

EECS150 - Digital Design

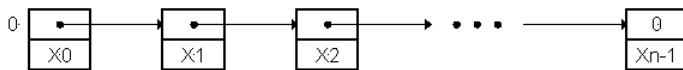
Lecture 25 - High-Level Design

(Part 2)

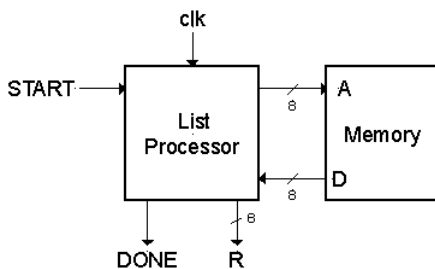
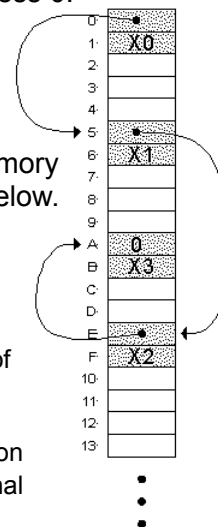
April 21, 2011
John Wawrzynek

List Processor Example

- Design a circuit that forms the sum of all the 2's complements integers stored in a linked-list structure starting at memory address 0:



- All integers and pointers are 8-bit. The link-list is stored in a memory block with an 8-bit address port and 8-bit data port, as shown below. The pointer from the last element in the list is 0.



I/Os:

- START resets to head of list and starts addition process.
- DONE signals completion
- R, Bus that holds the final result

List Example Resource Scheduling

- In this case, first spread out, then pack.

Memory	next ₁			x ₁		
adder		numa ₁			sum ₁	

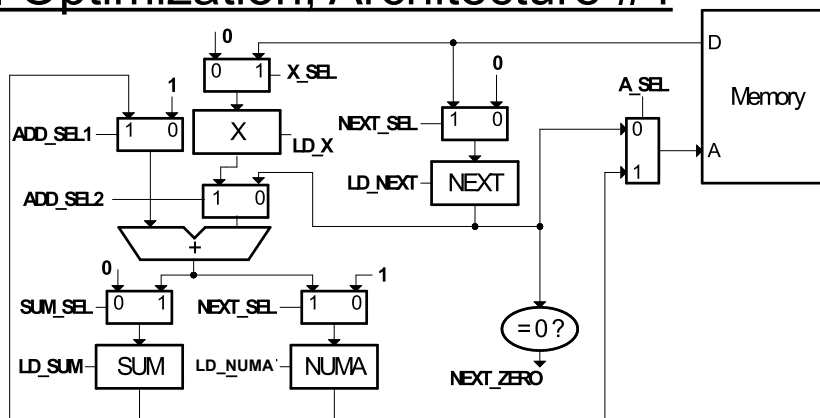
Memory	next ₁		next ₂	x ₁	next ₃	x ₂	next ₄	x ₃	
adder		numa ₁		numa ₂	sum ₁	numa ₃	sum ₂	numa ₄	sum ₃

- $X \leftarrow \text{Memory}[\text{NUMA}], \text{NUMA} \leftarrow \text{NEXT} + 1;$
- $\text{NEXT} \leftarrow \text{Memory}[\text{NEXT}], \text{SUM} \leftarrow \text{SUM} + X;$

- Three different loop iterations active at once.
- Short cycle time (no dependencies within a cycle)
- full utilization (only 2 cycles per result)
- Initialization: $x=0, \text{numa}=1, \text{sum}=0, \text{next}=\text{memory}[0]$
- Extra control states (out of the loop)
 - one to initialize next
 - one to finish off. 2 cycles after $\text{next}=0$.

5. Optimization, Architecture #4

- Datapath:



- Incremental cost:
 - Addition of another register & mux, adder mux, and control.
- Performance: find max time of the four actions

- $X \leftarrow \text{Memory}[\text{NUMA}], \text{NUMA} \leftarrow \text{NEXT} + 1;$
 $0.5 + 1 + 10 + 1 + 1 + 0.5 = 14\text{ns}$
 same for all $\Rightarrow T > 14\text{ns}, F < 71\text{MHz}$
- $\text{NEXT} \leftarrow \text{Memory}[\text{NEXT}], \text{SUM} \leftarrow \text{SUM} + X;$

Other Optimizations

- Node alignment restriction:
 - If the application of the list processor allows us to restrict the placement of nodes in memory so that they are aligned on even multiples of 2 bytes.
 - NUMA addition can be eliminated.
 - Controller supplies “0” for low-bit of memory address for NEXT, and “1” for X.
 - Furthermore, if we could use a memory with a 16-bit wide output, then could fetch entire node in one cycle:

$\{NEXT, X\} \leftarrow \text{Memory}[NEXT], \text{SUM} \leftarrow \text{SUM} + X;$

\Rightarrow execution time cut in half (half as many cycles)

List Processor Conclusions

- Through careful optimization:
 - clock frequency increased from 32MHz to 71MHz
 - little cost increase.
- “Scheduling” was used to overlap and to maximize use of resources.
- Questions:
 - Consider the design process we went through:
 - Could a computer program go from RTL description to circuits automatically?
 - Could a computer program derive the optimizations that we did?
 - It is the goal of “High-Level Synthesis” to do similar transformations and automatic mappings. “C-to-gates” compilers are an example.

Modulo Scheduling

- Review of list processor scheduling:

Memory	next ₁			x ₁		
adder		numa ₁			sum ₁	

- How did we know to “spread” out the schedule of one iteration to allow efficient packing?

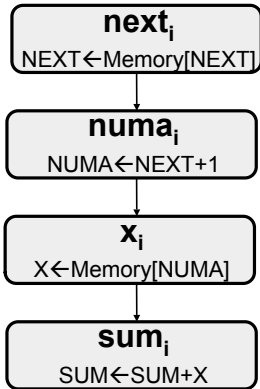
Memory	next ₁		next ₂	x ₁	next ₃	x ₂	next ₄	x ₃	
adder		numa ₁		numa ₂	sum ₁	numa ₃	sum ₂	numa ₄	sum ₃

- The goal of *modulo scheduling* is to find the schedule for one *characteristic section* of the computation. This is the part the control loops over.
- The entire schedule can then be derived, by repeating the characteristic section or repeating it with some pieces omitted.

Modulo Scheduling Procedure

1. Calculate *minimal length of characteristic section*.
 The maximum number of cycles that any one resource is used during one iteration of the computation (assuming a resource can only be used once per cycle).
2. Schedule one iteration of the computation on the characteristic section wrapping around when necessary. Each time the computation wraps around, decrease the iteration subscript by one.
3. If iteration will not fit on minimal length section, increase section by one and try again.

Modulo Scheduling List Processor



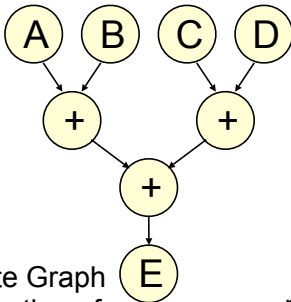
- Assuming a single adder and a single ported memory. Minimal schedule section length = 2. Because both memory and adder are used for 2 cycles during one iteration.

memory	next _i		
adder		numa _i	
memory	next _i	X _{i-1}	wrap-around, decrease subscript
adder		numa _i	
memory	next _i	X _{i-1}	wrap-around, decrease subscript
adder	sum _{i-2}	numa _i	

- Finished schedule for 4 iterations:

Memory	next ₁		next ₂	x ₁	next ₃	x ₂	next ₄	x ₃	
adder		numa ₁		numa ₂	sum ₁	numa ₃	sum ₂	numa ₄	sum ₃

Another Scheduling Example



Compute Graph
(one iteration of a
repeating calculation)

- Assume A, B, C, D, E stored in a dual port memory.
- Assume a single adder.
- Minimal schedule section length = 3. (Both memory and adder are used for 3 cycles during one iteration.)

memory port 1	load A	load C	
memory port 2	load B	load D	store E
adder	E =	A + B	C + D

Repeating schedule:

load A	load C		load A	load C		load A	load C	
load B	load D	store E	load B	load D	store E	load B	load D	store E
E =	A + B	C + D	E =	A + B	C + D	E =	A + B	C + D

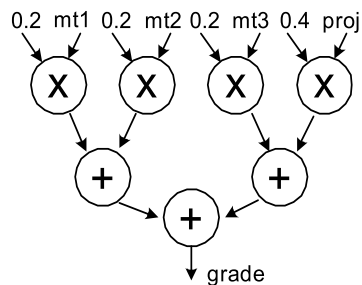
Parallelism

*Parallelism is the act of doing more than one thing at a time.
Optimization in hardware design often involves using
parallelism to trade between cost and performance.*

- Example, Student final grade calculation:

```
read mt1, mt2, mt3, project;  
grade = 0.2 × mt1 + 0.2 × mt2  
        + 0.2 × mt3 + 0.4 × project;  
write grade;
```

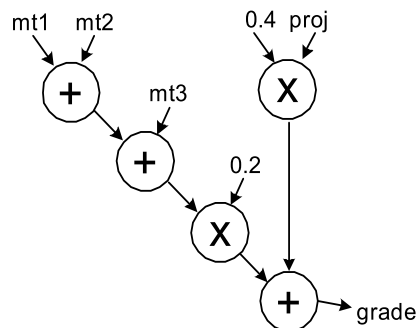
- High performance hardware implementation:



As many operations as possible are done in parallel.

Parallelism

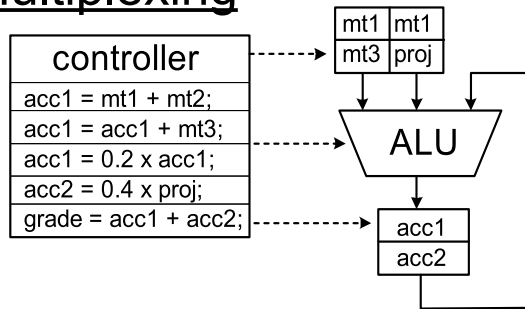
- Is there a lower cost hardware implementation? Different tree organization?
- Can factor out multiply by 0.2:



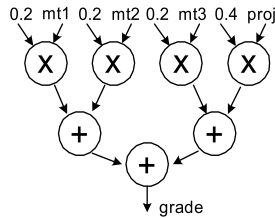
- How about sharing operators (multipliers and adders)?

Time-Multiplexing

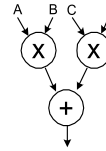
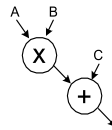
- *Time multiplex* single ALU for all adds and multiplies:
- Attempts to minimize cost at the expense of time.
 - Need to add extra register, muxes, control.



- If we adopt above approach, we can then consider the combinational hardware circuit diagram as an *abstract computation-graph*.



Using other primitives, other coverings are possible.



- This time-multiplexing “covers” the computation graph by performing the action of each node one at a time. (Sort of *emulates* it.)

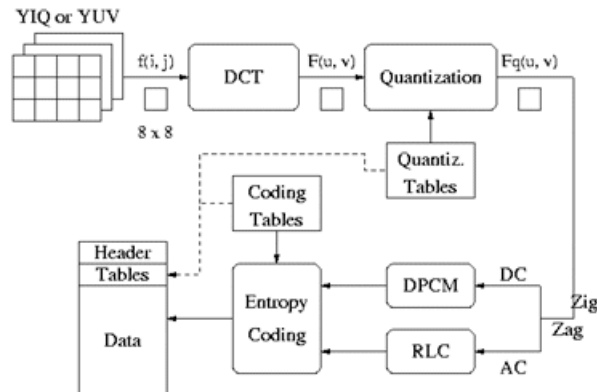
HW versus SW

- This **time-multiplexed ALU** approach is very similar to what a conventional software version would accomplish:


```
add r2, r1, r3
add r2, r2, r4
mult r2, r4, r5
      ⋮
```
- CPUs time-multiplex function units (ALUs, etc.)
- This model matches our tendency to express computation sequentially - even though most computations naturally contain parallelism.
- Our programming languages also strengthen a sequential tendency.
- In hardware we have the ability to exploit problem parallelism - gives us a “knob” to tradeoff performance & cost.
- Maybe best to express computations as abstract computations graphs (rather than “programs”) - should lead to wider range of implementations.
- *Note: modern processors spend much of their cost budget attempting to restore execution parallelism: “super-scalar execution”.*

Exploiting Parallelism in HW

- Example: Video Codec



- Separate algorithm blocks implemented in separate HW blocks, or HW is time-multiplexed.
- Entire operation is pipelined (with possible pipelining within the blocks).
- “Loop unrolling used within blocks” or for entire computation.

Spring 2009

EECS150 - Lec27-hld3

Page 15

Optimizing Iterative Computations

- Hardware implementations of computations almost always involves looping. Why?
- Is this true with software?
- Are there programs without loops?
 - Maybe in “through away” code.
- We probably would not bother building such a thing into hardware, would we?
 - (FPGA may change this.)
- Fact is, our computations are closely tied to loops. Almost all our HW includes some looping mechanism.
- What do we use looping for?

Spring 2009

EECS150 - Lec27-hld3

Page 16

Optimizing Iterative Computations

Types of loops:

1) Looping over input data (streaming):

- ex: MP3 player, video compressor, music synthesizer.

2) Looping over memory data

- ex: vector inner product, matrix multiply, list-processing

- 1) & 2) are really very similar. 1) is often turned into 2) by buffering up input data, and processing "offline". Even for "online" processing, buffers are used to smooth out temporary rate mismatches.

3) CPUs are one big loop.

- Instruction fetch \Rightarrow execute \Rightarrow Instruction fetch \Rightarrow execute \Rightarrow ...
- but change their personality with each iteration.

4) Others?

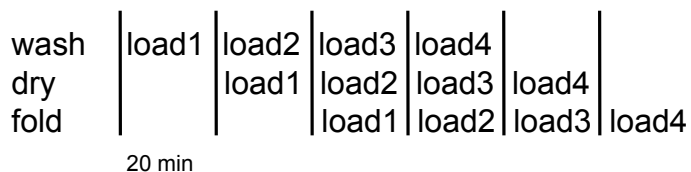
*Loops offer opportunity for parallelism
by executing more than one iteration at once,
using [parallel iteration execution &/or pipelining](#)*

Pipelining Principle

- With looping usually we are less interested in the latency of one iteration and more in the loop execution rate, or throughput.
- These can be different due to [parallel iteration execution &/or pipelining](#).
- Pipelining review from CS61C:

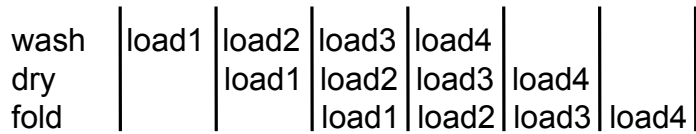
Analog to washing clothes:

step 1: wash (20 minutes)
 step 2: dry (20 minutes)
 step 3: fold (20 minutes)
 60 minutes x 4 loads \Rightarrow 4 hours



overlapped \Rightarrow 2 hours

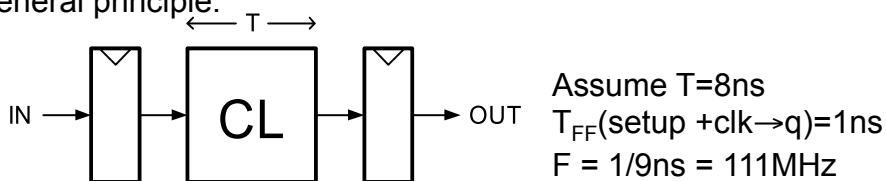
Pipelining



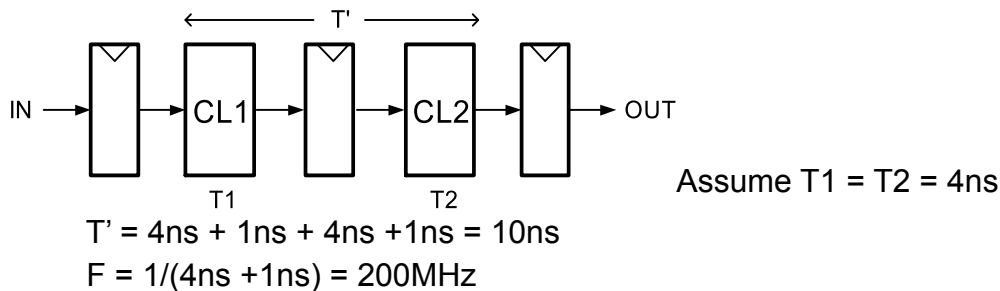
- In the limit, as we increase the number of loads, the average time per load approaches 20 minutes.
- The latency (time from start to end) for one load = 60 min.
- The throughput = 3 loads/hour
- The pipelined throughput \approx # of pipe stages x un-pipelined throughput.

Pipelining

- General principle:



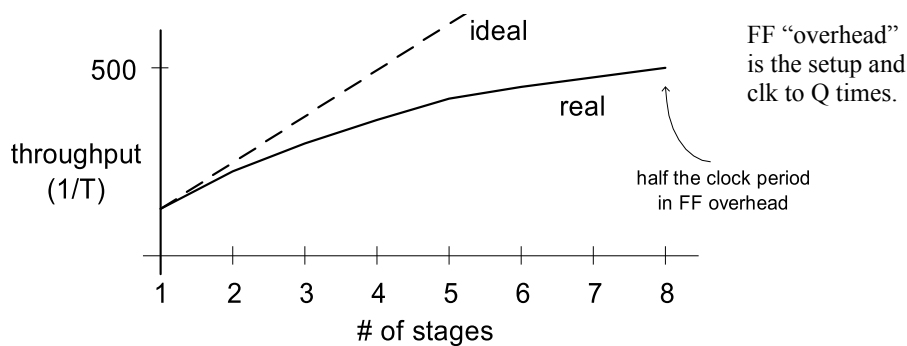
- Cut the CL block into pieces (stages) and separate with registers:



- CL block produces a new result every 5ns instead of every 9ns.

Limits on Pipelining

- Without FF overhead, throughput improvement \propto # of stages.
- After many stages are added FF overhead begins to dominate:



- Other limiters to effective pipelining:
 - clock skew contributes to clock overhead
 - unequal stages
 - FFs dominate *cost*
 - clock distribution power consumption
 - [feedback \(dependencies between loop iterations\)](#)

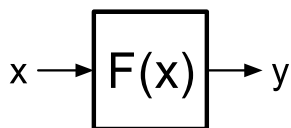
Spring 2009

EECS150 - Lec27-hld3

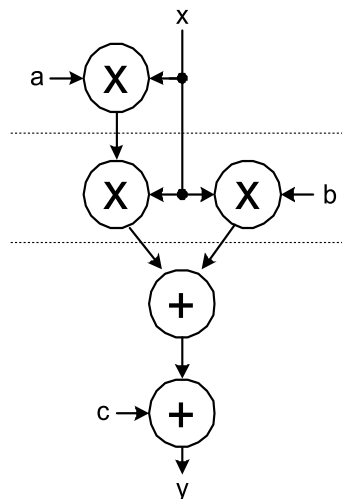
Page 21

Pipelining Example

- $F(x) = y_i = a x_i^2 + b x_i + c$



- Computation graph:



- x and y are assumed to be "streams"
- Divide into 3 (nearly) equal stages.
- Insert pipeline registers at dashed lines.
- Can we pipeline basic operators?

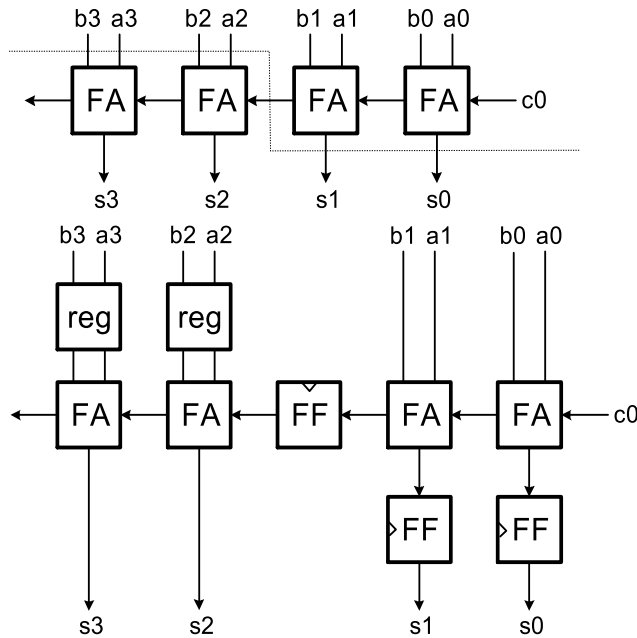
Spring 2009

EECS150 - Lec27-hld3

Page 22

Example: Pipelined Adder

- Possible, but usually not done.
(arithmetic units can often be made sufficiently fast without internal pipelining)



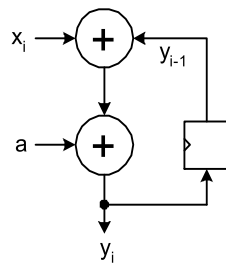
Pipelining Loops with Feedback

“Loop carry dependency”

- Example 1: $y_i = y_{i-1} + x_i + a$

unpipelined version:

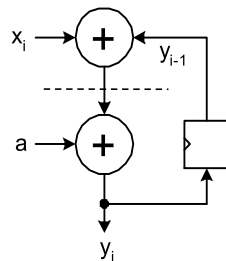
add ₁	$x_i + y_{i-1}$	$x_{i+1} + y_i$
add ₂	y_i	y_{i+1}



Can we “cut” the feedback and overlap iterations?

Try putting a register after add1:

add ₁	$x_i + y_{i-1}$	$x_{i+1} + y_i$
add ₂	y_i	y_{i+1}



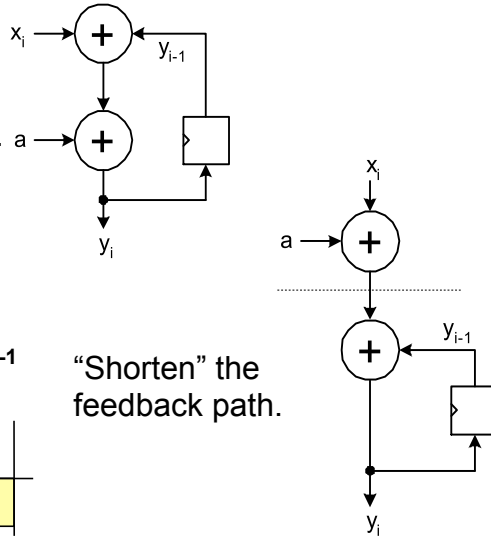
- Can't overlap the iterations because of the dependency.
- The extra register doesn't help the situation (actually hurts).
- In general, can't pipeline feedback loops.

Pipelining Loops with Feedback

“Loop carry dependency”

However, we can overlap the “non-feedback” part of the iterations:

Add is associative and commutative. Therefore we can reorder the computation to shorten the delay of the feedback path:



$$y_i = (y_{i-1} + x_i) + a = (a + x_i) + y_{i-1}$$

“Shorten” the feedback path.

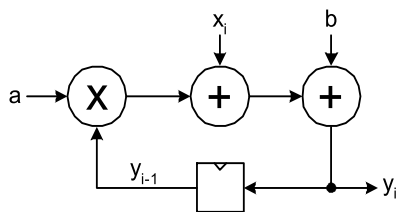
add ₁	x _i +a	x _{i+1} +a	x _{i+2} +a
add ₂		y _i	y _{i+1}

- Pipelining is limited to 2 stages.

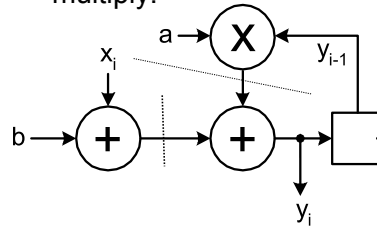
Pipelining Loops with Feedback

- Example 2:

$$y_i = a y_{i-1} + x_i + b$$



- Reorder to shorten the feedback loop and try putting register after multiply:



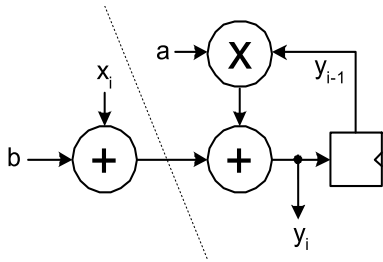
add ₁	x _i +b		x _{i+1} +b		x _{i+2} +b
mult	ay _{i-1}		ay _i		ay _{i+1}
add ₂		y _i		y _{i+1}	

- Still need 2 cycles/iteration

Pipelining Loops with Feedback

- Example 2:

$$y_i = a y_{i-1} + x_i + b$$



- Once again, adding register doesn't help. Best solution is to overlap non-feedback part with feedback part.
- Therefore critical path includes a multiply in series with add.
- Can overlap first add with multiply/add operation.
- Only 1 cycle/iteration. Higher performance solution (than 2 cycle version).

add ₁	x _i +b	x _{i+1} +b	x _{i+2} +b	
mult		ay _{i-1}	ay _i	ay _{i+1}
add ₂		y _i	y _{i+1}	y _{i+2}

- Alternative is to move register to after multiple, but same critical path.

“C-slow” Technique

- Another approach to increasing throughput in the presence of feedback: try to fill in “holes” in the chart with another (independent) computation:

add ₁	x _i +b		x _{i+1} +b		x _{i+2} +b	
mult	ay _{i-1}		ay _i		ay _{i+1}	
add ₂		y _i		y _{i+1}		y _{i+2}

If we have a second similar computation, can interleave it with the first:

$$x^1 \rightarrow \boxed{F^1} \rightarrow y^1 = a^1 y_{i-1}^1 + x_i^1 + b^1$$

Use muxes to direct each stream.
Time multiplex one piece of HW for both stream.

$$x^2 \rightarrow \boxed{F^2} \rightarrow y^2 = a^2 y_{i-1}^2 + x_i^2 + b^2$$

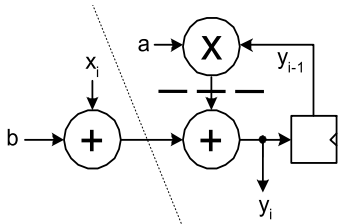
Each produces 1 result / 2 cycles.

- Here the feedback depth=2 cycles (we say C=2).
- Each loop has throughput of F_{clk}/C . But the aggregate throughput is F_{clk} .
- With this technique we could pipeline even deeper, assuming we could supply C independent streams.

“C-slow” Technique

- Essentially this means we go ahead and cut feedback path:
 - C computations (in this case C=2) can use the pipeline simultaneously.
 - Must be independent.
 - Input MUX interleaves input streams.
 - Each stream runs at half the pipeline frequency.
 - Pipeline achieves full throughput.
- This makes operations in adjacent pipeline stages independent and allows full cycle for each:

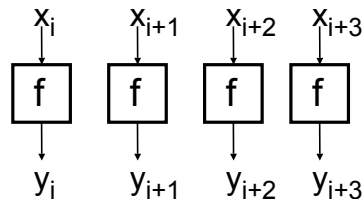
Multithreaded Processors use this.



add ₁	x+b	x+b	x+b	x+b	x+b	x+b
mult	ay	ay	ay	ay	ay	ay
add ₂	y	y	y	y	y	y

Beyond Pipelining - SIMD Parallelism

- An obvious way to exploit more parallelism from loops is to make multiple instances of the loop execution data-path and run them in parallel, sharing the same controller.
- For P instances, throughput improves by a factor of P.
- example: $y_i = f(x_i)$



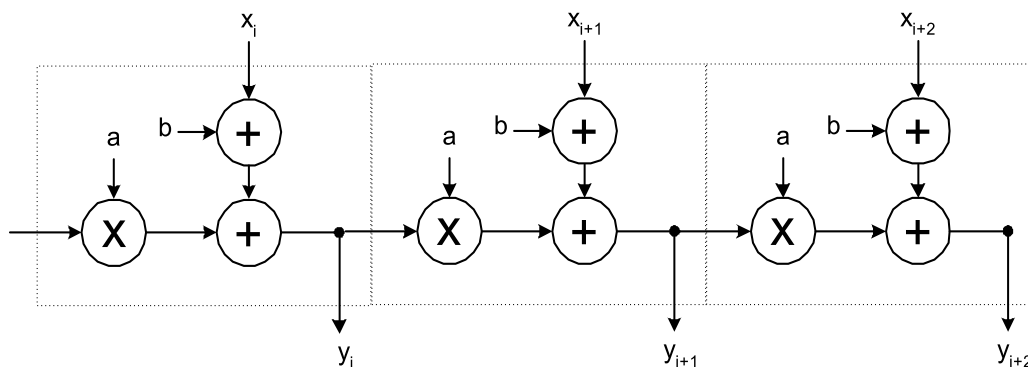
Usually called SIMD parallelism. Single Instruction Multiple Data

- Assumes the next 4 x values available at once. The validity of this assumption depends on the ratio of f repeat rate to input rate (or memory bandwidth).
- Cost $\propto P$. Usually, much higher than for pipelining. However, potentially provides a high speedup. Often applied after pipelining.
- Limited, once again, by loop carry dependencies. Feedback translates to dependencies between parallel data-paths.
- **Vector processors use this technique.**

SIMD Parallelism with Feedback

- Example, from earlier:

$$y_i = a y_{i-1} + x_i + b$$



- In this example end up with “carry ripple” situation.
- Could employ look-ahead / parallel-prefix optimization techniques to speed up propagation.
- As with pipelining, this technique is most effective in the absence of a loop carry dependence.