# ESE534:
# Computer Organization

Day 11:  March 1, 2010
Instruction Space Modeling

Penn

---

## Last Time

- Instruction Requirements
- Instruction Space

---

## Architecture Instruction Taxonomy

| Control Threads (PCs) | | | | | |
|---|---|---|---|---|---|
| | *pinsts* per Control Thread | | | | |
| | | Instruction Depth | | | |
| | | | Granularity | | |
| | | | | | Architecture/Examples |
| 0 | 0 | 0 | n/a | | Hardwired Functional Unit (*e.g.* ECC/EDC Unit, FP MPY) |
| | | | 1 | | FPGA |
| | $n$ | 1 | $w$ | | Reconfigurable ALUs |
| | | | | $n_v \cdot 1$ | Bitwise SIMD |
| | 1 | $c$ | $w$ | | Traditional Processors |
| 1 | | | | $n_v \cdot w$ | Vector Processors |
| | | $c$ | 1 | | DPGA |
| | $n$ | 8 | 16 | | PADDI |
| | | $c$ | $w$ | | 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) |

---

## Architecture Taxonomy

| PCs | Pints/PC | depth | width | Architecture |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | FPGA |
| 1 | 1 | 1024 | 32 | Scalar Processor (RISC) |
| 1 | N | D | W | VLIW (superscalar) |
| 1 | 1 | Small | W*N | SIMD, GPU, Vector |
| N | 1 | D | W | MIMD |
| 4 | 4 | 2048 | 64 | Quad core |

---

## Today

- Model Architecture from Instruction Parameters
  - implied costs
  - gross application characteristics

---

## 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 is a solution
    - Big, slow, hot, energy hungry….
  - 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
- We must deal with
  - estimation from components
  - technology differences
  - few area-time points

## Modeling Instruction Effects

- Restrictions from "ideal"
  - +save area and energy
  - limit usability (yield) of PE
    - May cost more energy, area in the end…

- Want to understand effects
  - area model [today] (energy model on HW5)
  - utilization/yield model

## Preclass

- Energies?
- 16-bit on 32-bit?
  - Sources of inefficiency?
- 8-bit operations per 16-bit operation?
- 16-bit on 8-bit?
  - Sources of inefficiency?

## Efficiency/Yield Intuition

- What happens when
  - Datapath is too wide?
  - Datapath is too narrow?
  - Instruction memory is too deep?
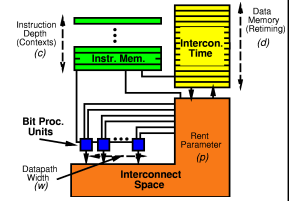  - Instruction memory is too shallow?
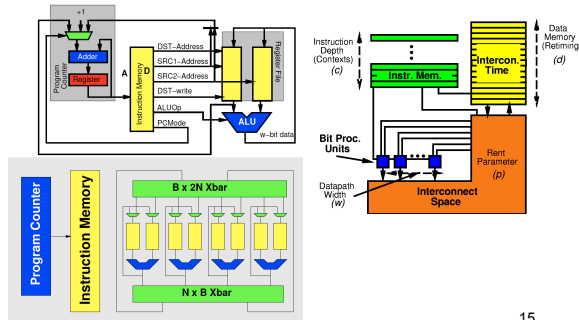
13

## Computing Device

- Composition
  - Bit Processing elements
  - Interconnect: space
  - Interconnect: time
  - Instruction Memory
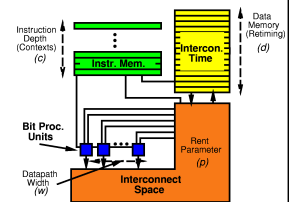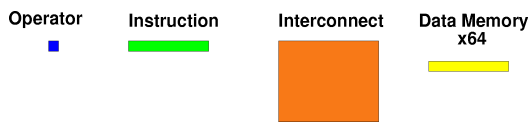
Tile together to build device

14

## Computing Device

15

## Computing Device

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

Tile together to build device

16

## Relative Sizes

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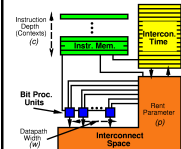- Instruction (w/ interconnect)    $80\text{K}\lambda^2$
- Memory bit (SRAM)       $1\text{-}2\text{K}\lambda^2$

| Operator | Instruction | Interconnect | Data Memory x64 |
|---|---|---|---|

17

## Model Area

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

18

3

## Architectures Fall in Space



Instruction depth (c) vs Word width (w)

VEGA, MIPS-X, T0 Vector, PADDI, DPGA, vDSP, FPGA

19

## 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$\lambda^2$** |

20

## Peak Densities from Model

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



Width (w), Density, Instr. Depth (c)
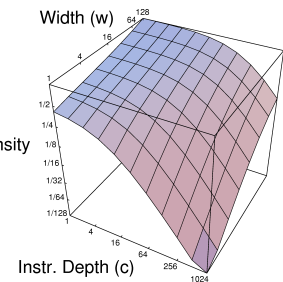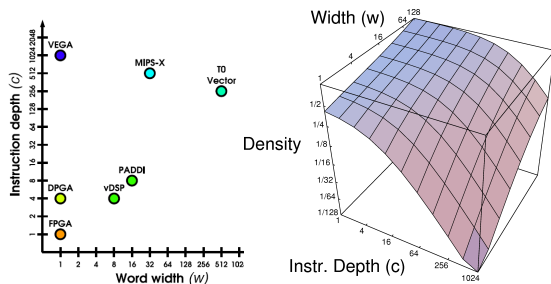
21

## Peak Densities from Model

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



Width (w), Density, Instr. Depth (c)

22

## Architectural parameters → Peak Densities



Instruction depth (c) vs Word width (w)

VEGA, MIPS-X, T0 Vector, PADDI, DPGA, vDSP, FPGA

Width (w), Density, Instr. Depth (c)

23

## Efficiency

- What do we really want to maximize?
  - Not peak, "guaranteed not to exceed" performance, but…
  - Useful work per unit silicon [per Joule]

- Yield Fraction / Area
- (or minimize (Area/Yielded performance) )

24

4

## Efficiency

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

## Width Mismatch Efficiency Calculation
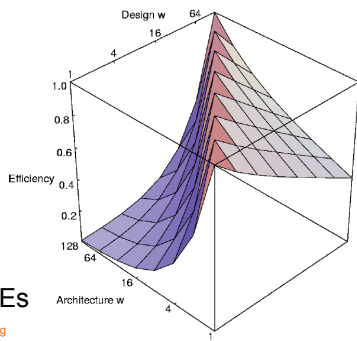
$$E = \frac{Area(Task-on-matched-Architecture)}{Area(Task-on-this-Architecture)}$$

$$E = \frac{W_{task} \times A_{bitelm|w=w_{task}}}{W_{arch} \times \left\lceil \dfrac{W_{task}}{W_{arch}} \right\rceil \times A_{bitelm|w=w_{arch}}}$$

## Efficiency: Width Mismatch



c=1,
16K PEs

## Efficiency for Preclass

$$E = \frac{Energy(Task-on-matched-Architecture)}{Energy(Task-on-this-Architecture)}$$

- Efficiency of 16-bit on 32-bit arch?
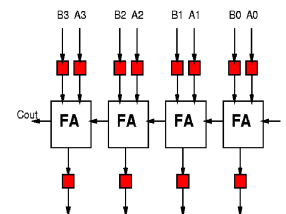- Efficiency of 16-bit on 8-bit arch?

## Application vs. Architecture

- $W_{task}$ vs. $W_{arch}$
- Path Length vs. Context Depth

## Path Length

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

5

## 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)?

*E.g.* Want meet 30ns real time rate with 1.5ns cycle time, can afford to issue 15 sequential ops.

31

## Context (Instruction) Depth

32

## Efficiency with fixed Width



Path Length

Context Depth

Efficiency

w=1, 16K PEs

33

## Ideal Efficiency (different model)



— c=1
— c=8

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

Applications need in different proportions.

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

## Breaking News (3/1/10)



— c=1
— c=8

Application Requirement

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

Note release from startup Tabula today:
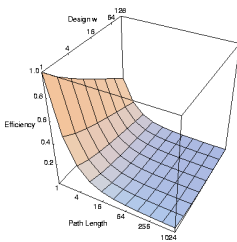http://www.eetimes.com/rss/showArticle.jhtml?articleID=223100915&cid=RSSfeed_eetimes_newsRS

35

## Robust Point depend on Width



w=1          w=8          w=64

36

6

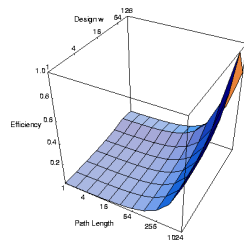## Processors and FPGAs
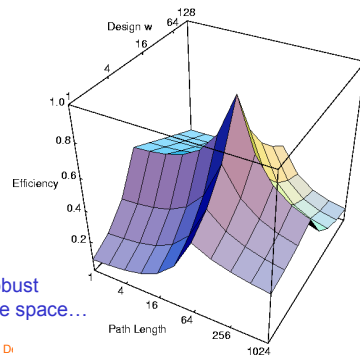### (architecture vs. two application axes)



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

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

37

## Intermediate Architecture

w=8
c=64
16K PEs



Hard to be robust across entire space…

38

## Caveats

- Model abstracts away many details that are important
  - interconnect (day 17--20)
  - control (day 23)
  - specialized functional units (day 14)
- Applications are a heterogeneous mix of characteristics

39

## 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?

40

## General Message

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

41

## Admin

- Should now have all background for HW5
  - Problem 2 similar (looking for robust point)
  - Different
    - Interconnect parameter
    - Energy efficiency
- Reading for Wednesday on Blackboard

42

7

## 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

43

## 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

44