

ESE534: Computer Organization

Day 7: February 8, 2010
VLSI Scaling



Today

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

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.
- Make direct comparison across technologies
 - E.g. to understand older designs
 - What comes from process vs. architecture

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

Preclass

- When will we have 32-core processors?

- Bits/DRAM chip in 2020?

Scaling

- **Premise:** features scale “uniformly”
 - everything gets better in a predictable manner
- **Parameters:**
 - λ (lambda) -- Mead and Conway (class)
 - F -- Half pitch – ITRS ($F=2\lambda$)
 - $1/\kappa$ – Dennard
 - S -- Bohr

Feature Size

λ is half the minimum feature size in a VLSI process

[minimum feature usually channel width]

7

Penn ESE534 Spring 2010 -- DeHon

Scaling

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

8

Penn ESE534 Spring 2010 -- DeHon

Scaling

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

9

Penn ESE534 Spring 2010 -- DeHon

Effects?

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

- Go through traditional / ideal
- ...then come back and ask what still makes sense today.

10

Penn ESE534 Spring 2010 -- DeHon

Area

- $\lambda \rightarrow \lambda/\kappa$
- $A = L * W$
- $A \rightarrow A/\kappa^2$

- 130nm \rightarrow 90nm
- 50% area
- 2x capacity same area

11

Penn ESE534 Spring 2010 -- DeHon

Area Perspective

Die size	Relative Logical Capacity
196mm	1
130mm	~10
70mm	~100
25mm	~1000
90nm	~10000
45nm	~100000
25nm	~900000

Penn E

Current

- Saturation Current

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

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

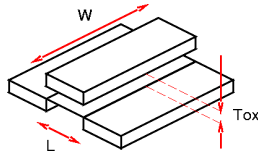
$$V_{TH} \rightarrow V_{TH}/\kappa$$

$$W \rightarrow W/\kappa$$

$$L \rightarrow L/\kappa$$

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

$$I_d \rightarrow I_d/\kappa$$

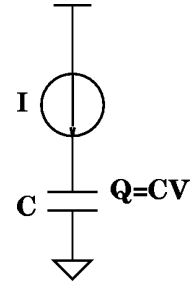


19

Penn ESE534 Spring 2010 -- DeHon

Gate Delay

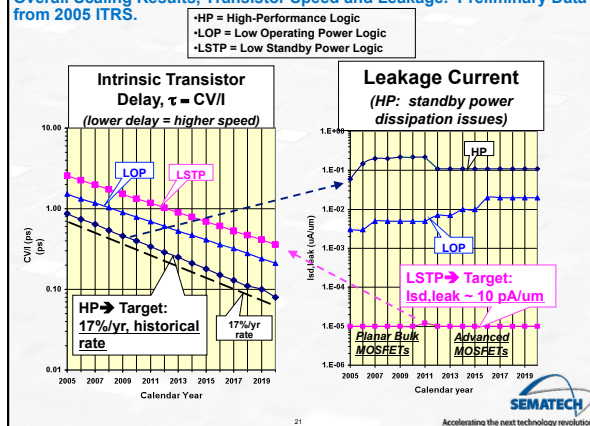
- $\tau_{gd} = Q/I = (CV)/I$
- $V \rightarrow V/\kappa$
- $I_d \rightarrow I_d/\kappa$
- $C \rightarrow C/\kappa$
- $\tau_{gd} \rightarrow \tau_{gd}/\kappa$



20

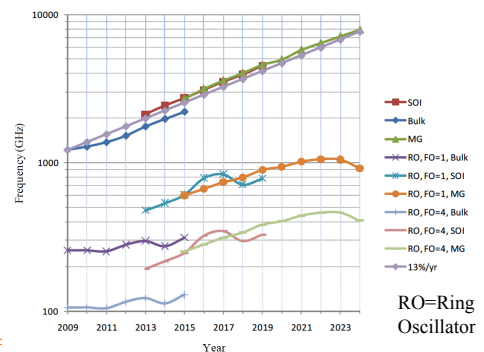
Penn ESE534 Spring 2010 -- DeHon

Overall Scaling Results, Transistor Speed and Leakage. Preliminary Data from 2005 ITRS.



Penn ESE534 Sp

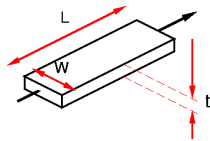
ITRS 2009 Transistor Speed



Penn ESE534 Sp

Resistance

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



23

Penn ESE534 Spring 2010 -- DeHon

Wire Delay

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

24

Penn ESE534 Spring 2010 -- DeHon

Power Dissipation (Static Load)

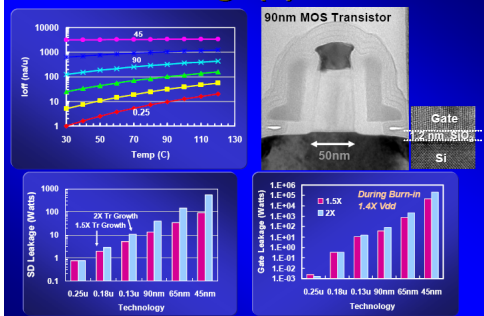
- Resistive Power
 - $P=V \cdot I$
 - $V \rightarrow V/\kappa$
 - $I_d \rightarrow I_d/\kappa$
 - $P \rightarrow P/\kappa^2$

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?
 - $\tau_{gd} \rightarrow \tau_{gd}/\kappa$
 - So: $f \rightarrow \kappa f$
 - $P \rightarrow P/\kappa^2$

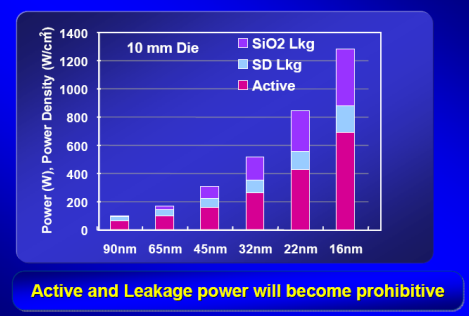
...and leakage

The Leakage(s)...



Intel on Leakage

Projected Power (unconstrained)



Active and Leakage power will become prohibitive

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]

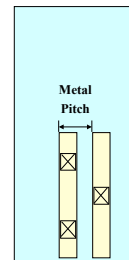


Source: 2001 ITRS - Exec. Summary, ORTC Figure
Penn ESE534 Spring 2010 -- DeHon

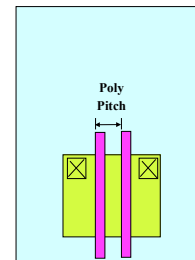
[from Andrew Kahng]

31

Half Pitch (= Pitch/2) Definition



(Typical
DRAM)



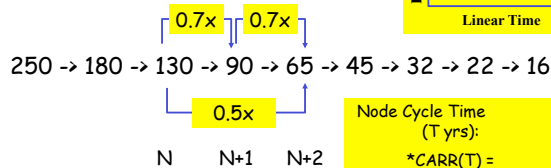
(Typical
MPU/ASIC)

Source: 2001 ITRS - Exec. Summary, ORTC Figure
Penn ESE534 Spring 2010 -- DeHon

[from Andrew Kahng]

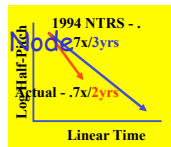
32

Scaling Calculator + Cycle Time:



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

Node Cycle Time
(T yrs):
*CARR(T) =
[[0.5]^(1/2T yrs)] - 1
CARR(3 yrs) = -10.9%
CARR(2 yrs) = -15.9%

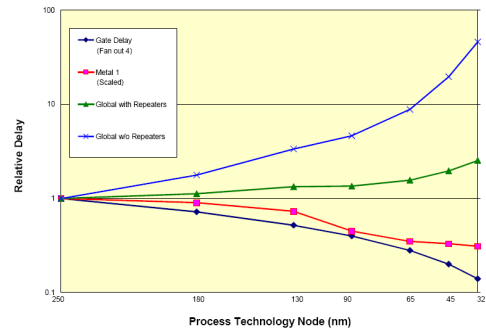


Source: 2001 ITRS - Exec. Summary, ORTC Figure
Penn ESE534 Spring 2010 -- DeHon

[from Andrew Kahng]

33

ITRS 2003,2005 Gate/Wire Scaling



Penn ESE534 Spring 2010 -- DeHon

34

What happens to delays?

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

Penn ESE534 Spring 2010 -- DeHon

35

Delays?

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

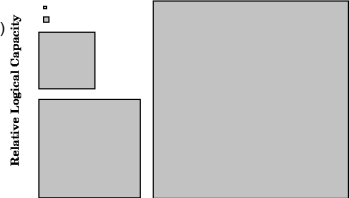
Penn ESE534 Spring 2010 -- DeHon

36

Delays

- Logical capacities growing
- Wirelengths?
 - No locality: $L \rightarrow \kappa$ (slower!)
 - Rent's Rule

- L proportional $\eta^{(p-0.5)}$
- $[p > 0.5]$



Penn ESE534 Spring 2010 -- DeHon

Compute Density

- Density = compute / (Area * Time)
- $\kappa^3 > \text{compute density scaling} > \kappa$
- κ^3 : gates dominate, $p < 0.5$
- κ^2 : moderate p , good fraction of gate delay
 - $[p$ from Rent's Rule again – more on Day18]
- κ : large p (wires dominate area and delay)

Penn ESE534 Spring 2010 -- DeHon

38

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 $\dots P/\kappa A$

Penn ESE534 Spring 2010 -- DeHon

39

Cheating...

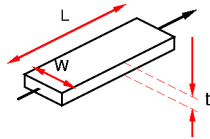
- Don't like some of the implications
 - High resistance wires
 - Higher capacitance
 - Atomic-scale dimensions
 - Quantum tunneling
 - Need for more wiring
 - Not scale speed fast enough
 - **Finite subthreshold slope (Wed.)**

Penn ESE534 Spring 2010 -- DeHon

40

Improving Resistance

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



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

Penn ESE534 Spring 2010 -- DeHon

41

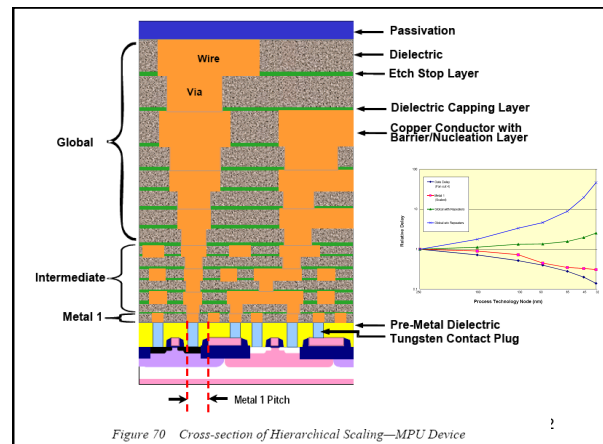


Figure 70 Cross-section of Hierarchical Scaling—MPU Device

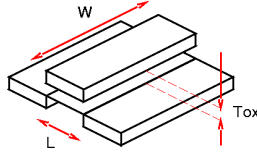
Capacitance and Leakage

- 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 (interconnect)

and Increase Dielectric to substitute for scaling T_{ox} (gate quantum tunneling) 43

ITRS 2009

Table PIDS3B Low Operating Power Technology Requirements

Grey cells delineate one of two time periods: either before initial production ramp has started for ultra-thin body fully depleted (UTB FD) SOI or multi-gate (MG) MOSFETs, or beyond when planar bulk or UTB FD MOSFETs have reached the limits of practical scaling (see the text and the table notes for further discussion).

Year of Production	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
WPU/SNK Metal 1 (M1) Pitch (nm) (continued)	54	45	38	32	27	24	21	18.9	16.9	15	13.4	11.9	10.6	9.5	8.4	7.5
Physical Logic for High Performance logic (nm)	29	27	24	22	20	18	17	15.3	14	12.8	11.7	10.7	9.7	8.9	8.1	7.4
Physical Logic for Low Operating Power (LOP) logic (nm) (1)	32	29	27	24	22	18	17	15.3	14	12.8	11.7	10.7	9.7	8.9	8.1	7.4
LOP Equivalent Oxide Thickness (nm) (2)	1	0.9	0.8	0.85	0.8	0.8	0.8	0.75	0.7							
Extended Planar Bulk																
UTB FD																
MG																
Gate poly depletion (nm) (3)																
Bulk	0.27	0.27	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Channel doping (E18 cm ⁻³) (4)																
Extended Planar Bulk	3	3.7	4.5	5	5.5	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
Minimum depth or body Thickness (nm) (5)																
Extended Planar Bulk (specimen)	14	13	11.5	10	9	8.2	6	6.1	4.7							
UTB FD (body)																
MG (body)																
LOP Equivalent Oxide Thickness (nm) (6)	1.64	1.52	1.23	1.18	1.14											
Extended Planar Bulk																
UTB FD																
MG																

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.

Dielectric	Dielectric constant (bulk)	Bandgap (eV)	Conduction band offset (eV)	Leakage current reduction w.r.t. SiO ₂	Thermal stability w.r.t. silicon (MEIS data)
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 ² -10 ³ ×	Silicate formation
Hafnium oxide (HfO ₂)	~20	6	1.5	10 ² -10 ³ ×	~950°C
Zirconium oxide (ZrO ₂)	~23	5.8	1.4	10 ² -10 ³ ×	~900°C
Strontium titanate (SrTiO ₃)	3.3		-0.1		
Zirconium silicate (ZrSiO ₄)	6*		1.5		
Hafnium silicate (HfSiO ₄)	6*		1.5		

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

Intel NYT Announcement

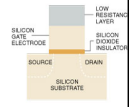
- Intel Says Chips Will Run Faster, Using Less Power

- NYT 1/27/07, John Markov
- Claim: "most significant change in the materials used to manufacture silicon chips since Intel pioneered the modern integrated-circuit transistor more than four decades ago"
- "Intel's advance was in part in finding a new insulator composed of an alloy of hafnium... will replace the use of silicon dioxide."

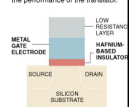
Small and Efficient

All microprocessor transistors become smaller, stopping undesired current leakage becomes more difficult. This leakage leads to shortened battery life. Intel's coming chips use a new insulator material to prevent this, reducing power consumption.

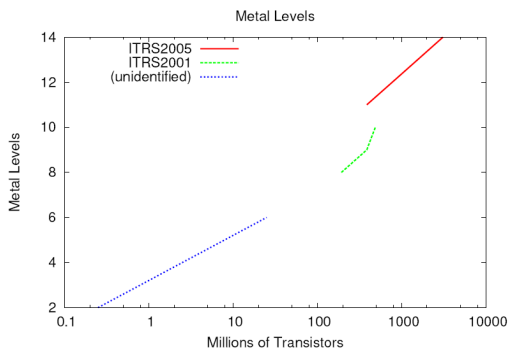
Current transistors use an extremely thin silicon dioxide insulator, which leads to current leakage. Transistors then decrease the leakage but reduce the electric charge passing through, impacting performance.



New transistors use a hafnium-based insulator and a metal gate electrode. Hafnium provides stronger electrical coupling, so the insulator can be made thicker to reduce leakage without degrading the performance of the transistor.



Wire Layers = More Wiring



Typical chip cross-section illustrating hierarchical scaling methodology

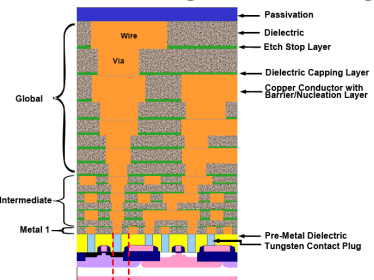


Figure 70 Cross-section of Hierarchical Scaling—MPU Device

Improving Gate Delay

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

- $V \rightarrow V/\kappa$

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

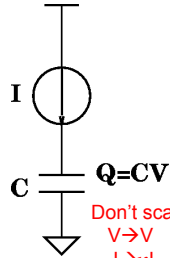
- $I_d \rightarrow I_d/\kappa$

- $C \rightarrow C/\kappa$

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

- Lower C.

- Don't scale V.



Don't scale V:
 $V \rightarrow V$
 $I \rightarrow I/\kappa$
 $\tau_{gd} \rightarrow \tau_{gd}/\kappa^2$

...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.

...And Power Density

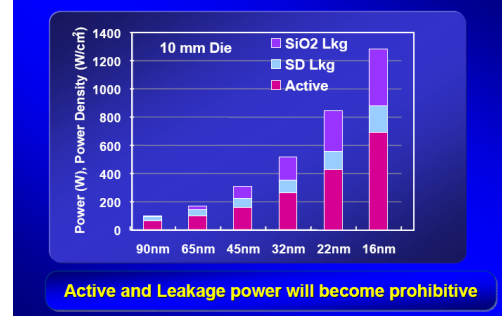
- $P \rightarrow P$ (increase frequency)
- $P \rightarrow P/\kappa$ (dynamic, same freq.)
- But... $A \rightarrow A/\kappa^2$
- $P/A \rightarrow \kappa P/A$... or ... $\kappa^2 P/A$

- **Power Density Increases**

...this is where some companies have gotten into trouble...

Intel on Leakage

Projected Power (unconstrained)



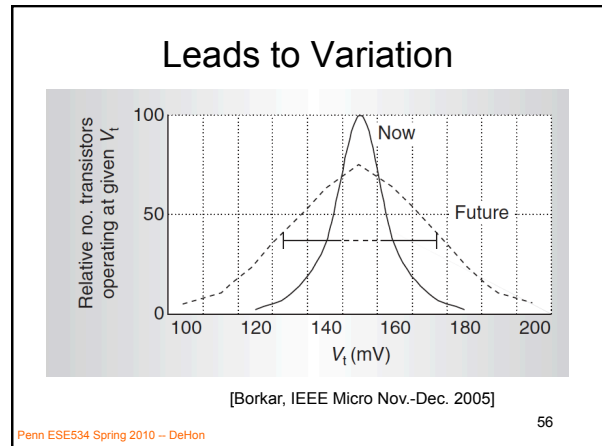
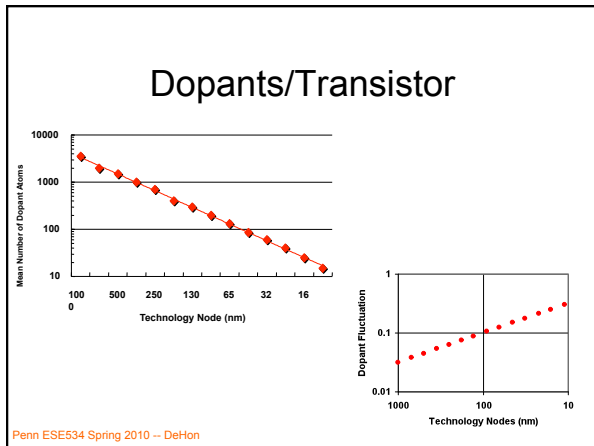
Active and Leakage power will become prohibitive

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



Electrons

Table 40b High-Performance Logic Technology Requirements—Long-term (continued)

Grey cells delineate one of two time periods, either before initial production ramp has started for ultra-thin body fully depleted (UTB FD) SOI or double-gate (DG) MOSFETs, or beyond when planar bulk or UTB FD MOSFETs have reached the limits of practical scaling (see the text and the table notes for further discussions).

Year of Production	2014	2015	2016	2017	2018	2019	2020
DRAM half-pitch (nm)	28	25	22	20	18	16	14
MPU ASIC Metal 1 Q11 half-pitch (nm)	28	25	22	20	18	16	14
MPU Physical Gate Length (nm)	11	10	9	8	7	6	6
R_{sd} Effective Parasitic series source/drain resistance [12]							
Planar Bulk ($\Omega\text{-}\mu\text{m}$)							
UTB FD ($\Omega\text{-}\mu\text{m}$)	75	75					
DG ($\Omega\text{-}\mu\text{m}$)	85	80	75	70	65	60	55
$C_{g, \text{total}}$ Total gate capacitance for calculation of CVT [14]							
Extended Planar Bulk (fF/ μm^2)							
UTB FD (fF/ μm^2)	4.22E-16	3.83E-16					
DG (fF/ μm^2)	3.80E-16	3.45E-16	3.45E-16	3.07E-16	2.68E-16	2.30E-16	1.92E-16
Extended Planar Bulk (fF/ μm^2)							
UTB FD (fF/ μm^2)	8.42E-16	5.03E-16					
DG (fF/ μm^2)	5.69E-16	5.25E-16	5.25E-16	4.87E-16	4.48E-16	4.10E-16	3.62E-16

$e = 1.6 \times 10^{-19} \text{ C}$

How many electrons?

$F = 14 \text{ nm}$ $C_{\text{min}} = 1.4\text{E-}2 \times 3.6\text{E-}16 = 5\text{E-}18 = 30e^{57}$

Penn ESE534 Spring 2010 -- DeHon

- ### 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
- [from Andrew Kahng]
- Penn ESE534 Spring 2010 -- DeHon

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	12-26	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	20-36	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	1.8-2.2	1.5	<1.5	<1.5

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

Penn ESE534 Spring 2010 -- DeHon [from Andrew Kahng]

ITRS 2009

Year of production	2009	2011	2013	2015	2017	2019	2021	2023	2025	2027	2029	2031	2033	2035
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	12-26	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								

Penn ESE534 Spring 2010 -- DeHon

Conventional Scaling

- Ends in your lifetime
- ...perhaps in your first few years out of school...
- 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

Finishing Up...

Admin

- Reading
 - Wed. on blackboard
 - None for next Monday
 - Next Wed. (will be) on blackboard

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
 - ...but not for much longer

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 ($1/\kappa$)
 - ...maybe a little less
- Wire delay not decrease, maybe increase
- Overall delay decrease less than ($1/\kappa$)