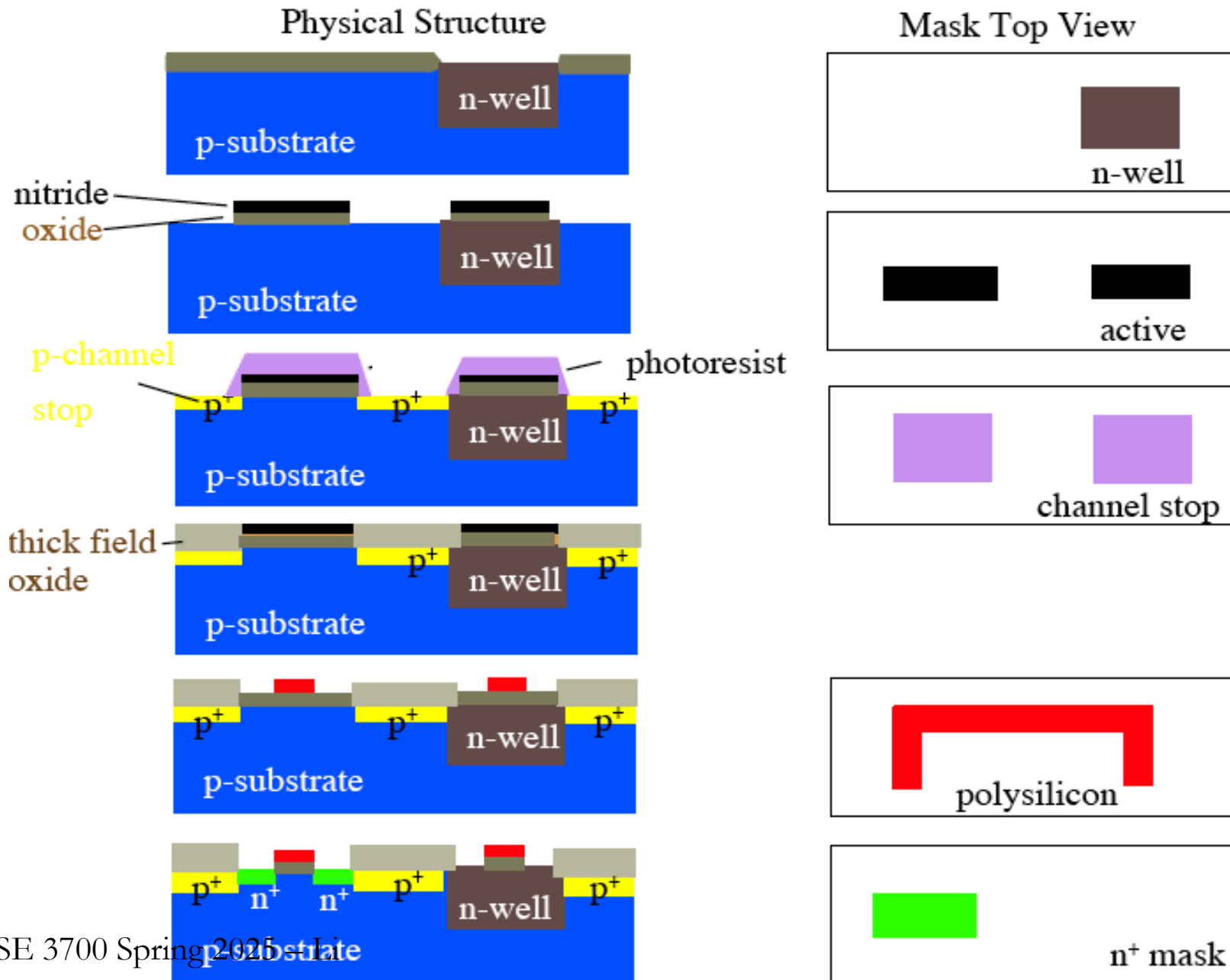
Layout to Circuit



- 2. Add connections
 - nMOS source to GND



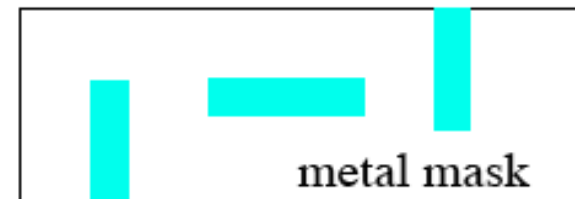
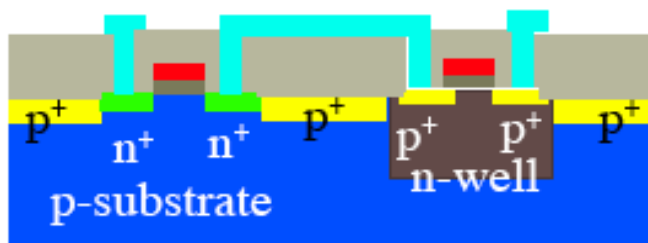
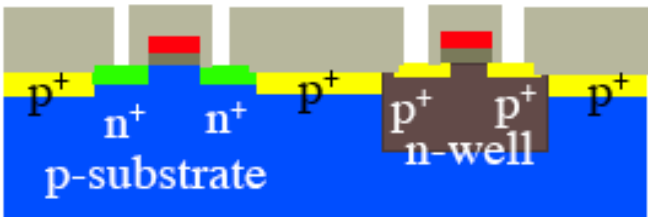
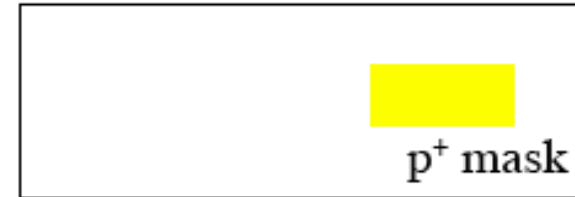
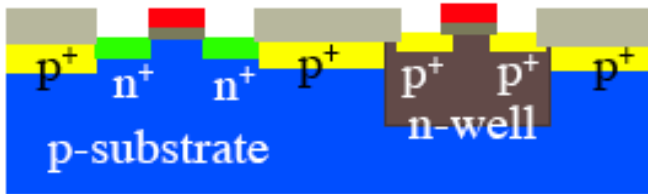
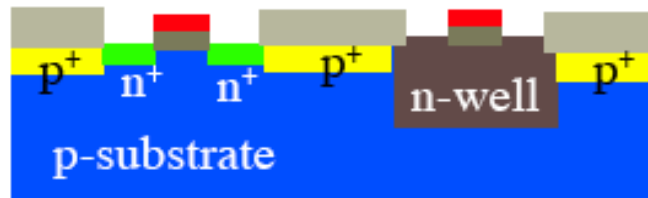
Typical N-Well CMOS Process



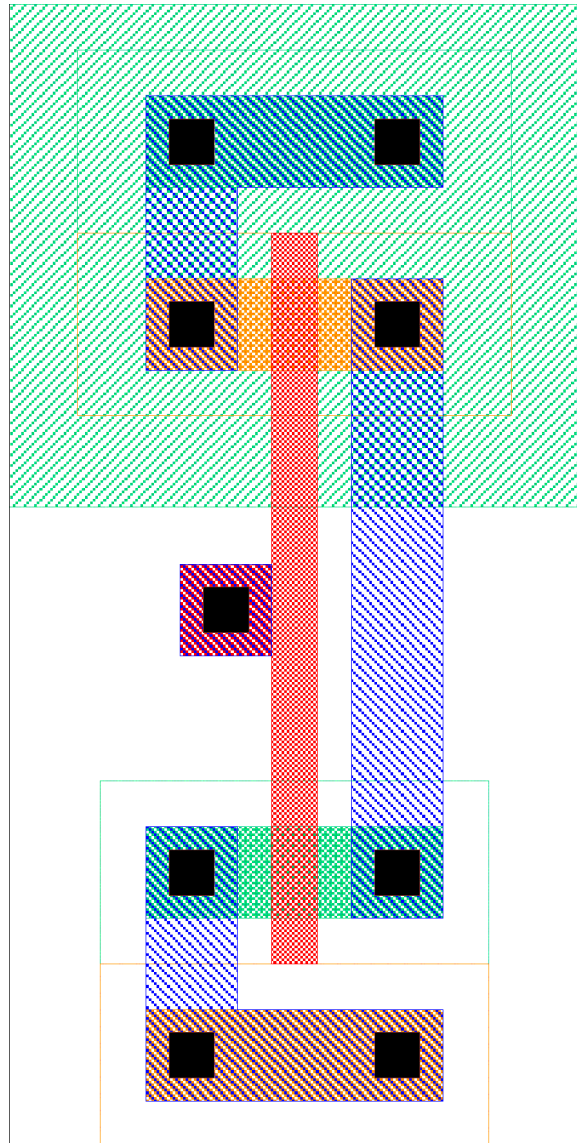
Typical N-Well CMOS Process

Physical Structure

Mask Top View



Design Rules



- Why not adjacent transistors?
 - Plenty of empty space
 - If area is money, pack in as much as possible
 - Shortens connections
- Recall: processing is imprecise
 - Margin of error for process variation

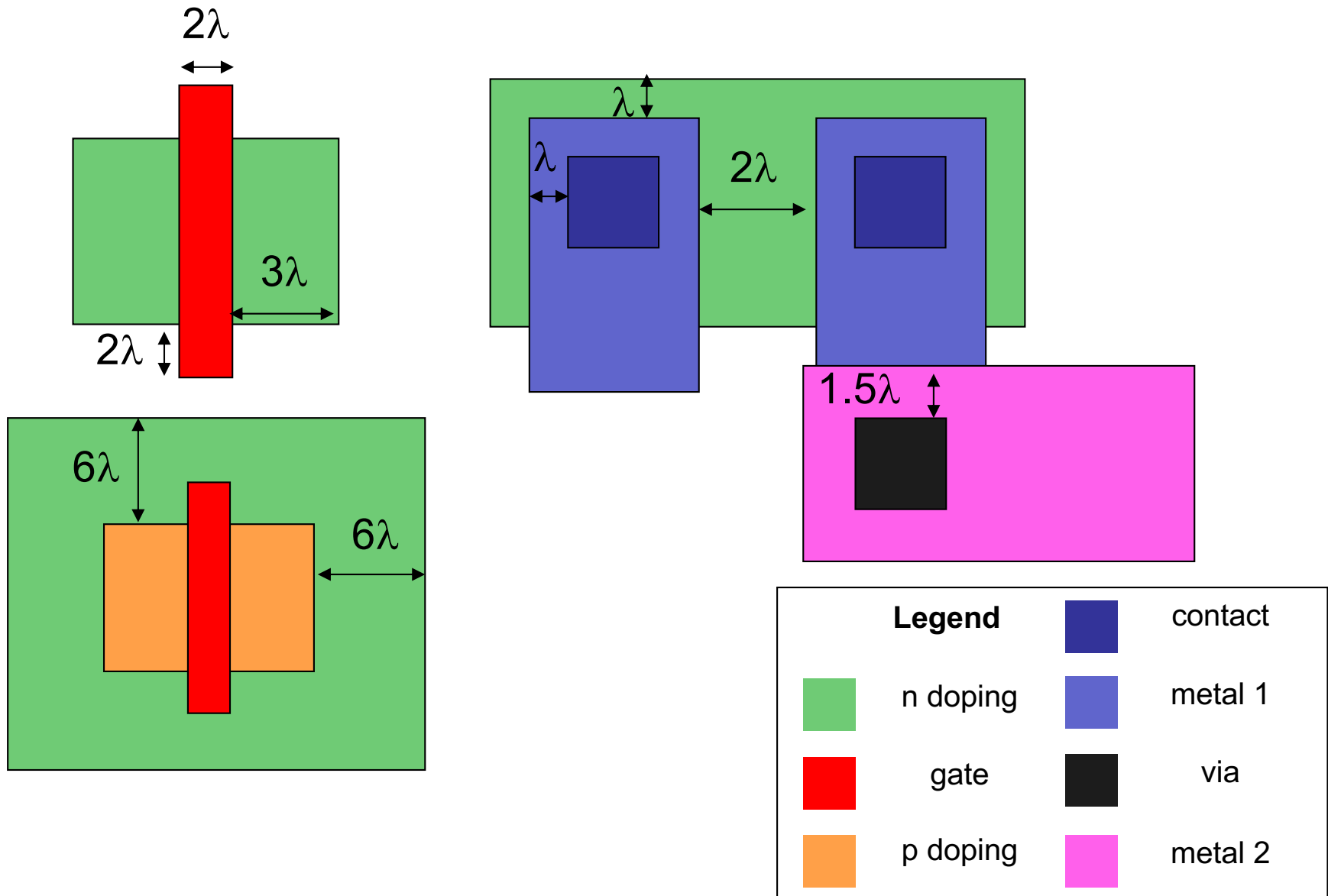


Design Rules

- ❑ Contract between process engineer & designer
 - Minimum width/spacing
 - Can be (often are) process specific

- ❑ Lambda rules: scalable design rules
 - In terms of $\lambda = 0.5 L_{\min}$ (L_{drawn})
 - Can migrate designs from similar process with lambda factor

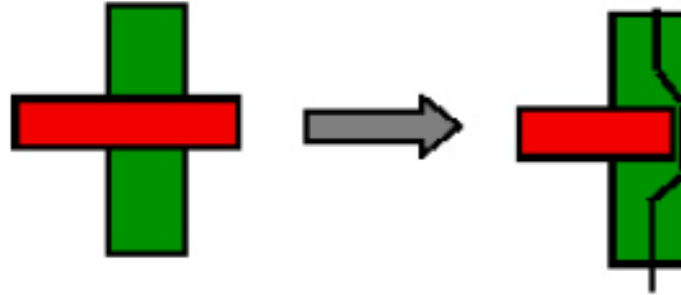
Design Rules: Some Examples



Potential Consequences of Design Rule Violations

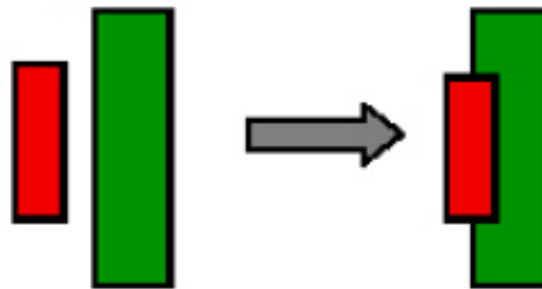
❑ Inter-Layer Design Rule Origins

Intended Transistor



Catastrophic Error – Unintended misalignment cause Source-Drain short circuit

Intended Unrelated Poly & Diffusion



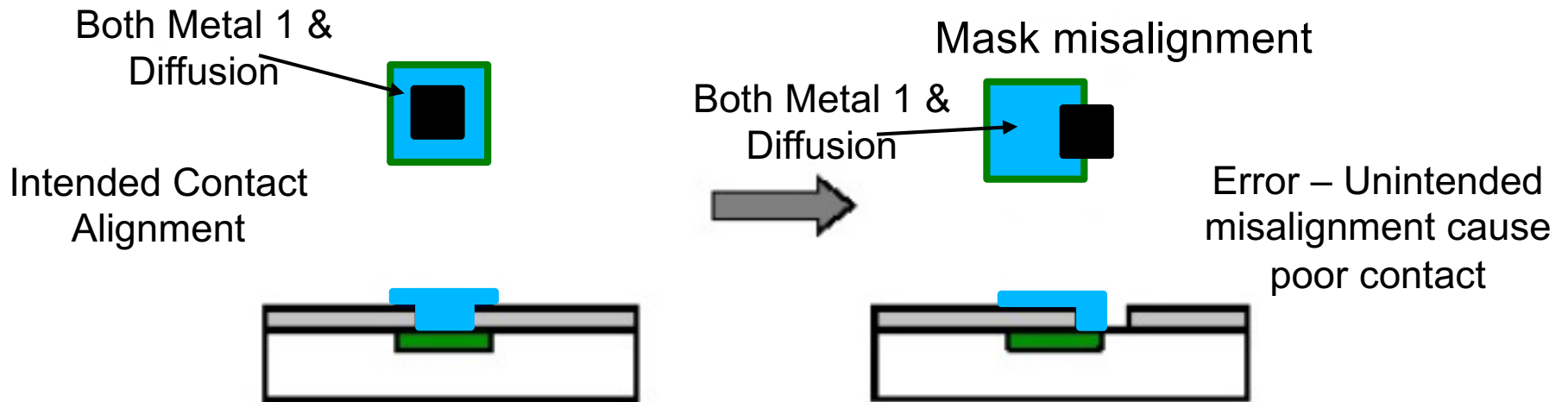
Catastrophic Error – Unintended overlap cause fabrication of a parasitic Transistor

Potential Consequences of Design Rule Violations

□ Inter-Layer Design Rule Origins

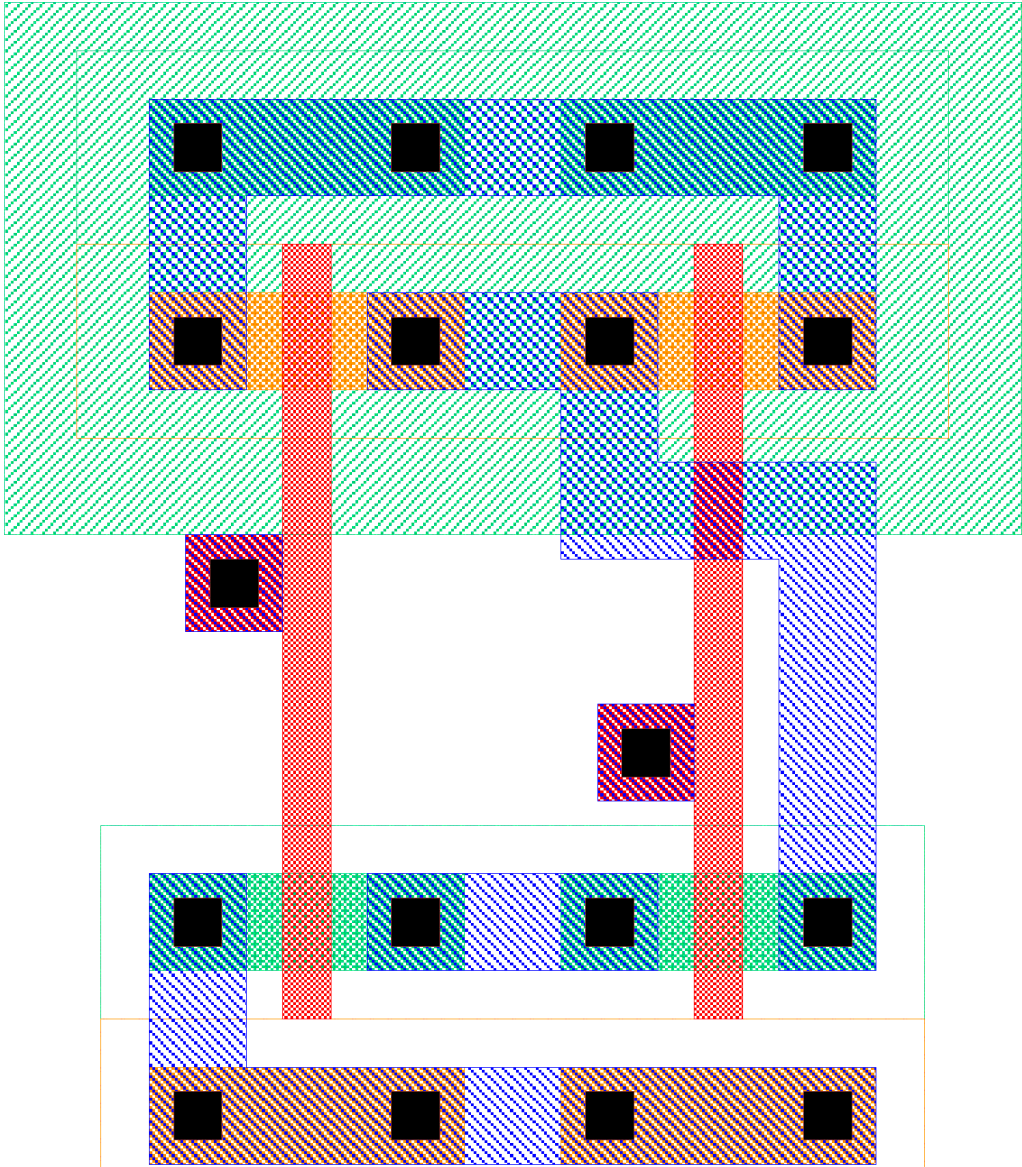
Contact and Via Masks

M1 contact to n-diffusion M1 contact to p-diffusion M1 contact to poly]	-> Contact Mask
Mn contact to Mn-1 for n = 2, 3,..		-> Via Mask



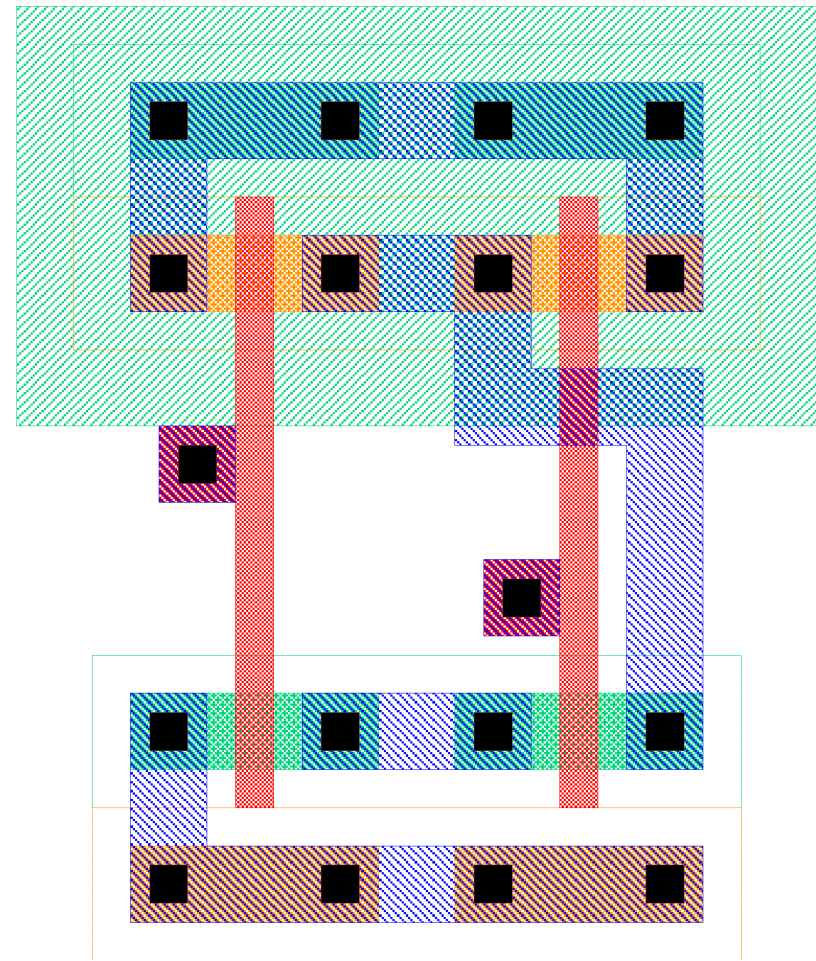


Layout #2 (preclass 2)



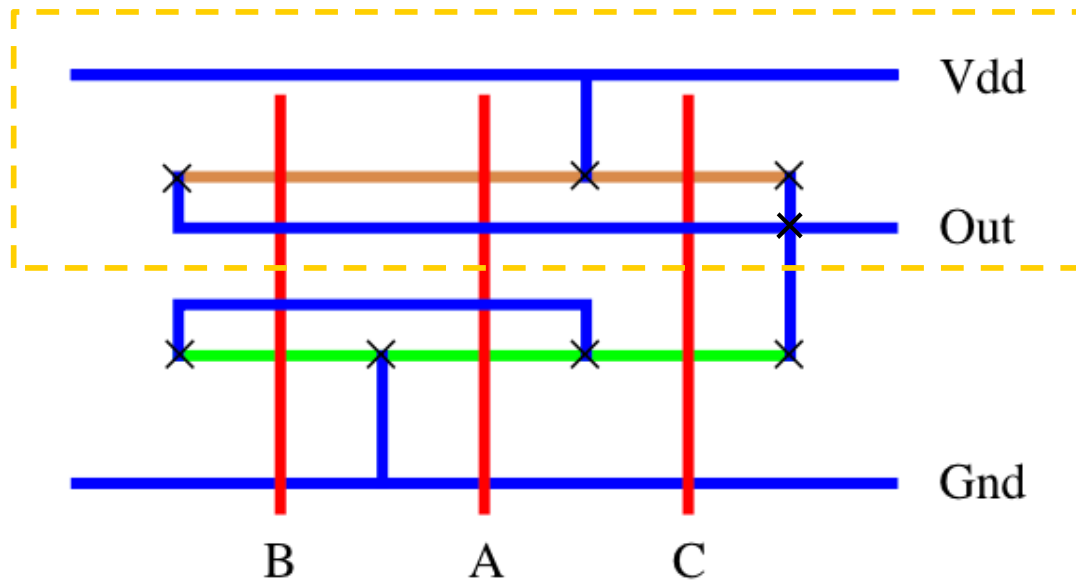
Layout #2 (preclass 2)

- How many transistors?
 - PMOS?
 - NMOS?
- How connected?
 - PMOS, NMOS?
- Inputs connected how?
- Outputs?
- What is it?



Symbolic Layout

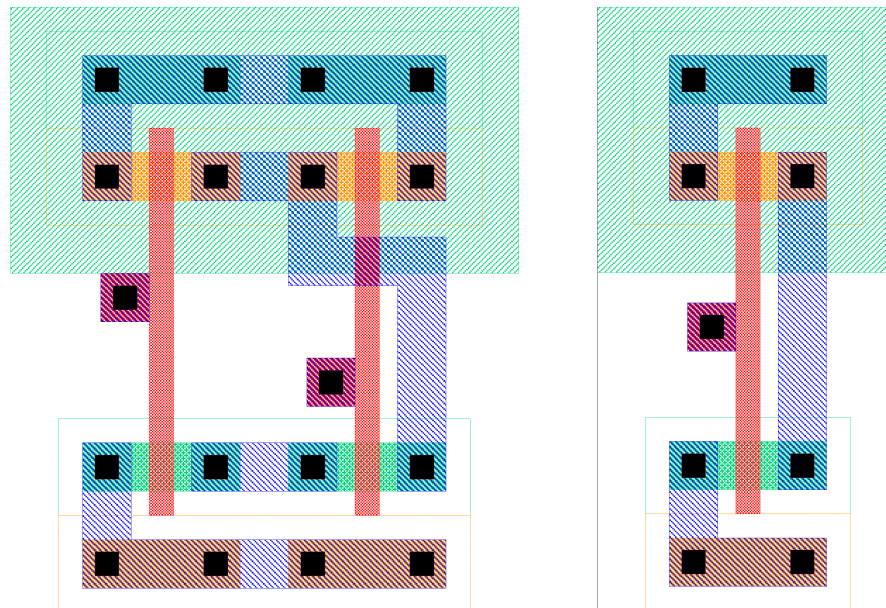
- Stick diagrams capture spatial relationships, but abstract away design rules



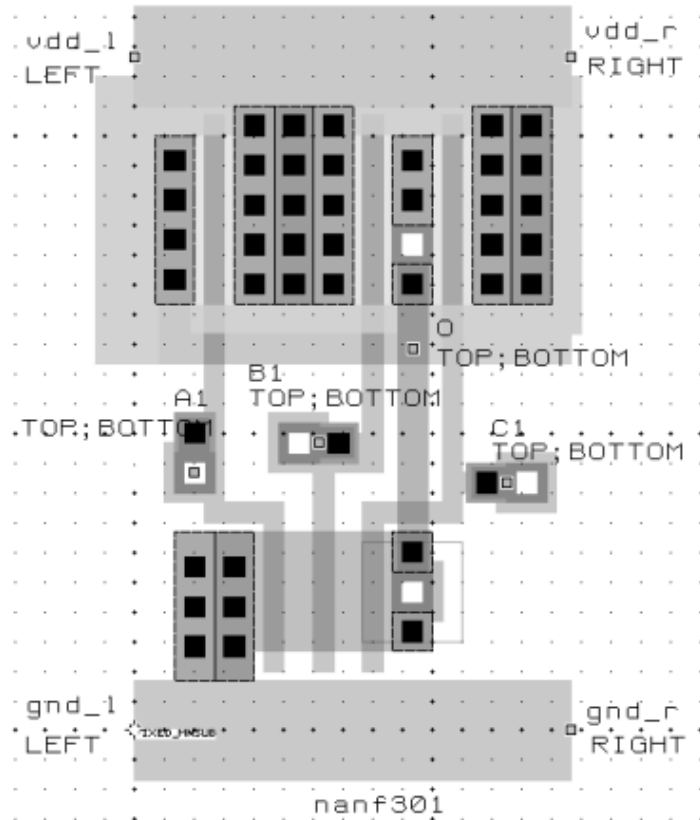
- What is the gate function?
 - How many NMOS? PMOS? D/S connections?
 - Draw schematic

Standard Cells

- Lay out gates so that heights match
 - Rows of adjacent cells
 - Standardized sizing of gate heights
- Motivation: automated place and route
 - EDA tools convert HDL to layout



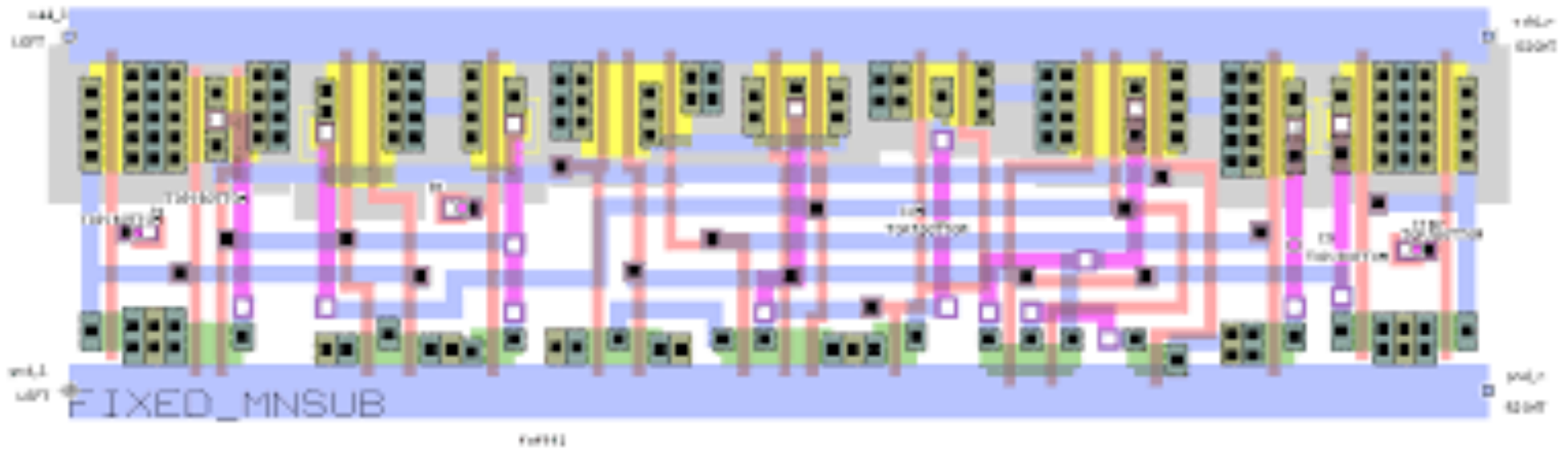
Standard Cells



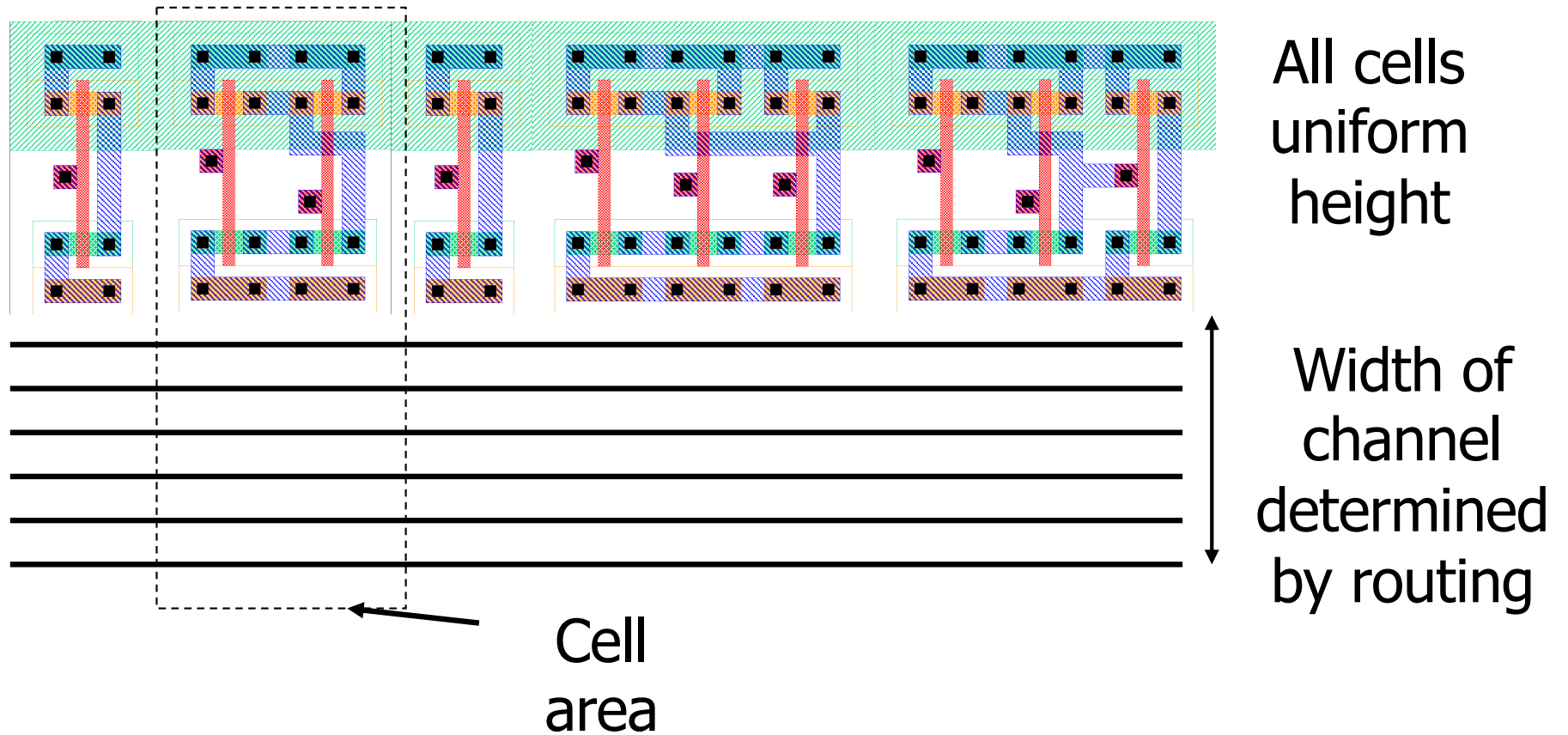
Fanout 4x	0.5 μm	1.0 μm	2.0 μm
<i>A1_tphl</i>	0.595	0.711	0.919
<i>A1_tplh</i>	0.692	0.933	1.360
<i>B1_tphl</i>	0.591	0.739	1.006
<i>B1_tplh</i>	0.620	0.825	1.1.81
<i>C1_tphl</i>	0.574	0.740	1.029
<i>C1_tplh</i>	0.554	0.728	1.026

3-input NAND cell
 (from Mississippi State Library)
 characterized for fanout of 4 and
 for three different technologies

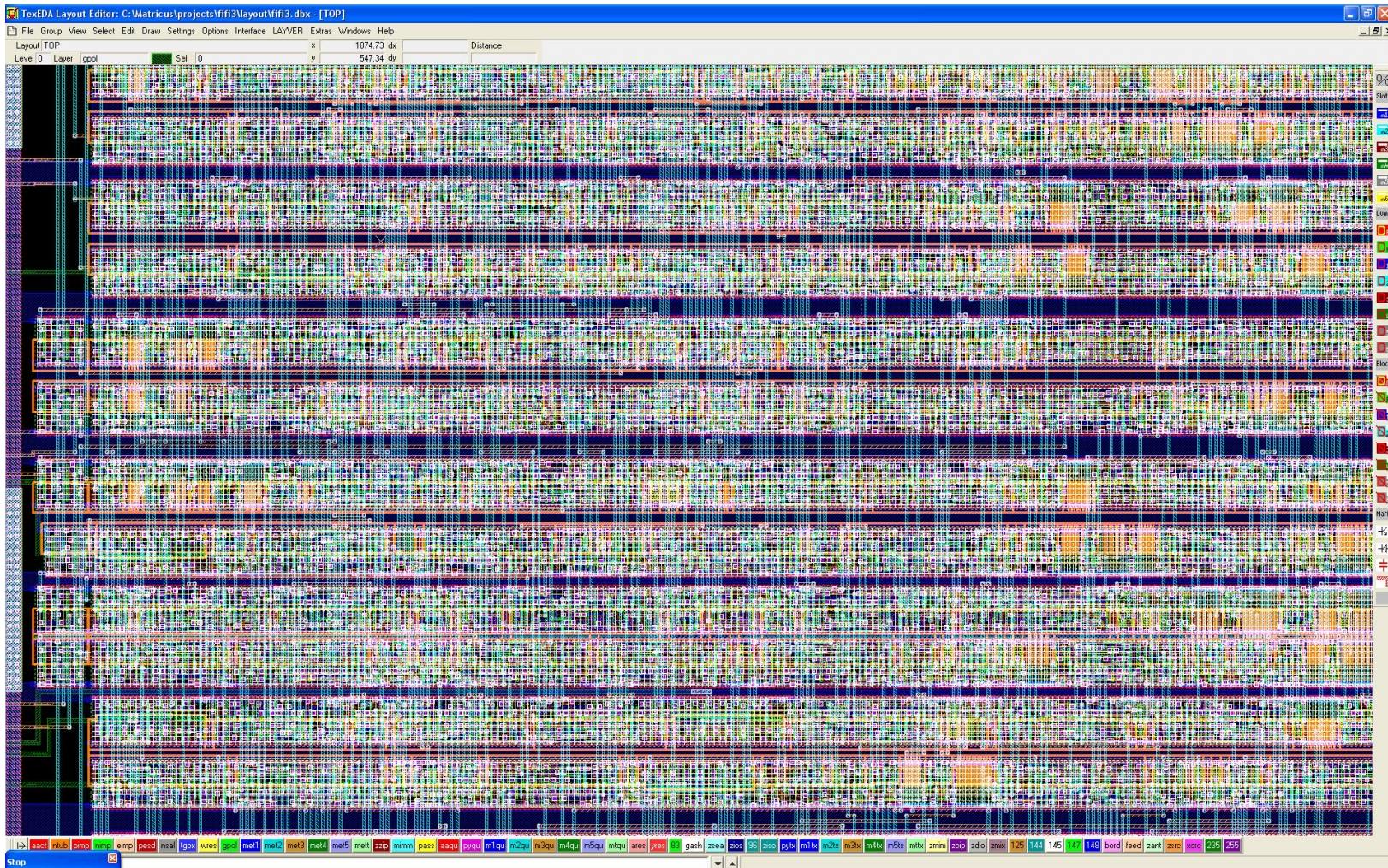
Standard Cell Layout Example



Standard Cell Area

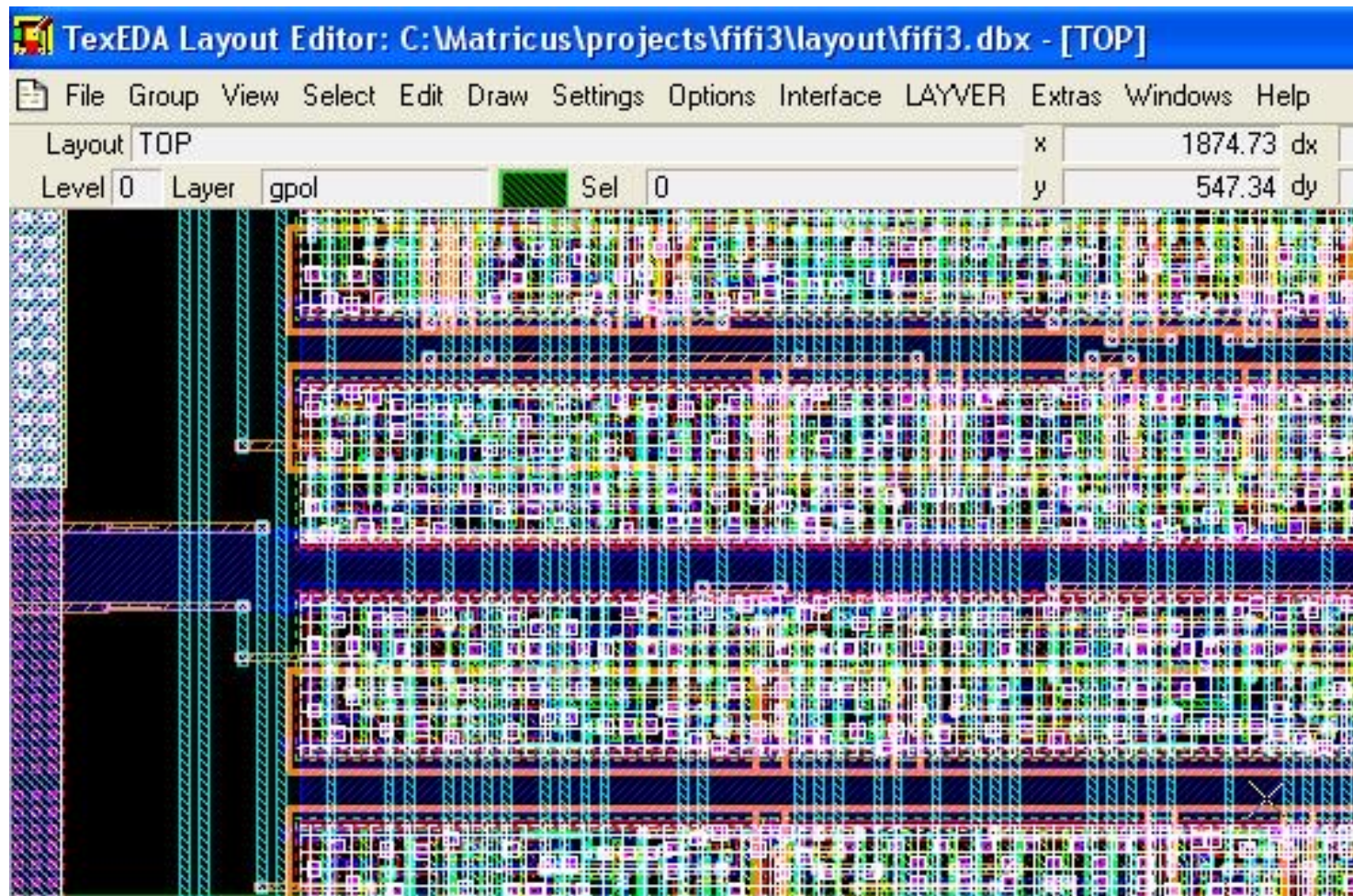


Standard Cell Layout Example



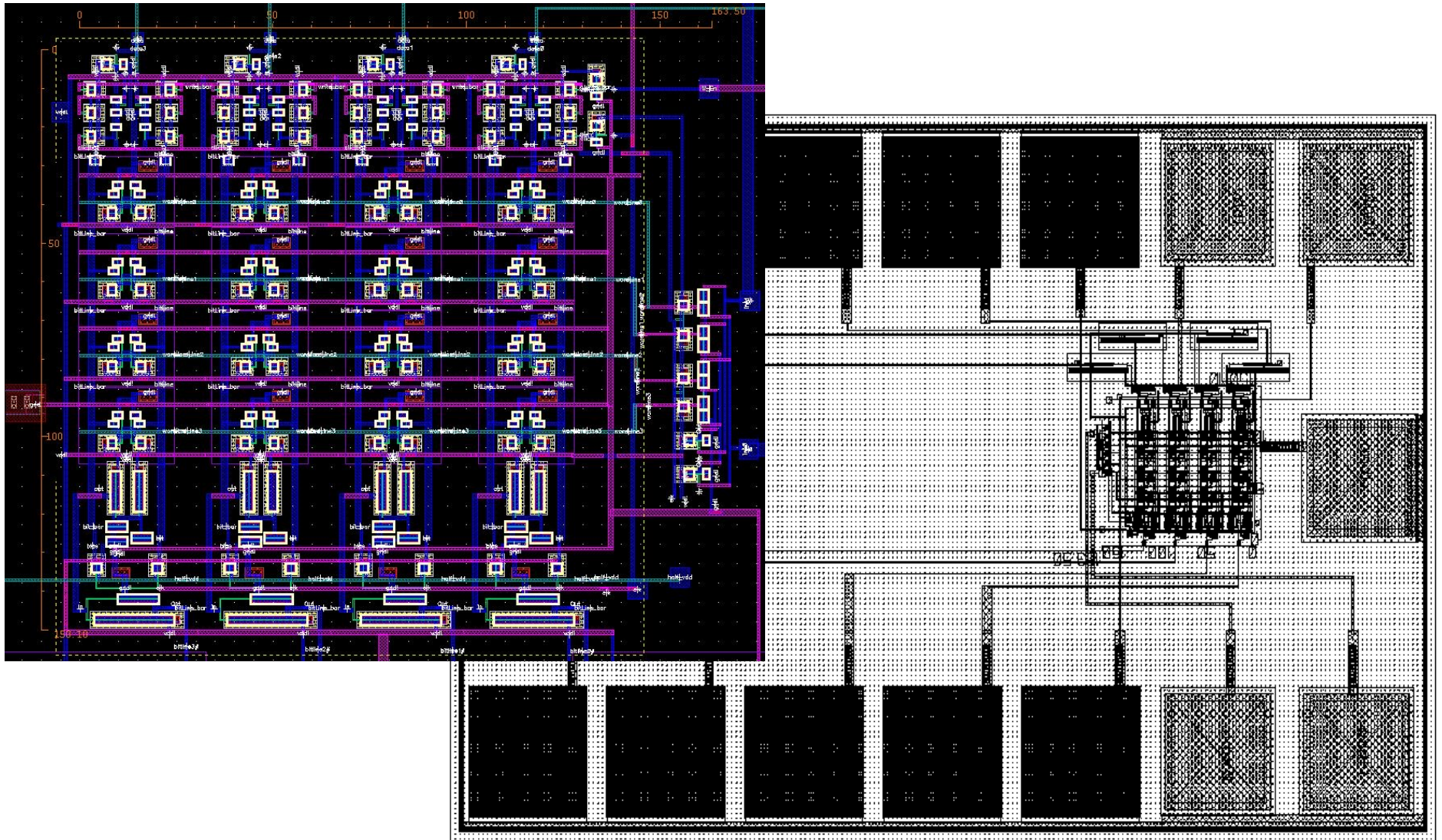
<http://www.laytools.com/images/StandardCells.jpg>

Standard Cell Layout Example



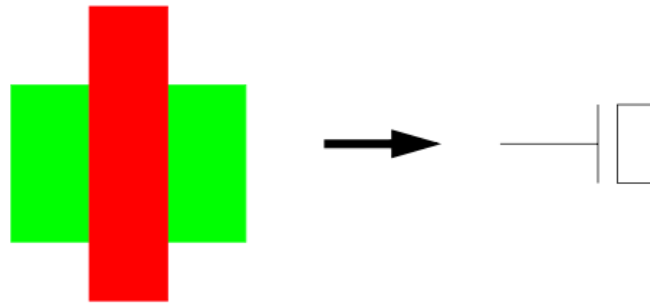
<http://www.laytools.com/images/StandardCells.jpg>

4x4 6T SRAM Memory



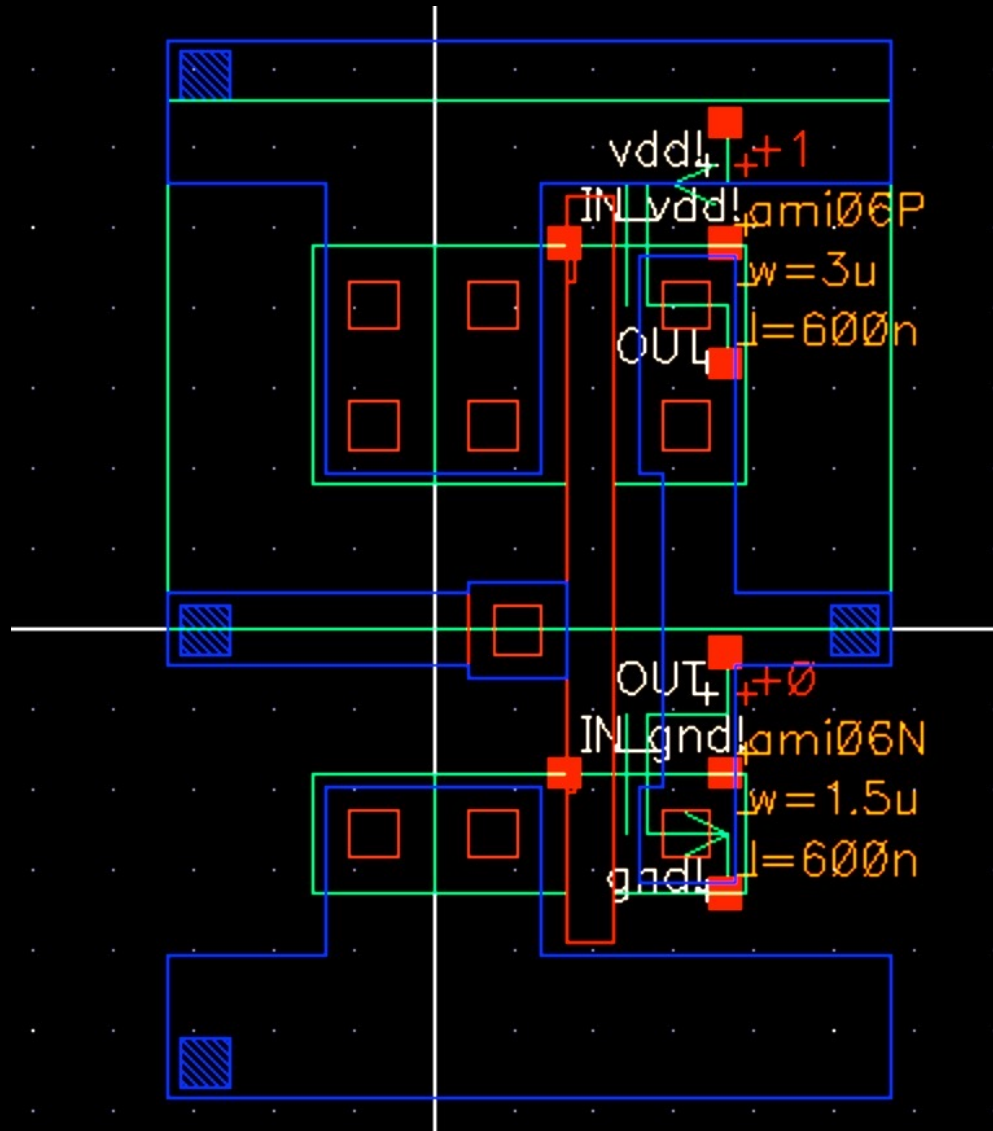
Circuit Extraction

- Circuit extraction extracts a schematic representation of a layout, including transistors, wires, and possibly wire and device resistance and capacitance.

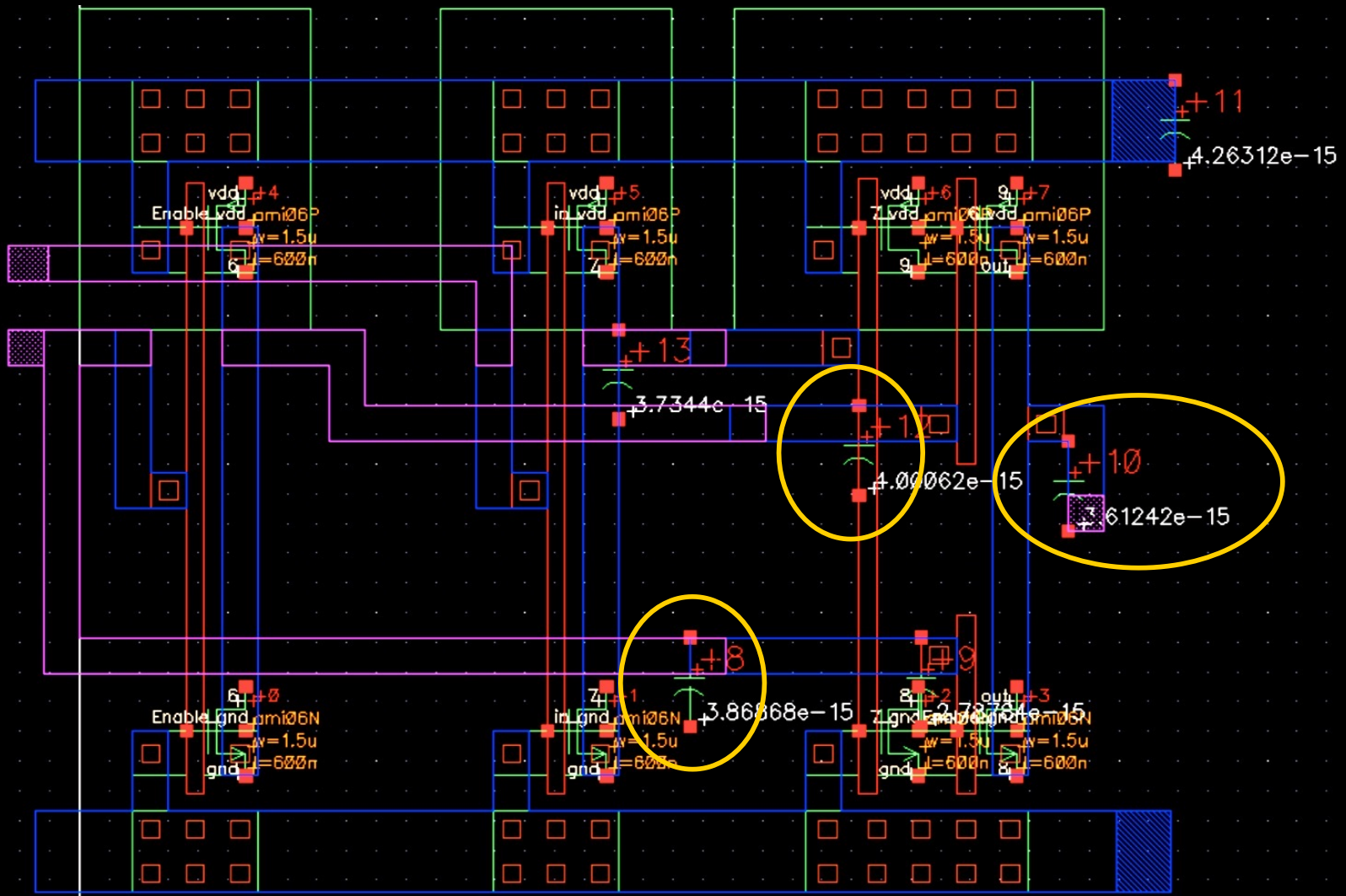


- Circuit extraction is used for LVS, and for spice simulation of layouts

Circuit Extraction

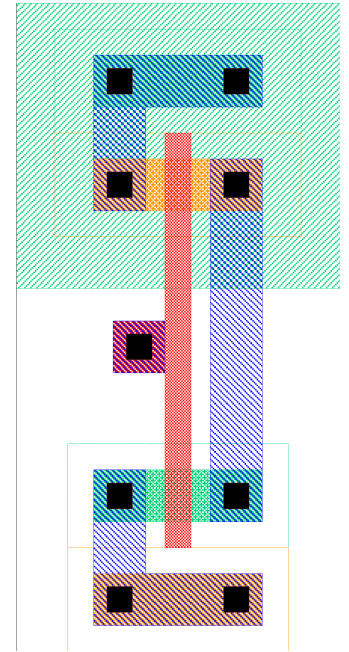


Circuit Extraction



Big Idea

- ❑ Layouts are physical realization of circuit
 - Geometry tradeoff
 - Can decrease spacing at the cost of yield
 - Design rules
- ❑ Can go from circuit to stick diagram/layout or stick diagram/layout to circuit by inspection
- ❑ Moderately predictable VLSI Scaling
 - unprecedented capacities/capability growth for engineered systems
 - ...but hits physical limit





Admin

- HW3 due Friday 2/14



Acknowledgement

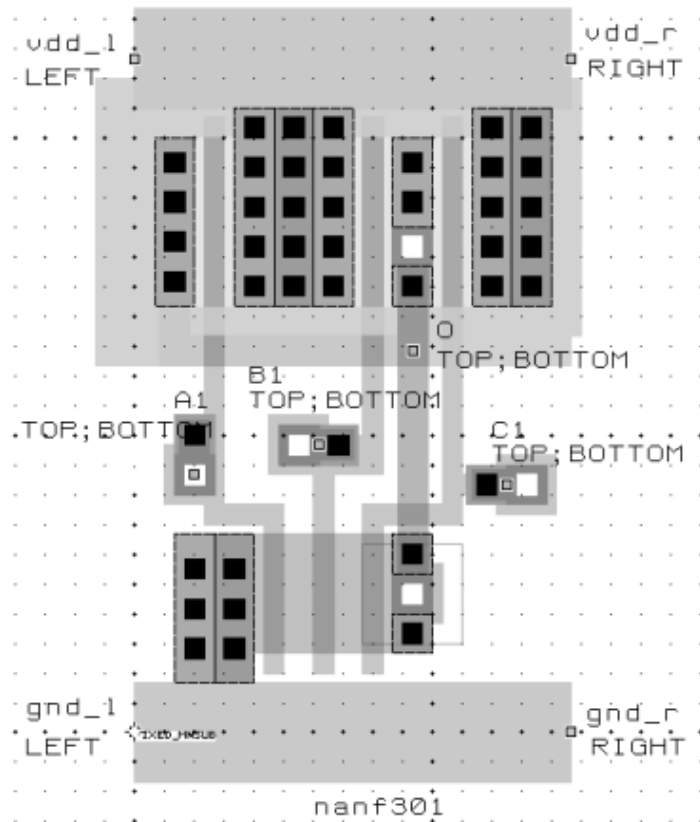
- ❑ Prof. André DeHon (University of Pennsylvania)
- ❑ Prof. Tania Khanna (University of Pennsylvania)

Scaling

(General Background Reading)



Standard Cells



Fanout 4x	0.5 μm	1.0 μm	2.0 μm
<i>A1_tphl</i>	0.595	0.711	0.919
<i>A1_tplh</i>	0.692	0.933	1.360
<i>B1_tphl</i>	0.591	0.739	1.006
<i>B1_tplh</i>	0.620	0.825	1.1.81
<i>C1_tphl</i>	0.574	0.740	1.029
<i>C1_tplh</i>	0.554	0.728	1.026

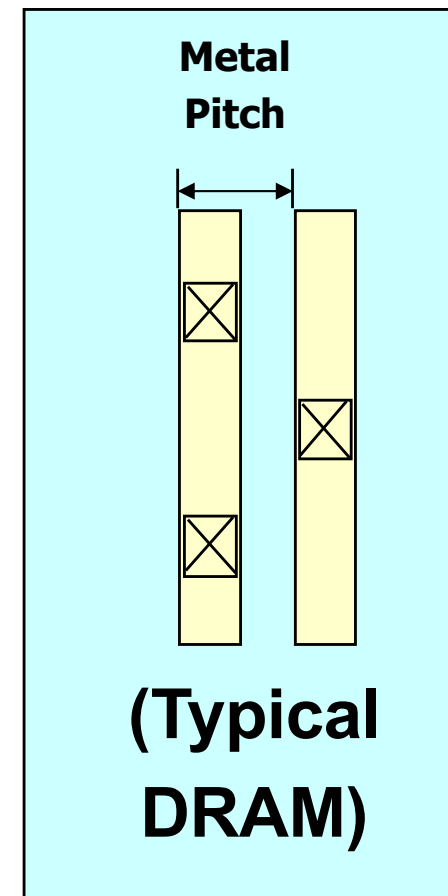
3-input NAND cell
 (from Mississippi State Library)
 characterized for fanout of 4 and
 for three different technologies

Scaling Technology

- ❑ **Premise:** features scale “uniformly”
 - everything gets better in a predictable manner

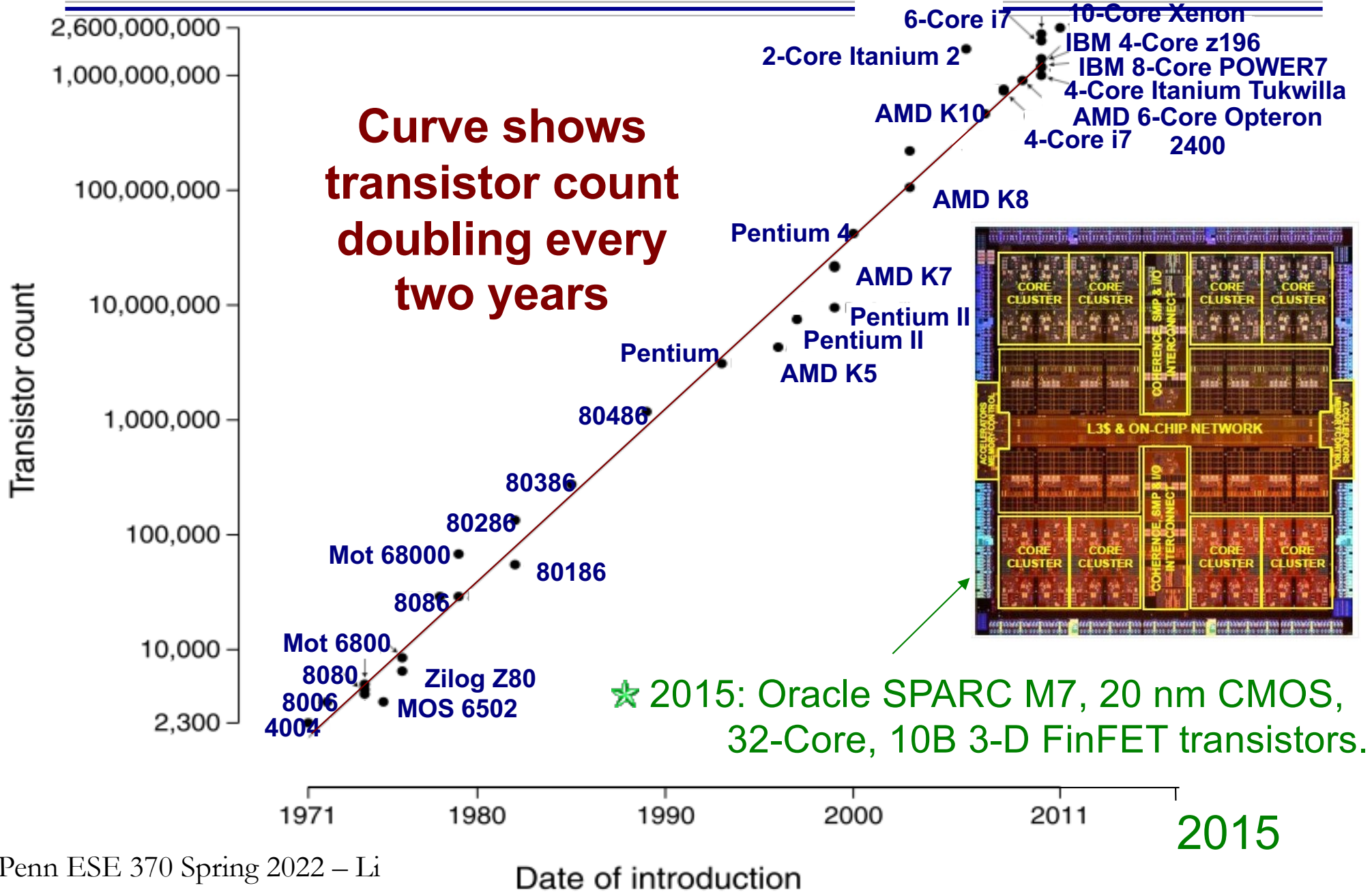
- ❑ **Parameters:**

- λ (lambda) -- Mead and Conway ($L=2\lambda$)
- F -- Half pitch – ITRS ($F=2\lambda=L$)
- S – scale factor – Rabaey
 - $F'=F/S$
 - $S>1$



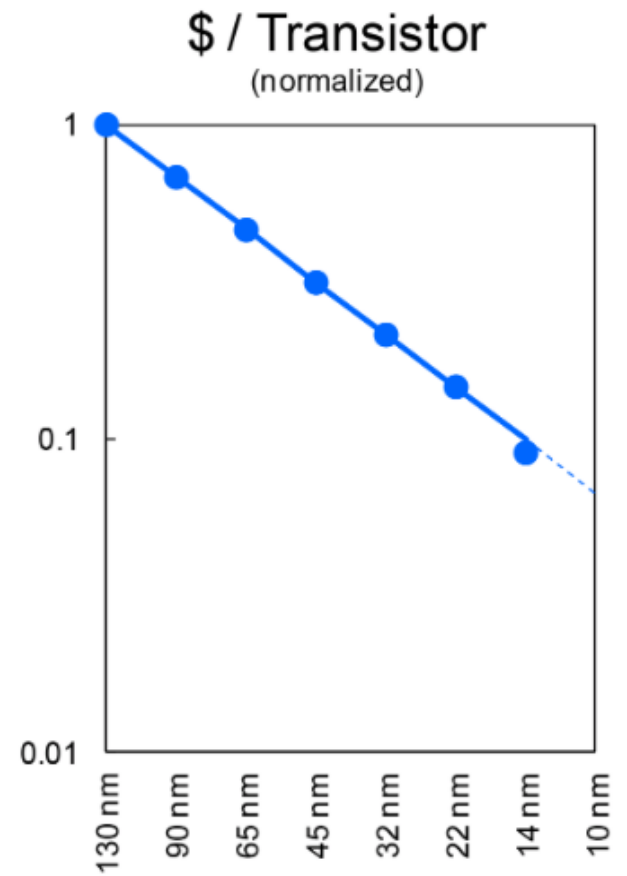
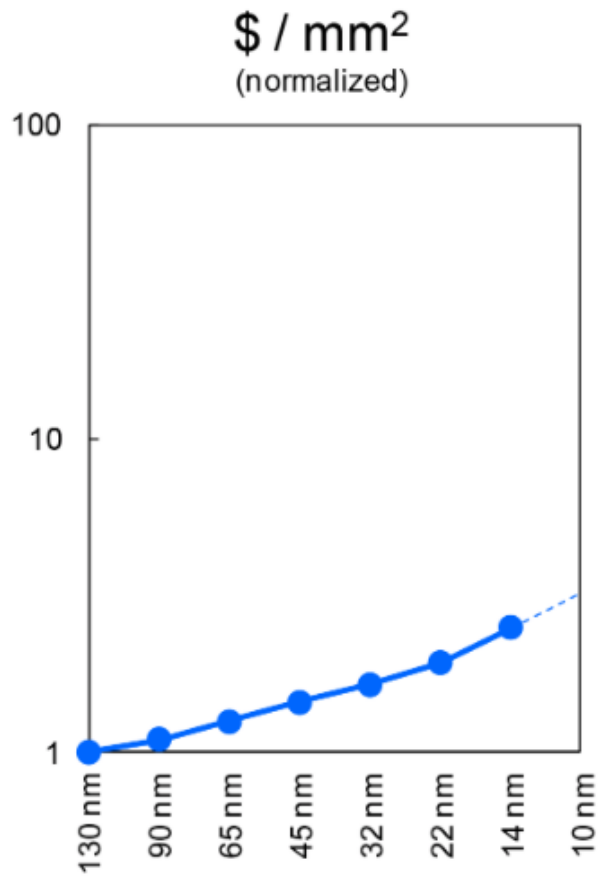
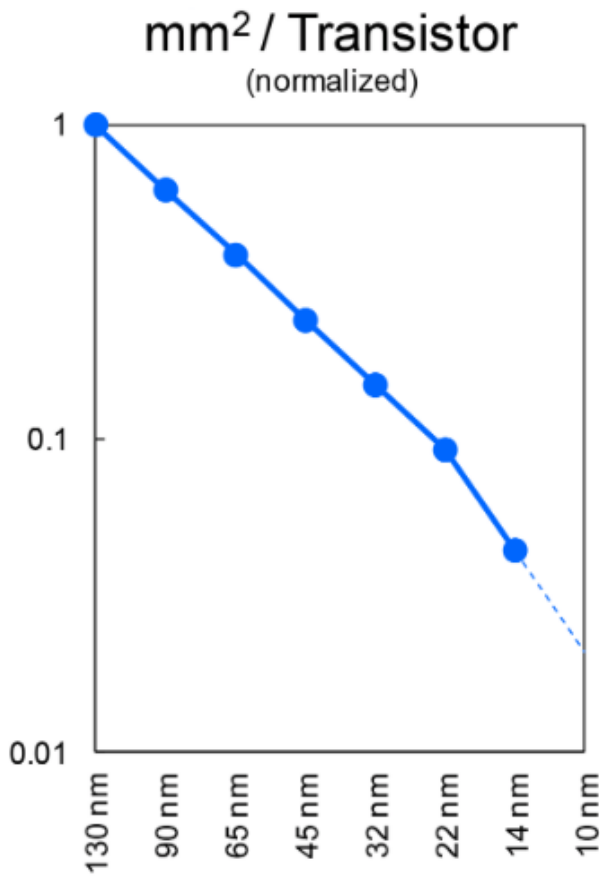
Microprocessor Trans Count 1971-2015

Curve shows transistor count doubling every two years





Intel Cost Scaling



<http://www.anandtech.com/show/8367/intels-14nm-technology-in-detail>

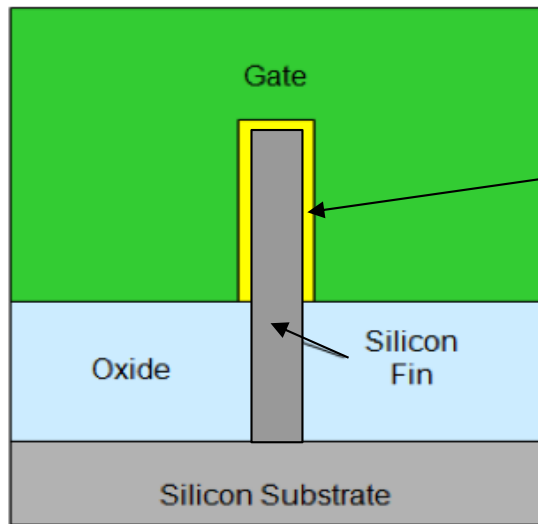
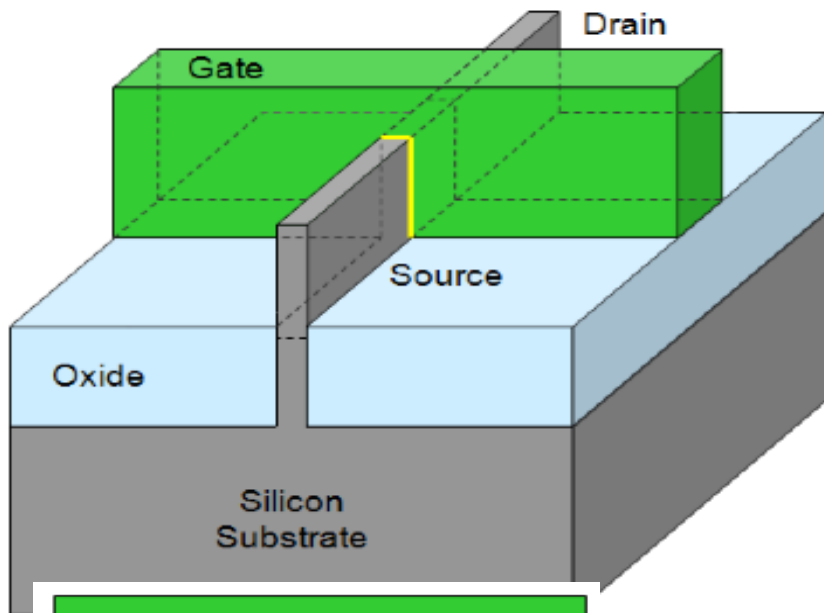


More Moore → Scaling

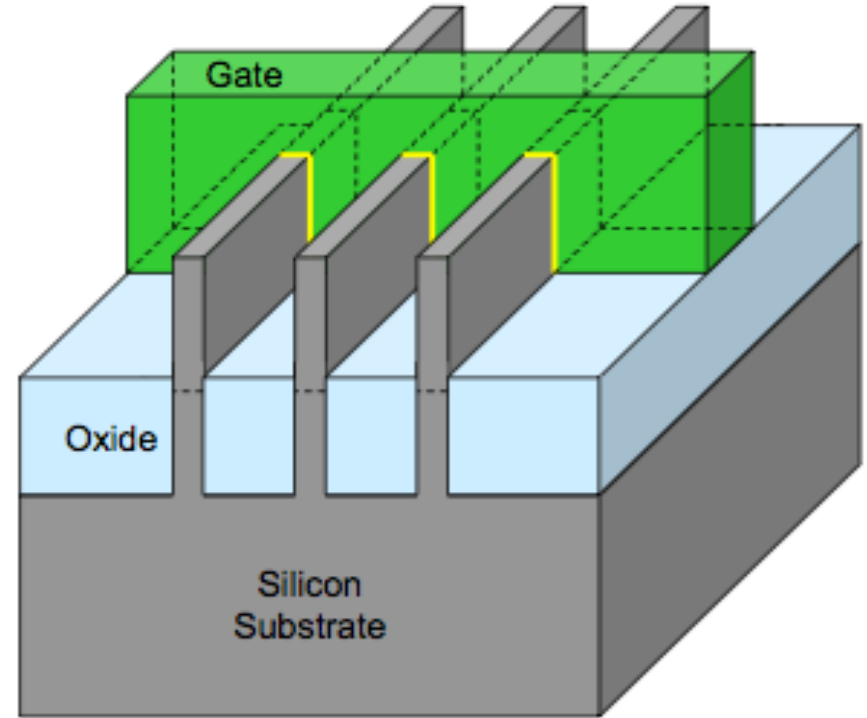
- ❑ Geometrical Scaling
 - continued shrinking of horizontal and vertical physical feature sizes

- ❑ Design Equivalent Scaling
 - design technologies that enable high performance, low power, high reliability, low cost, and high design productivity even if neither geometrical nor equivalent scaling can be used

22nm 3D FinFET Transistor



High-k
gate
dielectric



Tri-Gate transistors with multiple fins connected together increases total drive strength for higher performance

http://download.intel.com/newsroom/kits/22nm/pdfs/22nm-Details_Presentation.pdf

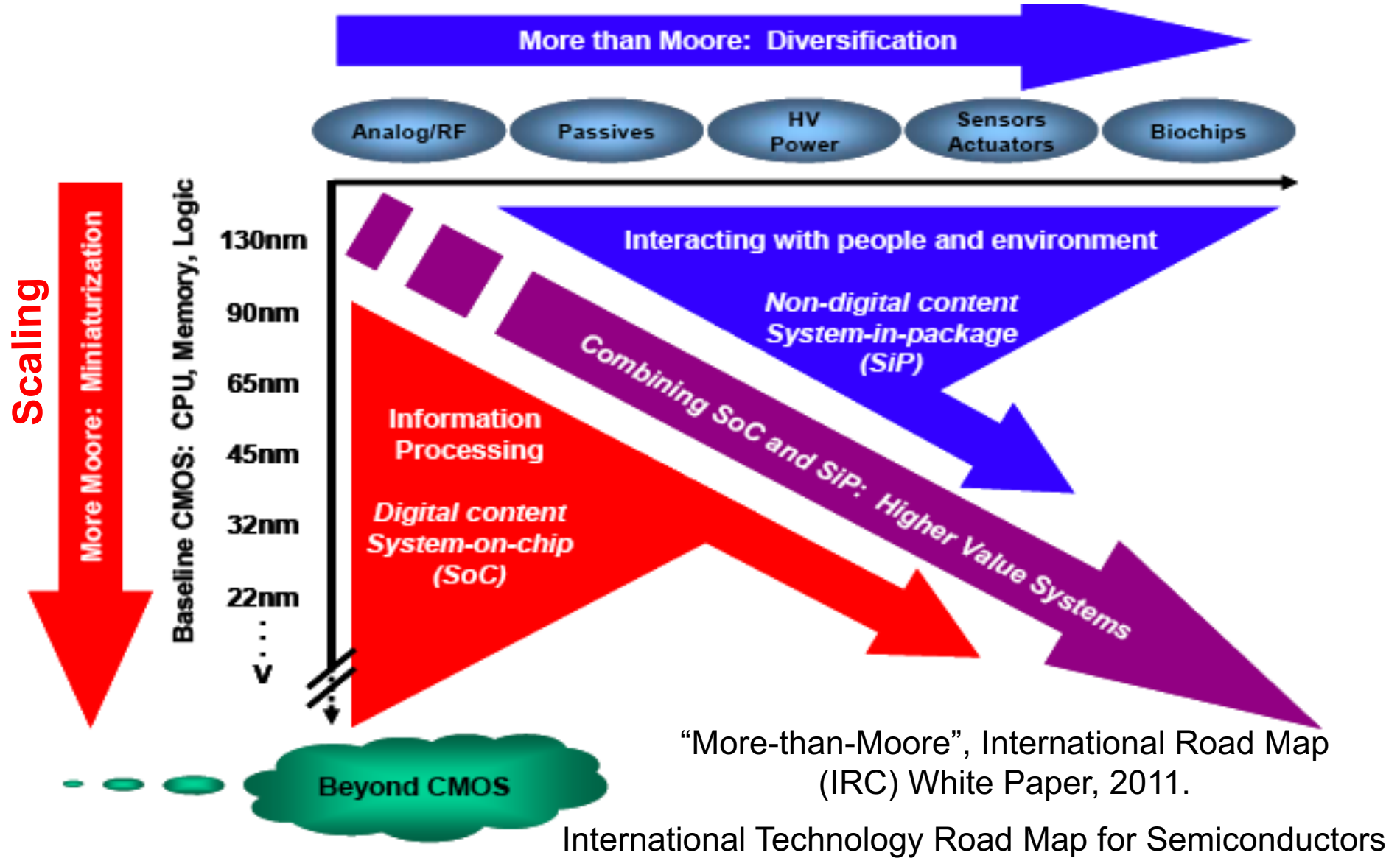


ITRS Roadmap

- ❑ International Technology Roadmap for Semiconductors
 - Try to predict where industry going
- ❑ ITRS 2.0 started in 2015 with new focus
 - System Integration, Heterogeneous Integration, Heterogeneous Components, Outside System Connectivity, More Moore, Beyond CMOS and Factory Integration.

- ❑ <http://www.itrs2.net/>

More-than-Moore





Question

- Scaling from 32nm \rightarrow 22nm, what is 1/S?
 - Scaling minimum gate length
 - And pitch distance

MOS Transistor **Scaling** - (1974 to present)

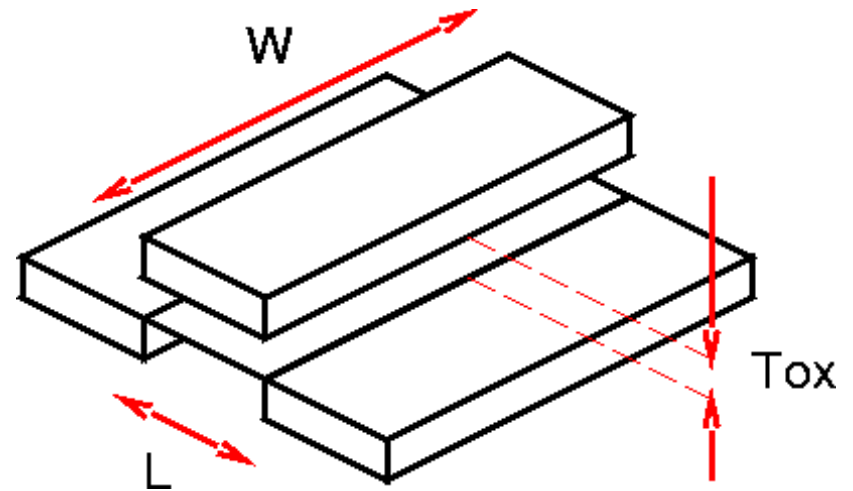
$1/S=0.7$
per technology node
[0.5x per 2 nodes]



**Source: 2001 ITRS - Exec. Summary, ORTC
Figure, Andrew Kahng**

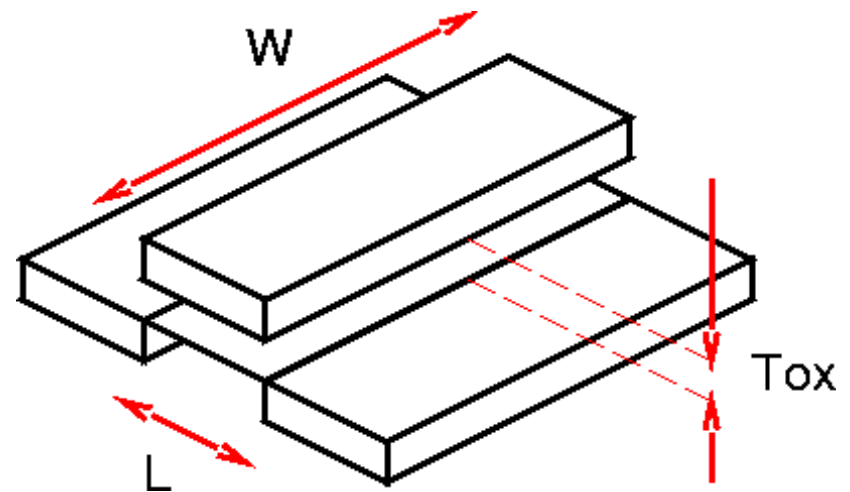
Scaling

- ❑ Channel Length (L)
- ❑ Channel Width (W)
- ❑ Oxide Thickness (t_{ox})
- ❑ Doping (N_a)
- ❑ Voltage (V_{DD}, V_t)



Full Scaling (Ideal Scaling)

- ❑ Channel Length (L) $1/S$
- ❑ Channel Width (W) $1/S$
- ❑ Oxide Thickness (t_{ox}) $1/S$
- ❑ Doping (N_a) S
- ❑ Voltage (V_{DD}, V_{t}) $1/S$



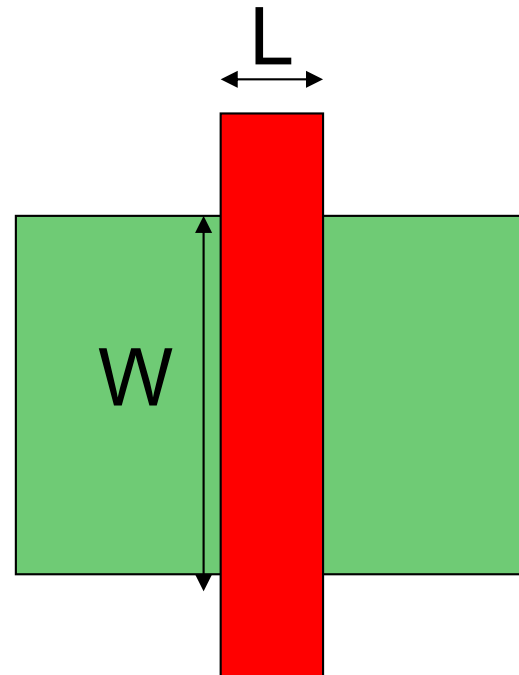


Effects on Physical Properties and Specs?

- ❑ Area
- ❑ Capacitance
 - C_{ox} and C_{gate}
- ❑ Resistance
- ❑ Current (I_d)
- ❑ Gate Delay (τ_{gd})
- ❑ Wire Delay (τ_{wire})
- ❑ Power
 - Same frequency
 - Scaled frequency
- ❑ Power Density
 - Same frequency
 - Scaled frequency

Area

- $\lambda' \rightarrow \lambda/S$
- Area impact?
- $A = L \times W$
-



Area

- $\lambda' \rightarrow \lambda/S$

- Area impact?

- $A = L \times W$

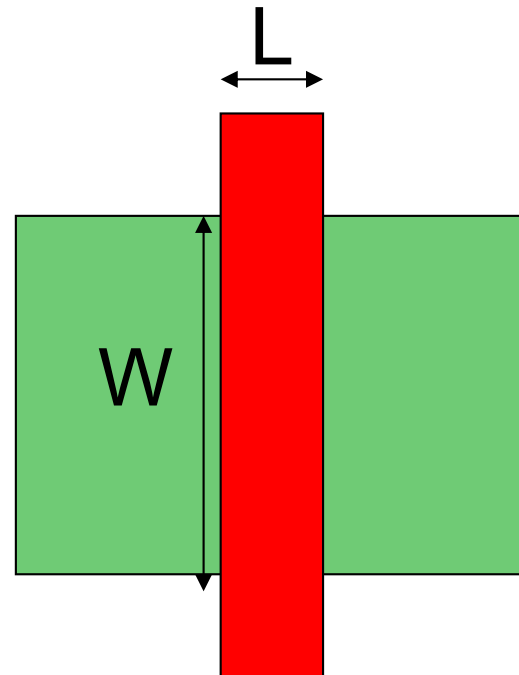
- $A' \rightarrow A/S^2$

- $32\text{nm} \rightarrow 22\text{nm}$

- 50% area

- $2 \times$ transistor capacity
for same area

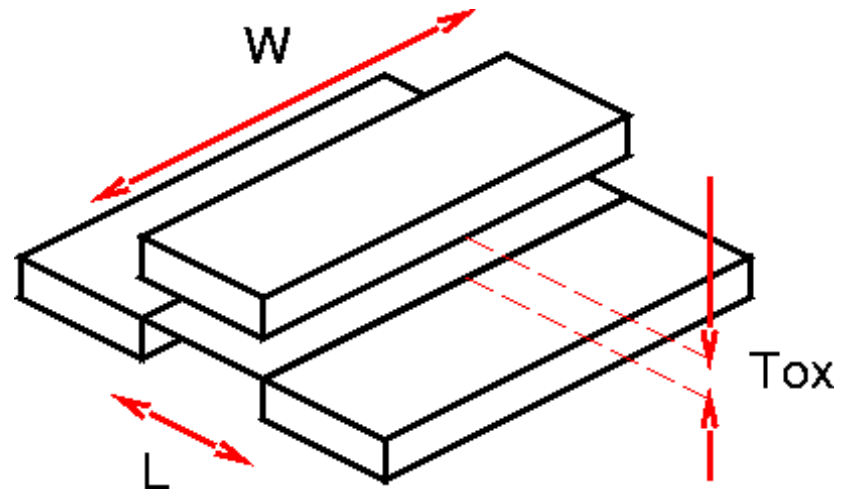
$$1/S=0.7$$



Capacitance

□ Capacitance per unit area scaling?

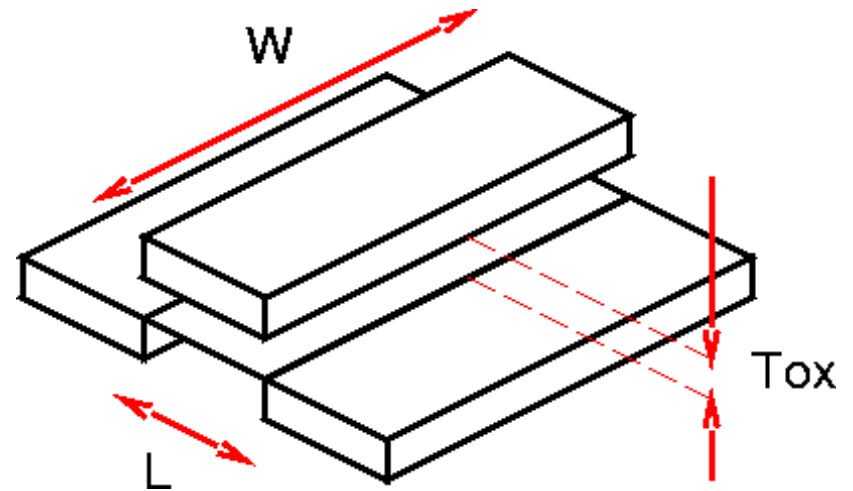
- $C_{\text{ox}} = \epsilon_{\text{SiO}_2} / t_{\text{ox}}$
- $t'_{\text{ox}} \rightarrow t_{\text{ox}} / S$



Capacitance

□ Capacitance per unit area scaling?

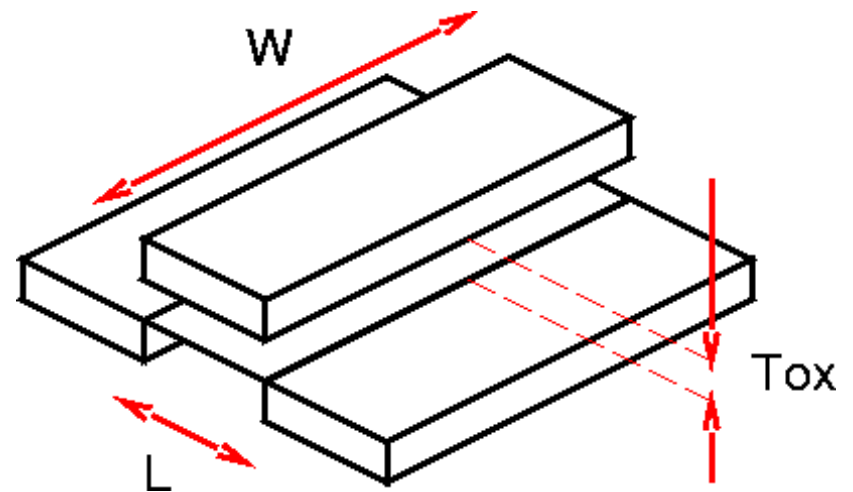
- $C_{\text{ox}} = \epsilon_{\text{SiO}_2} / t_{\text{ox}}$
- $t'_{\text{ox}} \rightarrow t_{\text{ox}} / S$
- $C'_{\text{ox}} \rightarrow C_{\text{ox}} \times S$



Capacitance

□ Gate Capacitance scaling?

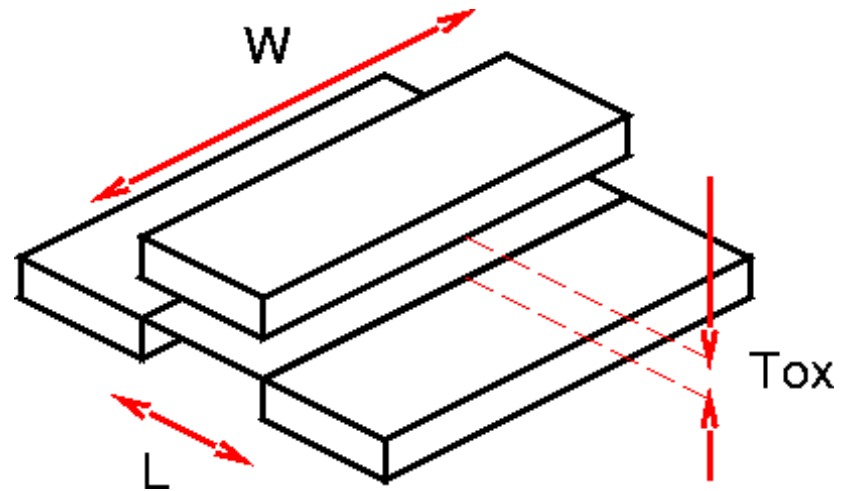
- $C_{\text{gate}} = A \times C_{\text{ox}}$
- $A' \rightarrow A/S^2$
- $C'_{\text{ox}} \rightarrow C_{\text{ox}} \times S$



Capacitance

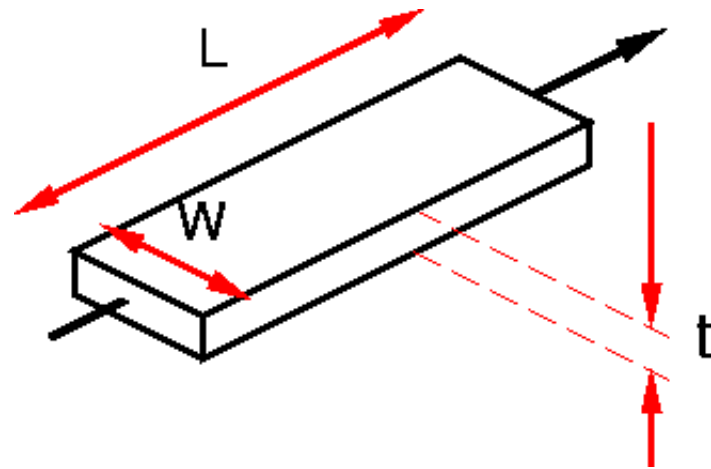
□ Gate Capacitance scaling?

- $C_{\text{gate}} = A \times C_{\text{ox}}$
- $A' \rightarrow A/S^2$
- $C'_{\text{ox}} \rightarrow C_{\text{ox}} \times S$
- $C'_{\text{gate}} \rightarrow C_{\text{gate}}/S$



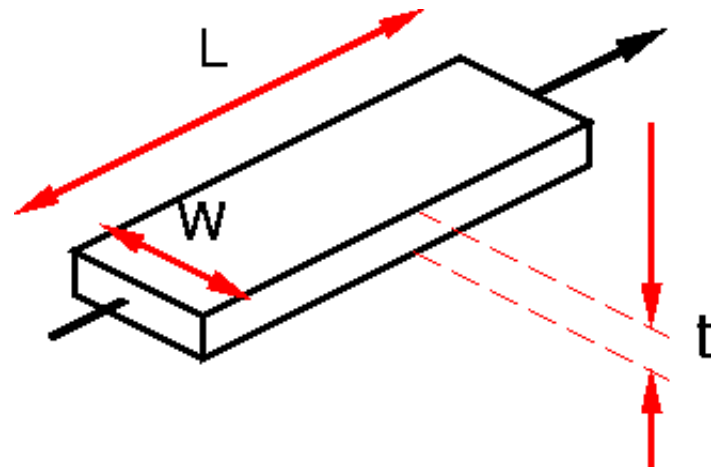
Wire Resistance

- Resistance scaling?
- $R = \rho L / (W * t)$
 - L, t remain similar (not scaled)
- $W \rightarrow W / S$



Wire Resistance

- Resistance scaling?
- $R = \rho L / (W * t)$
 - L, t remain similar (not scaled)
- $W \rightarrow W / S$
- $R' \rightarrow R \times S$



Current

- ❑ Which Voltages matters here? ($V_{gs}, V_{ds}, V_{th} \dots$)
- ❑ Transistor charging looks like voltage-controlled current source
- ❑ Saturation Current scaling?

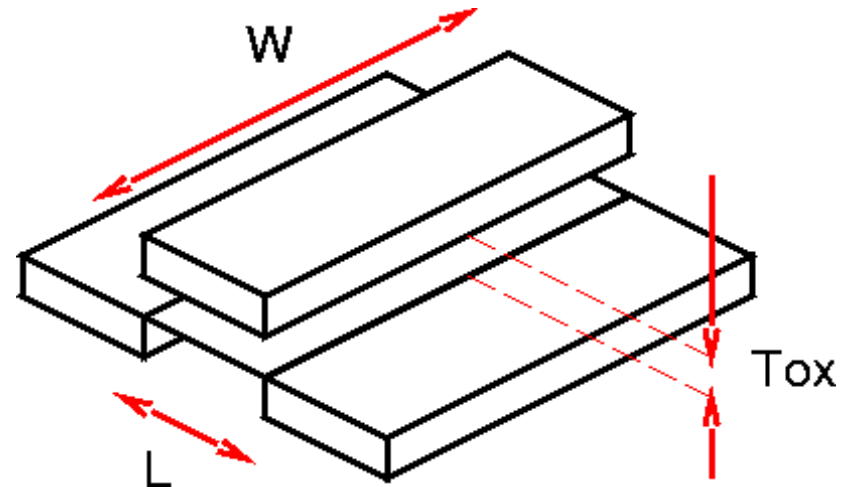
$$I_d = (\mu C_{OX} / 2) (W / L) (V_{gs} - V_{TH})^2$$

$$V_{gs}, V_{TH}: V' \rightarrow V / S$$

$$W' \rightarrow W / S$$

$$L' \rightarrow L / S$$

$$C'_{ox} \rightarrow C_{ox} \times S$$



Current

- Which Voltages matters here? ($V_{gs}, V_{ds}, V_{th} \dots$)
- Transistor charging looks like voltage-controlled current source
- Saturation Current scaling?

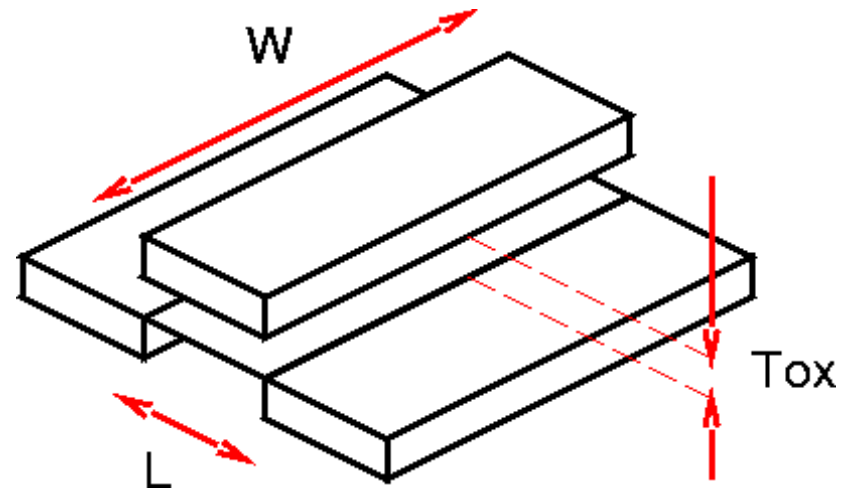
$$I_d = (\mu C_{OX} / 2) (W / L) (V_{gs} - V_{TH})^2$$

$$V_{gs}, V_{TH}: V' \rightarrow V / S$$

$$W' \rightarrow W / S$$

$$L' \rightarrow L / S$$

$$C'_{ox} \rightarrow C_{ox} \times S$$



$$I'_d = (\mu C_{OX} S / 2) ((W / S) / (L / S)) (V_{gs} / S - V_{TH} / S)^2$$

Current

- ❑ Which Voltages matters here? ($V_{gs}, V_{ds}, V_{th} \dots$)
- ❑ Transistor charging looks like voltage-controlled current source
- ❑ Saturation Current scaling?

$$I_d = (\mu C_{OX} / 2) (W / L) (V_{gs} - V_{TH})^2$$

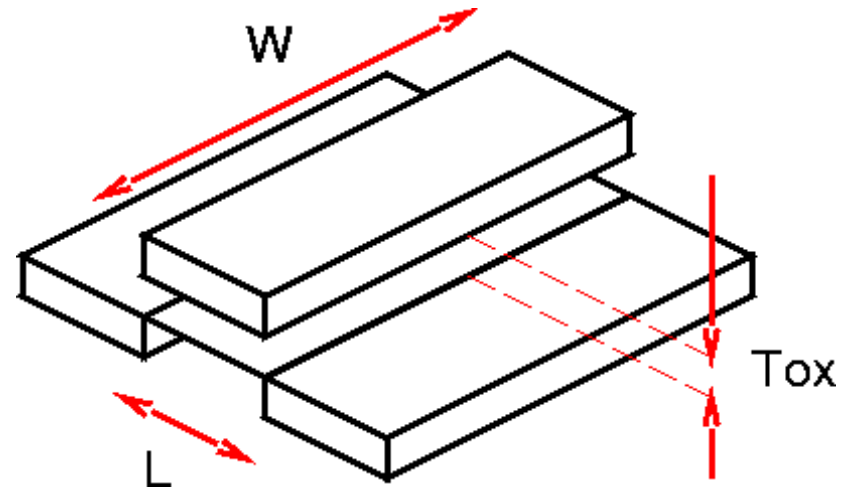
$$V_{gs}, V_{TH}: V' \rightarrow V / S$$

$$W' \rightarrow W / S$$

$$L' \rightarrow L / S$$

$$C'_{ox} \rightarrow C_{ox} \times S$$

$$I'_d \rightarrow I_d / S$$



Current

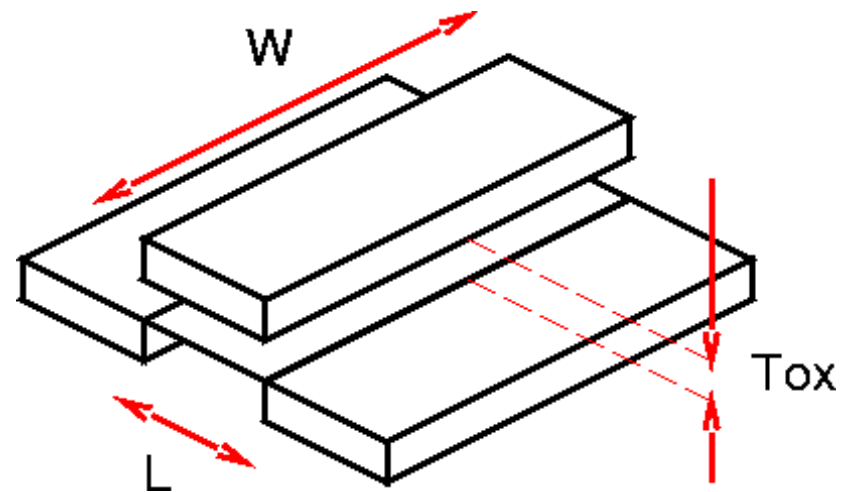
- Velocity Saturation Current scaling?

$$V_{gs}, V_{TH}: V' \rightarrow V/S$$

$$L' \rightarrow L/S$$

$$W' \rightarrow W/S$$

$$C'_{ox} \rightarrow C_{ox}S$$



Current

□ Velocity Saturation Current scaling?

$$V_{gs}, V_{TH}: V' \rightarrow V/S$$

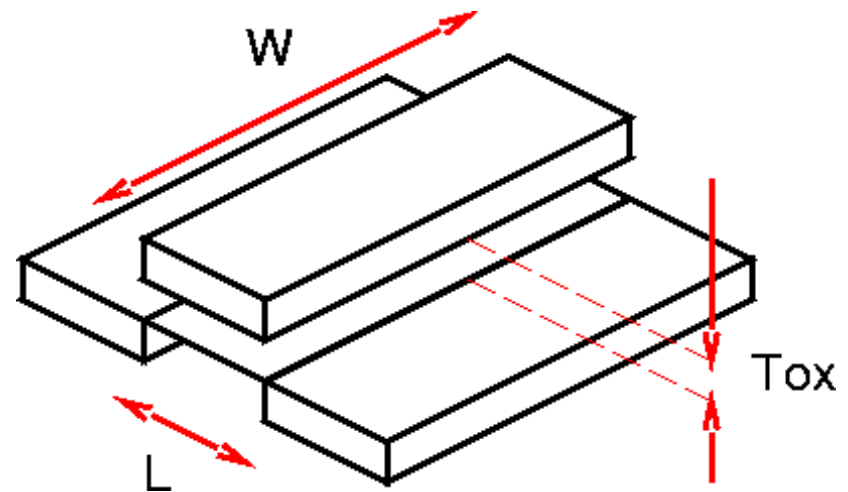
$$L' \rightarrow L/S$$

$$W' \rightarrow W/S$$

$$C'_{ox} \rightarrow C_{ox} S$$

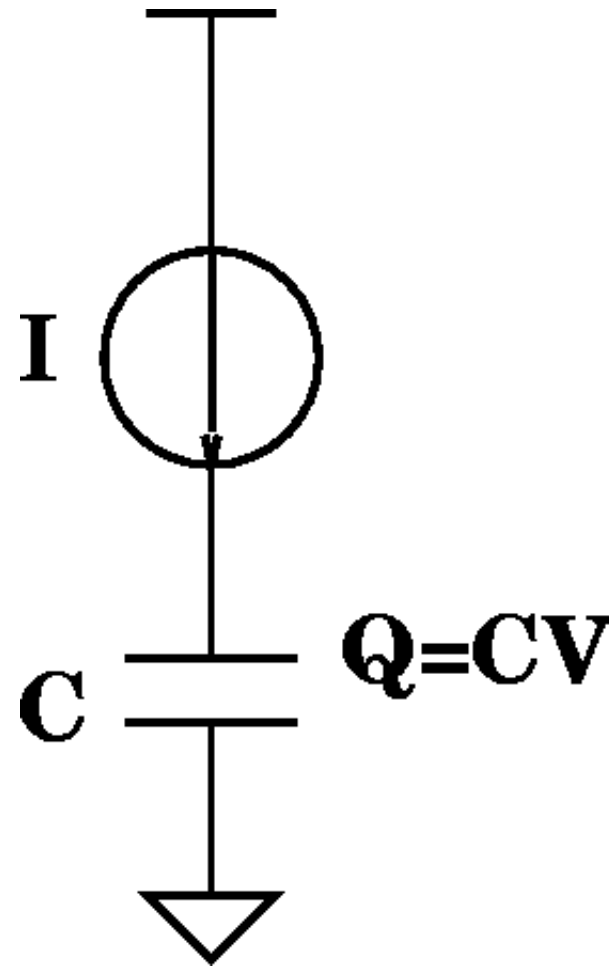
$$I'_d \rightarrow I_d/S$$

$$I_{DS} \approx v_{sat} C_{OX} W \left(V_{GS} - V_{TH} - \frac{V_{DSAT}}{2} \right)$$



Gate Delay

- Gate Delay scaling?
- $\tau_{gd} = Q/I = (CV)/I$
- $V' \rightarrow V/S$
- $I'_d \rightarrow I_d/S$
- $C'_g \rightarrow C_g/S$



Gate Delay

- Gate Delay scaling?

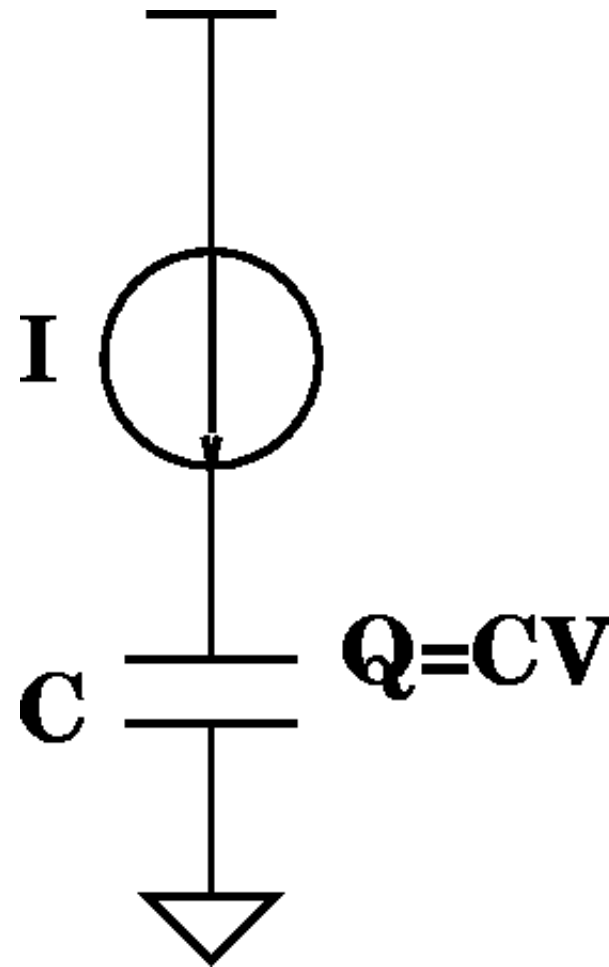
- $\tau_{gd} = Q/I = (CV)/I$

- $V' \rightarrow V/S$

- $I'_d \rightarrow I_d/S$

- $C'_g \rightarrow C_g/S$

- $\tau'_{gd} \rightarrow \tau_{gd}/S$





Wire Delay

- Wire delay scaling?
- $\tau_{\text{wire}} = R \times C$
- $R' \rightarrow R \times S$
- $C' \rightarrow C / S$
- Again assuming (logical) wire lengths remain constant



Wire Delay

- Wire delay scaling?
- $\tau_{\text{wire}} = R \times C$
- $R' \rightarrow R \times S$
- $C' \rightarrow C / S$
- $\tau'_{\text{wire}} \rightarrow \tau_{\text{wire}}$
- Again assuming (logical) wire lengths remain constant



Power Dissipation (Dynamic)

□ Capacitive (Dis)charging scaling?

□ $P = (1/2)CV^2f$

□ $V' \rightarrow V/S$

□ $C' \rightarrow C/S$



Power Dissipation (Dynamic)

□ Capacitive (Dis)charging scaling?

□ $P = (1/2)CV^2f$

□ $V' \rightarrow V/S$

□ $C' \rightarrow C/S$

□ $P' \rightarrow P/S^3$

Power Dissipation (Dynamic)

□ Capacitive (Dis)charging scaling?

$$\square P = (1/2)CV^2f$$

$$\square V' \rightarrow V/S$$

$$\square C' \rightarrow C/S$$

$$\square P' \rightarrow P/S^3$$

□ Increase Frequency?

$$\square \tau_{gd} \rightarrow \tau_{gd}/S$$

$$\square \text{So: } f \rightarrow f \times S$$

Power Dissipation (Dynamic)

□ Capacitive (Dis)charging scaling?

$$\square P = (1/2)CV^2f$$

$$\square V' \rightarrow V/S$$

$$\square C' \rightarrow C/S$$

$$\square P' \rightarrow P/S^3$$

□ Increase Frequency?

$$\square \tau_{gd} \rightarrow \tau_{gd}/S$$

$$\square \text{So: } f \rightarrow f \times S$$

$$\square P \rightarrow P/S^2$$



Effects?

- Area $1/S^2$
- Capacitance (C_{ox}, C_g) $S, 1/S$
- Resistance S
- Threshold (V_{th}) $1/S$
- Current (I_d) $1/S$
- Gate Delay (τ_{gd}) $1/S$
- Wire Delay (τ_{wire}) 1
- Power $1/S^3, 1/S^2$

$1/S=0.7$



Power Density

- ❑ $P' \rightarrow P/S^2$ (increased frequency)
- ❑ $P' \rightarrow P/S^3$ (same frequency)
- ❑ $A' \rightarrow A/S^2$

- ❑ Power Density: P/A two cases?



Power Density

- $P' \rightarrow P/S^2$ (increased frequency)
- $P' \rightarrow P/S^3$ (same frequency)
- $A' \rightarrow A/S^2$

- Power Density: P/A two cases?
 - $P/A \rightarrow P/A$ increase freq.
 - $P/A \rightarrow (P/A)/S$ same freq.

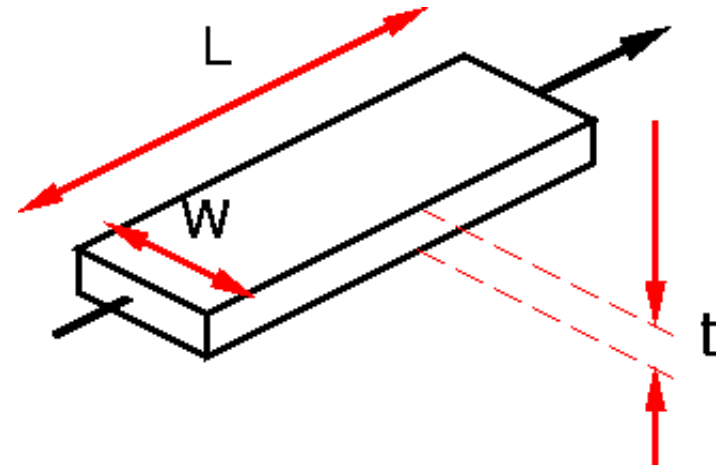


But...

- ❑ Don't like some of the implications
 - ❑ High resistance wires
 - ❑ Higher gate oxide capacitance with atomic-scale dimensions
 - ❑ Quantum tunneling
 - ❑ Need for more wiring
 - ❑ Not scale speed fast enough

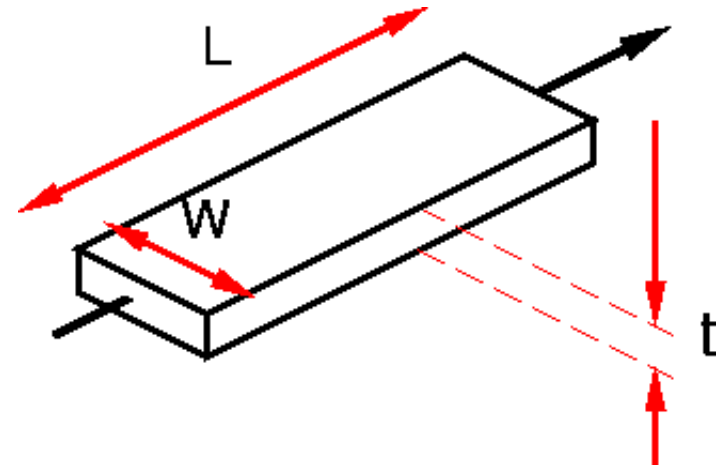
Improving Resistance

- $R = \rho L / (W \times t)$
- $W' \rightarrow W/S$
 - L, t similar
- $R' \rightarrow R \times S$



Improving Resistance

- $R = \rho L / (W \times t)$
- $W' \rightarrow W/S$
 - L, t similar
- $R' \rightarrow R \times S$



What might we do?

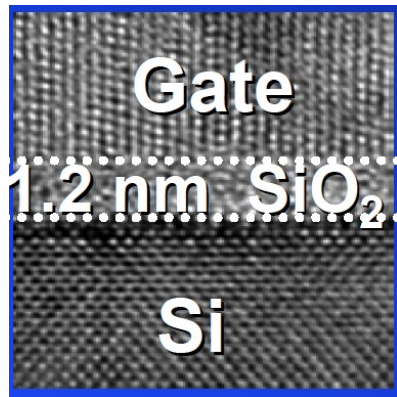
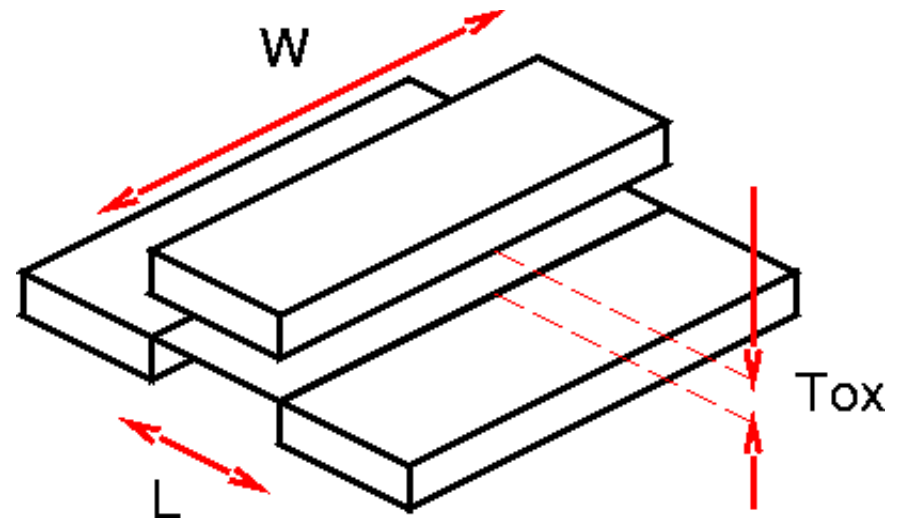
Decrease ρ (copper) – introduced 1997

<http://www.ibm.com/ibm100/us/en/icons/copperchip/>

Capacitance and Leakage

□ Capacitance per unit area

- $C_{ox} = \epsilon_{SiO_2} / t_{ox}$
- $t'_{ox} \rightarrow t_{ox} / S$
- $C'_{ox} \rightarrow C_{ox} \times S$



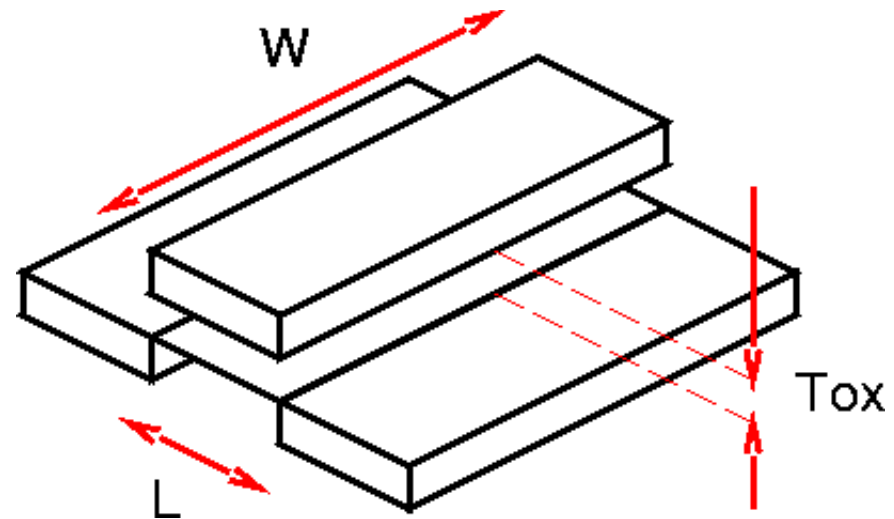
What's wrong with $t_{ox} = 1.2\text{nm}$?

source: Borkar/Micro 2004

Capacitance and Leakage

□ Capacitance per unit area

- $C_{\text{ox}} = \epsilon_{\text{SiO}_2} / t_{\text{ox}}$
- $t'_{\text{ox}} \rightarrow t_{\text{ox}} / S$
- $C'_{\text{ox}} \rightarrow C_{\text{ox}} \times S$



What might we do?

Reduce dielectric constant, ϵ , and not scale thickness to mimic t_{ox} scaling.

High-K dielectric Survey

Table 2 Selected material and electrical properties of high-*k* gate dielectrics. Data compiled from Robertson [25], Gusev et al. [20], Hubbard and Schlom [19], and other sources.

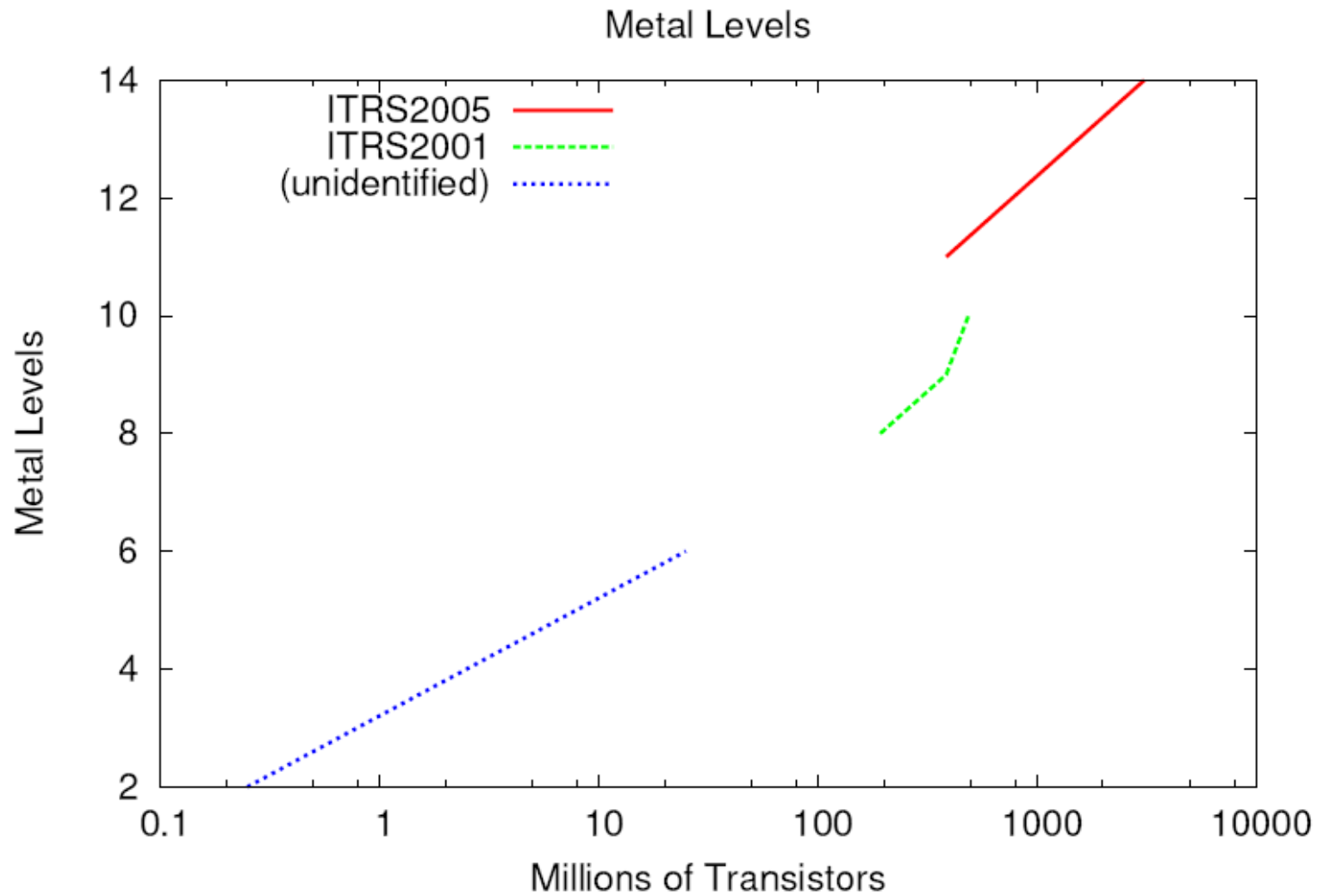
<i>Dielectric</i>	<i>Dielectric constant (bulk)</i>	<i>Bandgap (eV)</i>	<i>Conduction band offset (eV)</i>	<i>Leakage current reduction w.r.t. SiO₂</i>	<i>Thermal stability w.r.t. silicon (MEIS data)</i>
Silicon dioxide (SiO ₂)	3.9	9	3.5	N/A	>1050°C
Silicon nitride (Si ₃ N ₄)	7	5.3	2.4		>1050°C
Aluminum oxide (Al ₂ O ₃)	~10	8.8	2.8	10 ² –10 ³ ×	~1000°C, RTA
Tantalum pentoxide (Ta ₂ O ₅)	25	4.4	0.36		Not thermodynamically stable with silicon
Lanthanum oxide (La ₂ O ₃)	~21	6*	2.3		
Gadolinium oxide (Gd ₂ O ₃)	~12				
Yttrium oxide (Y ₂ O ₃)	~15	6	2.3	10 ⁴ –10 ⁵ ×	Silicate formation
Hafnium oxide (HfO ₂)	~20	6	1.5	10 ⁴ –10 ⁵ ×	~950°C
Zirconium oxide (ZrO ₂)	~23	5.8	1.4	10 ⁴ –10 ⁵ ×	~900°C
Strontium titanate (SrTiO ₃)		3.3	–0.1		
Zirconium silicate (ZrSiO ₄)		6*	1.5		
Hafnium silicate (HfSiO ₄)		6*	1.5		

*Estimated value.

Wong/IBM J. of R&D, V46N2/3P133—168, 2002



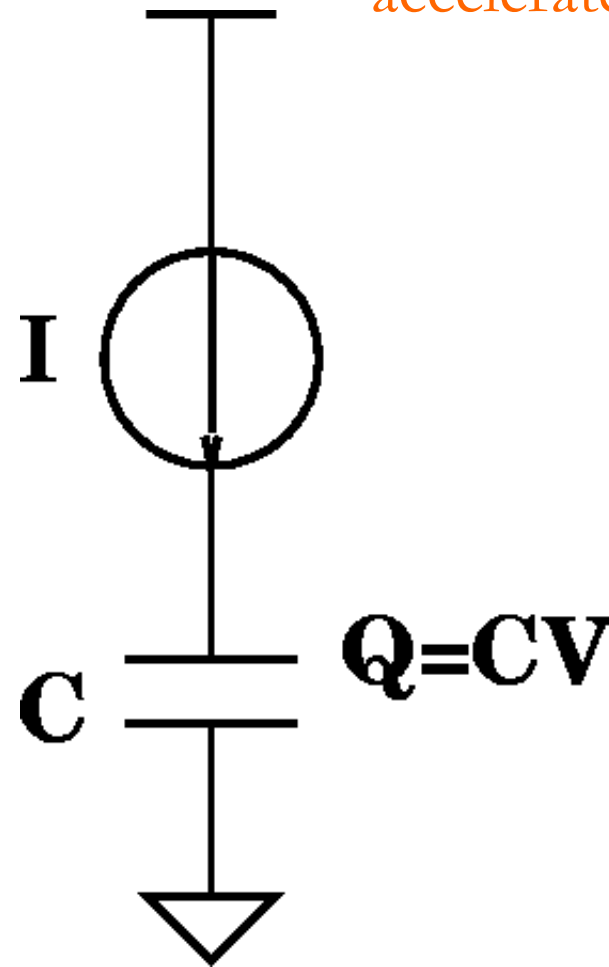
Wire Layers = More Wiring



Gate Delay

- $\tau_{gd} = Q/I = (CV)/I$
- $V' \rightarrow V/S$
- $I'_d \rightarrow I_d/S$
- $C'_g \rightarrow C_g/S$
- $\tau'_{gd} \rightarrow \tau_{gd}/S$

How might we
accelerate speed up?



Gate Delay

- $\tau_{gd} = Q/I = (CV)/I$

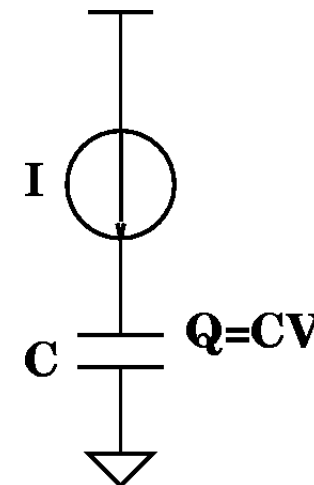
- $V' \rightarrow V$

- $I'_d = (\mu C_{OX} S/2) ((W/S)/(L/S)) (V_{gs} - V_{TH})^2$

- $I'_d \rightarrow I_d \times S$

How might we
accelerate speed up?

Don't scale V!



Gate Delay

- $\tau_{gd} = Q/I = (CV)/I$

- $V' \rightarrow V$

- $I'_d = (\mu C_{OX} S/2) ((W/S)/(L/S)) (V_{gs} - V_{TH})^2$

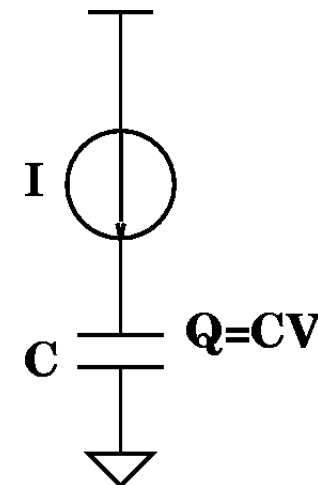
- $I'_d \rightarrow I_d \times S$

- $C'_g \rightarrow C_g/S$

- $\tau'_{gd} \rightarrow \tau_{gd}/S^2$

How might we
accelerate speed up?

Don't scale V!





But... Power Dissipation (Dynamic)

- Capacitive (Dis)charging
 - $P = (1/2)CV^2f$
 - $V' \rightarrow V$
 - $C' \rightarrow C/S$
 - $P' \rightarrow P/S$

But... Power Dissipation (Dynamic)

□ Capacitive (Dis)charging

- $P = (1/2)CV^2f$
- $V' \rightarrow V$
- $C' \rightarrow C/S$
- $P' \rightarrow P/S$

□ Increase Frequency?

- $\tau'_{gd} \rightarrow \tau_{gd}/S^2$
- $f' \rightarrow f \times S^2$
- $P' \rightarrow P \times S$

If don't scale V , power dissipation doesn't scale down!



...And Power Density

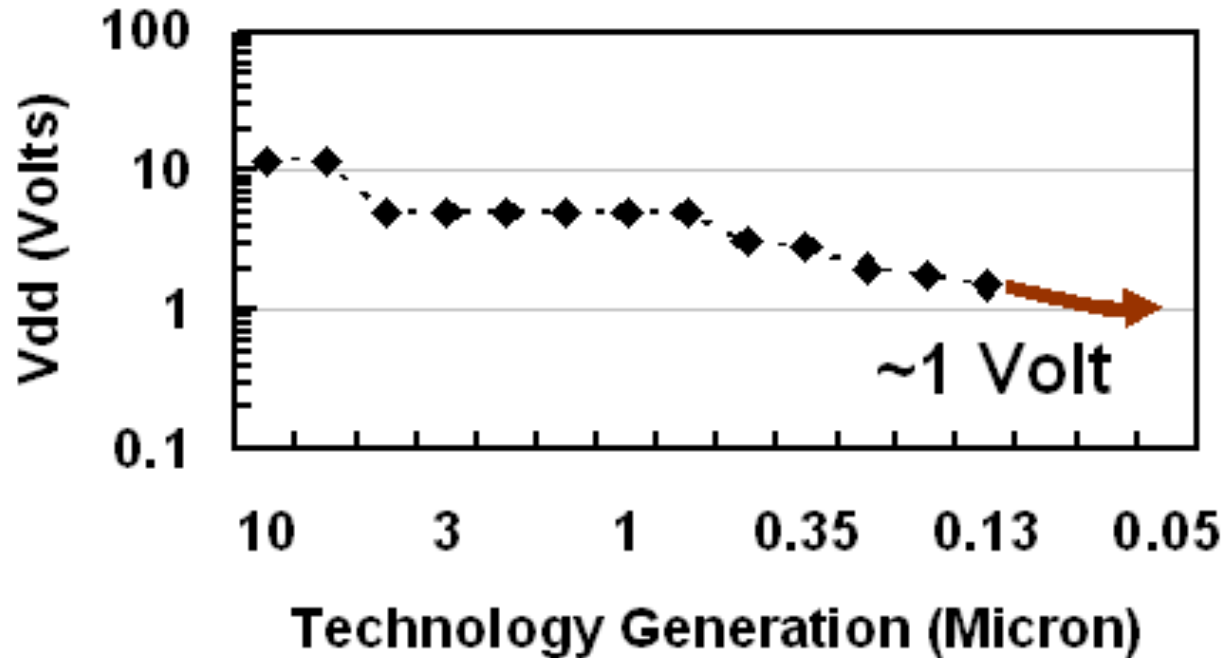
- ❑ $P' \rightarrow P \times S$ (increase frequency)
- ❑ $A' \rightarrow A/S^2$
- ❑ What happens to power density?



...And Power Density

- ❑ $P' \rightarrow P \times S$ (increase frequency)
- ❑ $A' \rightarrow A/S^2$
- ❑ What happens to power density?
- ❑ $P/A \rightarrow S^3 \times P$
- ❑ Power Density Increases!

Historical Voltage Scaling

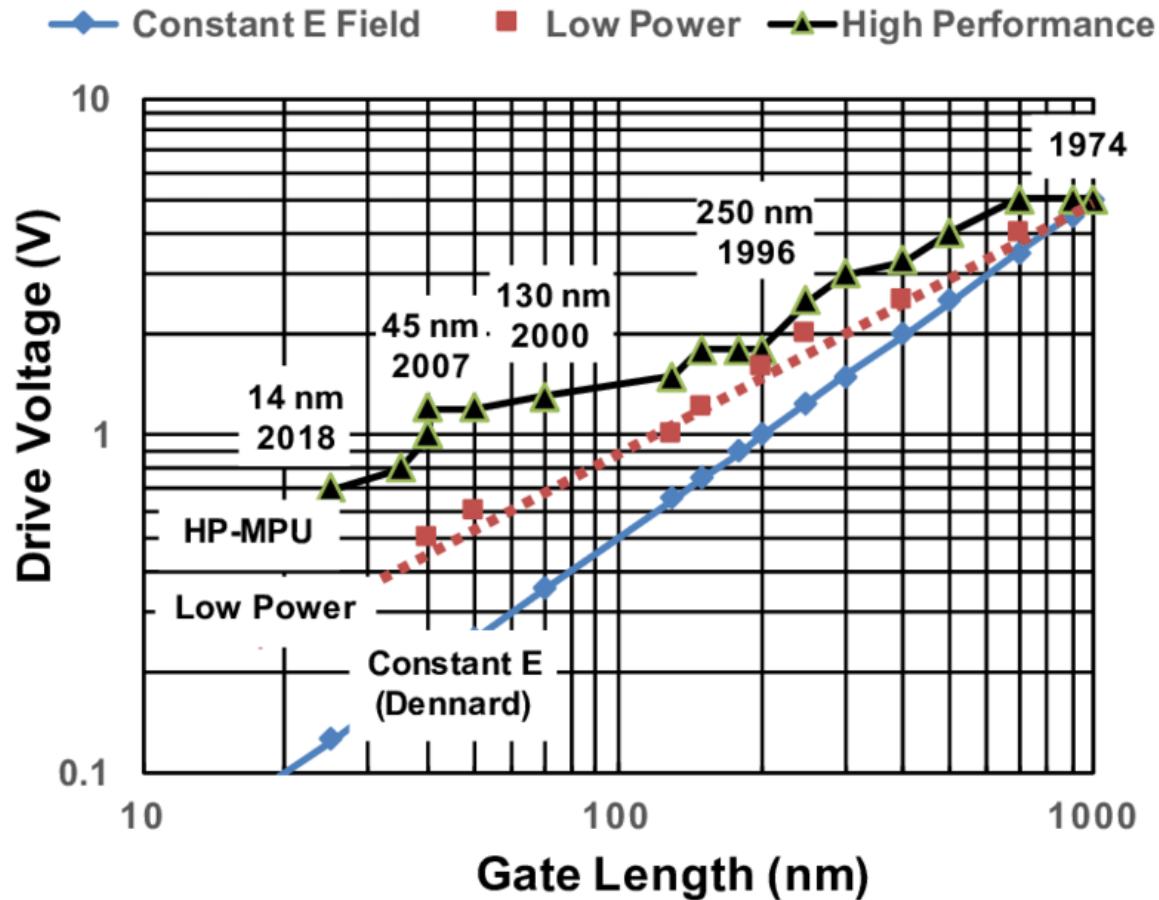


<http://software.intel.com/en-us/articles/gigascale-integration-challenges-and-opportunities/>

- ❑ Frequency impact?
- ❑ Power Density impact?

Historical Voltage Scaling

Drive Voltage Scaling



<http://software>

□ Freq

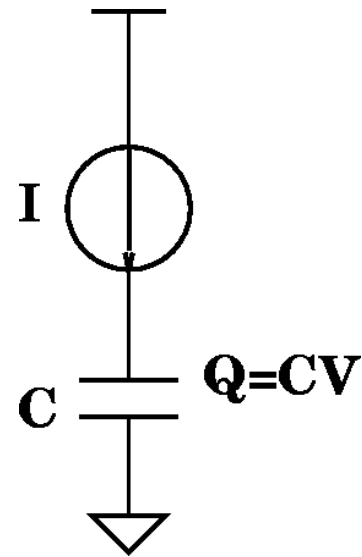
□ Pow

opportunities/

Scale V separately with Factor U , ($U < S$)

□ $\tau_{gd} = Q/I = (CV)/I$

□ $V' \rightarrow V/U$



Scale V separately with Factor U

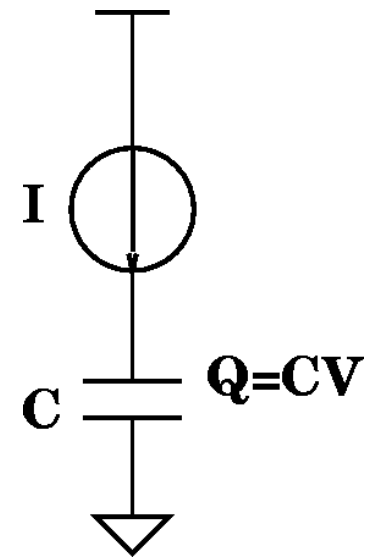
- $\tau_{gd} = Q/I = (CV)/I$

- $V' \rightarrow V/U$

- $I'_d = (\mu C_{OX} S/2) ((W/S)/(L/S)) (V_{gs}/U - V_{TH}/U)^2$

- $I'_d \rightarrow S/U^2 \times I_d$

- $C' \rightarrow C/S$



Scale V separately with Factor U

- $\tau_{gd} = Q/I = (CV)/I$

- $V' \rightarrow V/U$

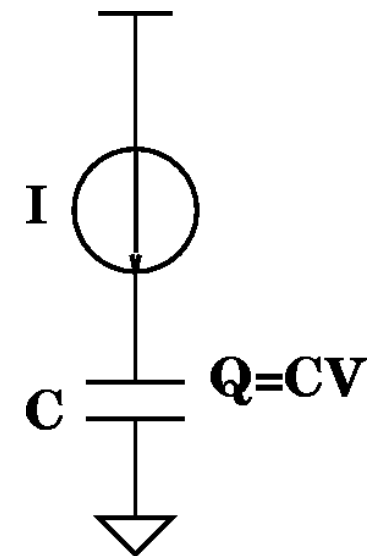
- $I'_d = (\mu C_{OX} S/2) ((W/S)/(L/S) (V_{gs}/U - V_{TH}/U)^2$

- $I'_d \rightarrow S/U^2 \times I_d$

- $C' \rightarrow C/S$

- $\tau'_{gd} \rightarrow ((1/(SU)) / (S/U^2)) \times \tau_{gd}$

- $\tau'_{gd} \rightarrow (U/S^2) \times \tau_{gd}$



Scale V separately with Factor U

$$\tau_{gd} = Q/I = (CV)/I$$

$$V' \rightarrow V/U$$

$$I'_d = (\mu C_{OX} S/2) ((W/S)/(L/S) (V_{gs}/U - V_{TH}/U)^2)$$

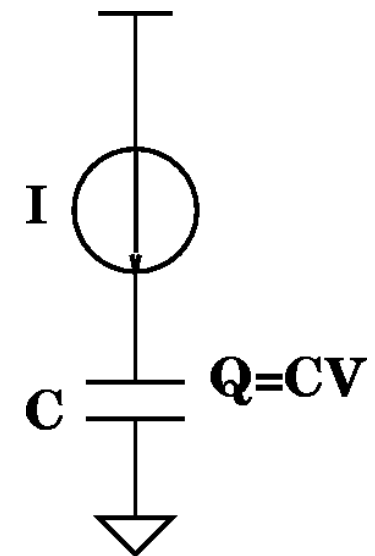
$$I'_d \rightarrow S/U^2 \times I_d$$

$$C' \rightarrow C/S$$

$$\tau'_{gd} \rightarrow ((1/(SU)) / (S/U^2)) \times \tau_{gd}$$

$$\tau'_{gd} \rightarrow (U/S^2) \times \tau_{gd}$$

$$f' \rightarrow (S^2/U) \times f$$



Scale V separately with Factor U

Ideal scale factors:

$$S=100$$

$$U=100$$

$$\tau=1/100$$

$$f_{\text{ideal}}=100$$

$$\tau_{\text{gd}} = Q/I = (CV)/I$$

$$V' \rightarrow V/U$$

$$I'_d = (\mu C_{\text{OX}} S/2) ((W/S)/(L/S) (V_{\text{gs}}/U - V_{\text{TH}}/U)^2$$

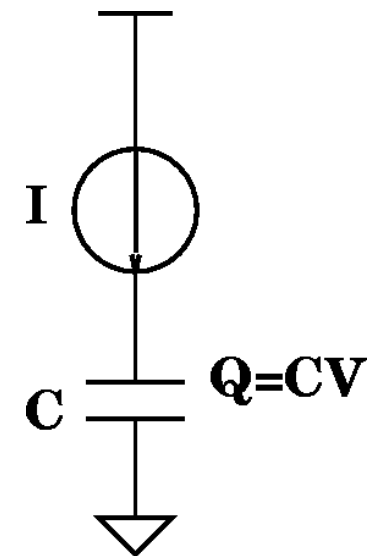
$$I'_d \rightarrow S/U^2 \times I_d$$

$$C' \rightarrow C/S$$

$$\tau'_{\text{gd}} \rightarrow ((1/(SU)) / (S/U^2)) \times \tau_{\text{gd}}$$

$$\tau'_{\text{gd}} \rightarrow (U/S^2) \times \tau_{\text{gd}}$$

$$f' \rightarrow (S^2/U) \times f$$





Question

- Assuming $V_{dd}=10V$ in a $10\mu m$ process and $V_{dd}=1V$ in a $100nm$ process, what are **S** and **U**? (assume everything else scales according to ideal scaling.)

- **S** =

- **U** =

Scale V separately with Factor U

Ideal scale factors:

$$S=100$$

$$U=100$$

$$\tau=1/100$$

$$f_{\text{ideal}}=100$$

$$\tau_{\text{gd}} = Q/I = (CV)/I$$

$$V' \rightarrow V/U$$

$$I'_d = (\mu C_{\text{OX}} S/2) ((W/S)/(L/S) (V_{\text{gs}}/U - V_{\text{TH}}/U)^2$$

$$I'_d \rightarrow S/U^2 \times I_d$$

$$C' \rightarrow C/S$$

$$\tau'_{\text{gd}} \rightarrow \tau_{\text{gd}} \rightarrow \left(\left(\frac{1}{(SU)} \right) / \left(\frac{S}{U^2} \right) \right)$$

$$\tau'_{\text{gd}} \rightarrow (U/S^2) \times \tau_{\text{gd}}$$

$$f' \rightarrow (S^2/U) \times f$$

Cheating factors:

$$S=100$$

$$U=10$$

How much faster are gates?

Scale V separately with Factor U

Ideal scale factors:

$$S=100$$

$$U=100$$

$$\tau=1/100$$

$$f_{\text{ideal}}=100$$

$$\tau_{\text{gd}} = Q/I = (CV)/I$$

$$V' \rightarrow V/U$$

$$I'_d = (\mu C_{\text{OX}} S/2) ((W/S)/(L/S) (V_{\text{gs}}/U - V_{\text{TH}}/U)^2)$$

$$I'_d \rightarrow S/U^2 \times I_d$$

$$C' \rightarrow C/S$$

$$\tau'_{\text{gd}} \rightarrow \tau_{\text{gd}} \rightarrow \left(\left(\frac{1}{(SU)} \right) / \left(\frac{S}{U^2} \right) \right)$$

$$\tau'_{\text{gd}} \rightarrow (U/S^2) \times \tau_{\text{gd}}$$

$$f' \rightarrow (S^2/U) \times f$$

Cheating factors:

$$S=100$$

$$U=10$$

$$\tau=1/1000$$

$$f_{\text{new}}=1000$$

How much faster are gates?

Scale V separately with Factor U

Ideal scale factors:

$$S=100$$

$$U=100$$

$$\tau=1/100$$

$$f_{ideal}=100$$

$$\tau_{gd} = Q/I = (CV)/I$$

$$V' \rightarrow V/U$$

$$I'_d = (\mu C_{OX} S/2) ((W/S)/(L/S) (V_{gs}/U - V_{TH}/U)^2$$

$$I'_d \rightarrow S/U^2 \times I_d$$

$$C' \rightarrow C/S$$

$$\tau'_{gd} \rightarrow ((1/(SU)) / (S/U^2)) \times \tau_{gd}$$

$$\tau'_{gd} \rightarrow (U/S^2) \times \tau_{gd}$$

$$f' \rightarrow (S^2/U) \times f$$

$$f_{new}/f_{ideal} = 10$$

Cheating factors:

$$S=100$$

$$U=10$$

$$\tau=1/1000$$

$$f_{new}=1000$$



Power Density Impact


- $P = (1/2)CV^2 f$
- $P \rightarrow (1/S) (1/U^2) (S^2/U) = S/U^3$
- $P/A \rightarrow (S/U^3) / (1/S^2) = S^3/U^3$



Power Density Impact

- $P = (1/2)CV^2 f$
- $P \rightarrow (1/S) (1/U^2) (S^2/U) = S/U^3$
- $P/A \rightarrow (S/U^3) / (1/S^2) = S^3/U^3$

- $U=10 \quad S=100$
- $P/A \rightarrow 1000 (P/A)$



Power Density Impact

- ❑ $P = (1/2)CV^2 f$
- ❑ $P \rightarrow (1/S) (1/U^2) (S^2/U) = S/U^3$
- ❑ $P/A \rightarrow (S/U^3) / (1/S^2) = S^3/U^3$

- ❑ $U=10 \quad S=100$
- ❑ $P/A \rightarrow 1000 (P/A)$

- ❑ **Compare with ideal scaling:**
- ❑ $P/A \rightarrow S^3 \times P$ (ideal scaling)
- ❑ $P/A \rightarrow 1,000,000 (P/A)$ (ideal scaling)

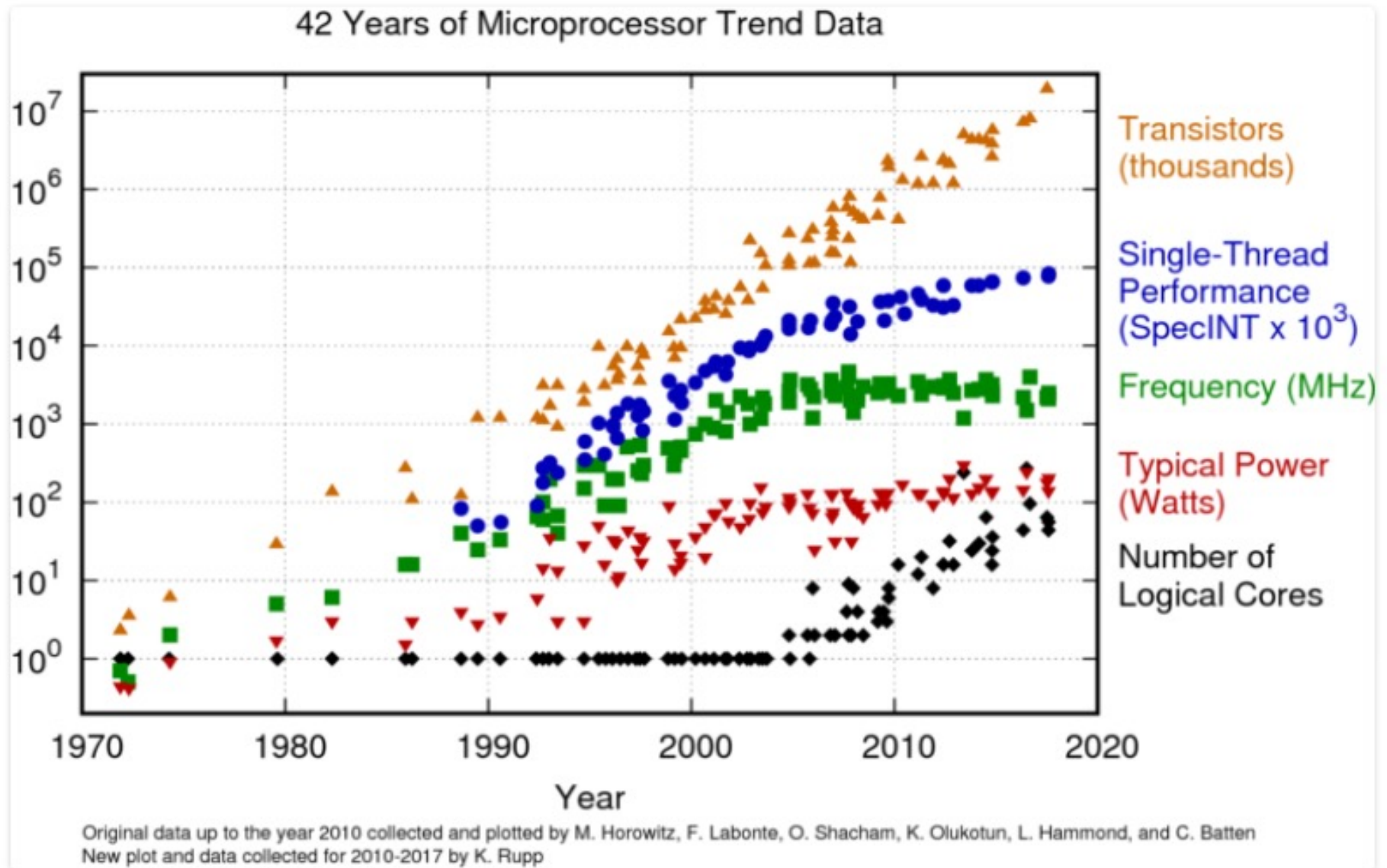
Scaling Methods

Table 3.8 Scaling scenarios for short-channel devices.

Parameter	Relation	Constant Field	General Scaling	Constant Voltage
W, L, t_{ox}		$1/S$	$1/S$	$1/S$
V_{DD}, V_T		$1/S$	$1/U$	1
N_{SUB}	V/W_{depl}^2	S	S^2/U	S^2
Area/Device	WL	$1/S^2$	$1/S^2$	$1/S^2$
C_{ox}	$1/t_{ox}$	S	S	S
C_{gate}	$C_{ox}WL$	$1/S$	$1/S$	$1/S$
k_n, k_p	$C_{ox}W/L$	S	S	S

$U < S$

42 Years of uP Trend Data



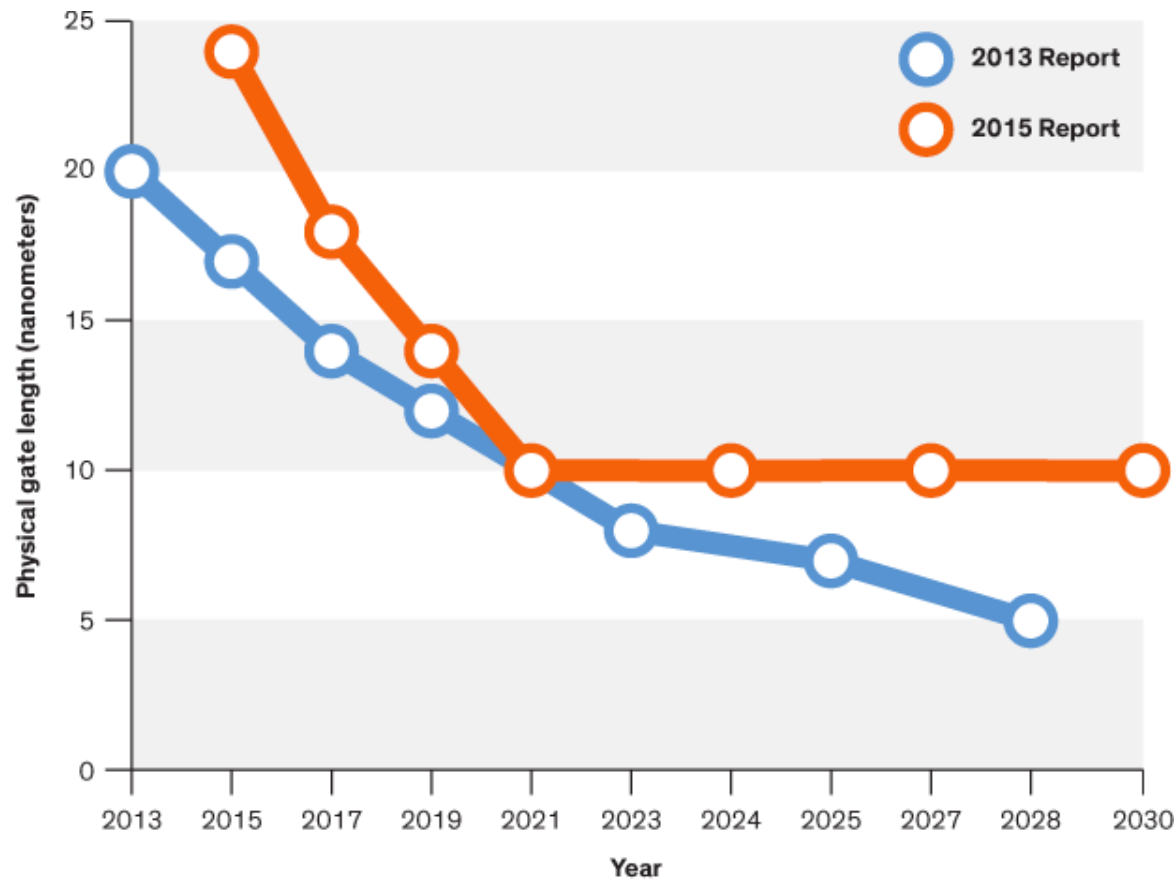


Conventional Scaling

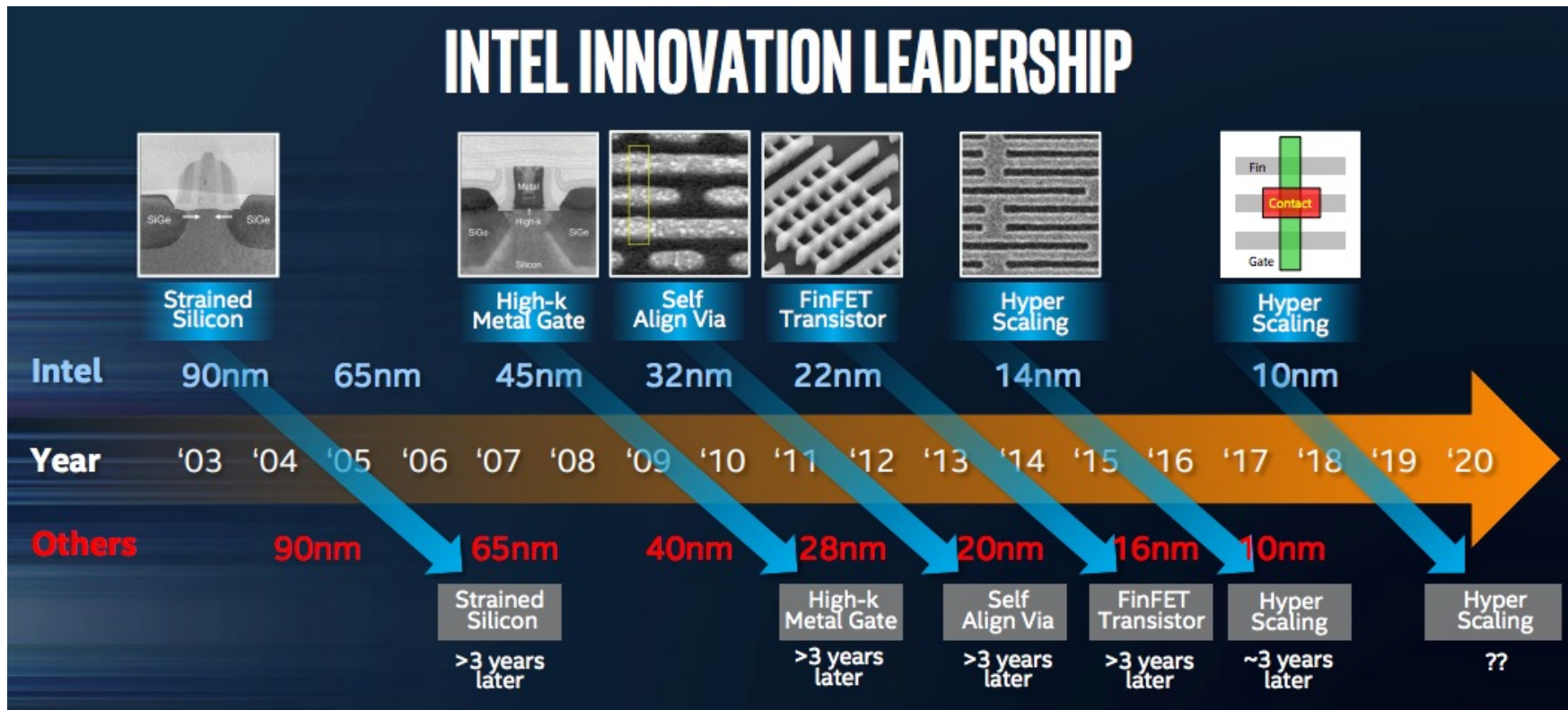
- ❑ Ends in your lifetime
- ❑ Perhaps already:
 - "Basically, this is the end of scaling."
 - May 2005, Bernard Meyerson, V.P. and chief technologist for IBM's systems and technology group

ITRS 2.0 Report 2015

- “After 2021, the report forecasts, it will no longer be economically desirable for companies to continue traditional transistor miniaturization in microprocessors.”

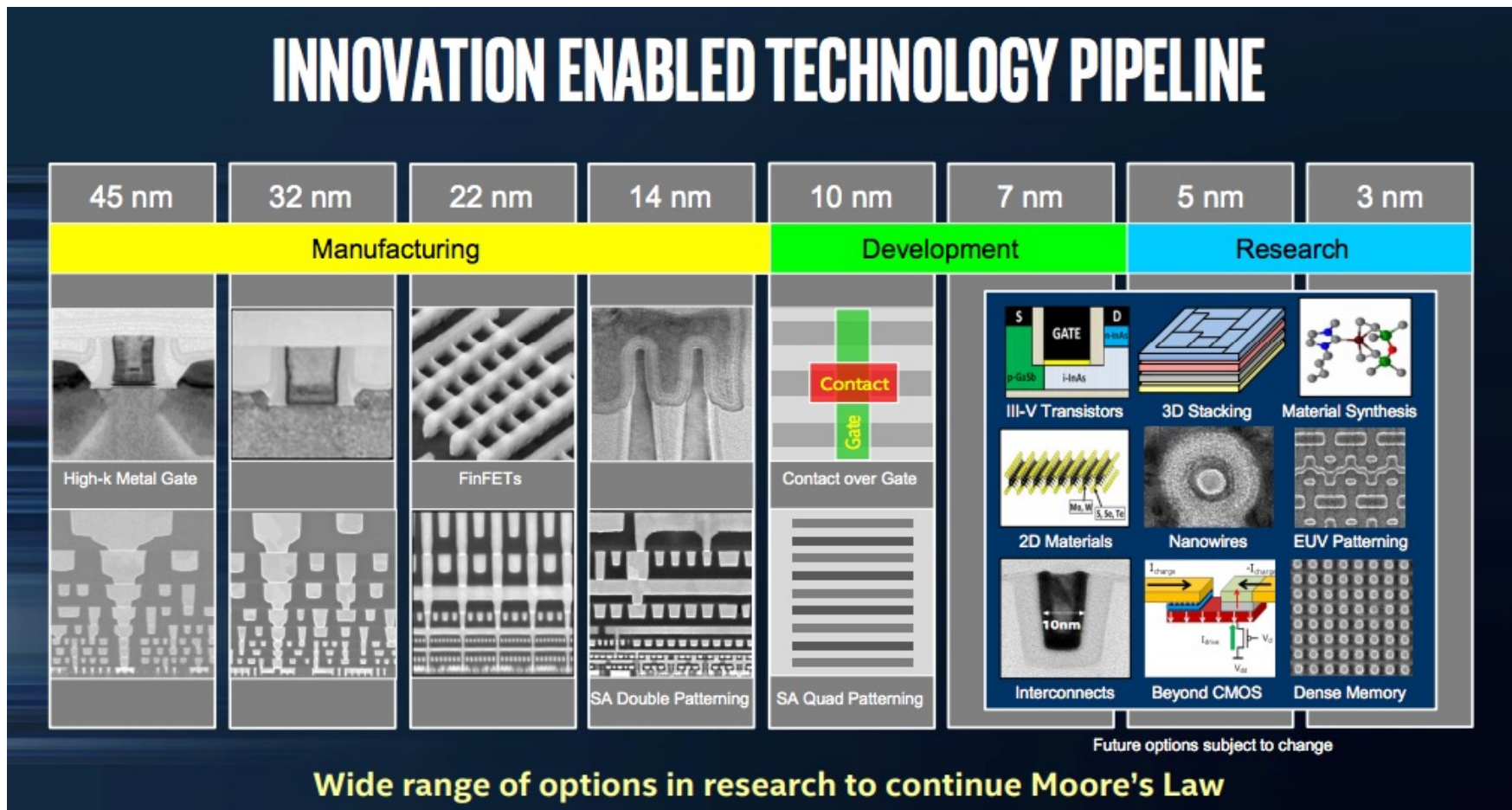


BUT...



Source: <https://newsroom.intel.com/newsroom/wp-content/uploads/sites/11/2017/09/mark-bohr-on-continuing-moores-law.pdf>

BUT...



Source: <https://newsroom.intel.com/newsroom/wp-content/uploads/sites/11/2017/09/mark-bohr-on-continuing-moores-law.pdf>