

CS184a: Computer Architecture (Structure and Organization)

Day 19: February 23, 2005
Retime 1: Transformations



Caltech CS184 Winter2005 -- DeHon

Previously

- Reviewed Pipelining
 - basic assignments on
- Saw spatial designs efficient
 - when reuse logic at maximum frequency
- Interconnect is dominant delay
 - and dominant area
 - heavy call to reuse to use efficiently

2

Caltech CS184 Winter2005 -- DeHon

Today

- Systematic transformation for retiming
 - preserve semantics (meaning)

3

Caltech CS184 Winter2005 -- DeHon

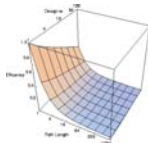
Motivation

4

Caltech CS184 Winter2005 -- DeHon

Motivation

- FPGAs (spatial computing)
 - run efficiently when all resources reused rapidly
 - cycle time minimized



- “Everything in the right place at the right time.”

5

Caltech CS184 Winter2005 -- DeHon

Motivating Questions

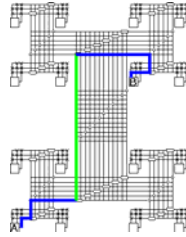
- Can I build a fixed-frequency (fixed clock) programmable architecture?
- Can I always make a design run at maximum clock rate?
- How do we systematically transform any computation to
 - Operate on fixed-frequency array?
 - Coordinate around mandatory registers in design?

6

Caltech CS184 Winter2005 -- DeHon

Interconnect Retiming

- Long Paths Slow
- Could limit cycle
- Add registers to long distance interconnect
 - At each switch?
 - In the middle of long wires?
- How justify these registers?

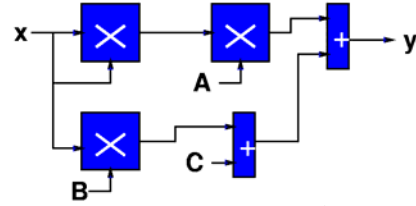


Caltech CS184 Winter2005 -- DeHon

7

Day 3

Spatial Quadratic



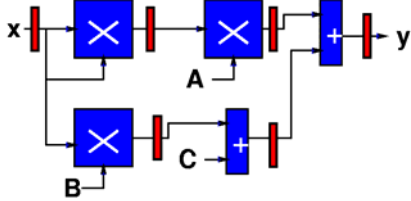
- How do we pipeline a design?

Caltech CS184 Winter2005 -- DeHon

8

Day 3

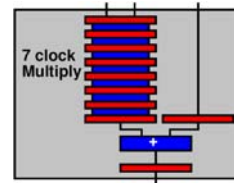
Pipelined Spatial Quadratic



Caltech CS184 Winter2005 -- DeHon

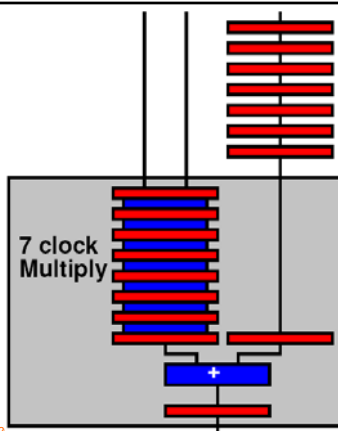
9

How do you use?



Caltech CS184 Winter2005 -- DeHon

10

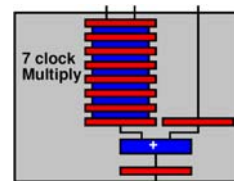


Caltech CS184 Winter2

11

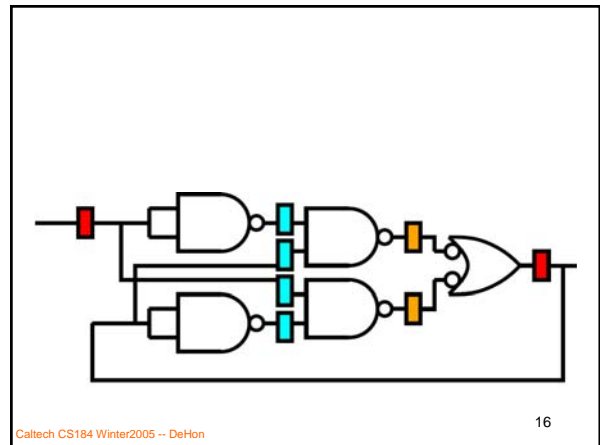
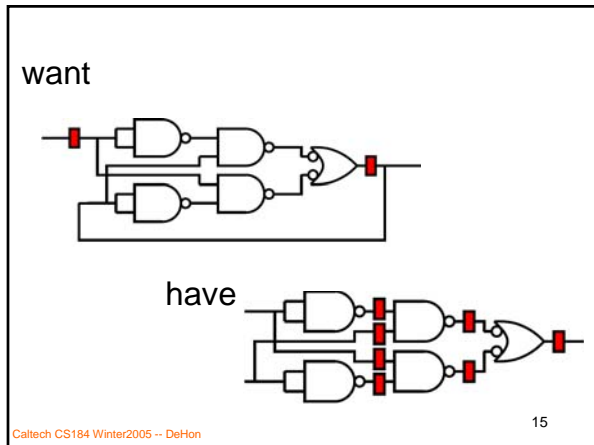
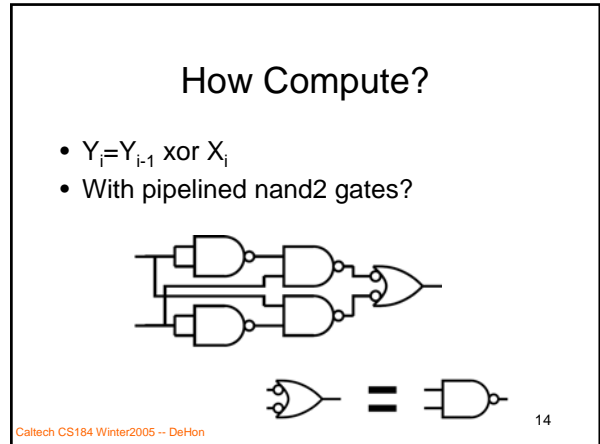
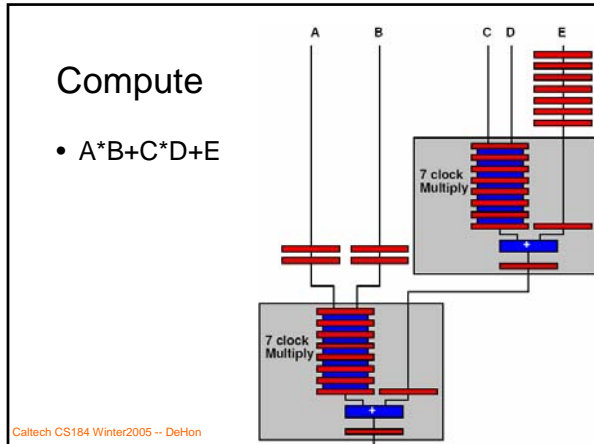
How do you use?

- To compute $A*B+C*D+E$



Caltech CS184 Winter2005 -- DeHon

12



Retiming Algorithm

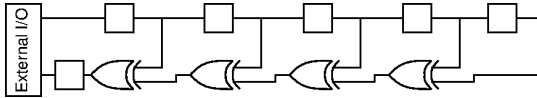
Caltech CS184 Winter2005 -- DeHon

Task

- Move registers to:
 - Preserve semantics
 - Minimize path length between registers
 - *i.e.* Make path length 1 for maximum throughput or reuse
 - ...while minimizing number of registers required

Caltech CS184 Winter2005 -- DeHon

Simple Example



Path Length (L) = 4

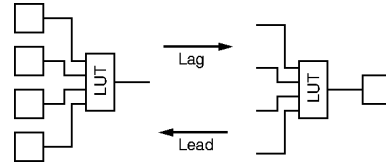
Can we do better?

Caltech CS184 Winter2005 -- DeHon

19

Legal Register Moves

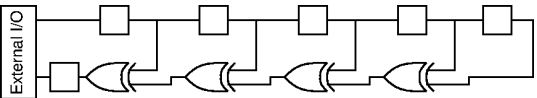
- Retiming Lag/Lead



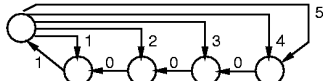
Caltech CS184 Winter2005 -- DeHon

20

Canonical Graph Representation



Observable I/O

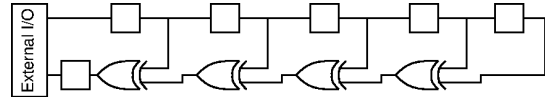


Separate arc for each path
Weight edges by number of registers
(weight nodes by delay through node)

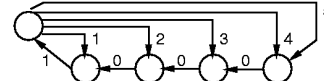
Caltech CS184 Winter2005 -- DeHon

21

Critical Path Length



Observable I/O

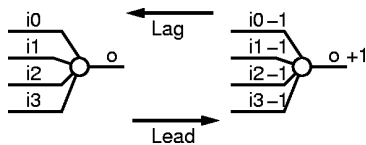


Critical Path: Length of longest path of zero weight nodes
Compute in $O(|E|)$ time by leveling network:

Topological sort, push path lengths forward until find register.

Caltech CS184 Winter2005 -- DeHon

Retiming Lag/Lead



Retiming: Assign a lag to every vertex

$$\text{weight}(e') = \text{weight}(e) + \text{lag}(\text{head}(e)) - \text{lag}(\text{tail}(e))$$

Caltech CS184 Winter2005 -- DeHon

23

Valid Retiming

- Retiming is valid as long as:
 - $\forall e$ in graph
 - $\text{weight}(e') = \text{weight}(e) + \text{lag}(\text{head}(e)) - \text{lag}(\text{tail}(e)) \geq 0$
- Assuming original circuit was a valid synchronous circuit, this guarantees:
 - non-negative register weights on all edges
 - no travel backward in time :-)
 - all cycles have strictly positive register counts
 - propagation delay on each vertex is non-negative (assumed 1 for today)

Caltech CS184 Winter2005 -- DeHon

24

Retiming Task

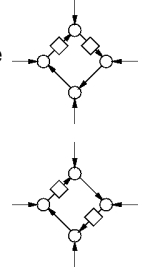
- Move registers \equiv assign lags to nodes
 - lags define all locally legal moves
- Preserving non-negative edge weights
 - (previous slide)
 - guarantees collection of lags remains consistent globally

Caltech CS184 Winter2005 -- DeHon

25

Retiming Transformation

- *N.B.*: unchanged by retiming
 - number of registers around a cycle
 - delay along a cycle
- Cycle of length P must have
 - at least P/c registers on it
 - to be retimeable to cycle c



Caltech CS184 Winter2005 -- DeHon

26

Optimal Retiming

- There is a retiming of
 - graph G
 - w/ clock cycle c
 - iff $G-1/c$ has no cycles with negative edge weights
- $G-\alpha \equiv$ subtract α from each edge weight

Caltech CS184 Winter2005 -- DeHon

27

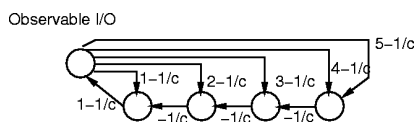
1/c Intuition

- Want to place a register every c delay units
- Each register adds one
- Each delay subtracts $1/c$
- As long as remains more positives than negatives around all cycles
 - can move registers to accommodate
 - Captures the $\text{regs}=P/c$ constraints

Caltech CS184 Winter2005 -- DeHon

28

$G-1/c$



Caltech CS184 Winter2005 -- DeHon

29

Compute Retiming

- $\text{Lag}(v) =$ shortest path to I/O in $G-1/c$
- Compute shortest paths in $O(|V||E|)$
 - Bellman-Ford
 - also use to detect negative weight cycles when c too small

Caltech CS184 Winter2005 -- DeHon

30

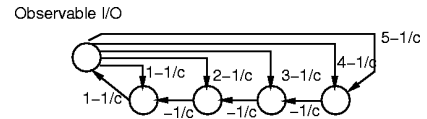
Bellman Ford

- For $k \leftarrow 0$ to N
 - $u_i \leftarrow -\infty$ (except $u_i=0$ for I/O)
- For $k \leftarrow 0$ to N
 - for $e_{ij} \in E$
 - $u_i \leftarrow \min(u_i, u_j + w(e_{ij}))$
- For $e_{ij} \in E$ //still update \rightarrow negative cycle
 - if $u_i > u_j + w(e_{ij})$
 - cycles detected

Caltech CS184 Winter2005 -- DeHon

31

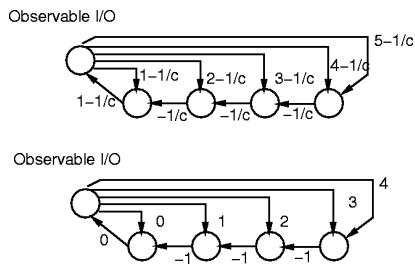
Apply to Example



Caltech CS184 Winter2005 -- DeHon

32

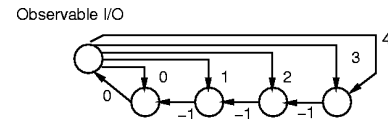
Try $c=1$



Caltech CS184 Winter2005 -- DeHon

33

Apply: Find Lags

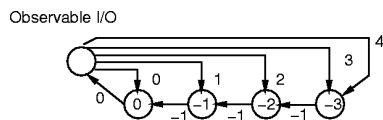


Negative weight cycles?
Shortest paths?

Caltech CS184 Winter2005 -- DeHon

34

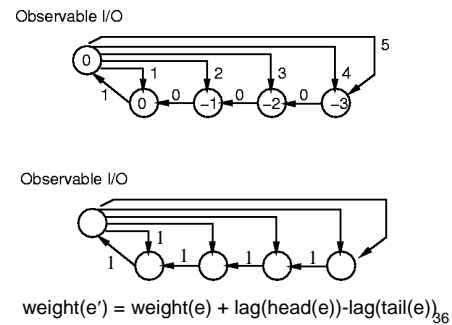
Apply: Lags



Caltech CS184 Winter2005 -- DeHon

35

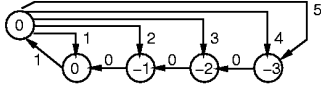
Apply: Move Registers



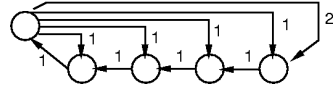
Caltech CS184 Winter2005 -- DeHon

Apply: Retimed

Observable I/O



Observable I/O

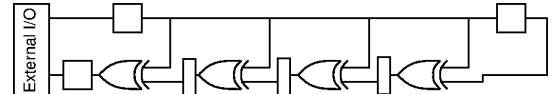
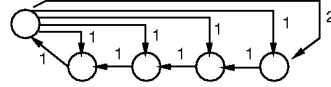


Caltech CS184 Winter2005 -- DeHon

37

Apply: Retimed Design

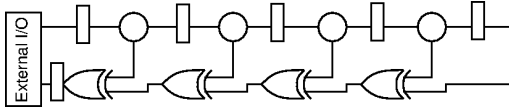
Observable I/O



Caltech CS184 Winter2005 -- DeHon

38

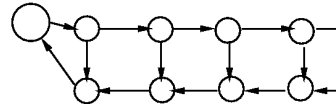
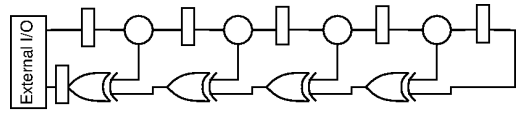
Revise Example (fanout delay)



Caltech CS184 Winter2005 -- DeHon

39

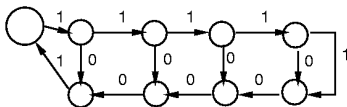
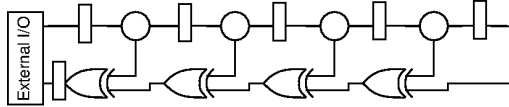
Revised: Graph



Caltech CS184 Winter2005 -- DeHon

40

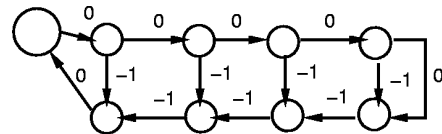
Revised: Graph



Caltech CS184 Winter2005 -- DeHon

41

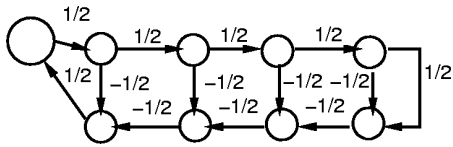
Revised: C=1?



Caltech CS184 Winter2005 -- DeHon

42

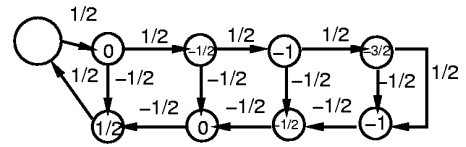
Revised: C=2?



Caltech CS184 Winter2005 -- DeHon

43

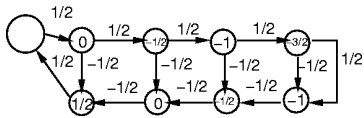
Revised: Lag



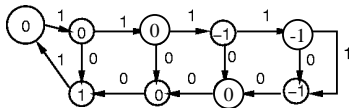
Caltech CS184 Winter2005 -- DeHon

44

Revised: Lag



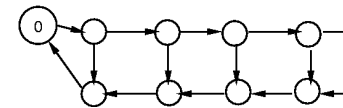
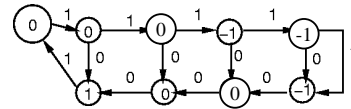
Take ceiling to convert to integer lags:



Caltech CS184 Winter2005 -- DeHon

45

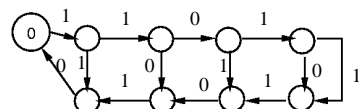
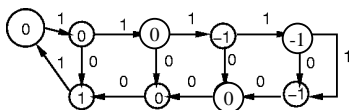
Revised: Apply Lag



Caltech CS184 Winter2005 -- DeHon

46

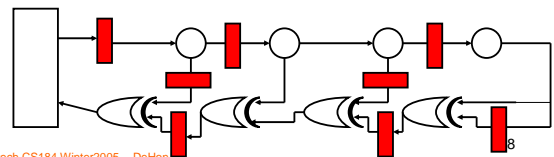
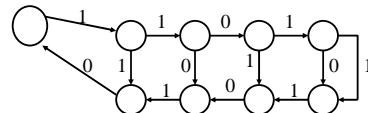
Revised: Apply Lag



Caltech CS184 Winter2005 -- DeHon

47

Revised: Retimed



Caltech CS184 Winter2005 -- DeHon

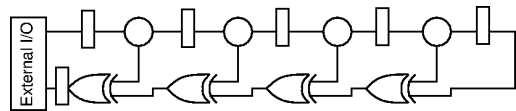
Pipelining

- We can use this retiming to pipeline
- Assume we have enough (infinite supply) registers at edge of circuit
- Retime them into circuit

Caltech CS184 Winter2005 -- DeHon

49

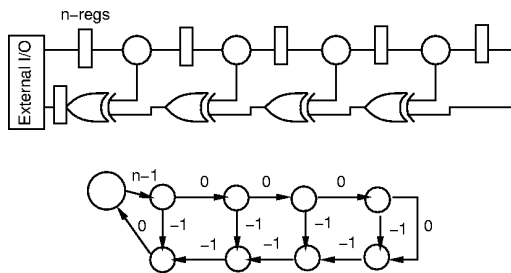
$C > 1 \implies$ Pipeline



Caltech CS184 Winter2005 -- DeHon

50

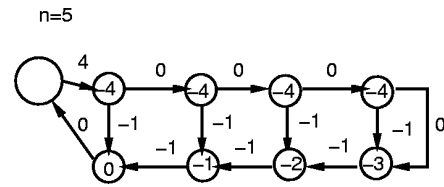
Add Registers



Caltech CS184 Winter2005 -- DeHon

51

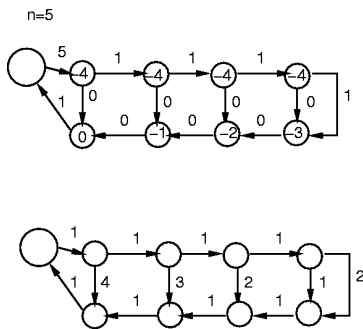
Pipeline Retiming: Lag



Caltech CS184 Winter2005 -- DeHon

52

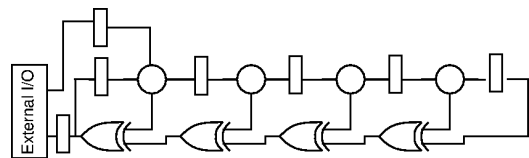
Pipelined Retimed



Caltech CS184 Winter2005 -- DeHon

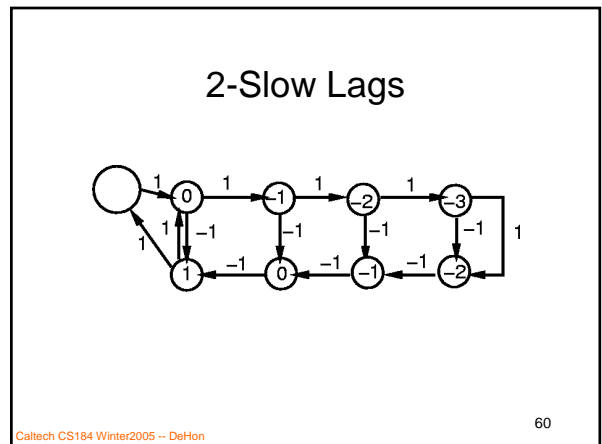
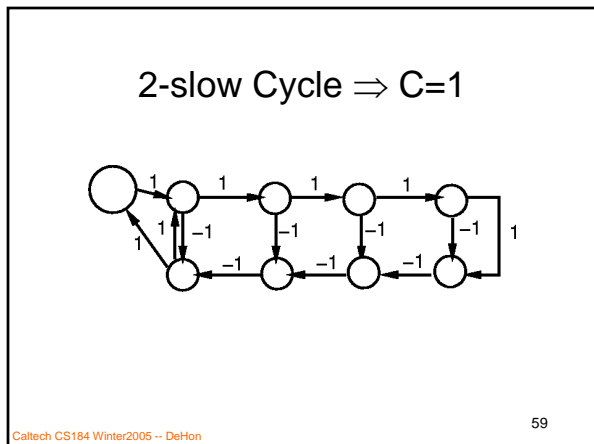
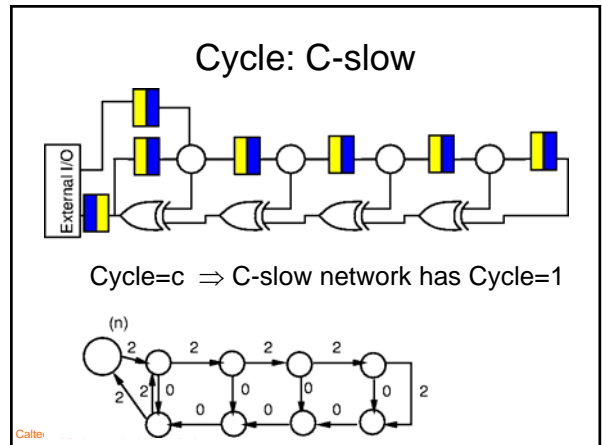
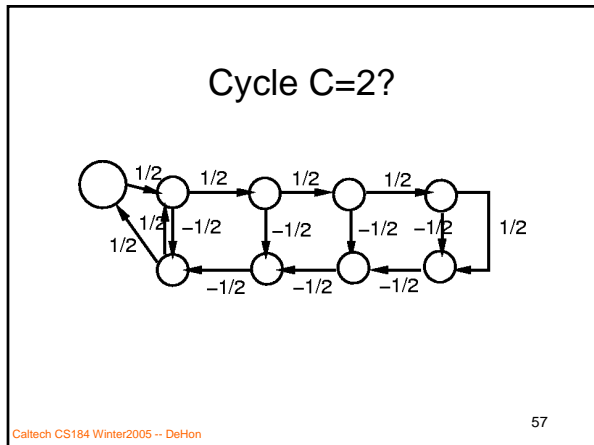
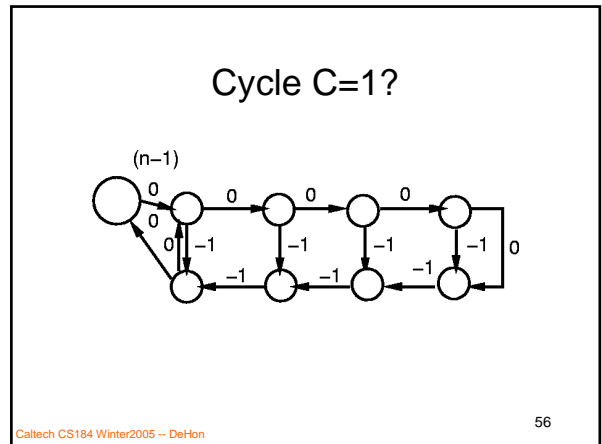
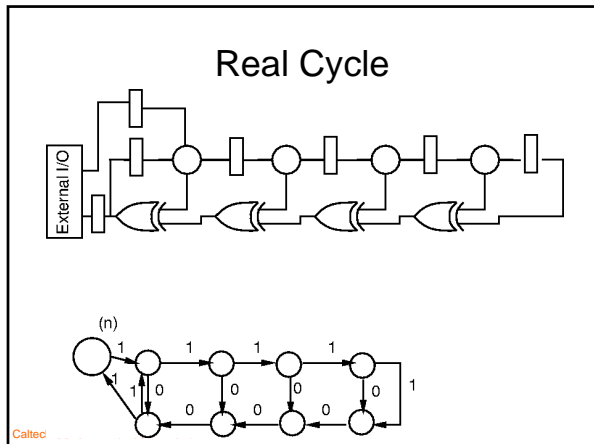
53

Real Cycle

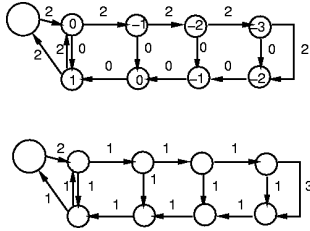


Caltech CS184 Winter2005 -- DeHon

54



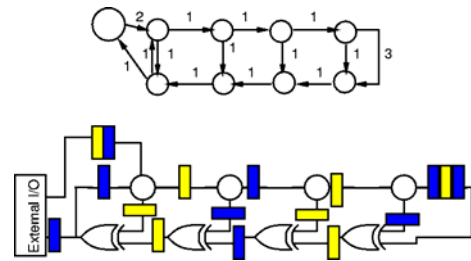
2-Slow Retime



Caltech CS184 Winter2005 -- DeHon

61

Retimed 2-Slow Cycle



Caltech CS184 Winter2005 -- DeHon

62

C-Slow applicable?

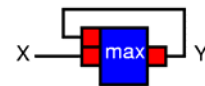
- Available parallelism
 - solve C identical, independent problems
 - Data-level parallelism
 - e.g. process packets (blocks) separately
 - e.g. independent regions in images
 - Commutative operators
 - e.g. max example

Caltech CS184 Winter2005 -- DeHon

63

Max Example

2-Slow design:



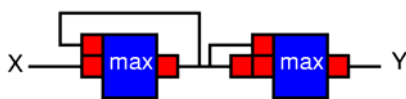
$X_2 X_2 X_1 X_1 X_0 X_0 \rightarrow Y_2 ? Y_1 ? Y_0 ?$

$B_2 A_2 B_1 A_1 B_0 A_0 \rightarrow Y A_2 Y B_1 Y A_1 Y B_0 Y A_0 ?$

Caltech CS184 Winter2005 -- DeHon

64

Max Example



Computes two interleaved streams: even max, odd max

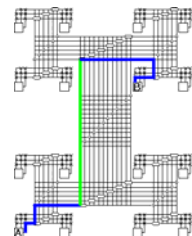
Computes final max of even and odd pairs

Caltech CS184 Winter2005 -- DeHon

65

HSRA Retiming

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



Caltech CS184 Winter2005 -- DeHon

66

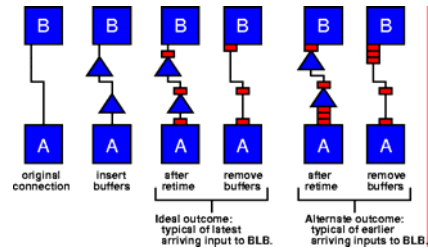
Accommodating HSRA Interconnect Delays

- 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 CS184 Winter2005 -- DeHon

67

Accommodating HSRA Interconnect Delays



Caltech CS184 Winter2005 -- DeHon

68

Big Ideas [MSB Ideas]

- Retiming transformations important to
 - minimize cycles
 - efficiently utilize spatial architectures
- Optimally solvable in $O(|V||E|)$ time
- Tells us
 - pipelining required
 - C-slow
 - where to move registers
- Can accommodate mandatory delays

Caltech CS184 Winter2005 -- DeHon

69