

# ESE680-002 (ESE534): Computer Organization

Day 24: April 11, 2007  
Specialization



Penn ESE680-002 Spring2007 -- DeHon

## Previously

- How to support bit processing operations
- How to compose any task
- Instantaneous << potential computation

Penn ESE680-002 Spring2007 -- DeHon

2

## Today

- What bit operations do I need to perform?
- Specialization
  - Binding Time
  - Specialization Time Models
  - Specialization Benefits
  - Expression

Penn ESE680-002 Spring2007 -- DeHon

3

## Quote

- The fastest instructions you can execute, are the ones you don't.

Penn ESE680-002 Spring2007 -- DeHon

4

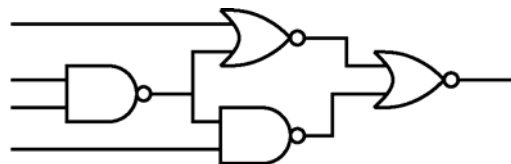
## Idea

- **Goal:** Minimize computation must perform
- Instantaneous computing requirements less than general case
- Some data known or predictable
  - compute minimum computational residue
- As know more data → reduce computation
- Dual of **generalization** we saw for local control

Penn ESE680-002 Spring2007 -- DeHon

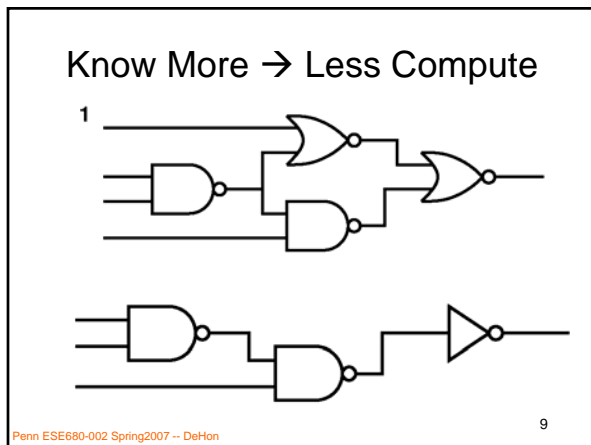
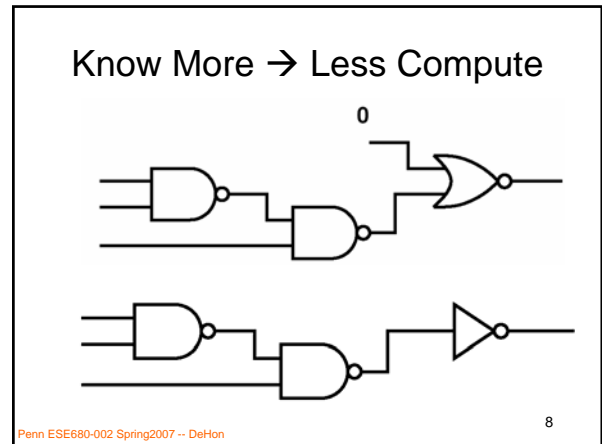
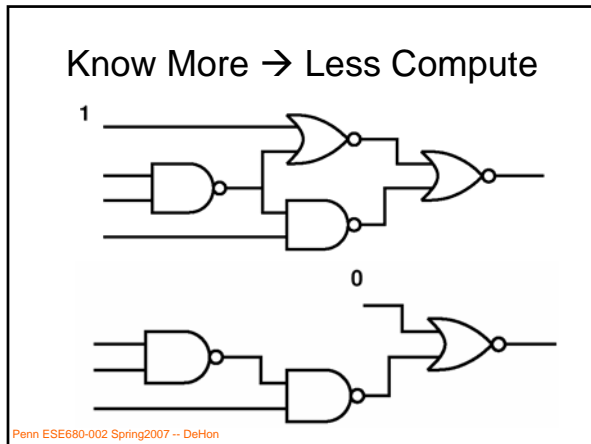
5

## Know More → Less Compute



Penn ESE680-002 Spring2007 -- DeHon

6



- ### Typical Optimization
- Once know another piece of information about a computation (data value, parameter, usage limit)
  - Fold into computation producing smaller computational residue
- Penn ESE680-002 Spring2007 -- DeHon

- ### Opportunity Exists
- Spatial unfolding of computation
    - can afford more specificity of operation
    - E.g. last assignment (FIR, IIR)
  - Fold (early) bound data into problem
  - Common/exceptional cases
- Penn ESE680-002 Spring2007 -- DeHon

- ### Opportunity
- Arises for programmables
    - can change their *instantaneous* implementation
    - don't have to cover all cases with a single configuration
    - can be heavily specialized
      - while still capable of solving entire problem
        - (all problems, all cases)
- Penn ESE680-002 Spring2007 -- DeHon

## Opportunity

- With bit level control
  - larger space of optimization than word level
- While true for both spatial and temporal programmables
  - bigger effect/benefits for spatial

## Multiply Example

Architecture	Feature Size ( $\mu\text{m}$ )	Area and Time	16x16		8x8	
			mpy /s	scale	mpy /s	scale
Custom 16x16	0.63 $\mu\text{m}$	2.6M $\mu\text{m}^2$ , 40 ns	9.6	9.6	9.6	9.6
Custom 8x8	0.80 $\mu\text{m}$	3.3M $\mu\text{m}^2$ , 4.3 ns			70	70
Gate-Array 16x16	0.75 $\mu\text{m}$	26M $\mu\text{m}^2$ , 30ns	1.3	1.3	1.3	1.3
FPGA (XC4K)	0.60 $\mu\text{m}$	1.25M $\mu\text{m}^2$ /CLB 316 CLBs, 26 ns 84 CLBs, 40 ns 220 CLBs, 12.1 ns 22 CLBs, 25 ns	0.097	0.24	0.30	1.5
16b DSP	0.65 $\mu\text{m}$	350M $\mu\text{m}^2$ , 50 ns	0.057	0.057	0.057	0.057
RISC (no multiplier)	0.75 $\mu\text{m}$	125M $\mu\text{m}^2$ , 66 ns/cycle two 16b operands – 44 cycles 16b constant – 7 cycles one 8b operand – 24 cycles 8b constant – 4 cycles	0.0028	0.017	0.0051	0.030

## Multiply Show

- Specialization in datapath width
- Specialization in data

## Benefits

Empirical Examples

## Benefit Examples

- UART
- Pattern match
- Less than
- Multiply revisited
  - more than just constant propagation
- ATR

## UART

- I8251 Intel (PC) standard UART
- Many operating modes
  - bits
  - parity
  - sync/async
- Run in same mode for length of connection

## UART FSMs

FSM	Fully Generic				Specialized	
	Speed Mapped CLBs	path	Area Mapped CLBs	path	CLBs	path
18251 processor i/o	11	3.5	11	3.5	6.5	2
	fast (any configuration) small (any configuration)				5.5	2.5
18251 transmitter	57.5	4.5	57.5	4.5	24	4
	Asynchronous, parity				27	4.5
	Asynchronous, no parity				31	4.5
	2 Sync chars, parity				31	4
18251 receiver	52.5	5.5	52.5	5.5	32.5	4.5
	Asynchronous, parity				36	4.5
	Asynchronous, no parity				29.5	4.5
	External Sync, parity				27	3.5
	Internal, 2 Sync chars, parity				28	4.5
	Internal, 1 Sync char, parity				28	4.5
	Internal, 1 Sync char, no parity				31.5	4.5

Penn ESE680-002

19

## UART Composite

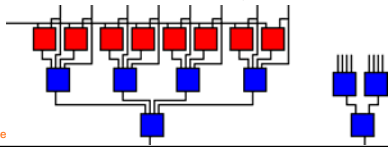
design	Fully Generic				Specialized	
	Speed Mapped CLBs	path	Area Mapped CLBs	path	CLBs	path
18251 core	358.5	8.5	348.5	10.5	216.5	7
	Async, 64 clks/bit, 8e2				201	6
	Async, 16 clks/bit, 8n1				141.5	4.5
	Async, 1 clks/bit, 5n1				165	4.5
	Sync, internal, 2 sync, 8o				136	5.5
	Sync, external, 5n					

Penn ESE680-002 Spring2007 -- DeHon

20

## Pattern Match

- Savings:
  - 2N bit input computation → N
  - if N variable, maybe trim unneeded portion
  - state elements store target
  - control load target



Pe

21

## Pattern Match

(size)	CLBs	path	CLBs	path	AT Ratio	
$a = b$	$b$ variable		$b$ constant			$w/state$
(8)	2.5 (+4)	2	1.5	2	0.60	0.23
(16)	5.5 (+8)	3	2.5	2	0.30	0.12
(32)	10.5 (+16)	3	5.5	3	0.52	0.21
(64)	21.5 (+32)	4	10.5	3	0.37	0.15

Penn ESE680-002 Spring2007 -- DeHon

22

## Less Than (Bounds check?)

- Area depend on target value
- But all targets less than generic comparison

Function (size)	Speed Mapped CLBs	path	Area Mapped CLBs	path	Speed Mapped CLBs	path	Area Mapped CLBs	path
$a < b$	$b$ variable				$b$ constant			
(8)	4	8	4	8	$< 2$	$< 2$	$< 1.5$	$< 3$
(16)	18.5	14	16.5	16	$< 6.5$	$< 3$	$< 3$	$< 5$
(32)	35	19	36	24	$< 13.5$	$< 4$	$< 6$	$< 11$
(64)	77.5	20	74.5	28	$< 30$	$< 5$	$< 14$	$< 16$

Penn ESE680-002 Spring2007 -- DeHon

23

## Multiply (revisited)

- Specialization can be more than constant propagation
- Naïve,
  - save product term generation
  - complexity number of 1's in constant input
- Can do better exploiting algebraic properties

Penn ESE680-002 Spring2007 -- DeHon

24

## Multiply

- Never really need more than  $\lfloor N/2 \rfloor$  one bits in constant
- If more than  $N/2$  ones:
  - invert c  $(2^{N+1}-1-c)$
  - (less than  $N/2$  ones)
  - multiply by x  $(2^{N+1}-1-c)x$
  - add x  $(2^{N+1}-c)x$
  - subtract from  $(2^{N+1})x = cx$

Penn ESE680-002 Spring2007 -- DeHon

25

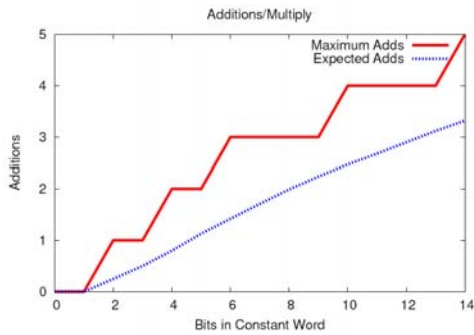
## Multiply

- At most  $\lfloor N/2 \rfloor + 2$  adds for any constant
- Exploiting common subexpressions can do better:
  - e.g.
    - $c=10101010$
    - $t1=x+x<<2$
    - $t2=t1<<5+t1<<1$

Penn ESE680-002 Spring2007 -- DeHon

26

## Multiply



Penn ESE680-002 Spring2007 -- DeHon

27

## Example: ATR

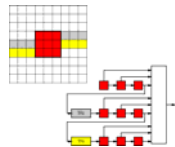
- Automatic Target Recognition
  - need to score image for a number of different patterns
    - different views of tanks, missiles, etc.
  - reduce target image to a binary template with don't cares
  - need to track many (e.g. 70-100) templates for each image region
  - templates themselves are sparse
    - small fraction of care pixels

Penn ESE680-002 Spring2007 -- DeHon

28

## Example: ATR

- $16 \times 16 \times 2 = 512$  flops to hold single target pattern
- $16 \times 16 = 256$  LUTs to compute match
- 256 score bits  $\rightarrow$  8b score  $\sim$  500 adder bits in tree
- more for retiming
- $\sim 800$  LUTs here
- Maybe fit 1 generic template in XC4010 (400 CLBs)?



Penn ESE680-002 Spring2007 -- DeHon

29

## Example: UCLA ATR

- UCLA
  - specialize to template
  - ignore don't care pixels
  - only build adder tree to care pixels
  - exploit common subexpressions
  - get 10 templates in a XC4010

[Villasenor et. al./FCCM'96]

Penn ESE680-002 Spring2007 -- DeHon

30

## Example: FIR Filtering

$$Y_i = w_1 x_i + w_2 x_{i+1} + \dots$$

Architecture	Feature Size ( $\lambda$ )	$\frac{TAPs}{\lambda^2}$
32b RISC	0.75 $\mu$ m	0.020
16b DSP	0.65 $\mu$ m	0.057
32b RISC/DSP	0.25 $\mu$ m	0.021
64b RISC	0.18 $\mu$ m	0.064
FPGA (XC4K) (Altera 8K)	0.60 $\mu$ m	1.9
Full Custom	0.30 $\mu$ m	3.6
	0.75 $\mu$ m	3.6
	0.60 $\mu$ m	3.5
	0.75 $\mu$ m	2.4
(fixed coefficient) (n.b. 16b samples)	0.60 $\mu$ m	56
		..

Penn ESE680-002 Spring2007 -- DeHon

## Usage Classes

Penn ESE680-002 Spring2007 -- DeHon

32

## Specialization Usage Classes

- Known binding time
- Dynamic binding, persistent use
  - apparent
  - empirical
- Common case

Penn ESE680-002 Spring2007 -- DeHon

33

## Known Binding Time

- Sum=0
- For  $l=0 \rightarrow N$ 
  - Sum += V[l]
- For  $l=0 \rightarrow N$ 
  - VN[l] = V[l] / Sum
- Scale(max,min,V)
  - for  $l=0 \rightarrow V.length$ 
    - tmp = (V[l] - min)
    - Vres[l] = tmp / (max - min)

Penn ESE680-002 Spring2007 -- DeHon

34

## Dynamic Binding Time

- cexp=0;
- For  $l=0 \rightarrow V.length$ 
  - if (V[l].exp != cexp)
    - cexp = V[l].exp;
  - Vres[l] = V[l].mant << cexp
- Thread 1:
  - a = src.read()
  - if (a.newavg())
    - avg = a.avg()
- Thread 2:
  - v = data.read()
  - out.write(v/avg)

Penn ESE680-002 Spring2007 -- DeHon

35

## Empirical Binding

- Have to check if value changed
  - Checking value  $O(N)$  area [pattern match]
  - Interesting because computations
    - can be  $O(2^N)$  [Day 9]
    - often greater area than pattern match
  - Also Rent's Rule:
    - Computation > linear in IO
    - $IO = C n^p \rightarrow n \propto IO^{(1/p)}$

Penn ESE680-002 Spring2007 -- DeHon

36

## Common/Uncommon Case

- For  $i=0 \rightarrow N$ 
  - If  $(V[i]==10)$ 
    - $SumSq+=V[i]*V[i];$
  - elseif  $(V[i]<10)$ 
    - $SumSq+=V[i]*V[i];$
  - else
    - $SumSq+=V[i]*V[i];$

Penn ESE680-002 Spring2007 -- DeHon

37

## Binding Times

- Pre-fabrication
- Application/algorithm selection
- Compilation
- Installation
- Program startup (load time)
- Instantiation (new ...)
- Epochs
- Procedure
- Loop

Penn ESE680-002 Spring2007 -- DeHon

38

## Exploitation Patterns

- Full Specialization (Partial Evaluation)
  - May have to run (synth?) p&r at runtime
- Worst-case footprint
  - e.g. multiplier worst-case, avg., this case
- Constructive Instance Generator
- Range specialization (wide-word datapath)
  - data width
- Template
  - e.g. pattern match – only fillin LUT prog.

Penn ESE680-002 Spring2007 -- DeHon

39

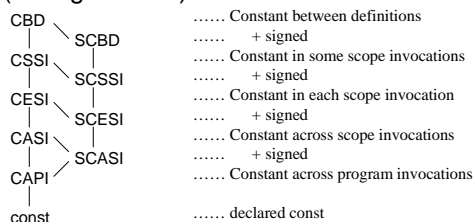
## Opportunity Example

Penn ESE680-002 Spring2007 -- DeHon

40

## Bit Constancy Lattice

- binding time for bits of variables (storage-based)



Penn ESE680-002 Spring2007 -- DeHon

[Experiment: Eylon Caspi/UCB] 41

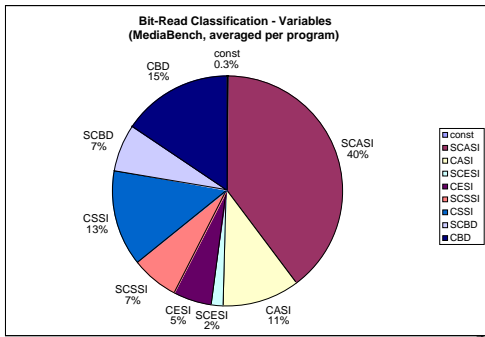
## Experiments

- Applications:
  - UCLA MediaBench: adpcm, epic, g721, gsm, jpeg, mesa, mpeg2 (not shown today - ghostscript, pegwit, pgp, rasta)
  - gzip, versatility, SPECint95 (parts)
- Compiler optimize  $\rightarrow$  instrument for profiling  $\rightarrow$  run
- analyze variable usage, ignore heap
  - heap-reads typically 0-10% of all bit-reads
  - 90-10 rule (variables) - ~90% of bit reads in 1-20% of bits

Penn ESE680-002 Spring2007 -- DeHon

[Experiment: Eylon Caspi/UCB] 42

## Empirical Bit-Reads Classification



Penn ESE680-002 Spring2007 -- DeHon

[Experiment: Eylon Caspi/UCB]

## Bit-Reads Classification

- regular across programs
  - SCASI, CASI, CBD stddev ~11%
- nearly no activity in variables declared const
- ~65% in constant + signed bits
  - trivially exploited

[Experiment: Eylon Caspi/UCB]

44

Penn ESE680-002 Spring2007 -- DeHon

## Constant Bit-Ranges

- 32b data paths are too wide
- 55% of all bit-reads are to sign-bits
- most CASI reads clustered in bit-ranges (10% of 11%)
- CASI+SCASI reads (50%) are positioned:
  - 2% low-order
  - 39% high-order
  - 8% whole-word constant
  - 1% elsewhere

Penn ESE680-002 Spring2007 -- DeHon

[Experiment: Eylon Caspi/UCB]

45

## Issue Roundup

Penn ESE680-002 Spring2007 -- DeHon

46

## Expression Patterns

- Generators
- Instantiation/Immutable computations
  - (disallow mutation once created)
- Special methods (only allow mutation with)
- Data Flow (binding time apparent)
- Control Flow
  - (explicitly separate common/uncommon case)
- Empirical discovery

Penn ESE680-002 Spring2007 -- DeHon

47

## Benefits

- Much of the benefits come from reduced area
  - reduced area
    - room for more spatial operation
    - maybe less interconnect delay
- Fully exploiting, full specialization
  - don't know how big a block is until see values
  - dynamic resource scheduling

Penn ESE680-002 Spring2007 -- DeHon

48



## Optimization Prospects

- Area-Time Tradeoff

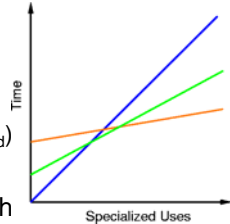
- $T_{\text{spl}} = T_{\text{sc}} + T_{\text{load}}$

- $AT_{\text{gen}} = A_{\text{gen}} \times T_{\text{gen}}$

- $AT_{\text{spl}} = A_{\text{spl}} \times (T_{\text{sc}} + T_{\text{load}})$

- If compute long enough

- $T_{\text{sc}} \gg T_{\text{load}} \rightarrow$  amortize out load



49

Penn ESE680-002 Spring2007 -- DeHon

## Storage

- Will have to store configurations somewhere
- LUT  $\sim 1M\lambda^2$
- Configuration 64+ bits
  - SRAM:  $80K\lambda^2$  (12-13 for parity)
  - Dense DRAM:  $6.4K\lambda^2$  (160 for parity)

50

Penn ESE680-002 Spring2007 -- DeHon

## Saving Instruction Storage

- Cache common, rest on alternate media
  - e.g. disk
- Compressed Descriptions
- Algorithmically composed descriptions
  - good for regular datapaths
  - think Kolmogorov complexity
- Compute values, fill in template
- Run-time configuration generation

51

Penn ESE680-002 Spring2007 -- DeHon

## Open

- How much opportunity exists in a given program?
- Can we measure entropy of programs?
  - How constant/predictable is the data compute on?
  - Maximum potential benefit if exploit?
  - Measure efficiency of architecture/implementation like measure efficiency of compressor?

52

Penn ESE680-002 Spring2007 -- DeHon

## Big Ideas

- Programmable advantage
  - Minimize work by specializing to instantaneous computing requirements
- Savings depends on functional complexity
  - but can be substantial for large blocks
  - close gap with custom?

53

Penn ESE680-002 Spring2007 -- DeHon

## Big Ideas

- Several models of structure
  - slow changing/early bound data, common case
- Several models of exploitation
  - template, range, bounds, full special

54

Penn ESE680-002 Spring2007 -- DeHon