# ESE680-002 (ESE534): Computer Organization

Day 9:  February 7, 2007
Instruction Space Modeling

---

# Last Time

- Instruction Requirements
- Instruction Space

2

---

# 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) |

3

---

# Today

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

4

---

# Quotes

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

5

---

# Modeling

- Why do we model?

6

1

## 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?
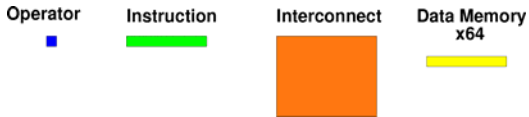  - Instruction memory is too shallow?

## Computing Device

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



Tile together to build device

12

## Relative Sizes

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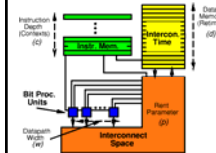- Memory bit (SRAM)             $1\text{-}2K\lambda^2$



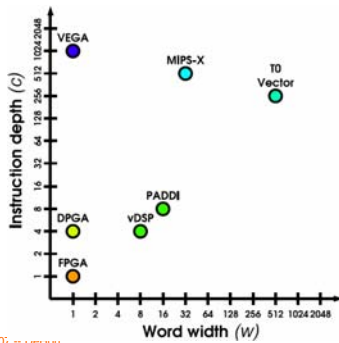| Operator | Instruction | Interconnect | Data Memory x64 |

13

## Model Area

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

14

## Architectures Fall in Space

15

## 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 | $2.1M\lambda^2$ |

16

## Peak Densities from Model

$$A_{bit\_elm} = A_{fixed} + \underbrace{N_{SW}(Np, w, p) \cdot A_{SW}}_{\textbf{interconnect}}$$
$$+ \underbrace{\left(\frac{c}{w}\right) \cdot n_{ibits} \cdot A_{mem\_cell}}_{\textbf{instruction memory}}$$
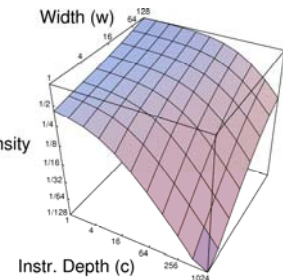$$+ \underbrace{d \cdot A_{mem\_cell}}_{\textbf{retiming memory}}$$
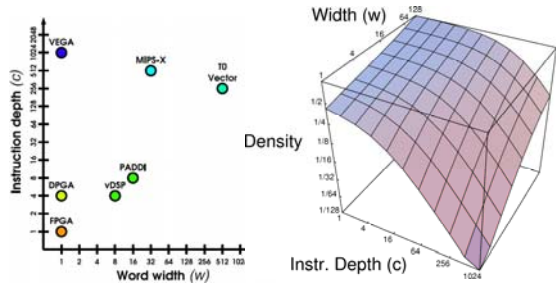
17

## Peak Densities from Model

- Only 2 of 4 parameters
  - small slice of space
  - 100× density across

- Large difference in peak densities
  - large design space!

18

3

## Peak Densities from Model

19

---

## Efficiency

- What do we want to maximize?
  - Useful work per unit silicon
  - (not potential/peak work)

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

20

---

## 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

21

---

## Width Mismatch Efficiency Calculation

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

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

22

---

## Efficiency:  Width Mismatch



c=1,
16K PEs

23

---

## Path Length

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

24

4

## 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)?

25

## Context Depth

26

## Efficiency with fixed Width



w=1,
16K PEs

27

## 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

29

## Processors and FPGAs



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

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

30

5

## Intermediate Architecture

w=8
c=64
16K PEs



Hard to be robust
  across entire space…

31

---

## Caveats

- Model abstracts away many details which are important
  - interconnect (day 13--18)
  - control      (day 24)
  - specialized functional units (next time)
- Applications are a heterogeneous mix of characteristics

32

---

## 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?

33

---

## General Message

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

34

---

## Admin

- Assignment 4 out today
  - Did push back due dates for 4 and 5
- Reading for Monday on web
  - Supplemental from this month TRCAD

35

---

## 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

36

---

# 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

37