

CS184a: Computer Architecture (Structure and Organization)

Day 20: March 3, 2003
Control



Caltech CS184 Winter2003 -- DeHon

Previously

- Looked broadly at instruction effects
- Looked at structural components of computation
 - interconnect
 - compute
 - retiming
- Looked at time-multiplexing

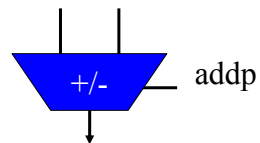
Caltech CS184 Winter2003 -- DeHon

Today

- Control
 - data-dependent operations
- Different forms
 - local
 - instruction selection
- Architectural Issues

Control

- **Control:** That point where the data affects the instruction stream (operation selection)
 - Typical manifestation
 - data dependent branching
 - if (a!=0) OpA else OpB, bne
 - data dependent state transitions
 - new => goto S0
 - else => stay
 - data dependent operation selection



Control

- **Viewpoint:** can have instruction stream sequence without control
 - *i.e.* static/data-independent progression through sequence of instructions is control free
 - $C_0 \rightarrow C_1 \rightarrow C_2 \rightarrow C_0 \rightarrow C_1 \rightarrow C_2 \rightarrow C_0 \rightarrow \dots$
 - Similarly, FSM w/ no data inputs
 - Recall could do multiply w/out branching...

Terminology (reminder)

- **Primitive Instruction (*pinst*)**
 - Collection of bits which tell a bit-processing element what to do
 - Includes:
 - select compute operation
 - input sources in space (interconnect)
 - input sources in time (retiming)
- **Configuration Context**
 - Collection of all bits (*pinsts*) which describe machine's behavior on one cycle

Why?

- Why do we need / want control?
- Static interconnect sufficient?
- Static sequencing?
- Static datapath operations?

Back to “Any” Computation

- Design must handle all potential inputs (computing scenarios)
- Requires sufficient generality
- However, computation for any given input may be *much* smaller than general case.

Two Control Options

1. Local control

- unify choices
 - build all options into spatial compute structure and select operation

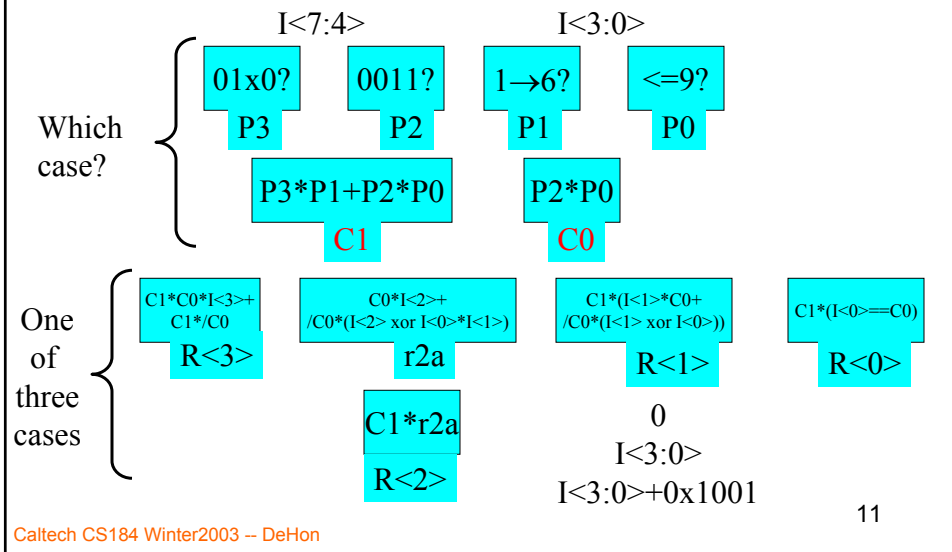
2. Instruction selection

- provide a different instruction (instruction sequence) for each option
- selection occurs when chose which instruction(s) to issue

Example: ASCII Hex→Binary

- If ($c \geq 0x30$ && $c \leq 0x39$)
 - $res = c - 0x30$ // $0x30 = '0'$
- elseif ($c \geq 0x41$ && $c \leq 0x46$)
 - $res = c - 0x41 + 10$ // $0x41 = 'A'$
- elseif ($c \geq 0x61$ && $c \leq 0x66$)
 - $res = c - 0x61 + 10$ // $0x61 = 'a'$
- else
 - $res = 0$

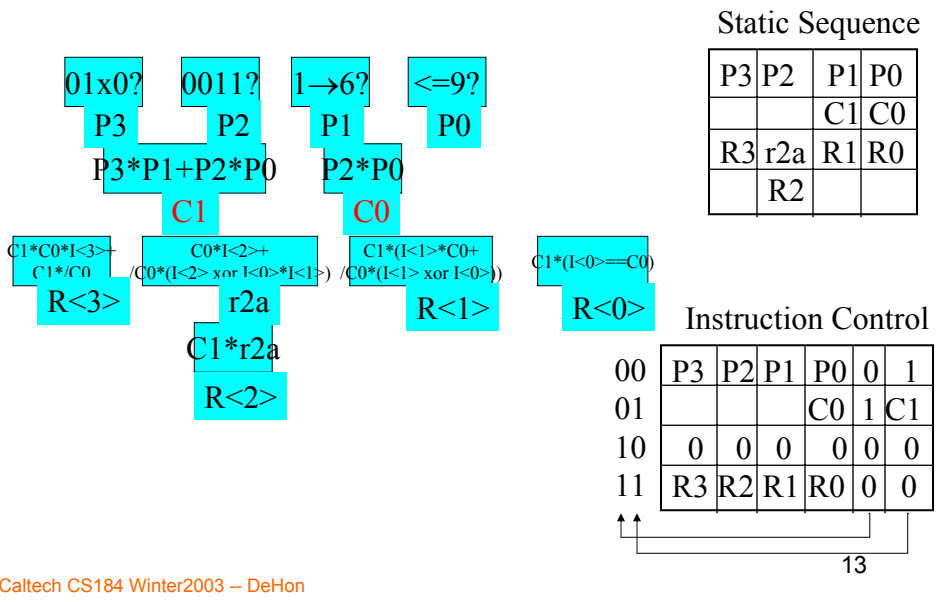
Local Control



Local Control

- LUTs used \neq LUT evaluations produced
- This example:
 - each case only requires 4 4-LUTs to produce $R\langle 3:0 \rangle$
 - takes 5 4-LUTs once unified
- \Rightarrow Counting LUTs not tell cycle-by-cycle LUT needs

Instruction Control



Static/Control

- Static sequence
 - 4 4-LUTs
 - depth 4
 - 4 context
 - maybe 3
 - shuffle r2a w/ C1, C0
 - execute 0 1 1 2 0 1 1 2
- Control
 - 6 4-LUTs
 - depth 3
 - 4 contexts
 - Example too simple to show big savings...

Recall from Multiplier

- Do selective add
 - vs. compute mask and add
 - computation in either case [bit 0, bit 1]
 - (even with branch op)
 - ...was less than unified case
- Also option to terminate early
 - Can optimize for **expected** case

Local vs. Instruction

- If can decide early enough
 - and afford schedule/reload
 - instruction select => less computation
- If load too expensive
 - local instruction
 - faster
 - maybe even less capacity (AT)

Slow Context Switch

- Instruction selection profitable only at coarse grain
 - Xilinx ms reconfiguration times
 - HSRA μ s reconfiguration times
 - still 1000s of cycles
- *E.g.* Video decoder [frame rate = 33ms]
 - if (packet==FRAME)
 - if (type==I-FRAME)
 - IF-context
 - else if (type==B-FRAME)
 - BF-context

Local vs. Instruction

- For multicontext device
 - fast (single) cycle switch
 - factor according to available contexts
- For conventional devices
 - factor only for gross differences
 - and early binding time

Optimization

- Basic Components
 - T_{load} -- config. time
 - T_{select} -- case compute time
 - T_{gen} -- generalized compute time
 - A_{select} -- case compute area
 - A_{gen} -- generalized compute area
- Minimize Capacity Consumed:
 - $AT_{local} = A_{gen} \times T_{gen}$
 - $AT_{select} =$
 - $A_{select} \times (T_{select} + T_{load})$
 - $T_{load} \rightarrow 0$ if can overlap w/ previous operation
 - know early enough
 - background load
 - have sufficient bandwidth to load

FSM Control Factoring Experiment

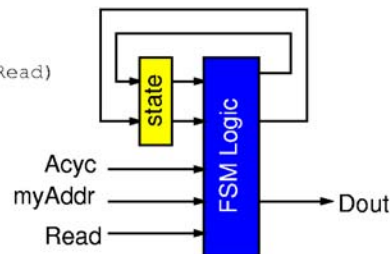
FSM Example

- FSM -- canonical “control” structure
 - captures many of these properties
 - can implement with deep multicontext
 - instruction selection
 - can implement as multilevel logic
 - unify, use local control
- Serve to build intuition

FSM Example (local control)

FSM Description

```
Idle (00):
  if (Acyc & myAddr & Read)
    goto Wait1
  else
    goto Idle
Wait1 (01):
  goto Data
Data (10):
  Assert Dout
  goto Wait2
Wait2 (11):
  goto Idle
```



FSM Logic

```
Dout = S1*/S0
NS0 = /S1*/S0*Acyc*myAddr*Read + S1*/S0
NS1 = /S1*S0 + S1*/S0
```

4 4-LUTs

2 LUT Delays

FSM Example (Instruction)

Context 0 (S1=0)

Dout = 0

NS0 = /S0*Acyc*myAddr*Read

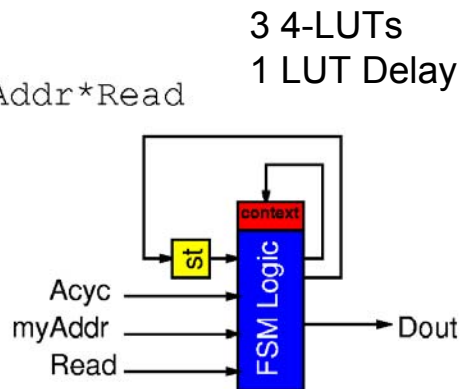
NS1 = S0

Context 1 (S1=1)

Dout = /S0

NS0 = /S0

NS1 = /S0



Full Partitioning Experiment

- Give each state its own context
- Optimize logic in state separately
- Tools
 - mustang, espresso, sis, Chortle
- Use:
 - one-hot encodings for single context
 - smallest/fastest
 - dense for multicontext
 - assume context select needs dense

Look at

- Assume stay in context for a number of LUT delays to evaluate logic/next state
- Pick delay from worst-case
- Assume single LUT-delay for context selection?
 - savings of 1 LUT-delay => comparable time
- Count LUTs in worst-case state

Full Partition (Area Target)

FSM	States	Single Context			Context per State			Ratio $\frac{A_{max}}{A_{avg}}$	Delta Levels
		Levels	N_{LUT}	Area (M λ^2)	Levels	N_{LUT}	Area (M λ^2)		
bbara	10	6	25	22.0	1	6	9.5	0.43	5
bbsse	16	4	50	43.9	3	12	24.6	0.56	1
bbtas	6	3	7	6.1	1	5	6.34	1.0	2
beecount	7	4	14	12.3	1	7	9.4	0.77	3
cse	16	6	83	72.9	2	15	30.7	0.42	4
dk14	7	4	58	50.9	1	8	10.8	0.21	3
dk15	4	12	25	22.0	1	7	7.8	0.35	11
dk16	27	5	80	70.2	1	8	23.2	0.33	4
dk17	8	6	19	16.7	1	6	8.5	0.51	5
dk512	15	2	20	17.6	1	7	13.8	0.79	1
donfile	24	2	46	40.4	1	6	16.0	0.40	1
ex1	20	7	120	105.4	2	26	61.4	0.58	5
ex4	14	7	21	18.4	1	13	24.6	1.33	6
ex6	8	5	57	50.0	1	11	15.7	0.31	4
keyb	19	7	112	98.3	4	14	32.0	0.32	3
mc	4	2	8	7.0	1	7	7.8	1.10	1
modulo12	12	6	12	10.5	1	5	8.7	0.82	5
planet	48	6	150	131.7	1	25	113.6	0.86	5
pma	24	6	82	72.0	2	15	40.1	0.56	4
s1	20	5	137	120.3	5	25	59.0	0.49	0
s1488	48	6	152	133.5	3	27	122.7	0.92	3
s1a	20	5	72	63.2	7	21	49.6	0.78	-2
s208	18	4	38	33.4	1	7	15.4	0.46	3
s27	6	2	5	4.4	1	4	5.1	1.20	1
s386	13	5	42	36.9	2	12	21.8	0.59	3
s420	18	3	40	35.1	1	7	15.4	0.44	2
s510	47	5	54	47.4	1	13	58.1	1.22	4
s8	5	4	12	10.5	1	4	4.7	0.45	3
s820	25	6	92	80.8	3	30	82.5	1.02	3
sand	32	7	178	156.3	5	30	98.9	0.63	2
sse	16	4	50	43.9	3	12	24.6	0.56	1
styr	30	7	186	163.3	4	21	65.9	0.40	3
tbk	32	8	340	298.5	6	33	108.8	0.36	2
Average								0.64	3

Full Partition (Delay Target)

FSM	States	Single Context			Context per State			Ratio $\frac{A_{single}}{A_{single}}$	Delta Levels
		Levels	N_{ALUT}	Area (M λ^2)	Levels	N_{ALUT}	Area (M λ^2)		
bbara	10	3	40	35.1	1	6	9.5	0.27	2
bbse	16	3	60	52.7	2	14	28.7	0.54	1
bbtas	6	2	9	7.9	1	5	6.3	0.80	1
beecount	7	2	19	16.7	1	7	9.4	0.57	1
cse	16	4	97	85.2	2	15	30.7	0.36	2
ck14	7	3	67	58.8	1	8	10.8	0.18	2
ck15	4	3	37	32.5	1	7	7.8	0.24	2
ck16	27	3	83	72.9	1	8	23.2	0.32	2
ck17	8	2	26	22.8	1	6	8.5	0.37	1
ck512	15	2	20	17.6	1	7	13.8	0.79	1
donfile	24	2	46	40.4	1	6	16.0	0.40	1
ex1	20	4	151	132.6	2	26	61.4	0.46	2
ex4	14	2	25	22.0	1	13	24.6	1.12	1
ex6	8	3	62	54.4	1	11	15.7	0.29	2
keyb	19	4	150	131.7	3	26	59.3	0.45	1
mc	4	2	8	7.0	1	7	7.8	1.10	1
modulo12	12	1	13	11.4	1	5	8.7	0.76	0
planet	48	4	172	151.0	1	25	113.6	0.75	3
pma	24	4	139	122.0	2	15	40.1	0.33	2
s1	20	4	195	171.2	3	30	70.8	0.41	1
s1488	48	4	183	160.7	2	28	127.2	0.79	2
s1a	20	3	107	93.9	4	30	70.8	0.75	-1
s208	18	3	40	35.1	1	7	15.4	0.44	2
s27	6	2	5	4.4	1	4	5.0	1.16	1
s386	13	4	54	47.4	2	12	21.8	0.46	2
s420	18	3	40	35.1	1	7	15.4	0.44	2
s510	47	3	76	66.7	1	13	58.1	0.87	2
s8	5	2	13	11.4	1	4	4.8	0.42	1
s820	25	3	137	120.3	3	30	82.5	0.69	0
sand	32	4	224	196.7	3	43	141.7	0.72	1
sse	16	3	60	52.7	2	14	28.7	0.54	1
styr	30	5	285	250.2	3	23	72.2	0.29	2
tbk	32	5	510	447.8	4	42	138.4	0.31	1
Average								0.56	1.36

27

Caltech CS184 Winter2003

Full Partitioning

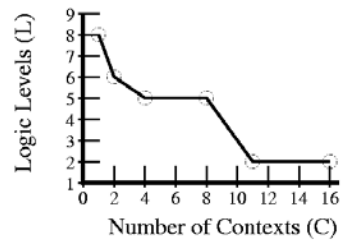
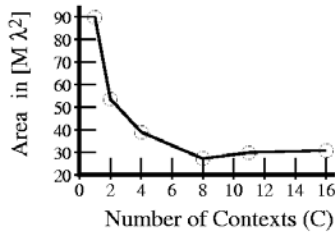
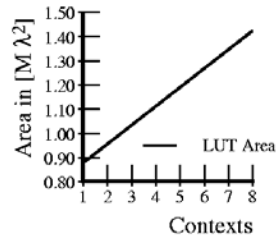
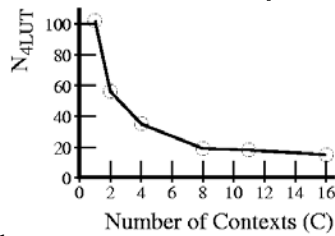
- Full partitioning comes out better
 - ~40% less area
- Note: full partition may not be optimal area case
 - e.g. intro example,
 - no reduction in area or time beyond 2-context implementation
 - 4-context (full partition) just more area
 - (additional contexts)

28

Caltech CS184 Winter2003 -- DeHon

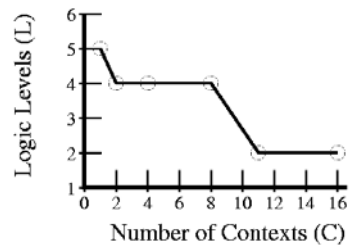
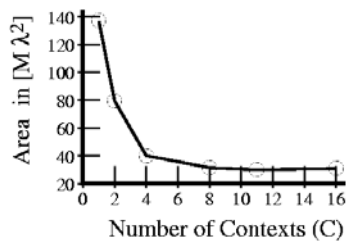
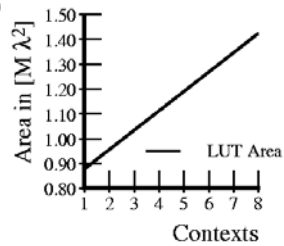
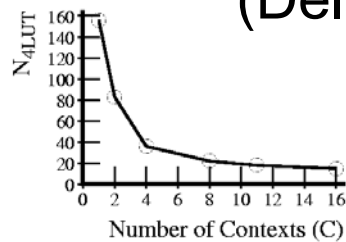
Partitioning versus Contexts (Area)

CSE benchmark



Caltech CS184

Partitioning versus Contexts (Delay)



Caltech CS184

Partitioning versus Contexts (Heuristic)

- Start with dense `mustang` state encodings
- Greedily pick state bit which produces
 - least greatest area split
 - least greatest delay split
- Repeat until have desired number of contexts

Partition to Fixed Number of Contexts

FSM	States	Best Single Context	Area Ratio by Number of Context						
			Dense Encodings						
			1	2	4	8	16	32	64
Area Target									
average ratio		1.00	1.51	0.86	0.63	0.56	0.70	1.09	1.92
average delta		0.00	-0.27	0.33	1.27	2.18	2.70	3.03	3.06
Delay Target									
average ratio		1.00	1.45	1.05	0.59	0.50	0.62	0.95	1.67
average delta		0.00	-0.91	-0.48	0.06	0.64	0.91	1.15	1.21

N.B. - more realistic, device has fixed number of contexts.

Extend Comparison to Memory

- Fully local => compute with LUTs
- Fully partitioned => lookup logic (context) in memory and compute logic
- How compare to fully memory?
 - Simply lookup result in table?

Memory FSM Compare (small)

FSM	states	ins	outs	Min area (Mλ ²)	Integral Addr. & Data Organization	Memory area (Mλ ²)	FPGA area (Mλ ²)	8-ctx DPGA area (Mλ ²)
bftas	6	2	2	0.1	2 ⁵ ×5	0.2	6.1	7.1
dk15	4	3	5	0.3	2 ⁵ ×7	0.3	21.9	10.0
dk17	8	2	3	0.2	2 ⁵ ×6	0.2	16.7	8.5
dk512	15	1	3	0.3	2 ⁵ ×7	0.3	17.6	10.0
mc	4	3	5	0.3	2 ⁵ ×7	0.3	7.0	10.0
modulo12	12	1	1	0.1	2 ⁵ ×5	0.2	10.5	7.1
beecount	7	3	4	0.5	2 ⁶ ×7	0.5	12.3	10.0
dk14	7	3	5	0.5	2 ⁶ ×8	0.6	50.9	11.4
dk16	27	2	3	1.0	2 ⁷ ×8	1.3	70.2	11.4
donfile	24	2	1	0.7	2 ⁷ ×6	0.9	40.4	8.5
s27	6	4	1	0.5	2 ⁷ ×4	0.6	4.4	5.7
s8	5	4	1	0.4	2 ⁷ ×4	0.6	10.5	5.7
bbara	10	4	2	1.2	2 ⁸ ×6	1.8	21.9	11.4
ex6	8	5	8	3.4	2 ⁸ ×11	3.4	50.0	15.7
ex4	14	6	9	14.0	2 ¹⁰ ×13	16.0	18.4	18.5
bbsse	16	7	7	27.0	2 ¹¹ ×11	27.0	43.9	21.4
cse	16	7	7	27.0	2 ¹¹ ×11	27.0	72.9	27.1
tbk	32	6	3	19.7	2 ¹¹ ×8	19.7	298.5	68.4

Memory FSM Compare (large)

FSM	states	ins	outs	Min area (Mλ ²)	Integral Addr. & Data Organization	Memory area (Mλ ²)	FPGA area (Mλ ²)	8-ctx DPGA area (Mλ ²)
sse	16	7	7	27.0	2 ¹¹ × 11	27.0	43.9	21.4
s386	13	7	7	22.9	2 ¹¹ × 11	27.0	36.9	18.5
keyb	19	7	2	20.4	2 ¹² × 7	34.4	98.3	31.3
planet	48	7	19	184.3	2 ¹³ × 25	245.8	131.7	54.1
pma	24	8	8	95.8	2 ¹³ × 13	127.8	72.0	34.2
s1	20	8	6	67.6	2 ¹³ × 11	108.1	120.3	62.7
s1a	20	8	6	67.6	2 ¹³ × 11	108.1	63.2	54.1
ex1	20	9	19	294.9	2 ¹⁴ × 24	471.9	105.4	55.5
s1488	48	8	19	368.6	2 ¹⁴ × 25	491.5	133.5	74.0
styr	30	9	10	276.5	2 ¹⁴ × 15	294.9	163.3	57.0
s208	18	11	2	309.7	2 ¹⁶ × 7	550.5	33.4	12.8
sand	32	11	9	1101.0	2 ¹⁶ × 14	1101.0	156.3	62.7
s820	25	18	19	188743.7	2 ²³ × 24	241591.9	80.8	64.1
s420	18	19	2	79272.3	2 ²⁴ × 7	140928.6	35.1	14.2
s510	47	19	7	384408.9	2 ²⁵ × 13	523449.1	47.4	25.6

Caltech CS184 Winter2003 -- DeHon

Memory FSM Compare (notes)

- Memory selected was “optimally” sized to problem
 - in practice, not get to pick memory allocation/organization for each FSM
 - no interconnect charged
- Memory operate in single cycle
 - but cycle slowing with inputs
- Smaller for <11 state+input bits
- Memory size not affected by CAD quality (FPGA/DPGA is)

Caltech CS184 Winter2003 -- DeHon

Control Granularity

Control Granularity

- What if we want to run multiple of these FSMs on the same component?
 - Local
 - Instruction

Consider

- Two network data ports
 - states: idle, first-datum, receiving, closing
 - data arrival uncorrelated between ports

Local Control Multi-FSM

- Not rely on instructions
- Each wired up independently
- Easy to have multiple FSMs
 - (units of control)

Instruction Control

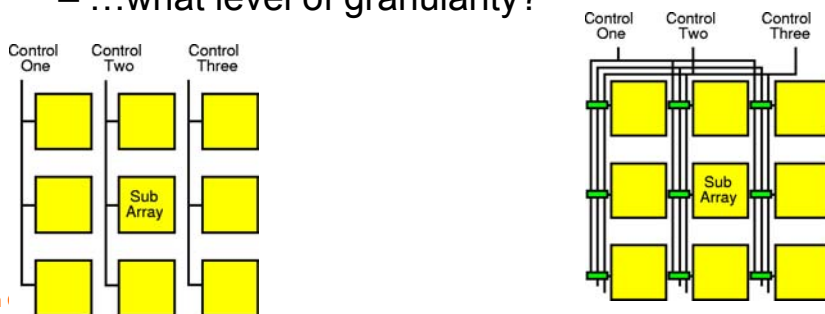
- If FSMs advance orthogonally
 - (really independent control)
 - context depth => product of states
 - for full partition
 - *i.e.* w/ single controller (PC)
 - must create product FSM
 - which may lead to state explosion
 - N FSMs, with S states => S^N product states
 - This example:
 - 4 states, 2 FSMs => 16 state composite FSM

Caltech CS184 Winter2003 -- DeHon

41

Architectural Questions

- How many pins/controller?
- Fixed or Configurable assignment of controllers to pins?
 - ...what level of granularity?



Caltech

Architectural Questions

- Effects of:
 - Too many controllers?
 - Too few controllers?
 - Fixed controller assignment?
 - Configurable controller assignment?

Architectural Questions

- Too many:
 - wasted space on extra controllers
 - synchronization?
- Too few:
 - product state space and/or underuse logic
- Fixed:
 - underuse logic if when region too big
- Configurable:
 - cost interconnect, slower distribution

Control and FPGAs

- Local/single instruction not rely on controller
- Potential strength of FPGA
- Easy to breakup capacity and deploy to orthogonal tasks
- How processor handle? Efficiency?

Control and FPGAs

- Data dependent selection
 - potentially fast w/ local control compared to uP
 - Can examine many bits and perform multi-way branch (output generation) in just a few LUT cycles
 - μ P requires sequence of operations

Architecture Instr. 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)
	n	1	1	FPGA
		w	w	Reconfigurable ALUs
1	1	c	$n_v \cdot 1$	Bitwise SIMD
		w	w	Traditional Processors
		w	$n_v \cdot w$	Vector Processors
	n	c	1	DPGA
		w	8	PADDI
		w	w	VLIW
m	n	1	1	HSRA/SCORE
	c	1	$n_v \cdot w$	MSIMD
	w	1	1	VEGA
m	1	w	8	PADDI-2
		w	w	MIMD (traditional)

Admin

- Final exercise
 - Suggest Meet Wednesday (no lecture W)
 - Also fillout EAS feedback sheets
- Feedback: retiming
 - How improve assignment?
 - Smaller pieces?
- Next term
 - Text: Hennessey and Patterson
 - Computer Architecture: A Quantitative Approach 3rd edition
 - Not at Caltech Bookstore, order now...Amazon

Big Ideas

[MSB Ideas]

- **Control**: where data effects instructions (operation)
- Two forms:
 - local control
 - all ops resident => fast selection
 - instruction selection
 - may allow us to reduce **instantaneous** work requirements
 - introduce issues
 - depth, granularity, instruction load time

Big Ideas

[MSB-1 Ideas]

- Intuition => looked at canonical FSM case
 - few context can reduce LUT requirements considerably (factor dissimilar logic)
 - similar logic more efficient in local control
 - overall, moderate contexts (e.g. 8)
 - exploits both properties
 - better than extremes
 - single context (all local control)
 - full partition
 - flat memory (except for smallest FSMs)