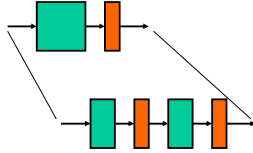


Pipelining

- Adding registers along a path
 - split combinational logic into multiple cycles
 - each cycle smaller than previously
 - increase throughput



Pipelining and Retiming 1

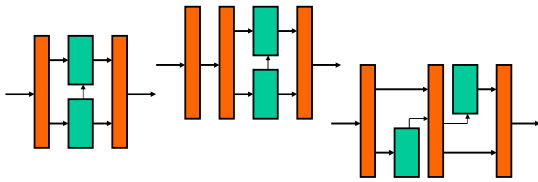
Pipelining

- Delay, d , of slowest combinational stage determines performance
- Throughput = $1/d$: rate at which outputs are produced
- Latency = $n \cdot d$: number of stages * clock period
- Pipelining increases circuit utilization
- Registers slow down data, synchronize data paths
- Wave-pipelining
 - no pipeline registers - waves of data flow through circuit
 - relies on equal-delay circuit paths - no short paths

Pipelining and Retiming 2

When and How to Pipeline?

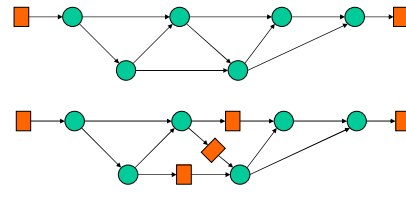
- Where is the best place to add registers?
 - splitting combinational logic
 - overhead of registers (propagation delay and setup time requirements)
- What about cycles in data path?
- Example: 16-bit adder, add 8-bits in each of two cycles



Pipelining and Retiming 3

Retiming

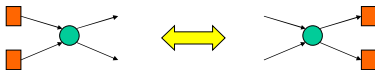
- Process of optimally distributing registers throughout a circuit
 - minimize the clock period
 - minimize the number of registers



Pipelining and Retiming 4

Retiming (cont'd)

- Fast optimal algorithm (Leiserson & Saxe 1983)
- Retiming rules:
 - remove one register from each input and add one to each output
 - remove one register from each output and add one to each input



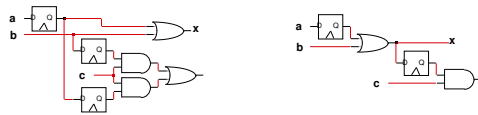
Pipelining and Retiming 5

Retiming examples

- Shortening critical paths



- Create simplification opportunities



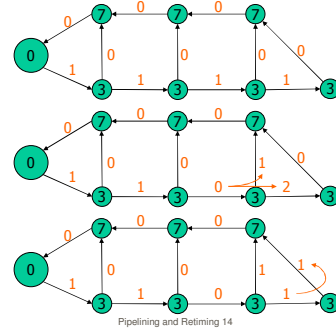
Pipelining and Retiming 6

Computing W and D

- $W[u,v]$ = number of registers on the minimum weight path from $u \rightarrow v$
 - Any retiming changes the weight of all paths by the same constant
 - i.e. Retiming cannot change which is the minimum weight path
- $D[u,v]$ = maximum delay over all paths with $W[u,v]$ registers
 - Retiming does not affect $D[u,v]$
- These matrices contain all the required register and delay information
 - If retiming removes all registers from the path $u \rightarrow v$, then $D[u,v]$ is the largest delay path that results

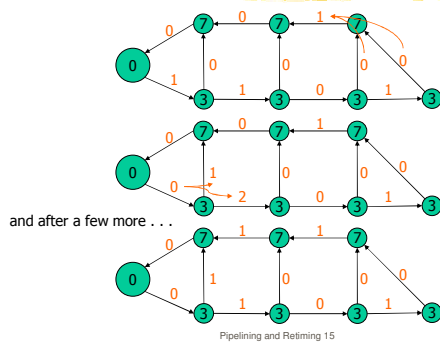
Pipelining and Retiming 13

Retiming: One Step at a Time



Pipelining and Retiming 14

Retiming: One Step at a Time (cont'd)



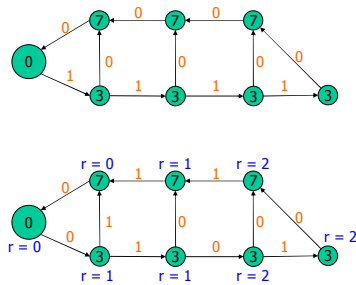
Pipelining and Retiming 15

Retiming: Problem Formulation

- $r(v)$: number of registers pushed through a node in the forward direction
 - $w_{\text{new}}(u, v) = w_{\text{old}}(u, v) + r(u) - r(v)$
- Problem statement
 - $r(v_h) = 0$ (*host is not retimed*)
 - $w_{\text{new}}(u, v) = w_{\text{old}}(u, v) + r(u) - r(v) \geq 0$, for all u, v
 - $r(u) - r(v) \geq -w_{\text{old}}(u, v)$ (*no negative registers!*)
 - For all $D[u,v] > T_{\text{clk}}$,
 - $w_{\text{new}}(u, v) = w_{\text{old}}(u, v) + r(u) - r(v) \geq 1$
 - $r(u) - r(v) \geq -w_{\text{old}}(u, v) + 1$ (*every long path has at least 1 reg*)
- Difference constraints like this can be solved by generating a graph that represents the constraints and using a shortest path algorithm like Bellman-Ford to find a set of $r(v)$ values that meets all the constraints
- The value of $r(v)$ returned by the algorithm can be used to generate the new positions of the registers in the retimed circuit

Pipelining and Retiming 16

Retimed Correlator



Pipelining and Retiming 17

Extensions to Retiming

- Host interface
 - add latency
 - multiple hosts
- Area considerations
 - limit number of registers
 - optimize logic across register boundaries
 - peripheral retiming
 - incremental retiming
 - pre-computation
- Generality
 - different propagation delays for different signals
 - widths of interconnections

Pipelining and Retiming 18

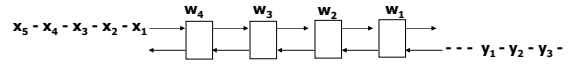
Systolic Arrays

- Set of identical processing elements
 - specialized or programmable
- Efficient nearest-neighbor interconnections (in 1-D, 2-D, other)
- SIMD-like
- Multiple data flows, converging to engage in computation

Analogy: data flowing through the system in a rhythmic fashion – from main memory through a series of processing elements and back to main memory

Example - Convolution

$$y_j = x_j w_1 + x_{j+1} w_2 + \dots + x_{j+n-1} w_n$$



$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

$$y_2 = x_2 w_1 + x_3 w_2 + x_4 w_3 + x_5 w_4$$

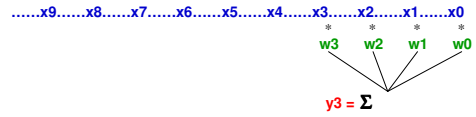
$$y_3 = x_3 w_1 + x_4 w_2 + x_5 w_3 + x_6 w_4$$

Example - Convolution

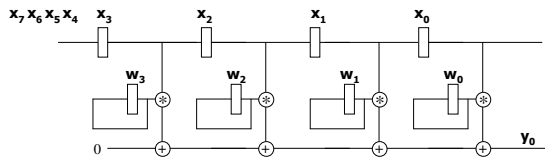
	w_4	w_3	w_2	w_1	
$x_6 - x_5 - x_4 - x_3 - x_2 -$	x_1				$y_1 - y_2 - y_3$
$x_6 - x_5 - x_4 - x_3 - x_2 -$	x_1				$y_1 - y_2 - y_3$
$x_6 - x_5 - x_4 - x_3 -$	x_2	x_1			$y_1 - y_2 - y_3$
$x_6 - x_5 - x_4 - x_3 -$	x_2	x_1			$y_1 - y_2 - y_3$
$x_6 - x_5 - x_4 -$	x_3	x_2	x_1		$y_2 - y_3$
$x_6 - x_5 - x_4 -$	x_3	x_2	x_1		$y_2 - y_3$
$x_6 - x_5 -$	x_4	x_3	x_2	x_1	y_3
$x_6 - x_5 -$	x_4	x_3	x_2	x_1	y_3

Convolution - Another Look

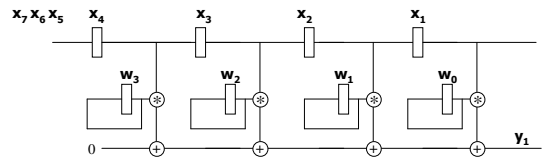
- Repeated vector product



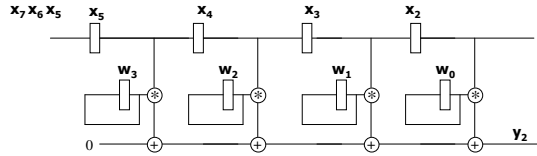
Convolution Example



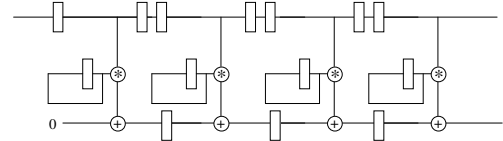
Convolution Example



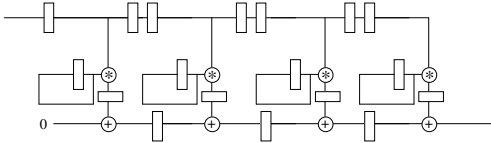
Convolution Example



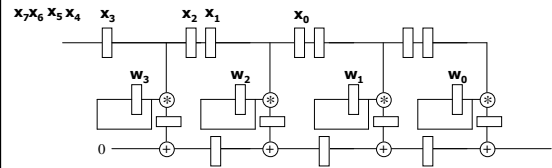
Pipelining and Retiming 25



Pipelining and Retiming 26



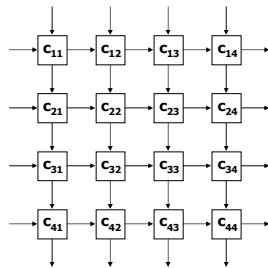
Pipelining and Retiming 27



Pipelining and Retiming 28

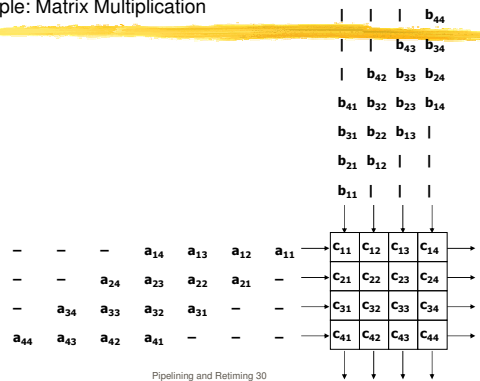
Example: Matrix Multiplication

$C = A \times B$ $c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$



Pipelining and Retiming 29

Example: Matrix Multiplication



Pipelining and Retiming 30

Systolic Algorithms

- 2D Convolution
 - Image processing
- FFT
- String matching
 - Dynamic programming
 - DNA comparison
- Matrix computations
 - LU decomposition
 - QR factorization

Pipelining and Retiming 31

Systolic Architectures

- Highly parallel
 - "fine-grained" parallelism
 - deep pipelining
- Local communication
 - wires are short - no global communication (except CLK)
 - linear array \rightarrow no clock skew
 - increasingly important as wire delays increase (relative to gate delays)
- Linear arrays
 - most systolic algorithms can be done with a linear array
 - include memory in each cell in the array
 - linear array a better match to I/O limitations
- Contrast to superscalar and vector architectures

Pipelining and Retiming 32

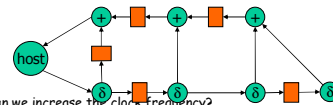
Systolic Computers

- Custom chips - early 1980's
- Warp (CMU) - 1987
 - linear array of 10 or more processing cells
 - optimized inter-cell communication for low-latency
 - pipelined cells and communication
 - conditional execution
 - compiler partitions problem into cells and generates microcode
- i-Warp (Intel) - 1990
 - successor to Warp
 - two-dimensional array
 - time-multiplexing of physical busses between cells
 - 32x32 array has 20Gflops peak performance
 - not a commercial success
- Currently confined to ASIC implementations

Pipelining and Retiming 33

Digital Correlator Revisited

- Optimally retimed circuit (clock cycle 13)

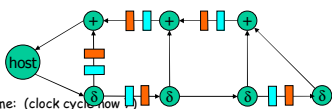


- How can we increase the clock frequency?
 - Work on multiple data sets at the same time

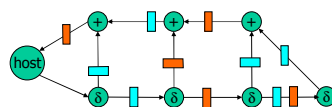
Pipelining and Retiming 34

C-slowng a Circuit

- Replace every register with C registers



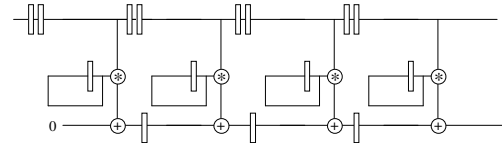
- Now retime: (clock cycle now 13)



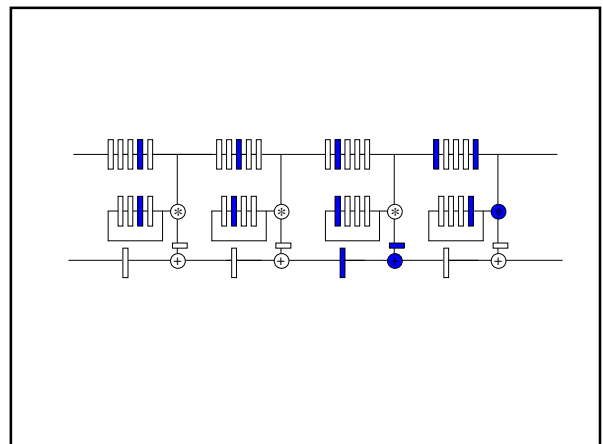
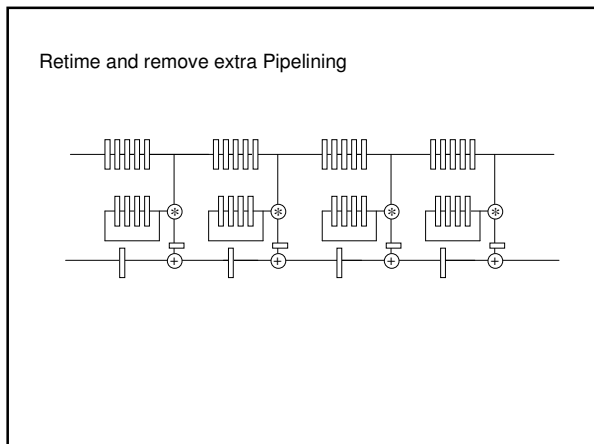
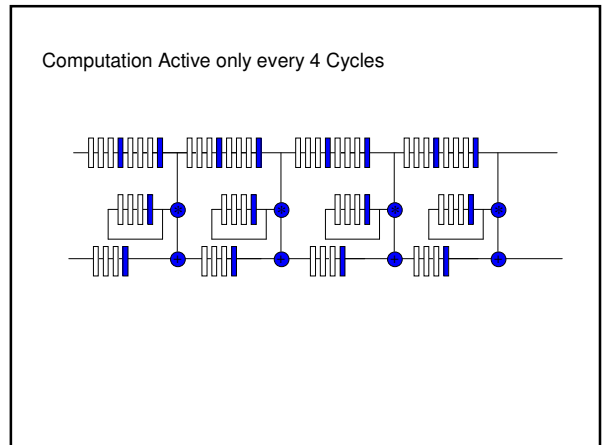
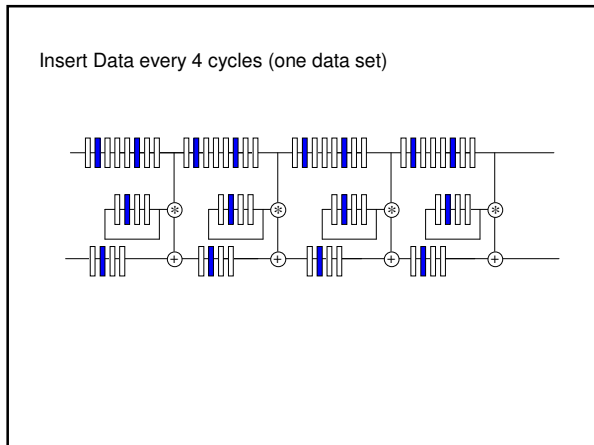
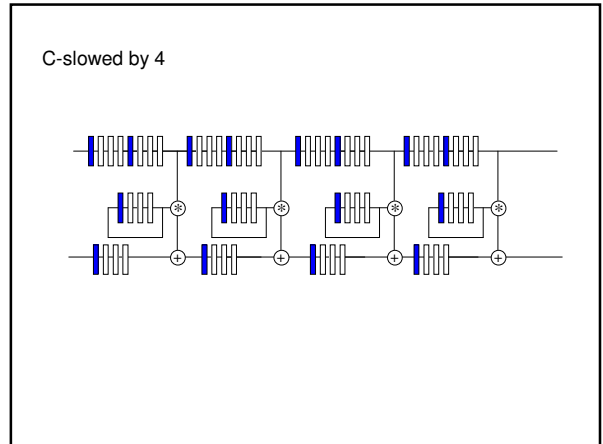
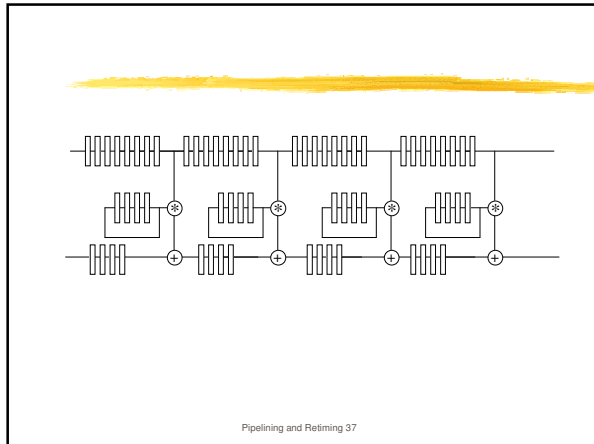
Pipelining and Retiming 35

C-slowng/Retiming for Resource Sharing

- Correlator circuit



Pipelining and Retiming 36



Computation spread over time

- Only need one multiplier and one adder
- We can use this method to schedule for any number of resources

