

CS184a: Computer Architecture (Structures and Organization)

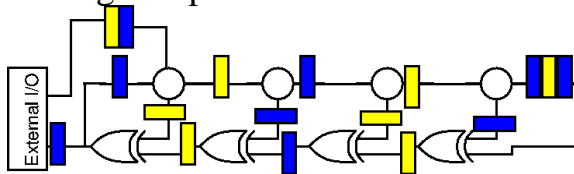
Day16: November 15, 2000
Retiming Structures

Caltech CS184a Fall2000 -- DeHon

1

Last Time

- Saw how to formulate and automate retiming:
 - start with network
 - calculate minimum achievable c
 - c = cycle delay (clock cycle)
 - make c -slow if want/need to make $c=1$
 - calculate new register placements and move



Caltech CS184a Fall2000 -- DeHon

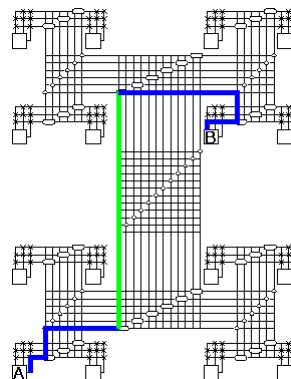
2

Today

- Systematic transformation for retiming
 - “justify” mandatory registers in design
- Retiming in the Large
- Retiming Requirements
- Retiming Structures

HSRA Retiming

- HSRA
 - adds mandatory pipelining to interconnect
- One additional twist
 - long, pipelined interconnect
 - \Rightarrow need more than one register on paths



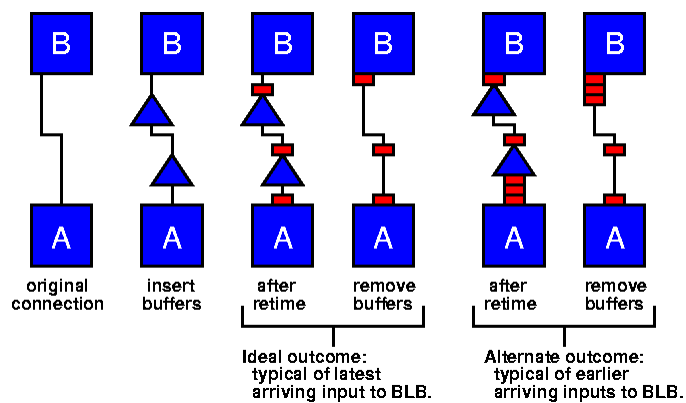
Accommodating HSRA Interconnect Delays

- Add buffers to LUT→LUT path to match interconnect register requirements
- Retime to $C=1$ as before
- Buffer chains force enough registers to cover interconnect delays

Caltech CS184a Fall2000 -- DeHon

5

Accommodating HSRA Interconnect Delays



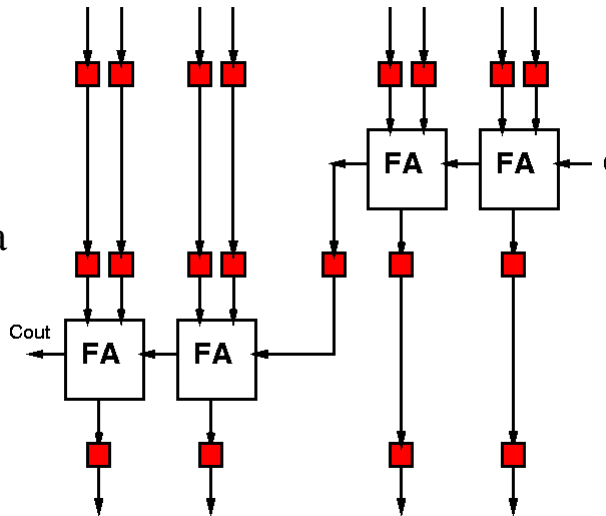
Caltech CS184a Fall2000 -- DeHon

6

Retiming in the Large

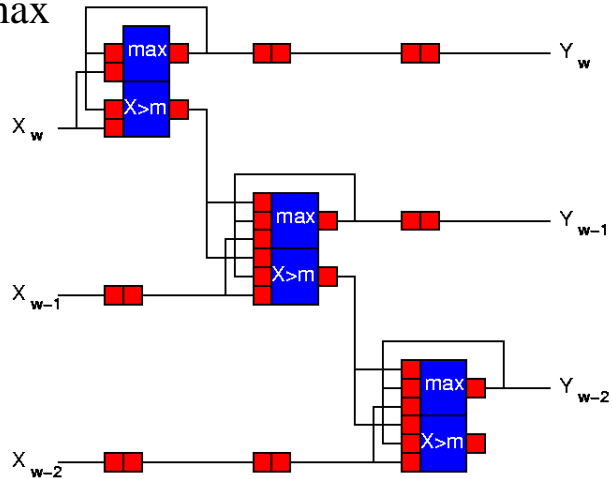
Align Data / Balance Paths

Day3:
registers
to align data



Systolic Data Alignment

- Bit-level max

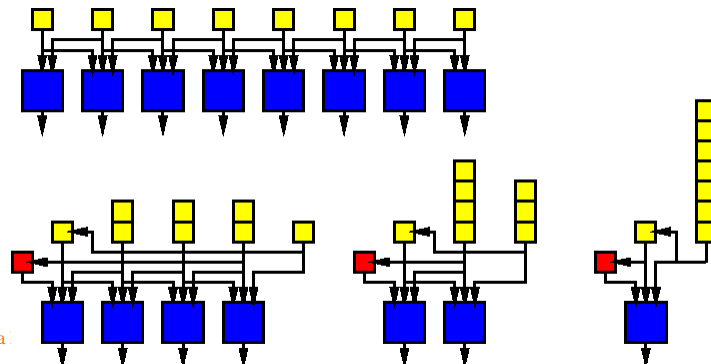


Caltech CS184a Fall2000 -- DeHon

9

Serialization

- Serialization
 - greater serialization => deeper retiming
 - total: same per compute: larger

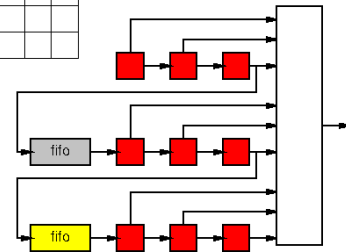
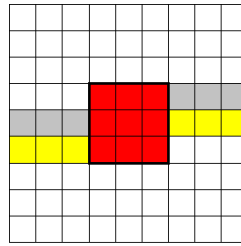


Caltech CS184a

10

Data Alignment

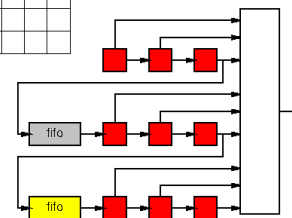
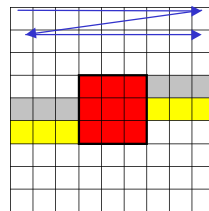
- For video (2D) processing
 - often work on local windows
 - retime scan lines
- E.g.
 - edge detect
 - smoothing
 - motion est.



Caltech CS184a Fall2000 -- DeHon

Image Processing

- See Data in raster scan order
 - adjacent, horizontal bits easy
 - adjacent, vertical bits
 - scan line apart

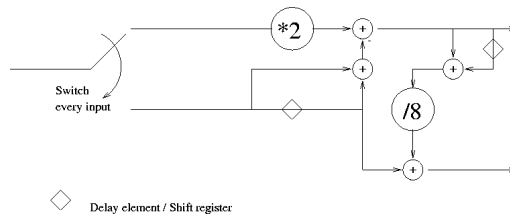


Caltech CS184a Fall2000 -- DeHon

12

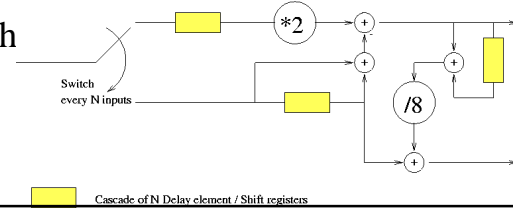
Wavelet

- Data stream for horizontal transform



- Data stream for vertical transform

– N =image width



Caltech CS184a Fall2000 -- DeHon

13

Retiming in the Large

- Aside from the local retiming for cycle optimization (last time)
- Many intrinsic needs to retime data for correct use of compute engine
 - some very deep
 - often arise from serialization

Caltech CS184a Fall2000 -- DeHon

14

Reminder: Temporal Interconnect

- Retiming \equiv Temporal Interconnect
- Function of *data* memory
 - perform retiming

Requirements not Unique

- Retiming requirements are not unique to the problem
- Depends on algorithm/implementation
- Behavioral transformations can alter significantly

Requirements Example

$$Q = A * B + C * D + E * F$$

- For $I \leftarrow 1$ to N
 - $t1[I] \leftarrow A[I] * B[I]$
- For $I \leftarrow 1$ to N
 - $t2[I] \leftarrow C[I] * D[I]$
- For $I \leftarrow 1$ to N
 - $t3[I] \leftarrow E[I] * F[I]$
- For $I \leftarrow 1$ to N
 - $t2[I] \leftarrow t1[I] + t2[I]$
- For $I \leftarrow 1$ to N
 - $Q[I] \leftarrow t2[I] + t3[I]$
- For $I \leftarrow 1$ to N
 - $t1 \leftarrow A[I] * B[I]$
 - $t2 \leftarrow C[I] * D[I]$
 - $t1 \leftarrow t1 + t2$
 - $t2 \leftarrow E[I] * F[I]$
 - $Q[I] \leftarrow t1 + t2$
- left $\Rightarrow 3N$ regs
- right $\Rightarrow 2$ regs

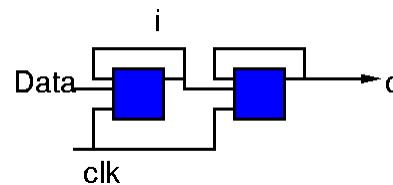
Retiming Structure and Requirements

Structures

- How do we implement programmable retiming?
- Concerns:
 - Area: λ^2/bit
 - Throughput: bandwidth (bits/time)
 - Latency important when do not know when we will need data item again

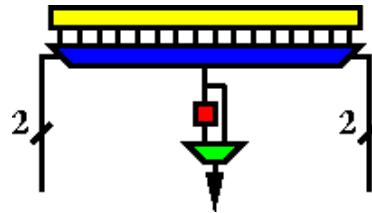
Just Logic Blocks

- Most primitive
 - build flip-flop out of logic blocks
 - $I \leftarrow D*/Clk + I*Clk$
 - $Q \leftarrow Q*/Clk + I*Clk$
 - Area: 2 LUTs (800K \rightarrow 1M λ^2 /LUT each)
 - Bandwidth: 1b/cycle



Optional Output

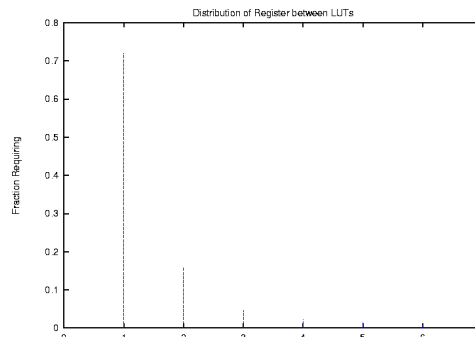
- Real flip-flop (optionally) on output



- flip-flop: $4-5K\lambda^2$
- Switch to select: $\sim 5K\lambda^2$
- Area: 1 LUT ($800K \rightarrow 1M\lambda^2/\text{LUT}$)
- Bandwidth: 1b/cycle

Output Flip-Flop Needs

- Pipeline and C-slow to LUT cycle
- Always need an output register



Number of Registers	1	2	3	4	5	6	7	8	9	10
Percentage	72	16	4.5	2.2	1.3	0.96	1.2	0.46	0.12	0.11

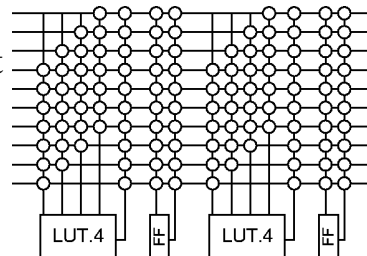
Table 3: Benchmark Wide Distribution of Registers Required between LUTs
Average Regs/LUT 1.7, some designs need 2--7x

Separate Flip-Flops

- Network flip flop w/ own interconnect

+ can deploy where needed

- requires more interconnect



- Assume routing goes as inputs

- 1/4 size of LUT

- Area: $200K\lambda^2$ each

- Bandwidth: 1b/cycle

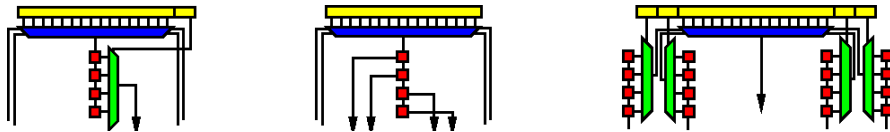
Deeper Options

- Interconnect / Flip-Flop is expensive
- How do we avoid?

Deeper

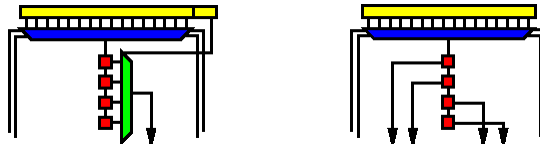
- Implication
 - don't need result on every cycle
 - number of regs > bits need to see each cycle
 - => lower bandwidth acceptable
 - => less interconnect

Deeper Retiming



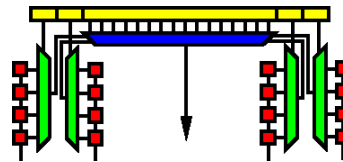
Output

- Single Output
 - Ok, if don't need other timings of signal
- Multiple Output
 - more routing

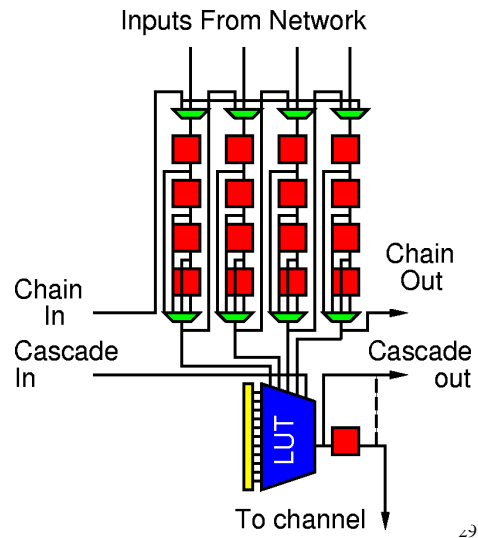
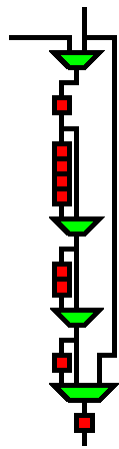


Input

- More registers ($K \times$)
 - $7-10K\lambda^2/\text{register}$
 - $4\text{-LUT} \Rightarrow 30-40K\lambda^2/\text{depth}$
- No more interconnect than unretimed
 - *open*: compare savings to additional reg. cost
 - Area: 1 LUT ($1M+d*40K\lambda^2$) get Kd regs
 - $d=4, 1.2M\lambda^2$
 - Bandwidth: 1b/cycle
 - $1/d$ th capacity



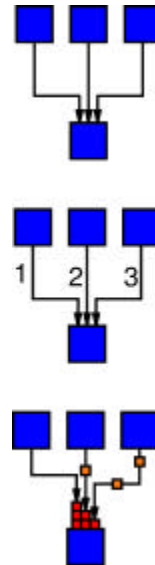
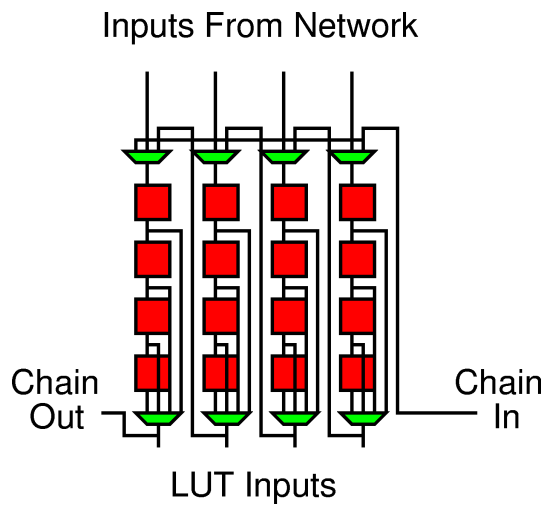
HSRA Input



Caltech CS184a Fall2000 -- DeHon

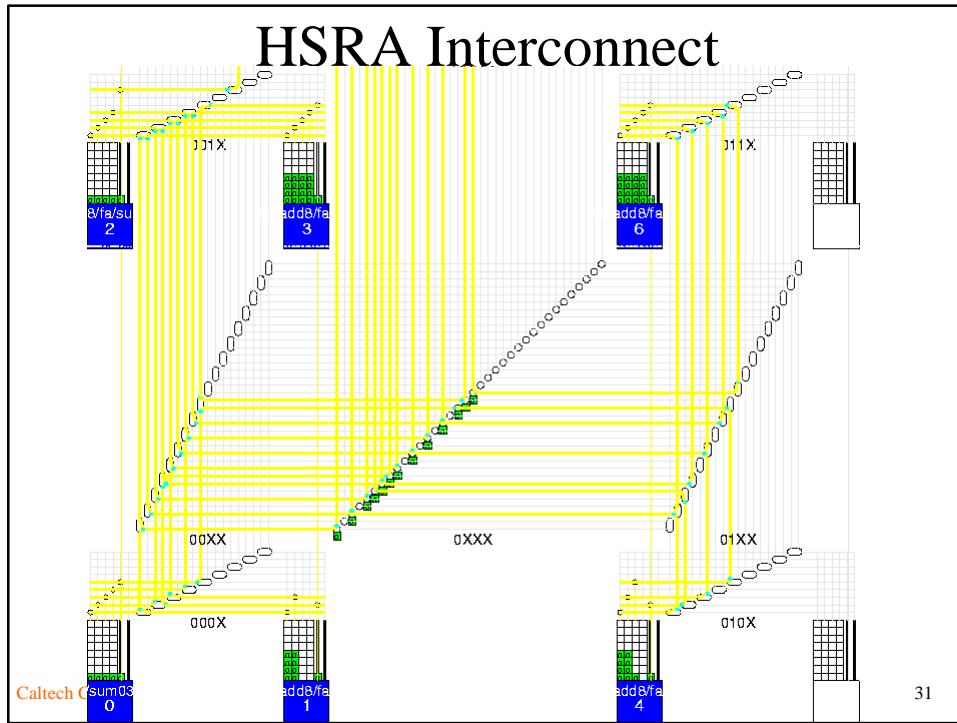
29

Input Retiming



Caltech CS184a Fall2000 -- DeHon

30



Flop Experiment #1

- Pipeline and retime to single LUT delay per cycle
 - MCNC benchmarks to 256 4-LUTs
 - no interconnect accounting

Number of Registers	1	2	3	4	5	6	7	8	9	10
Percentage	72	16	4.5	2.2	1.3	0.96	1.2	0.46	0.12	0.11

– average 1.7 registers/LUT (some circuits 2--7)

Flop Experiment #2

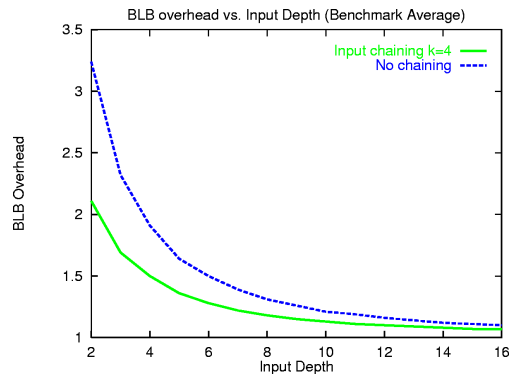
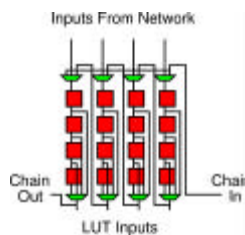
- Pipeline and retime to HSRA cycle
 - place on HSRA
 - single LUT or interconnect timing domain
 - same MCNC benchmarks

Number of Registers	1	2	3	4	5	6	7	8	9	10	>10
Percentage	60	6.9	5.9	3.8	4.3	2.7	2.6	1.9	1.5	1.2	9.2

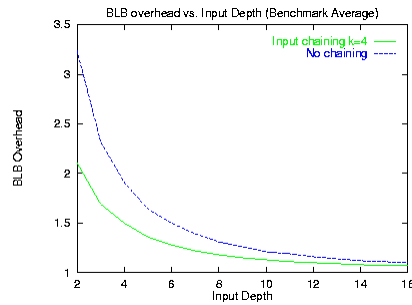
- average 4.7 registers/LUT

Input Depth Optimization

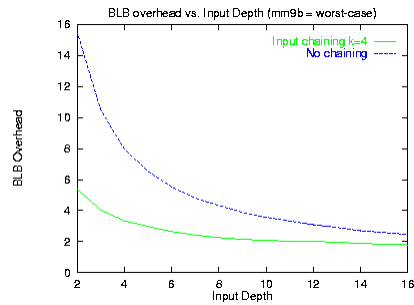
- Real design, fixed input retiming depth
 - truncate deeper and allocate additional logic blocks



Extra Blocks (limited input depth)



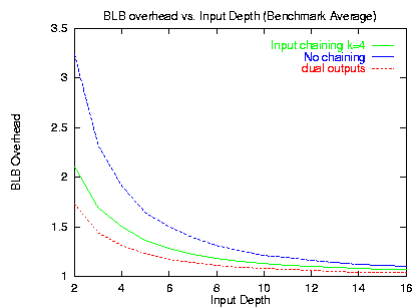
Average



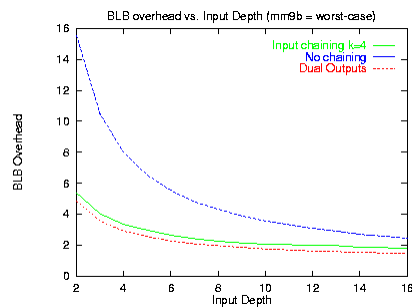
Worst Case Benchmark

With Chained Dual Output

[can use one BLB as 2 retiming-only chains]

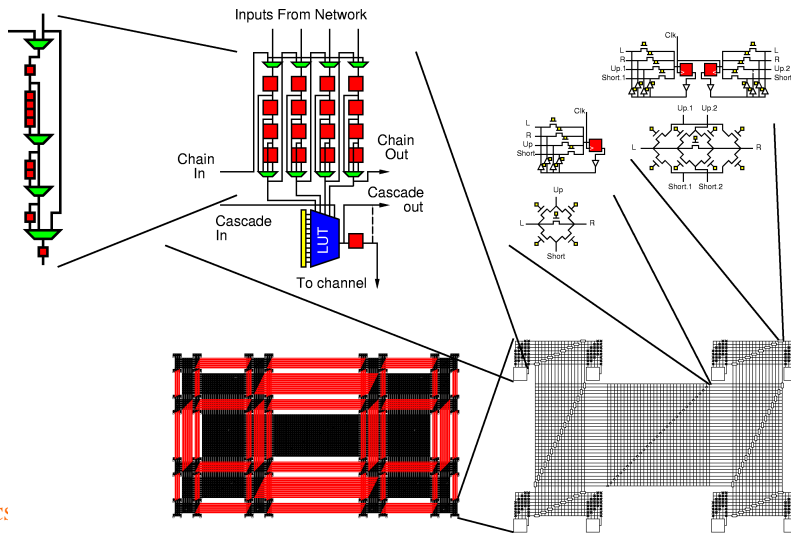


Average



Worst Case Benchmark

HSRA Architecture

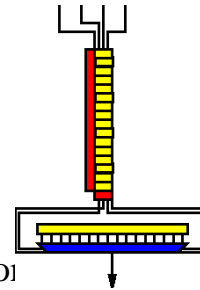


Caltech C!

37

Register File

- From MIPS-X
 - $1K\lambda^2/\text{bit} + 500\lambda^2/\text{port}$
 - $\text{Area}(\text{RF}) = (d+6)(W+6)(1K\lambda^2 + \text{port})$
- $w \gg 6, d \gg 6 \quad I+o=2 \Rightarrow 2K\lambda^2/\text{bit}$
- $w=1, d \gg 6 \quad I=0=4 \Rightarrow 35K\lambda^2/\text{bit}$
 - comparable to input chain
- More efficient for wide-word cases

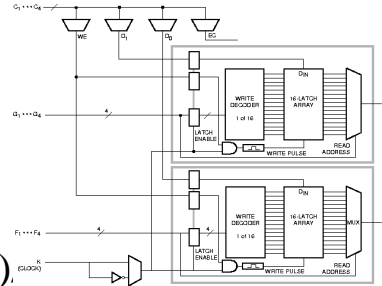


Caltech CS184a Fall2000 -- DeHon

38

Xilinx CLB

- Xilinx 4K CLB
 - as memory
 - works like RF
- Area: $1/2$ CLB ($640K\lambda^2$),
 - but need 4 CLBs to control
- Bandwidth: $1b/2$ cycle ($1/2$ CLB)
 - $1/16$ th capacity



Memory Blocks

- SRAM bit $\approx 1200\lambda^2$ (large arrays)
- DRAM bit $\approx 100\lambda^2$ (large arrays)
- Bandwidth: W bits / 2 cycles
 - usually single read/write
 - $1/2^A$ th capacity

Disk Drive

- Cheaper per bit than DRAM/Flash
 - (not MOS, no λ^2) 😊
- Bandwidth: 10-20Mb/s
 - For 4ns array cycle
 - 1b/12.5 cycles @20Mb/s

Hierarchy/Structure Summary

- “Memory Hierarchy” arises from area/bandwidth tradeoffs
 - Smaller/cheaper to store words/blocks
 - (saves routing and control)
 - Smaller/cheaper to handle long retiming in larger arrays (reduce interconnect)
 - High bandwidth out of registers/shallow memories

λ^2	DRAM	SRAM	RF bit	FF/RF	RF \times 1	XC	In FF	net FF	FF/LUT
bw/cap.	100	1200	2K	5K	40K	40K	75K	200K	800K
	$1/10^7$	$1/10^5-10^3$		$1/100$	$1/100$	$1/16$	$1/4$	$1/1$	$1/1$

Big Ideas [MSB Ideas]

- Can systematically justify registers in architecture (interconnect, FU pipeline)

Big Ideas [MSB Ideas]

- Tasks have a wide variety of retiming distances
- Retiming requirements affected by high-level decisions/strategy in solving task
- Wide variety of retiming costs
 - $100 \lambda^2 \rightarrow 1M\lambda^2$
- Routing and I/O bandwidth
 - big factors in costs
- Gives rise to memory (retiming) hierarchy