

ESE534: Computer Organization

Day 10: February 24, 2010
Computing Requirements and
Instruction Space



Today

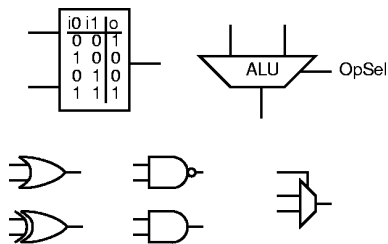
- Computing Requirements
- Instructions
 - Requirements
 - Taxonomy

Computing Requirements (review)

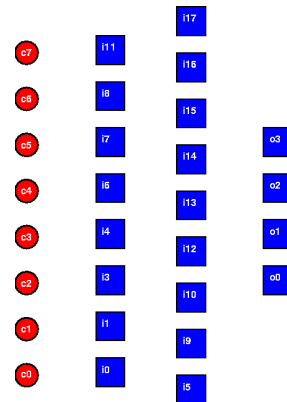
Requirements

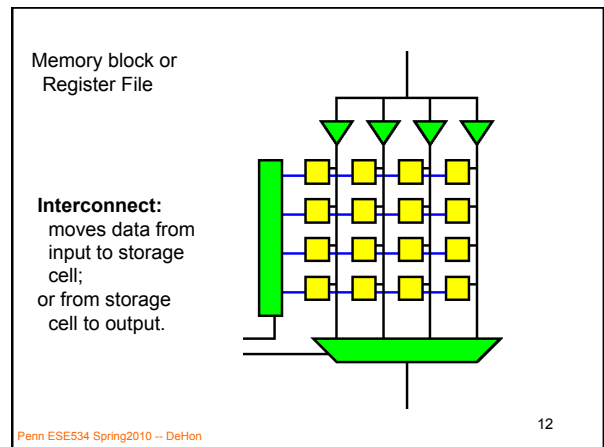
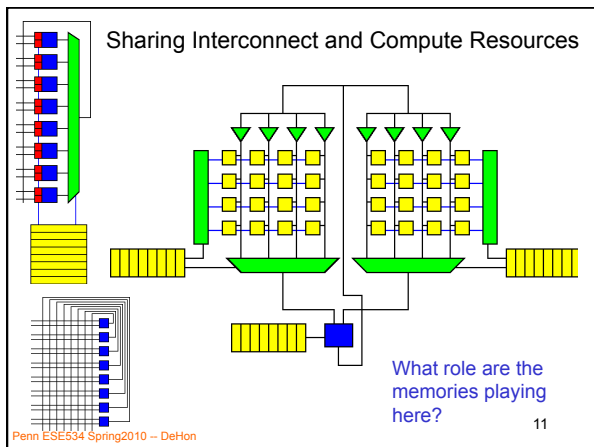
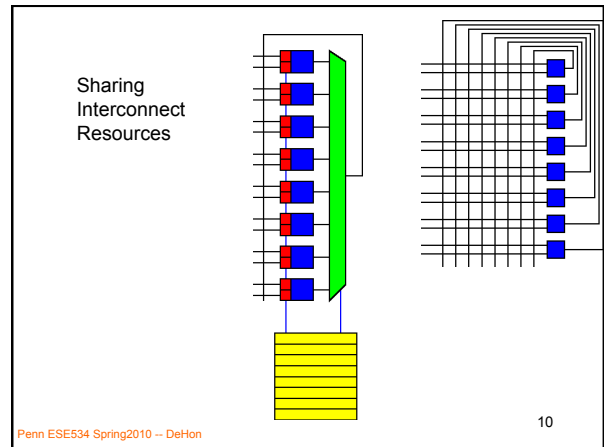
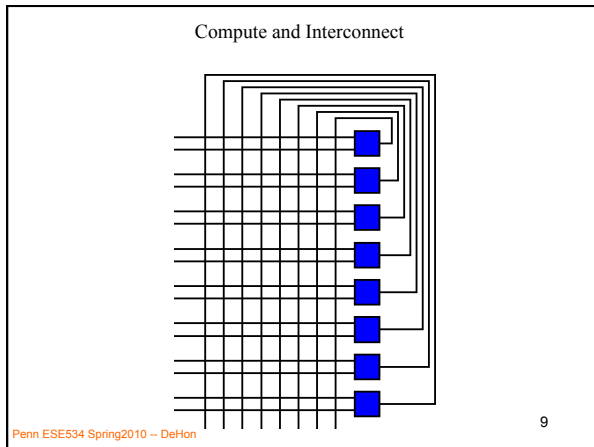
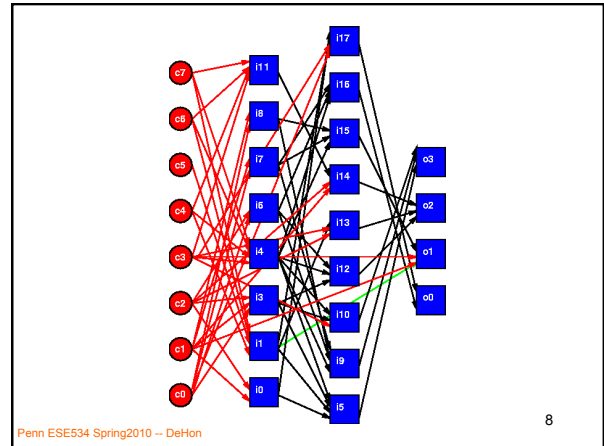
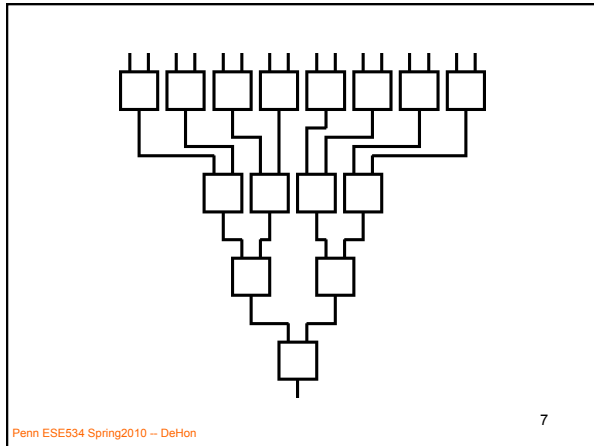
- In order to build a **general-purpose** (*programmable*) computing device, we absolutely must have?

-
-
-
-
-
-

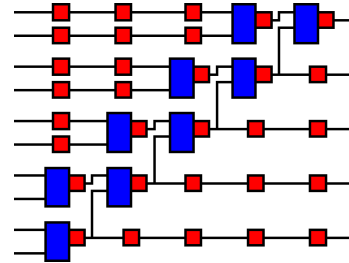
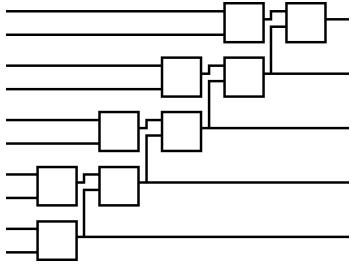


Primitive compute
elements enough?





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



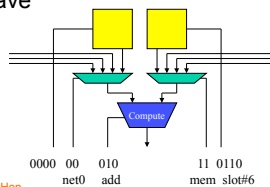
Requirements

- 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
 - Control (e.g. Program Counter)

Instruction Taxonomy

- 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

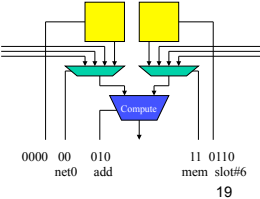
Terminology

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

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

- Includes:

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

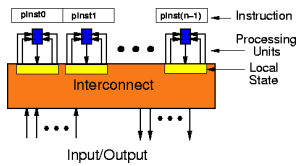


Preclass

- How big is *pinst* for preclass?
 - (problem 1)

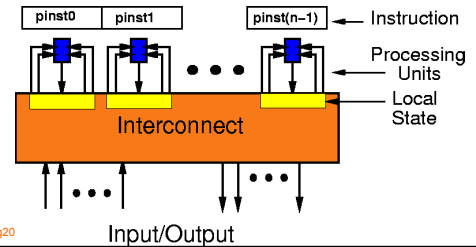
Computational Array Model

- 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



“Ideal” Instruction Distribution

- Why don't we do this?

Preclass

- How many total instruction bits?
 - (preclass 2)
- How many pins on a chip?

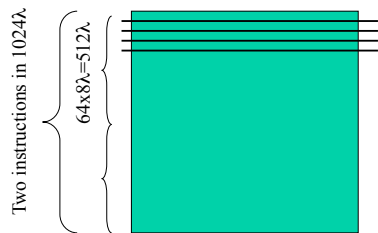
Preclass

- How wide is instruction distribution?
- For $\lambda=20\text{nm}$?

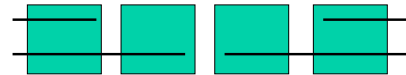
“Ideal” Instruction Distribution

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

Instruction Distribution

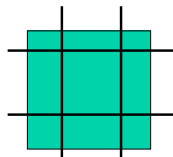


Instruction Distribution



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) = $2x$
- Two complete metal layers = $2x$
- \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

Penn ESE534 Spring2010 -- DeHon

31

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

Penn ESE534 Spring2010 -- DeHon

32

Avoid Instruction BW Saturation?

- How might we avoid this?

Penn ESE534 Spring2010 -- DeHon

33

Instruction Memory Requirements

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

Penn ESE534 Spring2010 -- DeHon

34

Instruction Pragmatics

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

Penn ESE534 Spring2010 -- DeHon

35

Typical Structure?

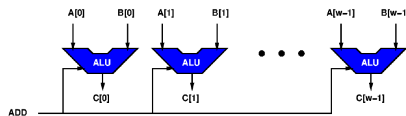
- What structure do we usually expect?

Penn ESE534 Spring2010 -- DeHon

36

Two Extremes

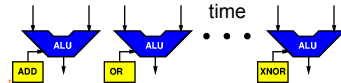
- SIMD Array (microprocessors)
 - Instruction/cycle
 - share instruction across array of PEs
 - uniform operation in space
 - operation variance in time
- SIMD = Single Instruction Multiple Data



Penn ESE534 Spring2010 -- DeHon

Two Extremes

- SIMD Array (microprocessors)
 - Instruction/cycle
 - share instruction across array of PEs
 - uniform operation in space
 - operation variance in time
- FPGA (Field-Programmable Gate Array)
 - Instruction/PE
 - assume temporal locality of instructions (same)
 - operation variance in space
 - uniform operations in time



Penn ESE534 Spring2010 -- DeHon

38

Placing Architectures

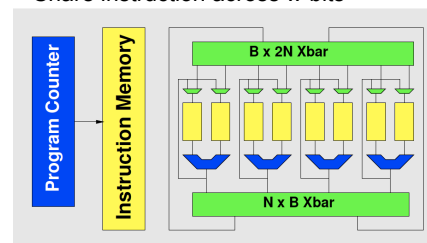
- What programmable architectures (organizations) are you familiar with?

Penn ESE534 Spring2010 -- DeHon

39

Hybrids

- VLIW = Very Long Instruction Word
 - Few *pinsts/cycle*
 - Share instruction across w bits

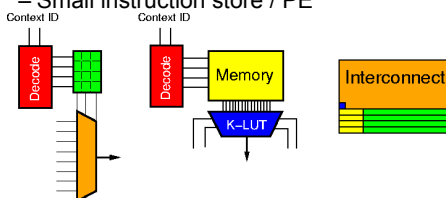


Penn ESE534 Spring2010 -- DeHon

40

Hybrids

- VLIW = Very Long Instruction Word
 - Few *pinsts/cycle*
 - Share instruction across w bits
- DPGA
 - Small instruction store / PE



Penn ESE534 Spring2010 -- DeHon

41

Architectural Differences

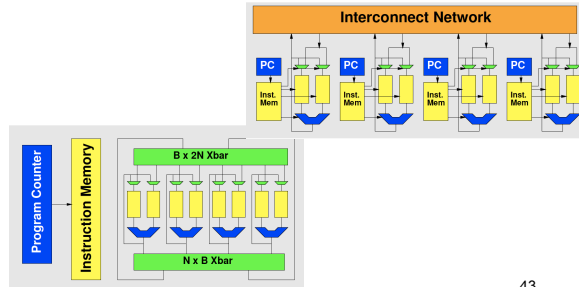
- What differentiates a VLIW from a multicore?
 - *E.g.*
 - 4-issue VLIW vs.
 - 4 single-issue processors

Penn ESE534 Spring2010 -- DeHon

42

Architectural Differences

- What differentiates a VLIW from a multicore?



Penn ESE534 Spring2010 -- DeHon

43

Gross Parameters

- Instruction sharing width
 - SIMD width
 - granularity
- Instruction depth
 - Instructions stored locally per compute element
- pins per control thread
 - E.g. VLIW width

Penn ESE534 Spring2010 -- DeHon

44

Architecture Instruction Taxonomy

Control Threads (PCs)		pins per Control Thread		Instruction Depth		Granularity		Architecture/Examples	
0	0	n/a						Hardwired Functional Unit (e.g. ECC/EDC Unit, FP MPY)	
	n	1	w					FPGA	
	n	1	w					Reconfigurable ALUs	
	1	c	w					Bitwise SIMD	
	1	c	w					Traditional Processors	
	1	c	w					Vector Processors	
	1	c	1					DPGA	
	n	8	16					PADDI	
	m	1	1					VLIW	Superscalar
	m	1	1					HSRA/SCORE	
	1	c	n _c · w					MSIMD	GPUs
	1	c	1					VEGA	
	m	1	8	16				PADDI-2	
	m	1	c	w				MIMD (traditional)	Multicore

Penn ESE534 Spring2010 -- DeHon

45

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

Penn ESE534 Spring2010 -- DeHon

46

Admin

- Reading on blackboard
- HW5
 - Should be able to do all of Problem 1 now
 - Day11/Monday relevant to Problem 2

Penn ESE534 Spring2010 -- DeHon

47

Big Ideas

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

Penn ESE534 Spring2010 -- DeHon

48