

ESE534: Computer Organization

Day 24: April 16, 2012
Specialization



Previously

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

Today

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

Quote

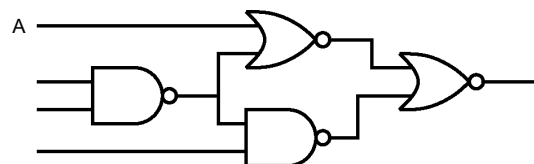
- The fastest instructions you can execute, are the ones you don't.

– ...and the least energy, too!

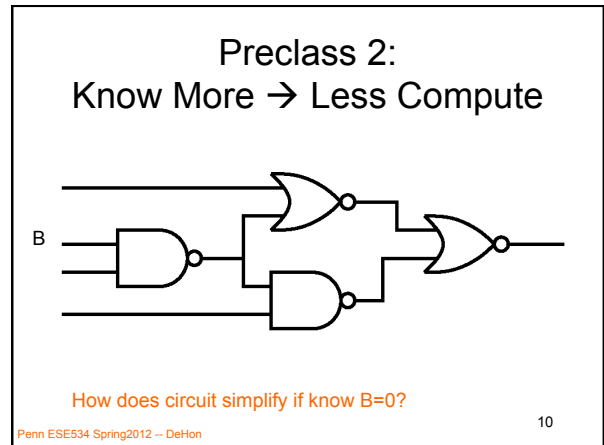
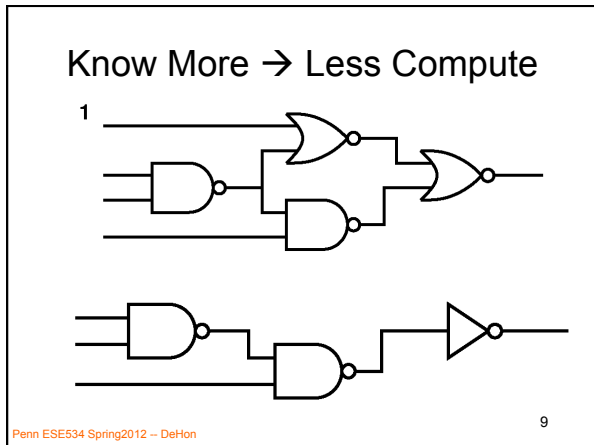
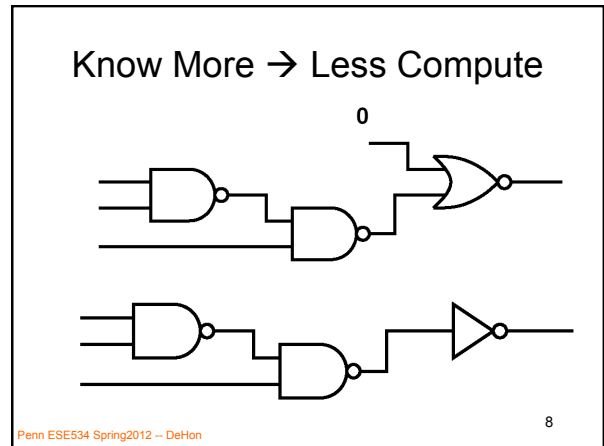
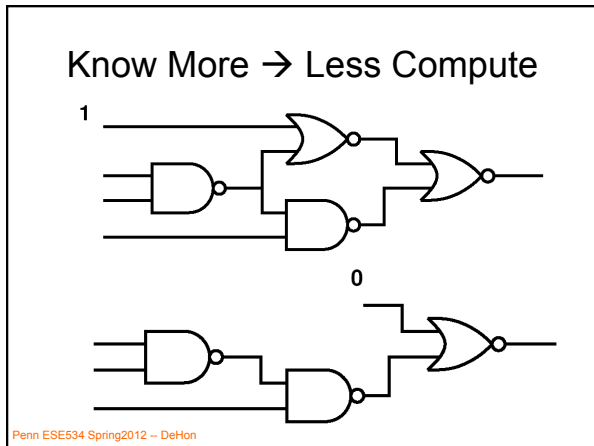
Idea

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

Preclass 1: Know More → Less Compute



How does circuit simplify if know A=1?



Possible Optimization

- Once know another piece of information about a computation (data value, parameter, usage limit)
- Fold into computation producing smaller computational residue

Penn ESE534 Spring2012 -- DeHon

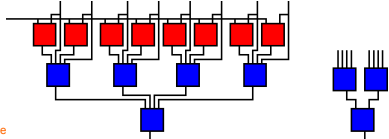
Preclass 3

- How many 4-LUTs for 8b-equality compare?
- How many 4-LUTs for 8b compare to constant?

Penn ESE534 Spring2012 -- DeHon

Pattern Match

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



Pe

13

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 ESE534 Spring2012 -- DeHon

14

Opportunity Exists

- Spatial unfolding of computation
 - can afford more specificity of operation
- Fold (early) bound data into problem
- Common/exceptional cases

Penn ESE534 Spring2012 -- DeHon

15

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 ESE534 Spring2012 -- DeHon

16

Preclass 4

```

7: int compare(char *target, char *potential)
8: {
9:     int i;
10:    char *p1=target;
11:    char *p2=potential;
12:    for (i=0;i<MATCH_LENGTH;i++)
13:    {
14:        if (*p1!=*p2) return(0);
15:        p1++;
16:        p2++;
17:    }
18:    return(1);
19: }
20:
21: int count_matches(FILE *fd, char *target)
22: {
23:    char *line=(char *)malloc(sizeof(char)*MAX_LINE_LENGTH);
24:    int matches=0;
25:    while (!feof(fd))
26:    {
27:        int cread=read_line(fd, line, MAX_LINE_LENGTH);
28:        while (cread-MATCH_LENGTH>0)
29:        {
30:            if (compare(target, line)>0)
31:                matches++;
32:            line++;
33:        }
34:        return(matches);
35:    }
36: }

```

MATCH_LENGTH bound?
 cread bound?
 Uses per binding?

Penn ESE534 Spring2012 -- DeHon

17

Preclass 4

```

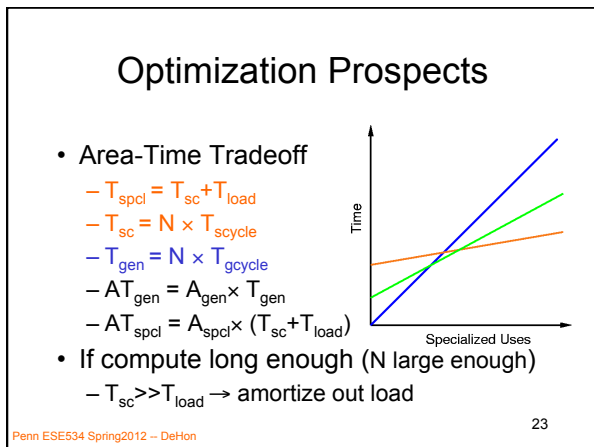
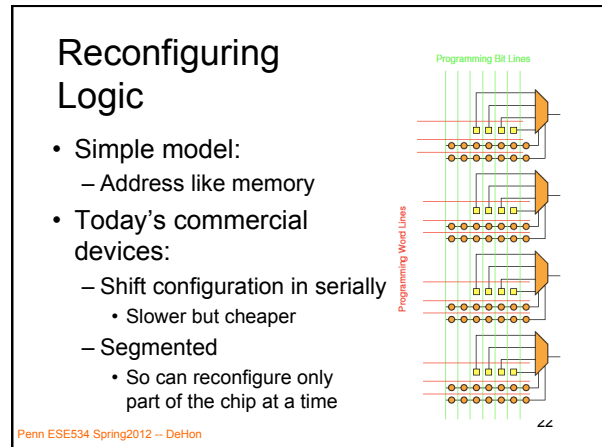
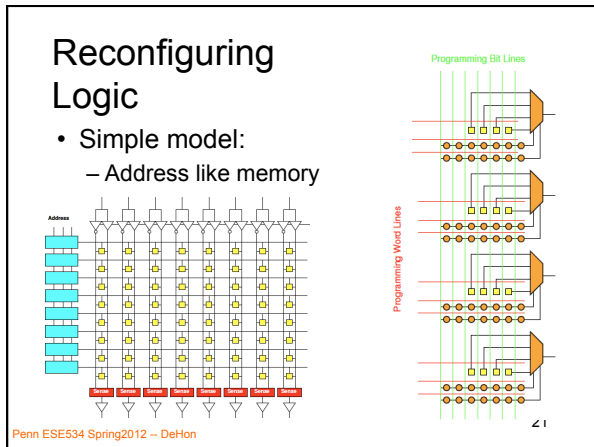
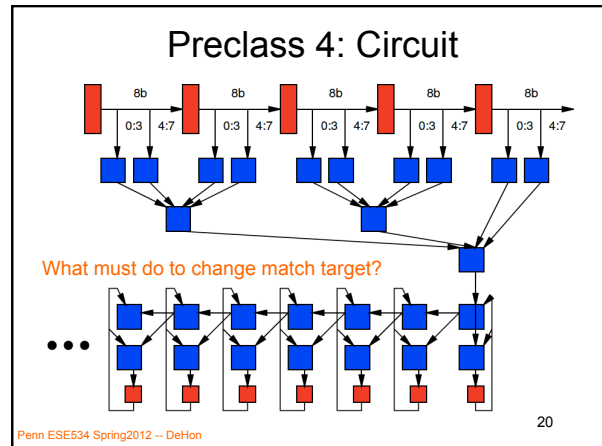
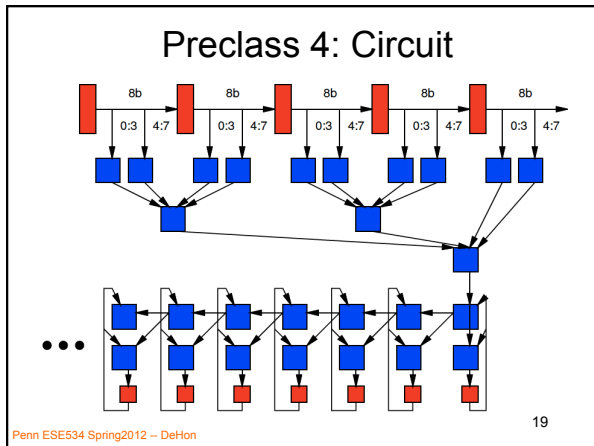
38: int main(int argc, char *argv[])
39: {
40:     FILE *fd;
41:     int cnt;
42:     char *target;
43:
44:     target=(char *)malloc(sizeof(char)*MATCH_LENGTH);
45:
46:     if (argc>=2)
47:     {
48:         strncpy(target, argv[1], MATCH_LENGTH);
49:         fd=fopen(argv[2], "r");
50:         cnt=count_matches(fd, target);
51:         fprintf(stdout, "Matches=%d\n", cnt);
52:     }
53: }

```

Target[0] bound?
 Uses per binding?

Penn ESE534 Spring2012 -- DeHon

18



Preclass 5

- $T_{load} = 100\mu s, T_{scycle} = 1ns$
- $T_{gload} = 0, T_{gcycle} = 2ns$
- Ratio T_{gtask}/T_{ctask} for $N = 10^6$?
- Breakeven N?

Penn ESE534 Spring2012 -- DeHon

24

Optimization Prospects

- Area-Time Tradeoff

- $T_{spcl} = T_{sc} + T_{load}$

- $T_{sc} = N \times T_{scycle}$

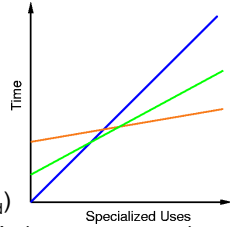
- $T_{gen} = N \times T_{gcycle}$

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

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

- If compute long enough (N large enough)

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



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 (Λ)	Area and Time	16 × 16		8 × 8	
			mpy ₁₆	scale ₁₆	mpy ₈	scale ₈
Custom 16 × 16	0.63 μm	2.6MΛ ² , 40 ns	9.6	9.6	9.6	9.6
Custom 8 × 8	0.80 μm	3.3MΛ ² , 4.3 ns			70	70
Gate-Array 16 × 16	0.75 μm	26MΛ ² , 30ns	1.3	1.3	1.3	1.3
FPGA (XC4K)	0.60 μm	1.25MΛ 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	
16b DSP	0.65 μm	350MΛ ² , 50 ns	0.057	0.057	0.057	1.5
RISC (no multiplier)	0.75 μm	125MΛ 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

- Less than
- Multiply revisited
 - more than just constant propagation
- ATR

Benefit Examples

- Less than
- Multiply revisited
 - more than just constant propagation
- ATR

Less Than (Bounds check?)

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

Function (size)	Speed Mapped CLBs	Area Mapped path	Speed Mapped CLBs	Area Mapped path	Speed Mapped CLBs	Area Mapped path	Speed Mapped CLBs	Area Mapped 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 ESE534 Spring2012 -- DeHon

31

Multiply

- How savings in a multiply by constant?
- Multiply by 80?
– 0101000
- Multiply by 255?

Penn ESE534 Spring2012 -- DeHon

32

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 ESE534 Spring2012 -- DeHon

33

Multiply

- Never really need more than $\lfloor N/2 \rfloor$ one bits in constant
- Example: multiply by 255:
 - $256x - x = 255x$
 - $t1 = x \ll 8$
 - $res = t1 - x$

Penn ESE534 Spring2012 -- DeHon

34

Multiply

- Never really need more than $\lfloor N/2 \rfloor$ one bits in constant
- If more than $N/2$ ones:
 - invert c $(2^{N+1}-1-c)$ 11111111-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 ESE534 Spring2012 -- DeHon

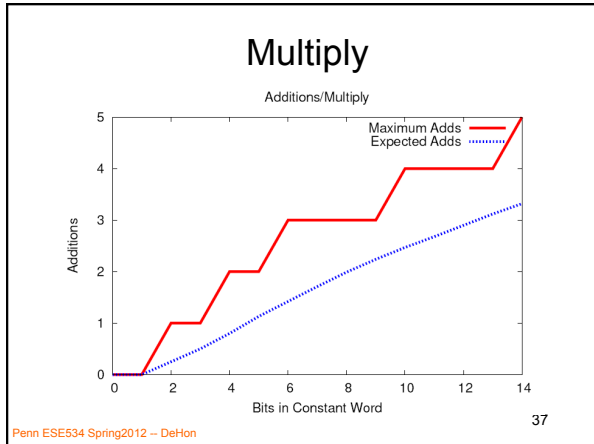
35

Multiply

- At most $\lfloor N/2 \rfloor + 2$ adds for any constant
- Exploiting common subexpressions can do better:
 - e.g.
 - $c = 10101010$
 - $t1 = x + x \ll 2$ (101x)
 - $t2 = t1 \ll 5 + t1 \ll 1$
 - \rightarrow 2 adds instead of 4

Penn ESE534 Spring2012 -- DeHon

36



Multiply Example

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

Penn ESE534 Spring2012 -- DeHon

Example: FIR Filtering

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

Application metric:
TAPs = filter taps
multiply accumulate

Architecture	Feature Size (λ)	TAPs /s
32b RISC	0.75 μ m	0.020
16b DSP	0.65 μ m	0.057
32b RISC/DSP	0.25 μ m	0.021
64b RISC	0.18 μ m	0.064
FPGA (XC4K)	0.60 μ m	1.9
(Altera 8K)	0.30 μ m	3.6
Full Custom	0.75 μ m	3.6
	0.60 μ m	3.5
	0.75 μ m	2.4
(fixed coefficient)	0.60 μ m	56
(n.b. 16b samples)		~

Penn ESE534 Spring2012 -- DeHon

- ### Opportunity Exists
- Spatial unfolding of computation
 - can afford more specificity of operation
- $$Y_i = w_1 x_i + w_2 x_{i+1} + \dots$$
- What opportunity do we lose if sequentializing on single multiplier?
- Penn ESE534 Spring2012 -- DeHon

- ### 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 ESE534 Spring2012 -- DeHon

- ### Example: ATR
- 16x16x2=512 flops to hold single target pattern
 - 16x16=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 ESE534 Spring2012 -- DeHon

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]

43

Penn ESE534 Spring2012 -- DeHon

Usage Classes

44

Penn ESE534 Spring2012 -- DeHon

Specialization Usage Classes

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

45

Penn ESE534 Spring2012 -- DeHon

Known Binding Time

- Sum=0 Scope/Procedure Invocation
- For I=0→N Scale(max,min,V)
 - Sum+=V[I] for I=0→V.length
- For I=0→N tmp=(V[I]-min)
 - VN[I]=V[I]/Sum Vres[I]=tmp/(max-min)

46

Penn ESE534 Spring2012 -- DeHon

Dynamic Binding Time

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

47

Penn ESE534 Spring2012 -- DeHon

Empirical Binding

- Have to check if value changed
 - Checking value O(N) area [pattern match]
 - Interesting because computations
 - can be O(2^N) [Day 12]
 - often greater area than pattern match
 - Also Rent's Rule:
 - Computation > linear in IO
 - IO=C n^P → n ∝ IO^(1/p)

48

Penn ESE534 Spring2012 -- DeHon

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
 - $SumSq += V[i] * V[i];$
- For $i=0 \rightarrow N$
 - If $(V[i] == 10)$
 - $SumSq += 100;$
 - elseif $(V[i] < 10)$
 - $SumSq += V[i] * V[i];$
 - else
 - $SumSq += V[i] * V[i];$

Potential Binding Times

- What are the potential binding times for values?
 - i.e. at what points might values be defined then held constant?

Binding Times

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

Exploitation Patterns

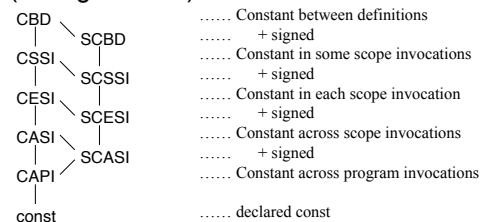
- Full Specialization (Partial Evaluation)
 - May have to run (synth?) p&r at runtime
- Worst-case footprint
 - e.g. multiplier worst-case, avg., this case
- Constructive Instance Generator
- Range specialization (wide-word datapath)
 - data width
- Template
 - e.g. pattern match – only fillin LUT prog.

Opportunity Example

(Lecture ended here)

Bit Constancy Lattice

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



Experiments

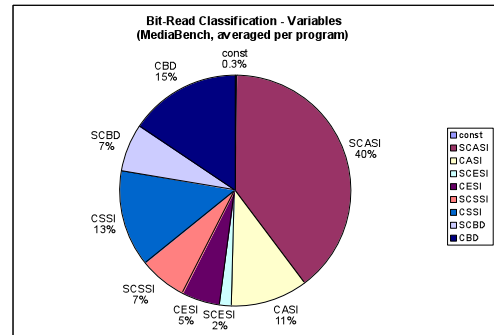
- Applications:
 - UCLA MediaBench: adpcm, epic, g721, gsm, jpeg, mesa, mpeg2 (not shown today - ghostscript, pegwit, pgp, rasta)
 - gzip, versatility, SPECint95 (parts)
- Compiler optimize → instrument for profiling → run
- analyze variable usage, ignore heap
 - heap-reads typically 0-10% of all bit-reads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

Penn ESE534 Spring2012 -- DeHon

[Experiment: Eylon Caspi/UCB]

55

Empirical Bit-Reads Classification



Penn ESE534 Spring2012 -- 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

Penn ESE534 Spring2012 -- DeHon

[Experiment: Eylon Caspi/UCB]

57

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 constant
 - 39% high-order constant
 - 8% whole-word constant
 - 1% elsewhere

Penn ESE534 Spring2012 -- DeHon

[Experiment: Eylon Caspi/UCB]

58

Issue Roundup

Penn ESE534 Spring2012 -- DeHon

59

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 ESE534 Spring2012 -- DeHon

60

Benefits

- Benefits come from reduced area & energy
 - reduced area → performance
 - room for more spatial operation
 - maybe less interconnect delay
- Challenge: Fully exploiting, full specialization
 - don't know how big a block is until see values
 - dynamic resource scheduling

Storage

- Will have to store configurations somewhere
- LUT ~ 250K F²
- Configuration 64+ bits
 - SRAM: 20KF² (12-13 for parity)
 - Dense DRAM: 1.6KF² (160 for parity)

Saving Instruction Storage

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

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?

Admin

- FM1 graded
- Next priority: FM2 feedback
- Final week of discussion period
 - Ends April 24th
- Reading for Wednesday online

Big Ideas [MSB]

- 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?

Big Ideas [MSB-1]

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