

# ESE680-002 (ESE534): Computer Organization

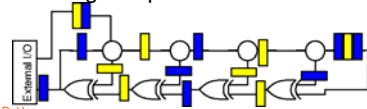
Day 20: March 28, 2007  
Retiming 2:  
Structures and Balance



Penn ESE680-002 Spring2007 -- DeHon

## Last Time

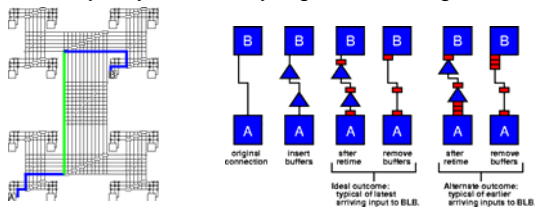
- Saw how to formulate and automate retiming:
  - start with network
  - calculate minimum achievable  $c$ 
    - $c$  = cycle delay (clock cycle)
  - make  $c$ -slow if want/need to make  $c=1$
  - calculate new register placements and move



Penn ESE680-002 Spring2007 -- DeHon

## Last Time

- Systematic transformation for retiming
  - “justify” mandatory registers in design



Penn ESE680-002 Spring2007 -- DeHon

3

## Today

- Retiming in the Large
- Retiming Requirements
- Retiming Structures

Penn ESE680-002 Spring2007 -- DeHon

4

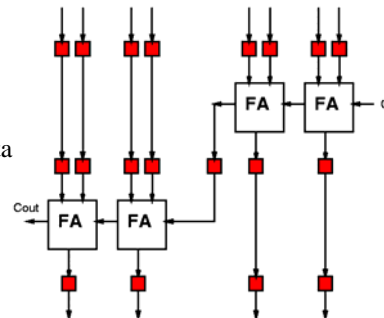
## Retiming in the Large

Penn ESE680-002 Spring2007 -- DeHon

5

## Align Data / Balance Paths

Day3:  
registers  
to align data

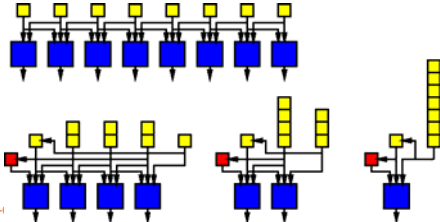


Penn ESE680-002 Spring2007 -- DeHon

6

## Serialization

- Serialization
  - greater serialization → deeper retiming
  - **total:** same     **per compute:** larger

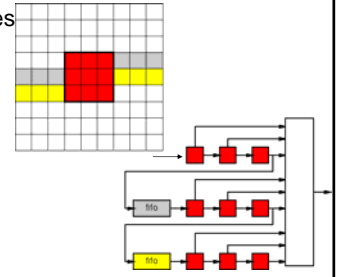


Penn ESE680-4

7

## Data Alignment

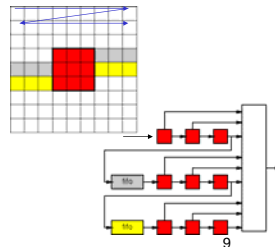
- For video (2D) processing
  - often work on local windows
  - retime scan lines
- *E.g.*
  - edge detect
  - smoothing
  - motion est.



Penn ESE680-002 Spring2007 -- DeHon

## Image Processing

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



Penn ESE680-002 Spring2007 -- DeHon

9

## 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

Penn ESE680-002 Spring2007 -- DeHon

10

## Reminder: Temporal Interconnect

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

Penn ESE680-002 Spring2007 -- DeHon

11

## Requirements not Unique

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

Penn ESE680-002 Spring2007 -- DeHon

12

## Requirements Example

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

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

Penn ESE680-002 Spring2007 -- DeHon

13

## Retiming Requirements

Penn ESE680-002 Spring2007 -- DeHon

14

## Flop Experiment #1

- Pipeline/C-slow/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)

Penn ESE680-002 Spring2007 -- DeHon

15

## 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

Penn ESE680-002 Spring2007 -- DeHon

16

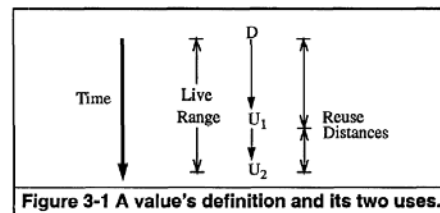
## Value Reuse Profiles

- What is the distribution of retiming distances needed?
  - Balance of retiming and compute
  - Fraction which need various depths
  - Like wire-length distributions....

Penn ESE680-002 Spring2007 -- DeHon

17

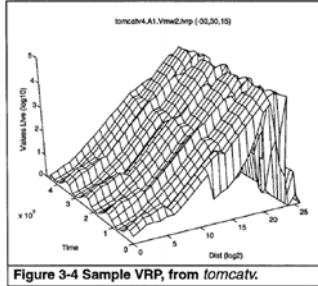
## Value Reuse Profiles



Penn ESE680-002 Spring2007 -- DeHon

[Huang&Shen/Micro 1995] 18

## Example Value Reuse Profile



Penn ESE680-002 Spring2007 -- DeHon

[Huang&Shen/Micro 1995]

19

## Interpreting VRP

- Huang and Shen data assume small number of Ops per cycle
- What happens if exploit more parallelism?
  - Values reused more frequently
  - Distances shorten

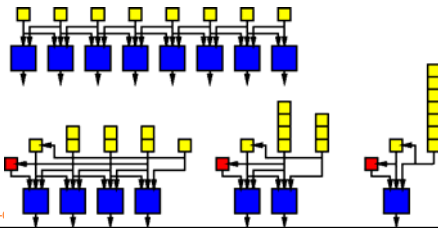
Penn ESE680-002 Spring2007 -- DeHon

20

Recall

## Serialization

- Serialization
  - greater serialization → deeper retiming
  - **total:** same    **per compute:** larger



Penn ESE680-002

21

## Idea

- Task, implemented with a given amount of parallelism
  - Will have a distribution of retiming requirements
  - May differ from task to task
  - May vary independently from compute/interconnect requirements
  - Another balance issue to watch
  - May need a canonical way to measure
    - Like Rent?

Penn ESE680-002 Spring2007 -- DeHon

22

## Retiming Structure

Penn ESE680-002 Spring2007 -- DeHon

23

## Structures

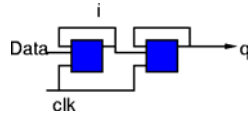
- How do we implement programmable retiming?
- Concerns:
  - Area:  $\lambda^2/\text{bit}$
  - Throughput: bandwidth (bits/time)
  - Latency important when do not know when we will need data item again

Penn ESE680-002 Spring2007 -- DeHon

24

## 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 → 1Mλ<sup>2</sup>/LUT each)
  - Bandwidth: 1b/cycle

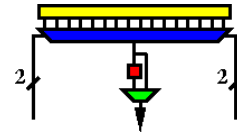


Penn ESE680-002 Spring2007 -- DeHon

25

## Optional Output

- Real flip-flop (optionally) on output



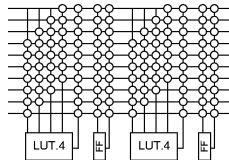
- flip-flop: 4-5Kλ<sup>2</sup>
- Switch to select: ~ 5Kλ<sup>2</sup>
- Area: 1 LUT (800K → 1Mλ<sup>2</sup>/LUT)
- Bandwidth: 1b/cycle

Penn ESE680-002 Spring2007 -- DeHon

26

## Separate Flip-Flops

- Network flip flop w/ own interconnect
  - + can deploy where needed
  - requires more interconnect
  - + Vary LUT/FF ratio
    - Arch. Parameter
  - Assume routing ∝ inputs
    - 1/4 size of LUT
  - Area: 200Kλ<sup>2</sup> each
  - Bandwidth: 1b/cycle



Penn ESE680-002 Spring2007 -- DeHon

27

## Deeper Options

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

Penn ESE680-002 Spring2007 -- DeHon

28

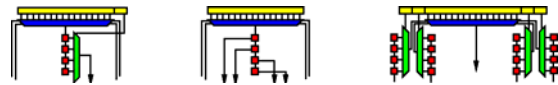
## Deeper

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

Penn ESE680-002 Spring2007 -- DeHon

29

## Deeper Retiming



Penn ESE680-002 Spring2007 -- DeHon

30

## Output

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

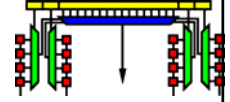


Penn ESE680-002 Spring2007 -- DeHon

31

## Input

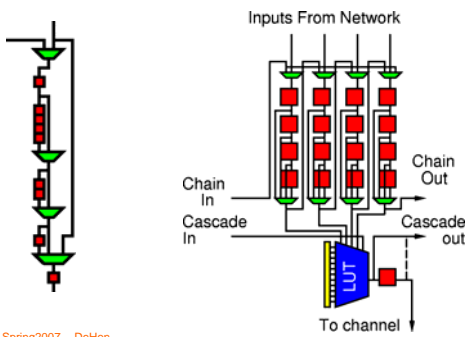
- More registers ( $K \times$ )
  - $7-10K\lambda^2/\text{register}$
  - $4\text{-LUT} \Rightarrow 30-40K\lambda^2/\text{depth}$
- No more interconnect than unretimed
  - **open**: compare savings to additional reg. cost
  - Area: 1 LUT ( $1M+d*40K\lambda^2$ ) get  $Kd$  regs
    - $d=4, 1.2M\lambda^2$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



Penn ESE680-002 Spring2007 -- DeHon

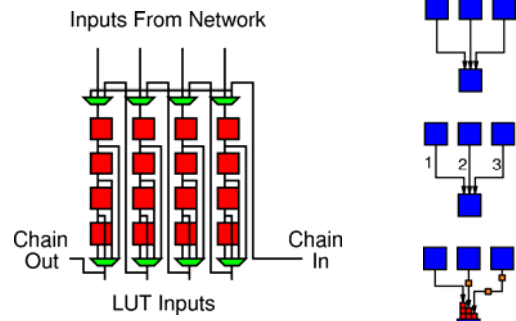
32

## HSRA Input



Penn ESE680-002 Spring2007 -- DeHon

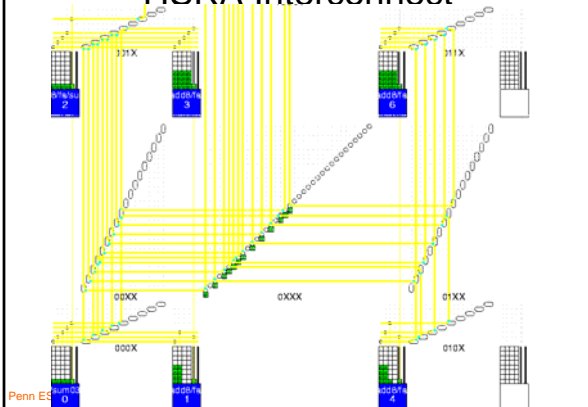
## Input Retiming



Penn ESE680-002 Spring2007 -- DeHon

34

## HSRA Interconnect



Penn ESE680-002 Spring2007 -- DeHon

Recall

## 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

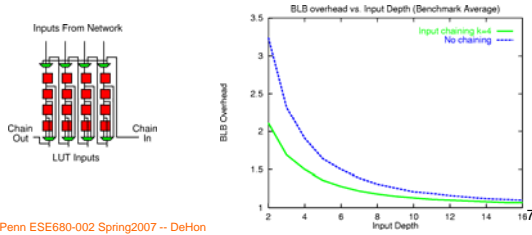
– average 4.7 registers/LUT

Penn ESE680-002 Spring2007 -- DeHon

36

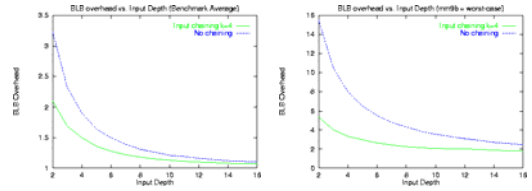
## Input Depth Optimization

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



Penn ESE680-002 Spring2007 -- DeHon

## Extra Blocks (limited input depth)



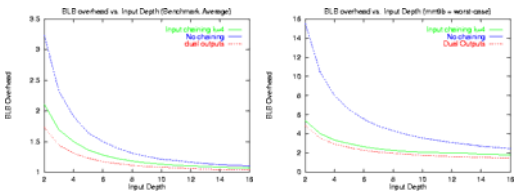
Average

Worst Case Benchmark

Penn ESE680-002 Spring2007 -- DeHon

38

## With Chained Dual Output [can use one BLB as 2 retiming-only chains]



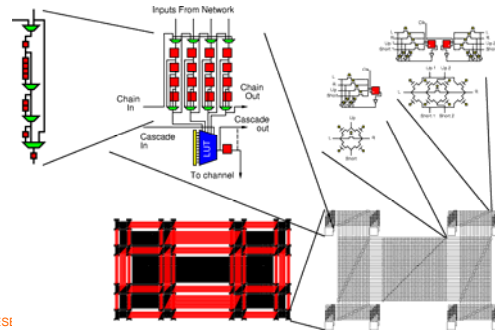
Average

Worst Case Benchmark

Penn ESE680-002 Spring2007 -- DeHon

39

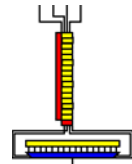
## HSRA Architecture



Penn ESI

## Register File

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

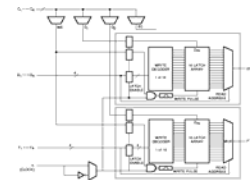


Penn ESE680-002 Spring2007 -- DeHon

41

## Xilinx CLB

- Xilinx 4K CLB
  - as memory
  - works like RF
- Area:  $1/2 \text{ CLB } (640K\lambda^2)/16 \approx 40K\lambda^2/\text{bit}$ 
  - but need 4 CLBs to control
- Bandwidth:  $1b/2 \text{ cycle } (1/2 \text{ CLB})$ 
  - $1/16 \text{ th capacity}$



Penn ESE680-002 Spring2007 -- DeHon

42

## Virtex SRL16

- Xilinx Virtex 4-LUT
  - Use as 16b shiftreg
- Area:  $\sim 1M\lambda^2/16 \approx 60K\lambda^2/b_{LU}$ 
  - Does not need CLBs to control
- Bandwidth: 1b/2 cycle (1/2 CLB)
  - 1/16 th capacity

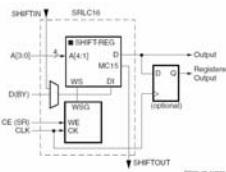


Figure 21: Shift Register Configurations

## 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  $\lambda^2$ )
- Bandwidth: 150MB/s
  - For 1ns array cycle
    - $\sim 1b/cycle @ 1.2Gb/s$

## 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

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

## Modern FPGAs

- Output Flop (depth 1)
- Use LUT as Shift Register (16)
- Embedded RAMs (16Kb)
- Interface off-chip DRAM ( $\sim 0.1-1Gb$ )
- No retiming in interconnect
  - ....yet

## Modern Processors

- DSPs have accumulator (depth 1)
- Inter-stage pipelines (depth 1)
  - Lots of pipelining in memory path...
- Reorder Buffer (4–32)
- Architected RF (16, 32, 128)
- Actual RF (256, 512...)
- L1 Cache ( $\sim 64Kb$ )
- L2 Cache ( $\sim 1Mb$ )
- L3 Cache (10-100Mb)
- Main Memory in DRAM ( $\sim 10-100Gb$ )



## Big Ideas [MSB Ideas]

- Tasks have a wide variety of retiming distances (depths)
- 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