

CS184a: Computer Architecture (Structures and Organization)

Day15: November 13, 2000
Retiming

Previously

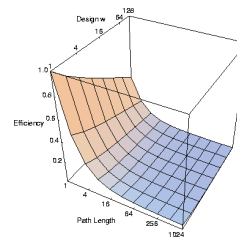
- Reviewed Pipelining
 - basic assignments on
- Saw spatial designs efficient
 - when reuse logic at maximum frequency
- Interconnect dominant delay
 - and dominant area
 - heavy call to reuse to use efficiently

Today

- Systematic transformation for retiming
 - maximize throughput
 - preserve semantics
 - “justify” mandatory registers in design

Motivation

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

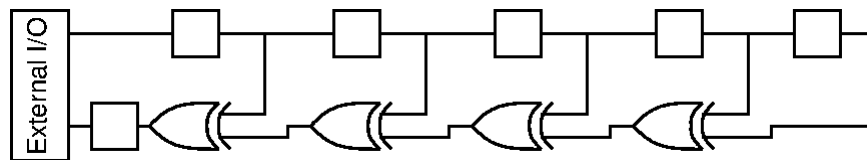


- “Everything in the right place at the right time.”

Task

- Move registers to:
 - preserve semantics
 - Minimize path length between registers
 - (make path length 1 for maximum throughput or reuse)
 - Maximize reuse rate
 - ...while minimizing number of registers required

Simple Example

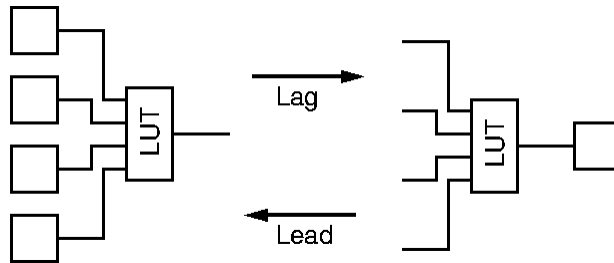


Path Length (L) = 4

Can we do better?

Legal Register Moves

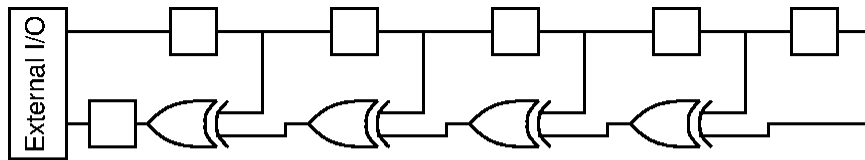
- Retiming Lag/Lead



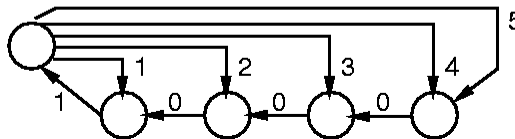
Caltech CS184a Fall2000 -- DeHon

7

Canonical Graph Representation



Observable I/O



Separate arch for each path

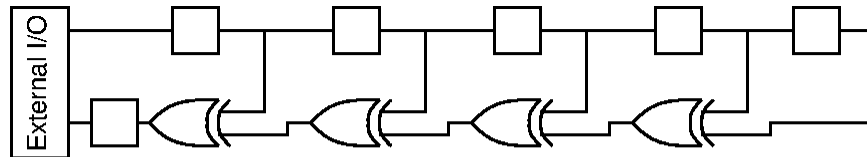
Weight edges by number of registers

(weight nodes by delay through node)

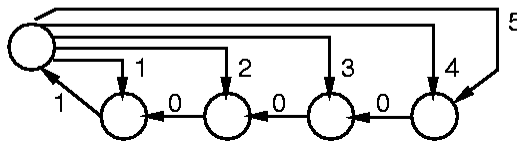
Caltech CS184a Fall2000 -- DeHon

8

Critical Path Length



Observable I/O

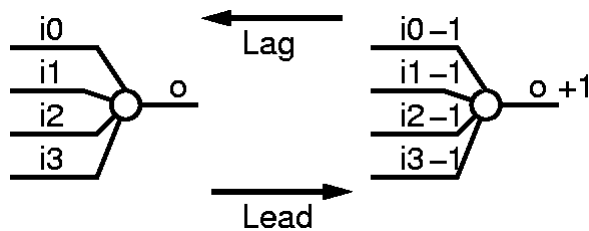


Critical Path: Length of longest path of zero weight nodes

Compute in $O(|E|)$ time by levelizing network:

Caltech CS184 Fall 2000 -- DeHon Topological sort, push path lengths forward until find register. ⁹

Retiming Lag/Lead



Retiming: Assign a lag to every vertex

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

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 CS184a Fall2000 -- DeHon

11

Retiming Task

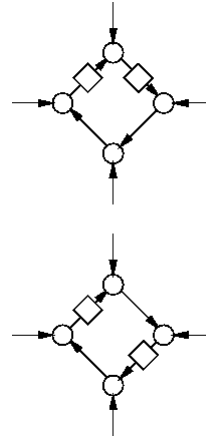
- Move registers \equiv assign lags to nodes
 - lags define all locally legal moves
- Preserving non-negative edge weights
 - (previous slide)
 - guarantees collection of lags remains consistent globally

Caltech CS184a Fall2000 -- DeHon

12

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

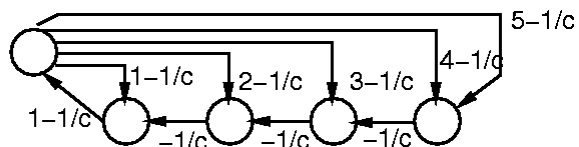


Optimal Retiming

- 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

$G-1/c$

Observable I/O



Compute Retiming

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

Bellman Ford

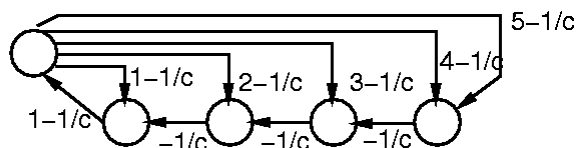
- For $k \leftarrow 0$ to N
 - $u_i \leftarrow \infty$ (except $u_i = 0$ for IO)
- For $k \leftarrow 0$ to N
 - for $e_{i,j} \in E$
 - $u_i \leftarrow \min(u_i, u_j + w(e_{i,j}))$
- for $e_{i,j} \in E$
 - if $u_i > u_j + w(e_{i,j})$
 - cycles detected

Caltech CS184a Fall2000 -- DeHon

17

Apply to Example

Observable I/O

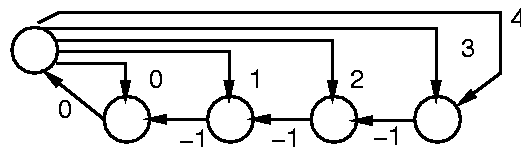


Caltech CS184a Fall2000 -- DeHon

18

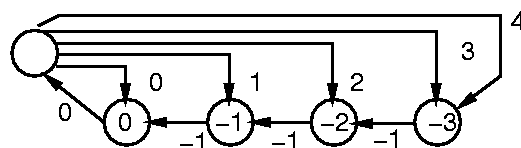
Apply: Find Lags

Observable I/O



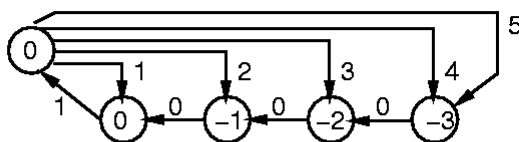
Apply: Lags

Observable I/O

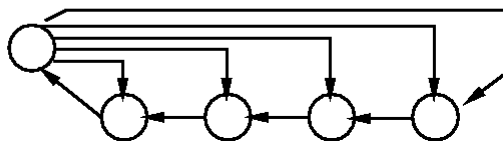


Apply: Move Registers

Observable I/O



Observable I/O



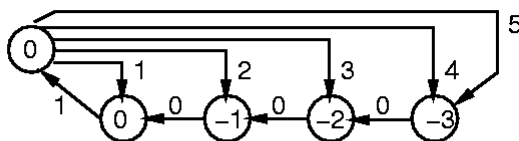
$$\text{weight}(e') = \text{weight}(e) + \text{lag}(\text{head}(e)) - \text{lag}(\text{tail}(e))$$

Caltech CS184a Fall2000 -- DeHon

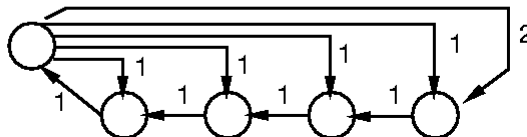
21

Apply: Retimed

Observable I/O



Observable I/O

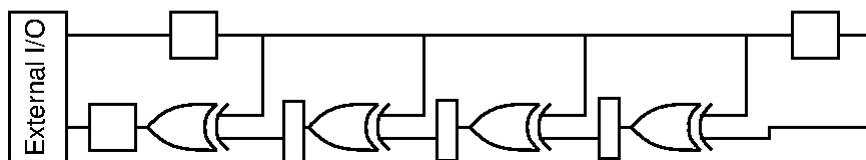
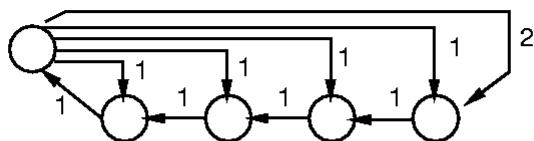


Caltech CS184a Fall2000 -- DeHon

22

Apply: Retimed Design

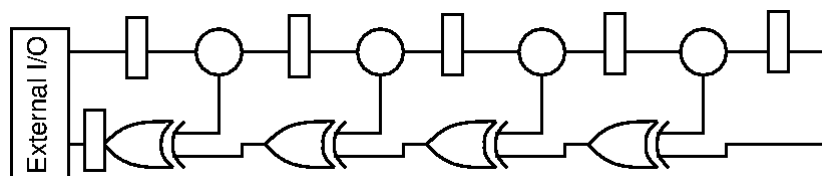
Observable I/O



Caltech CS184a Fall2000 -- DeHon

23

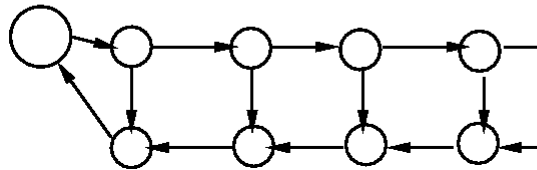
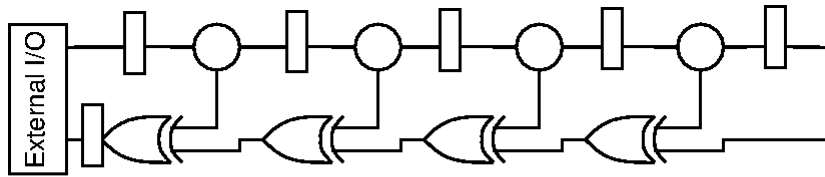
Revise Example (fanout delay)



Caltech CS184a Fall2000 -- DeHon

24

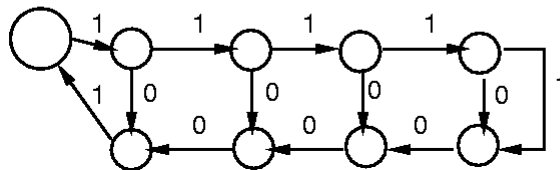
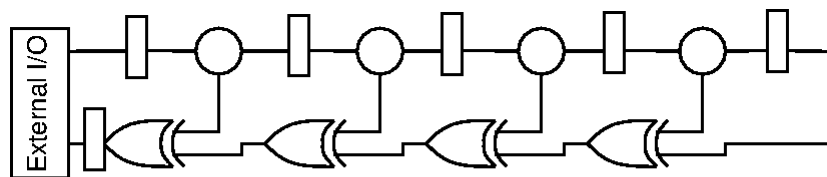
Revised: Graph



Caltech CS184a Fall2000 -- DeHon

25

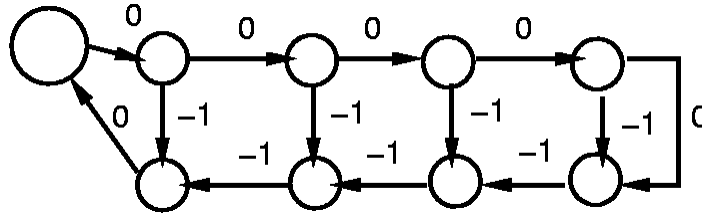
Revised: Graph



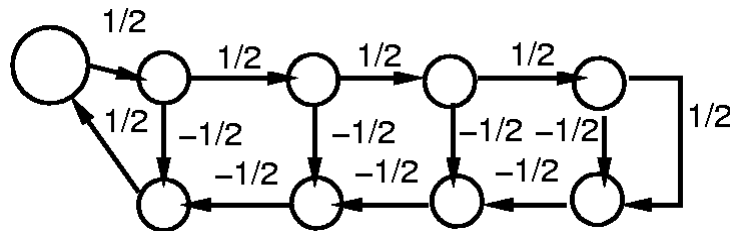
Caltech CS184a Fall2000 -- DeHon

26

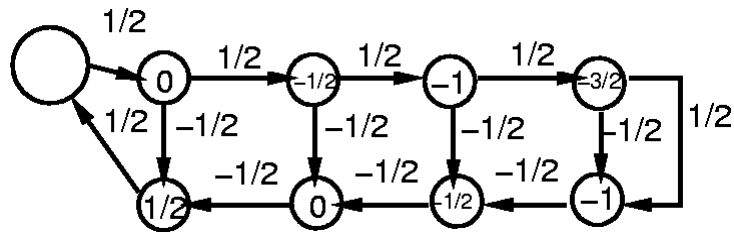
Revised: $C=1$?



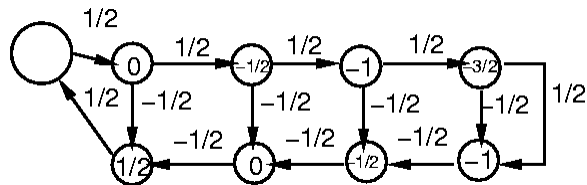
Revised: $C=2$?



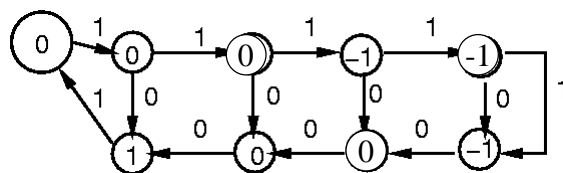
Revised: Lag



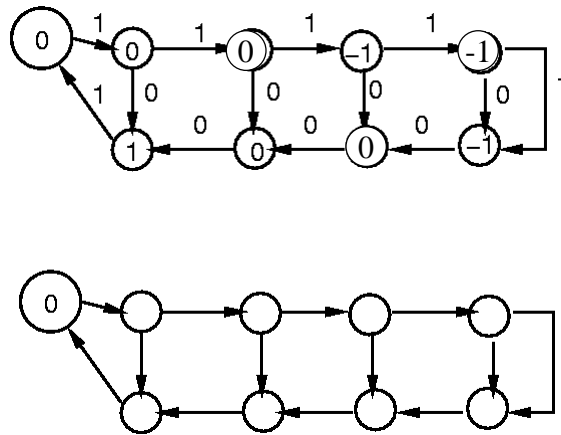
Revised: Lag



Take ceiling to convert to integer lags:



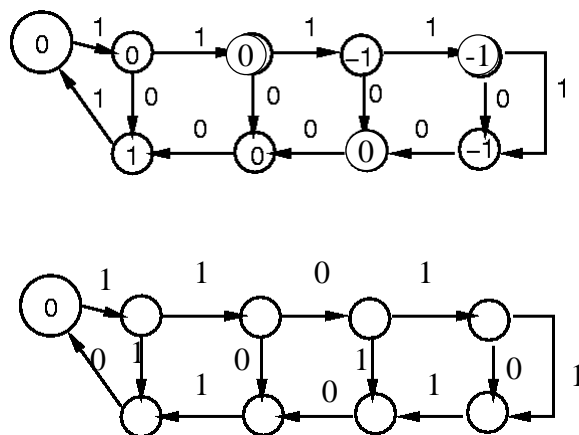
Revised: Apply Lag



Caltech CS184a Fall2000 -- DeHon

31

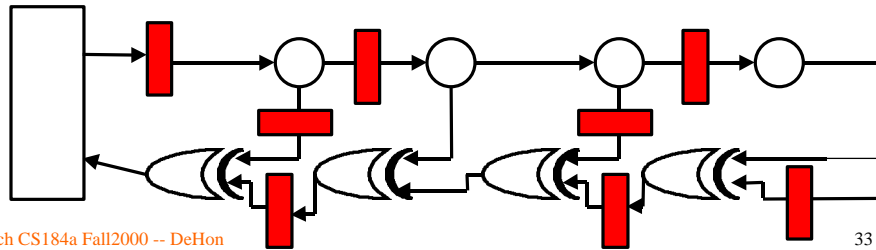
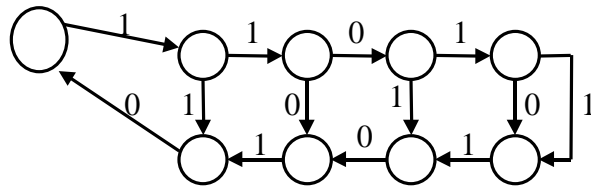
Revised: Apply Lag



Caltech CS184a Fall2000 -- DeHon

32

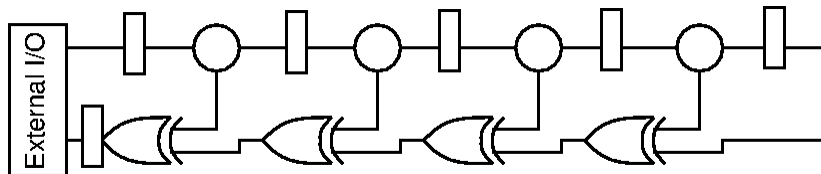
Revised: Retimed



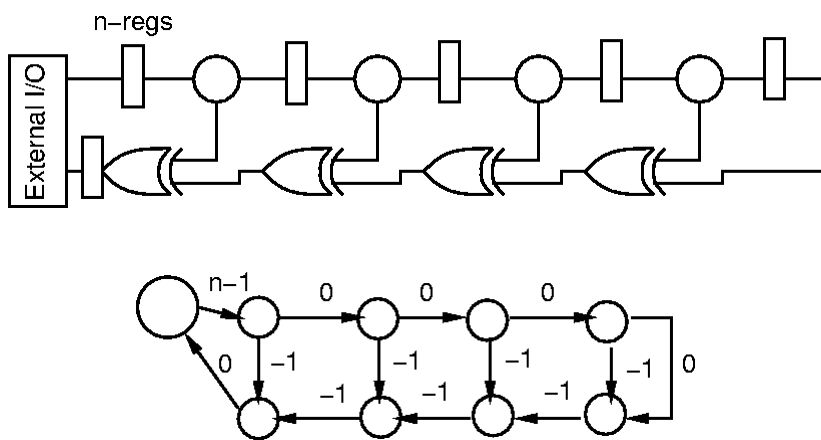
Pipelining

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

$C > 1 \implies$ Pipeline

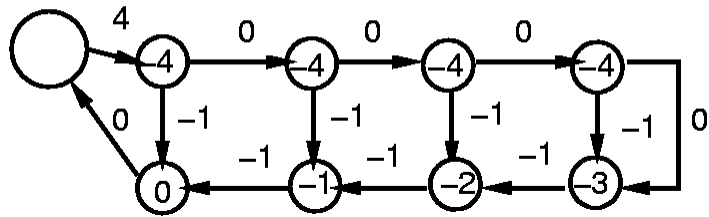


Add Registers



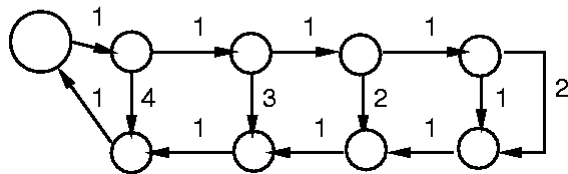
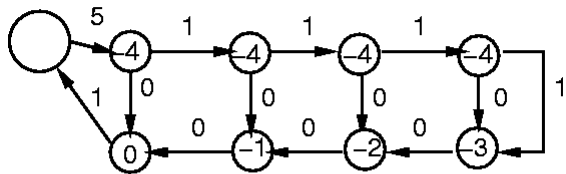
Pipeline Retiming: Lag

n=5

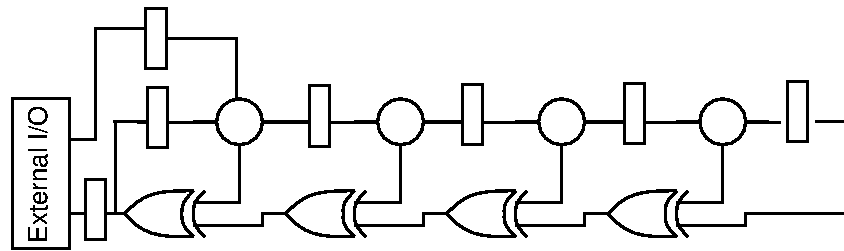


Pipelined Retimed

n=5



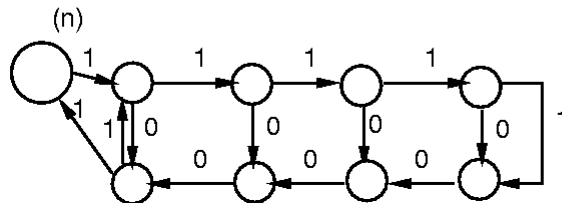
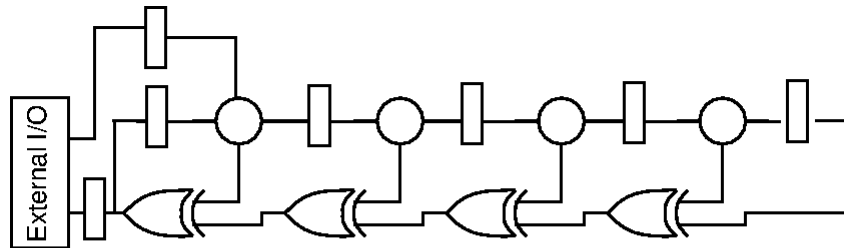
Real Cycle



Caltech CS184a Fall2000 -- DeHon

39

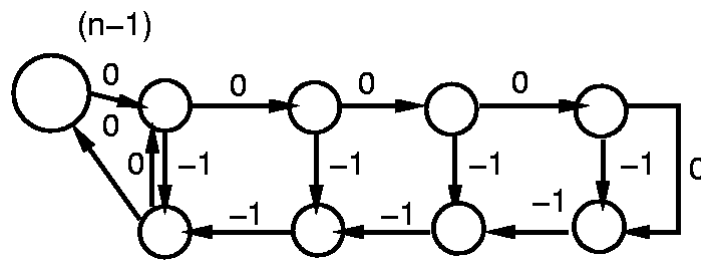
Real Cycle



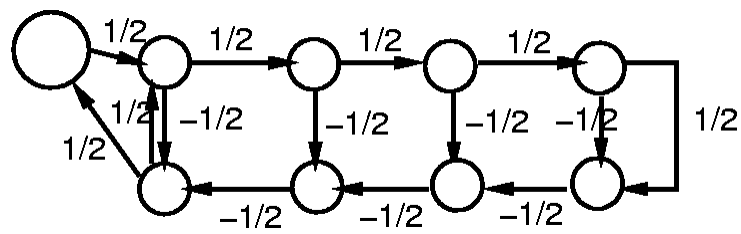
Caltech

40

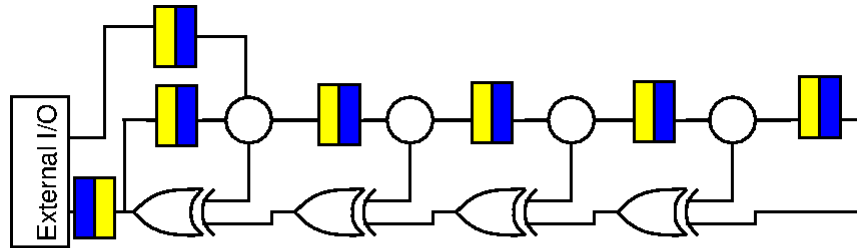
Cycle $C=1$?



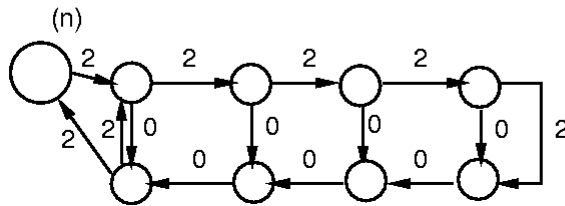
Cycle $C=2$?



Cycle: C-slow



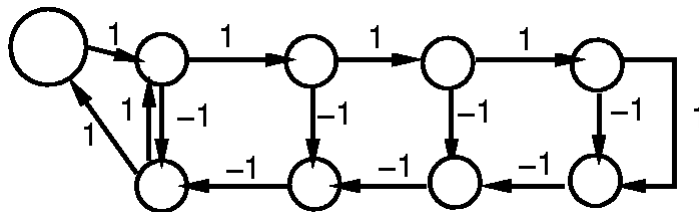
Cycle=c \Rightarrow C-slow network has Cycle=1



Caltec

43

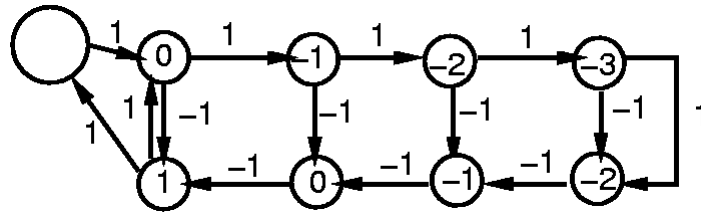
2-slow Cycle \Rightarrow C=1



Caltech CS184a Fall2000 -- DeHon

44

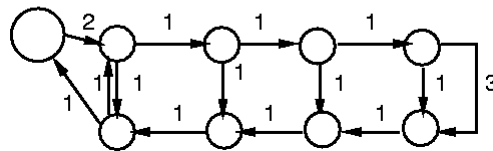
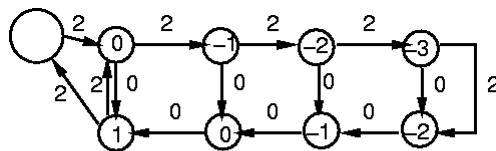
2-Slow Lags



Caltech CS184a Fall2000 -- DeHon

45

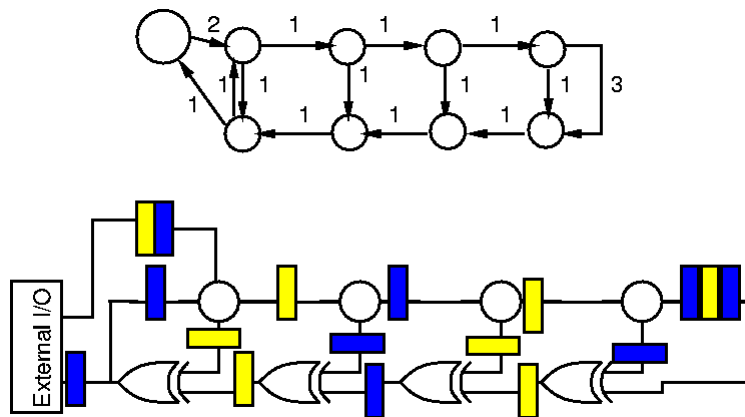
2-Slow Retime



Caltech CS184a Fall2000 -- DeHon

46

Retimed 2-Slow Cycle



Caltech CS184a Fall2000 -- DeHon

47

C-Slow applicable?

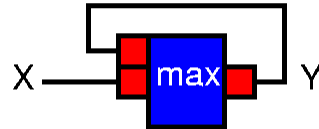
- Available parallelism
 - solve C identical, independent problems
 - e.g. process packets (blocks) separately
 - e.g. independent regions in images
- Commutative operators
 - e.g. max example

Caltech CS184a Fall2000 -- DeHon

48

Max Example

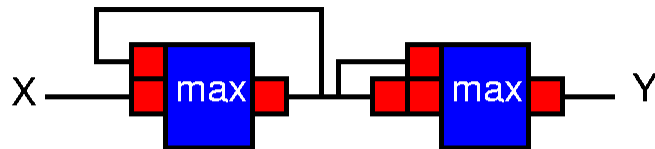
2-Slow design:



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

$B_2 A_2 B_1 A_1 B_0 A_0 \rightarrow YA_2 YB_1 YA_1 YB_0 YA_0 ?$

Max Example



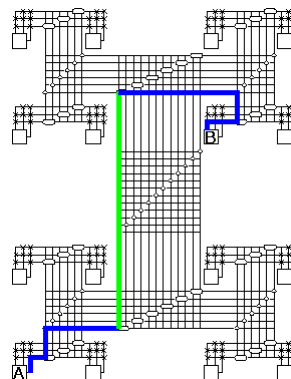
Computes two
interleaved streams:
even max, odd max

Computes final
max of even and
odd pairs

Monday Lecture Stopped Here

HSRA Retiming

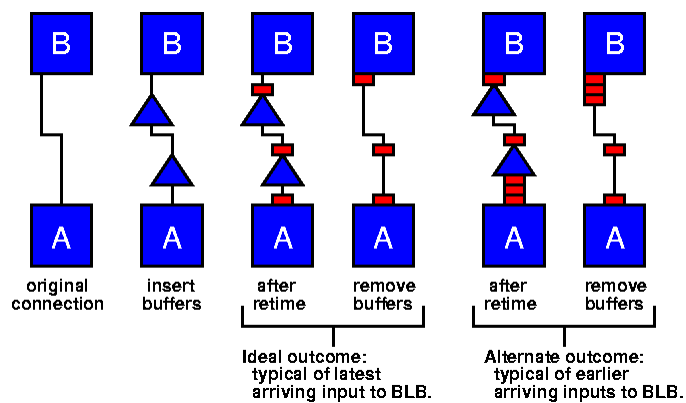
- HSRA
 - adds mandatory pipelining to interconnect
- One additional twist
 - long, pipelined interconnect
 - \Rightarrow need more than one register on paths



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

Accommodating HSRA Interconnect Delays



Big Ideas

[MSB Ideas]

- Retiming 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