

CS184a: Computer Architecture (Structure and Organization)

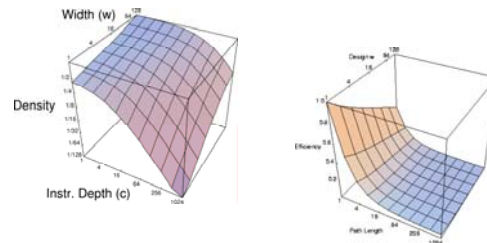
Day 10: January 28, 2005
Empirical Comparisons



Caltech CS184 Winter2005 -- DeHon

Last Time

- Instruction Space Modeling



Caltech CS184 Winter2005 -- DeHon

2

Today

- Empirical Data
 - Processors
 - FPGAs
 - Custom
 - Gate Array
 - Std. Cell
 - Full
 - Tasks

Caltech CS184 Winter2005 -- DeHon

3

Empirical Comparisons

Caltech CS184 Winter2005 -- DeHon

4

Empirical

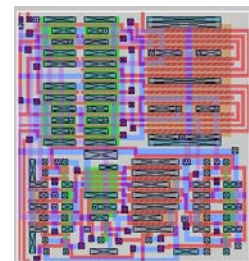
- Ground modeling in some concretes
- Start sorting out
 - custom vs. configurable
 - spatial configurable vs. temporal

Caltech CS184 Winter2005 -- DeHon

5

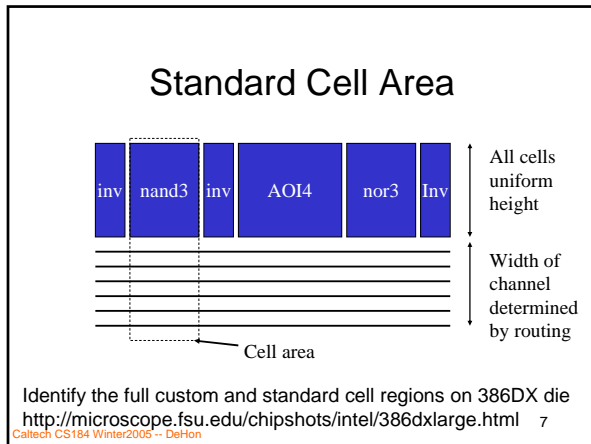
Full Custom

- Get to define all layers
- Use any geometry you like
- Only rules are process design rules
- CS181



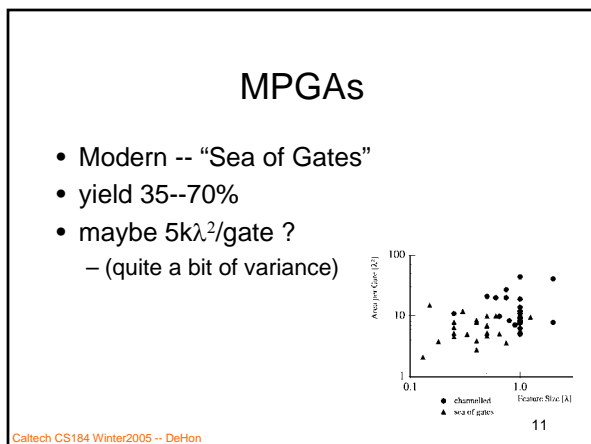
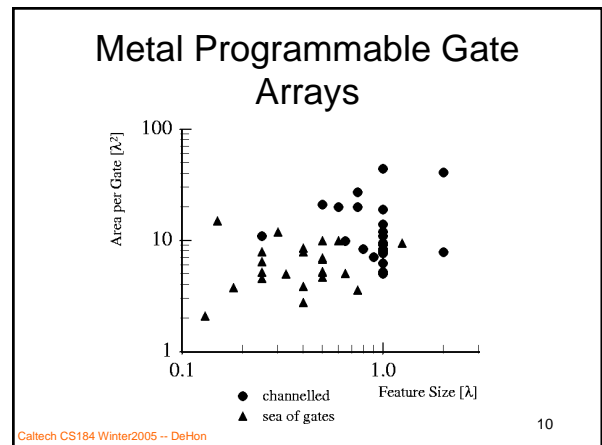
Caltech CS184 Winter2005 -- DeHon

6



- ### MPGA
- Metal Programmable Gate Array
 - Gates pre-placed (poly, diffusion)
 - Only get to define metal connections
 - Cheap – only have to pay for metal mask(s)
- Caltech CS184 Winter2005 -- DeHon 8

- ### MPGA vs. Custom?
- | | |
|--|--|
| <ul style="list-style-type: none"> • AMI CICC'83 <ul style="list-style-type: none"> – MPGA 1.0 – Std-Cell 0.7 – Custom 0.5 • Toshiba DSP <ul style="list-style-type: none"> – Custom 0.3 • Mosaid RAM <ul style="list-style-type: none"> – Custom 0.2 | <ul style="list-style-type: none"> • GE CICC'86 <ul style="list-style-type: none"> – MPGA 1.0 – Std-Cell 0.4--0.7 <ul style="list-style-type: none"> • FF/counter 0.7 • FullAdder 0.4 • RAM 0.2 <p>MPGA = Metal Programmable Gate Array (traditional Gate Array)</p> |
|--|--|
- Caltech CS184 Winter2005 -- DeHon 9



FPGA Table

Year	Design	Organization	Max	λ	λ^2 area	cycle
1986	Xilinx 2K	CLB (4-LUT)	100	1μ	500K	20 ns
1988	Xilinx 3K	CLB (2x4-LUT)	320	0.6μ	1.3M	13 ns
1992	Xilinx 4K	CLB (2x4-LUT +)	1024	0.6μ	1.25M	7 ns
1995	Xilinx 5K	CLB (4x4-LUTS)	484	0.3μ	2.25M	6 ns
1995	Altera 8K	LE (4-LUT)	1296	0.3μ	920K	7.5 ns
1995	ORCA 2C	PLC (4x4-LUT)	900	0.3μ	4.3M	7 ns
1998	HSRA	BLB (5-LUT/2x4-LUT ?)	–	0.2μ	2M	4 ns
	Model	4-LUT	2K	–	800K	–
	Model	4-LUT	16K	–	1M	–

Caltech CS184 Winter2005 -- DeHon 12

Modern FPGAs

- APEX 20K1500E
 - 52K LEs
 - 0.18 μ m
 - 24mm \times 22mm
- XC2V1000
 - 10.44mm \times 9.90mm [source: Chipworks]
 - 0.15 μ m
 - 11,520 4-LUTs
 - 1.5M λ^2 /4-LUT

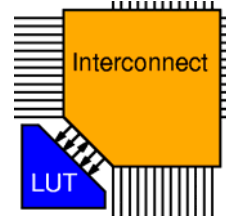
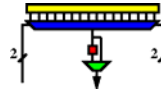
[Both also have RAM in cited area]

Caltech CS184 Winter2005 -- DeHon

13

Conventional FPGA Tile

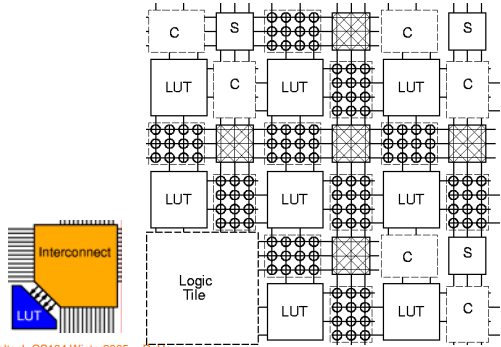
K-LUT (typical k=4)
w/ optional
output Flip-Flop



Caltech CS184 Winter2005 -- DeHon

14

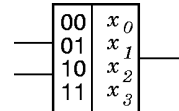
Toronto FPGA Model



Caltech CS184 Winter2005 -- DeHon

15

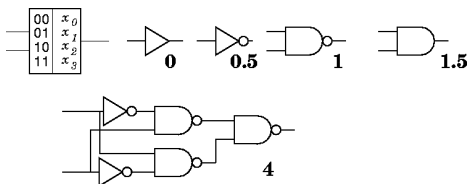
How many gates?



Caltech CS184 Winter2005 -- DeHon

16

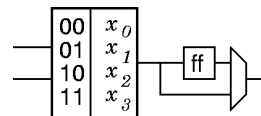
"gates" in 2-LUT



Caltech CS184 Winter2005 -- DeHon

17

Now how many?



Caltech CS184 Winter2005 -- DeHon

18

Which gives:
More usable gates?
More gates/unit area?

Caltech CS184 Winter2005 -- DeHon 19

Gates Required?

Depth=3, Depth=2048?

Caltech CS184 Winter2005 -- DeHon 20

Gate metric for FPGAs?

- Day8: several components for computations
 - compute element
 - interconnect:
 - space
 - time
 - instructions
- Not all applications need in same **balance**
- Assigning a single “capacity” number to device is an oversimplification

Caltech CS184 Winter2005 -- DeHon 21

MPGA vs. FPGA

- MPGA (SOG GA)
 - $5K\lambda^2/\text{gate}$
 - 35-70% usable (50%)
 - $7-17K\lambda^2/\text{gate net}$
- Xilinx XC4K
 - $1.25M\lambda^2/\text{CLB}$
 - 17-48 gates (26?)
 - $26-73K\lambda^2/\text{gate net}$
- Ratio: 2--10 (5)

Adding ~2x Custom/MPGA,
Custom/FPGA ~10x

Caltech CS184 Winter2005 -- DeHon 22

MPGA vs. FPGA

- MPGA (SOG GA)
 - $\lambda=0.6\mu$
 - $\tau_{gd}\sim 1\text{ns}$
- Xilinx XC4K
 - $\lambda=0.6\mu$
 - 1-7 gates in 7ns
 - 2-3 gates typical
- Ratio: 1--7 (2.5)

Caltech CS184 Winter2005 -- DeHon 23

Processors vs. FPGAs

Caltech CS184 Winter2005 -- DeHon 24

Processors and FPGAs

Metric: $\frac{4 \text{ input gate-evaluations}}{\lambda^2 \cdot s}$

Processor: $\frac{2 \times N_{ALU} \times w_{ALU}}{A_{proc} \times t_{cycle}}$ **FPGA:** $\frac{N_{ALUT}}{A_{array} \times t_{cycle}}$

Component Example

- Single die in 0.35μm
 - XC4085XL-09 3,136 CLBs 4.6ns
682 Bit Ops/ns
 - Alpha 1996 2×64b ALUs 2.3ns
55.7 Bit Ops/ns

[1 "bit op" = 2 gate evaluations]

Processors and FPGAs

Year	Design	Organization	λ	λ^2 area	cycle	$\frac{ge's}{\lambda^2 \cdot s}$
Microprocessors						
1984	MIPS	1 × 32	1.5μt	15M	250ns	17
1987	MIPS-X	1 × 32	1.0μt	68M	50ns	19
1994	MIPS	1 × 32	0.28μt	1.7G	2ns	19
1992	Alpha	1 × 64	0.38μt	1.7G	5ns	15
1995	Alpha	2 × 64	0.25μt	4.8G	3.3ns	18
1996	Alpha	2 × 64	0.18μt	6.8G	2.3ns	17
Reconfigurable ALUs						
1992	PADDI	8 × 16	0.6μt	126M	40ns	50
1995	PADDI-2	48 × 16	0.5μt	515M	20ns	150
FPGAs						
1986	Xilinx 2K	1 CLB (4 LUT)	1.0μt	500K	20ns	100
1988	Xilinx 3K	64 CLBs (2 4-LUT)	0.6μt	83M	13ns	120
1992	Xilinx 4K	49 CLBs (2 4-LUT)	0.6μt	61M	7ns	230
1995	Xilinx 5K	49 CLBs (4 4-LUT)	0.3μt	110M	6ns	290

Raw Density Summary

- Area
 - MPGA 2-3x Custom
 - FPGA 5x MPGA
- Area-Time
 - Gate Array 6-10x Custom
 - FPGA 15-20x Gate Array
 - Processor 10x FPGA

Raw Density Caveats

- Processor/FPGA may solve more specialized problem
- Problems have different resource balance requirements
 - ...can lead to low yield of raw density

Homework

- Day behind
- Current assignment
 - Involves cascades PLAs

Task Comparisons

Broadening Picture

- Compare larger computations
- For comparison
 - throughput density metric: results/area-time
 - normalize out area-time point selection
 - high throughput density
 - most in fixed area
 - least area to satisfy fixed throughput target

Multiply

Architecture	Feature Size (λ)	Area and Time	16×16		8×8	
			mpy λ^2/s	scale λ^2	mpy λ^2/s	scale λ^2
Custom 16×16	$0.63 \mu m$	$2.6M\lambda^2$, 40 ns	9.6	9.6	9.6	9.6
Custom 8×8	$0.80 \mu m$	$3.3M\lambda^2$, 4.3 ns			70	70
Gate-Array 16×16	$0.75 \mu m$	$26M\lambda^2$, 30ns	1.3	1.3	1.3	1.3
FPGA (XC4K)	$0.60 \mu m$	$1.25M\lambda^2/CLB$ 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	$0.65 \mu m$	$350M\lambda^2$, 60 ns	0.057	0.057	0.057	0.057
RISC (no multiplier)	$0.75 \mu m$	$125M\lambda^2$, 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

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)	$TAPs \lambda^2/s$
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)	$0.60 \mu m$	1.9
(Altera 8K)	$0.30 \mu m$	3.6
Full Custom	$0.75 \mu m$	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

IIR/Biquad

Architecture	Feature Size (λ)	Area and Time	16b TAPs λ^2/s	10b TAPs λ^2/s
16b DSP	$0.60 \mu m$	$200M\lambda^2$, 500 ns/biquad	0.010	0.010
FPGA (XC4K)	$0.60 \mu m$	60 CLBs, 320 ns/biquad 43 CLBs, 200 ns/biquad	0.044	0.093
Full Custom	$0.90 \mu m$	$68M\lambda^2$, 11.8 ns/4 biquads		5.0

Simplest IIR: $Y_i = A \times X_i + B \times Y_{i-1}$

DES Keysearch

Architecture	Feature Size (λ)	Area	Keys/Second	Keys λ^2/s
DES IC	$1.5 \mu m$	$11.1M\lambda^2$	310K	0.028
FPGA (Altera 8K)	$0.30 \mu m$	$81188 (930M\lambda^2)$	800K	0.00086
RISC	$0.30 \mu m$	$1.8G\lambda^2$	41K	0.000023

<<http://www.cs.berkeley.edu/~iang/isaac/hardware/>>

DNA Sequence Match

- **Problem:** “cost” of transform $S_1 \rightarrow S_2$
- **Given:** cost of insertion, deletion, substitution
- **Relevance:** similarity of DNA sequences
 - evolutionary similarity
 - structure predict function
- **Typically:** new sequence compared to large database

Caltech CS184 Winter2005 -- DeHon

37

DNA Sequence Match

Architecture	Feature Size (λ)	Area	Cell Updates per Second	$\frac{CU}{\lambda^2 s}$
Custom FPGA	2.0 μm	270M λ^2	500M	1.9
(SPLASH 2)	0.60 μm	43G λ^2	3,000M	0.070
(SPLASH) RISC	0.60 μm	33G λ^2	370M	0.012
(SparcStation 1)	0.75 μm	273M λ^2	0.87M	0.0032
(SparcStation 10)	0.40 μm	1.6G λ^2	1.2M	0.00075

N.B. includes memory area for SPLASH

Caltech CS184 Winter2005 -- DeHon

38

Floating-Point Add (single prec.)

Architecture	Reference	λ	area and time	32b FP ADDs/s
Custom Macrocells				
32b FP ALU	JSSC92	0.6 μm	32M λ^2 , 30 ns	1.0
32b Adder	EUROASIC92	0.4 μm	49M λ^2 , 60 ns	0.34
32/64b FPU	CICC92	0.25 μm	980M λ^2 , 2 per 12.5 ns	0.16
Coprocessor				
32b FP Processor	ISSCC88	0.75 μm	220M λ^2 , 55 ns	0.083
32b FP Processor	SSC89	0.65 μm	400M λ^2 , 23 ns	0.11
Processor FP support				
64b RISC w/FPU	ISSCC90	0.4 μm	1.4G λ^2 , 2x25 ns cycle	0.014
64b RISC w/FPU	ISSCC92	0.38 μm	1.7G λ^2 , 5 ns cycle	0.12
64b RISC w/FPU	ISSCC95	0.25 μm	4.8G λ^2 , 3.3ns cycle	0.063
FPGA				
	FCCM96	0.3 μm	500 Flex8K LEs (920K λ^2 /LE), 150 ns	0.014
	FCCM98	0.25 μm	315 XC4K LEs (1.25M λ^2 /CLB), 33 ns	0.077
Processor no FP support				
32b RISC, no FPU	ASPLoS85	1.5 μm	15M λ^2 , 16 μs	0.0042

Caltech CS184 Winter2005 -- DeHon

39

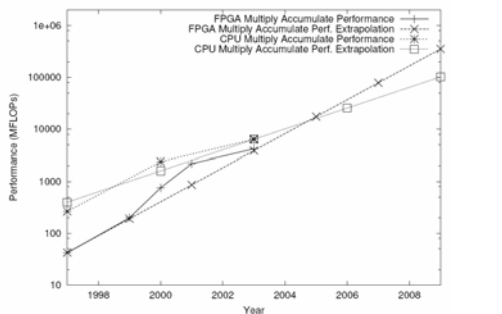
Floating-Point Mpy (single prec.)

Architecture	Reference	λ	area and time	32b FP MPYs/s
Custom Macrocells				
32b mpy	JSSC92	0.6 μm	43M λ^2 , 30 ns	0.78
32b mpy	EUROASIC92	0.4 μm	81M λ^2 , 42 ns	0.29
32/64b FPU	CICC92	0.25 μm	980M λ^2 , 2 per 12.5 ns	0.16
Coprocessor				
dedicated 32b mpy	JSSC84	1 μm	33M λ^2 , 79 ns	0.38
dedicated 32b mpy	ISSCC87	0.6 μm	160M λ^2 , 67 ns	0.093
32b FP Processor	ISSCC88	0.75 μm	200M λ^2 , 55 ns	0.083
32b FP Processor	JSSC89	0.65 μm	400M λ^2 , 23 ns	0.11
Processor FP support				
64b RISC w/FPU	ISSCC90	0.4 μm	1.4G λ^2 2x25 ns cycle	0.014
64b RISC w/FPU	ISSCC92	0.38 μm	1.7G λ^2 , 5 ns cycle	0.12
64b RISC w/FPU	ISSCC95	0.25 μm	4.8G λ^2 , 3.3ns cycle	0.063
FPGA				
	FCCM96	0.3 μm	344 Flex8K LEs (920K λ^2 /LE), 755 ns	0.0042
	FCCM98	0.25 μm	265 CLB (1.25M λ^2 /CLB), 200 ns	0.0161
Processor no FP support				
32b RISC, no FPU	ASPLoS85	1.5 μm	15M λ^2 , 7 μs	0.0096

Caltech CS184 Winter2005 -- DeHon

40

FPGA vs. Processor FP (Double Precision FP MAC)



Caltech CS184 Winter2005 -- DeHon

[Underwood/FPGA'2004] 41

Degrade from Peak

Caltech CS184 Winter2005 -- DeHon

42

Degrade from Peak: FPGAs

- Long path length → not run at cycle
- Limited throughput requirement
 - bottlenecks elsewhere limit throughput req.
- Insufficient interconnect
- Insufficient retiming resources (bandwidth)

Caltech CS184 Winter2005 -- DeHon

43

Degrade from Peak: Processors

- Ops w/ no gate evaluations (interconnect)
- Ops use limited word width
- Stalls waiting for retimed data

$$E(\text{Functional Density}) = \frac{\text{Gate Evaluations}}{\text{Datapath Bit}} \times \frac{\text{Datapath Bits}}{\text{pinst}} \times \frac{\text{pinsts}}{\text{Issue Slot}} \times \frac{1}{\text{Clock Cycle} \times \text{area} \times t_{\text{cycle}}}$$

Caltech CS184 Winter2005 -- DeHon

44

Degrade from Peak: Custom/MPGA

- Solve more general problem than required
 - (more gates than really need)
- Long path length
- Limited throughput requirement
- Not needed or applicable to a problem

Caltech CS184 Winter2005 -- DeHon

45

Degrade Notes

- We'll cover these issues in more detail as we get into them later in the course

Caltech CS184 Winter2005 -- DeHon

46

Big Ideas [MSB Ideas]

- Raw densities:
custom:ga:fpga:processor
 - 1:5:100:1000
 - close gap with specialization

Caltech CS184 Winter2005 -- DeHon

47