

ESE534: Computer Organization

Day 11: February 20, 2012
Instruction Space Modeling



Last Time

- Instruction Requirements
- Instruction Space

Architecture Taxonomy

PCs	Pints/PC	depth	width	Architecture
0	N	1	1	FPGA
1	N (48,640)	8	1	Tabula ABAX (A1EC04)
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
16	1 (4?)	2048	64	16-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

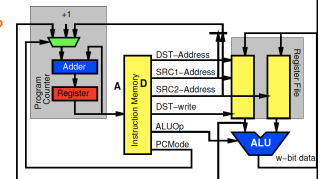
- Restrictions from “ideal”
 - + save area and energy
 - limit usability (yield) of PE
 - May cost more energy, area in the end...
- Want to understand effects
 - area model [today] (energy model on HW5)
 - utilization/yield model

Preclass

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

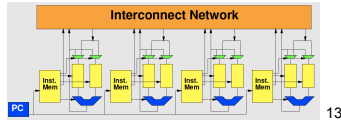
Efficiency/Yield Intuition

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



Efficiency/Yield Intuition

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

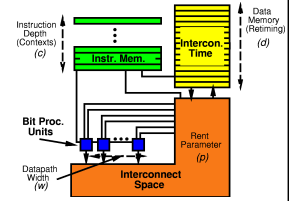


Penn ESE534 Spring 2012 -- DeHon

13

Computing Device

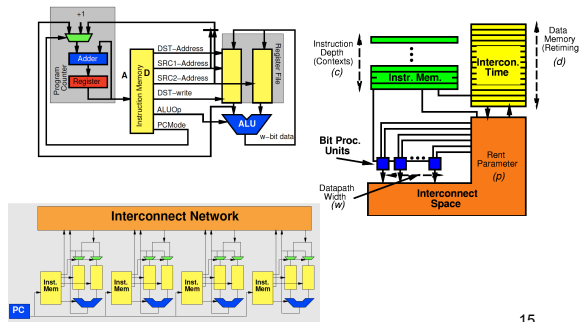
- Composition
 - Bit Processing elements
 - Interconnect: space
 - Interconnect: time
 - Instruction Memory



Tile together to build device

14

Computing Device

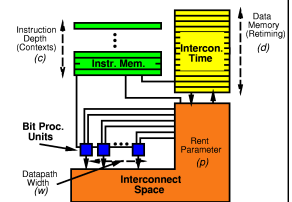


Penn ESE534 Spring 2012 -- DeHon

15

Computing Device

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

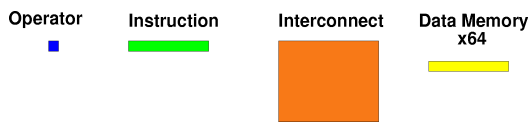


Tile together to build device

16

Relative Sizes

- Bit Operator 3-5KF²
- Bit Operator Interconnect 200K-250KF²
- Instruction (w/ interconnect) 20KF²
- Memory bit (SRAM) 250-500F²

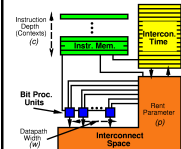


Penn ESE534 Spring 2012 -- DeHon

17

Model Area

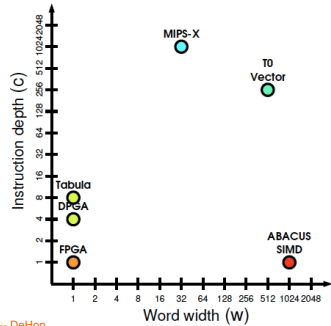
$$A_{bit_elm} = A_{fixed} + \frac{N_{SW}(N_p, w, p)}{interconnect} \cdot A_{SW} + \left(\frac{c}{w}\right) \cdot n_{ibits} \cdot A_{mem_cell} + \frac{d \cdot A_{mem_cell}}{retiming\ memory}$$



Penn ESE534 Spring 2012 -- DeHon

18

Architectures Fall in Space



Penn ESE534 Spring 2012 -- DeHon

19

Calibrate Model

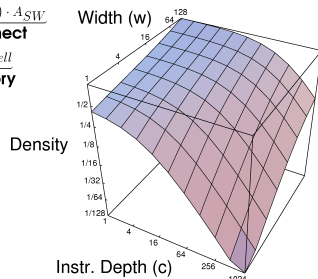
FPGA	model $w = 1, d = c = 1, k = 4$	$880K\lambda^2$
	Xilinx 4K	$630K\lambda^2$
	Altera 8K	$930K\lambda^2$
SIMD	model $w = 1000, c = 0, d = 64, k = 3$	$170K\lambda^2$
	Abacus	$190K\lambda^2$
Processor	model $w = 32, d = 32, c = 1024, k = 2$	$2.6M\lambda^2$
	MIPS-X	$2.1M\lambda^2$

Penn ESE534 Spring 2012 -- DeHon

20

Peak Densities from Model

$$A_{bit_elm} = A_{fixed} + \frac{N_{SW}(N_p, w, d) \cdot ASW}{interconnect} + \left(\frac{c}{w}\right) \cdot n_{ibits} \cdot A_{mem_cell} + \frac{d \cdot A_{mem_cell}}{refining\ memory}$$

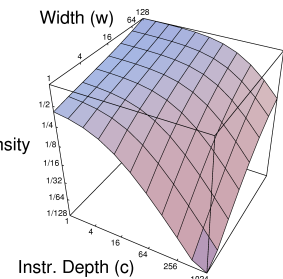


Penn ESE534 Spring 2012 -- DeHon

21

Peak Densities from Model

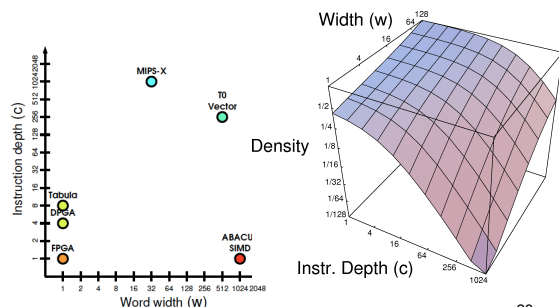
- Only 2 of 4 parameters
 - small slice of space
 - 100x density across
- Large difference in peak densities
 - large design space!



Penn ESE534 Spring 2012 -- DeHon

22

Architectural parameters → Peak Densities



Penn ESE534 Spring 2012 -- DeHon

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))

Penn ESE534 Spring 2012 -- DeHon

24

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

Penn ESE534 Spring 2012 -- DeHon

25

Width Mismatch Efficiency Calculation

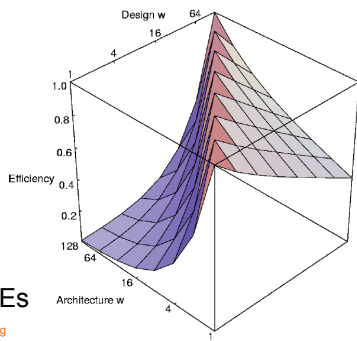
$$E = \frac{\text{Area}(\text{Task-on-matched-Architecture})}{\text{Area}(\text{Task-on-this-Architecture})}$$

$$E = \frac{W_{task} \times A_{bitelm|w=w_{task}}}{W_{arch} \times \left[\frac{W_{task}}{W_{arch}} \right] \times A_{bitelm|w=w_{arch}}}$$

Penn ESE534 Spring 2012 -- DeHon

26

Efficiency: Width Mismatch



Penn ESE534 Spring

27

Efficiency for Preclass

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

- Preclass 6 table

Penn ESE534 Spring 2012 -- DeHon

28

Application vs. Architecture

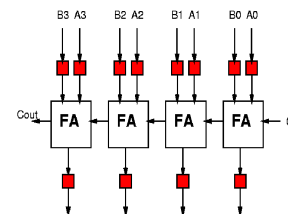
- W_{task} vs. W_{arch}
- Path Length vs. Context Depth

Penn ESE534 Spring 2012 -- DeHon

29

Path Length

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



Penn ESE534 Spring 2012 -- DeHon

30

Reuse

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

Penn ESE534 Spring 2012 -- DeHon

Context (Instruction) Depth

Penn ESE534 Spring 2012 -- DeHon

32

Efficiency with fixed Width

$w=1,$
16K PEs

Penn ESE534 Spring 2012 -- DeHon

33

Ideal Efficiency (different model)

Two resources here:

- active processing elements
- operation description/state

Applications need in different proportions.

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

Penn ESE534 Spring 2012 -- DeHon

Robust Point

- What is Energy Robust Point for preclass model?

Penn ESE534 Spring 2012 -- DeHon

35

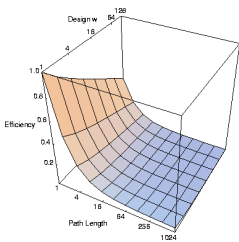
Robust Point depend on Width

$w=1$ $w=8$ $w=64$

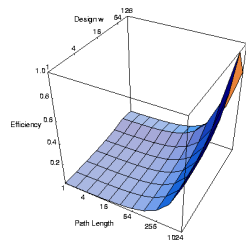
Penn ESE534 Spring 2012 -- DeHon

36

Processors and FPGAs (architecture vs. two application axes)



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



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

Penn ESE534 Spring 2012 -- DeHon

37

Application Needs

- What are common application datawidths?
- What are common application path lengths?

Penn ESE534 Spring 2012 -- DeHon

38

Examples

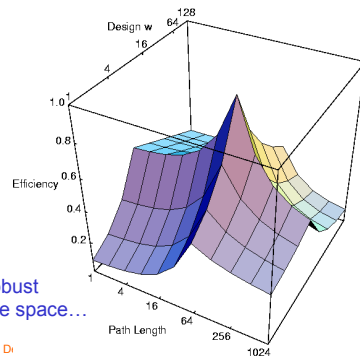
Application	Wapp	Lcritpath	Lpath	Notes
Conway LIFE	1	1	1	Run as fast as possible
Entropy Code	1	1-10	100	100ns memory interface
Video	8	1-6	24	1GHz for 1024x1024 x30 frames/s
Audio	16	1-10	20,000	44KHz for 1GHz
FDTD	35	1-5	1-5	

Penn ESE534 Spring 2012 -- DeHon

39

Intermediate Architecture

$w=8$
 $c=64$
16K PEs



Hard to be robust across entire space...

Penn ESE534 Spring 2012 -- Di

40

Caveats

- Model abstracts away many details that are important
 - interconnect (day 15--18)
 - control (day 22)
 - specialized functional units (day 14)
- Applications are a heterogeneous mix of characteristics

Penn ESE534 Spring 2012 -- DeHon

41

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?

Penn ESE534 Spring 2012 -- DeHon

42

General Message

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

Penn ESE534 Spring 2012 -- DeHon

43

Admin

- Should now have all background for HW5
 - Problem 2 similar (looking for robust point)
 - Different: Interconnect parameter, Energy
- No class Wednesday
- No office hours Tuesday
- Next class Monday
 - Reading online
- HW 6 out
 - 1 and 2 due Friday 2/3
 - Should be able to do 1 now

Penn ESE534 Spring 2012 -- DeHon

44

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

Penn ESE534 Spring 2012 -- DeHon

45

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

Penn ESE534 Spring 2012 -- DeHon

46