# Introduction to CMOS VLSI Design

Lecture 21: Scaling and Economics

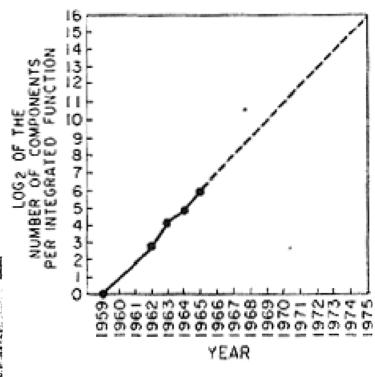
#### Outline

- □ Scaling
  - Transistors
  - Interconnect
  - Future Challenges
- □ VLSI Economics

#### Moore's Law

- ☐ In 1965, Gordon Moore predicted the exponential growth of the number of transistors on an IC
- Transistor count doubled every year since invention
- □ Predicted > 65,000 transistors by 1975!
- ☐ Growth limited by power





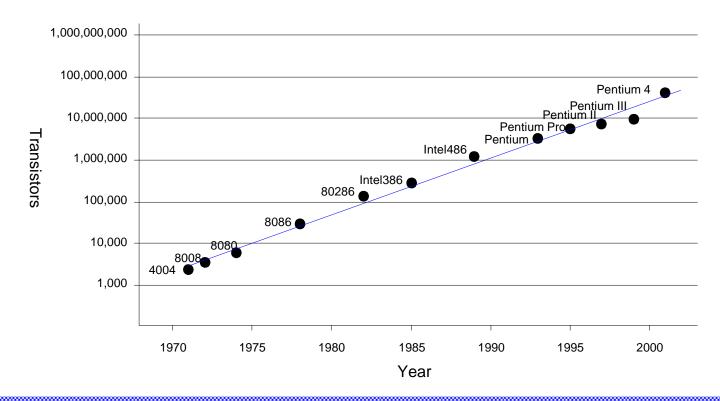
21: Scaling and Economics

**CMOS VLSI Design** 

Slide 3

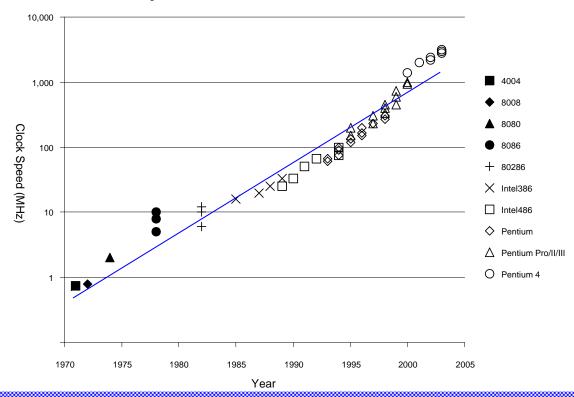
#### More Moore

Transistor counts have doubled every 26 months for the past three decades.



#### Speed Improvement

- ☐ Clock frequencies have also increased exponentially
  - A corollary of Moore's Law



# Why?

☐ Why more transistors per IC?

■ Why faster computers?

# Why?

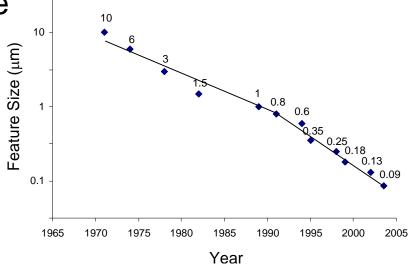
- Why more transistors per IC?
  - Smaller transistors
  - Larger dice
- Why faster computers?

## Why?

- Why more transistors per IC?
  - Smaller transistors
  - Larger dice
- Why faster computers?
  - Smaller, faster transistors
  - Better microarchitecture (more IPC-instruction per cyle)
  - Fewer gate delays per cycle

## Scaling

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



#### Scaling Assumptions

- What changes between technology nodes?
- Constant Field Scaling
  - All dimensions  $(x, y, z \Rightarrow W, L, t_{ox})$
  - Voltage (V<sub>DD</sub>)
  - Doping levels
- Lateral Scaling
  - Only gate length L
  - Often done as a quick gate shrink (S = 1.05)

Table 4.15 Influence of scaling on MOS device characteristics				
Parameter	Sensitivity	Constant Field	Lateral	
Scaling	Parameters			
Length: $L$				
Width: $W$				
Gate oxide thickness: $t_{ox}$				
Supply voltage: $V_{DD}$				
Threshold voltage: $V_{tn}$ , $V_{tp}$				
Substrate doping: $N_{\!\scriptscriptstyle\mathcal{A}}$				
Device Cl	haracteristics			
β				
Current: $I_{ds}$				
Resistance: R				
Gate capacitance: C			•	
Gate delay: τ				
Clock frequency: f			:	
Dynamic power dissipation (per gate): P				
Chip area: $A$				
Power density				
Current density				

Table 4.15 Influence of scaling on MOS device characteristics				
Parameter	Sensitivity	Constant Field	Lateral	
Sc	aling Parameters			
Length: $L$		1/S	1/ <i>S</i>	
Width: W		1/S	1	
Gate oxide thickness: $t_{ m ox}$		1/S	1	
Supply voltage: $V_{DD}$		1/S	1	
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1	
Substrate doping: $N_A$		S	1	
Dev	ice Characteristics			
β				
Current: $I_{ds}$		,	,	
Resistance: R		;	:	
Gate capacitance: C		<del></del>	<del>- 1</del>	
Gate delay: τ		,		
Clock frequency: f		:	!	
Dynamic power dissipation (per gate	e): P			
Chip area: A			'	
Power density				
Current density		,	'	

Table 4.15 Influence of scaling	g on MOS device	characteris	tics
Parameter	Sensitivity	Constant Field	Lateral
Sca	ling Parameters		
Length: $L$		1/S	1/ <i>S</i>
Width: W		1/S	1
Gate oxide thickness: $t_{ m ox}$		1/S	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Devid	e Characteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$			1
Resistance: R			
Gate capacitance: C		1	;
Gate delay: τ			
Clock frequency: f		:	:
Dynamic power dissipation (per gate):	: P		
Chip area: A			'
Power density			
Current density		'	'

Table 4.15 Influence of scaling of	n MOS device o	characteris	tics
Parameter	Sensitivity	Constant Field	Lateral
Scalin	g Parameters		
Length: $L$		1/ <i>S</i>	1/S
Width: W		1/S	1
Gate oxide thickness: $t_{ox}$		1/ <i>S</i>	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device (	Characteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big( V_{DD} - V_t \big)^2$	1/S	S
Resistance: R			
Gate capacitance: C		1	•
Gate delay: τ			
Clock frequency: f		!	:
Dynamic power dissipation (per gate): P			
Chip area: A		1	'
Power density			
Current density		1	'

Table 4.15 Influence of scaling of	n MOS device	characteris	tics
Parameter	Sensitivity	Constant Field	Lateral
Scalin	g Parameters		•
Length: $L$		1/S	1/S
Width: $W$		1/S	1
Gate oxide thickness: $t_{\rm ox}$		1/S	1
Supply voltage: $V_{DD}$		1/ <i>S</i>	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/ <i>S</i>	1
Substrate doping: $N_A$		S	1
Device (	Characteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big( V_{DD} - V_t \big)^2$	1/S	S
Resistance: R	$rac{V_{DD}}{I_{ds}}$	1	1/S
Gate capacitance: C			
Gate delay: τ		-	
Clock frequency: f		!	!
Dynamic power dissipation (per gate): P			
Chip area: A		1	'
Power density			
Current density		'	

Parameter	Sensitivity	Constant	Lateral
		Field	
-	Parameters		
Length: $L$		1/S	1/S
Width: W		1/ <i>S</i>	1
Gate oxide thickness: $t_{ox}$		1/ <i>S</i>	1
Supply voltage: $V_{DD}$		1/ <i>S</i>	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/ <i>S</i>	1
Substrate doping: $N_A$		S	1
Device C	haracteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{di}$	$\beta \big( V_{DD} - V_t \big)^2$	1/S	S
Resistance: R	$rac{{V_{DD}}}{{I_{ds}}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ			'
Clock frequency: f		!	!
Dynamic power dissipation (per gate): P			
Chip area: A		1	'
Power density			
Current density		'	

Parameter	Sensitivity	Constant	Lateral
i didiliotoi	Constitution	Field	Lateral
Scaling	Parameters		·
Length: $L$		1/ <i>S</i>	1/ <i>S</i>
Width: W		1/ <i>S</i>	1
Gate oxide thickness: $t_{ox}$		1/ <i>S</i>	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device C	haracteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big( V_{DD} - V_t \big)^2$	1/S	S
Resistance: R	$rac{V_{DD}}{I_{ds}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ	RC	1/ <i>S</i>	$1/S^{2}$
Clock frequency: f		!	:
Dynamic power dissipation (per gate): P			
Chip area: A		•	
Power density			
Current density			'

Table 4.15 Influence of scaling or	n MOS device o	haracterist	ics
Parameter	Sensitivity	Constant Field	Lateral
Scaling	Parameters	•	'
Length: $L$		1/ <i>S</i>	1/S
Width: W		1/ <i>S</i>	1
Gate oxide thickness: $t_{\rm ox}$		1/ <i>S</i>	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device C	haracteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big(V_{DD}-V_{t}\big)^{2}$	1/S	S
Resistance: R	$rac{V_{DD}}{I_{ds}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ	RC	1/S	$1/S^{2}$
Clock frequency: f	1/τ	S	$S^2$
Dynamic power dissipation (per gate): P		1	
Chip area: A		1	•
Power density			
Current density		<del>'</del>	,

Table 4.15 Influence of scaling or	n MOS device	characteris	tics
Parameter	Sensitivity	Constant Field	Lateral
Scaling	Parameters		
Length: $L$		1/ <i>S</i>	1/S
Width: W		1/S	1
Gate oxide thickness: $t_{\rm ox}$		1/S	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device C	haracteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big(V_{DD}-V_{t}\big)^{2}$	1/S	S
Resistance: R	$rac{V_{DD}}{I_{ds}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ	RC	1/S	$1/S^2$
Clock frequency: f	1/τ	S	$S^2$
Dynamic power dissipation (per gate): P	$CV^2f$	$1/S^2$	S
Chip area: A	·	-	
Power density			
Current density		'	'

Table 4.15 Influence of scaling or	n MOS device o	haracterist	tics
Parameter	Sensitivity	Constant Field	Lateral
Scaling	Parameters		
Length: $L$		1/S	1/ <i>S</i>
Width: W		1/S	1
Gate oxide thickness: $t_{\rm ox}$		1/ <i>S</i>	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device C	haracteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big(V_{DD}-V_{t}\big)^{2}$	1/S	S
Resistance: R	$rac{V_{DD}}{I_{ds}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ	RC	1/S	$1/S^2$
Clock frequency: f	1/τ	S	$S^2$
Dynamic power dissipation (per gate): P	$CV^2f$	$1/S^2$	S
Chip area: A	-	$1/S^2$	1
Power density			
Current density		1	'

Table 4.15 Influence of scaling or			
Parameter	Sensitivity	Constant Field	Lateral
Scaling	Parameters		•
Length: $L$		1/S	1/S
Width: W		1/S	1
Gate oxide thickness: $t_{\rm ox}$		1/S	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device C	haracteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big(V_{DD}-V_{t}\big)^{2}$	1/S	S
Resistance: R	$rac{V_{DD}}{I_{ds}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ	RC	1/S	$1/S^{2}$
Clock frequency: f	1/τ	S	$S^2$
Dynamic power dissipation (per gate): P	$CV^2f$	$1/S^2$	S
Chip area: A		$1/S^2$	1
Power density	P/A	1	S
Current density		-	

Parameter	Sensitivity	Constant Field	Lateral
Scali	ng Parameters		
Length: $L$		1/ <i>S</i>	1/ <i>S</i>
Width: W		1/S	1
Gate oxide thickness: $t_{ m ox}$		1/S	1
Supply voltage: $V_{DD}$		1/S	1
Threshold voltage: $V_{tn}$ , $V_{tp}$		1/S	1
Substrate doping: $N_A$		S	1
Device	Characteristics		
β	$\frac{W}{L} \frac{1}{t_{\text{ox}}}$	S	S
Current: $I_{ds}$	$\beta \big( V_{DD} - V_t \big)^2$	1/S	S
Resistance: R	$rac{{V_{DD}}}{{I_{ds}}}$	1	1/S
Gate capacitance: C	$\frac{WL}{t_{ m ox}}$	1/S	1/S
Gate delay: τ	RC	1/8	$1/S^2$
Clock frequency: f	1/τ	S	$S^2$
Dynamic power dissipation (per gate): I	$CV^2f$	1/52	S
Chip area: A		1/82	1
Power density	P/A	1	S
Current density	$I_d/A$	S	S

#### Observations

- Gate capacitance per micron is nearly independent of process
- But ON resistance \* micron improves with process (W/L increase)
- ☐ Gates get faster with scaling (good)
- Dynamic power goes down with scaling (good)
- ☐ Current density goes up with scaling (bad)
- □ Velocity saturation makes lateral scaling unsustainable

#### Example

- Gate capacitance is typically about 2 fF/μm
- ☐ The FO4 inverter delay in the TT corner for a process of feature size *f* (in nm) is about 0.5*f* ps
- $\Box$  Estimate the ON resistance of a unit (4/2 λ) transistor.

First letter refers to the NMOS corner, and the second letter refers to the PMOS corner. In this naming convention, three corners exist: typical, fast and slow. Fast and slow corners exhibit carrier mobilities that are higher and lower than normal, respectively. For example, a corner designated as FS denotes fast NFETs and slow PFETs. There are therefore five possible corners: typical-typical (TT), fast-fast (FF), slow-slow (SS), fast-slow (FS), and slow-fast (SF).

#### Solution

- Gate capacitance is typically about 2 fF/μm
- ☐ The FO4 inverter delay in the TT corner for a process of feature size *f* (in nm) is about 0.5*f* ps
- $\Box$  Estimate the ON resistance of a unit (4/2  $\lambda$ ) transistor.
- $\Box$  FO4 = 5 τ = 15 RC
- $\square$  RC = (0.5*f*) / 15 = (*f*/30) ps/nm
- $\Box$  If W = 2*f* (double the min. feature size), R = 8.33 k $\Omega$ 
  - Unit resistance is roughly independent of f



#### Scaling Assumptions

- Wire thickness
  - Hold constant vs. reduce in thickness
- □ Wire length
  - Local / scaled interconnect
  - Global interconnect
    - Die size scaled by D<sub>c</sub> ≈ 1.1



Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Scaling I	Parameters		
Width: พ			
Spacing: s		<b>T</b>	
Thickness: t		7	
Interlayer oxide height: <i>h</i>		T	
Characteristics Per Unit Length			
Wire resistance per unit length: $R_{w}$		,	
Fringing capacitance per unit length: $C_{\it wf}$		-	
Parallel plate capacitance per unit length: $C_{\it wp}$		-	+
Total wire capacitance per unit length: $C_w$		1	1
Unrepeated RC constant per unit length: $t_{wu}$			
Repeated wire RC delay per unit length: $t_w$ , (assuming constant field scaling of gates in Table 4.15)			1.
Crosstalk noise		1	+



Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Scaling F	Parameters			
Width: w			1/S	
Spacing: s			1/S	
Thickness: t		1/S	1	
Interlayer oxide height: h			1/ <i>S</i>	
Characteristics Per Unit Length				
Wire resistance per unit length: $R_{\scriptscriptstyle  exttt{w}}$		,		
Fringing capacitance per unit length: $C_{\it wf}$		+	+	
Parallel plate capacitance per unit length: $C_{\it wp}$		+		
Total wire capacitance per unit length: $C_w$		1	1	
Unrepeated RC constant per unit length: $t_{wu}$				
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)				
Crosstalk noise		+	1	



Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Scaling P	arameters	•		
Width: w		1/S		
Spacing: s			1/S	
Thickness: t		1/S	1	
Interlayer oxide height: <i>b</i>			1/ <i>S</i>	
Characteristics Per Unit Length				
Wire resistance per unit length: $R_{\scriptscriptstyle  extbf{w}}$	<u>1</u> wt	S <sup>2</sup>	S	
Fringing capacitance per unit length: $C_{\it wf}$		+		
Parallel plate capacitance per unit length: $C_{w\!p}$		+	+	
Total wire capacitance per unit length: $C_{w}$		1	1	
Unrepeated RC constant per unit length: $t_{ww}$				
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)				
Crosstalk noise		1	+	



Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Scaling Pa	arameters			
Width: ω		1/S		
Spacing: s			1/S	
Thickness: t		1/S	1	
Interlayer oxide height: <i>b</i>		1	1/ <i>S</i>	
Characteristics Per Unit Length				
Wire resistance per unit length: $R_{\scriptscriptstyle  ext{w}}$	<u>1</u> wt	S²	S	
Fringing capacitance per unit length: $C_{w\!f}$	<u>t</u> s	1	S	
Parallel plate capacitance per unit length: $C_{\!\scriptscriptstyle \it wp}$				
Total wire capacitance per unit length: $C_w$		1	1	
Unrepeated RC constant per unit length: $t_{ww}$				
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)				
Crosstalk noise		+	+	



Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Scaling P	arameters			
Width: w		1/ <i>S</i>		
Spacing: s		1	1/S	
Thickness: t		1/S	1	
Interlayer oxide height: <i>b</i>		1	L/S	
Characteristics Per Unit Length				
Wire resistance per unit length: $R_{\scriptscriptstyle  ext{w}}$	$\frac{1}{wt}$	S²	S	
Fringing capacitance per unit length: $C_{w\!f}$	<u>t</u> s	1	S	
Parallel plate capacitance per unit length: $C_{w\!p}$	$\frac{w}{b}$	1	1	
Total wire capacitance per unit length: $C_w$			+	
Unrepeated RC constant per unit length: t <sub>ww</sub>				
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)				
Crosstalk noise			1	



Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Scaling Pa	arameters		
Width: w		1/S	
Spacing: s		:	1/ <i>S</i>
Thickness: t		1/S	1
Interlayer oxide height: <i>h</i>		:	1/ <i>S</i>
Characteristics Per Unit Length			
Wire resistance per unit length: $R_{\scriptscriptstyle  ext{w}}$	$\frac{1}{wt}$	$S^2$	S
Fringing capacitance per unit length: $C_{w\!f}$	<u>t</u> s	1	S
Parallel plate capacitance per unit length: $C_{\ensuremath{\it up}}$	$\frac{w}{b}$	1	1
Total wire capacitance per unit length: $C_w$	$C_{wf}$ + $C_{wp}$	1	between 1,
Unrepeated RC constant per unit length: $t_{wu}$			
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)			
Crosstalk noise		1	1



Parameter	Sensitivity	Reduced	Constant
		Thickness	Thickness
Scaling Pa	arameters		
Width: w		1/S	
Spacing: s		_	/S
Thickness: t		1/S	1
Interlayer oxide height: $b$		1	/S
Characteristics Per Unit Length			
Wire resistance per unit length: $R_w$	$\frac{1}{wt}$	$S^2$	S
Fringing capacitance per unit length: $C_{\it wf}$	<u>t</u> s	1	S
Parallel plate capacitance per unit length: $C_{w\!p}$	$\frac{w}{h}$	1	1
Total wire capacitance per unit length: $C_w$	$C_{wf} + C_{wp}$	1	between 1, S
Unrepeated RC constant per unit length: $t_{ww}$	$R_wC_w$	$S^2$	between S, S <sup>2</sup>
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)			
Crosstalk noise		1	i



Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Scaling Pa	arameters		
Width: $w$		1/S	
Spacing: s		1/S	
Thickness: t		1/S	1
Interlayer oxide height: h		1	/S
Characteristics Per Unit Length			
Wire resistance per unit length: $R_w$	$\frac{1}{wt}$	$S^2$	S
Fringing capacitance per unit length: $C_{w\!f}$	<u>t</u> s	1	S
Parallel plate capacitance per unit length: $C_{w\!p}$	$\frac{w}{b}$	1	1
Total wire capacitance per unit length: $C_{w}$	$C_{wf}$ + $C_{wp}$	1	between 1, S
Unrepeated RC constant per unit length: $t_{wu}$	$R_wC_w$	S <sup>2</sup>	between $S$ , $S^2$
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)	$\sqrt{RCR_wC_w}$	<b>√</b> S	between 1, $\sqrt{S}$
Crosstalk noise			



Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Scaling P	arameters		
Width: w		1/S	
Spacing: s			1/S
Thickness: t		1/S	1
Interlayer oxide height: $b$			1/ <i>S</i>
Characteristics Per Unit Length			
Wire resistance per unit length: $R_{\scriptscriptstyle  ext{w}}$	$\frac{1}{wt}$	S²	S
Fringing capacitance per unit length: $C_{\it wf}$	<u>t</u> s	1	S
Parallel plate capacitance per unit length: $C_{w\!p}$	$\frac{w}{b}$	1	1
Total wire capacitance per unit length: $C_w$	$C_{wf}$ + $C_{wp}$	1	between 1,
Unrepeated RC constant per unit length: $t_{ww}$	$R_w C_w$	$S^2$	between S, S <sup>2</sup>
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)	$\sqrt{RCR_wC_w}$	$\sqrt{S}$	between 1, $\sqrt{S}$
Crosstalk noise	<u>t</u> s	1	S



# Interconnect Delay

Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Sca	ling Parameters		
Width: w			1/ <i>S</i>
Spacing: s			1/ <i>S</i>
Thickness: t		1/S	1
Interlayer oxide height: h		1/S	
Local/Scaled Interconnect Characterist	tics	+	-
Length: /			
Unrepeated wire RC delay			'
Repeated wire delay			
Global Interconnect Characteristics			
Length: /			
Unrepeated wire RC delay		·	
Repeated wire delay		<del>-                                    </del>	



Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Sca	aling Parameters		
Width: w			1/S
Spacing: s			1/S
Thickness: t		1/S	1
Interlayer oxide height: $h$			1/ <i>S</i>
Local/Scaled Interconnect Characteris	tics		
Length: /		1/S	
Unrepeated wire RC delay		1	1
Repeated wire delay			
Global Interconnect Characteristics		+	
Length: /			
Unrepeated wire RC delay			
Repeated wire delay		+	+



Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Sca	aling Parameters		-	
Width: w			1/ <i>S</i>	
Spacing: s			1/ <i>S</i>	
Thickness: t		1/S	1	
Interlayer oxide height: $h$			1/S	
Local/Scaled Interconnect Characteris	tics	•	•	
Length: /		1/S		
Unrepeated wire RC delay	$l^2 t_{wu}$	1	between 1/S, 1	
Repeated wire delay		:		
Global Interconnect Characteristics			-	
Length: /				
Unrepeated wire RC delay				
Repeated wire delay		<del>'</del>	1	



Parameter	Sensitivity	Reduced Thickness	Constant Thickness
Sca	ling Parameters	THICKIESS	THICKIESS
Width: ω			1/ <i>S</i>
Spacing: s			1/ <i>S</i>
Thickness: t		1/S	1
Interlayer oxide height: h			1/ <i>S</i>
Local/Scaled Interconnect Characterist	ics		
Length: /			1/ <i>S</i>
Unrepeated wire RC delay	$l^2 t_{wu}$	1	between 1/S, 1
Repeated wire delay	lt <sub>wr</sub>	$\sqrt{1/S}$	between $1/S$ , $\sqrt{1/S}$
Global Interconnect Characteristics			
Length: /			
Unrepeated wire RC delay			
Repeated wire delay		1	1



	g on interconnec			
Parameter	Sensitivity	Reduced Thickness	Constant Thickness	
Sca	ling Parameters	THICKHESS	THICKITESS	
Width: w			1/S	
Spacing: s			1/S	
Thickness: t		1/S	1	
Interlayer oxide height: h			1/ <i>S</i>	
Local/Scaled Interconnect Characterist	ics			
Length: /			1/ <i>S</i>	
Unrepeated wire RC delay	$l^2t_{wu}$	1	between $1/S$ , 1	
Repeated wire delay	lt <sub>wr</sub>	$\sqrt{1/S}$ between $1/S$ , $\sqrt{1}$		
Global Interconnect Characteristics				
Length: /		$D_{\epsilon}$		
Unrepeated wire RC delay		•		
Repeated wire delay		-		



Parameter	Sensitivity	Reduced Thickness	Constant Thickness		
Sca	aling Parameters		•		
Width: ω			1/S		
Spacing: s			1/S		
Thickness: t		1/S	1		
Interlayer oxide height: $h$			1/ <i>S</i>		
Local/Scaled Interconnect Characteris	tics				
Length: /			1/S		
Unrepeated wire RC delay	$l^2 t_{wu}$	1	between 1/S, 1		
Repeated wire delay	$lt_{wr}$	$\sqrt{1/S}$ between 1/S, $\sqrt{1}$			
Global Interconnect Characteristics					
Length: /		$D_{\epsilon}$			
Unrepeated wire RC delay	$l^2 t_{wu}$	$S^2D_c^2$	between $SD_c^2$ , $S^2D_c^2$		
Repeated wire delay		<del>-  </del>			



Parameter	Sensitivity	Reduced Thickness	Constant Thickness		
S	caling Parameters				
Width: w			1/ <i>S</i>		
Spacing: s			1/ <i>S</i>		
Thickness: t		1/S	1		
Interlayer oxide height: $h$			1/S		
Local/Scaled Interconnect Characteri	stics		-		
Length: /			1/S		
Unrepeated wire RC delay	$l^2t_{wu}$	1	between 1/S, 1		
Repeated wire delay	lt <sub>wr</sub>	$\sqrt{1/S}$	between $1/S$ , $\sqrt{1/S}$		
Global Interconnect Characteristics					
Length: /		$D_{\epsilon}$			
Unrepeated wire RC delay	$l^2t_{wu}$	$S^2D_c^2$	between $SD_c^2$ , $S^2D_c^2$		
Repeated wire delay	lt <sub>wr</sub>	$D_c \sqrt{S}$	between $D$		

#### Observations

- □ Capacitance per micron is remaining constant
  - About 0.2 fF/μm
  - Roughly 1/10 of gate capacitance
- ☐ Local wires are getting faster
  - Not quite tracking transistor improvement
  - But not a major problem
- ☐ Global wires are getting slower
  - No longer possible to cross chip in one cycle

### **ITRS**

- □ Semiconductor Industry Association forecast
  - Intl. Technology Roadmap for Semiconductors

Table 4.17 Predictions from the 2002 ITRS							
Year	2001	2004	2007	2010	2013	2016	
Feature size (nm)	130	90	65	45	32	22	
$V_{DD}\left(\mathbf{V}\right)$	1.1-1.2	1-1.2	0.7-1.1	0.6-1.0	0.5-0.9	0.4-0.9	
Millions of transistors/die	193	385	773	1564	3092	6184	
Wiring levels	8-10	9–13	10-14	10-14	11-15	11-15	
Intermediate wire pitch (nm)	450	275	195	135	95	65	
Interconnect dielectric	3-3.6	2.6-3.1	2.3-2.7	2.1	1.9	1.8	
constant							
I/O signals	1024	1024	1024	1280	1408	1472	
Clock rate (MHz)	1684	3990	6739	11511	19348	28751	
FO4 delays/cycle	13.7	8.4	6.8	5.8	4.8	4.7	
Maximum power (W)	130	160	190	218	251	288	
DRAM capacity (Gbits)	0.5	1	4	8	32	64	

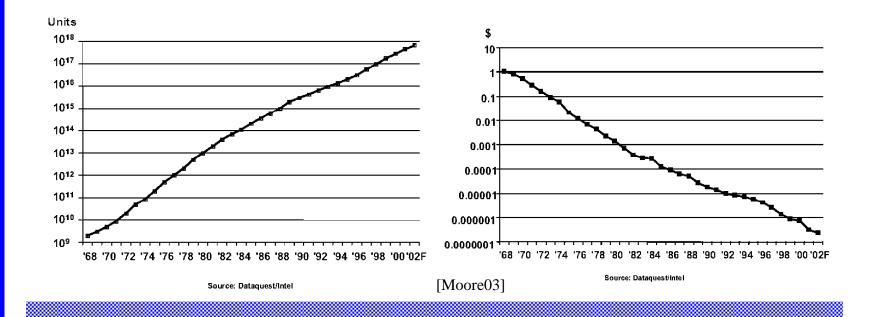
## Scaling Implications

- ☐ Improved Performance
- ☐ Improved Cost
- ☐ Interconnect Woes
- Power Woes
- Productivity Challenges
- Physical Limits

## Cost Improvement

- ☐ In 2003, \$0.01 bought you 100,000 transistors
  - Moore's Law is still going strong

21: Scaling and Economics

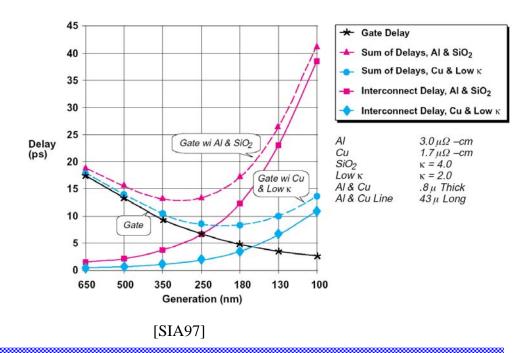


**CMOS VLSI Design** 

Slide 46

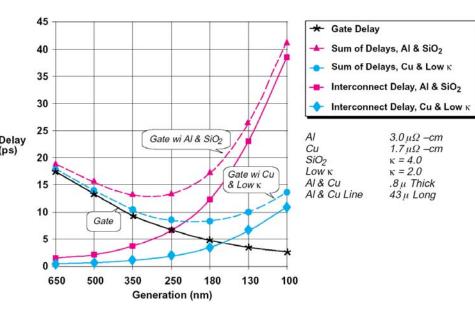
### Interconnect Woes

- ☐ SIA made a gloomy forecast in 1997
  - Delay would reach minimum at 250 180 nm,
     then get worse because of wires
- □ But...



#### Interconnect Woes

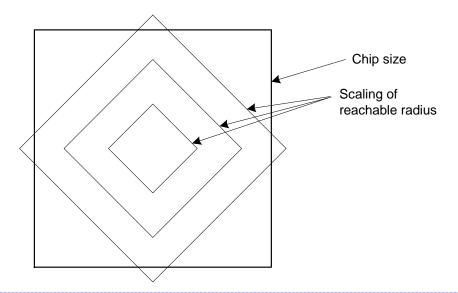
- ☐ SIA made a gloomy forecast in 1997
  - Delay would reach minimum at 250 180 nm,
     then get worse because of wires
- **□** But...
  - Misleading scale
  - Global wires
- ☐ 100 kgate blocks ok (ps)





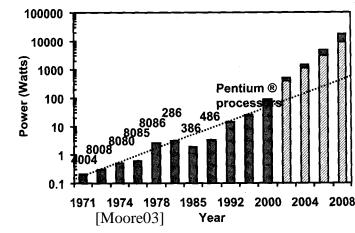
#### Reachable Radius

- We can't send a signal across a large fast chip in one cycle anymore
- ☐ But the microarchitect can plan around this
  - Just as off-chip memory latencies were tolerated



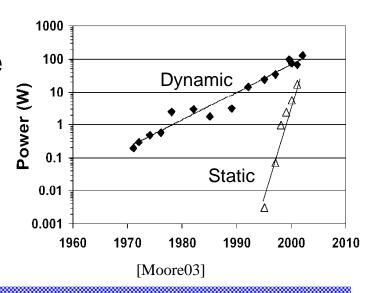
### Dynamic Power

- □ Intel VP Patrick Gelsinger (ISSCC 2001)
  - If scaling continues at present pace, by 2005, high speed processors would have power density of nuclear reactor, by 2010, a rocket nozzle, and by 2015, surface of sun.
  - "Business as usual will not work in the future."
- ☐ Intel stock dropped 8% on the next day
- But attention to power is increasing



#### Static Power

- □ V<sub>DD</sub> decreases
  - Save dynamic power
  - Protect thin gate oxides and short channels
  - No point in high value because of velocity sat.
- □ V<sub>t</sub> must decrease to maintain device performance
- ☐ But this causes exponential increase in OFF leakage
- Major future challenge



### Productivity

- Transistor count is increasing faster than designer productivity (gates / week)
  - Bigger design teams
    - Up to 500 for a high-end microprocessor
  - More expensive design cost
  - Pressure to raise productivity
    - Rely on synthesis, IP blocks
  - Need for good engineering managers

### Physical Limits

- ☐ Will Moore's Law run out of steam?
  - Can't build transistors smaller than an atom...
- Many reasons have been predicted for end of scaling
  - Dynamic power
  - Subthreshold leakage, tunneling
  - Short channel effects
  - Fabrication costs
  - Electromigration
  - Interconnect delay
- Rumors of demise have been exaggerated

#### VLSI Economics

- ☐ Selling price S<sub>total</sub>
  - $-S_{total} = C_{total} / (1-m)$
- $\Box$  m = profit margin
- $\Box$  C<sub>total</sub> = total cost
  - Nonrecurring engineering cost (NRE)
  - Recurring cost
  - Fixed cost

#### NRE

- ☐ Engineering cost
  - Depends on size of design team
  - Include benefits, training, computers
  - CAD tools:
    - Digital front end: \$10K
    - Analog front end: \$100K
    - Digital back end: \$1M
- Prototype manufacturing
  - Mask costs: \$500k 1M in 130 nm process
  - Test fixture and package tooling

### Recurring Costs

- Fabrication
  - Wafer cost / (Dice per wafer \* Yield)
  - Wafer cost: \$500 \$3000
  - Dice per wafer:  $N = \pi \left[ \frac{r^2}{A} \frac{2r}{\sqrt{2A}} \right]$
  - Yield:  $Y = e^{-AD}$ 
    - For small A, Y ≈ 1, cost proportional to area
    - For large A,  $Y \rightarrow 0$ , cost increases exponentially
- Packaging
- □ Test

### Fixed Costs

- □ Data sheets and application notes
- Marketing and advertising
- ☐ Yield analysis

### Example

- ☐ You want to start a company to build a wireless communications chip. How much venture capital must you raise?
- □ Because you are smarter than everyone else, you can get away with a small team in just two years:
  - Seven digital designers
  - Three analog designers
  - Five support personnel

### Solution

- Digital designers:
  - salary
  - overhead
  - computer
  - CAD tools
  - Total:
- Analog designers
  - salary
  - overhead
  - computer
  - CAD tools
  - Total:

- ☐ Support staff
  - salary
  - overhead
  - computer
  - Total:
- Fabrication
  - Back-end tools:
  - Masks:
  - Total:
- Summary

### Solution

- Digital designers:
  - \$70k salary
  - \$30k overhead
  - \$10k computer
  - \$10k CAD tools
  - Total: \$120k \* 7 = \$840k ☐ Fabrication
- Analog designers
  - \$100k salary
  - \$30k overhead
  - \$10k computer
  - \$100k CAD tools
  - Total: \$240k \* 3 = \$720k

- Support staff
  - \$45k salary
  - \$20k overhead
  - \$5k computer
  - Total: \$70k \* 5 = \$350k
- - Back-end tools: \$1M
  - Masks: \$1M
  - Total: \$2M / year
- Summary
  - 2 years @ \$3.91M / year
  - + \$8M design & prototype

### Cost Breakdown

- New chip design is fairly capital-intensive
- Maybe you can do it for less?

