# Parallelism

**Parallelism:** execute **multiple operations simultaneously**

**Some Types**

--- **Instruction-Level Parallelism ===>** execute **multiple instructions** from **same job** **(ILP)**

--- **Data Parallelism ===>** operate on **multiple data** items from **same job (SIMD, MIMD, SPMD)**

--- **Thread-Level Parallelism ===>** execute **multiple jobs** but **same program**

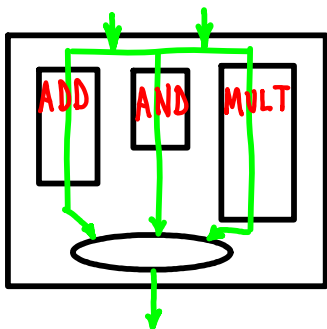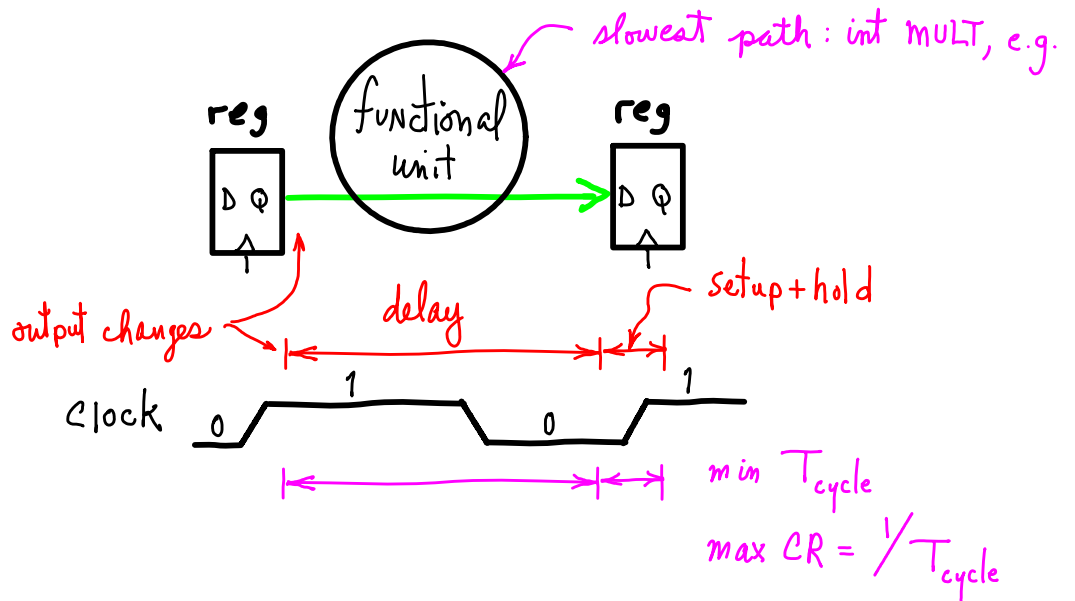--- **Task-Level Parallelism ===>** execute **multiple jobs, multiple programs**

# ILP, Pipeline

$$P_{erf} = \frac{n \; instructions}{Time} = \frac{n}{n \; CPI \; (1/CR)} = CR/CPI \qquad CR\uparrow \implies perf \uparrow$$

*slowest path : int MULT, e.g.*

**CR limited by slowest path.**

--- **Other, faster units** could be **clocked faster**.

--- **But cannot** because of slow path.

--- Faster units **waiting/idle**.

reg  **functional unit**  reg

D Q  →  D Q

*output changes*   *delay*   *Setup + hold*

Clock  0  1  0  1

$$min \; T_{cycle}$$
$$max \; CR = 1/T_{cycle}$$

ALU

ADD  AND  MULT

— Only 1 unit used
— every instruction takes same time
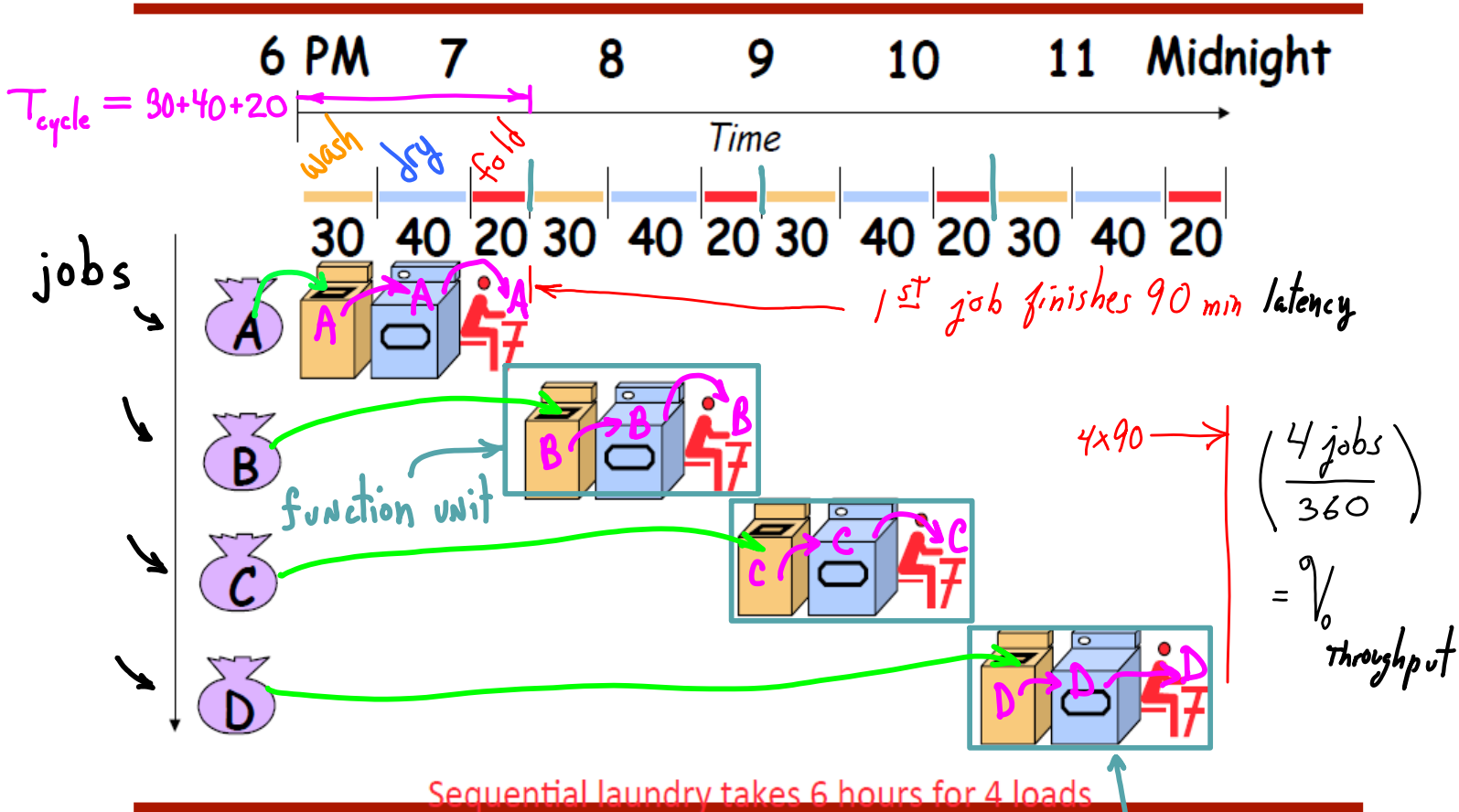
Amdahl's Law says **attack the common case.**
**Every instruction faster** w/ increased **CR**.
--- **MULT** chopped into **k fast pieces**.
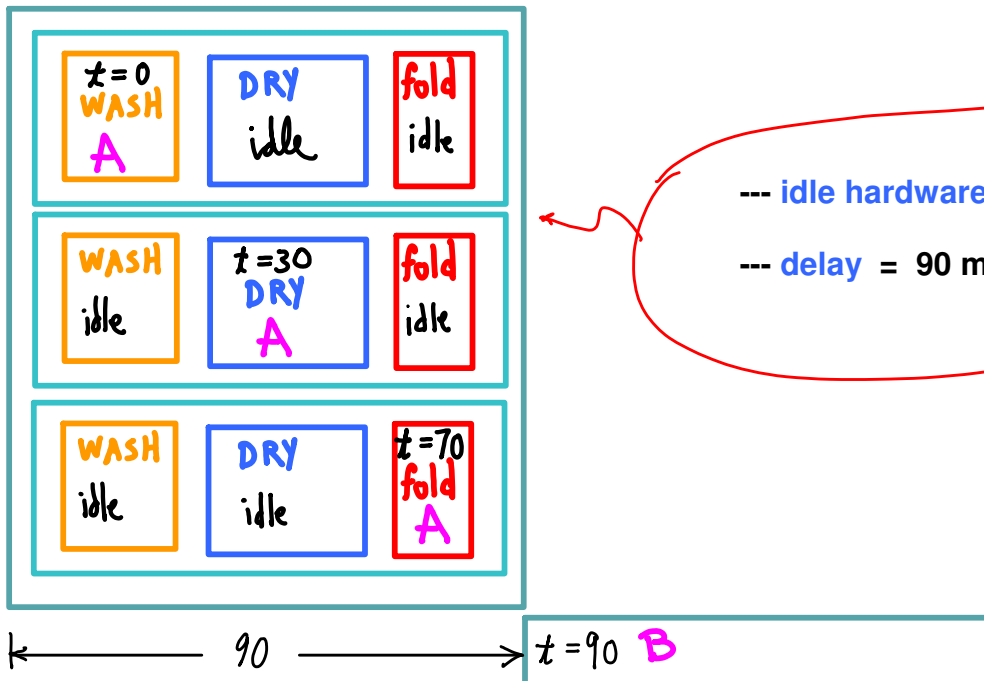--- **More cycles** for MULT, **but** faster **CR**.

*also, see if we can get k pieces of MULT to run in Parallel*
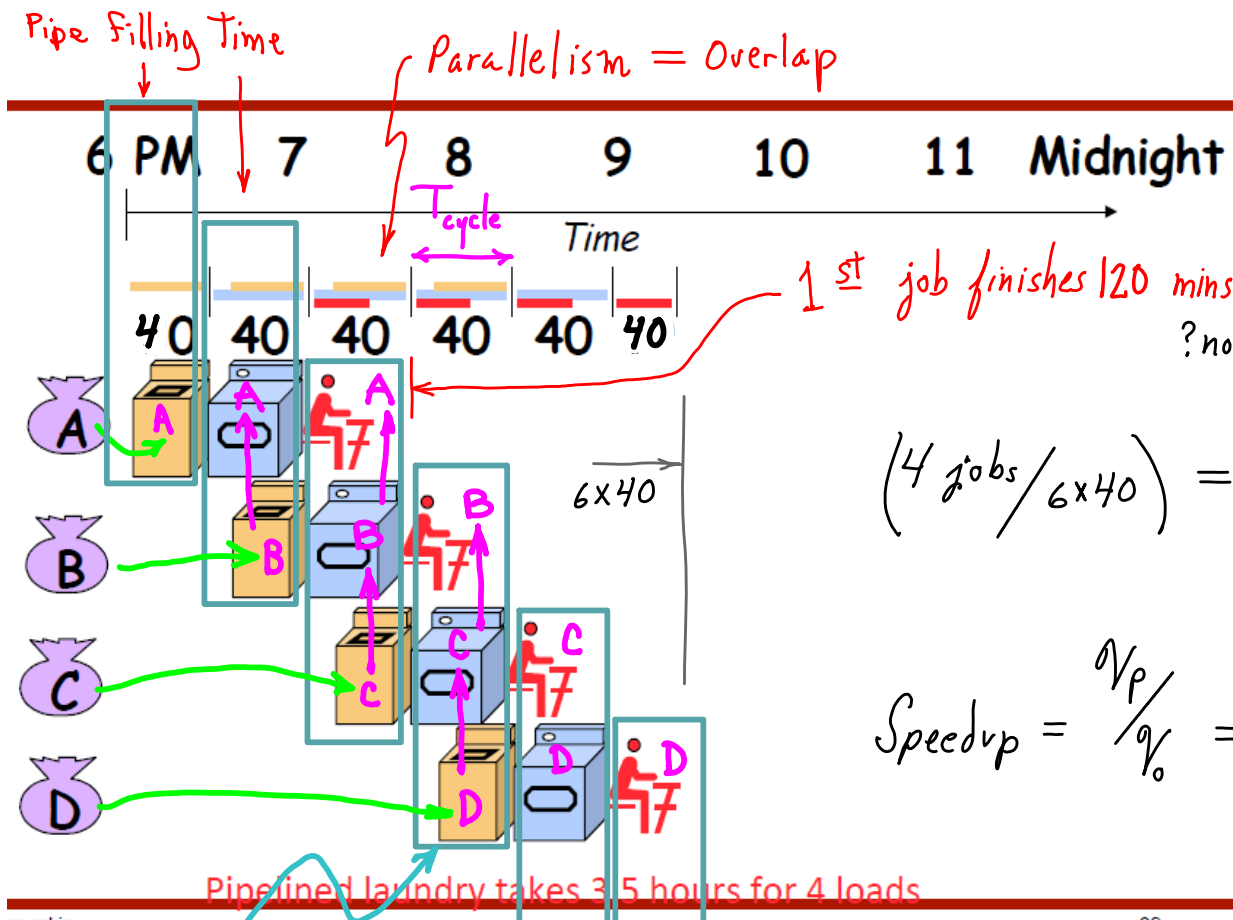
# Sequential: 1-piece functional unit

$T_{cycle} = 30+40+20$

6 PM   7   8   9   10   11   Midnight

*Time*

wash   dry   fold

30  40  20  30  40  20  30  40  20  30  40  20

jobs

A

B

*function unit*

C

D

1st job finishes 90 min latency

$4 \times 90$ → $\left(\dfrac{4 \text{ jobs}}{360}\right) = \%$ throughput

Sequential laundry takes 6 hours for 4 loads

functional unit
at $t \in [3 \times 90, 4 \times 90]$

Laundry unit = [ wash + dry + fold ]

| | | |
|---|---|---|
| $t=0$ WASH A | DRY idle | fold idle |
| WASH idle | $t=30$ DRY A | fold idle |
| WASH idle | DRY idle | $t=70$ fold A |

--- **idle hardware** in **LAUNDRY** unit

--- **delay = 90 min**

|← 90 →| $t=90$ B

# Pipelined: *pipelined* functional unit

Pipe Filling Time

Parallelism = Overlap

6 PM  7  8  9  10  11  Midnight

$T_{cycle}$  Time

40  40  40  40  40  40

1$^{st}$ job finishes 120 mins **latency**

? no improvement ?

A  B  C  D

$6 \times 40$

$\left( 4 \text{ jobs} / 6 \times 40 \right) = V_p$ throughput

$Speedup = \dfrac{V_p}{V_0} = \dfrac{360}{240} = \dfrac{3}{2}$

The pipeline:
3 devices
in parallel

Pipelined laundry takes 3.5 hours for 4 loads

zyrakis                                                    39

pipe draining time

$\max V_p = \dfrac{1 \text{ job completes}}{40}$

$= 4 \text{ jobs} / 160$

$S' = \dfrac{360}{160} = 2^+$

$t \in [0, 40]$     $t \in [80, 120]$

no idle units

| fold | | fold |
| idle | ... | A |
| DRY | | DRY |
| idle | | B |
| WASH | | WASH |
| A | | C |

|← 40 →|

$\dfrac{CR_{new}}{CR_{old}} = \dfrac{T_{cycle}^{old}}{T_{cycle}^{new}} = \dfrac{90}{40} = 2^+$

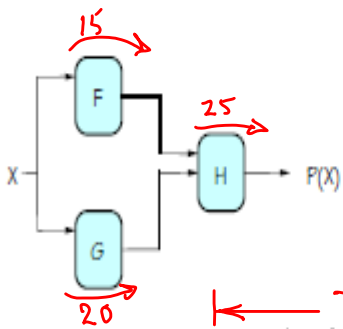Other instructions $S' = 2^+$

$\% \text{ overhead} = \dfrac{\text{fill time} + \text{drain time}}{\text{total time}}$

decreases as
total time increases

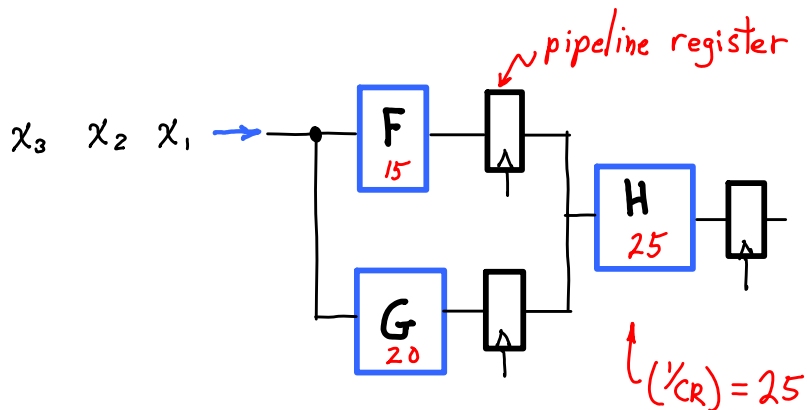# General pipelines



15

25

X → [F] [H] → P(X)

[G]

20

For combinational logic:
latency = $t_{PD}$,
throughput = $1/t_{PD} \approx \left( \dfrac{1 \text{ job}}{T \text{ sec}} \right)$

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

|← T →|

X

F(X)  **15**  ← idle

G(X)  **5**  ← idle

P(X)  **25**

idle ↗

←——— 45 ———→

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

latency = 45

$\gamma = \dfrac{1 \text{ job}}{45}$

$X_3 \; X_2 \; X_1 \rightarrow$

pipeline register

[F 15] [ ]

[ ] [H 25] [ ] → $f(X_3) \; f(X_2) \; f(X_1)$

[G 20] [ ]

$(1/CR) = 25$

1 job exits per 25

latency = 25 + 25 = 50

$\gamma_p = \dfrac{1 \text{ job}}{25}$

## different Pipeline diagrams

Clock cycle →

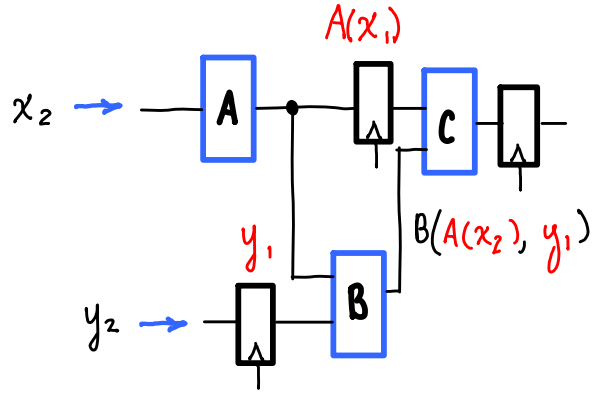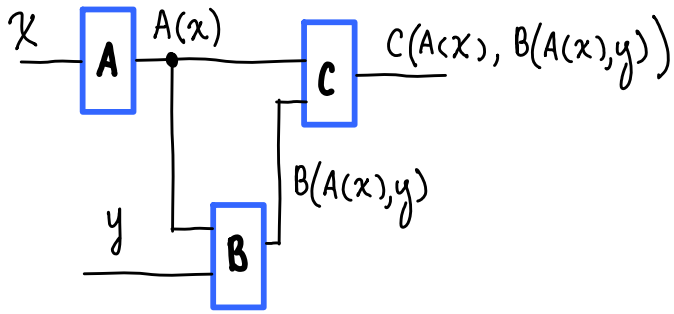| Pipeline stages | i | i+1 | i+2 | i+3 | |
|---|---|---|---|---|---|
| Input | $X_i$ | $X_{i+1}$ | $X_{i+2}$ | $X_{i+3}$ | ... |
| F Reg | | $F(X_i)$ | $F(X_{i+1})$ | $F(X_{i+2})$ | |
| G Reg | | $G(X_i)$ | $G(X_{i+1})$ | $G(X_{i+2})$ | |
| H Reg | | | $H(X_i)$ | $H(X_{i+1})$ | $H(X_{i+2})$ |

pipe

The results associated with a particular set of input data moves diagonally through the diagram, progressing through one pipeline stage each clock cycle.
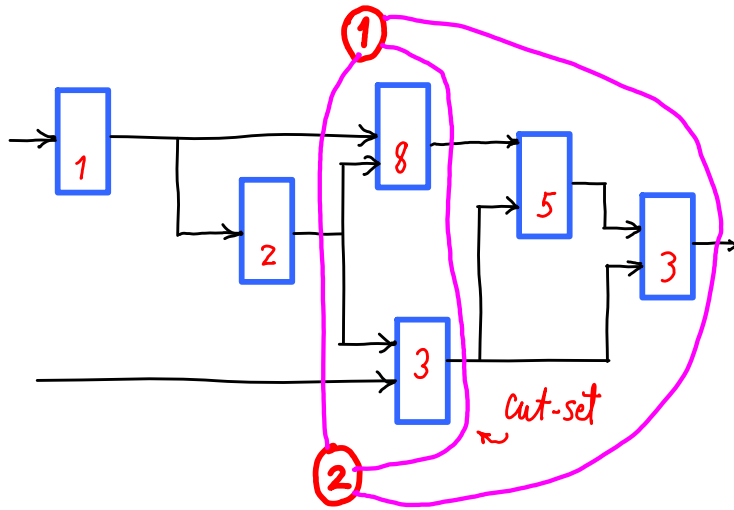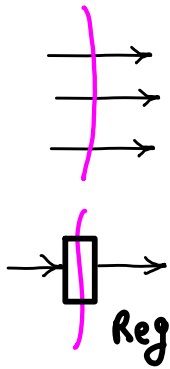
job1    job 2    job 3

# Bad pipelining

$x$ → **A** → $A(x)$ → **C** → $C(A(x), B(A(x),y))$

$B(A(x),y)$

$y$ → **B**

---

$x_2$ → **A** → $A(x_1)$ → □ → **C** → □

$y_1$

$y_2$ → □ → **B** → $B(A(x_2), y_1)$

## what's wrong? How to fix?

---

# General Method

**Cut all paths** in **same direction** **bisect graph**

**===> cut set**

For each cut **add Reg** to path

**Reg**



(1) (2) cut-set

**balance stage delays**

---

## Inductive proof

signals in cut-set are correct



STABLE INPUTS → stable outputs

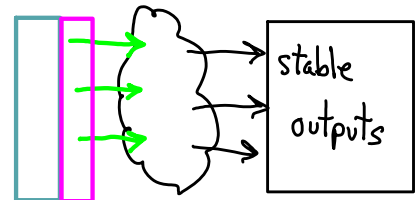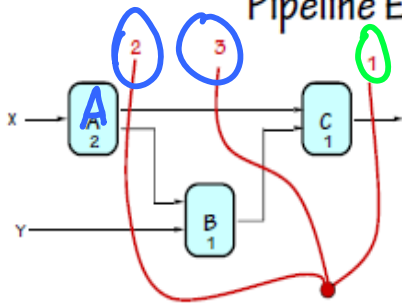inputs to FF are correct  flip flop

STABLE INPUTS

Before Tick

flip flop  outputs are correct

stable outputs

After Tick

# 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.
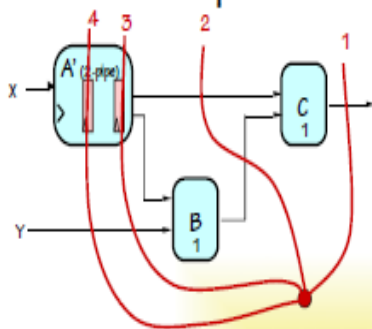- Back-to-back registers are often required to keep pipeline well-formed.

| | LATENCY | THROUGHPUT |
|---|---|---|
| 0-pipe: | $\frac{2+1+1}{4}$ | 1/4 |
| 1-pipe: | | 1/4 |
| 2-pipe: | $\frac{2+2}{4}$ | 1/2 |
| 3-pipe: | $\frac{2+2+2}{6}$ | 1/2 |

$CR = 1/4$

$CR = 1/2$   1, 2-pipe:

$CR = 1/2$   1, 2, 3-pipe:
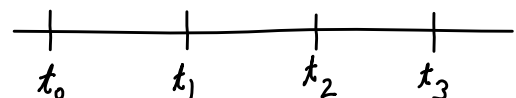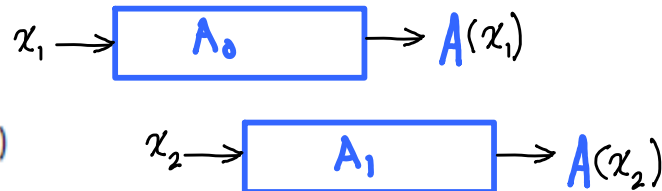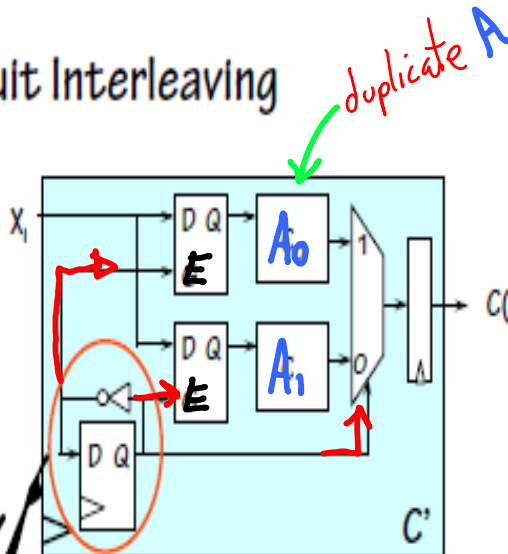
# Pipelined Components



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

4-stage pipeline, thruput=1

## Can't split A in two w/ T=1 ? Fake it!

# Circuit Interleaving

duplicate A

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

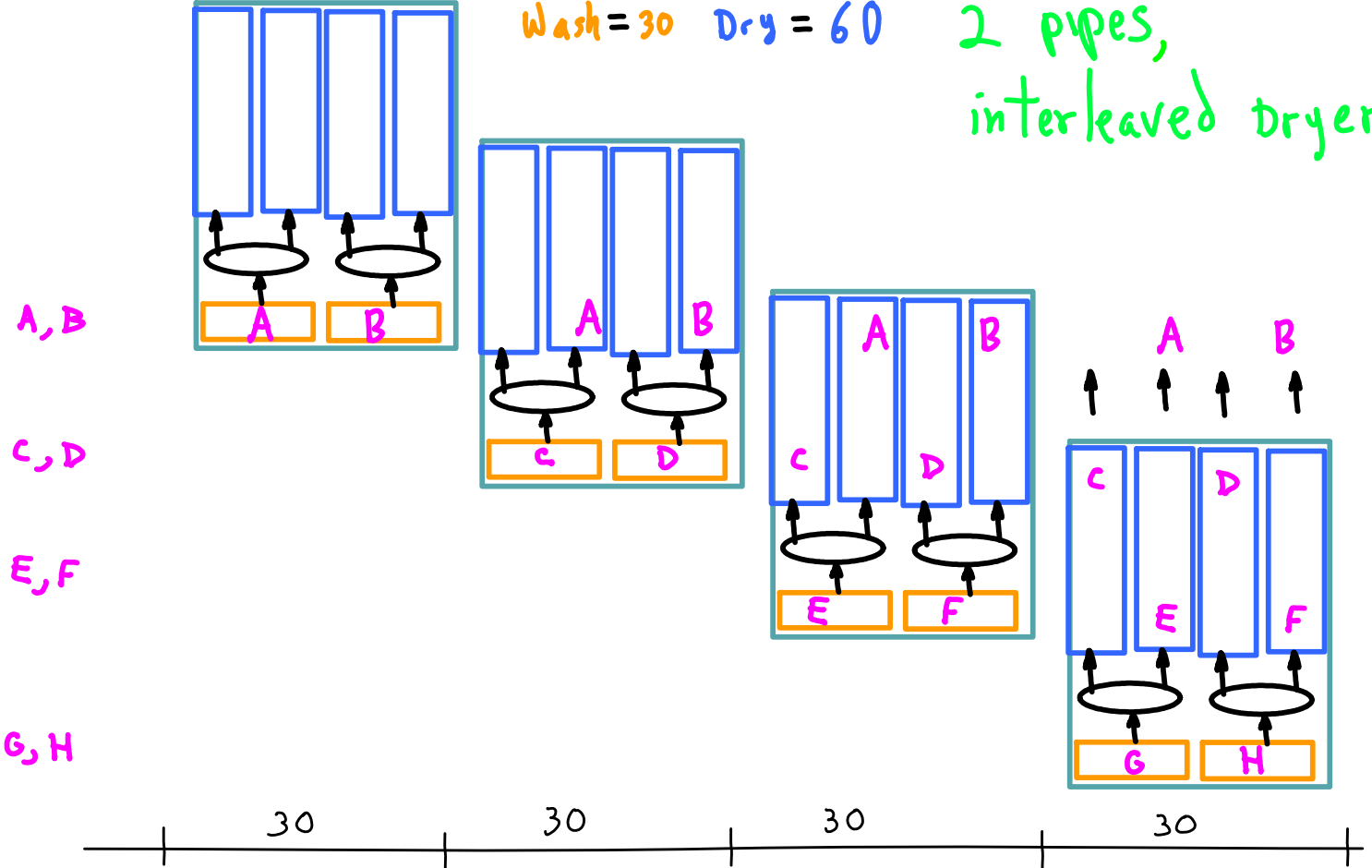This is a simple 2-state FSM that alternates between 0 and 1 on each clock



$C(X_{i-2})$

$C'$

$x_1 \rightarrow \boxed{A_0} \rightarrow A(x_1)$

$x_2 \rightarrow \boxed{A_1} \rightarrow A(x_2)$

$t_0 \quad t_1 \quad t_2 \quad t_3$

# Parallel pipes + interleave



Wash = 30   Dry = 60

2 pipes,
interleaved dryers

A, B

C, D

E, F

G, H

30   30   30   30

latency = 90

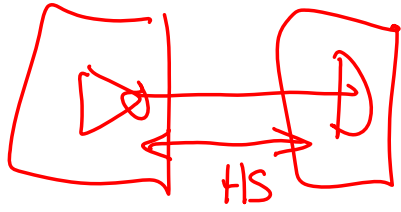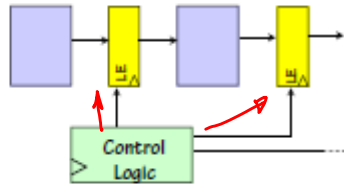$$S = \frac{2 \text{ jobs}/30}{1/90} = 6$$

$$\overline{CPI} = ?$$

$T_{cycle} = 30$

$$1 \text{ cycle}/2 \text{ instr} = \tfrac{1}{2}$$

Provided pipe is full.
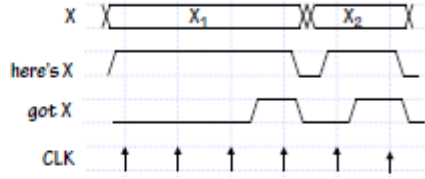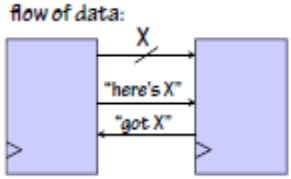
← superscalar

# Control Structure Alternatives

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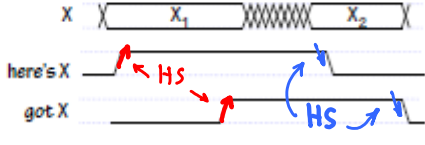Control signals (e.g., load enables)
From FSM controller

Control Logic
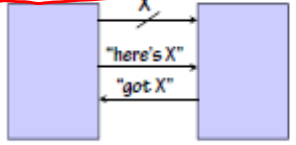
HS

Synchronous, locally-timed:
Local circuitry, "handshake" controls
flow of data:

X
"here's X"
"got X"

X — $X_1$ — $X_2$
here's X
got X
CLK

I/O busses w/ clk
both use same clock

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

X
"here's X"
"got X"

X — $X_1$ — $X_2$
here's X — HS
got X — HS

independent clocks
slow I/O device
w/ HS protocols

( or no clocks,
Self-timed )
( fast
Logic )

6.004 - Spring 2009    3/5/09    L08

# 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!

|  | Synchronous | Asynchronous |
|---|---|---|
| Globally Timed | Centralized clocked FSM generates all control signals. | 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). |

Global
Local

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

$A_1$
c
$A_2$    use c

6.004 - Spring 2009    3/5/09    L08 - Pipelining 26

**Data Dependency Graphs** and **Pipelining**

--- **Entire algorithm** is **dependency graph**
--- **Data-Flow: nodes fire when inputs ready**
--- **Multiple jobs**
--- **High throughput**
--- Is parallelism **as high as possible**?
--- Does **timing depend on data**?

$w_2$  $x_2$  $y_2$  $z_2$  $job_2$
$w_1$  $x_1$  $y_1$  $z_1$  $job_1$   $A_1$

pipe cut

x    z

c    This data depends on

$A_2$

$c = x + z$

RESULTS

**Making 4n-bit multipliers from n-bit ones: 2 "induction steps"**

| $a_3$ | $a_2$ | $a_1$ | $a_0$ |
|---|---|---|---|

$\times$

| $b_3$ | $b_2$ | $b_1$ | $b_0$ |
|---|---|---|---|



$a_0 b_3$
$a_1 b_3$ $a_0 b_2$
$a_1 b_2$
$a_2 b_3$ $a_0 b_1$
$a_3 b_3$ $a_2 b_2$ $a_1 b_1$ $a_0 b_0$
$a_3 b_2$ $a_1 b_0$
$a_2 b_1$
$a_3 b_1$ $a_2 b_0$
$a_3 b_0$

+

$m_3$ $m_2$ $m_1$ $m_0$

*add by columns*

$a_0 b_1$

$m_1$  $m_0$

## Design of 1-bit multiplier "Brick":

*add*



$b_3$
$b_2$
$a_0 b_3$ $b_1$
$a_1 b_3$ $a_0 b_2$ $b_0$
$a_2 b_3$ $a_1 b_2$ $a_0 b_1$
$a_3 b_3$ $a_2 b_2$ $a_1 b_1$ $a_0 b_0$ $a_0$
$a_3 b_2$ $a_2 b_1$ $a_1 b_0$ $a_1$
$a_3 b_1$ $a_2 b_0$ $a_2$
$a_3 b_0$ $a_3$

$\leftarrow$ *partial product*

**Array Layout:**
- operand bits bused diagonally
- Carry bits propagate right-to-left
- Sum bits propagate down

**Brick design:**
- AND gate forms 1x1 product
- 2-bit sum propagates from top to bottom
- Carry propagates to left

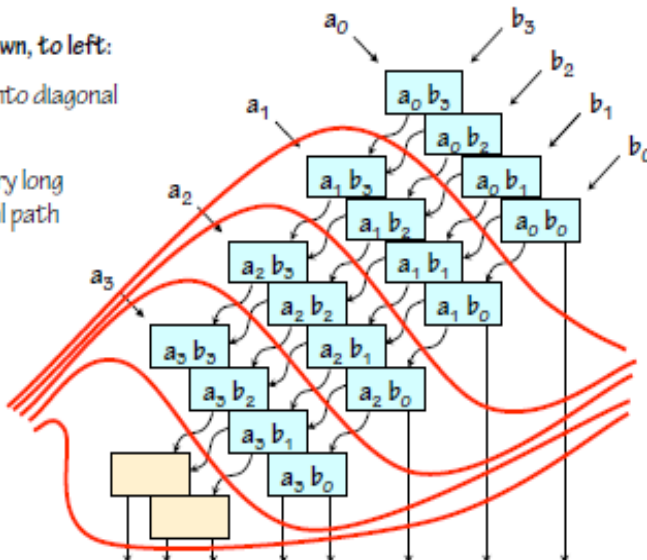Wastes some gates... but consider (say) optimized 4x4-bit brick!



$S_{k+1}$     $S_k$     $b_i$

$C_{k+2} \leftarrow$  FA  FA  $\leftarrow C_k$

$a_j$

$S'_{k+1}$  $S'_k$

*1-bit multiply $\times$*
*2-bit result*

*Column sums*

## Breaking O(n) combinational paths

**LONG PATHS go down, to left:**
- Break array into diagonal slices
- Segment every long combinational path



$a_0$     $b_3$
$a_1$     $b_2$
$a_0 b_3$     $b_1$
$a_0 b_2$     $b_0$
$a_2$  $a_1 b_3$  $a_0 b_1$
$a_1 b_2$  $a_0 b_0$
$a_3$  $a_2 b_3$  $a_1 b_1$
$a_2 b_2$  $a_1 b_0$
$a_3 b_3$  $a_2 b_1$
$a_3 b_2$  $a_2 b_0$
$a_3 b_1$
$a_3 b_0$

*cut both ways?*
*$\Rightarrow$ systolic array*

**GOAL: $\Theta(n)$ stages;  $\Theta(1)$ clock period!**

# 1-cycle MIPS processor
--- **Harvard Architecture** (two memories)
--- **clocking PC initiates cycle**



## Single Cycle Processor Performance
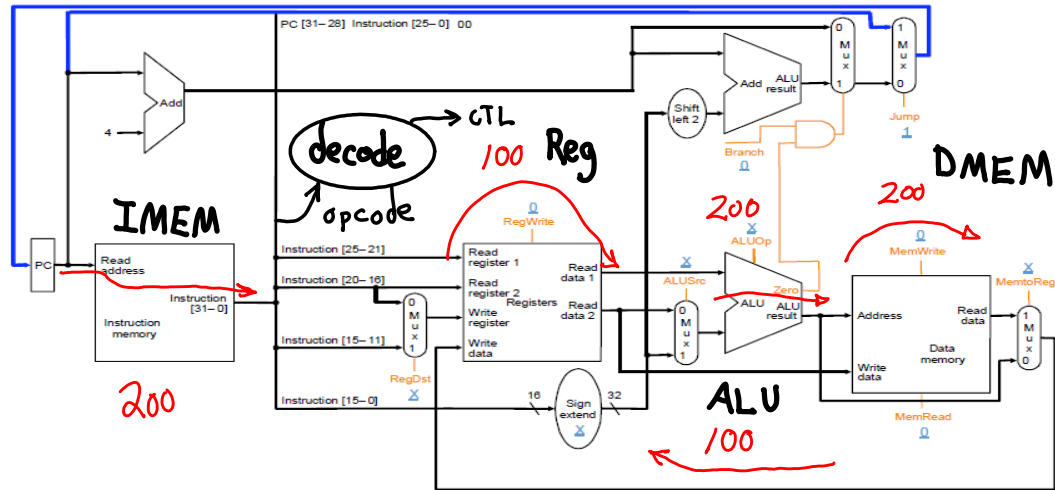
- Functional unit delay
  - Memory: 200ps
  - ALU and adders: 200ps
  - Register file: 100 ps

$$ps = 10^{-12} sec$$

| Instruction Class | Instruction memory | Register read | ALU operation | Data memory | Register write | Total |
|---|---|---|---|---|---|---|
| R-type | 200 | 100 | 200 | | 100 | 600 |
| load | 200 | 100 | 200 | 200 | 100 | 800 |
| store | 200 | 100 | 200 | 200 | | 700 |
| branch | 200 | 100 | 200 | | | 500 |
| jump | 200 | | | | | 200 |

max delay = $T_{clock}$

$$\frac{1}{T_{clock}} = \frac{1}{0.8\,ns} = 1.25\,GHz$$

- CPU clock cycle = 800 ps = 0.8ns (1.25GHz)

## what if we let clock trigger delay by opcode?

BR 560 ps
LD 800 ps

- Instruction Mix
  - 45% ALU
  - 25% loads
  - 10% stores
  - 15% branches
  - 5% jumps

| Instruction Class | Instruction memory | Register read | ALU operation | Data memory | Register write | Total |
|---|---|---|---|---|---|---|
| R-type | 200 | 100 | 200 | | 100 | 600 |
| load | 200 | 100 | 200 | 200 | 100 | 800 |
| store | 200 | 100 | 200 | 200 | | 700 |
| branch | 200 | 100 | 200 | | | 500 |
| jump | 200 | | | | | 200 |

ps → 0.6 ns
0.8
0.7
0.5
0.2

⟹ **1.6 GHz**

- CPU clock cycle = 0.6x45% + 0.8x25% + 0.7x10% + 0.5x15% + 0.2x5%
  = 0.625 ns (1.6GHz)

speedup? $S_{new-old} = \frac{1.6}{1.25} = 1.28$

# Pipelining Load *lw*

- Load instruction takes 5 stages
  - Five independent functional units work on each stage
    - Each functional unit used only once
  - Another load can start as soon as 1st finishes IF stage
  - Each load still takes 5 cycles to complete
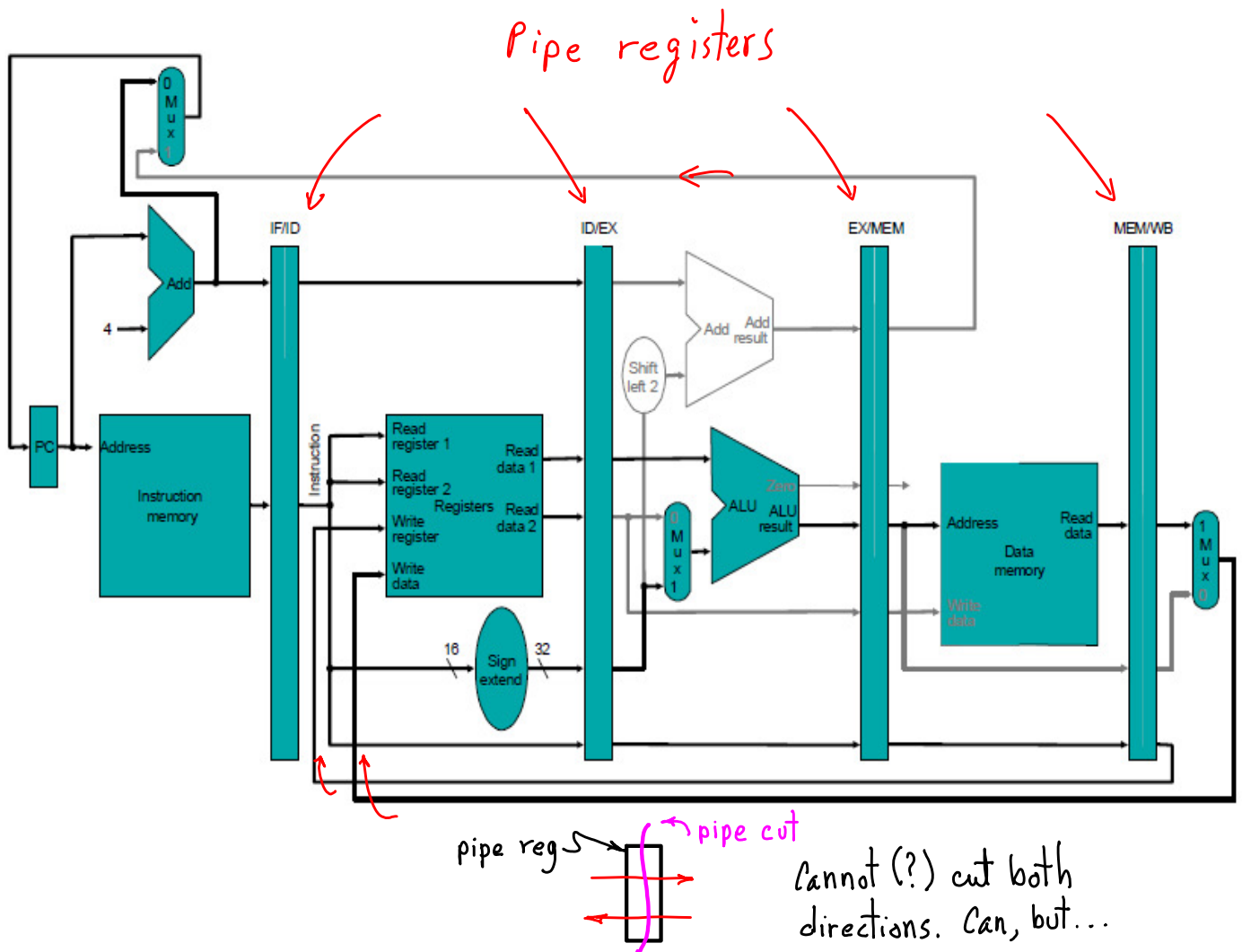  - The *throughput*, however, is much higher

*when all stages busy*

*1 instr. exits per cycle*

$$\overline{CPI} = \frac{1 \text{ instr.}}{1 \text{ cycle}}$$

*latency = 5 cycles*



| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6 | Cycle 7 |
|---|---|---|---|---|---|---|---|
| Clock | | | | | | | |
| JOB 1st lw | IF | RF/ID | EX | MEM | WB | | |
| 2nd lw | | IF | RF/ID | EX | MEM | WB | |
| 3rd lw | | | IF | RF/ID | EX | MEM | WB |

*Pipe registers*



*pipe regs*   *pipe cut*

*Cannot (?) cut both directions. Can, but...*

# Max out pipeling?

$$\text{delay} \rightarrow \frac{\text{delay}}{n} \implies CR \rightarrow n \times CR \qquad (n \rightarrow \infty?)$$

NO:

— $T_{cycle} > (\text{setup} + \text{hold}) \implies CR < 1/(\text{setup} + \text{hold})$

— small operations harder to divide into pieces

— n stages $\implies$ (n fill/drain overhead) $\times \left( T_{cycle} = 1/n \times CR \right)$

looks like n cancels, but keeping pipe full gets harder
→ pay penalty more often

— $n \times CR \implies n \times \text{Power}$

$\approx 20$ stages limit for CPU (But elsewhere?)