

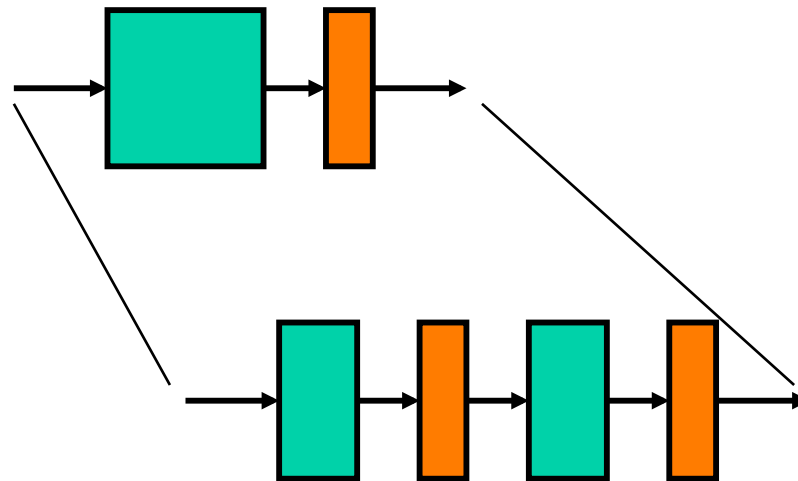
REVIEW



- Placement - Greedy vs. Simulated Annealing
- Routing - Shortest Path (Dijkstra), A*, Pathfinder
- Transmission lines - high frequency loss, dispersion, ringing, overshoot
- Differential signaling - advantages, eye diagrams
- Crosstalk - inductive, capacitive

PIPELINING

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



PIPELINING

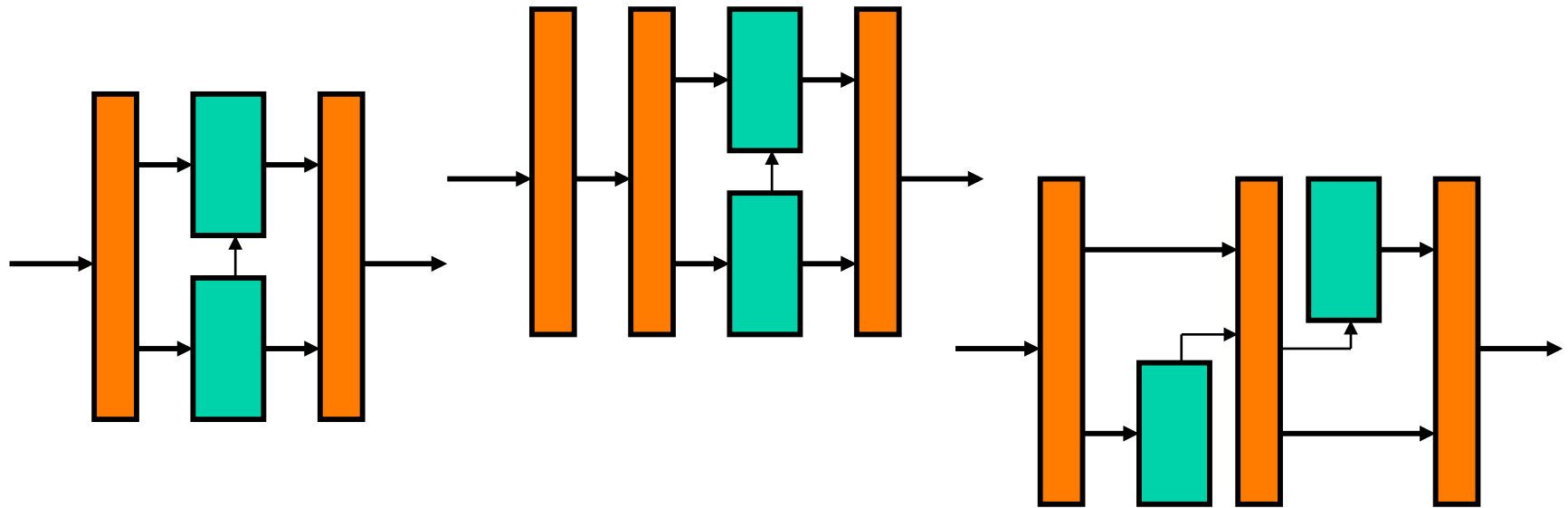


- Delay, d , of slowest combinational stage determines performance
- Throughput = $1/d$: rate at which outputs are produced
- Latency = $n \cdot d$: number of stages * clock period
- Pipelining increases circuit utilization
- Registers slow down data, synchronize data paths

- Wave-pipelining
 - no pipeline registers - waves of data flow through circuit
 - relies on equal-delay circuit paths - no short paths

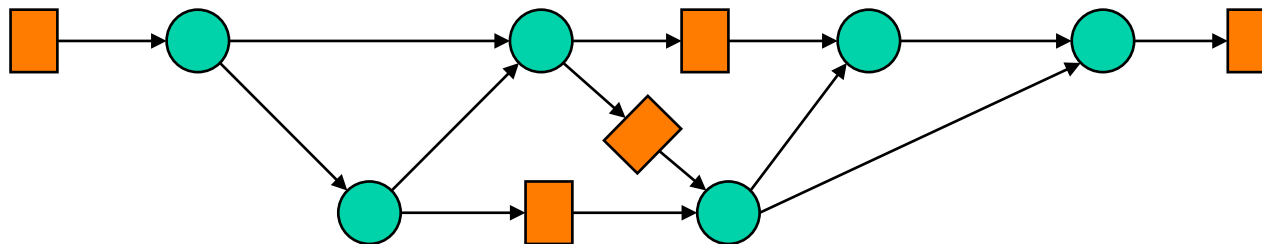
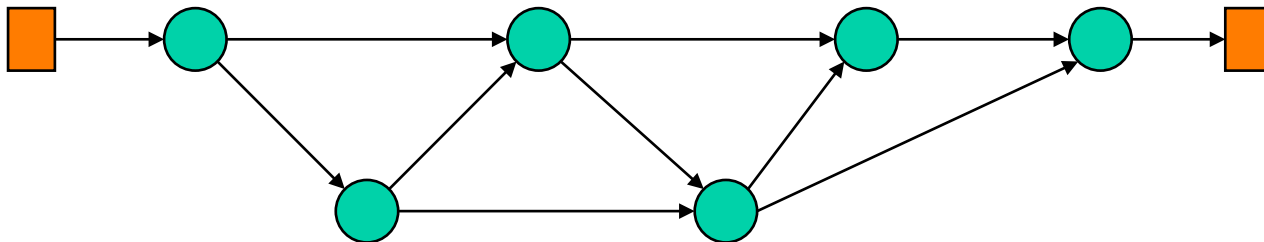
WHEN AND HOW TO PIPELINE?

- 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



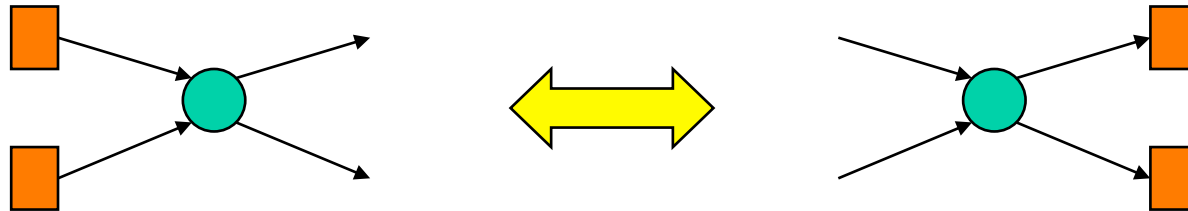
RETIMING

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



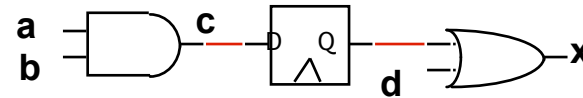
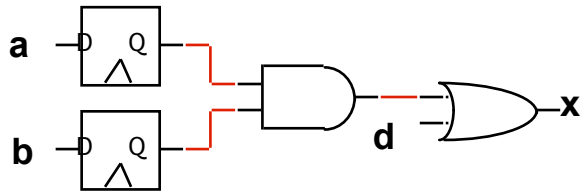
RETIMING (CONT'D)

- 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

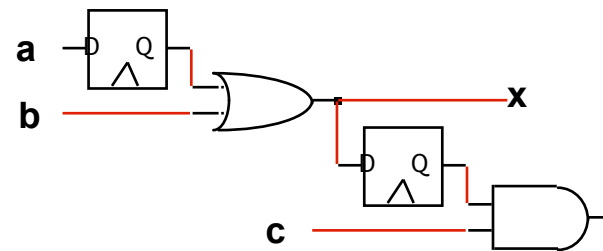
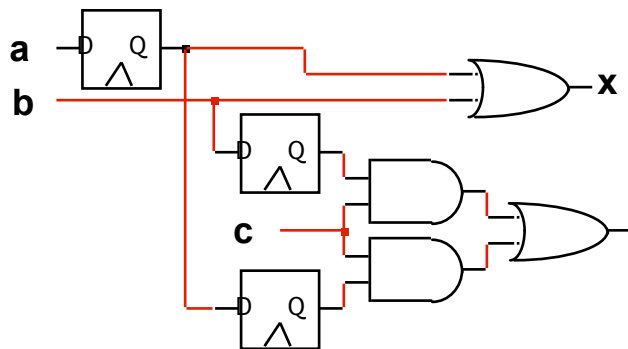


RETIMING EXAMPLES

- Shortening critical paths

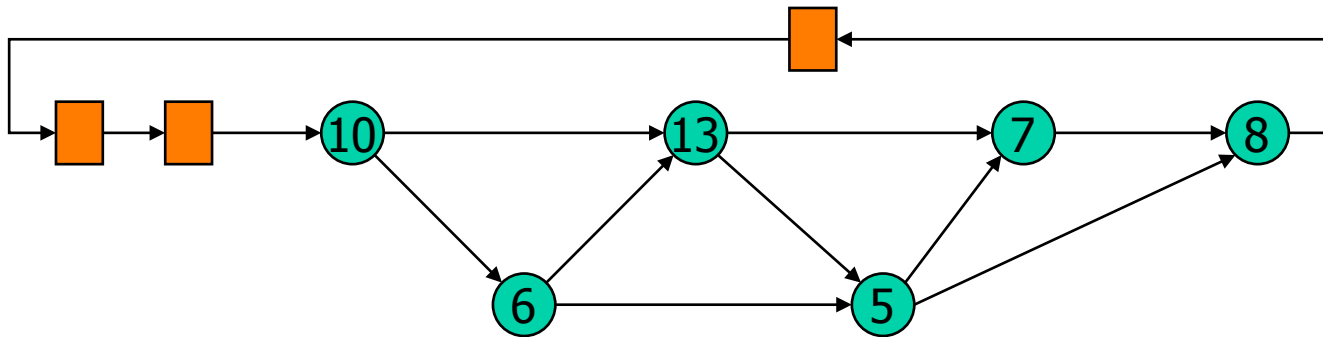


- Create simplification opportunities



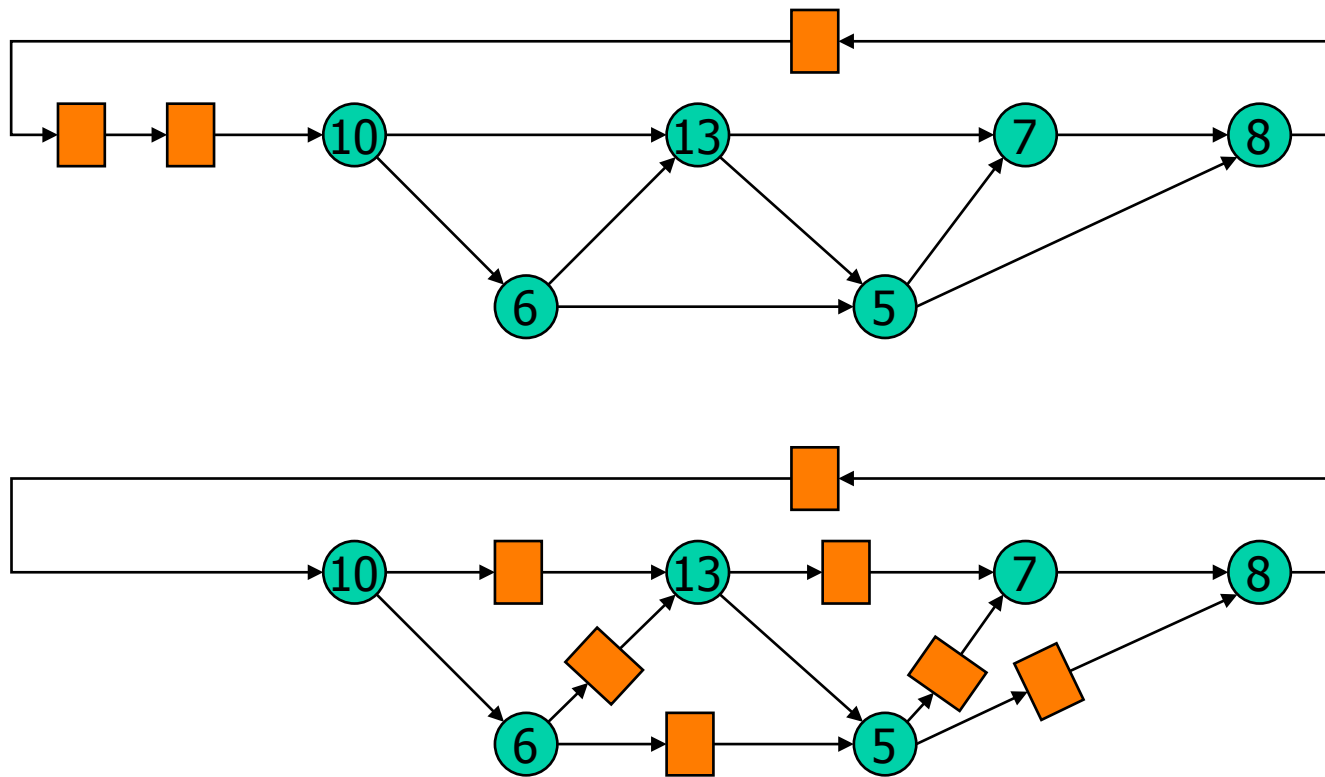
OPTIMAL PIPELINING

- Add registers - use retiming to find optimal location



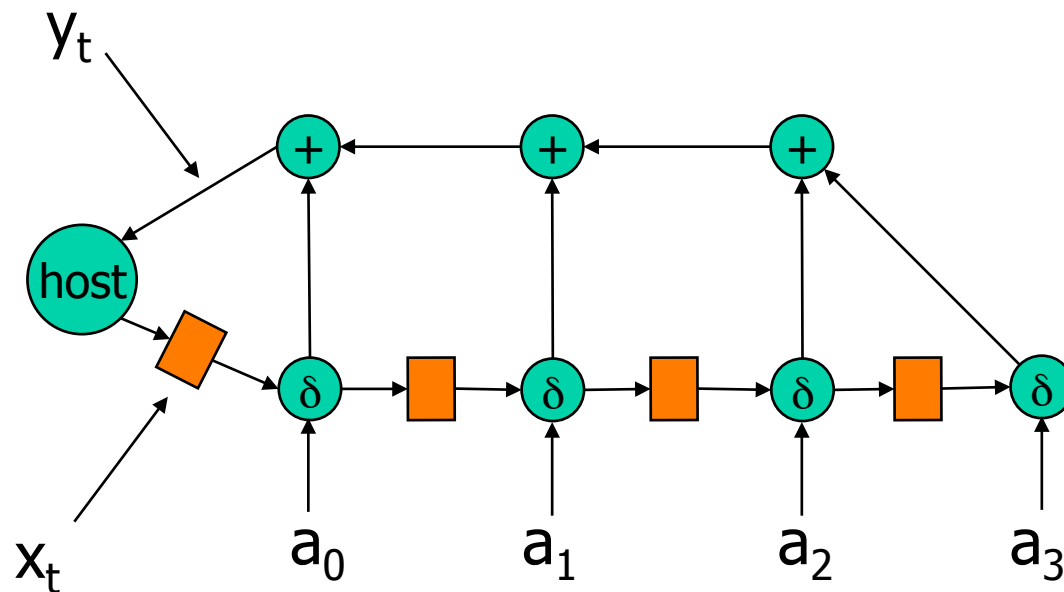
OPTIMAL PIPELINING

- Add registers - use retiming to find optimal location



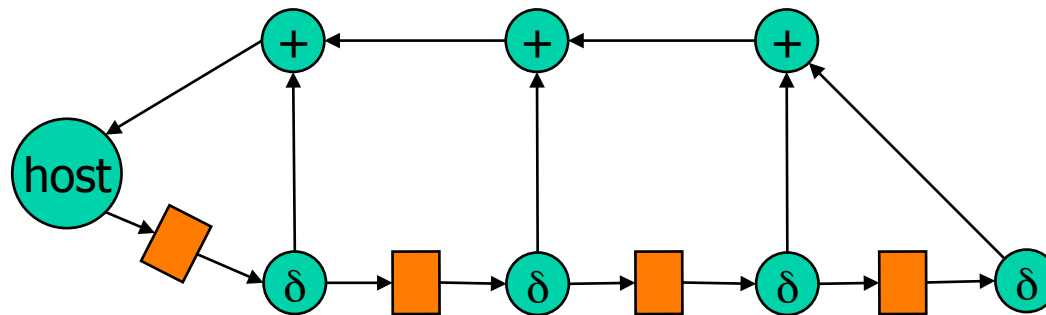
EXAMPLE - DIGITAL CORRELATOR

- $y_t = \delta(x_t, a_0) + \delta(x_{t-1}, a_1) + \delta(x_{t-2}, a_2) + \delta(x_{t-3}, a_3)$
- $\delta(x_t, a_0) = 0$ if $x \neq a$, 1 otherwise (and passes x along to the right)



EXAMPLE - DIGITAL CORRELATOR (CONT'D)

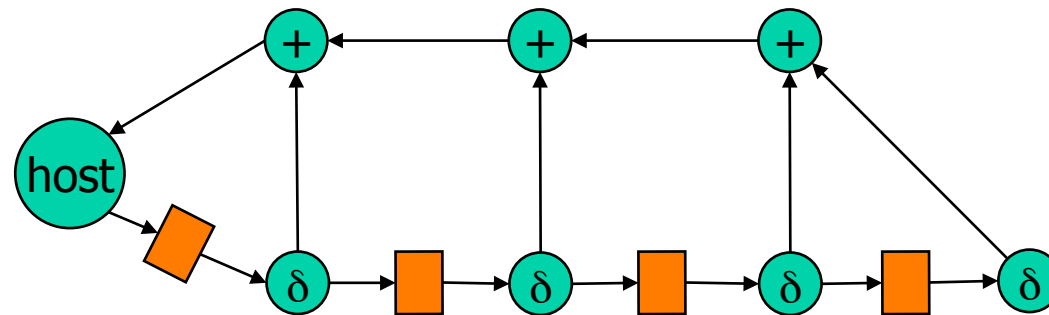
- Delays: adder, 7; comparator, 3; host, 0



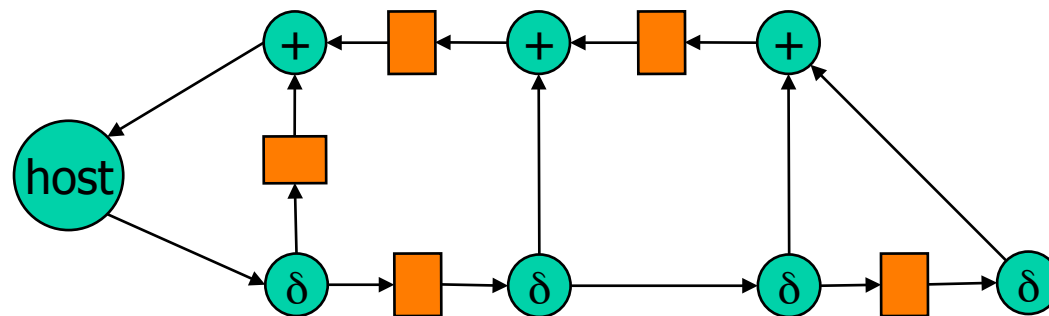
cycle time =

EXAMPLE - DIGITAL CORRELATOR (CONT'D)

- Delays: adder, 7; comparator, 3; host, 0



cycle time = 24



cycle time = 13

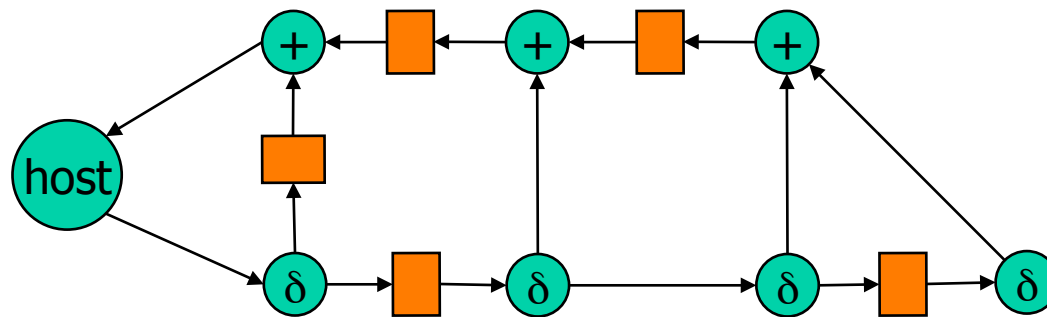
EXTENSIONS TO RETIMING



- 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

DIGITAL CORRELATOR REVISITED

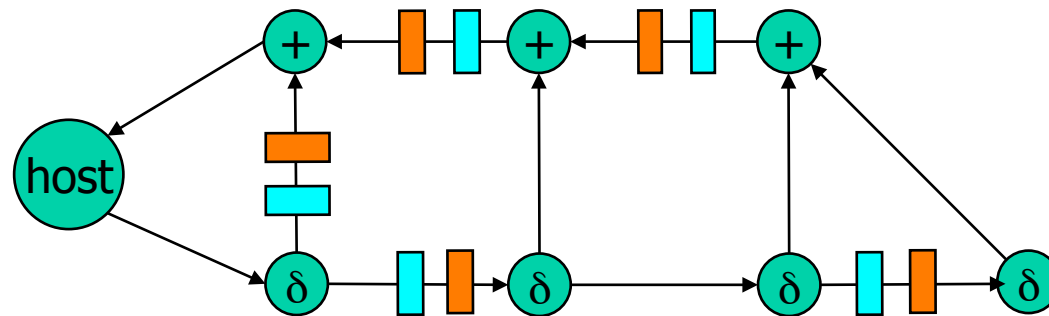
- Optimally retimed circuit (clock cycle 13)



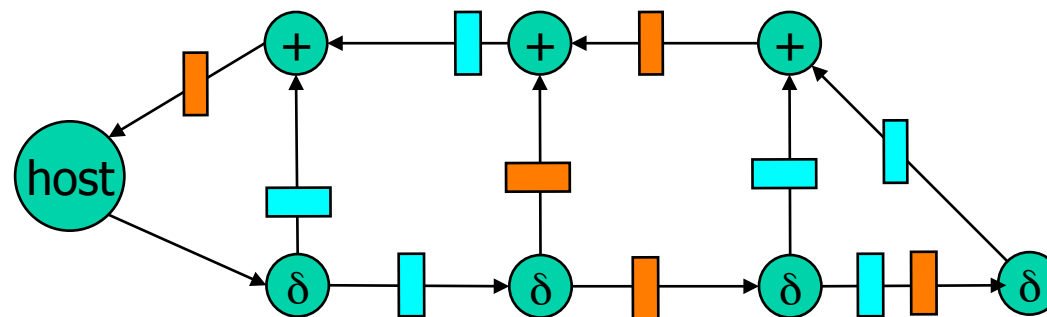
- How can we increase the clock frequency?
 - Work on multiple data sets at the same time

C-SLOW'ING A CIRCUIT

- Replace every register with C registers

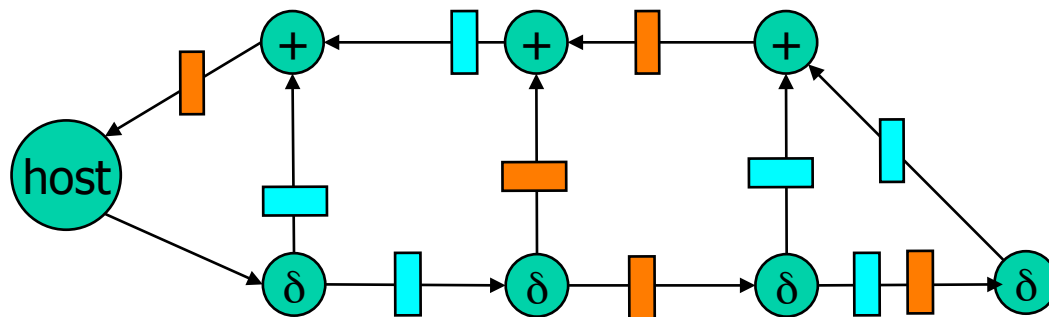


- Now retime: (clock cycle now 7)



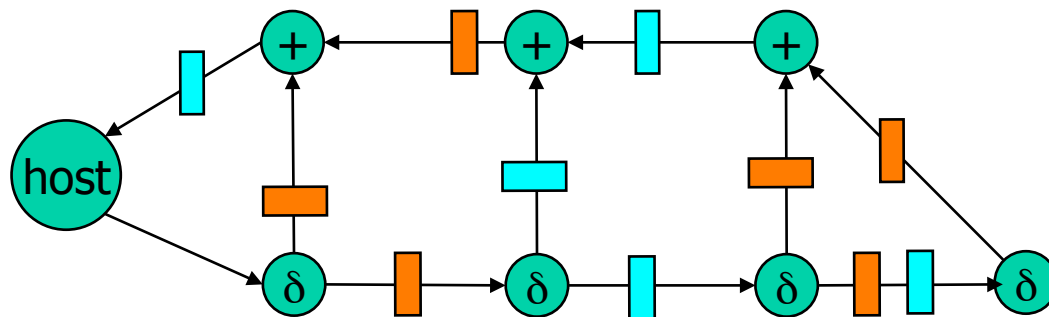
C-SLOW'ING A CIRCUIT = MULTI-THREADING

- In this case there are two threads (blue and orange)
 - Host alternates between the two threads
- Input blue data, remove orange results



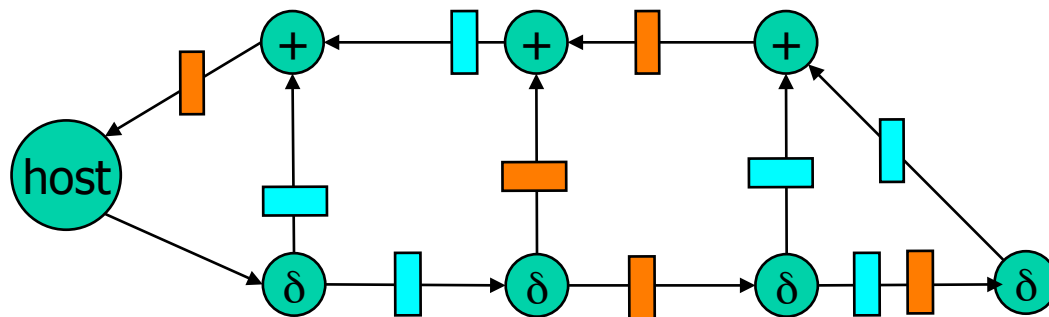
C-SLOW'ING A CIRCUIT = MULTI-THREADING

- In this case there are two threads (blue and orange)
 - Host alternates between the two threads
- Input orange data, remove blue results



C-SLOW'ING A CIRCUIT = MULTI-THREADING

- In this case there are two threads (blue and orange)
 - Host alternates between the two threads
- Input blue data, remove orange results
 - Throughput of each thread (1/14) is almost what it was before (1/13)!



PIPELINING PROCESSORS



- Pipelined processors are complex because of feedback loops (cycles)
 - Forwarding
 - Branch prediction
 - Long latency ops (e.g. cache miss) cause stalls
- Solution - *C*-slowing!
 - Start with non-pipelined processor
 - Program counter, register file
 - *C*-slow by *N*
 - *N* program counters, register files
 - Pipeline datapath, ignoring hazards!

MULTI-THREADED PROCESSOR



- Pipelined, c-slotted processor is "multi-threaded"
 - Executing N different instruction streams simultaneously
 - Each executes an instruction every N cycles
 - Allows op feedback latency of N cycles
 - e. g. cache read miss
- Tera Computer MTA (now Cray) is multi-threaded
 - $N = 1024$
 - There is no cache!
 - Remote memory access < 1024 cycles (2 usec)
- Requires huge parallelism (N threads!)

MULTI-THREADED PROCESSORS



- Multi-threaded processors are very simple
 - No stalls, no forwarding, etc.
- Great for FPGAs
 - Muxes are expensive and slow!
 - Registers are almost free
- Typical $N = 4$ for multi-threaded FPGA processors
 - Almost 4x performance increase
 - Higher utilization (no stalls) (2x)
 - Higher frequency (2x)

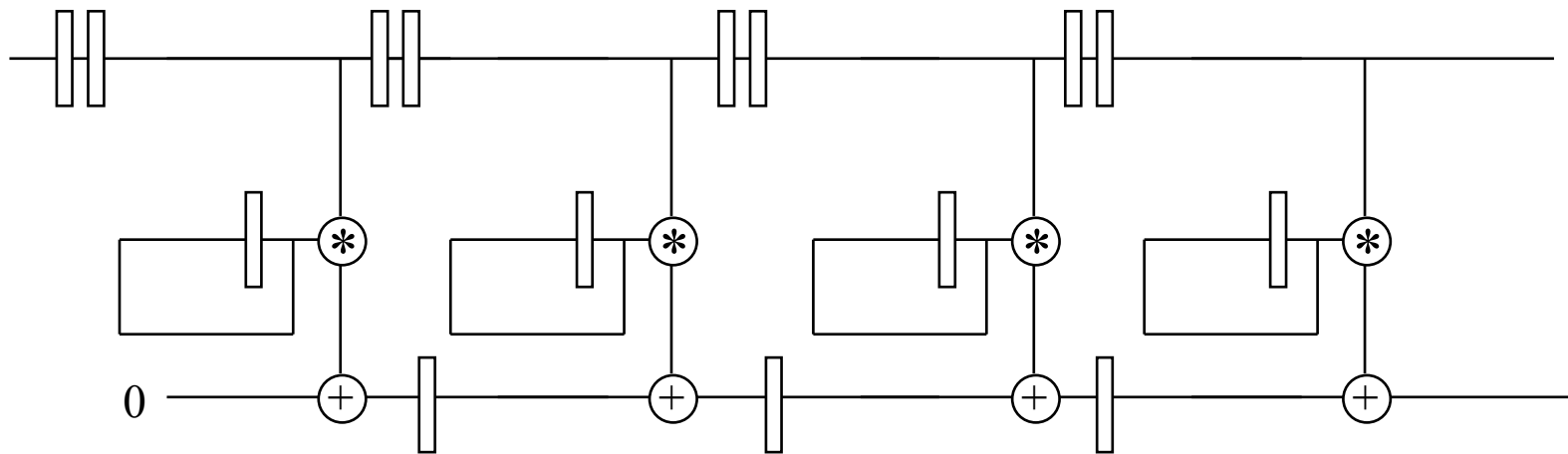
SMT PROCESSORS

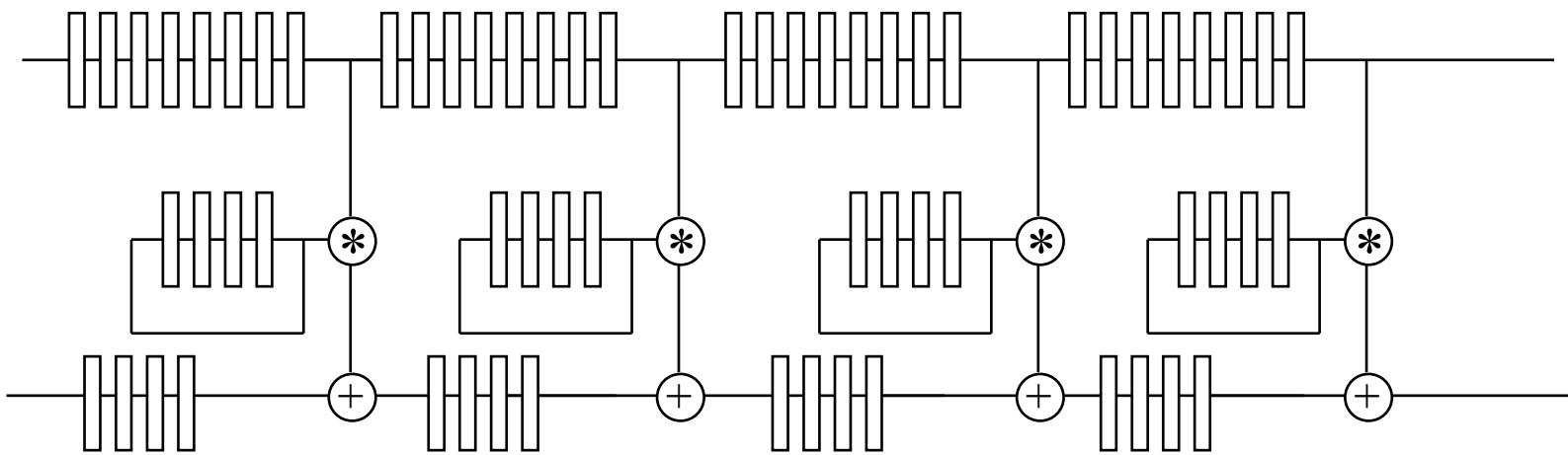


- Dynamic multi-threaded, super-scalar, out-of-order processors
 - Hugely complicated
- Take CSE 471 !!

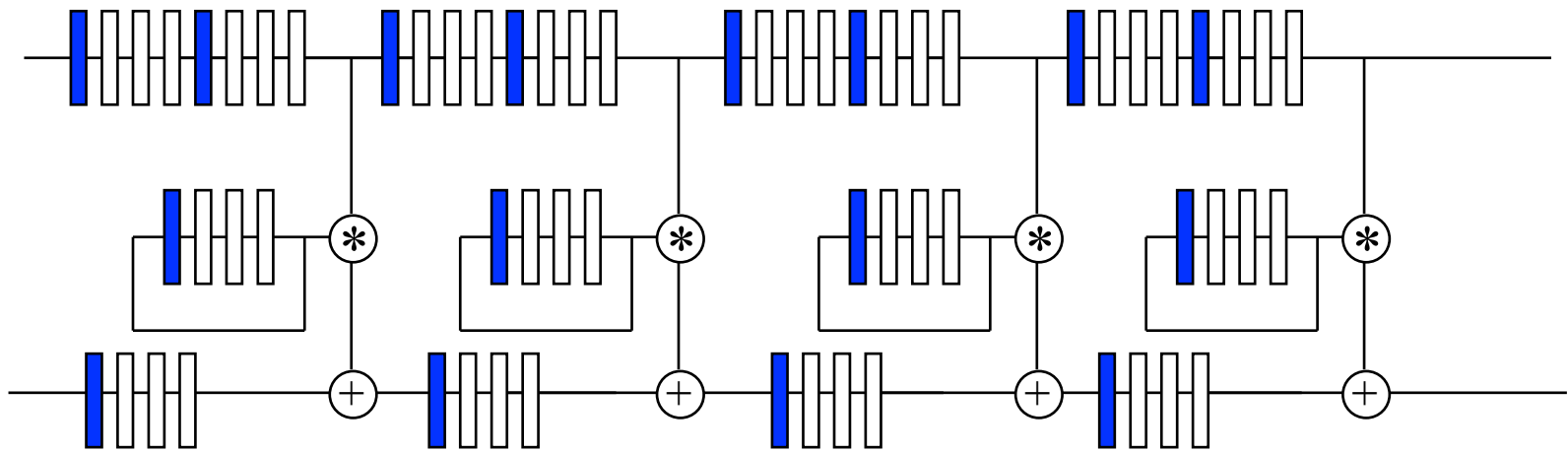
C-SLOWING/RETIMING FOR RESOURCE SHARING

- Correlator circuit

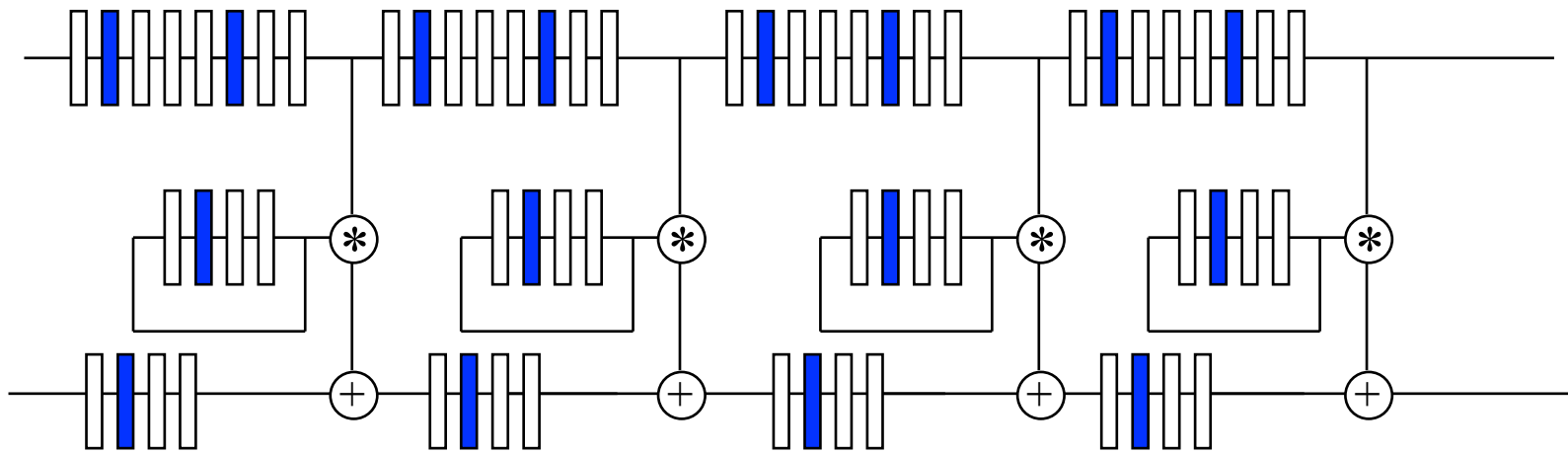


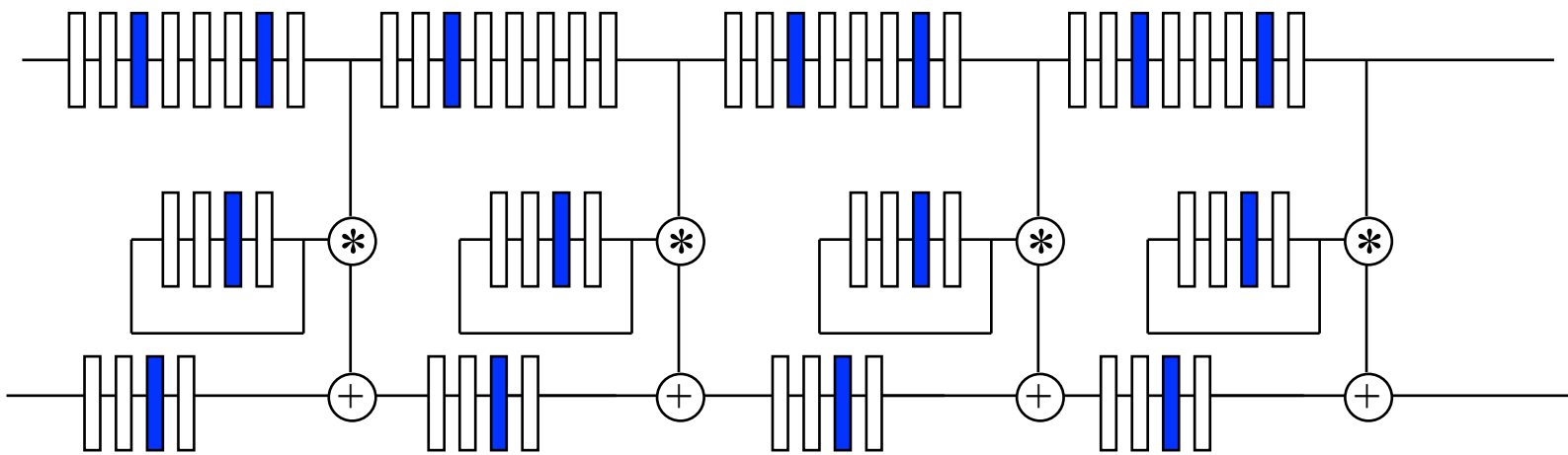


C-SLOWED BY 4

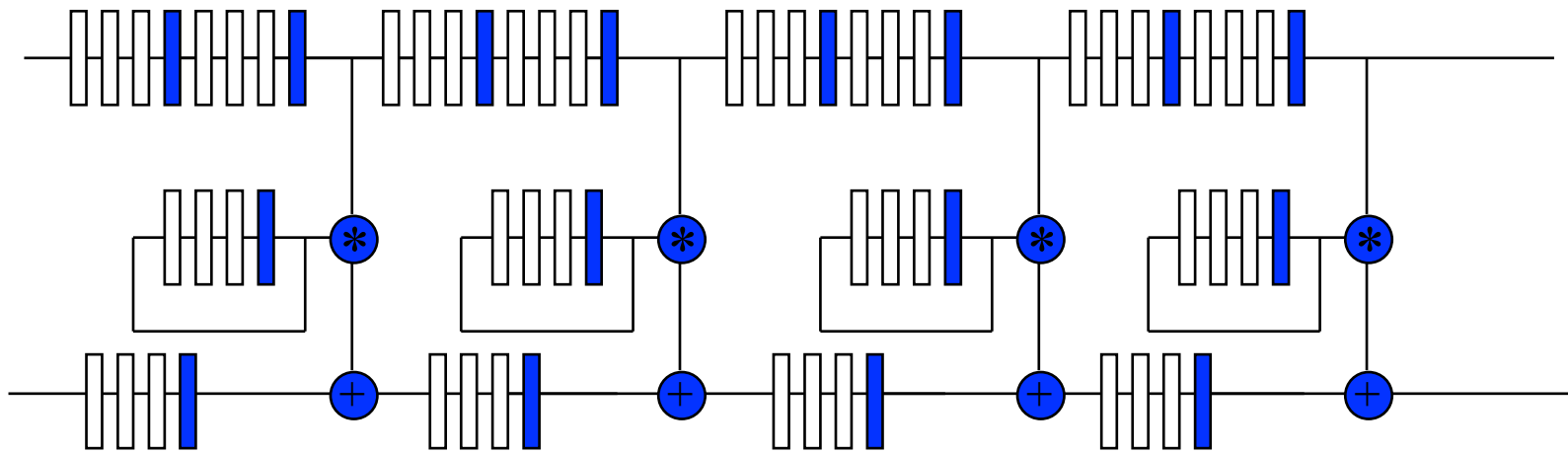


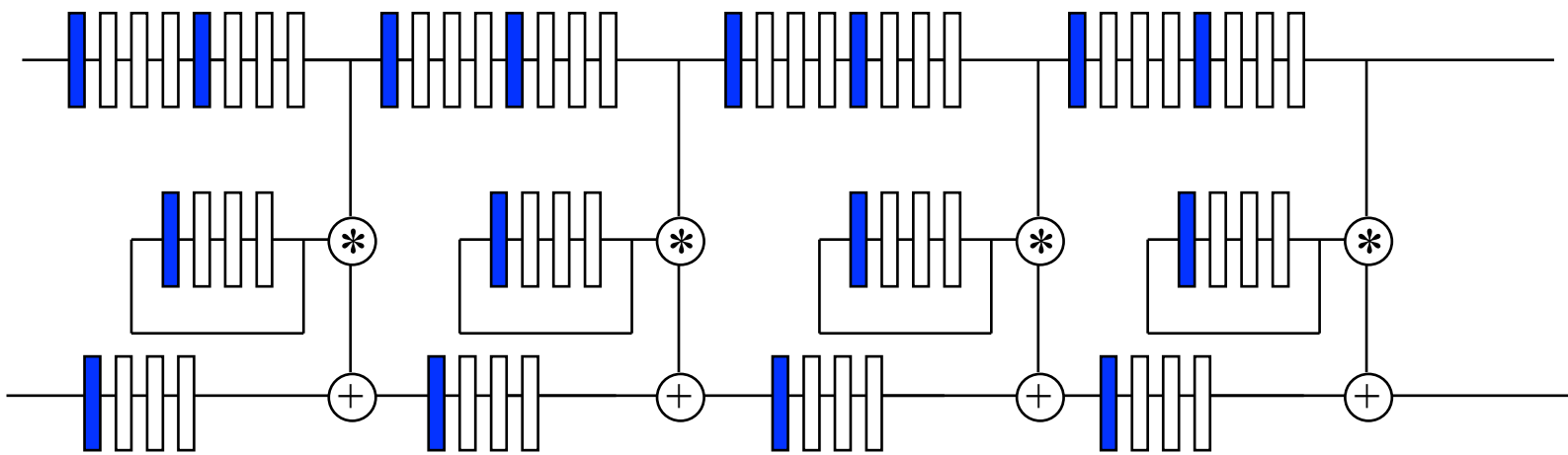
INSERT DATA EVERY 4 CYCLES (ONE DATA SET)

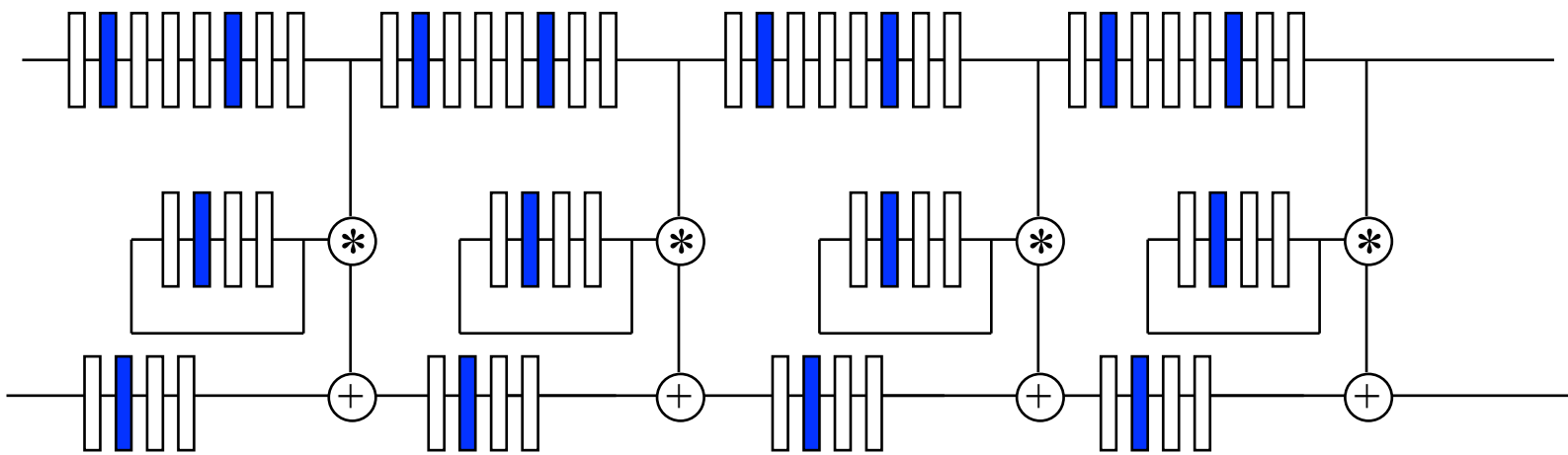


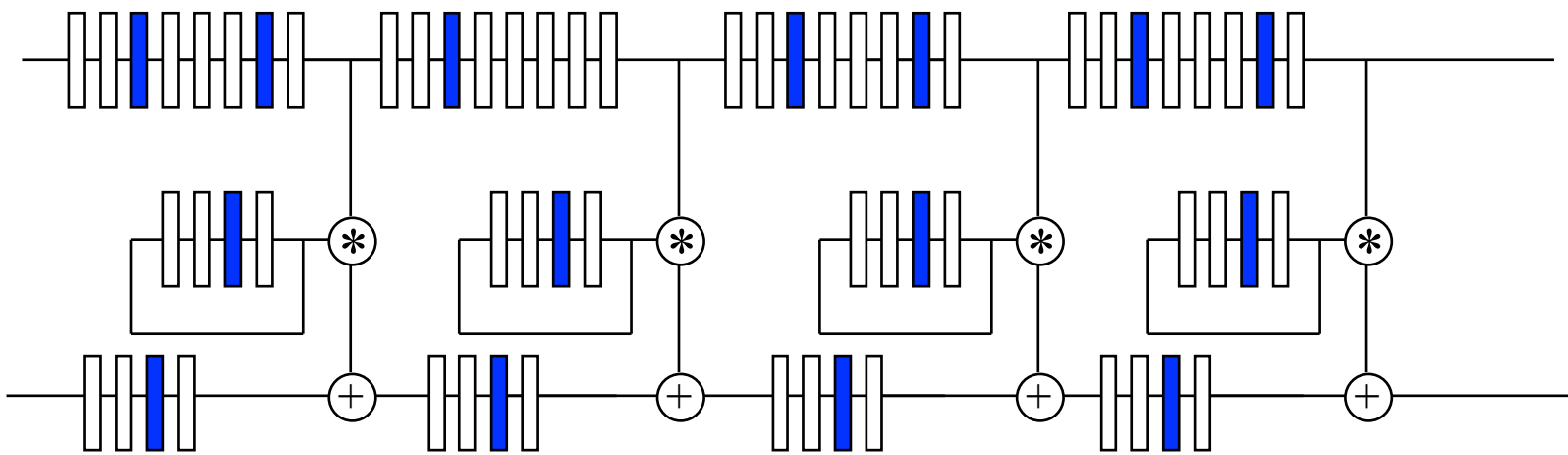


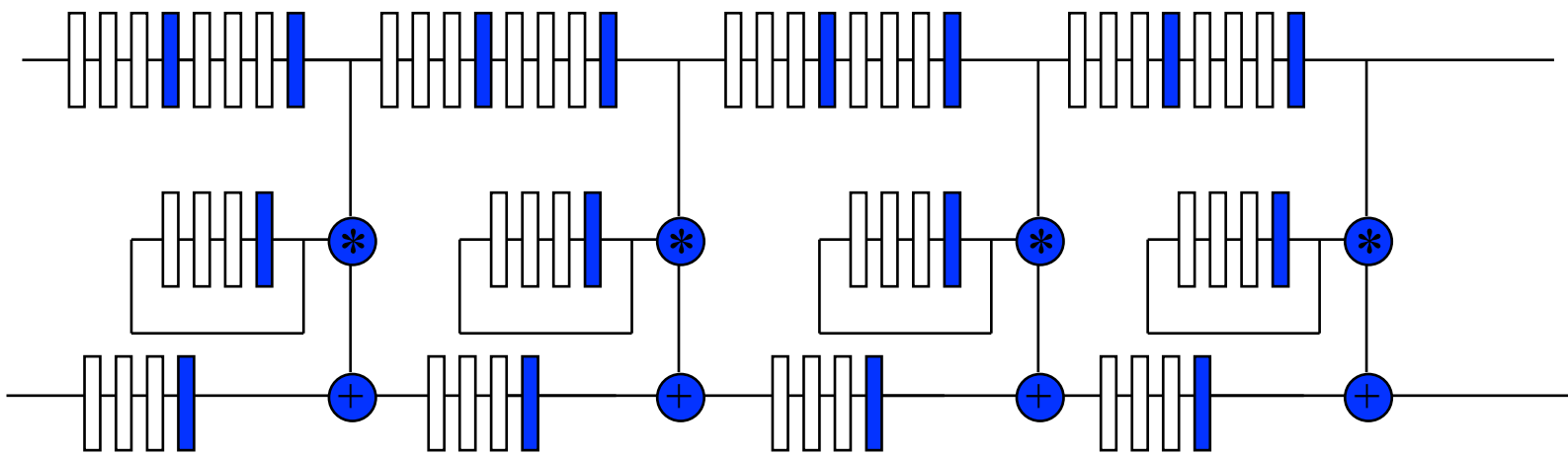
COMPUTATION ACTIVE ONLY EVERY 4 CYCLES



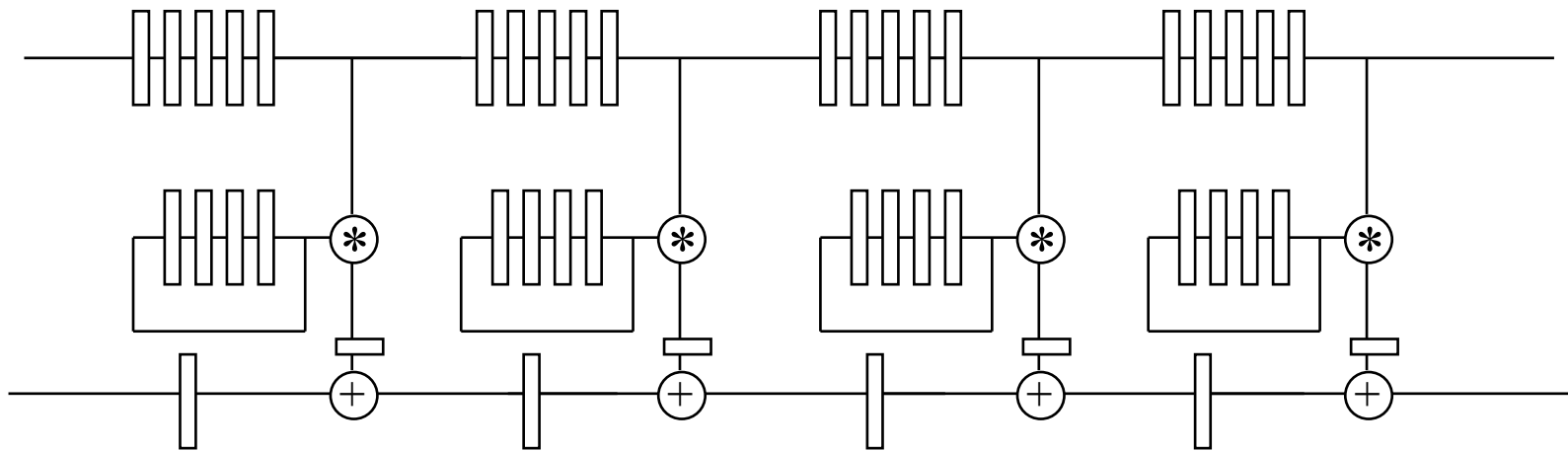


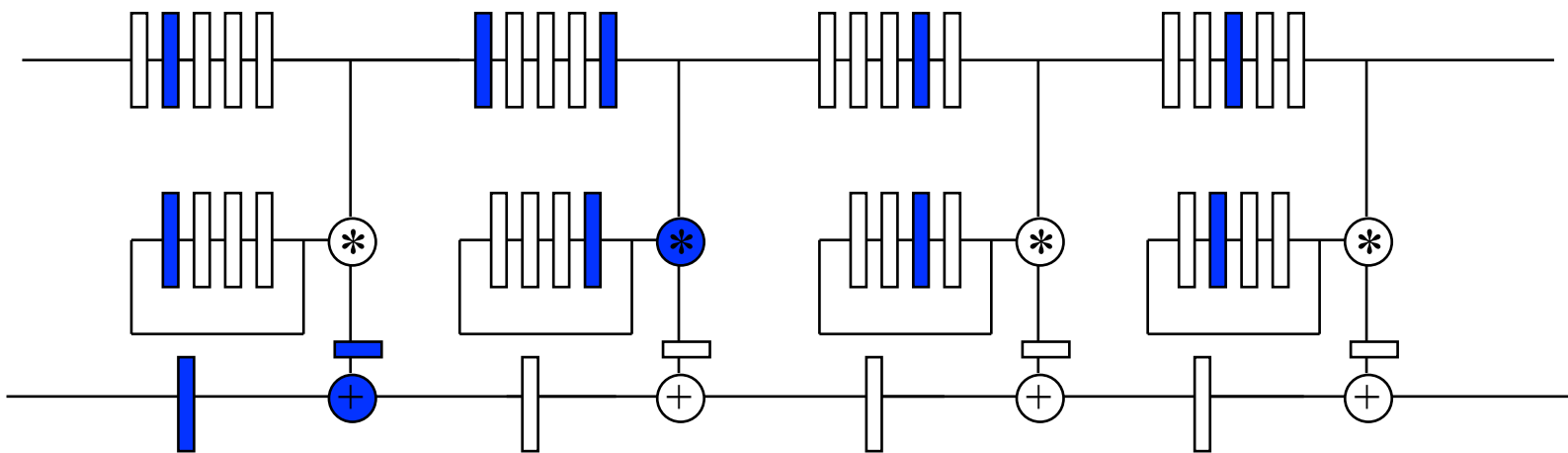


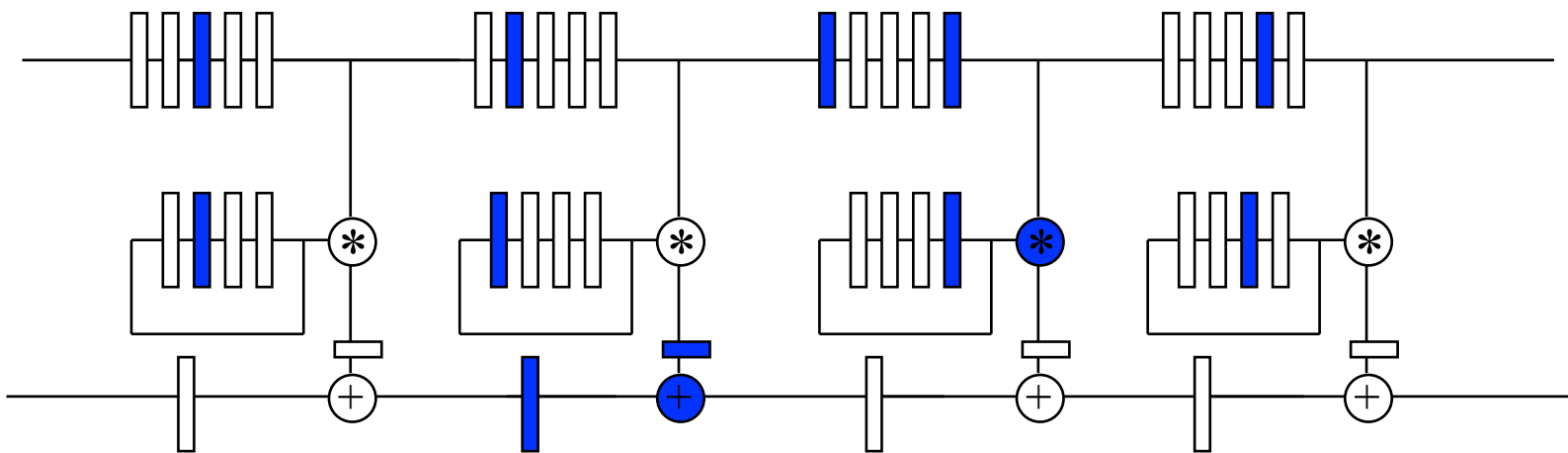


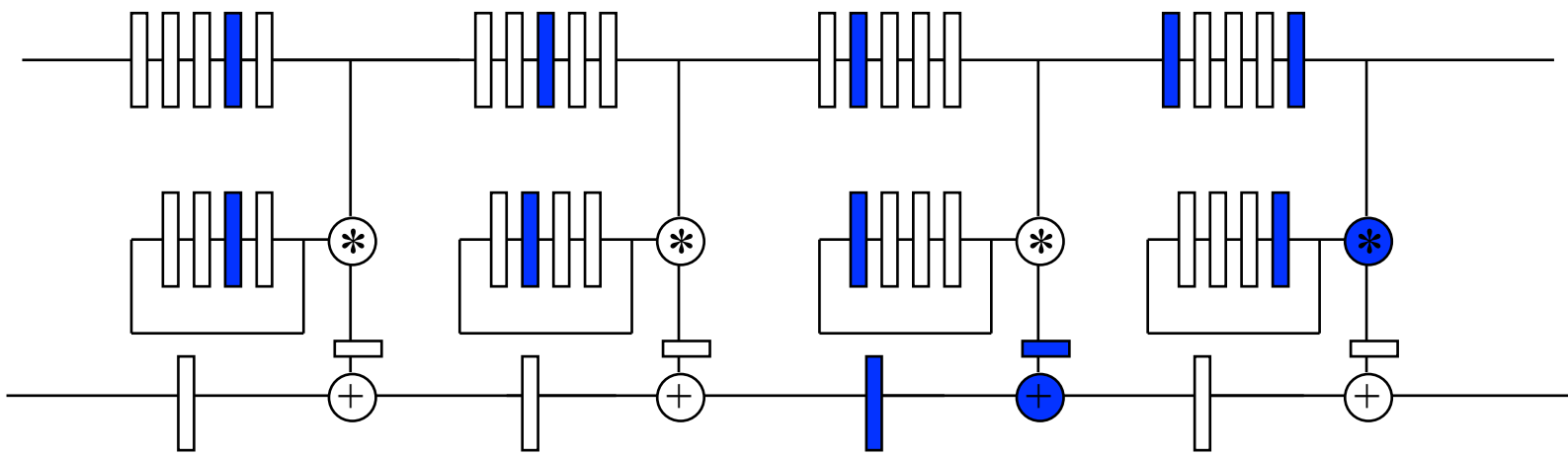


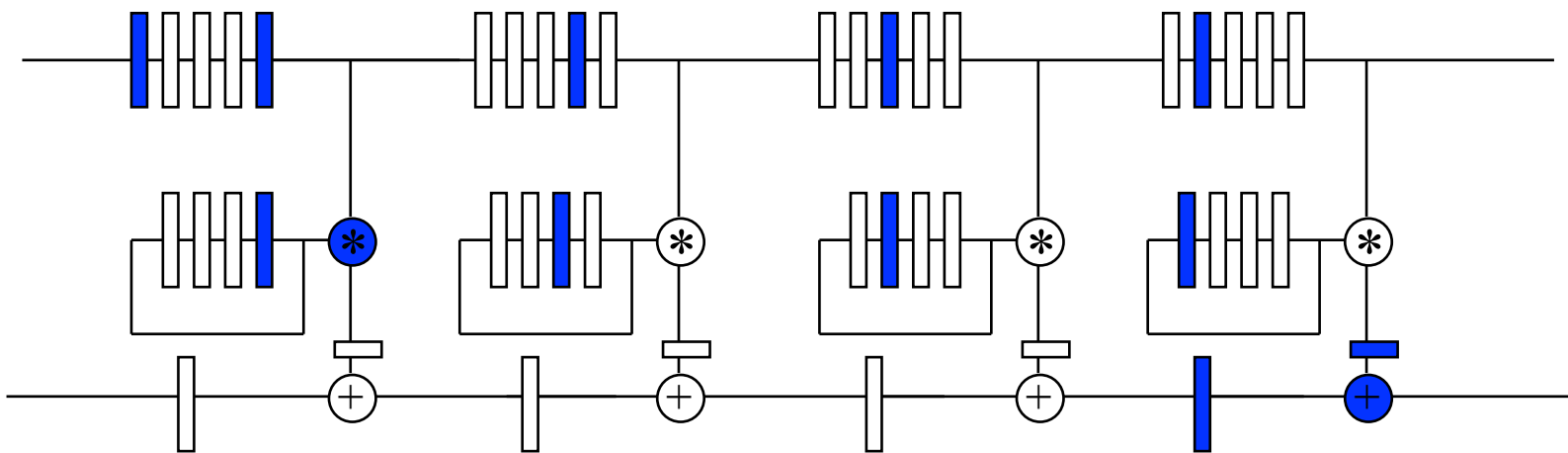
RETIME AND REMOVE EXTRA PIPELINING

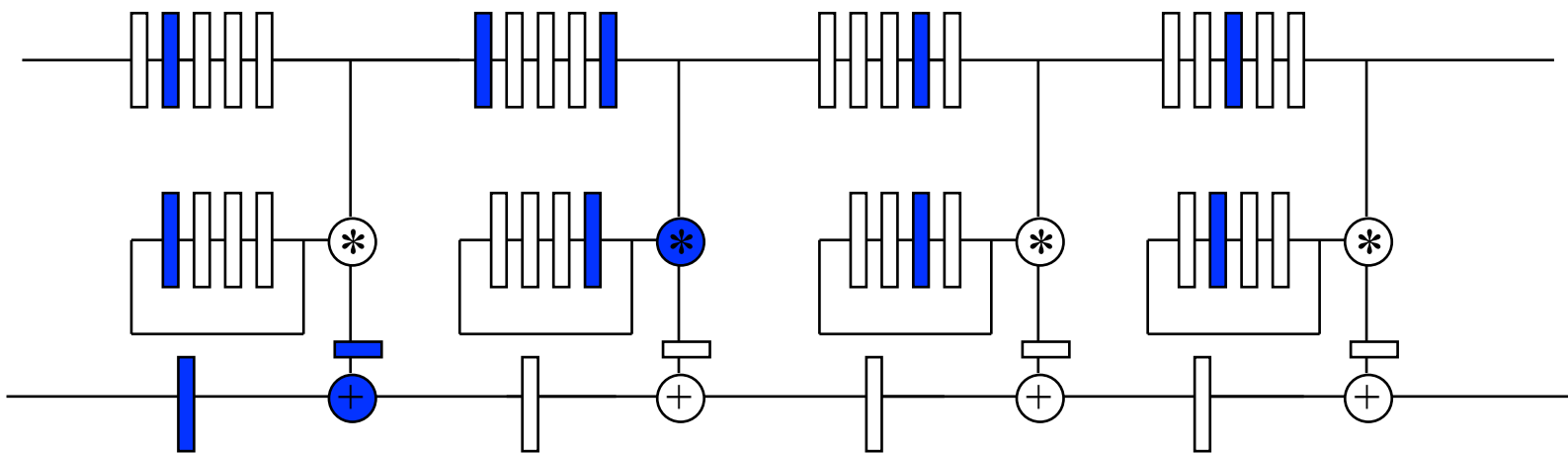


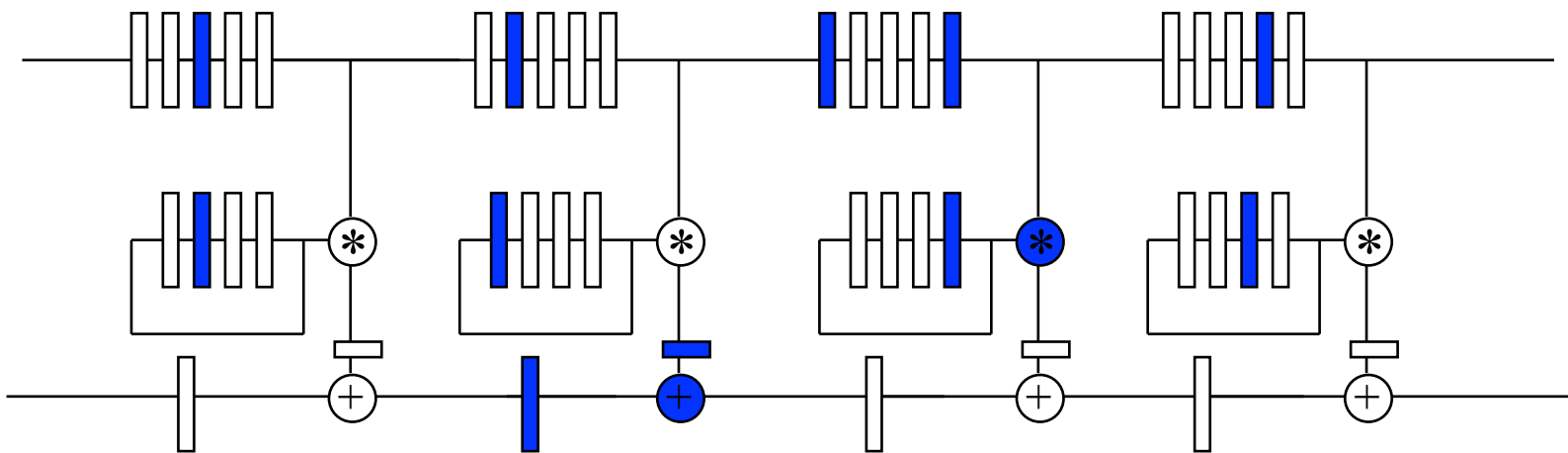


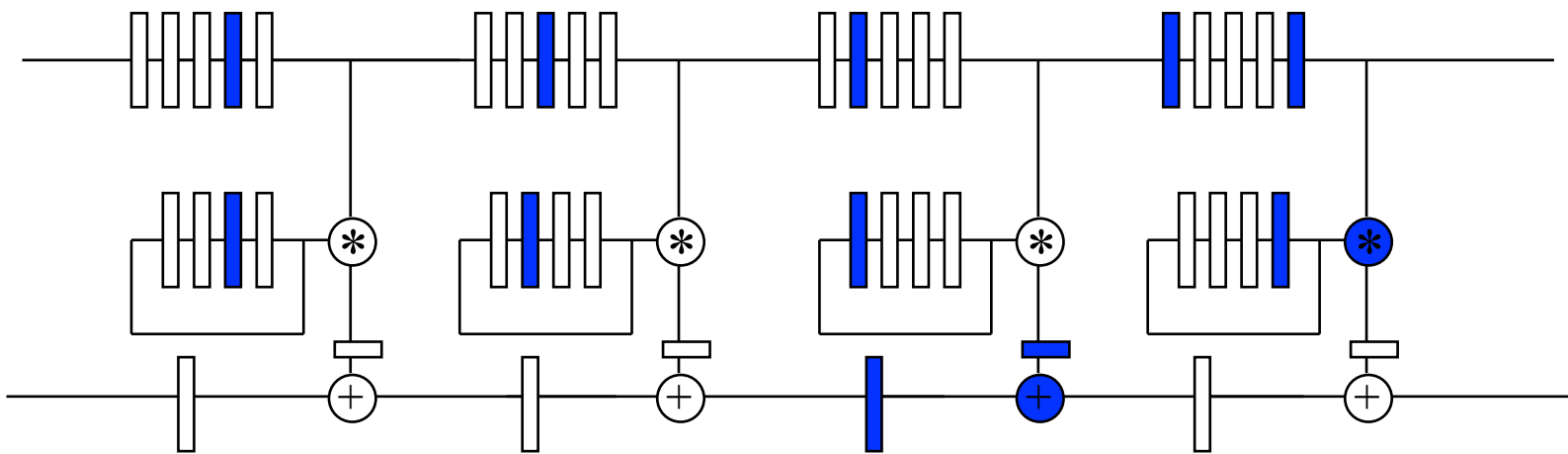












SYSTOLIC ARRAYS



- 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

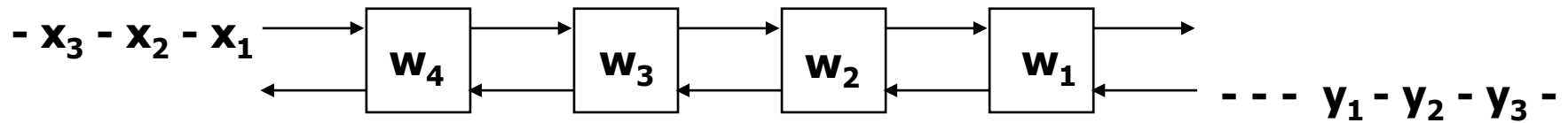
■ $y_j = x_j w_1 + x_{j+1} w_2 + \dots + x_{j+n-1} w_n$

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

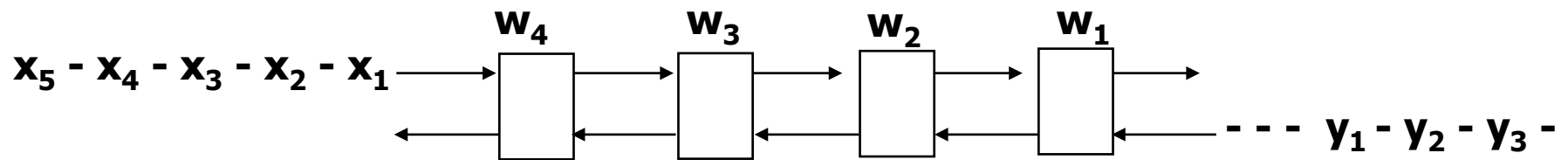
$$y_2 = x_2 w_1 + x_3 w_2 + x_4 w_3 + x_5 w_4$$

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

.....



EXAMPLE - CONVOLUTION

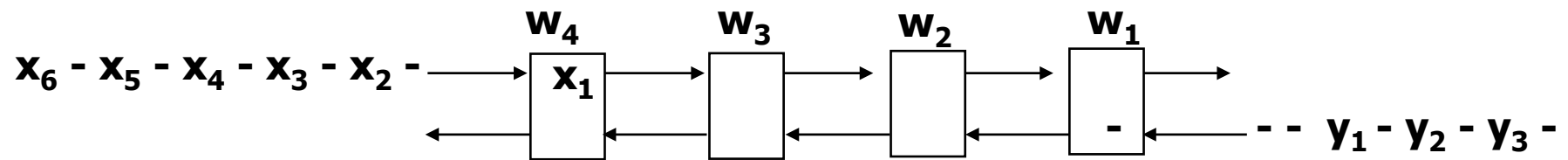


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

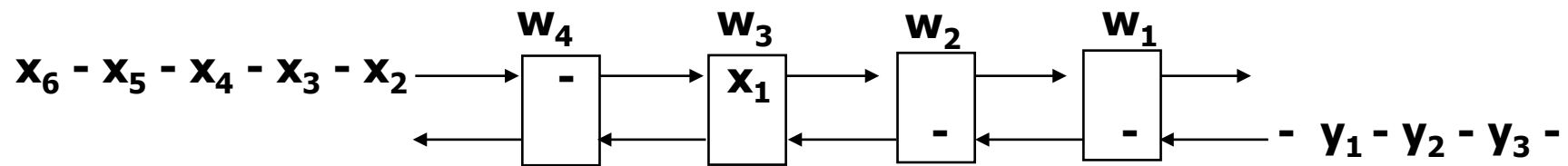


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

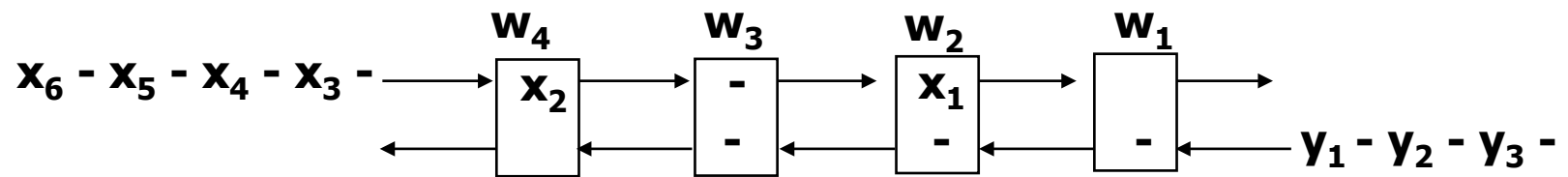


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

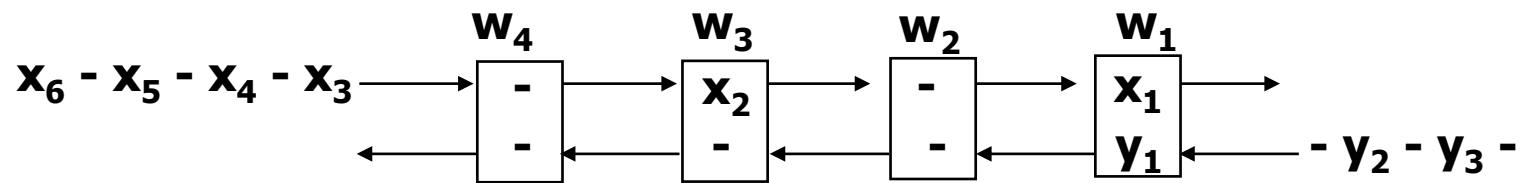


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

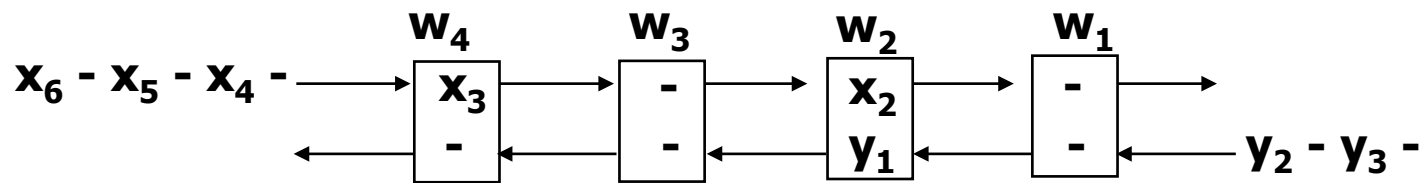


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

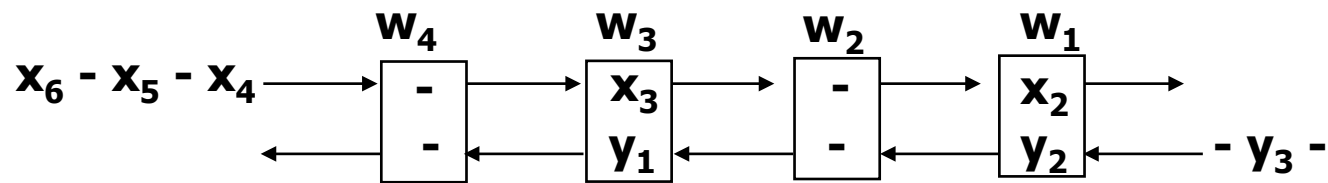


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

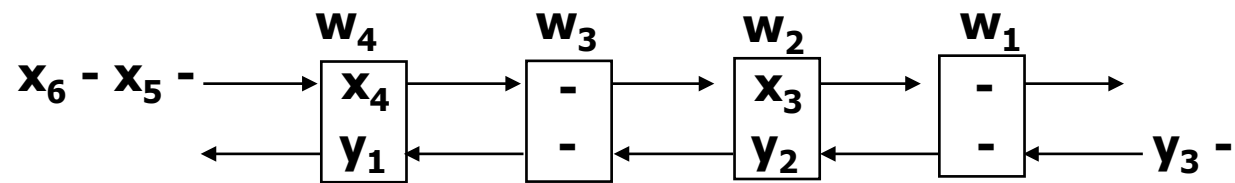


$$y_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4$$

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

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

EXAMPLE - CONVOLUTION

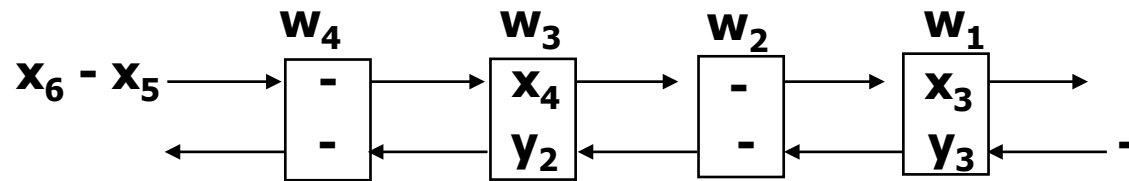


$$y_1 = x_1w_1 + x_2w_2 + x_3w_3 + x_4w_4$$

$$y_2 = x_2w_1 + x_3w_2 + x_4w_3 + x_5w_4$$

$$y_3 = x_3w_1 + x_4w_2 + x_5w_3 + x_6w_4$$

EXAMPLE - CONVOLUTION

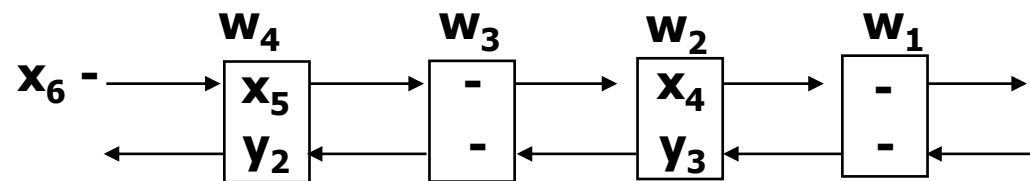


$$y_1 = x_1w_1 + x_2w_2 + x_3w_3 + x_4w_4$$

$$y_2 = x_2w_1 + x_3w_2 + x_4w_3 + x_5w_4$$

$$y_3 = x_3w_1 + x_4w_2 + x_5w_3 + x_6w_4$$

EXAMPLE - CONVOLUTION



$$y_1 = x_1w_1 + x_2w_2 + x_3w_3 + x_4w_4$$

$$y_2 = x_2w_1 + x_3w_2 + x_4w_3 + x_5w_4$$

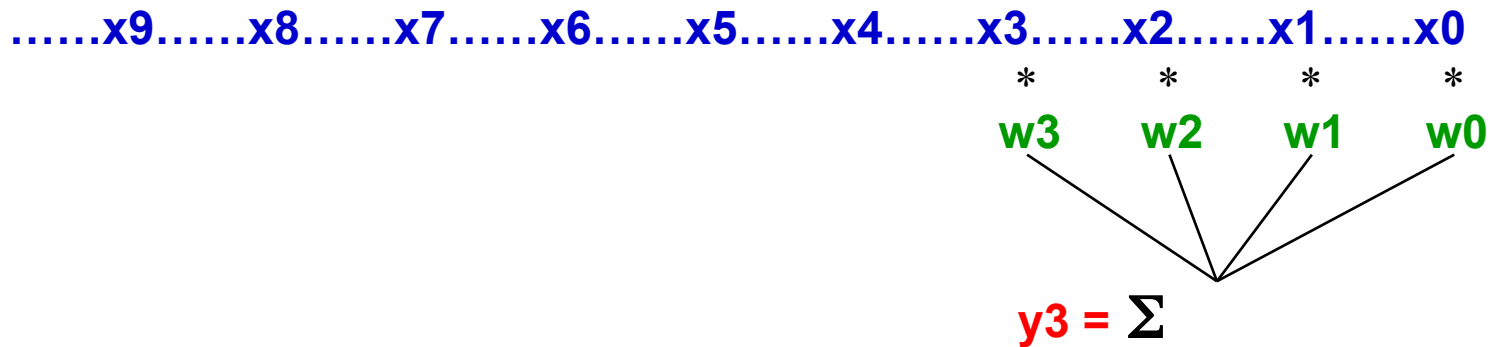
$$y_3 = x_3w_1 + x_4w_2 + x_5w_3 + x_6w_4$$

EXAMPLE - CONVOLUTION

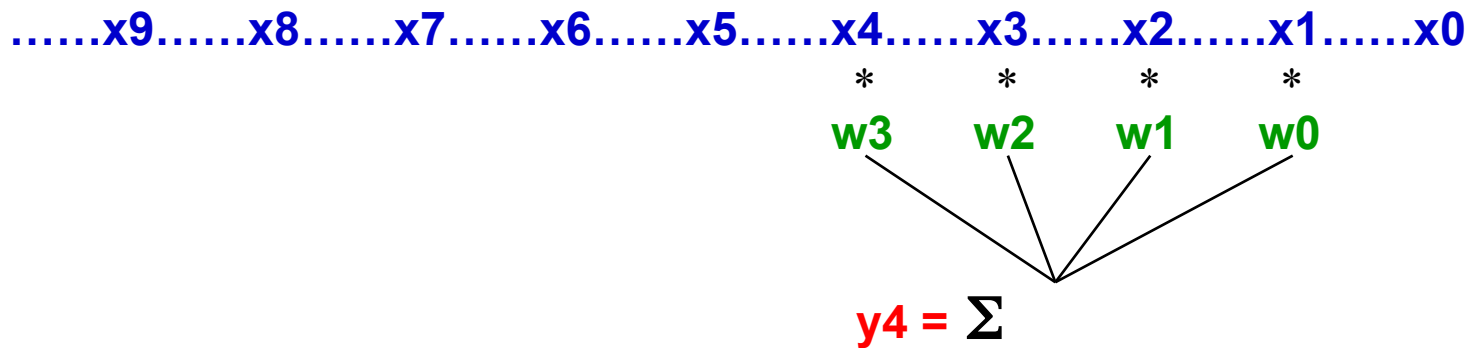
| | W₄ | W₃ | W₂ | W₁ | |
|--|----------------------|----------------------|----------------------|----------------------|--|
| x₆ - x₅ - x₄ - x₃ - x₂ - | x₁ | | | - | - - y₁ - y₂ - y₃ |
| x₆ - x₅ - x₄ - x₃ - x₂ | - | x₁ | | - | - y₁ - y₂ - y₃ |
| x₆ - x₅ - x₄ - x₃ - | x₂ | - | x₁ | - | y₁ - y₂ - y₃ |
| x₆ - x₅ - x₄ - x₃ | - | x₂ | - | x₁ | - y₂ - y₃ |
| x₆ - x₅ - x₄ - | x₃ | - | x₂ | - | y₂ - y₃ |
| x₆ - x₅ - x₄ | - | x₃ | - | x₂ | - y₃ |
| x₆ - x₅ - | x₄ | - | x₃ | - | y₃ |
| x₆ - x₅ | - | x₄ | - | x₃ | |

CONVOLUTION - ANOTHER LOOK

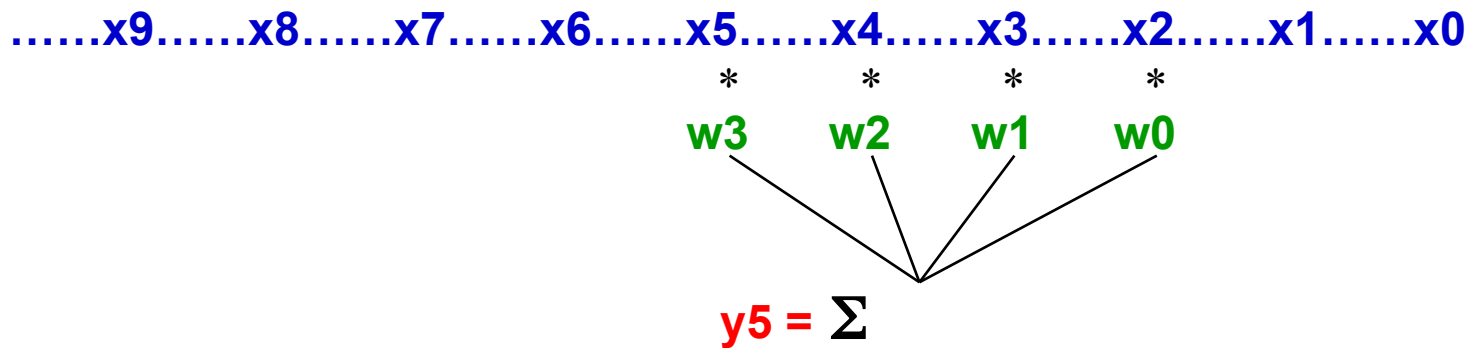
- Repeated vector product



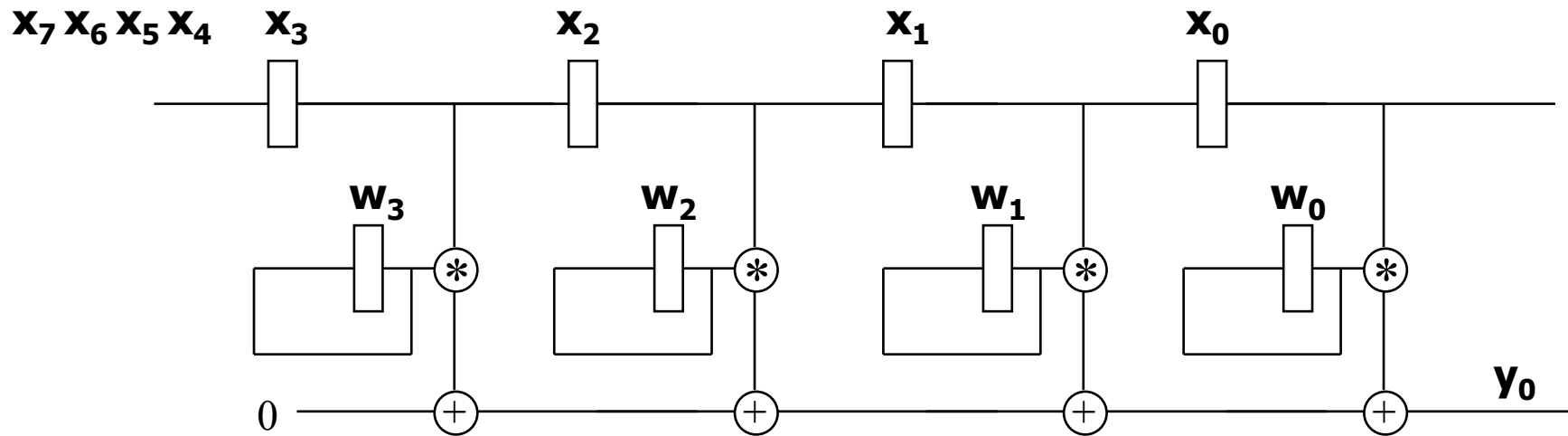
EXAMPLE: CONVOLUTION



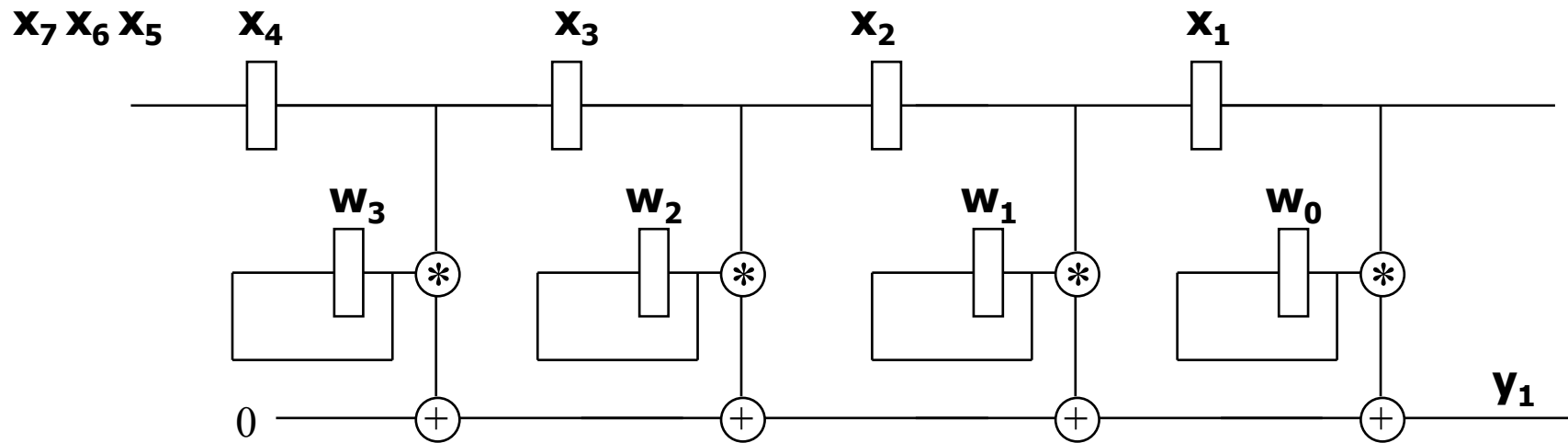
EXAMPLE: CONVOLUTION



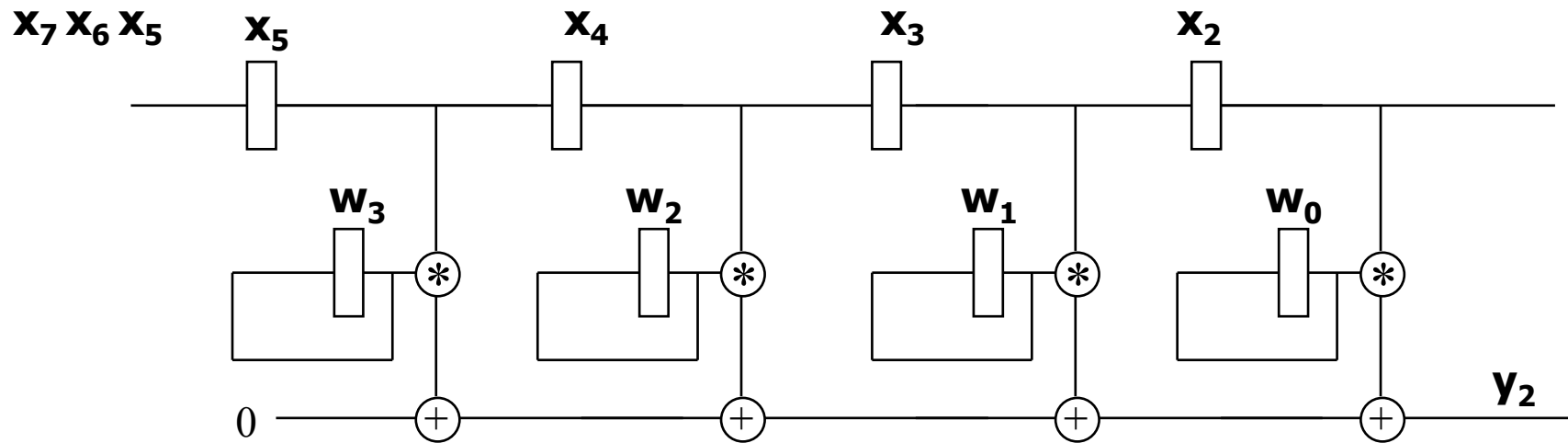
CONVOLUTION EXAMPLE

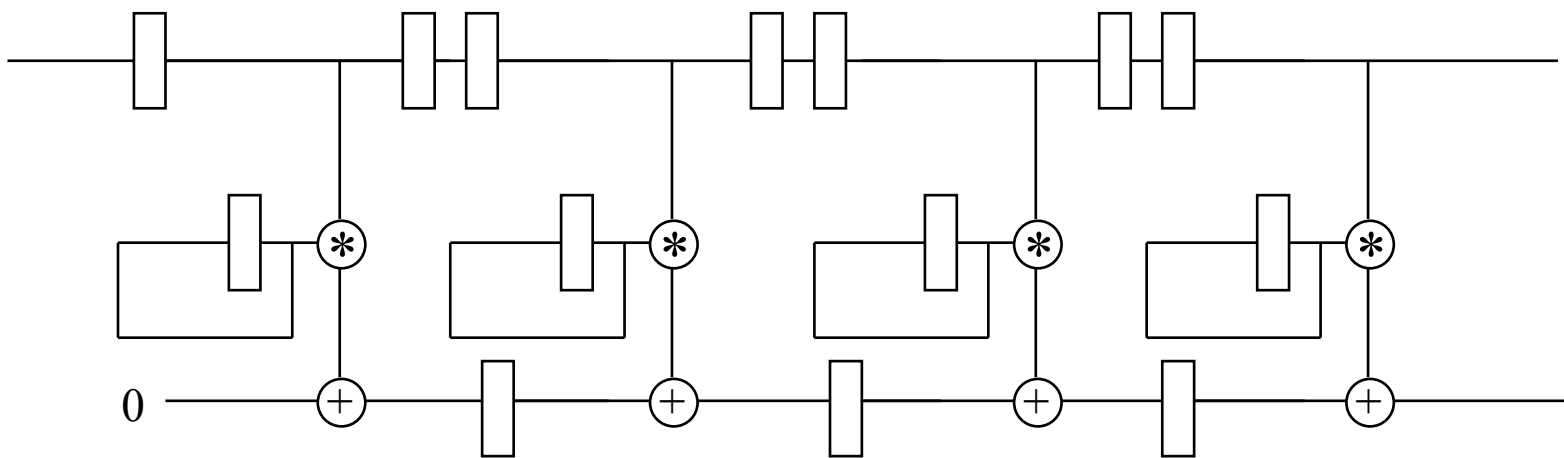


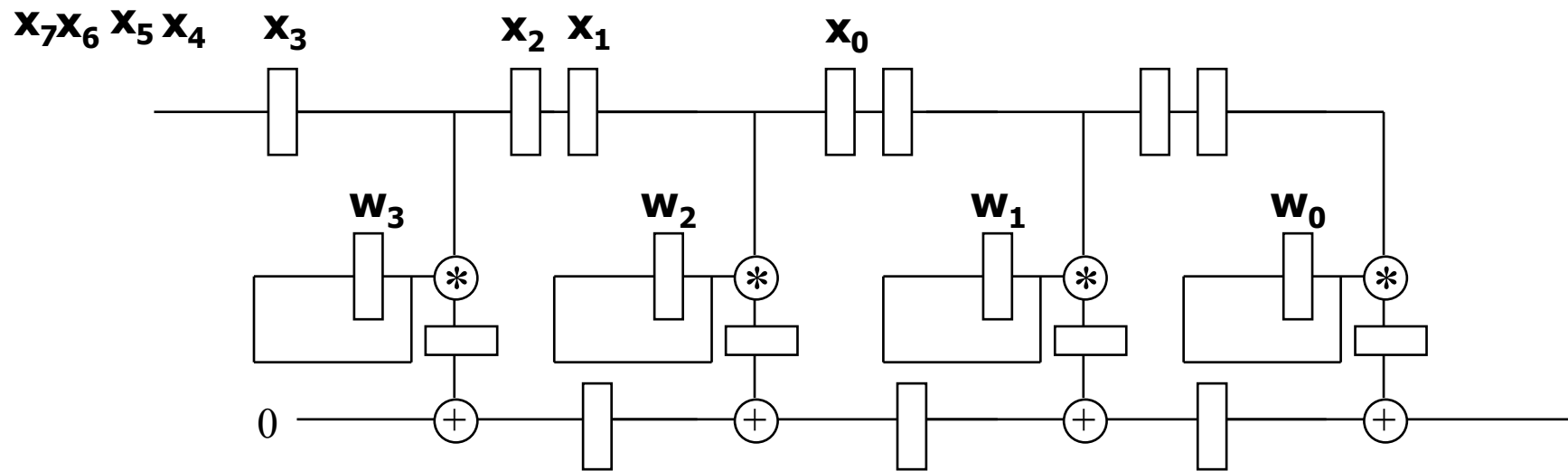
CONVOLUTION EXAMPLE

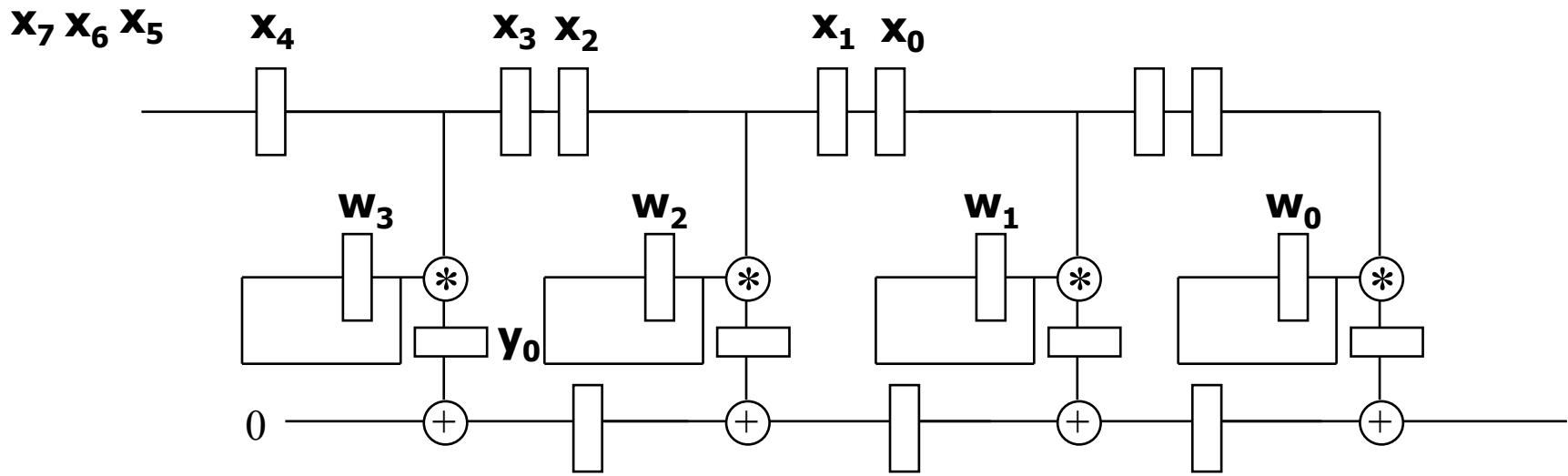


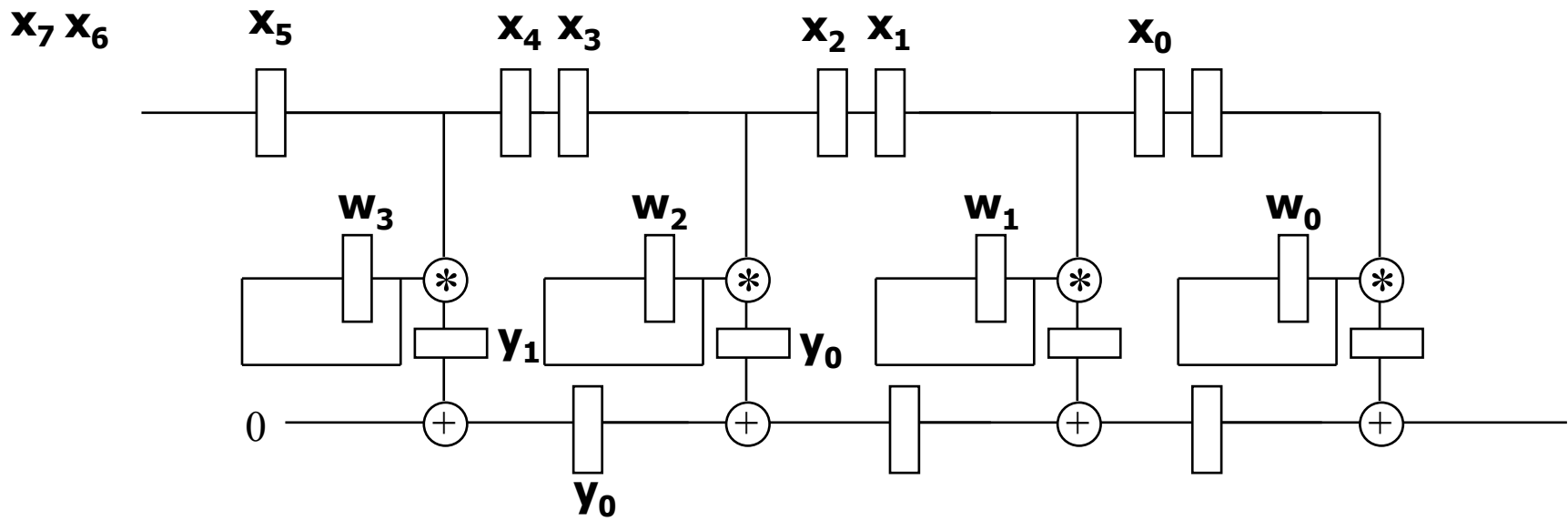
CONVOLUTION EXAMPLE

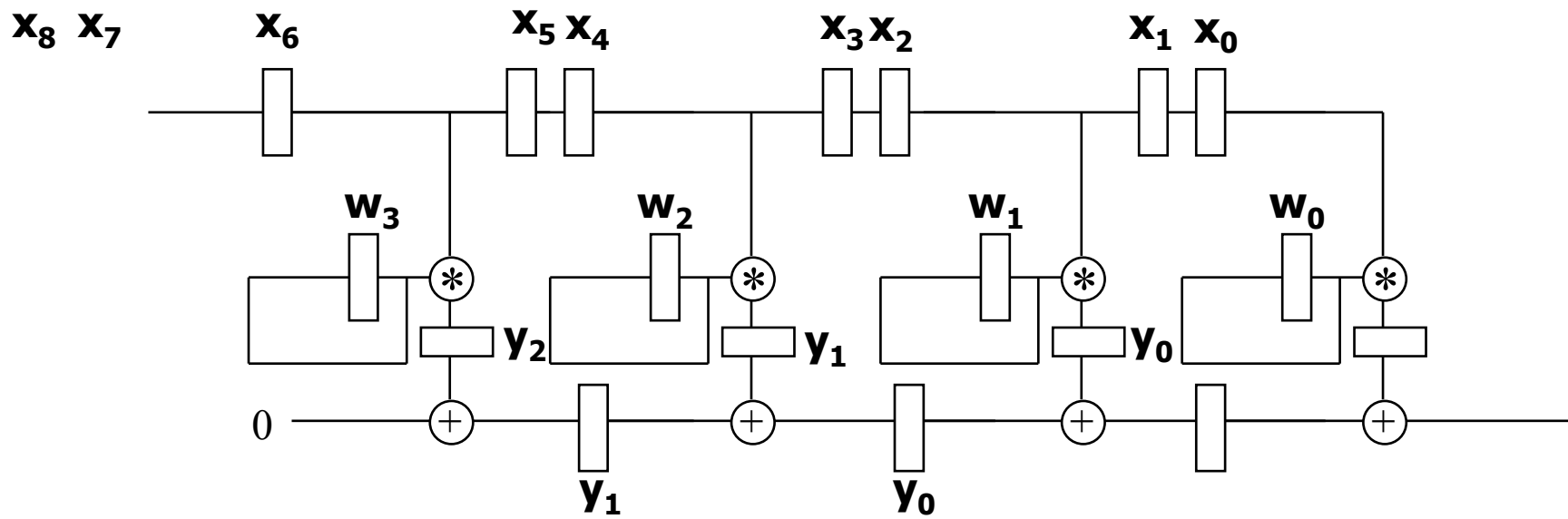


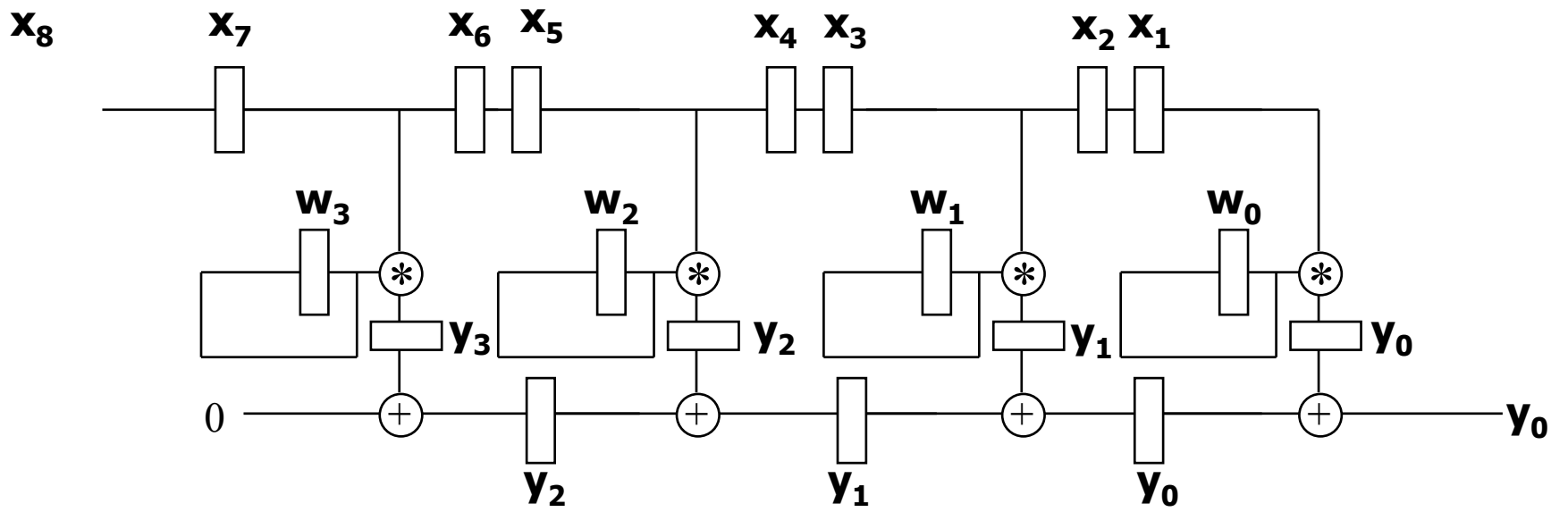


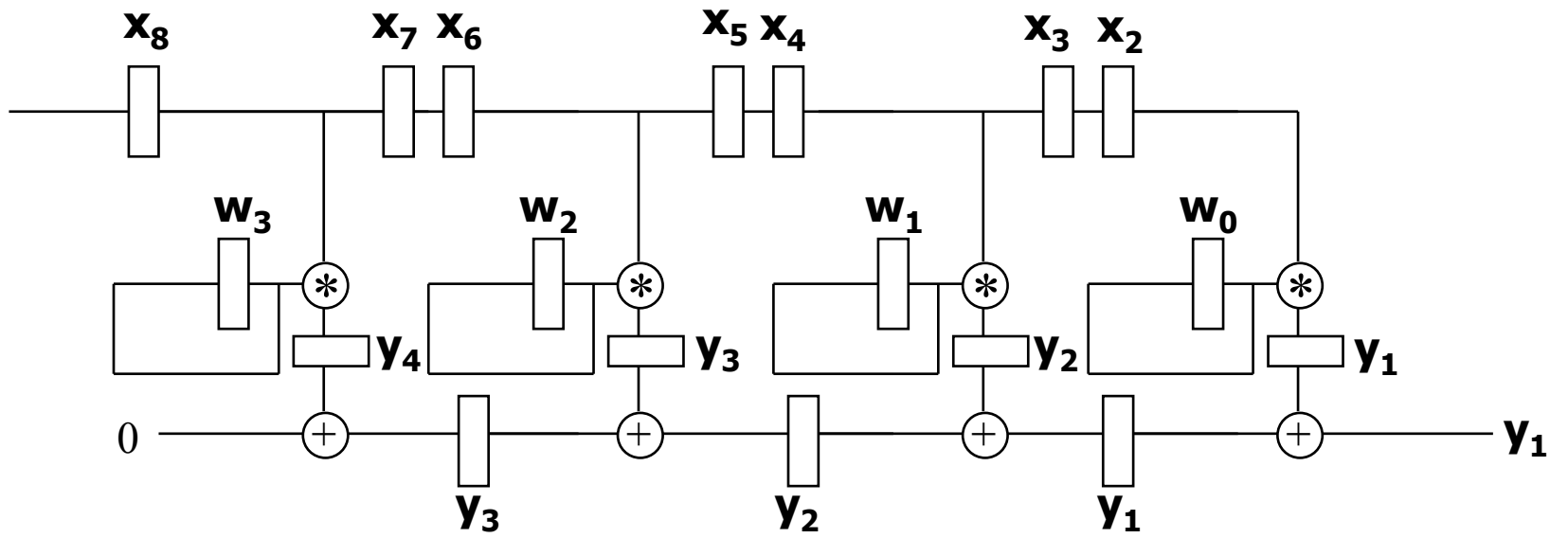






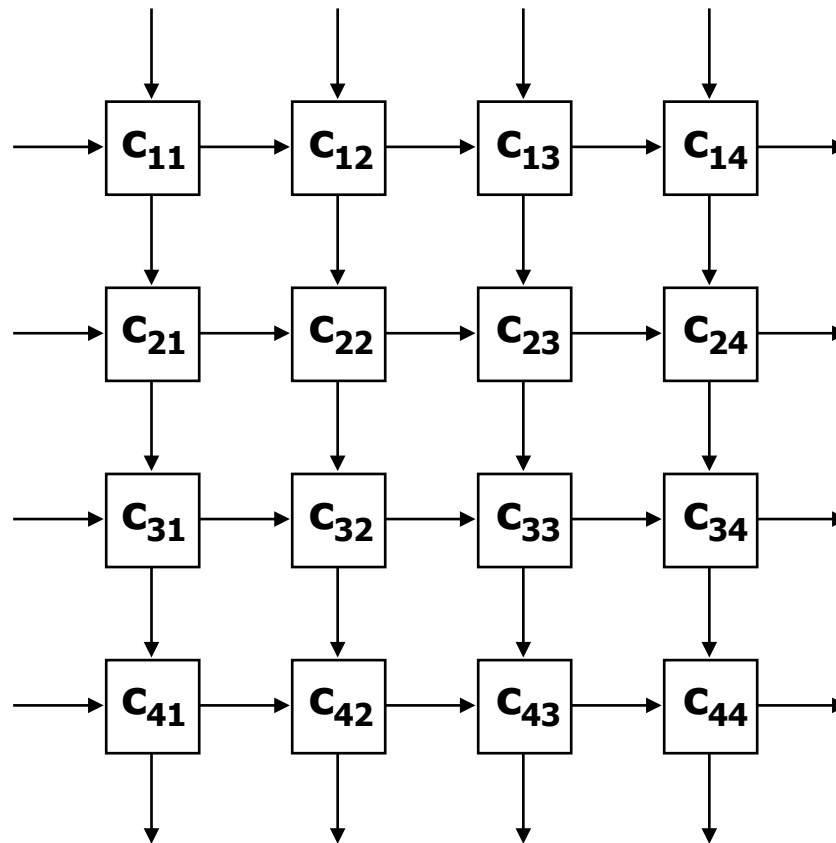




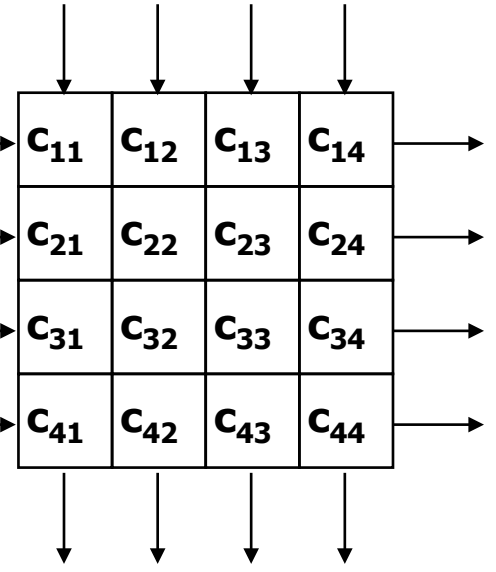
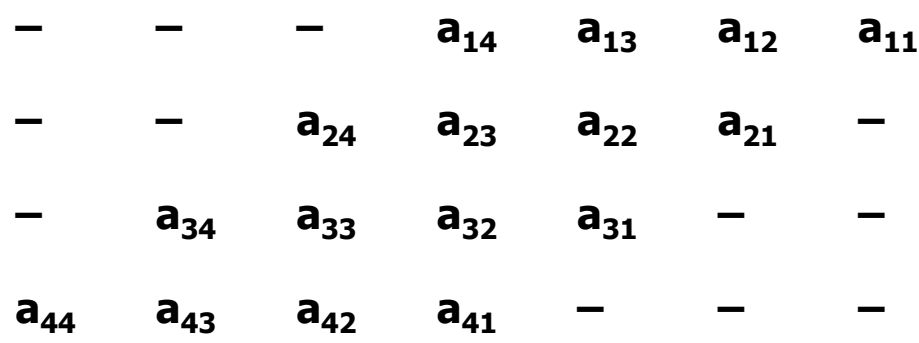
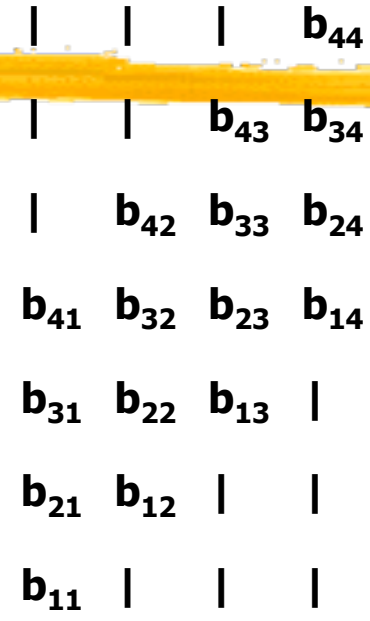


EXAMPLE: MATRIX MULTIPLICATION

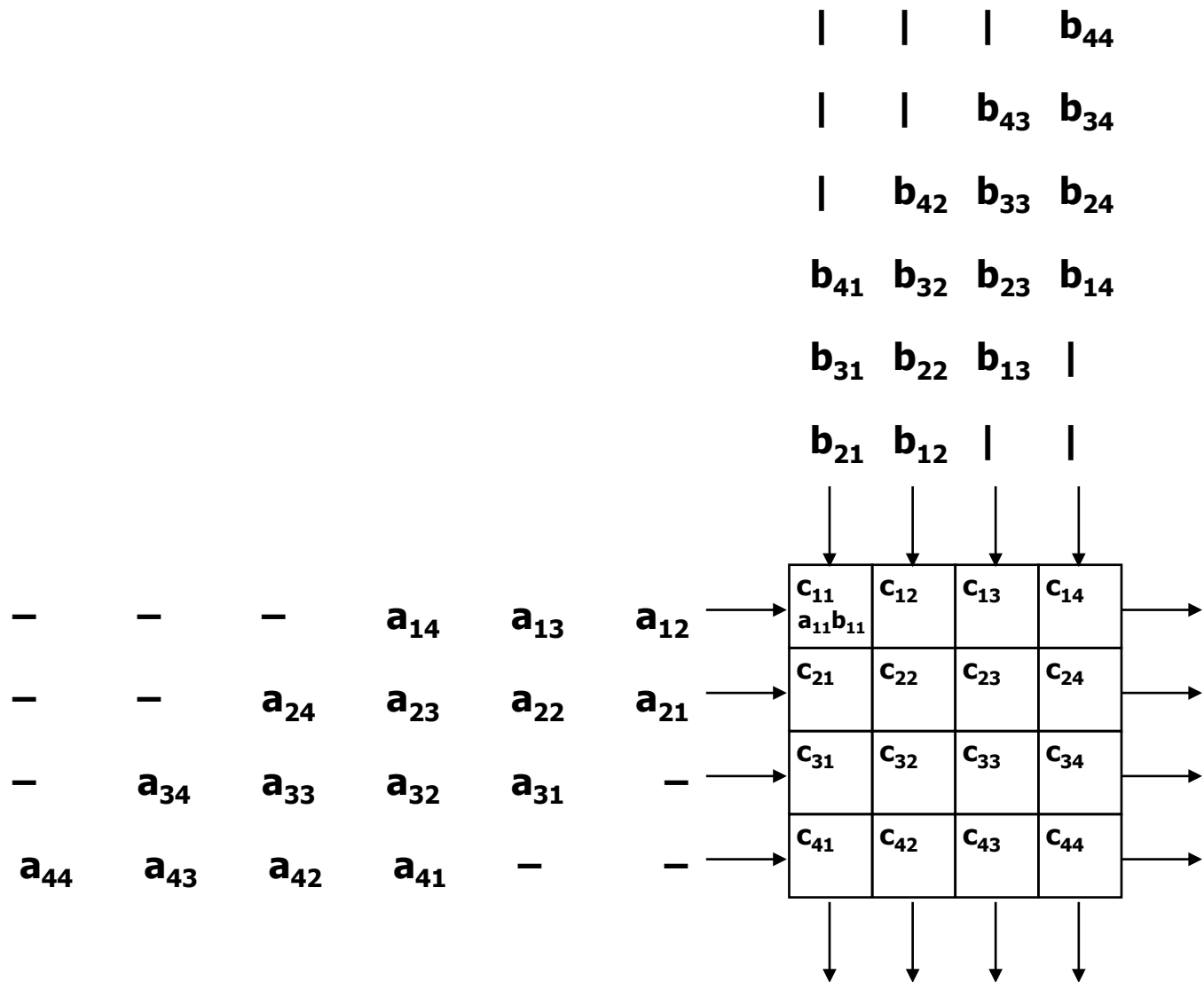
■ $C = A \times B$ $c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$



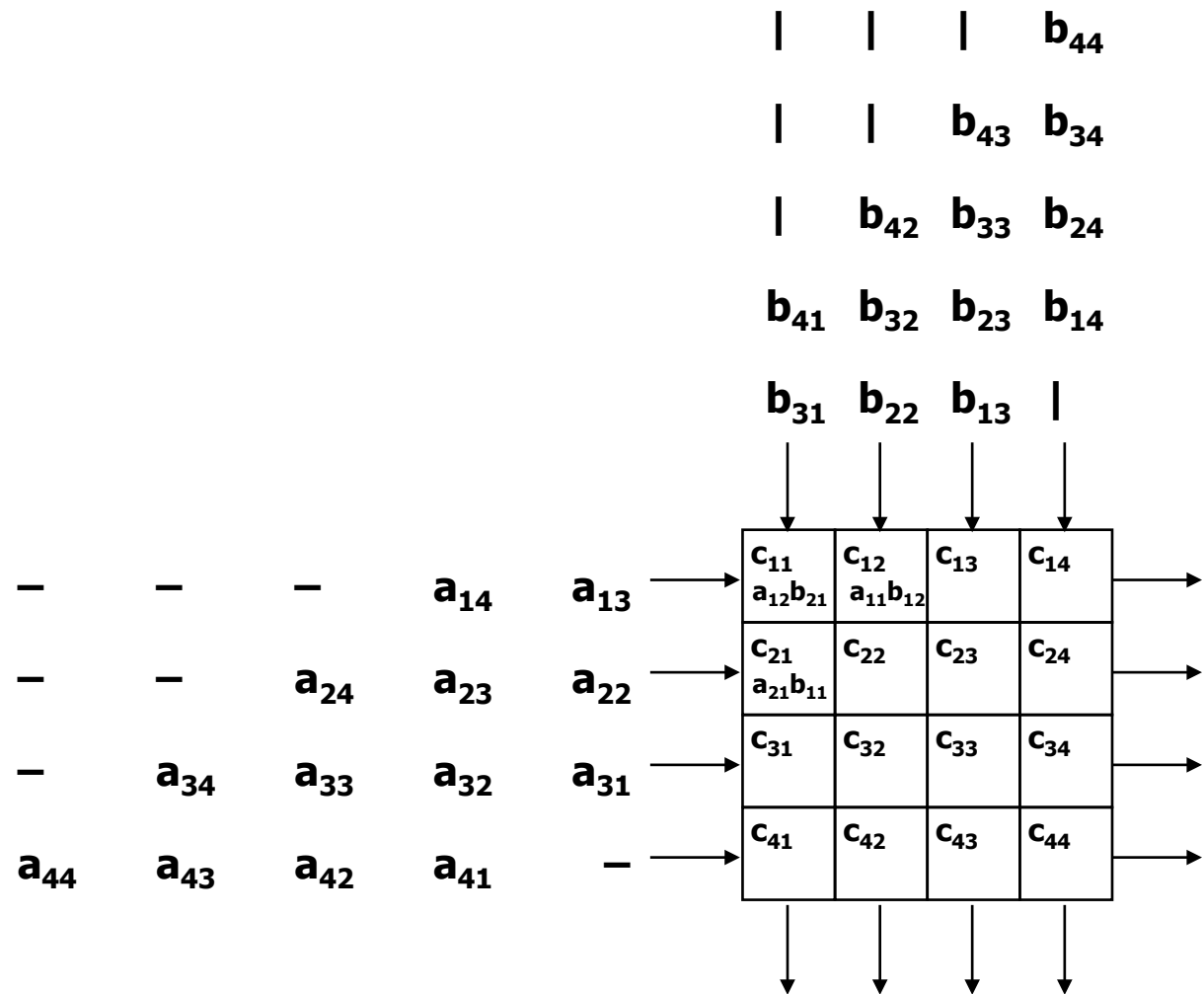
EXAMPLE: MATRIX MULTIPLICATION



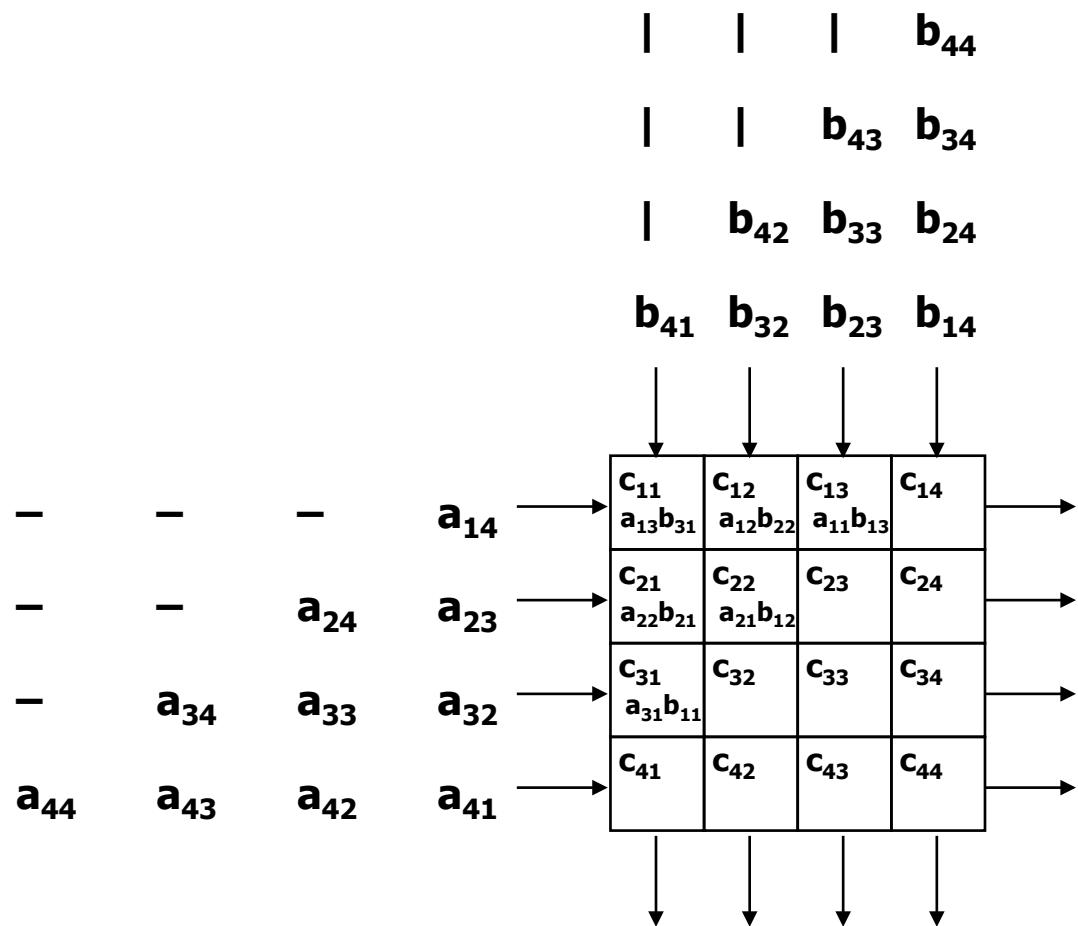
EXAMPLE: MATRIX MULTIPLICATION



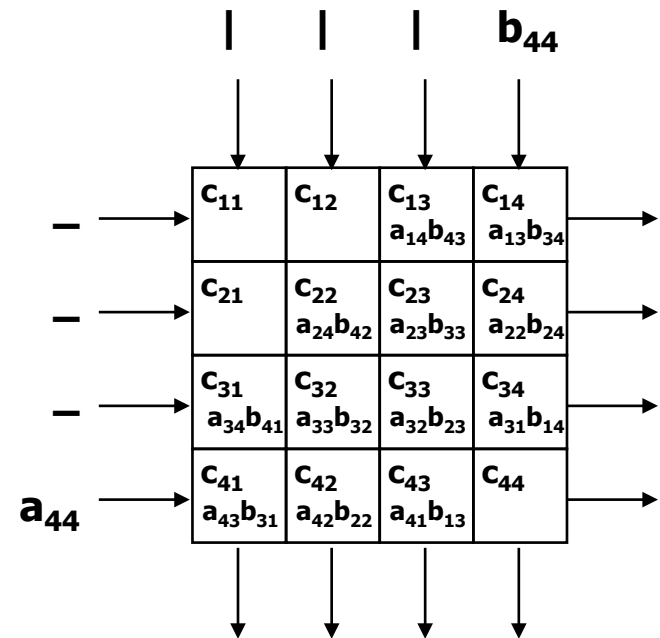
EXAMPLE: MATRIX MULTIPLICATION



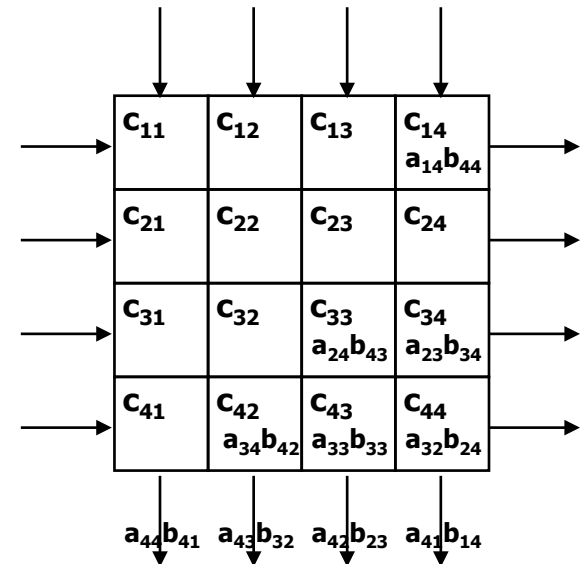
EXAMPLE: MATRIX MULTIPLICATION



EXAMPLE: MATRIX MULTIPLICATION



EXAMPLE: MATRIX MULTIPLICATION



SYSTOLIC ALGORITHMS



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

SYSTOLIC ARCHITECTURES



- Highly parallel
 - “fine-grained” parallelism
 - deep pipelining
- Local communication
 - wires are short - no global communication (except CLK)
 - linear array → 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

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