

Pipelining

what Seymour Cray taught the laundry industry

Quiz 2 (next week) will cover materials through Tuesday's lecture. Lab 3 is due tonight.

Forget 6.004... lets solve a "Real Problem"

INPUT:
dirty laundry



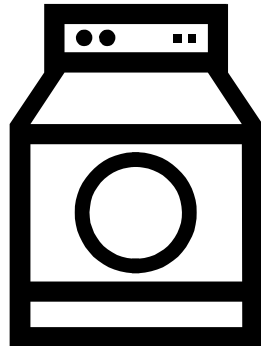
OUTPUT:
6 more weeks



Device: Washer

Function: Fill, Agitate, Spin

Washer_{PD} = 30 mins



Device: Dryer

Function: Heat, Spin

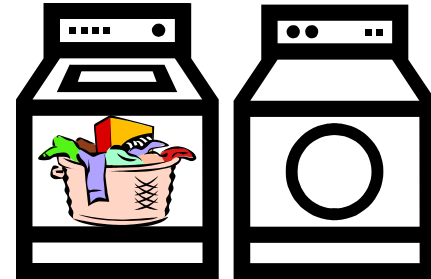
Dryer_{PD} = 60 mins

One load at a time

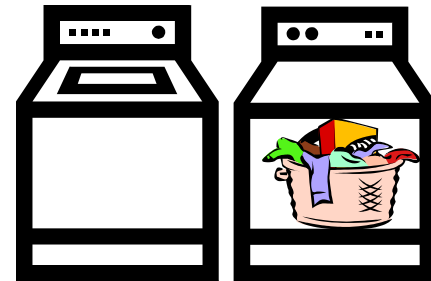
Everyone knows that the real reason that MIT students put off doing laundry so long is not because they procrastinate, are lazy, or even have better things to do.

The fact is, doing one load at a time is not smart.

Step 1:



Step 2:

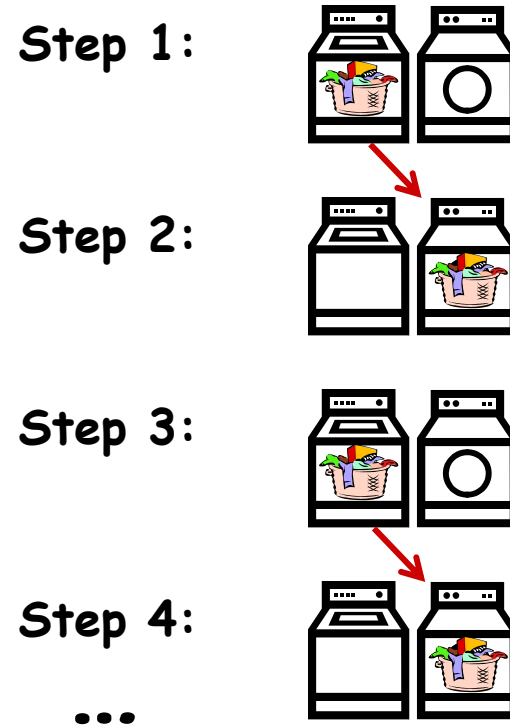


$$\begin{aligned} \text{Total} &= \text{Washer}_{PD} + \text{Dryer}_{PD} \\ &= \underline{\quad 90 \quad} \text{ mins} \end{aligned}$$

Doing N loads of laundry

Here's how they do laundry at Harvard, the "combinational" way.

(Of course, this is just an urban legend. No one at Harvard actually *does* laundry. The butlers all arrive on Wednesday morning, pick up the dirty laundry and return it all pressed and starched in time for afternoon tea)



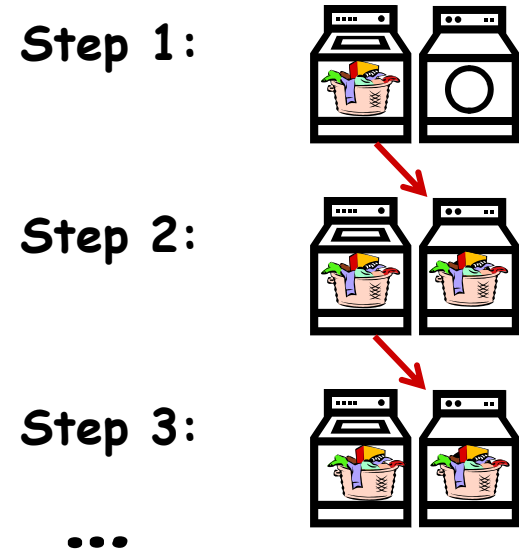
$$\begin{aligned} \text{Total} &= N * (\text{Washer}_{PD} + \text{Dryer}_{PD}) \\ &= \underline{\quad N * 90 \quad} \text{ mins} \end{aligned}$$

Doing N Loads... the MIT way

MIT students “pipeline”
the laundry process.

That’s why we wait!

Actually, it’s more like $N \cdot 60 + 30$ if we account for the startup transient correctly. When doing pipeline analysis, we’re mostly interested in the “steady state” where we assume we have an infinite supply of inputs.



$$\begin{aligned} \text{Total} &= N * \text{Max}(\text{Washer}_{PD}, \text{Dryer}_{PD}) \\ &= \underline{N \cdot 60} \text{ mins} \end{aligned}$$

Some definitions

Latency:

The delay from when an input is established until the output associated with that input becomes valid.

(Harvard Laundry = _____ mins)

(MIT Laundry = _____ mins)
120

Assuming that the wash is started as soon as possible and waits (wet) in the washer until dryer is available.

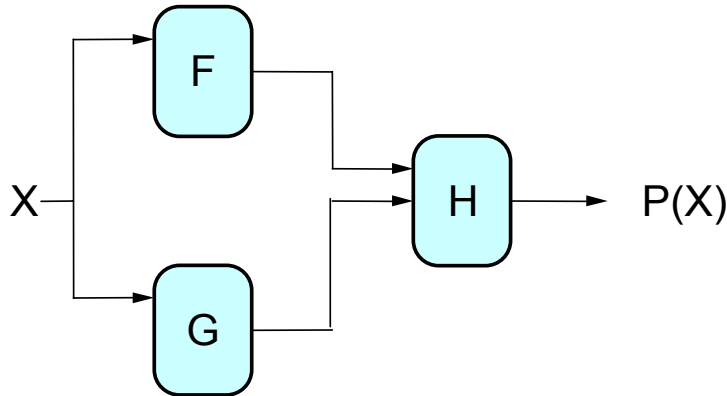
Throughput:

The rate of which inputs or outputs are processed.

(Harvard Laundry = _____ outputs/min)

(MIT Laundry = _____ outputs/min)
1/90
1/60

Okay, back to circuits...

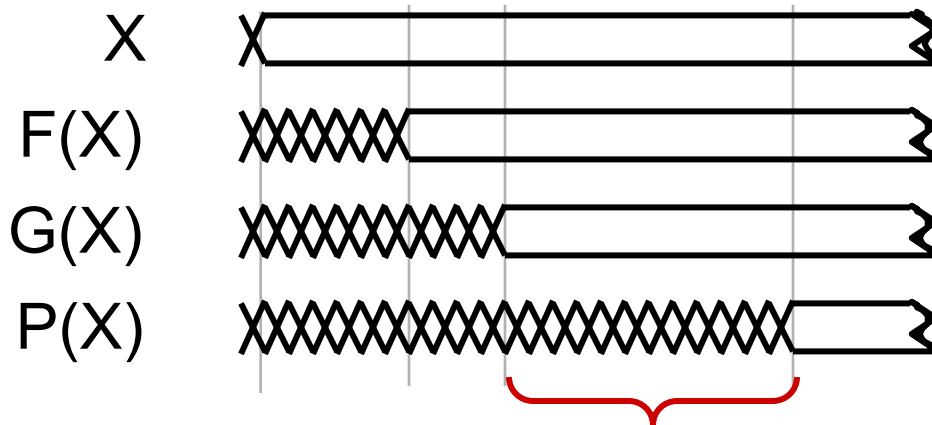


For combinational logic:

$$\text{latency} = t_{PD},$$

$$\text{throughput} = 1/t_{PD}.$$

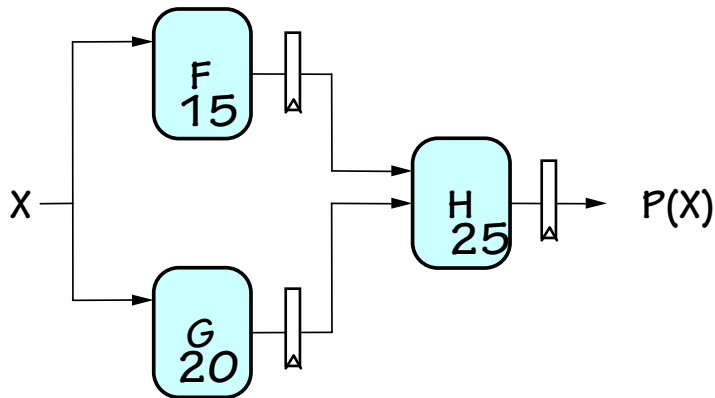
We can't get the answer faster, but are we making effective use of our hardware at all times?



F & G are "idle", just holding their outputs stable while H performs its computation

Pipelined Circuits

use registers to hold H's input stable!

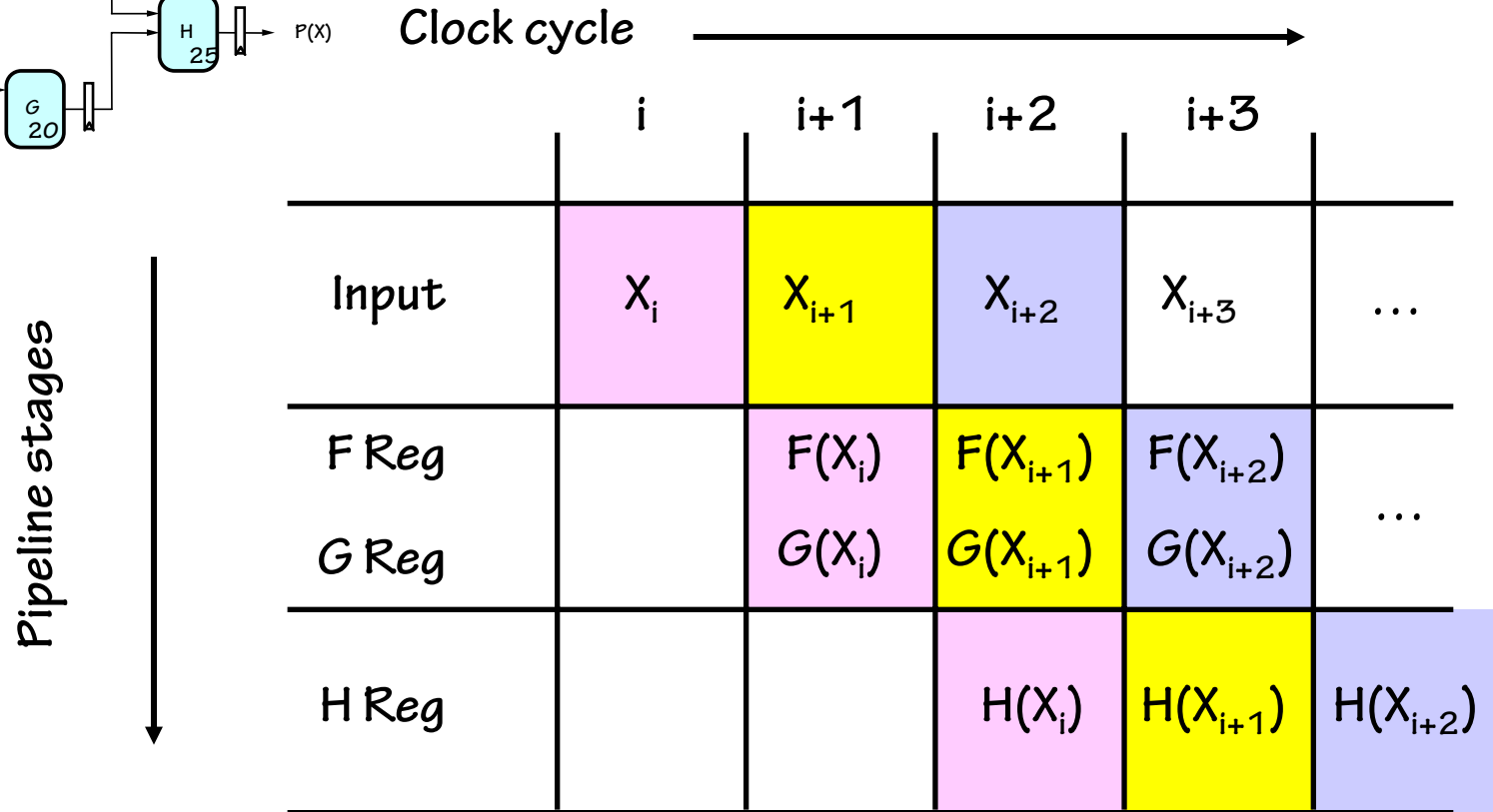
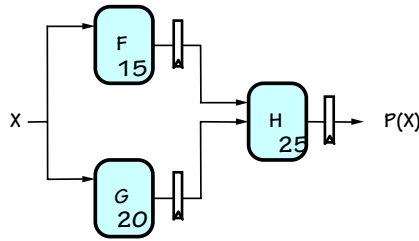


Now F & G can be working on input X_{i+1} while H is performing its computation on X_i . We've created a 2-stage *pipeline*: if we have a valid input X during clock cycle j, P(X) is valid during clock j+2.

Suppose F, G, H have propagation delays of 15, 20, 25 ns and we are using ideal zero-delay registers:

	<u>latency</u>	<u>throughput</u>
unpipelined	45	1/45
2-stage pipeline	<u>50</u>	<u>1/25</u>
	worse	better

Pipeline diagrams



The results associated with a particular set of input data moves *diagonally* through the diagram, progressing through one pipeline stage each clock cycle.

Pipeline Conventions

DEFINITION:

a *K-Stage Pipeline* (“K-pipeline”) is an acyclic circuit having exactly K registers on every path from an input to an output.

a COMBINATIONAL CIRCUIT is thus an 0-stage pipeline.

CONVENTION:

Every pipeline stage, hence every K-Stage pipeline, has a register on its *OUTPUT* (not on its input).

ALWAYS:

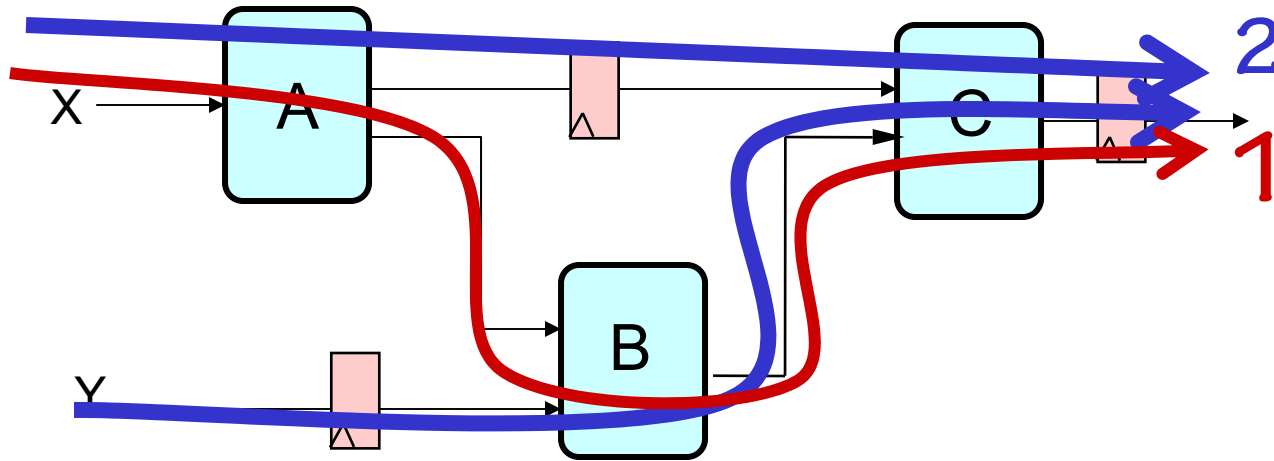
The CLOCK common to all registers must have a period sufficient to cover propagation over combinational paths PLUS (input) register t_{PD} PLUS (output) register t_{SETUP} .

The LATENCY of a K-pipeline is K times the period of the clock common to all registers.

The THROUGHPUT of a K-pipeline is the frequency of the clock.

Ill-formed pipelines

Consider a BAD job of pipelining:



For what value of K is the following circuit a K-Pipeline? ANS: _____

none

Problem:

Successive inputs get mixed: e.g., $B(A(X_{i+1}), Y_i)$. This happened because some paths from inputs to outputs had 2 registers, and some had only 1!

Can this happen on a well-formed K pipeline?

A pipelining methodology

Step 1:

Draw a line that crosses every output in the circuit, and mark the endpoints as terminal points.

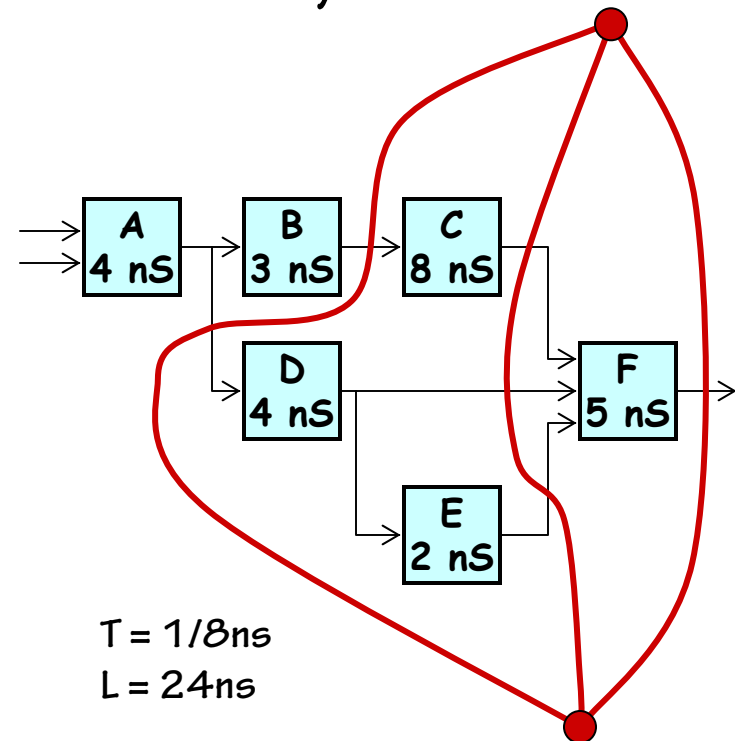
Step 2:

Continue to draw new lines between the terminal points across various circuit connections, ensuring that every connection crosses each line in the same direction. These lines demarcate *pipeline stages*.

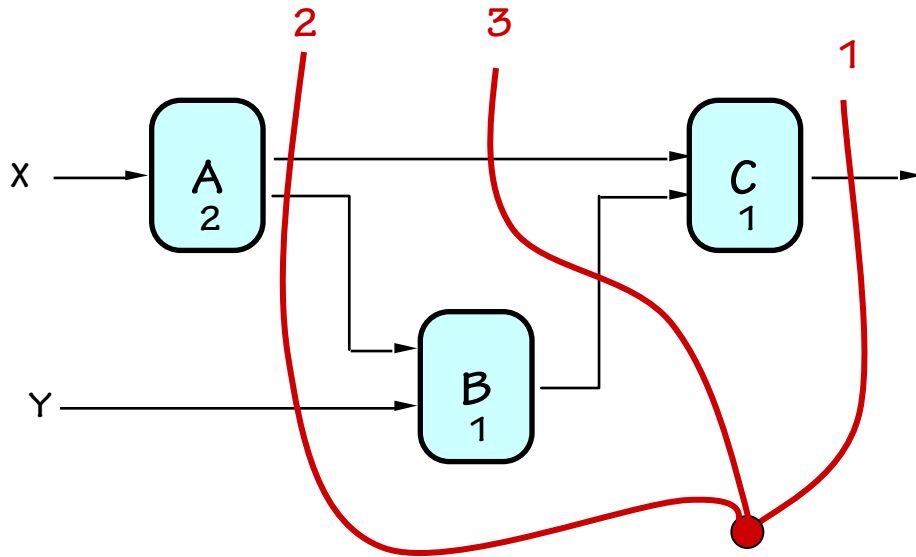
Adding a pipeline register at every point where a separating line crosses a connection will always generate a valid pipeline.

STRATEGY:

Focus your attention on placing pipelining registers around the slowest circuit elements (BOTTLENECKS).



Pipeline Example



OBSERVATIONS:

- 1-pipeline improves neither L or T.
- T improved by breaking long combinational paths, allowing faster clock.
- Too many stages cost L, don't improve T.
- Back-to-back registers are often required to keep pipeline well-formed.

	LATENCY	THROUGHPUT
0-pipe:	4	1/4
1-pipe:	4	1/4
2-pipe:	4	1/2
3-pipe:	6	1/2

Pipelining Summary

Advantages:

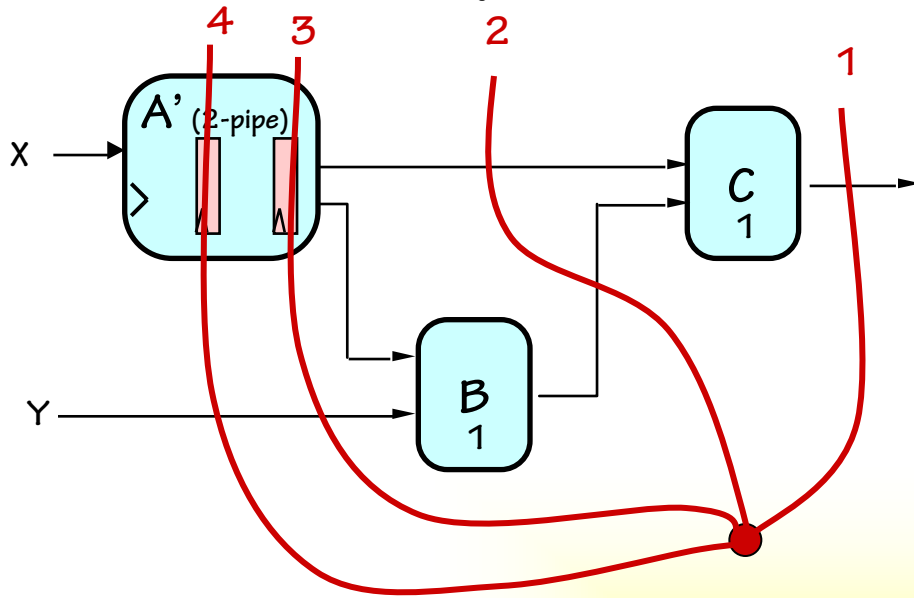
- Allows us to increase throughput, by breaking up long combinational paths and (hence) increasing clock frequency

Disadvantages:

- May increase latency...
- Only as good as the weakest link: slowest step constrains system throughput.

Isn't there a way around this "weak link" problem?

Pipelined Components



4-stage pipeline, throughput=1

but... but...

How can one pipeline
a clothes dryer???

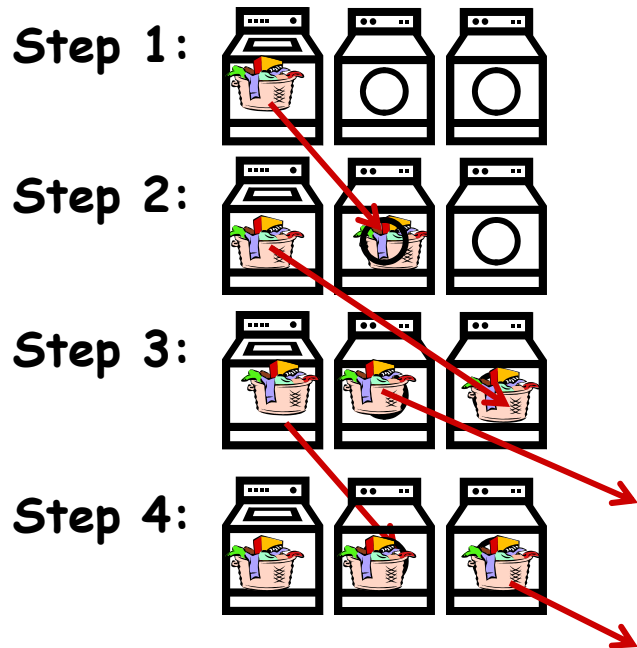
Pipelined systems can be hierarchical:

- Replacing a slow combinational component with a k-pipe version may increase clock frequency
- Must account for new pipeline stages in our plan

How do 6.004 Aces do Laundry?

They work around the bottleneck.

First, they find a place with twice as many dryers as washers.



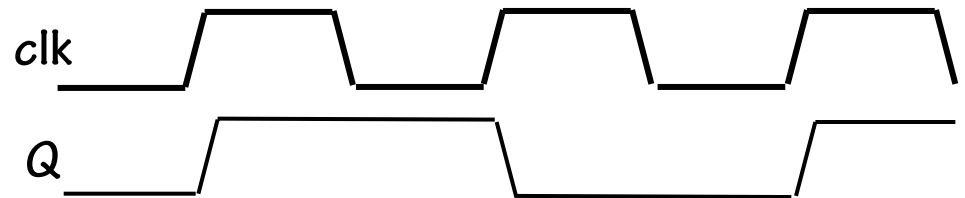
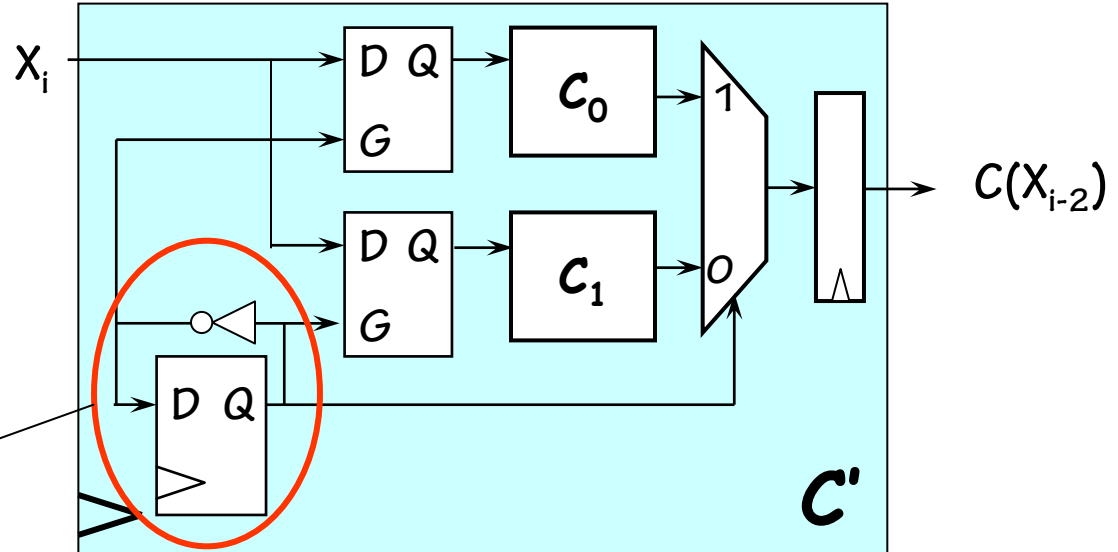
$$\text{Throughput} = \frac{\quad}{1/30} \text{ loads/min}$$

$$\text{Latency} = \frac{\quad}{90} \text{ mins/load}$$

Circuit Interleaving

We can simulate a pipelined version of a slow component by replicating the critical element and alternate inputs between the various copies.

This is a simple 2-state FSM that alternates between 0 and 1 on each clock

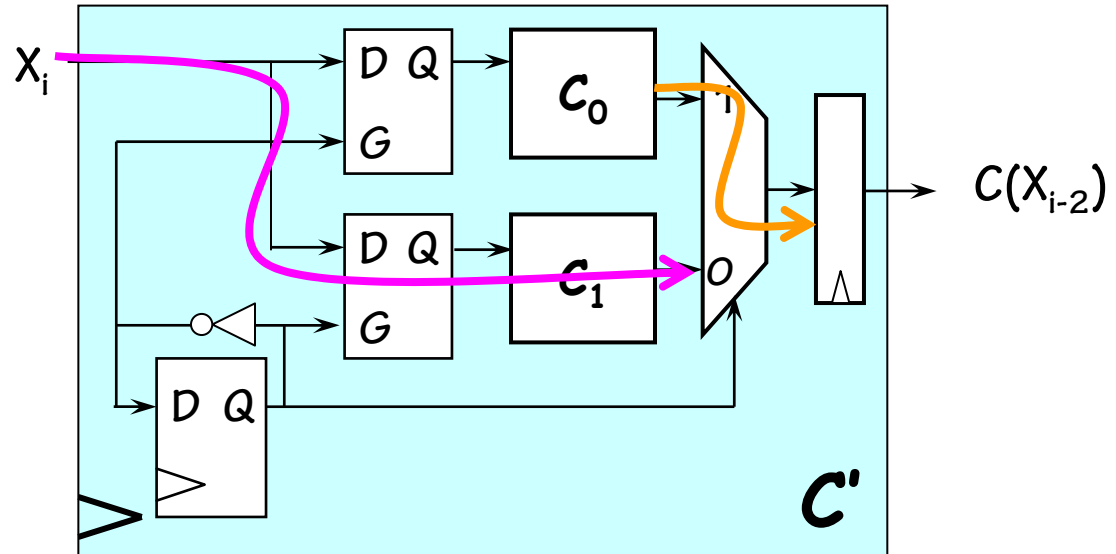


Circuit Interleaving

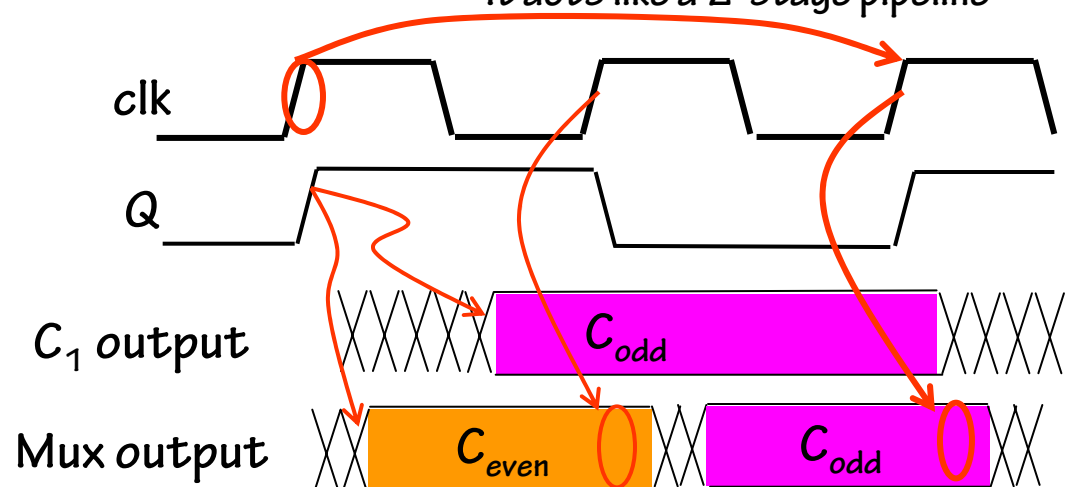
We can simulate a pipelined version of a slow component by replicating the critical element and alternate inputs between the various copies.

When Q is 1 the lower path is combinational (the latch is open), yet the output of the upper path will be enabled onto the input of the output register ready for the NEXT clock edge.

Meanwhile, the other latch maintains the input from the last clock.



"It acts like a 2-stage pipeline"

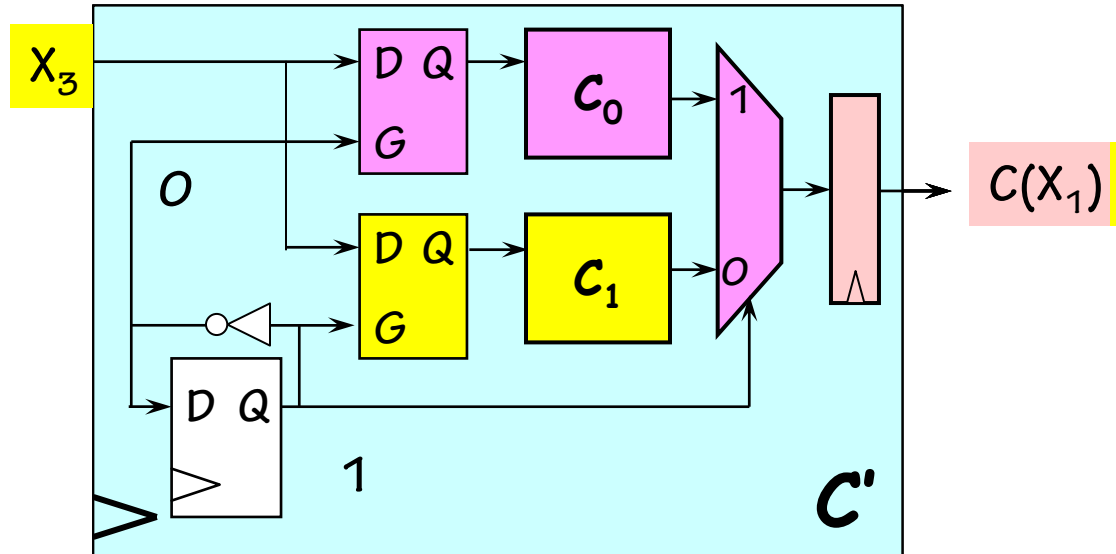


Circuit Interleaving

2-Clock Martinizing

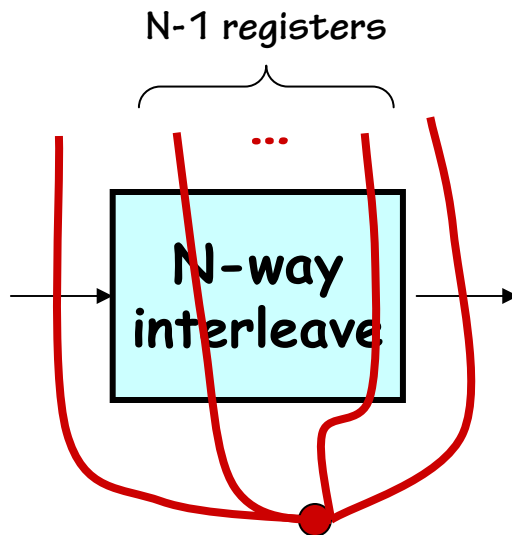
“In by t_i , out by t_{i+2} ”

N-way interleaving
is equivalent to
N pipeline Stages...



Latency = 2 clocks

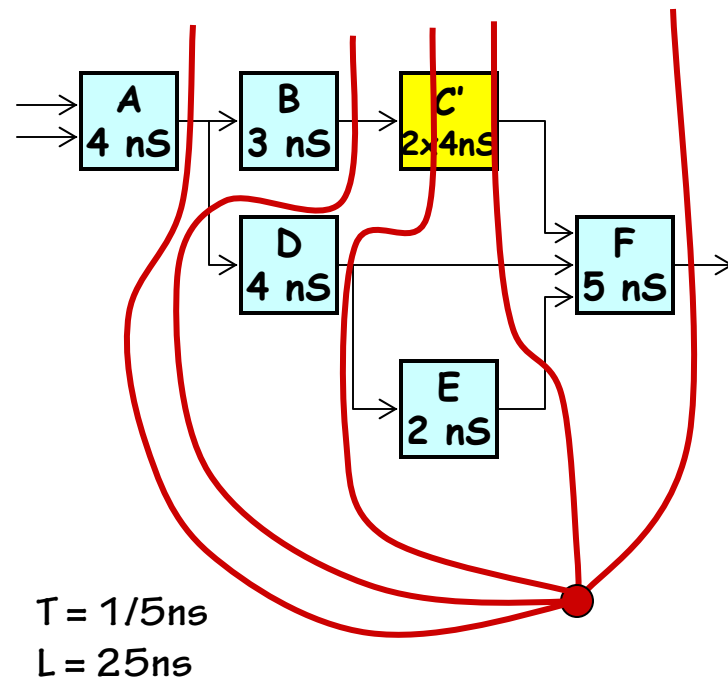
- Clock period 0: X_0 presented at input, propagates thru upper latch, C_0 .
- Clock period 1: X_1 presented at input, propagates thru lower latch, C_1 . $C_0(X_0)$ propagates to register inputs.
- Clock period 2: X_2 presented at input, propagates thru upper latch, C . $C_0(X_0)$ loaded into register, appears at output.



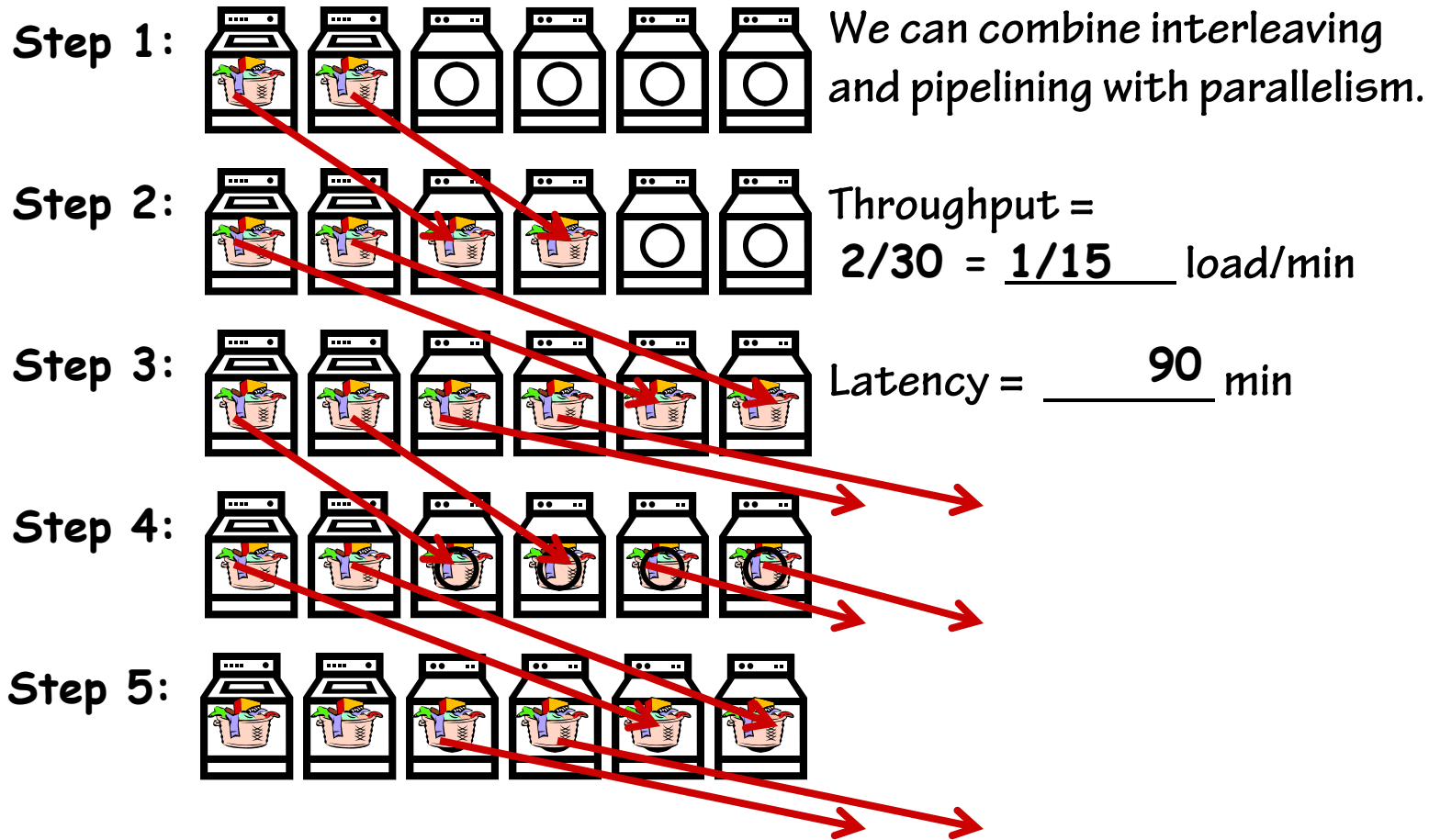
Combining techniques

We can combine interleaving and pipelining. Here, C' interleaves two C elements with a propagation delay of 8 nS . The resulting C' circuit has a throughput of $1/4\text{ nS}$, and latency of 8 nS . This can be considered as an extra pipelining stage that passes through the middle of the C' module. One of our separation lines must pass through this pipeline stage.

By combining interleaving with pipelining we move the bottleneck from the C element to the F element.



And a little parallelism...



Control Structure Approaches

Synchronous

ALL computation “events” occur at active edges of a periodic clock: time is divided into fixed-size discrete intervals.

RIGID

Globally Timed

Timing dictated by centralized FSM according to a fixed schedule.

Asynchronous

Events -- eg the loading of a register -- can happen at arbitrary times.

Laid Back

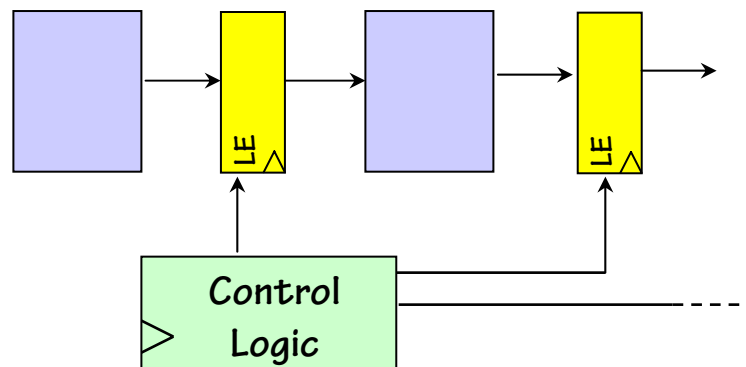
Locally Timed

Each module takes a START signal, generates a FINISHED signal. Timing is dynamic, data dependent.



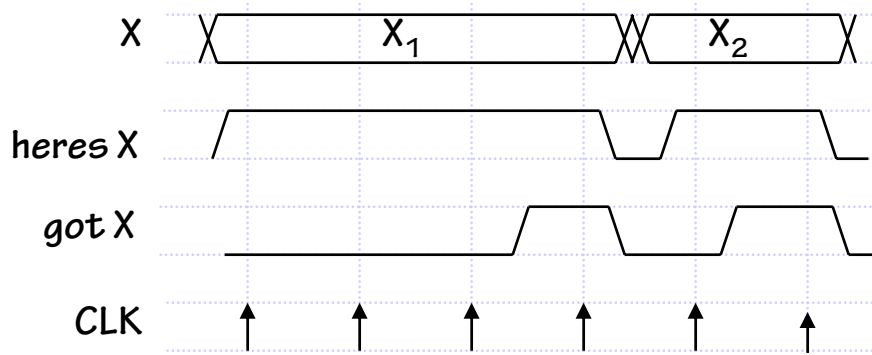
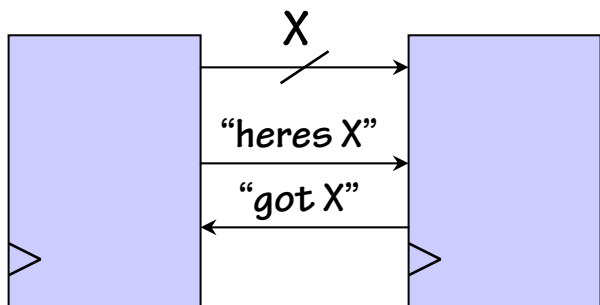
Control Structure Alternatives

Synchronous, globally-timed:
Control signals (e.g., load enables)
From FSM controller

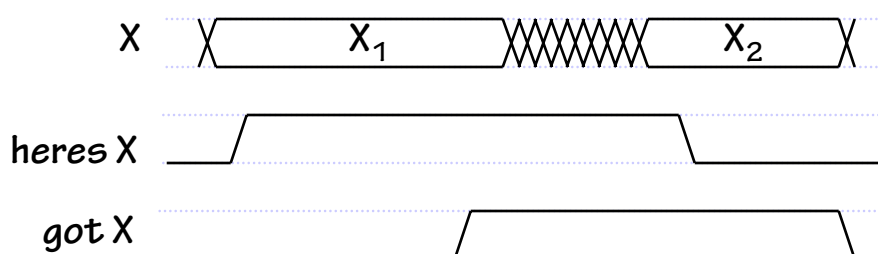
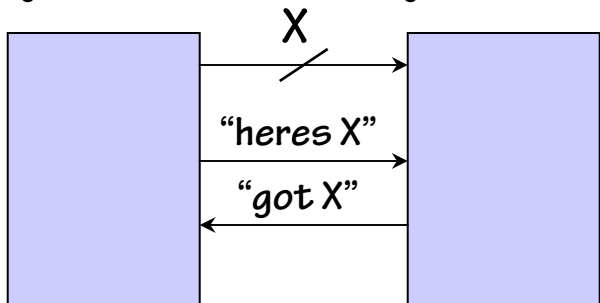


Synchronous, locally-timed:
Local circuitry, "handshake" controls

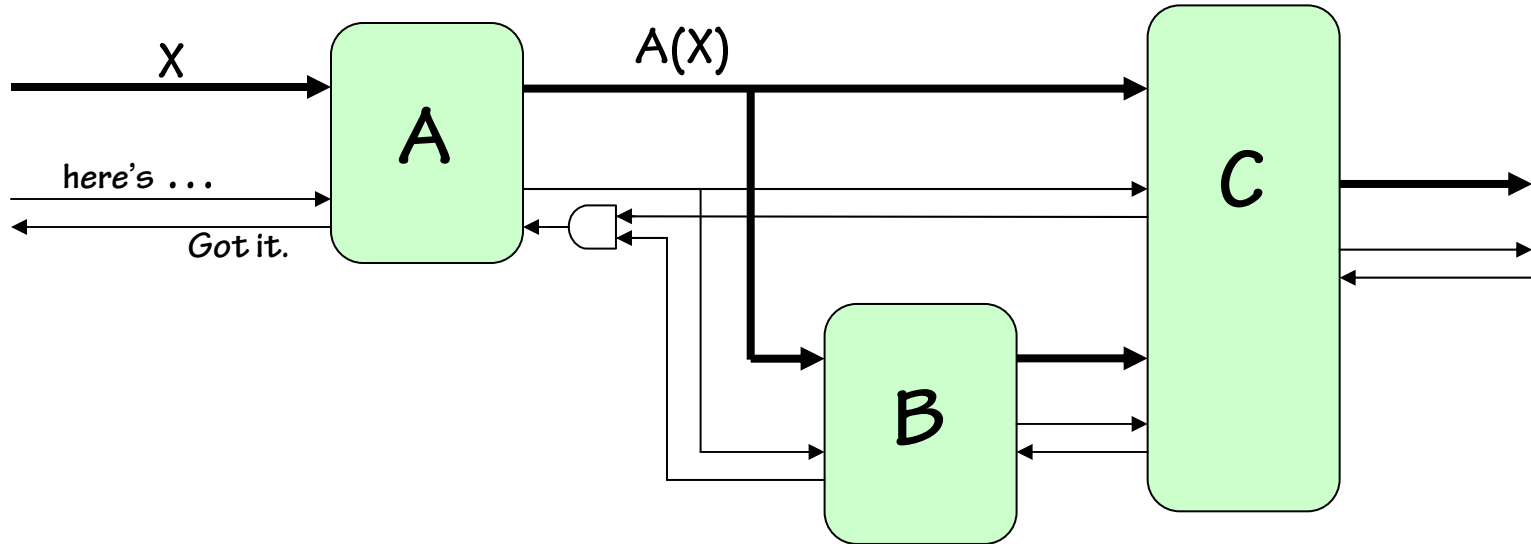
flow of data:



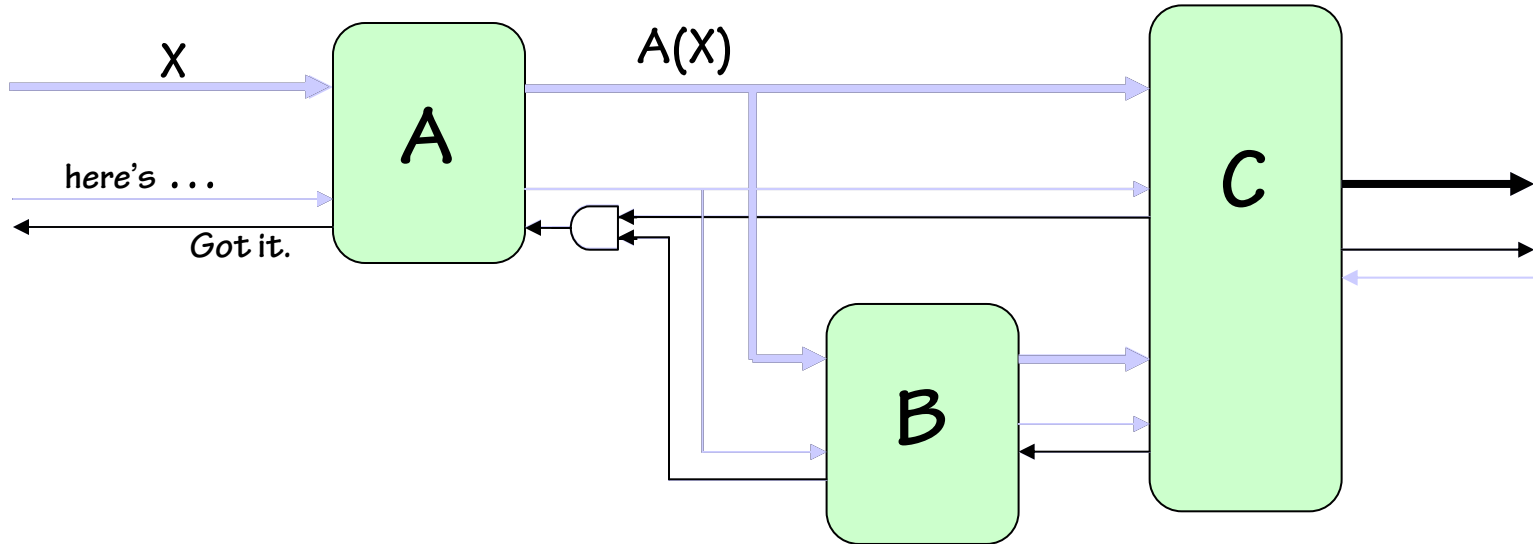
Asynchronous, locally-timed system using *transition signaling*:



Self-timed Example



Self-timed Example



Elegant, timing-independent design:

- Each component specifies its own time constraints
- Local adaptation to special cases (eg, multiplication by 0)
- Module performance improvements automatically exploited
- Can be made asynchronous (no clock at all!) or synchronous

Control Structure Taxonomy

Easy to design but fixed-sized interval can be wasteful (no data-dependencies in timing)

Large systems lead to very complicated timing generators... just say no!

Globally Timed

Synchronous

Centralized clocked FSM generates all control signals.

Asynchronous

Central control unit tailors current time slice to current tasks.

Locally Timed

Start and Finish signals generated by each major subsystem, synchronously with global clock.

Each subsystem takes asynchronous Start, generates asynchronous Finish (perhaps using local clock).

The best way to build large systems that have independently-timed components.

The "next big idea" for the last several decades: a lot of design work to do in general, but extra work is worth it in special cases

Summary

- Latency (L) = time it takes for given input to arrive at output
- Throughput (T) = rate at each new outputs appear
- For combinational circuits: $L = t_{PD}$ of circuit, $T = 1/L$
- For K -pipelines ($K > 0$):
 - always have register on output(s)
 - K registers on every path from input to output
 - Inputs available shortly after clock i , outputs available shortly after clock $(i+K)$
 - $T = 1/(t_{PD,REG} + t_{PD}$ of slowest pipeline stage $+ t_{SETUP})$
 - more throughput \rightarrow split slowest pipeline stage(s)
 - use replication/interleaving if no further splits possible
 - $L = K / T$
 - pipelined latency \geq combinational latency