# CS184a:
# Computer Architecture
# (Structure and Organization)

Day 6:  January 19, 2005
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.

---

# 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
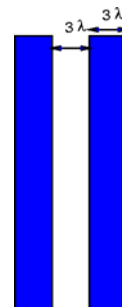  - lose out to someone who can stay on curve

---

# Scaling

- **Premise:** features scale "uniformly"
  - everything gets better in a predictable manner

- **Parameters:**
  - $\lambda$ (lambda) -- Mead and Conway (class)
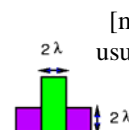  - S -- Bohr
  - $1/\kappa$ -- Dennard

---

# Feature Size



$\lambda$ is half the minimum feature size in a VLSI process
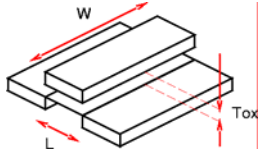
[minimum feature usually channel width]

---

1

## Scaling

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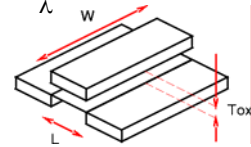- Oxide Thickness ($T_{ox}$)
- Doping ($N_a$)
- Voltage (V)

7

## Scaling

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

8

## Effects?

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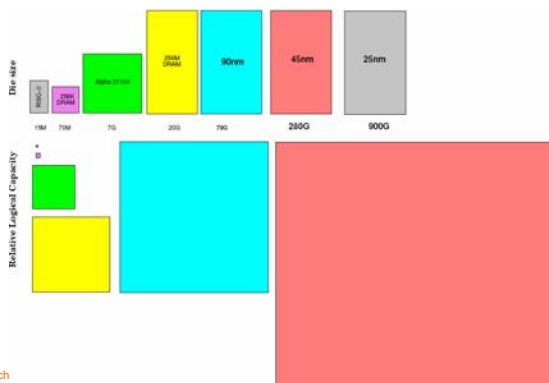- Wire Delay ($\tau_{wire}$)
- Power

9

## Area

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

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

10

## Area Perspective

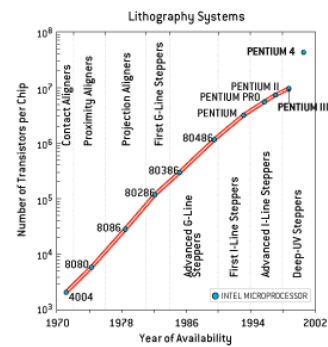## Capacity Scaling from Intel

12

2

## Capacitance

- Capacitance per unit area
  - $C_{ox} = \varepsilon_{SiO_2}/T_{ox}$
  - $T_{ox} \rightarrow T_{ox}/\kappa$
  - $C_{ox} \rightarrow \kappa C_{ox}$

## Capacitance

- Gate Capacitance
  - $C_{gate} = A*C_{ox}$
  - $A \rightarrow A/\kappa^2$
  - $C_{ox} \rightarrow \kappa C_{ox}$
  - $C_{gate} \rightarrow C_{gate}/\kappa$

## Threshold Voltage

**Before:**

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

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

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

**After:**

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

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

## Threshold Voltage

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

## Current

- Saturation Current
  $I_d = (\mu C_{OX}/2)(W/L)(V_{gs} - V_{TH})^2$
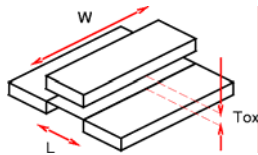
  $V_{gs=}V \rightarrow V/\kappa$
  $V_{TH} \rightarrow V_{TH}/\kappa$
  $W \rightarrow W/\kappa$
  $L \rightarrow L/\kappa$
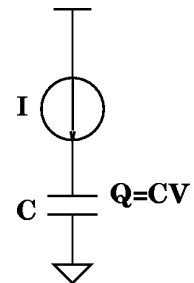  $C_{ox} \rightarrow \kappa C_{ox}$

  $I_d \rightarrow I_d/\kappa$

## Gate Delay

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



$I$

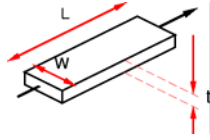$C$　$Q=CV$

## Resistance

- $R = \rho L / (W \cdot t)$

- $W \rightarrow W/\kappa$
- L, t similar

- $R \rightarrow \kappa R$

19

## 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)

20

## Power Dissipation (Static Load)

- Resistive Power
  - $P = V \cdot I$

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

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

21

## Power Dissipation (Dynamic)

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

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

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

- Increase Frequency?

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

  - So: $f \rightarrow \kappa f$ ?

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

22

## Effects?

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

23

## ITRS Roadmap

- Semiconductor Industry rides this scaling curve
- Try to predict where industry going
  - (requirements…self fulfilling prophecy)

- http://public.itrs.net

24

4

## Slide 25

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

# S=0.7
## [0.5x per 2 nodes]



Pitch  Gate

[from Andrew Kahng]  25

## Slide 26

### Half Pitch (= Pitch/2) Definition



Metal Pitch

Poly Pitch

**(Typical DRAM)**

**(Typical MPU/ASIC)**

[from Andrew Kahng]  26

## Slide 27

# Scaling Calculator +
## Node Cycle Time:

**1994 NTRS - .7x/3yrs**

**Actual - .7x/2yrs**

Log Half-Pitch

Linear Time

0.7x  0.7x

250 -> 180 -> 130 -> 90 -> 65 -> 45 -> 32 -> 22 -> 16

0.5x

N    N+1   N+2

* 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%

[from Andrew Kahng]  27

## Slide 28

### ITRS Roadmap Acceleration Continues…Gate Length

[from Andrew Kahng]  28

## Slide 29

# ITRS 2003 Gate/Wire Scaling

29

## Slide 30

# What happens to delays?

- If delays in gates/switching?

- If delays in interconnect?

- Logical interconnect lengths?

30

## Delays?

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

## Delays

- Logical capacities growing
- Wirelengths?
  - No locallity: $L \rightarrow \kappa$    (slower!)
  - Rent's Rule
    - $L \rightarrow n^{(p-0.5)}$
    - [p>0.5]

## Compute Density

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

## Power Density

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

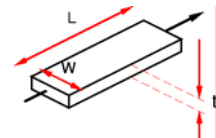- $P/A \rightarrow P/A$ … or … $P/\kappa A$

## Cheating…

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

## Improving Resistance

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

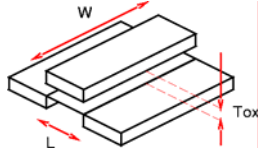  - Don't scale t quite as fast.
  - Decrease $\rho$  (copper)

## Capacitance and Leakage

- Capacitance per unit area
  - $C_{ox} = \varepsilon_{SiO_2}/T_{ox}$
  - $T_{ox} \rightarrow T_{ox}/\kappa$
  - $C_{ox} \rightarrow \kappa\, C_{ox}$



W

Tox

L

Reduce Dielectric Constant $\varepsilon$ (interconnect)

or Substitute for scaling $T_{ox}$ (gate quantum tunneling)

37

## Threshold Voltage

**Before:**

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

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

$$\textbf{adjust } V_{s\text{-sub}} \textbf{ so } \left(\phi_s + V_{s\text{-sub}}\right) \rightarrow \frac{\left(\phi_s + V_{s\text{-sub}}\right)}{\kappa}$$

**After:**

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

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

## ITRS 2003



Table 47a   High-performance Logic Technology Requirements—Near-term

Table 81a

39

## 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.

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

40

## Wire Layers = More Wiring

## Typical chip cross-section illustrating hierarchical scaling methodology



[from Andrew Kahng]

42

7

## Improving Gate Delay

- $\tau_{gd}=Q/I=(CV)/I$
- $V \rightarrow V/\kappa$
- $I_d=(\mu C_{OX}/2)(W/L)(V_{gs}-V_{TH})^2$
- $I_d \rightarrow I_d/\kappa$
- $C \rightarrow C/\kappa$
- $\tau_{gd} \rightarrow \tau_{gd}/\kappa$

**I**

**C** ⊥ **Q=CV**

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

- Lower C.
- Don't scale V.

## …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.)
- $A \rightarrow A/\kappa^2$

- $P/A \rightarrow \kappa P/A$ … or … $\kappa^2 P/A$
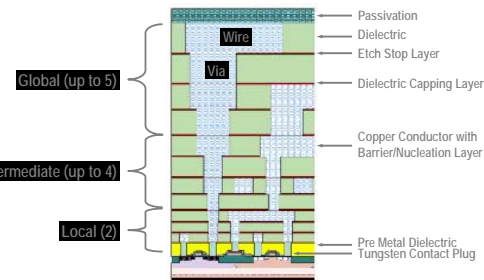
- **Power Density Increases**

  …this is where some companies have gotten into trouble…

## 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 <u>no known solution</u>

- Alternate definition:  Red Brick = something that REQUIRES billions of dollars in R&D investment

[from Andrew Kahng]

## 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 | 0.5-0.8 |
| 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

[from Andrew Kahng]

49

---

## Conventional Scaling

- Ends in your lifetime
- …perhaps in your first few years after grad school…

50

---

# Finishing Up...

51

---

## 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

52

---

## Big Ideas
## [MSB-1 Ideas]

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

53