

# CS184a: Computer Architecture (Structure and Organization)

Day 8: January 24, 2005  
Computing Requirements and  
Instruction Space



## Previously

- Fixed and Programmable Computation
- Area-Time-Energy Tradeoffs
- VLSI Scaling

## Today

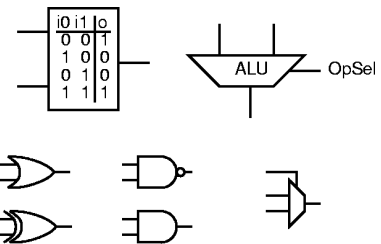
- Computing Requirements
- Instructions
  - Requirements
  - Taxonomy
  - Model Architecture [if time permits]
    - implied costs
    - gross application characteristics

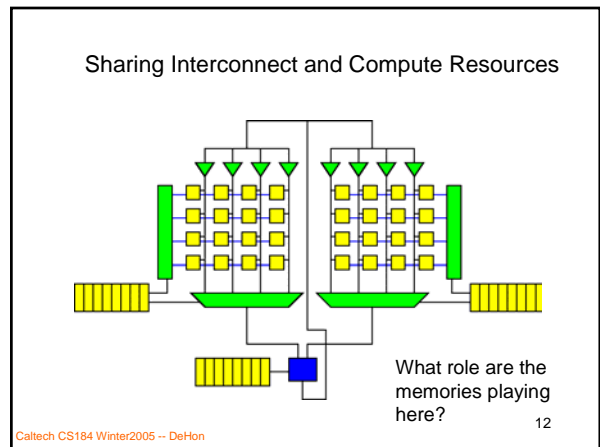
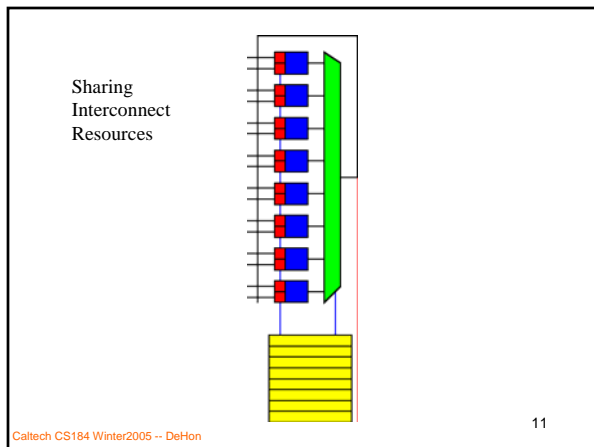
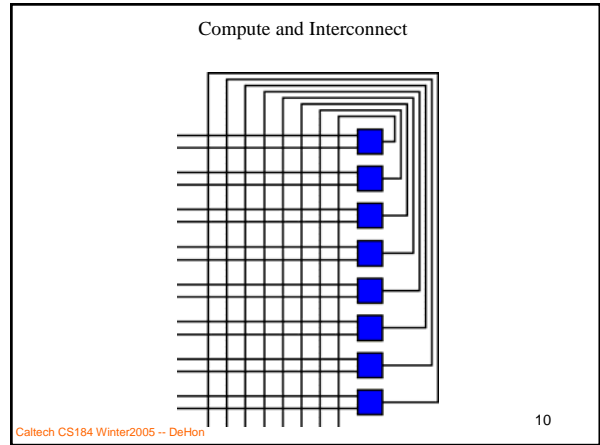
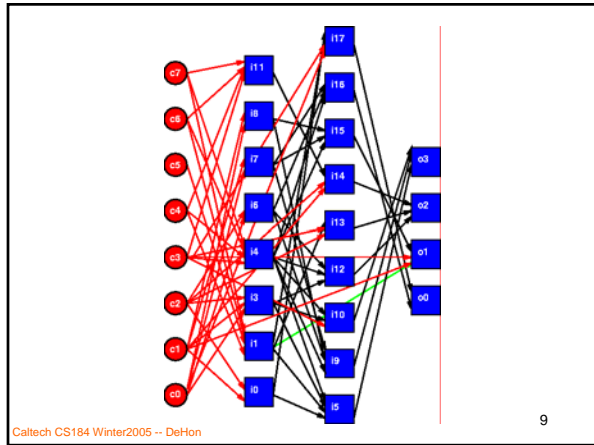
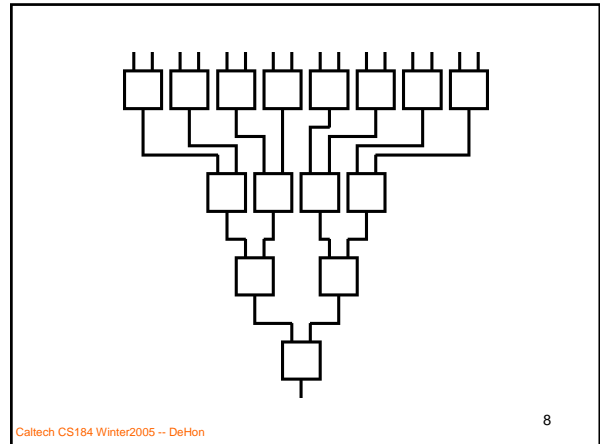
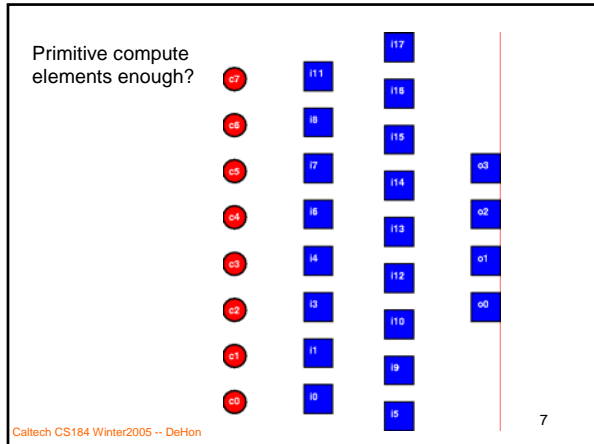
## Computing Requirements (review)

## Requirements

- In order to build a **general-purpose** (*programmable*) computing device, we absolutely must have?

- 
- 
- 
- 
- 
- 





Memory block or Register File

**Interconnect:**  
moves data from input to storage cell;  
or from storage cell to output.

Caltech CS184 Winter2005 -- DeHon 13

What do I need to be able to use this circuit properly?  
(reuse it on different data?)

Caltech CS184 Winter2005 -- DeHon 14

Caltech CS184 Winter2005 -- DeHon 15

### Requirements

- In order to build a **general-purpose (programmable)** computing device, we absolutely must have?
  - Compute elements
  - Interconnect: space
  - Interconnect: time (retiming)
  - Interconnect: external (IO)
  - Instructions

Caltech CS184 Winter2005 -- DeHon 16

### Instruction Taxonomy

Caltech CS184 Winter2005 -- DeHon 17

### Instructions

- Distinguishing feature of programmable architectures?
  - *Instructions* -- bits which tell the device how to behave

Caltech CS184 Winter2005 -- DeHon 18

## Focus on Instructions

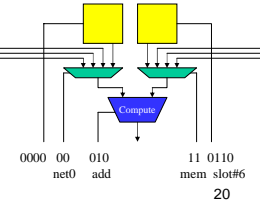
- Instruction organization has a large effect on:
  - size or compactness of an architecture
  - realm of efficient utilization for an architecture

Caltech CS184 Winter2005 -- DeHon

19

## Terminology

- **Primitive Instruction (*pinst*)**
  - Collection of bits which tell a single bit-processing element what to do
  - Includes:
    - select **compute** operation
    - input sources in space
      - (**interconnect**)
    - input sources in time
      - (**retiming**)

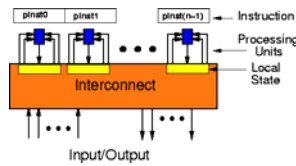


Caltech CS184 Winter2005 -- DeHon

20

## Computational Array Model

- Collection of computing elements
  - compute operator
  - local storage/retiming
- Interconnect
- Instruction

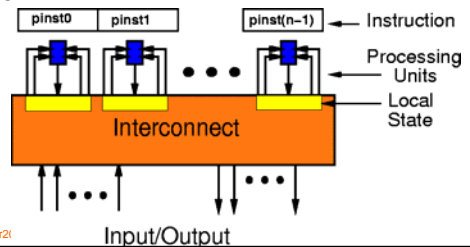


Caltech CS184 Winter2005 -- DeHon

21

## “Ideal” Instruction Control

- Issue a new instruction to every computational bit operator on every cycle



Caltech CS184 Winter21

## “Ideal” Instruction Distribution

- Why don't we do this?

Caltech CS184 Winter2005 -- DeHon

23

## “Ideal” Instruction Distribution

- **Problem:** Instruction bandwidth (and storage area) quickly dominates everything else
  - Compute Block  $\sim 1M\lambda^2$  ( $1K\lambda \times 1K\lambda$ )
  - Instruction  $\sim 64$  bits
  - Wire Pitch  $\sim 8\lambda$
  - Memory bit  $\sim 1.2K\lambda^2$

Caltech CS184 Winter2005 -- DeHon

24

### Instruction Distribution

Two instructions in  $1024\lambda$ .

$64 \times 8\lambda = 512\lambda$

Caltech CS184 Winter2005 -- DeHon 25

### Instruction Distribution

Distribute from both sides =  $2x$

Caltech CS184 Winter2005 -- DeHon 26

### Instruction Distribution

Distribute X and Y =  $2x$

Caltech CS184 Winter2005 -- DeHon 27

### Instruction Distribution

- Room to distribute 2 instructions across PE per metal layer ( $1024 = 2 \times 8 \times 64$ )
- Feed top and bottom (left and right) =  $2x$
- Two complete metal layers =  $2x$
- $\Rightarrow$  8 instructions / PE Side

Caltech CS184 Winter2005 -- DeHon 28

### Instruction Distribution

- Maximum of 8 instructions per PE side
- Saturate wire channels at  $8 \times \sqrt{N} = N$
- $\Rightarrow$  at 64 PE
  - beyond this:
    - instruction distribution dominates area
- Instruction consumption goes with area
- Instruction bandwidth goes with perimeter

Caltech CS184 Winter2005 -- DeHon 29

### Instruction Distribution

- Beyond 64 PE, instruction bandwidth dictates PE size

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

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

- As we build larger arrays
  - $\Rightarrow$  processing elements become less dense

Caltech CS184 Winter2005 -- DeHon 30

## Instruction Memory Requirements

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

Caltech CS184 Winter2005 -- DeHon

31

## Instruction Pragmatics

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

Caltech CS184 Winter2005 -- DeHon

32

## Typical Structure?

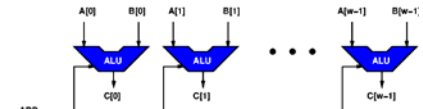
- What structure do we usually expect?

Caltech CS184 Winter2005 -- DeHon

33

## Two Extremes

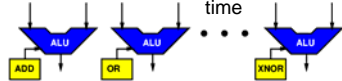
- SIMD Array (microprocessors)
  - Instruction/cycle
  - share instruction across array of PEs
  - uniform operation in space
  - operation variance in time



Caltech CS184 Winter2005 -- DeHon

## Two Extremes

- SIMD Array (microprocessors)
  - Instruction/cycle
  - share instruction across array of PEs
  - 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 CS184 Winter2005 -- DeHon

35

## Placing Architectures

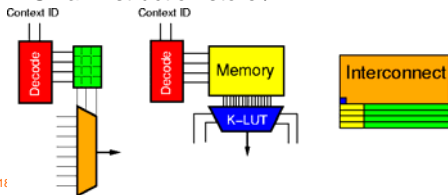
- What programmable architectures (organizations) are you familiar with?

Caltech CS184 Winter2005 -- DeHon

36

## Hybrids

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



Caltech CS184

37

## Architecture Instruction Taxonomy

Control Threads (PCs)		Instruction Depth		Granularity		Architecture/Examples	
$m$	$c$	$n$	$w$	$n_c$	$w_c$		
0	0	n/a				Hardwired Functional Unit (e.g. ECC/EDC Unit, FP MPV)	
				1	1	FPGA	
				$n$	$w$	Reconfigurable ALUs	
				$n_c$	$w_c$	Bitwise SIMD	
	1	$c$	$w$			Traditional Processors	
		$c$	$w$			Vector Processors	
		$n$	8	16		DPGA	
		$c$	$w$			VLIW	
$m$	$n$	1	1			HSRA/SCORE	
	1	$c$	$n_c$	$w$		MSIMD	
		$c$	1			VEGA	
$m$	1	8	16			PADDI-2	
	$c$	$w$				MIMD (traditional)	

Caltech CS184 Winter2005 -- DeHon

38

## Instruction Message

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

Caltech CS184 Winter2005 -- DeHon

39

## Quotes

- *If it can't be expressed in figures, it is not science; it is opinion.* -- Lazarus Long

Caltech CS184 Winter2005 -- DeHon

40

## Modeling

- Why do we model?

Caltech CS184 Winter2005 -- DeHon

41

## Motivation

- Need to understand
  - How costly (big) is a solution
  - How compare to alternatives
  - Cost and benefit of flexibility

Caltech CS184 Winter2005 -- DeHon

42

## What we really want:

- Complete implementation of our application
- For each architectural alternatives
  - In same implementation technology
  - w/ multiple area-time points

## Reality

- Seldom get it packaged that nicely
  - much work to do so
  - technology keeps moving
- Deal with
  - estimation from components
  - technology differences
  - few area-time points

## Modeling Instruction Effects

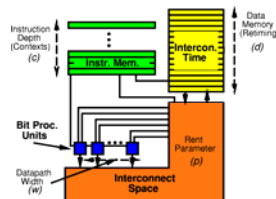
- Restrictions from “ideal” save area
- Restriction from “ideal” limits usability (yield) of PE
  - area model
  - utilization/yield model
- Want to understand effects
  - area model
  - utilization/yield model

## Efficiency/Yield Intuition

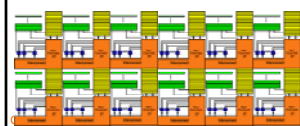
- What happens when
  - Datapath is too wide?
  - Datapath is too narrow?
  - Instruction memory is too deep?
  - Instruction memory is too shallow?

## Computing Device

- Composition
  - Bit Processing elements
  - Interconnect: space
  - Interconnect: time
  - Instruction Memory

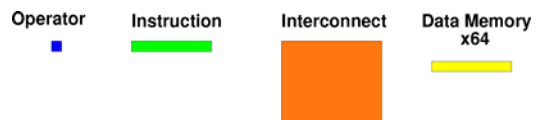


Tile together to build device



## Relative Sizes

- Bit Operator  $10\text{-}20K\lambda^2$
- Bit Operator Interconnect  $500K\text{-}1M\lambda^2$
- Instruction (w/ interconnect)  $80K\lambda^2$
- Memory bit (SRAM)  $1\text{-}2K\lambda^2$





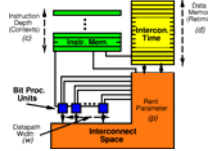
## Model Area

$$A_{bit\_elm} = A_{fixed} + \underbrace{N_{SW}(N_p, w, p) \cdot A_{SW}}_{\text{interconnect}}$$

$$+ \left(\frac{c}{w}\right) \cdot n_{ibits} \cdot A_{mem\_cell}$$

$$+ d \cdot A_{mem\_cell}$$

instruction memory  
retiming memory



Caltech CS184 Winter2005 -- DeHon

49

## Calibrate Model

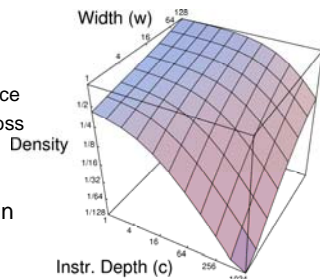
FPGA	model $w = 1, d = c = 1, k = 4$	$880K\lambda^2$
	Xilinx 4K	$630K\lambda^2$
	Altera 8K	$930K\lambda^2$
SIMD	model $w = 1000, c = 0, d = 64, k = 3$	$170K\lambda^2$
	Abacus	$190K\lambda^2$
Processor	model $w = 32, d = 32, c = 1024, k = 2$	$2.6M\lambda^2$
	MIPS-X	$2.1M\lambda^2$

Caltech CS184 Winter2005 -- DeHon

50

## Peak Densities from Model

- Only 2 of 4 parameters
  - small slice of space
  - 100x density across
- Large difference in peak densities
  - large design space!



Caltech CS184 Winter2005 -- DeHon

51

## Efficiency

- What do we want to maximize?
  - Useful work per unit silicon
  - (not potential/peak work)
- Yield Fraction / Area
- (or minimize (Area/Yield) )

Caltech CS184 Winter2005 -- DeHon

52

## Efficiency

- For comparison, look at relative efficiency to ideal.
- Ideal = architecture exactly matched to application requirements
- Efficiency =  $A_{ideal}/A_{arch}$
- $A_{arch} = \text{Area Op/Yield}$

Caltech CS184 Winter2005 -- DeHon

53

## Efficiency Calculation

$$\text{Efficiency} = \frac{A_{matched\_arch}}{A_{arch}}$$

E.g.

If  $w_{task} > w_{arch}$ :

$$\text{Efficiency} = \frac{w_{task} \times A_{bit\_elm}|_{w=w_{task}}}{\left[\frac{w_{task}}{w_{arch}}\right] \times w_{arch} \times A_{bit\_elm}|_{w=w_{arch}}}$$

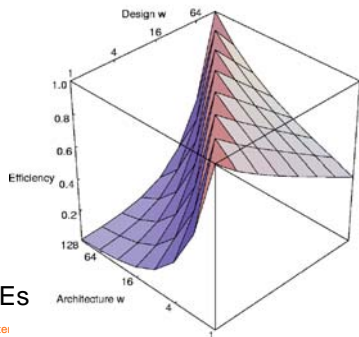
If  $w_{task} < w_{arch}$ :

$$\text{Efficiency} = \frac{w_{task} \times A_{bit\_elm}|_{w=w_{task}}}{w_{arch} \times A_{bit\_elm}|_{w=w_{arch}}}$$

Caltech CS184 Winter2005 -- DeHon

54

## Efficiency: Width Mismatch



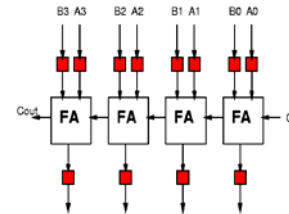
$c=1$ ,  
16K PEs

Caltech CS184 Winter

55

## Path Length

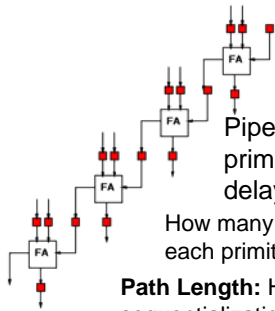
- How many primitive-operator delays before can perform next operation?
  - Reuse the resource



Caltech CS184 Winter2005 -- DeHon

56

## Reuse



Pipeline and reuse at primitive-operator delay level.

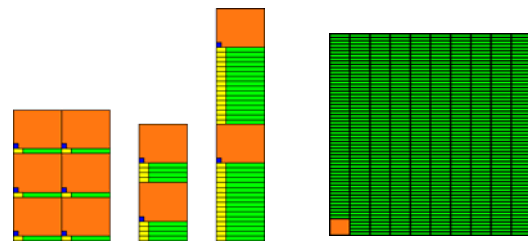
How many times can I reuse each primitive operator?

**Path Length:** How much sequentialization is allowed (required)?

57

Caltech CS184 Winter2005 -- DeHon

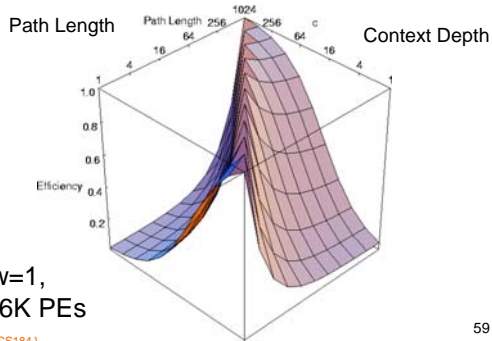
## Context Depth



Caltech CS184 Winter2005 -- DeHon

58

## Efficiency with fixed Width

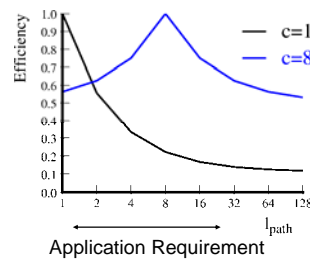


$w=1$ ,  
16K PEs

Caltech CS184 Winter2005 -- DeHon

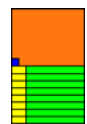
59

## Ideal Efficiency (different model)



- Two resources here:
- active processing elements
  - operation description/state

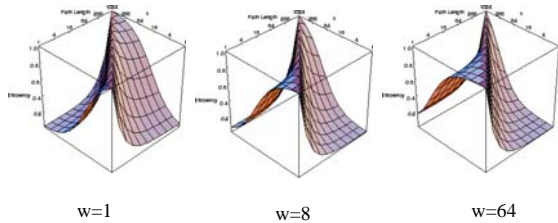
Applications need in different proportions.



Robust point:  $c \cdot A_{ctx} = A_{base}$

Caltech CS184 Winter2005 -- DeHon

## Robust Point depend on Width



w=1

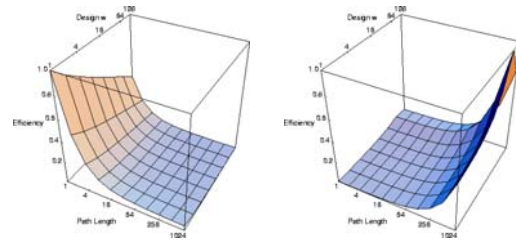
w=8

w=64

Caltech CS184 Winter2005 -- DeHon

61

## Processors and FPGAs



FPGA  
c=d=1, w=1, k=4

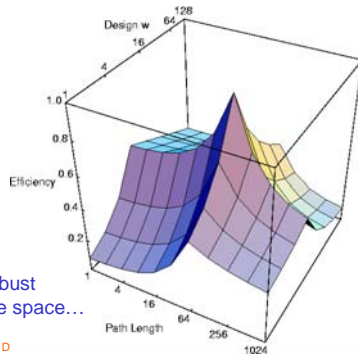
"Processor"  
c=d=1024, w=64, k=2

Caltech CS184 Winter2005 -- DeHon

62

## Intermediate Architecture

w=8  
c=64  
16K PEs



Hard to be robust  
across entire space...

Caltech CS184 Winter2005 -- D

63

## Caveats

- Model abstracts away many details which are important
  - interconnect (day 12--17)
  - control (day 21)
  - specialized functional units (next time)
- Applications are a heterogeneous mix of characteristics

Caltech CS184 Winter2005 -- DeHon

64

## Modeling Message

- Architecture space is **huge**
- Easy to be very inefficient
- Hard to pick one point robust across entire space
- Why we have so many architectures?

Caltech CS184 Winter2005 -- DeHon

65

## General Message

- Parameterize architectures
- Look at continuum
  - costs
  - benefits
- Often have competing effects
  - leads to maxima/minima

Caltech CS184 Winter2005 -- DeHon

66

## Big Ideas [MSB Ideas]

- Basic elements of a programmable computation
  - Compute
  - Interconnect
    - (space and time, outside system [IO])
  - Instructions
- Instruction resources can be significant
  - dominant/limiting resource

Caltech CS184 Winter2005 -- DeHon

67

## Big Ideas [MSB Ideas]

- Applications typically have structure
- Exploit this structure to reduce resource requirements
- Architecture is about understanding and exploiting structure and costs to reduce requirements

Caltech CS184 Winter2005 -- DeHon

68

## Big Ideas [MSB-1 Ideas]

- Two key functions of memory
  - retiming
  - instructions
    - description of computation

Caltech CS184 Winter2005 -- DeHon

69

## Big Ideas [MSB Ideas]

- Instruction organization induces a design space (taxonomy) for programmable architectures
- Arch. structure and application requirements mismatch  $\Rightarrow$  inefficiencies
- Model  $\Rightarrow$  visualize efficiency trends
- Architecture space is huge
  - can be very inefficient
  - need to learn to navigate

Caltech CS184 Winter2005 -- DeHon

70