

CS184a: Computer Architecture (Structures and Organization)

Day6: October 11, 2000
Instruction Taxonomy
VLSI Scaling

Last Time

- Computing requirements
- Instruction requirements
- Structure

Today

- Instruction Taxonomy
- VLSI Scaling

Instruction Distribution

- Beyond 64 PE, instruction bandwidth dictates PE size

$$\frac{\sqrt{PE_{\text{area}}} \times 4 \times \sqrt{N}}{(64 \times 8\lambda)} = N$$

$$PE_{\text{area}} = 16K\lambda^2 \times N$$

- Build larger arrays
⇒ processing elements become less dense

Instruction Memory Requirements

- **Idea:** put instruction memory in array
- **Problem:** Instruction memory can quickly dominate area, too
 - Memory Area = $64 \times 1.2K\lambda^2/\text{instruction}$
 - $PE_{\text{area}} = 1M\lambda^2 + (\text{Instructions}) \times 80K\lambda^2$

Instruction Pragmatics

- Instruction requirements *could* dominate array size.
- Standard architecture trick:
 - Look for structure to exploit in “typical computations”

Two Extremes

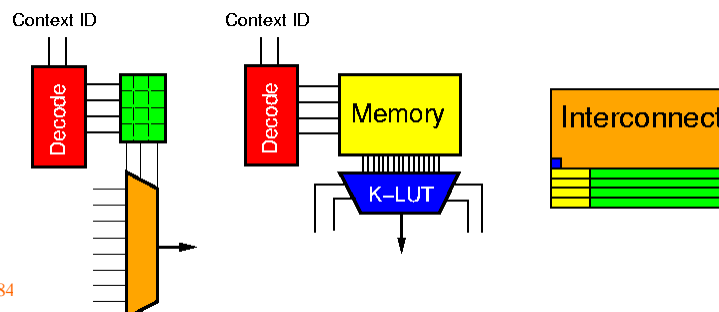
- SIMD Array (microprocessors)
 - Instruction/cycle
 - share instruction across array of PE s
 - uniform operation in space
 - operation variance in time
- FPGA
 - Instruction/PE
 - assume temporal locality of instructions (same)
 - operation variance in space
 - uniform operations in time

Caltech CS184a Fall2000 -- DeHon

7

Hybrids

- VLIW (SuperScalar)
 - Few *pinsts*/cycle
 - Share instruction across w bits
- DPGA
 - Small instruction store / PE



Caltech CS184

8

Architecture Instruction Taxonomy

Control Threads (PCs)					
<i>pinsts</i> per Control Thread					
Instruction Depth					
Granularity					
Architecture/Examples					
0	0	0	n/a	Hardwired Functional Unit (e.g. ECC/EDC Unit, FP MPY)	
	n	1	1	FPGA	
		w	w	Reconfigurable ALUs	
1	1	c	$n_v \cdot 1$	Bitwise SIMD	
			w	Traditional Processors	
			$n_v \cdot w$	Vector Processors	
	n	c	1	DPGA	
			8	PADDI	
			w	VLIW	
m	n	1	1	HSRA/SCORE	
		1	c	$n_v \cdot w$	MSIMD
		c	1	VEGA	
m	1	8	16	PADDI-2	
		c	w	MIMD (traditional)	

Caltech CS184a Fall2000 --

9

Instruction Message

- Architectures fall out of:
 - general model too expensive
 - look for structure in common problems
 - exploit structure to reduce resource requirements
- Architectures can be viewed in a unified design space

Caltech CS184a Fall2000 -- DeHon

10

VLSI Scaling

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.

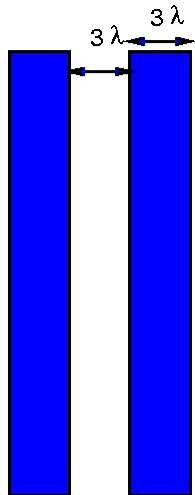
Why Care

- 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

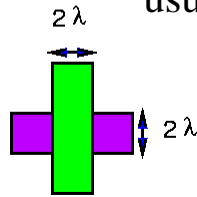
- **Premise:** features scale “uniformly”
 - everything gets better in a predictable manner
- **Parameters:**
 - λ (lambda) -- Mead and Conway (class)
 - S -- Bohr
 - $1/\kappa$ -- Dennard

Feature Size



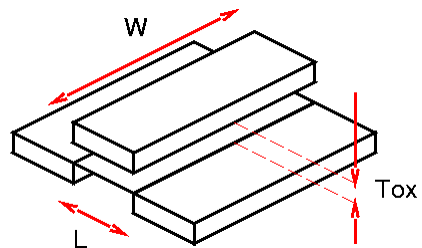
λ is half the minimum feature size in a VLSI process

[minimum feature usually channel width]



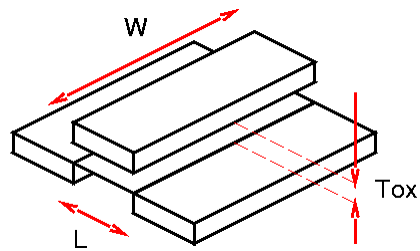
Scaling

- 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

$$\lambda \rightarrow \lambda/\kappa$$

$$A = L * W$$

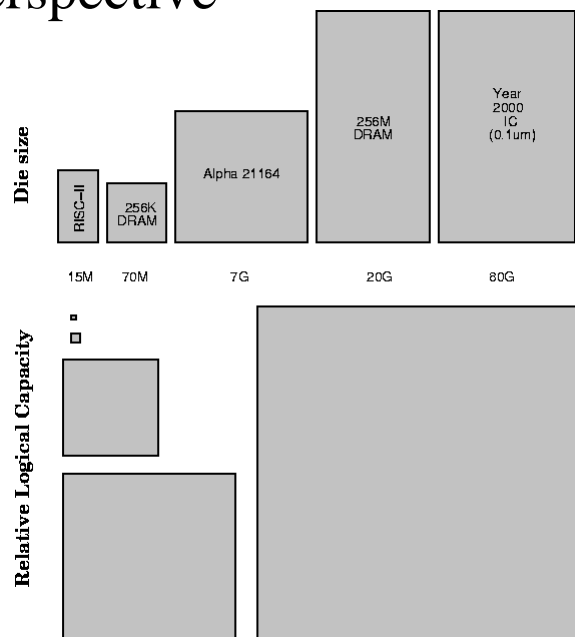
$$A \rightarrow A/\kappa^2$$

$$0.35\mu\text{m} \rightarrow 0.25\mu\text{m}$$

- 50% area
- 2x capacity same area

Area Perspective

[2000 tech.]
 18mm×18mm
 0.18μm
 60G λ²



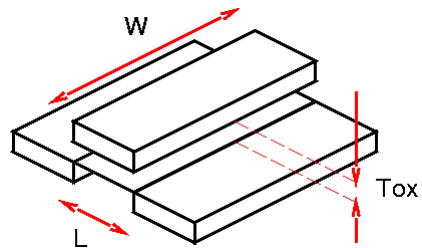
Capacitance

- Capacitance per unit area

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

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

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



Capacitance

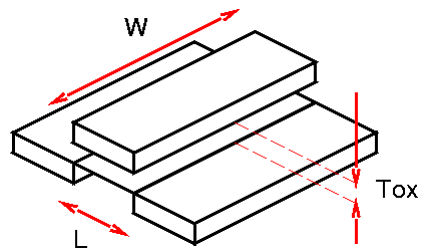
- Gate Capacitance

- $C_{gate} = A * C_{ox}$

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

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

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



Threshold Voltage

Before:

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

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

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

After:

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

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

Threshold Voltage

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

Current

- Saturation Current

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

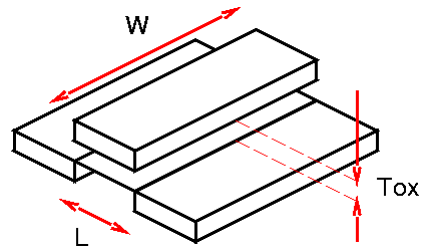
- $V_{gs} \rightarrow V/\kappa$

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

- $W \rightarrow W/\kappa$

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

- $I_d \rightarrow I_d/\kappa$



Gate Delay

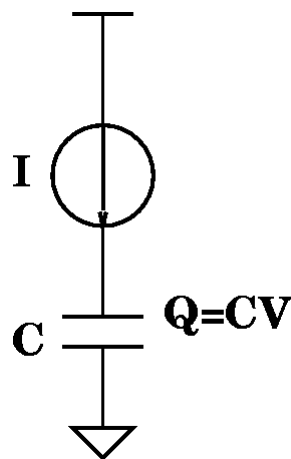
$$\tau_{gd} = Q/I = (CV)/I$$

- $V \rightarrow V/\kappa$

- $I_d \rightarrow I_d/\kappa$

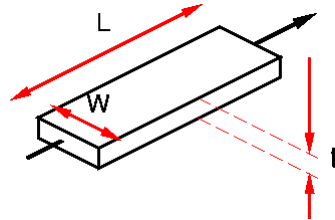
- $C \rightarrow C/\kappa$

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



Resistance

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



Wire Delay

$$\tau_{\text{wire}} = R_L C$$

- $R \rightarrow \kappa R$
- $C \rightarrow C / \kappa$

$$\tau_{\text{wire}} \rightarrow \tau_{\text{wire}}$$

- ...assuming (logical) wire lengths remain constant...

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 $1/\kappa^2 \rightarrow 1/\kappa^3$

Delays?

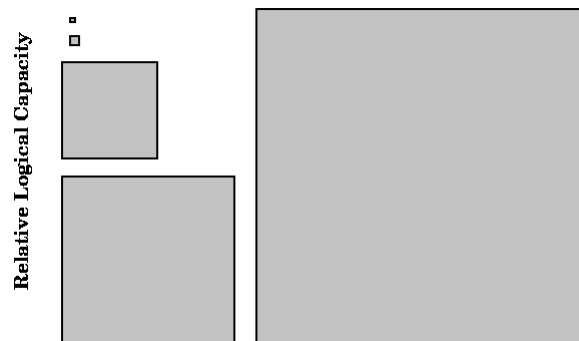
- If delays in gates/switching?
- If delays in interconnect?
- Logical interconnect lengths?

Delays?

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

Delays

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



Capacity

- Rent: $IO=C*N^p$
- $p>0.5$
- $A= C*N^{2p}$
- Logical Area $\rightarrow \kappa^2$
 $\kappa^2 A= C*N_2^{2p}$
 $\kappa^2 N^{2p} = N_2^{2p}$
 $N_2 = \kappa^{(1/p)} N$
- Sanity Check
 - $p=1$
 - $N_2 = \kappa N$
 - $p\sim 0.5$
 - $N_2 \sim \kappa^2 N$

Compute Density

- Density = compute / (Area * Time)

$\kappa^3 >$ compute density scaling $> \kappa$

κ^3 : gates dominate, $p < 0.5$

κ^2 : moderate p , good fraction of gate delay

κ : large p (wires dominate area and delay)

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 $\dots P/\kappa A$

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

Finishing Up...

Big Ideas [MSB Ideas]

- Instruction organization induces a design space (taxonomy) for programmable architectures
- 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 ($1/\kappa$)