

# CS184a: Computer Architecture (Structure and Organization)

Day 7: January 24, 2003  
Instruction Space



Caltech CS184 Winter2003 -- DeHon

## Previously

- Computing Requirements
  - Compute
  - Interconnect
    - space
    - time
  - Instructions
- VLSI Scaling

Caltech CS184 Winter2003 -- DeHon

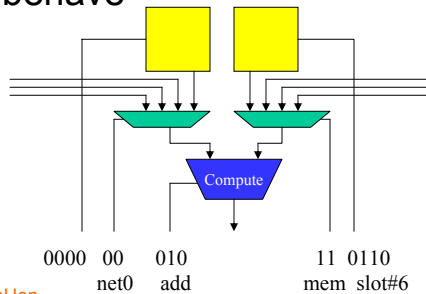
# Today

- Instructions
  - Requirements
  - Taxonomy
  - Model Architecture
    - implied costs
    - gross application characteristics

# Instruction Taxonomy

# Instructions

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



# Focus on Instructions

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

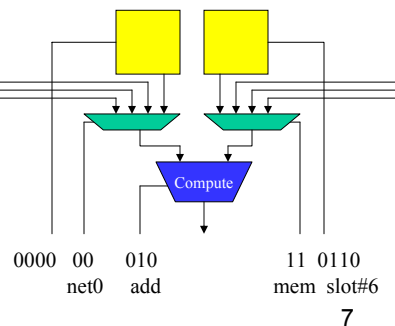
# Terminology

- **Primitive Instruction (*pinst*)**

- Collection of bits which tell a single bit-processing element what to do

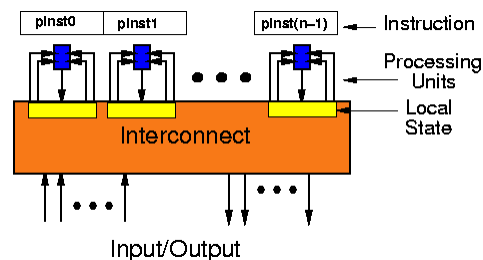
- Includes:

- select **compute** operation
- input sources in space
  - (**interconnect**)
- input sources in time
  - (**retiming**)



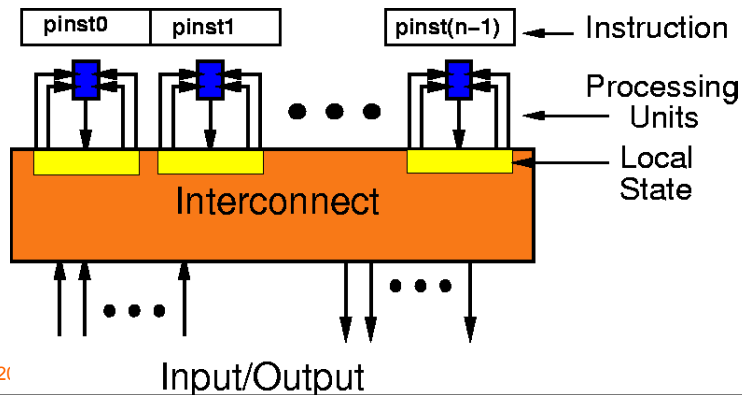
# Computational Array Model

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



## “Ideal” Instruction Control

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



Caltech CS184 Winter20

## “Ideal” Instruction Distribution

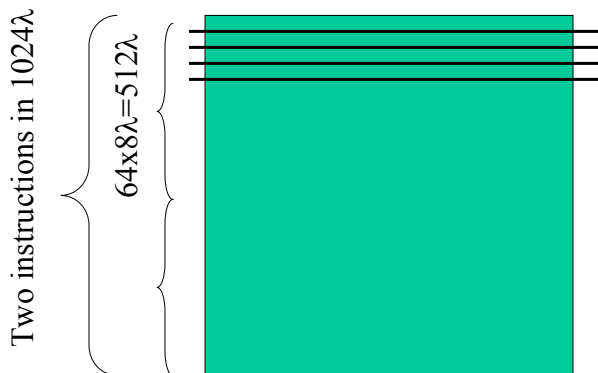
- Why don't we do this?

Caltech CS184 Winter2003 -- DeHon

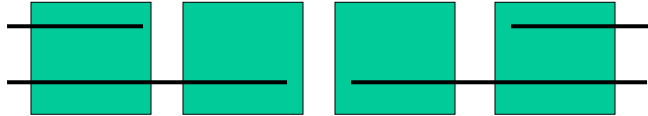
# “Ideal” Instruction Distribution

- **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  $\sim 1.2K\lambda^2$

# Instruction Distribution

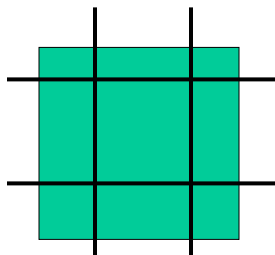


# Instruction Distribution



Distribute from both sides =  $2x$

# Instruction Distribution



Distribute X and Y =  $2x$

## 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) =  $2 \times$
- Two complete metal layers =  $2 \times$
  
- $\Rightarrow$  8 instructions / PE Side

## 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



## Instruction Distribution

- Beyond 64 PE, instruction bandwidth dictates PE size

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

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

- As we 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

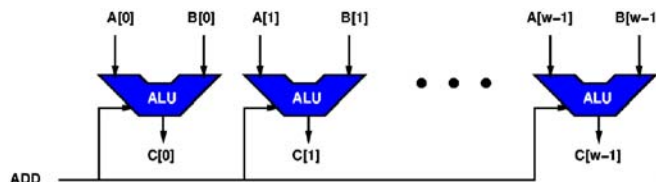
- Instruction requirements *could* dominate array size.
- Standard architecture trick:
  - Look for **structure** to exploit in “typical computations”

# Typical Structure?

- What structure do we usually expect?

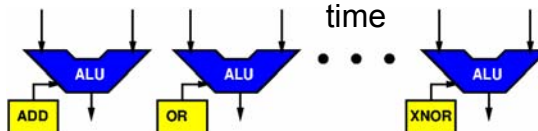
# Two Extremes

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



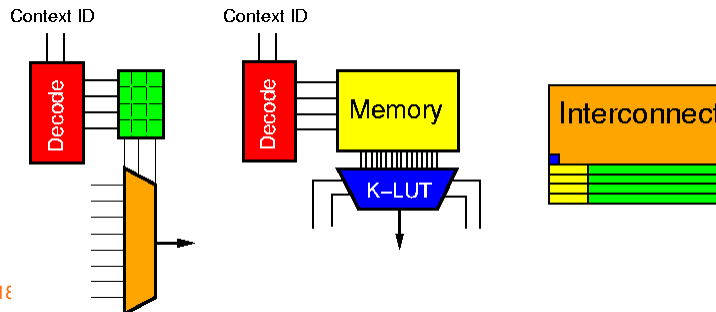
# Two Extremes

- |  |  |
|--|--|
| <ul style="list-style-type: none"> <li>• SIMD Array (microprocessors)                     <ul style="list-style-type: none"> <li>– Instruction/cycle</li> <li>– share instruction across array of PEs</li> <li>– uniform operation in space</li> <li>– operation variance in time</li> </ul> </li> </ul> | <ul style="list-style-type: none"> <li>• FPGA                     <ul style="list-style-type: none"> <li>– Instruction/PE</li> <li>– assume temporal locality of instructions (same)</li> <li>– operation variance in space</li> <li>– uniform operations in time</li> </ul> </li> </ul> |
|--|--|



# Hybrids

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



# Architecture Instruction Taxonomy

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

# 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

# Quotes

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

# Modeling

- Why do we model?

# Motivation

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

## 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
  
- Want to understand effects
  - area model
  - utilization/yield model

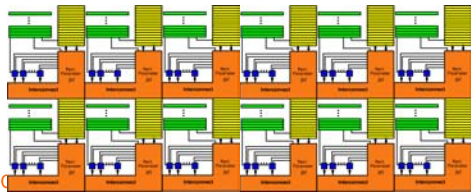
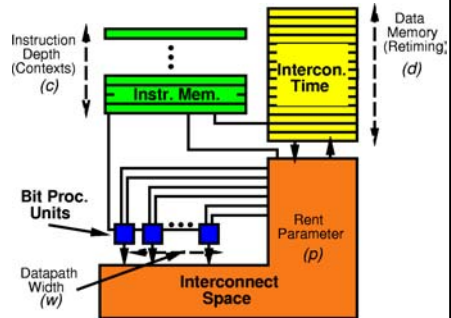
## Efficiency/Yield Intuition

- 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

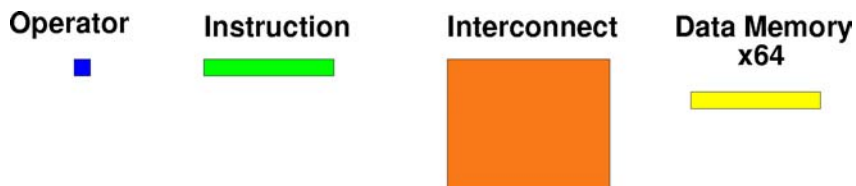


Tile together to build device

33

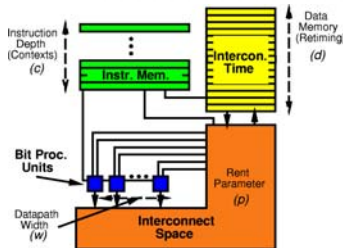
## Relative Sizes

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



## Model Area

$$A_{bit\_elm} = A_{fixed} + \underbrace{N_{SW}(N_p, w, p) \cdot A_{SW}}_{\text{interconnect}} + \underbrace{\left(\frac{c}{w}\right) \cdot n_{ibits} \cdot A_{mem\_cell}}_{\text{instruction memory}} + \underbrace{d \cdot A_{mem\_cell}}_{\text{retiming memory}}$$



Caltech CS184 Winter2003 -- DeHon

35

## Calibrate Model

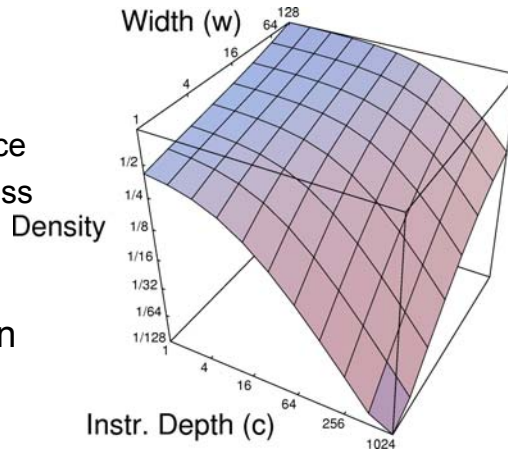
|                        |   |                         |
|------------------------|---|-------------------------|
| <b>FPGA</b>            | <b>model</b> $w = 1, d = c = 1, k = 4$        | <b>880K</b> $\lambda^2$ |
|                        | <b>Xilinx 4K</b>                              | <b>630K</b>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

Caltech CS184 Winter2003 -- DeHon

36

## Peak Densities from Model

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



Caltech CS184 Winter2003 -- DeHon

37

## 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 Winter2003 -- DeHon

38

# Efficiency

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

# Efficiency Calculation

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

E.g.

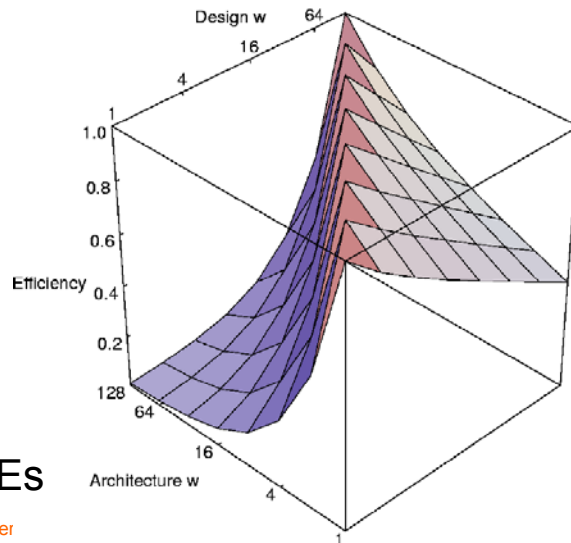
If  $w_{task} > w_{arch}$ :

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

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

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

# Efficiency: Width Mismatch



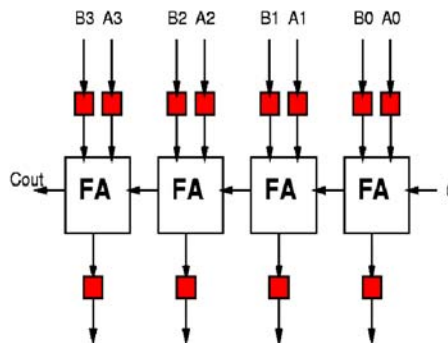
$c=1,$   
16K PEs

Caltech CS184 Winter

41

# Path Length

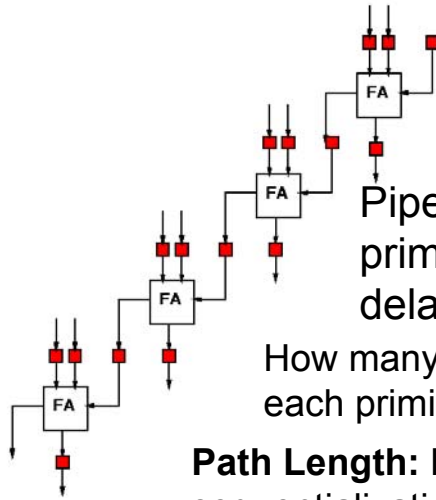
- How many primitive-operator delays before can perform next operation?
  - Reuse the resource



Caltech CS184 Winter2003 -- DeHon

42

# Reuse



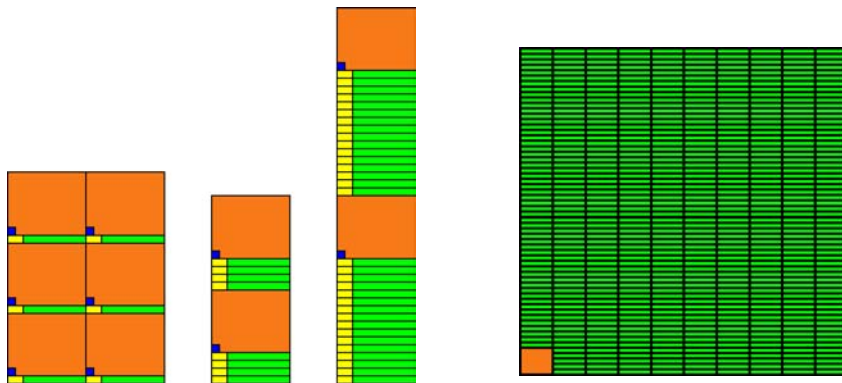
Pipeline and reuse at primitive-operator delay level.

How many times can I reuse each primitive operator?

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

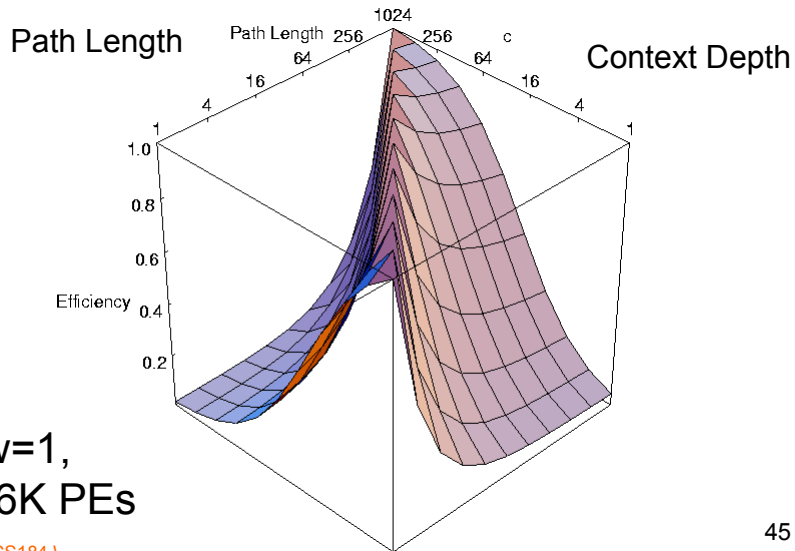
43

# Context Depth



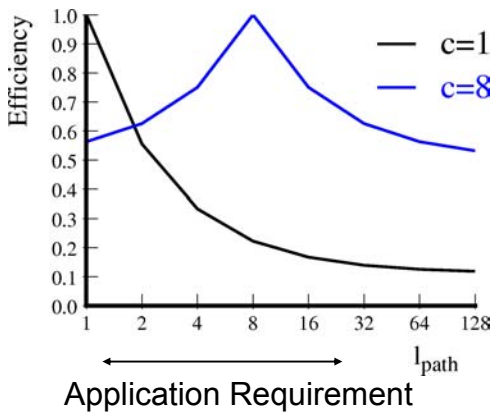
44

# Efficiency with fixed Width



45

# Ideal Efficiency (different model)



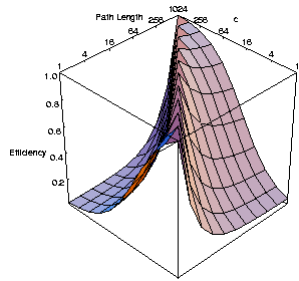
- Two resources here:
- active processing elements
  - operation description/state

Applications need in different proportions.

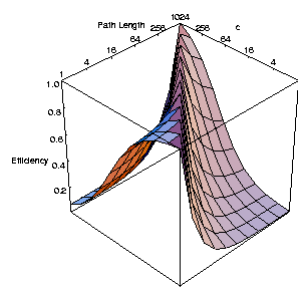


Robust point:  $c \cdot A_{ctx} = A_{base}$

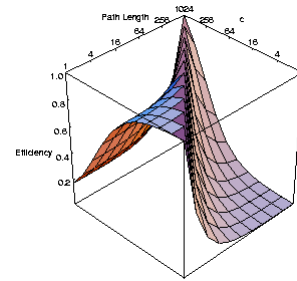
# Robust Point depend on Width



w=1

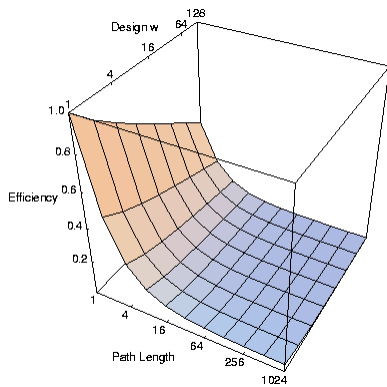


w=8

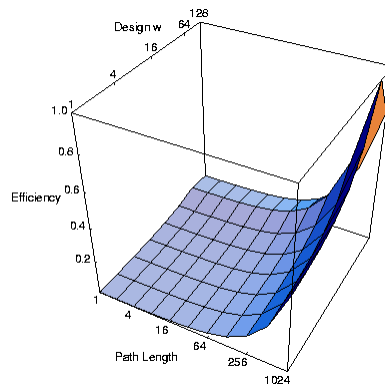


w=64

# Processors and FPGAs



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

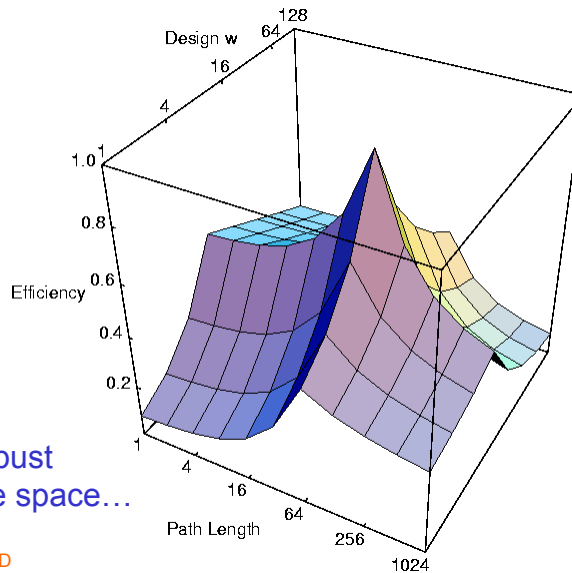


"Processor"  
c=d=1024, w=64, k=2



# Intermediate Architecture

w=8  
c=64  
16K PEs



Hard to be robust  
across entire space...

Caltech CS184 Winter2003 -- D

49

## Caveats

- Model abstracts away many details which are important
  - interconnect (day 11--16)
  - control (day 20)
  - specialized functional units (next time)
- Applications are a heterogeneous mix of characteristics

Caltech CS184 Winter2003 -- DeHon

50

## 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?

## General Message

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

## Big Ideas [MSB Ideas]

- Instruction resources can be significant
  - dominant/limiting resource
- Applications typically have structure
- Exploit this structure to reduce resource requirements
- Architecture is about understanding and exploiting structure and costs to reduce requirements

## 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