

CS184a: Computer Architecture (Structure and Organization)

Day 6: January 22, 2003
VLSI Scaling



Caltech CS184 Winter2003 -- DeHon

Today

- VLSI Scaling Rules
- Effects
- Historical/predicted scaling
- Variations (cheating)
- Limits

Caltech CS184 Winter2003 -- DeHon

Why Care?

- In this game, we must be able to predict the future
- Rapid technology advance
- Reason about changes and trends
- re-evaluate prior solutions given technology at time X.

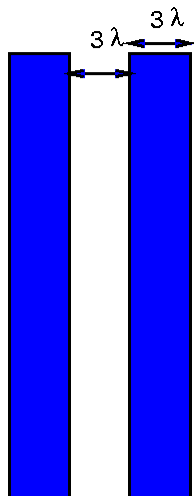
Why Care

- Cannot compare against what competitor does today
 - but what they can do at time you can ship
- Careful not to fall off curve
 - lose out to someone who can stay on curve

Scaling

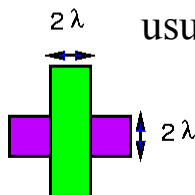
- **Premise:** features scale “uniformly”
 - everything gets better in a predictable manner
- **Parameters:**
 - λ (lambda) -- Mead and Conway (class)
 - S -- Bohr
 - $1/\kappa$ -- Dennard

Feature Size



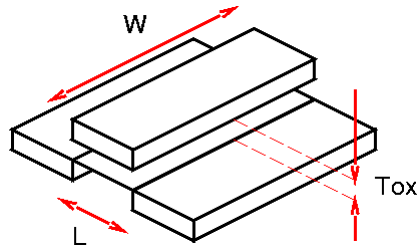
λ is half the minimum feature size in a VLSI process

[minimum feature usually channel width]



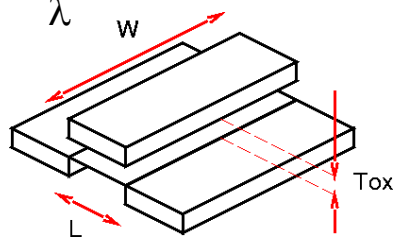
Scaling

- Channel Length (L)
- Channel Width (W)
- Oxide Thickness (T_{ox})
- Doping (N_a)
- Voltage (V)



Scaling

- Channel Length (L) λ
- Channel Width (W) λ
- Oxide Thickness (T_{ox}) λ
- Doping (N_a) $1/\lambda$
- Voltage (V) λ



Effects?

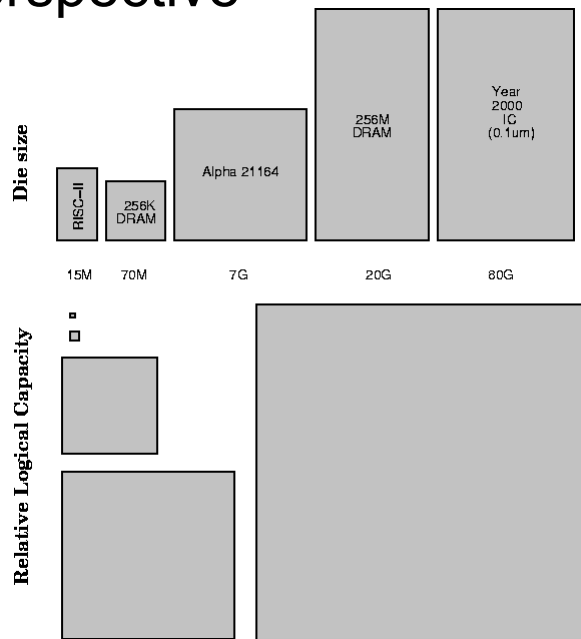
- Area
- Capacitance
- Resistance
- Threshold (V_{th})
- Current (I_d)
- Gate Delay (τ_{gd})
- Wire Delay (τ_{wire})
- Power

Area

- $\lambda \rightarrow \lambda/\kappa$
- $A = L * W$
- $A \rightarrow A/\kappa^2$
- $0.35\mu m \rightarrow 0.25\mu m$
- 50% area
- 2x capacity same area

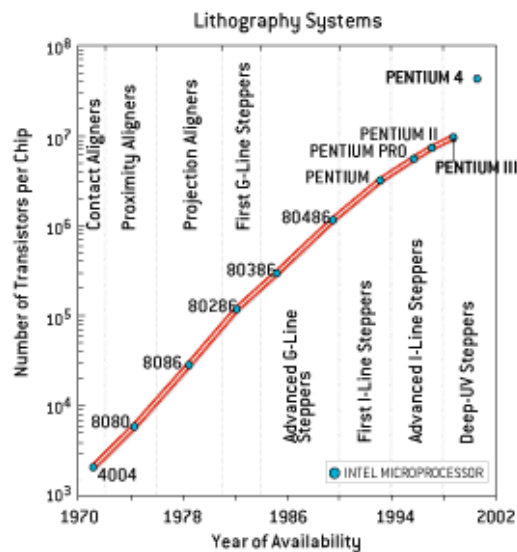
Area Perspective

[2000 tech.]
18mm×18mm
0.18μm
60G λ²



Caltech CS184 Winter2003 -- DeHon

Capacity Scaling from Intel



SOURCES: VLSI Research, Inc.; Integrated Circuit Engineering Corporation; Intel

Caltech CS184 Winter2003 -- DeHon

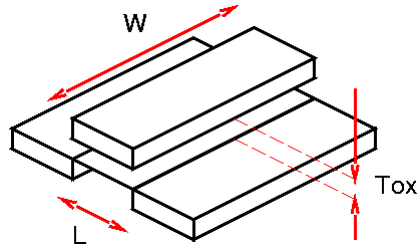
Capacitance

- Capacitance per unit area

$$- C_{ox} = \epsilon_{SiO_2} / T_{ox}$$

$$- T_{ox} \rightarrow T_{ox} / \kappa$$

$$- C_{ox} \rightarrow \kappa C_{ox}$$



Capacitance

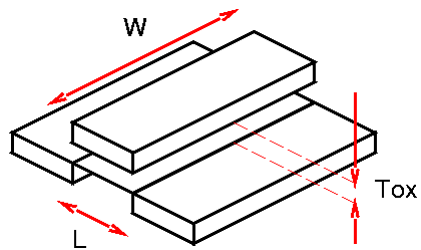
- Gate Capacitance

$$\blacksquare C_{gate} = A * C_{ox}$$

$$\blacksquare A \rightarrow A / \kappa^2$$

$$\blacksquare C_{ox} \rightarrow \kappa C_{ox}$$

$$\blacksquare C_{gate} \rightarrow C_{gate} / \kappa$$



Threshold Voltage

Before:

$$V_{th} = \frac{1}{C_{OX}} \left(-Q_{eff} + \left(2\epsilon_{Si} q N_a (\phi_s + V_{s-sub}) \right)^{1/2} \right) + (W_f + \phi_s)$$

$$(W_f + \phi_s) \approx 0$$

adjust V_{s-sub} **so** $(\phi_s + V_{s-sub}) \rightarrow \frac{(\phi_s + V_{s-sub})}{\kappa}$

After:

$$V'_{th} = \frac{1}{\kappa C_{OX}} \left(-Q_{eff} + \left(2\epsilon_{Si} q \kappa N_a \frac{(\phi_s + V_{s-sub})}{\kappa} \right)^{1/2} \right)$$

$$V'_{th} \approx \frac{V_{th}}{\kappa}$$

Caltech CS184 Winter2003 -- DeHon

Threshold Voltage

- $V_{TH} \rightarrow V_{TH} / \kappa$

Caltech CS184 Winter2003 -- DeHon

Current

- Saturation Current

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

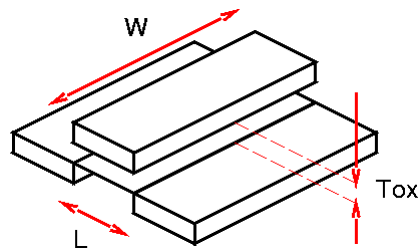
- $V_{gs} = V \rightarrow V/\kappa$

- $V_{TH} \rightarrow V_{TH}/\kappa$

- $W \rightarrow W/\kappa$

- $C_{OX} \rightarrow \kappa C_{OX}$

- $I_d \rightarrow I_d/\kappa$



Gate Delay

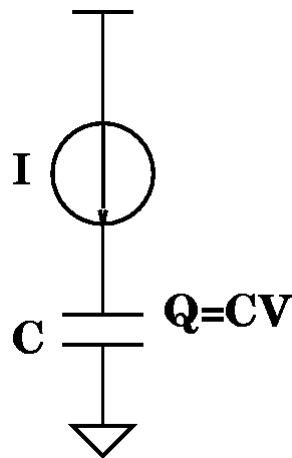
- $\tau_{gd} = Q/I = (CV)/I$

- $V \rightarrow V/\kappa$

- $I_d \rightarrow I_d/\kappa$

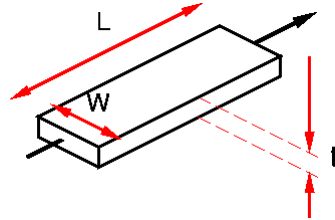
- $C \rightarrow C/\kappa$

- $\tau_{gd} \rightarrow \tau_{gd}/\kappa$



Resistance

- $R = \rho L / (W \cdot t)$
- $W \rightarrow W / \kappa$
- L, t similar
- $R \rightarrow \kappa R$



Wire Delay

- $\tau_{\text{wire}} = R \times C$
- $R \rightarrow \kappa R$
- $C \rightarrow C / \kappa$
- $\tau_{\text{wire}} \rightarrow \tau_{\text{wire}}$
- ...assuming (logical) wire lengths remain constant...
- Assume short wire or buffered wire
- (unbuffered wire ultimately scales as length squared)

Power Dissipation (Static)

- Resistive Power

$$- P = V \cdot I$$

$$- V \rightarrow V / \kappa$$

$$- I_d \rightarrow I_d / \kappa$$

$$- P \rightarrow P / \kappa^2$$

Power Dissipation (Dynamic)

- Capacitive (Dis)charging
 - $P = (1/2)CV^2f$

$$▪ V \rightarrow V / \kappa$$

$$▪ C \rightarrow C / \kappa$$

$$▪ P \rightarrow P / \kappa^3$$

- Increase Frequency?

$$▪ f \rightarrow \kappa f ?$$

$$▪ P \rightarrow P / \kappa^2$$

Effects?

- Area $1/\kappa^2$
- Capacitance $1/\kappa$
- Resistance κ
- Threshold (V_{th}) $1/\kappa$
- Current (I_d) $1/\kappa$
- Gate Delay (τ_{gd}) $1/\kappa$
- Wire Delay (τ_{wire}) 1
- Power $1/\kappa^2 \rightarrow 1/\kappa^3$

ITRS Roadmap

- Semiconductor Industry rides this scaling curve
- Try to predict where industry going
 - (requirements...self fulfilling prophecy)
- <http://public.itrs.net>

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

$$S=0.7$$

[0.5x per 2 nodes]

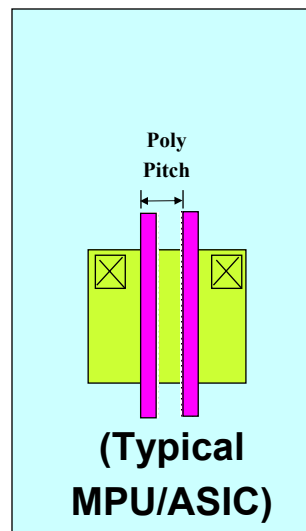
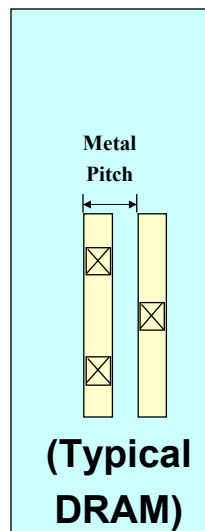


Source: 2001 ITRS - Exec. Summary, ORTC Figure
Caltech CS184 Winter2003 -- DeHon

[from Andrew Kahng]

25

Half Pitch (= Pitch/2) Definition

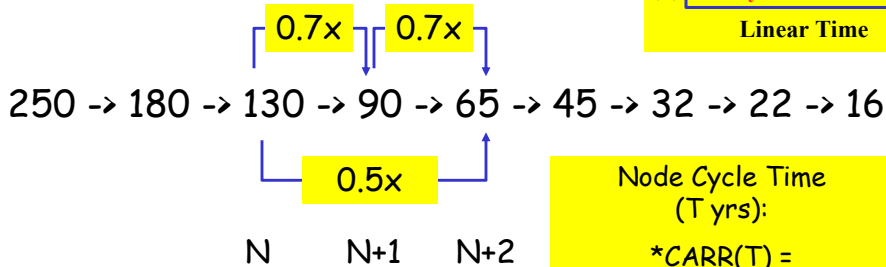


Source: 2001 ITRS - Exec. Summary, ORTC Figure
Caltech CS184 Winter2003 -- DeHon

[from Andrew Kahng]

26

Scaling Calculator + Node Cycle Time:



* CARR(T) = Compound Annual
Reduction Rate
(@ cycle time period, T)

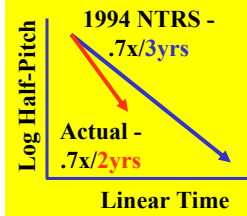
Node Cycle Time
(T yrs):

$$*CARR(T) =$$

$$[(0.5)^{(1/2T \text{ yrs})}] - 1$$

$$CARR(3 \text{ yrs}) = -10.9\%$$

$$CARR(2 \text{ yrs}) = -15.9\%$$

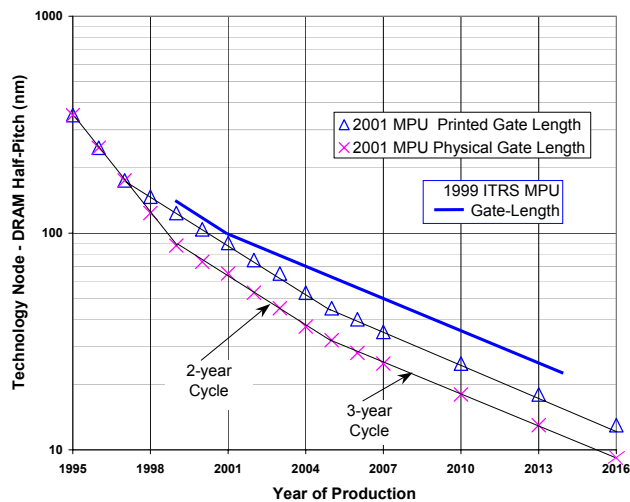


Source: 2001 ITRS - Exec. Summary, ORTC Figure
Caltech CS184 Winter2003 -- DeHon

[from Andrew Kahng]

27

ITRS Roadmap Acceleration Continues...Gate Length



Source: 2001 ITRS - Exec. Summary, ORTC Figure
Caltech CS184 Winter2003 -- DeHon

[from Andrew Kahng]

28

Delays?

- If delays in gates/switching?
- If delays in interconnect?
- Logical interconnect lengths?

Delays?

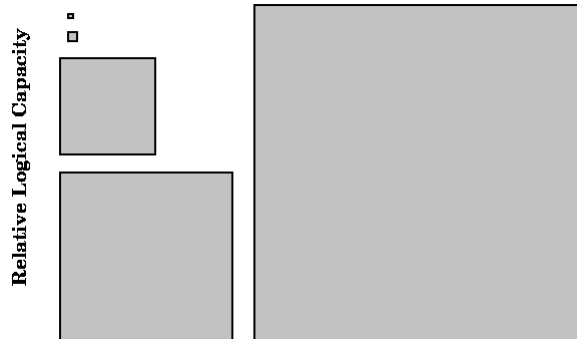
- If delays in gates/switching?
 - Delay reduce with $1/\kappa [\lambda]$

Delays

- Logical capacities growing
- Wirelengths?
 - No locality $\rightarrow \kappa$ (slower!)

- Rent's Rule

- $L \rightarrow n^{(p-0.5)}$
- $[p > 0.5]$



Caltech CS184 Winter2003 -- DeHon

Compute Density

- Density = compute / (Area * Time)
- κ^3 : compute density scaling $> \kappa$
- κ^3 : gates dominate, $p < 0.5$
- κ^2 : moderate p , good fraction of gate delay
 - $[p$ from Rent's Rule again – more on Day12]
- κ : large p (wires dominate area and delay)

Caltech CS184 Winter2003 -- DeHon

Power Density

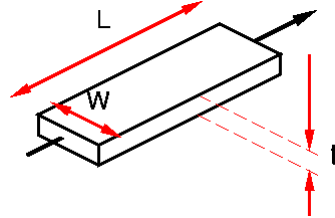
- $P \rightarrow P/\kappa^2$ (static, or increase frequency)
- $P \rightarrow P/\kappa^3$ (dynamic, same freq.)
- $A \rightarrow A/\kappa^2$
- $P/A \rightarrow P/A \dots$ or $\dots P/\kappa A$

Cheating...

- Don't like some of the implications
 - High resistance wires
 - Higher capacitance
 - Need for more wiring
 - Not scale speed fast enough

Improving Resistance

- $R = \rho L / (W \cdot t)$
- $W \rightarrow W / \kappa$
- L, t similar
- $R \rightarrow \kappa R$



- Don't scale t quite as fast.
- Decrease ρ (copper)

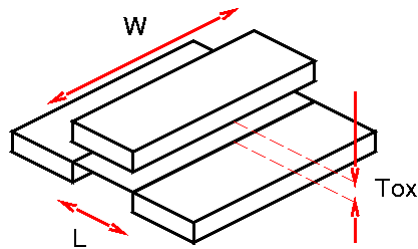
Improving Capacitance

- Capacitance per unit area

$$- C_{ox} = \epsilon_{SiO_2} / T_{ox}$$

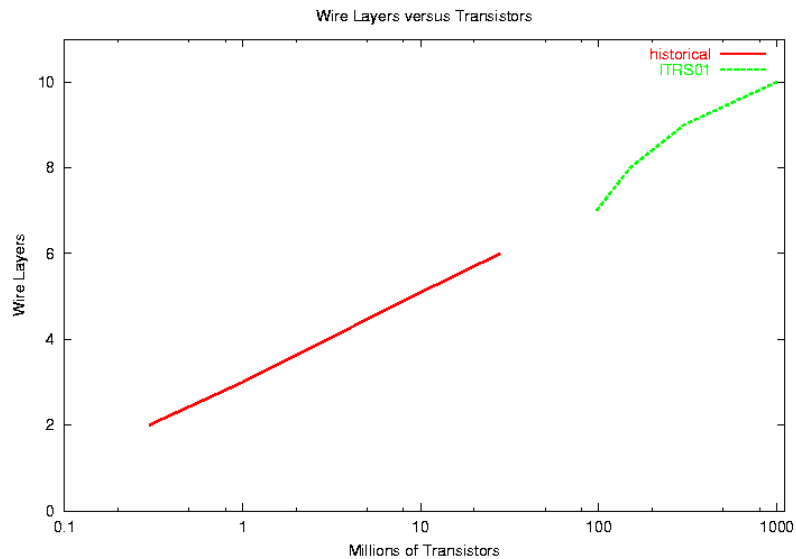
$$- T_{ox} \rightarrow T_{ox} / \kappa$$

$$- C_{ox} \rightarrow \kappa C_{ox}$$

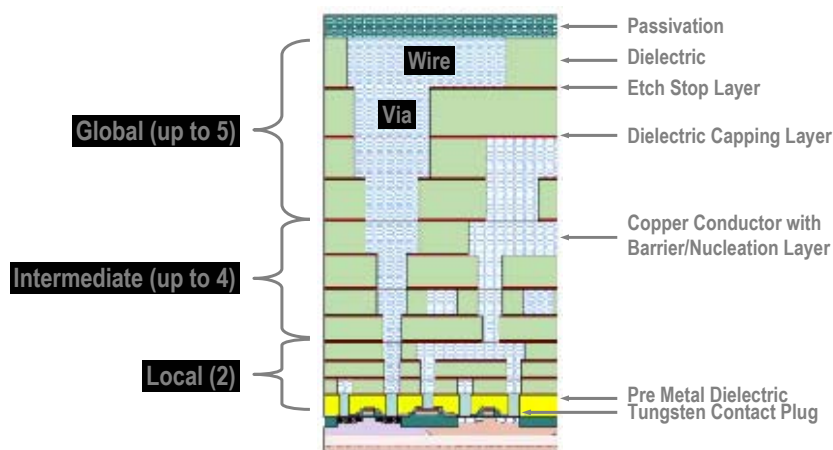


Reduce Dielectric Constant ϵ

Wire Layers = More Wiring



Typical chip cross-section illustrating hierarchical scaling methodology



[from Andrew Kahng]

Improving Gate Delay

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

- $V \rightarrow V/\kappa$

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

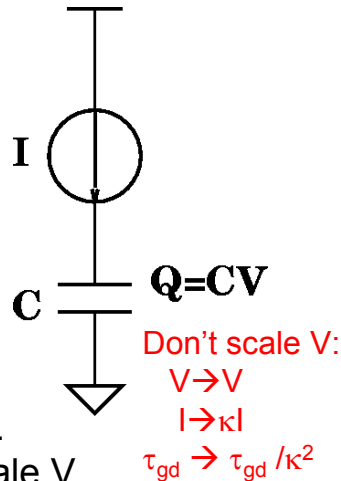
- $I_d \rightarrow I_d/\kappa$

- $C \rightarrow C/\kappa$

- $\tau_{gd} \rightarrow \tau_{gd}/\kappa$

- Lower C.

- Don't scale V.



39

Caltech CS184 Winter2003 -- DeHon

...But Power Dissipation (Dynamic)

- Capacitive (Dis)charging

- $P = (1/2)CV^2f$

- $V \rightarrow V/\kappa$

- $C \rightarrow C/\kappa$

- $P \rightarrow P/\kappa^3$

- Increase Frequency?

- $f \rightarrow \kappa f$?

- $P \rightarrow P/\kappa^2$

If not scale V, power dissipation not scale.

40

Caltech CS184 Winter2003 -- DeHon

...And Power Density

- $P \rightarrow P$ (increase frequency)
- $P \rightarrow > P/\kappa$ (dynamic, same freq.)
- $A \rightarrow A/\kappa^2$
- $P/A \rightarrow \kappa P/A \dots$ or $\dots \kappa^2 P/A$
- Power Density Increases

Physical Limits

- Doping?
- Features?

Physical Limits

- Depended on
 - bulk effects
 - doping
 - current (many electrons)
 - mean free path in conductor
 - localized to conductors
- Eventually
 - single electrons, atoms
 - distances close enough to allow tunneling

What Is A “Red Brick” ?

- Red Brick = ITRS Technology Requirement with no known solution
- Alternate definition: Red Brick = something that REQUIRES billions of dollars in R&D investment

The “Red Brick Wall” - 2001 ITRS vs 1999

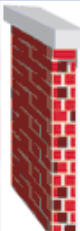
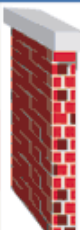
Table 1. 2001 Status of Red Brick Wall							
Year of production	2001	2003	2005		2007	2010	2016
DRAM half-pitch (nm)	130	100	80		65	45	22
Overlay accuracy (nm)	46	35	28		23	18	9
MPU gate length (nm)	90	65	45		35	25	13
CD control (nm)	8	5.5	3.9		3.1	2.2	1.1
T _{ox} (equivalent) (nm)	1.3-1.6	1.1-1.6	0.8-1.3		0.6-1.1	0.5-0.8	0.4-0.5
Junction depth (nm)	48-95	33-66	24-47		18-37	13-26	7-13
Metal cladding thickness (nm)	16	12	9		7	5	2.5
Intermetal dielectric constant, k	3.0-3.6	3.0-3.6	2.6-3.1		2.3-2.7	2.1	1.8

Table 2. 1999 Status of Red Brick Wall							
Year of production	1999	2002	2005		2008	2011	2014
DRAM half-pitch (nm)	180	130	100		70	50	35
Overlay accuracy (nm)	65	45	35		25	20	15
MPU gate length (nm)	140	85-90	65		45	30-32	20-22
CD control (nm)	14	9	6		4	3	2
T _{ox} (equivalent) (nm)	1.9-2.5	1.5-1.9	1.0-1.5		0.8-1.2	0.6-0.8	0.5-0.6
Junction depth (nm)	42-70	25-43	20-33		16-26	11-19	8-13
Metal cladding thickness (nm)	17	13	10		0	0	0
Intermetal dielectric constant, k	3.5-4.0	2.7-3.56	1.6-2.2		1.5	<1.5	<1.5

Source: Semiconductor International - <http://www.e-insite.net/semiconductor/index.asp?layout=article&articleId=CA187876>

Caltech CS184 Winter2003 -- DeHon [from Andrew Kahng]

45

Finishing Up...

Big Ideas [MSB Ideas]

- Moderately predictable VLSI Scaling
 - unprecedented capacities/capability growth for engineered systems
 - **change**
 - be prepared to exploit
 - account for in comparing across time

Big Ideas [MSB-1 Ideas]

- Uniform scaling reasonably accurate for past couple of decades
- Area increase κ^2
 - Real capacity maybe a little less?
- Gate delay decreases ($1/\kappa$)
- Wire delay not decrease, maybe increase
- Overall delay decrease less than ($1/\kappa$)