

Chapter 6

PIPELINING & SUPERSCALAR TECHNIQUES

• LINEAR PIPELINE

- Cascade of processing stages to perform a fixed function over a stream of data flowing from Stage 1 to stage n
- Applications : instruction execution
arithmetic computation
memory-access functions

Models → Asynchronous

↓ → Synchronous

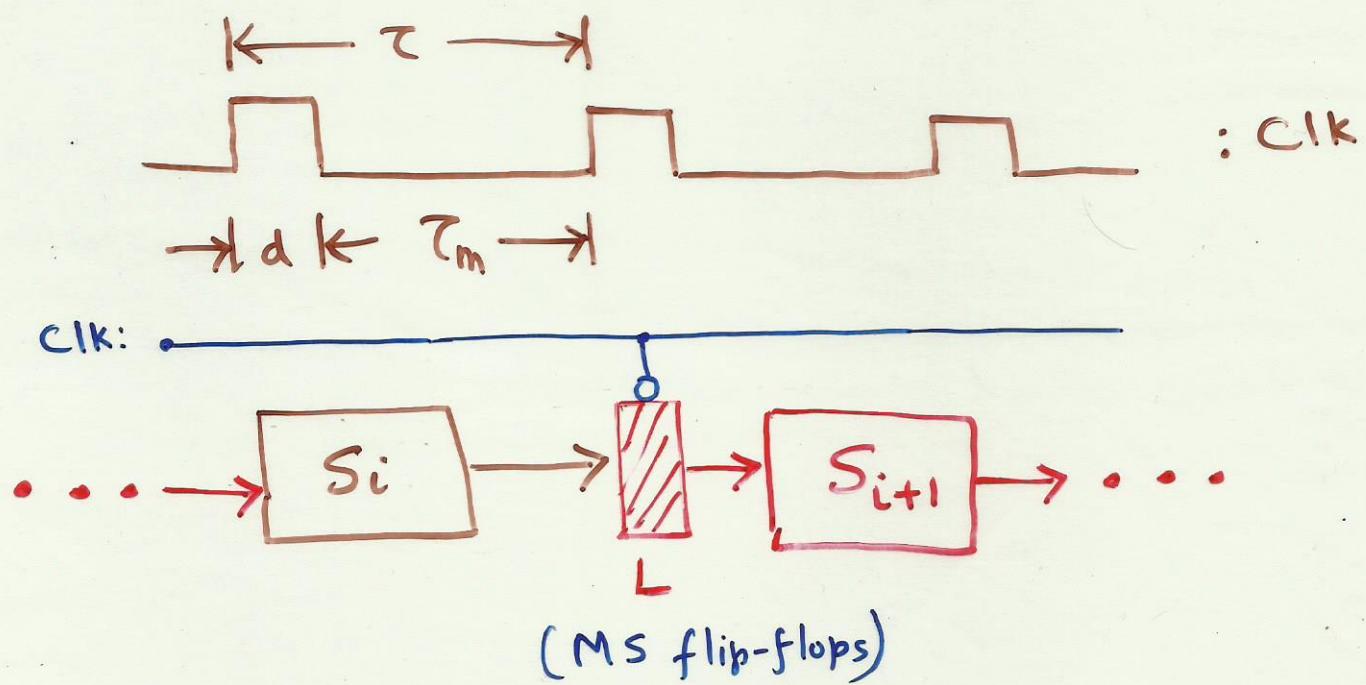
→ depends on the flow control along the pipeline stages.

Asynchronous Model

- Adjacent stages use handshaking protocol (Ready/ack signals)
- Wormhole routing ;
- Variable throughput is possible
- delay per stage is different

[fig. 6.1 (a), pg: 266]

Synchronous Model



- Upon arrival of a clock pulse, all the latches transfer data to the next stage simultaneously

- All stages have more-or-less equal delays

Pipeline utilization is studied by using a reservation table.

	1	2	3	4
stages	s_1	s_2	s_3	s_4
	X		X	
		X		X
			X	
				X

$t \rightarrow$ CLK cycles

- After 4 clock cycles an o/p appears from stage s_4 .

T_m : maximum stage delay : $\max_{1 \leq i \leq k} \{ T_i \}$

$$T = T_m + d \rightarrow ①$$

In general, $T_m \gg d$.

$$\text{Pipeline frequency } (f) = (\gamma/c) \rightarrow \textcircled{2} \quad (4)$$

Thus, f represents the max. throughput.
 However, the actual throughput may be less than f , as it depends on the initiation of the successive tasks.

Clock skewing

$$d + t_{\max} + s \leq \underline{\underline{T}} \leq T_m + t_{\min} - s$$

i.e., t_{\max} : time delay of the longest logic path within a stage

t_{\min} : shortest logic path delay

Thus to avoid any race, we need

$$T_m > t_{\max} + s \quad [s: \text{offset time}]$$

$$\& d \leq t_{\min} - s$$

- if $s=0$ (ideal case) $t_{\max} = T_m$ & $t_{\min} = d$

(5)

If we have k stages & n tasks

$n > k$; we need

k cycles to complete the first job

+ $(n-1)$ cycles to complete the other jobs

$$\Rightarrow T_k = [k + (n-1)] \tau \rightarrow ③$$

$$\text{Speed up} = \frac{T(\text{non-pipelined System})}{T(\text{pipelined system})}$$

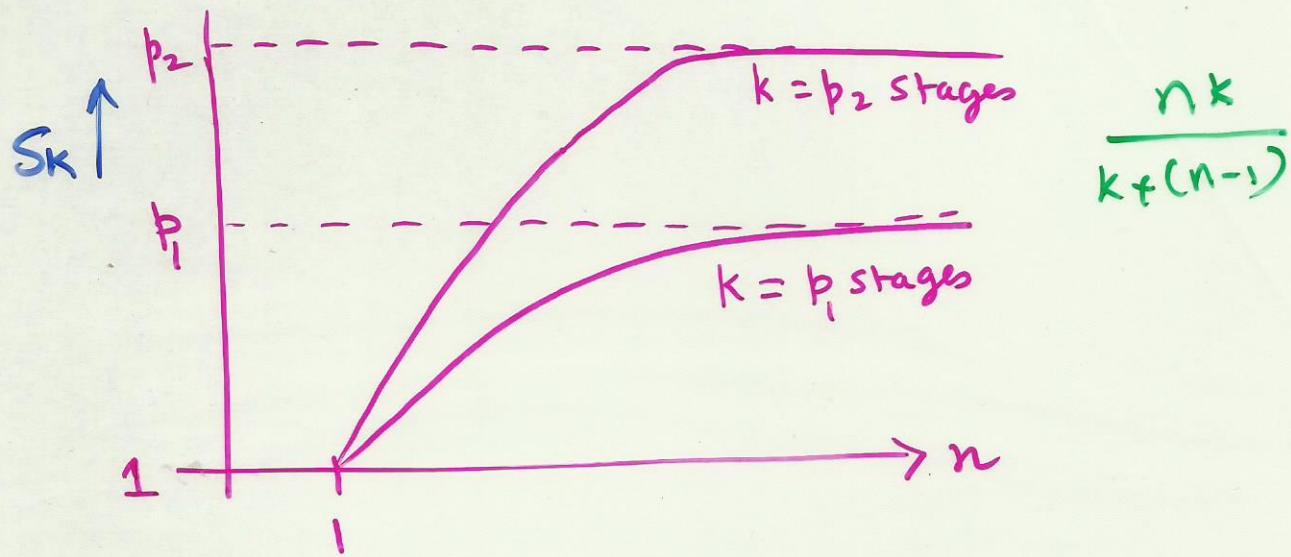
$$= \frac{n \cdot k \cdot \tau}{[k + (n-1)] \tau} \rightarrow \left\{ \begin{array}{l} \text{equivalent function} \\ \text{non-pipelined} \\ \text{processor with} \\ \text{a delay } (k\tau). \end{array} \right.$$

$$S_k = \frac{n k}{k + (n-1)} \rightarrow ④$$

\rightarrow [See fig 6.2(a), S_k vs n , for some fixed k]

(6)

Thus, as $n \rightarrow \infty$, $S_k \rightarrow k$. This can be seen from the fig.



Can we determine an optimal # of stages to be used?

micro pipelining \rightarrow finest level of pipelining

$$2 \leq k \leq 15$$

macro pipelining \rightarrow at the "processor level"

Optimal choice \rightarrow maximize

$$\left(\frac{\text{Performance}}{\text{Cost}} \right) [\text{PCR}]$$

(7)

t : total time required for a non-pipelined sequential pgm

With k stages, the required clock period is

$$b = \left(\frac{t}{k}\right) + d, \quad d: \text{latch delay}$$

$$\Rightarrow \text{Throughput (max)} \ f = 1/b$$

$$= \frac{1}{\left[\frac{t}{k} + d\right]}$$

Total pipeline cost is estimated by: $(c + kh)$,

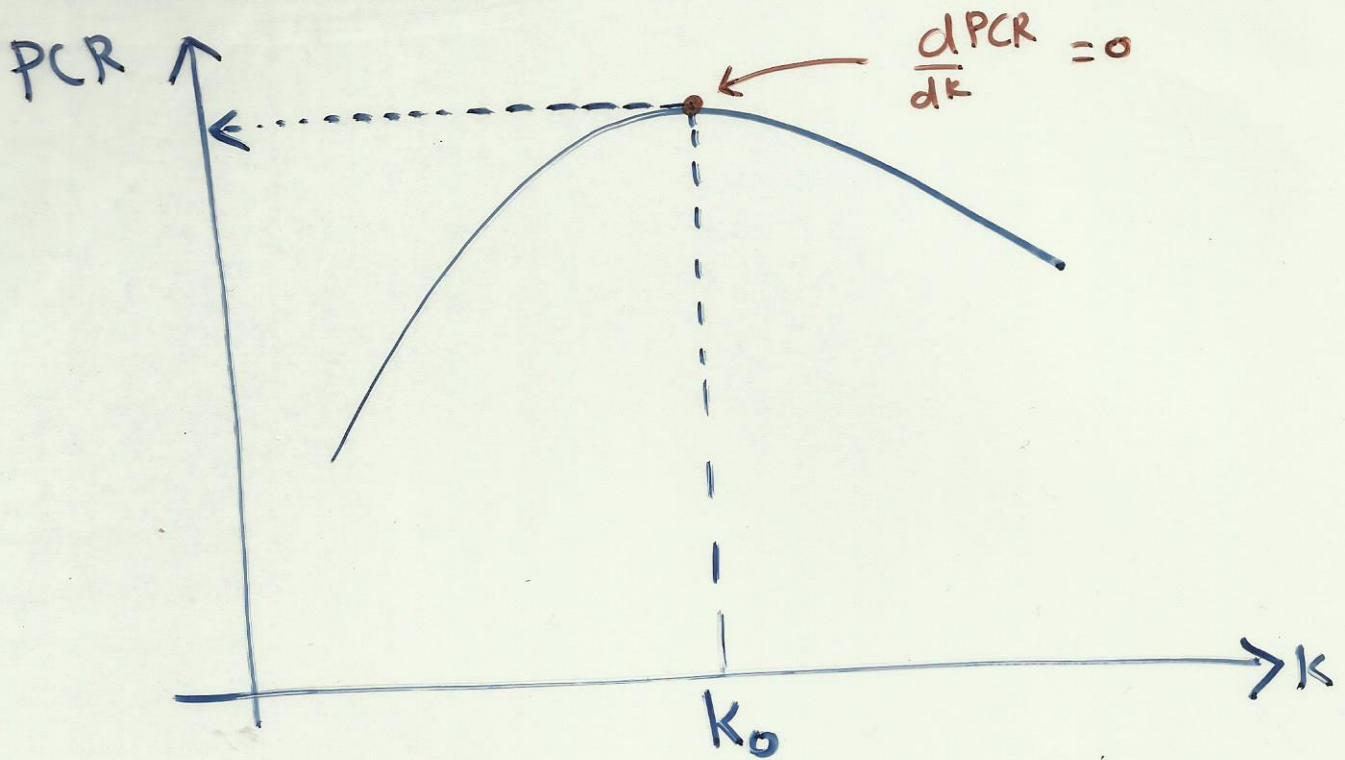
where, c : Cost of all logic stages

h : Cost of each latch

Thus, $\text{PCR} = \left(\frac{f}{c + kh} \right)$

$$\text{PCR} = \frac{1}{\left[\frac{t}{k} + d\right] [c + kh]}$$

→ ⑤



K_0 (Optimal # of stages)

$$\rightarrow K_0 = \sqrt{\frac{t \cdot c}{d \cdot h}} \rightarrow ⑥$$

Efficiency (E_K) = $\left(\frac{S_K}{K} \right) = \frac{n}{K + (n-1)}$

[Speed-up per stage]

→ ⑦

As, $n \rightarrow \infty$, $E_K \rightarrow 1$, as expected

Also, $E_K \geq \frac{1}{K}$ [Sets a lower bound on E_K , when $n=1$]

(9)

Pipeline Throughput

H_k : # of tasks performed per unit time

$$= \frac{n}{[k + (n-1)\tau]}$$

$$H_k = \frac{n \cdot f}{k + (n-1)\tau} \rightarrow ⑧$$

Thus,

$$H_k = E_k \cdot f = E_k / \tau = \frac{S_k}{k \tau}$$

Problems: 6.1, 6.4,

Non-linear Pipeline Processors

(10)

- Allows feedforward & feedback connections in addition to the streamline connections.

[see fig 6.3 (a), pg: 27]

With the feedforward & feedback connections, scheduling of tasks on a pipeline becomes a difficult problem.

→ basic question : how to maximize the throughput of the pipeline?

Equivalently,

with what minimum time gap, I can initiate tasks?

(11)

TABLE
1.

s_1	X				X	X	
s_2		X	X				
s_3			X	X		X	
	1	2	3	4	5	6	7

 $\rightarrow t$

→ initial state will be given. In this fig, this is the way in which data flows through the pipeline.

Reservation table → timing diagram; also referred to as a time-space graph

Each function evaluation is specified by a reservation table.

- Fig 6.3 (b) → we will consider this table for analysis throughout.
(pg: 271)

Some Terminology!

- Evaluation time : # of columns in a RT
- A checkmark (x) in a cell (i,j) means stage i is used at the j^{th} CLK cycle.
- Latency : # of time units (CLK cycles) between 2 initiations of a pipeline
is the latency between them.]
- Any attempt to use a stage by 2 initiations at the same time results in a Collision
 Collision \rightarrow Resource conflicts ; must be avoided

- Latencies that cause collision are called **forbidden latencies**. (FL)

Refer to Table 1.

	1	2	3	4	5	6	7	8	9	10
1	X		X			X		XX		X
2		X		XX		X				
3		X		XX		XX		XX		X

first initiation $t=1$ } latency of 2
 Second initiation $t=3$

Thus, 2 is a FL

Similarly, 5 is also a FL

→ Verify: $FL = \{2, 4, 5, 7\}$

→ Permissible latencies
 $P = \{1, 3, 6, 10\}$



(4)

- How to find/determine the forbidden latencies, given a RT ?

Check the distance between any two checkmarks in the same row of the RT.

Refer to Table 1.

$$\begin{aligned}
 \text{row 1 : } \text{FL} &\rightarrow \{2, 5, 7\} & t_1, t_6, t_8 \\
 && \underbrace{}_{5} \quad \underbrace{}_2 \\
 \text{row 2 : } \text{FL} &\rightarrow \{2\} \quad (t_2, t_4) \\
 \text{row 3 : } \text{FL} &\rightarrow \{2, 4\} \quad t_3, t_5, t_7 \\
 && \underbrace{}_2 \quad \underbrace{}_4 \\
 \Rightarrow \text{FL} &\rightarrow \{2, 4, 5, 7\}
 \end{aligned}$$

- Latency Sequence : sequence of permissible (non-forbidden) latencies between successive initiations

(15)

- Latency cycle : latency sequence which repeats the same subsequence indefinitely.

Refer to Table. 1.

$$LC(1,8) = 1, 8, 1, 8, \dots \dots$$

$$\Rightarrow t_1 \rightarrow \text{job 1}$$

$$t_2 \rightarrow \text{job 2}$$

$$t_{10} \rightarrow \text{job 3}$$

$$t_{11} \rightarrow \text{job 4}$$

⋮
⋮

[See fig 6.5 (a)]

i.e., successive initiations of new tasks are separated by one cycle and 8 cycles alternately.

Similarly, $LC(3) = 3, 3, 3, \dots \dots$

$$LC(6) = 6, 6, 6, \dots \dots$$

$$\frac{\text{Average Latency}}{(ALC(\dots))} = \frac{\sum_i l_i}{\# \text{ of latencies}}$$

$$ALC(1,8) = \frac{1+8}{2}, ALC(3) = 3, ALC(6) = 6$$

(16)

- Given these, how do I effect a collision free scheduling?

formal treatment

Collision vectors

- For a RT with n columns,
the maximum FL (m) $\leq n-1$.
- The permissible latency p should be as small as possible, i.e.,
 $1 \leq p \leq m-1$
- $C = (C_m, C_{m-1}, C_{m-2}, \dots, C_2, C_1)$
is a m -tuple vector that captures the $FL \leq p$ latencies
- $$\rightarrow \begin{cases} C_i = 1, & \text{if latency } i \text{ causes collision} \\ = 0, & \text{if latency } i \text{ is permissible} \end{cases}$$

(17)

Note that $m = 1$ always, as it corresponds to the max. FL.

Refer to Table 1. $n = 8$, $m \leq 7$

$$FL = \{2, 4, 5, 7\}$$

$$P = \{1, 6, 3, 0\}$$

$m = 7 \Rightarrow 7$ bit collision vector

$$C = C_7 C_6 C_5 C_4 C_3 C_2 C_1$$

$$1 \textcolor{red}{0} \quad 1 \quad 1 \textcolor{red}{0} \textcolor{blue}{1} \textcolor{red}{0}$$

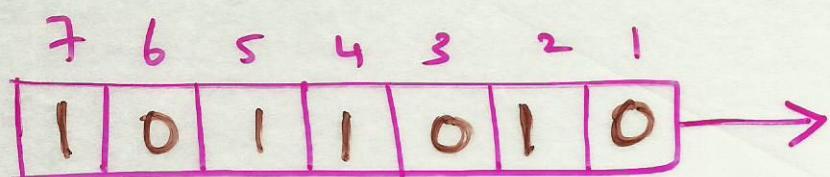
Thus, the collision vector (to evaluate a function X) for a given RT is

$$\boxed{C_X = 1011010}$$

O/p function

$$GY = \dots \dots \dots$$

(18)



m-bit shift register

At time 1, we have an initial collision vector

Each 1-bit shift corresponds to an increase in the latency by 1.

Thus, starting from time 1, when a 0-bit emerges from the right end

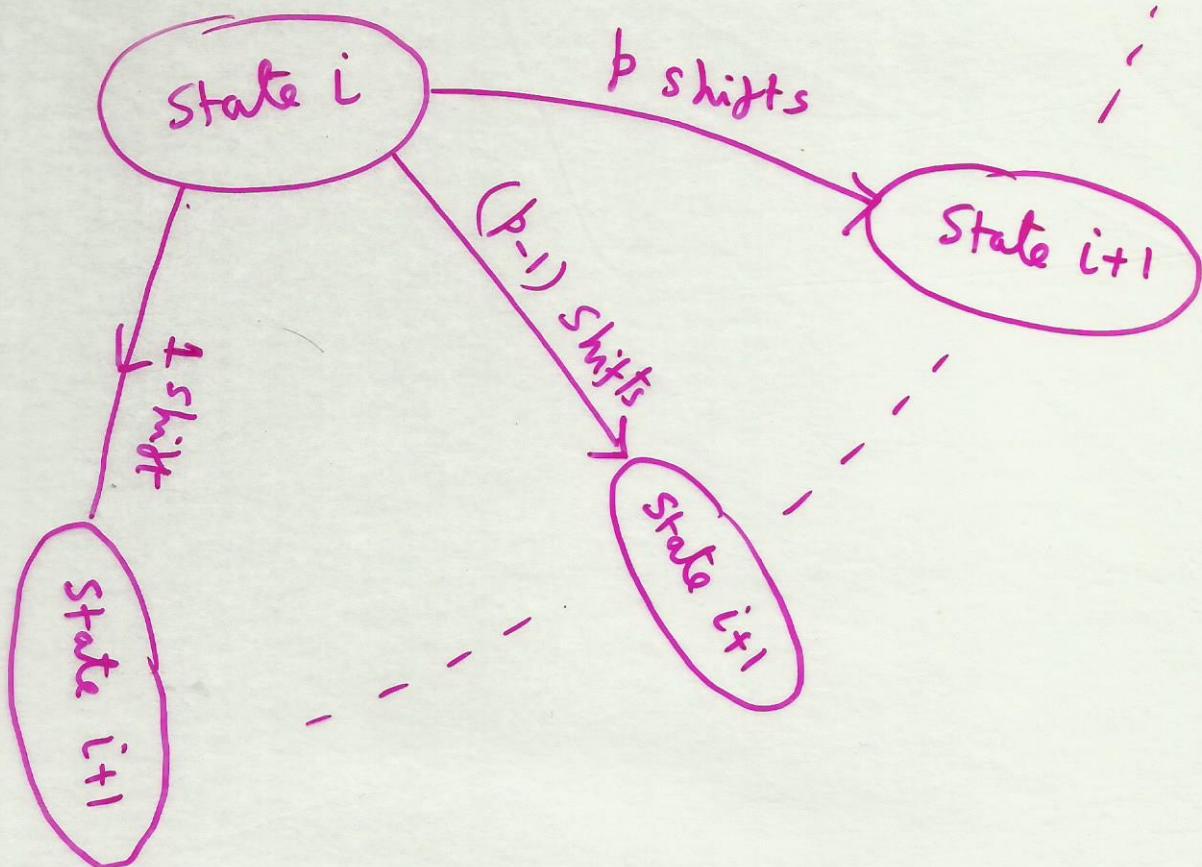
after b shifts, this means

b is a permissible latency

Verify: t_1, t_3, t_6 & $\{1, 3, 6\}$ are Per. lat.

- As we shift to the right, a logical '0' enters from the left end.

(19)



- next state after p shifts is obtained

by $C_x \xrightarrow{\text{bitwise OR operation}} \text{OR (shifted contents in SR)}$

Jhvs,

$$C_x = (1011010) \xrightarrow{\text{RTS}} \begin{matrix} 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \end{matrix}$$

after 1 right shift gives

$\rightarrow 0101101$ (shifted contents)

1011010 (C_x)

$\overline{1111111}$

next state

(20)

$$C_x = (1011010)$$

1 right shift : 0101101

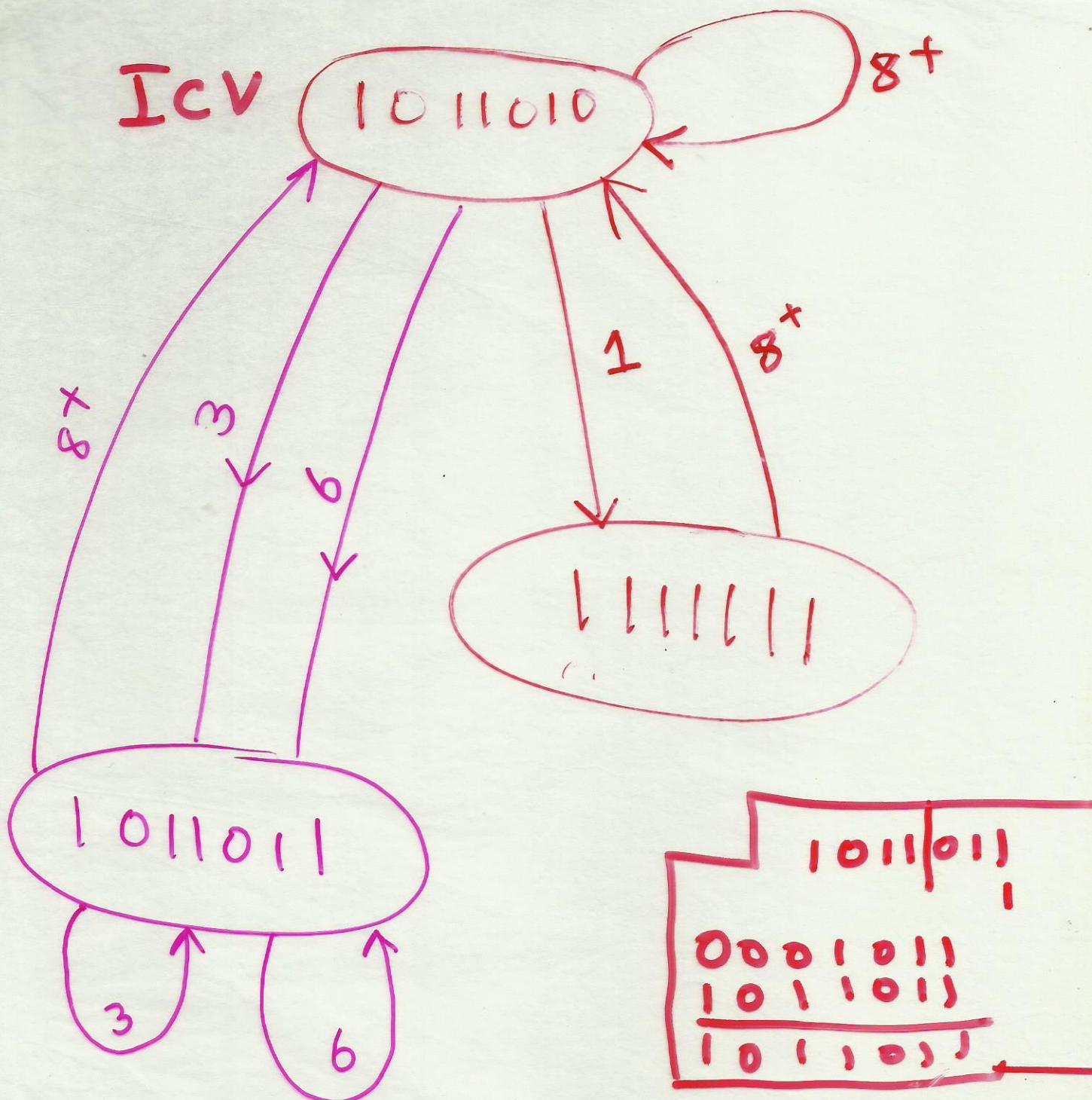
2nd r.sh : 0010110 ✓

3rd r.sh : 0001011 ←

next state after 3 shifts

$$\begin{array}{r}
 & 0001011 \\
 (\text{OR}) & \underline{1011010} \leftarrow C_x \\
 & \underline{\checkmark 1011011} \checkmark
 \end{array}$$

- With C_x , if you shift 8 times, you will reach C_x
- In general, with $(m+1)$ shifts or greater all transitions are redirected to the initial state (denoted as m^+)



State diagram for table 1 (See fig 6.6(b))

- Once the initial collision vector is known State diagram can be uniquely determined.

Thus,

- 0's & 1's in the current state, say at time t , indicate the P & F L latencies , at time t .
- bitwise ORing of the shifted contents of the present state with the initial collision vector is meant to prevent collisions from future initiations starting at time $t+1$ onwards.
- Collision free scheduling \Rightarrow choose that permissible latency that has minimum average latency (MAL)

Greedy cycles

There are infinitely many latency cycles one can trace from the SD.

(1,8)

(1,8,6,8)

(3)

(6)

(3,6,3)

:

We look for "Simple Cycles"



Latency cycle in which each state appears only once (including constant cycles)

(1,8), (3,8), (6,8), (3),(6),
(8)

→ $(1, 8, 6, 8)$ is not simple

Since state (1011010) is traversed twice

Some Simple cycles are Greedy.

- Greedy cycle is one whose edges are all made with minimum latencies from their respective starting states.

→ refer to the SD (slide # 21) [Fig 6.6(b)]
 greedy cycles : $(1, 8)$ & 3

- Greedy cycles must first be simple cycles
- Avg. latency of ~~all~~ a greedy cycle must be lower than those of simple cycles

$$\rightarrow (1, 8) \rightarrow \frac{1+8}{2} = 4.5$$

$$4.5 < \text{ALC}(6, 8)$$

↗ simple cycle

The greedy cycle (3) has a constant latency & if equal to the minimal average latency (MAL), in

This example.

- * At least one of the greedy cycles will lead to MAL

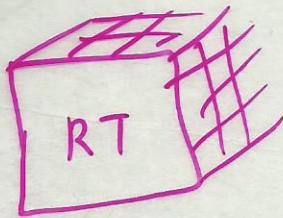
Thus,

Collision free scheduling \Rightarrow find all the greedy cycles from the set of simple cycles.

\rightarrow Greedy cycle yielding the MAL is the final solution!

So, The procedure is :

given :



find

FL & P sets

draw
SD



form the CV

identify simple cycles



identify greedy cycles



extract MAL !

Now, Can we improve the throughput ?

If so, how ?

→ The purpose is to achieve an optimal latency cycle, which is absolutely the shortest.

Some Observations

[1972 by Shar]

Max. # of
checkmarks
in any row

$$\leq \text{MAL} \leq \text{average latency of any greedy cycle}$$

and avg. latency of any greedy ~~cycle~~ cycle $\leq \# \text{ of } 1's \text{ in the initial CV} + 1$

facts : (after going through the proof's)

- Optimal latency cycle $\in \{\text{greedy cycles}\}$
- However, a greedy cycle is not sufficient to guarantee the optimality of MAL.
It is the lower bound that guarantees the optimality.

[out e.g: $\text{MAL} = 3$ meets the lower bound]

From our Table 1 (6.6 b)
 & SD ←

- Upper bound on the MAL = $4+1=5$
 - To optimize the MAL, we need to hit the lower bound by modifying the RT.
-

approach: reduce the max. # of checkmarks in any row, but RT must preserve the function!

∴ Idea is to introduce non-compute delays

⇒ modified SD! → results in greedy cycles that meet the low. bound.

Example 6.3 (page 277- 279)

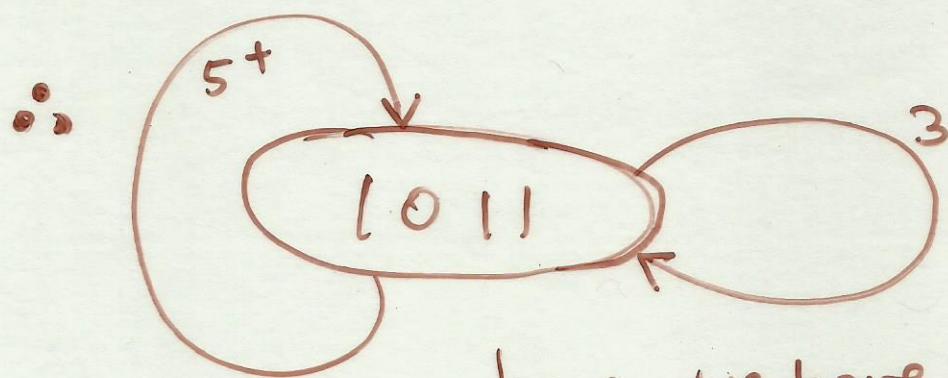
Consider

	1	2	3	4	5
S_1	X				X
S_2		X		X	
S_3			X	X	

$$FL = \{4, 2, 1\} \quad m = 4$$

$$C = C_4 C_3 C_2 C_1$$

$$\therefore C = (1011)$$



3 shifts

0001	(OR)
1011	
1011	

here we have a
greedy cycle of latency

$$3 = MAL$$

But in row 2, we have 2 checkmarks.

\Rightarrow lower bound is 2

\Rightarrow $MAL = 3$ is not optimal!

So, we go about inserting some non-computing delays.

(D_1)
 → a delay after S_3 will delay both X_1 & X_2 operations one cycle beyond

Time 4

	1	2	3	4	5	6
S_1	X			X		
S_2		X	X			
S_3		X	X			

	1	2	3	4	5	6
S_1	X			O	x_1	
S_2	X		X			
S_3		X	O	x_1		
D_1				D_1		
D_2						

	1	2	3	4	5	6	7
S_1	X			O	O	x_2	
S_2		X	X				
S_3		X	O	x_1			
D_1				D_1			
D_2							

→ t

→ Delay elements are inserted as extra stages.

→ Note: All that we need to take care is that S_1 starts after S_2 & S_3 finish their computation. This fact must be preserved! (see the original RT & modified RT now!)

→ modified table has $3+2=5$ rows
 $5+2=7$ columns

Now determine the CV for this table.

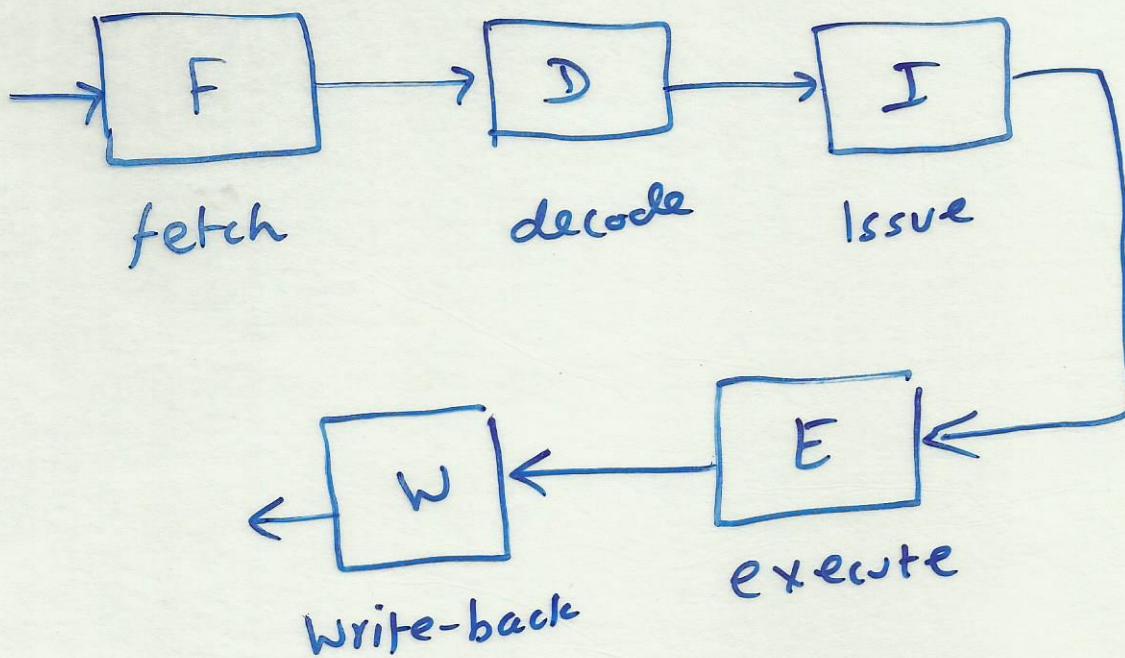
100010

See the SD (pg: 278) ←

verify the MAL } $\Rightarrow \frac{1+3}{2} = 2$
 given by }
 greedy cycle (1,3) }
 lower bound for the
 MAL .

Instruction Pipeline Design

- We will see some typical pipelines used in CISC & RISC scalar processors.



F: fetches instructions from cache, usually one per cycle

D: reveals the instruction function to be performed & identifies the resources needed

I: Reserves resources (GPRs, buses, etc)
+ Operands are read from the registers.

E : Execution stage

W : Writing back the results into the registers

* Memory reading/storing operations are usually treated as a part of E stage.

in-order execution: If an instruction is blocked from issuing due to any data/resource dependence/conflicts, all the instructions that follow are blocked
[see fig 6.9 (b)]

Reordered issuing:

Issue all the required load operations in the beginning.

→ results in an improved time performance.

(33) b



Example of an out-of-order execution

(Delay due to data dependency).

I1)	ADD R1, R2, R3	$C(R_3) \leftarrow C(R_1) + C(R_2)$
I2)	MUL R3, R4, R5	$C(R_5) \leftarrow C(R_3) * C(R_4)$
I3)	SUB R7, R2, R6	$C(R_6) \leftarrow C(R_7) - C(R_2)$
I4)	INC R3	$C(R_3) \leftarrow C(R_3) + 1.$

	1	2	3	4	5	6	7	8	9
I1	FI	DE	EX	MEM	SR				
I2		FI	DE	X	X	EX	MEM	SR	
I3			FI	DE	EX	MEM	SR		
I4				FI	DE	X	X	MEM	SR

. MUL (I2) should not be executed in cycle 4 as the value of R_3 it needs will not be stored in the reg. file by I1 & thus not available.

→ data dependency, called data hazard.

• But, I₃ does not need R₃ & hence,
it can proceed without waiting (!!)

• I₄ also needs R₃. It can EX at clock
cycle 6 as the result is available for use.

• But, ALU is being used in clock cycle 6
⇒ it has to wait till cycle 7 to increment R₃.

* Out-of-order execution. Sometimes unacceptable!

Soln: lock the pipeline. Thus,

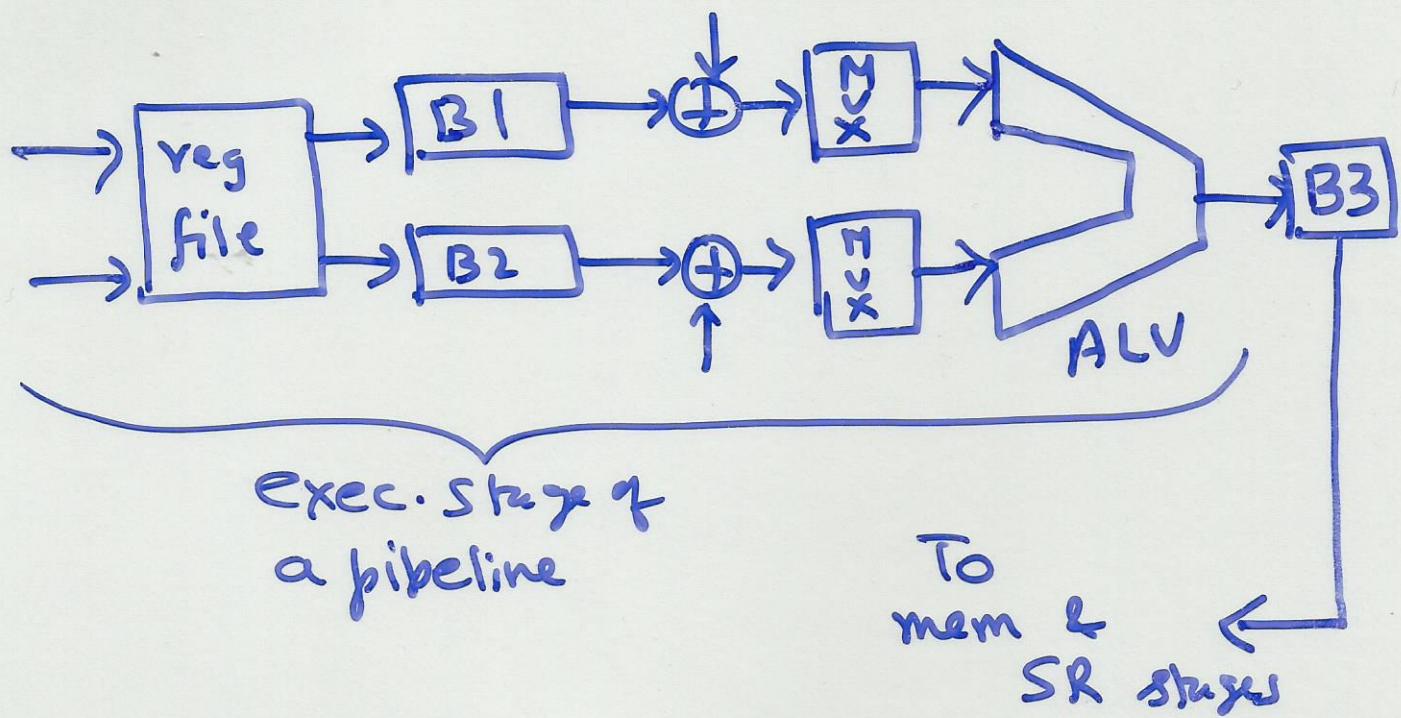
	1	2	3	4	5	6	7	8	9	10
I ₁	FI	DE	Ex	MEM	SR					
I ₂	FI	DE	X	X	Ex	MEM	SR			
I ₃	FI	X	X	DE	Ex	MEM	SR			
I ₄	X	X	FI	DE	Ex	MEM	SR			

finish time: 10 units

Here, the order (in-order) is preserved, but takes
1 cycle more!

How to avoid pipeline delays due to data dependency?

H/w Solution — register forwarding



- Key idea - instead of waiting till the SR stage, a special path is provided from B3 to ALU so that the result of the current inst. is made available to the next instr. This is referred to as reg. forwarding
- std feature in all pipelined processors

→ Calculation of timing is between
 the start of the first execution
 to the start of the last execution
 This is to eliminate the "start-up" delays
 in the pipeline.

(in your
book, ex.
fig 6.9)

Typically:

{ load/store (E stage) : 4 clk cycles
 floating pt +, x : 3 clk cycles

These are typical CISC values!

→ Pay attention to fig 6.9.

→ Reading assignment → Example 6.4.
 [MIPS R4000 IP]

Delays in pipeline due to resource constraints

Suppose if one common memory is used for both data & instructions, then :

	1	2	3	4	5	6	7	8
→ i	F	D	E	M	R/W			
i+1		F	D	E	M	R/W		
i+2			F	D	E	M	R/W	
→ i+3				F	D	E	M	R/W

Either i or $i+3$ can be executed \Rightarrow
 "pipeline stalling"

Loss of Speed-up ?

Let f : fraction of the total # of instructions
 that are floating pt instructions that
 take $(n+k)$ cycles.

m : # of instructions to be executed

n : # of stages in the pipeline.

Time to execute m instructions with no pipelining

$$\text{Time taken with no pipelining} = m(1-f) \cdot n + m \cdot f(n+k)$$

Time taken with pipelining:

$$\text{first instruction} \rightarrow n + (m-1)\{(1-f) + kf\}$$

\Rightarrow Speed-up =

$$\frac{m(1-f)n + m \cdot f(n+k)}{n + (m-1)\{(1-f) + kf\}}$$

$$= \frac{(n+kf)}{\left(\frac{n}{m}\right) + \left(\frac{m-1}{m}\right)(1-f) + kf} \quad \leftarrow$$

$$\approx \frac{n}{1-f+kf} \quad \begin{array}{l} m \gg n \\ n \gg kf \end{array}$$

So, if 10% of instructions are fl. pt. instructions, with each fl. point instr. takes 2 extra clk cycles,

$$k=2, f=0.1$$

$$\text{Speed-up} = \frac{n}{0.9 + 0.2} = \underline{\underline{0.909n}}$$

Superscalar Performance

Compare the relative performance of a superscalar processor with that of a scalar base machine. [see pg: 159 for defn.]

N : # independent instructions ~~(parallel pipelines)~~

for base scalar machine

$$T(1,1) = k + (N-1) \text{ (base cycles)}$$

$$T(m,1) = k + \frac{N-m}{m}$$

(m : # of pipelines)

k : time required to execute the first m instructions

$\left(\frac{N-m}{m}\right) \leftarrow$ remaining # of instructions,
 $m \frac{\text{instructions}}{\text{per cycle}}$ with m pipelines

$$\therefore S(m,1) = \frac{T(1,1)}{T(m,1)}$$

(Speed-up)

$$= \frac{n+k-1}{\frac{n}{m} + k-1}$$

$$= \frac{m(n+k-1)}{n+m(k-1)} \quad \leftarrow$$

\Rightarrow as $n \rightarrow \infty$, $S(m,1) \rightarrow m$, as expected

Superpipelined Design:

for a superpipelined machine of degree n
with k stages in the pipeline

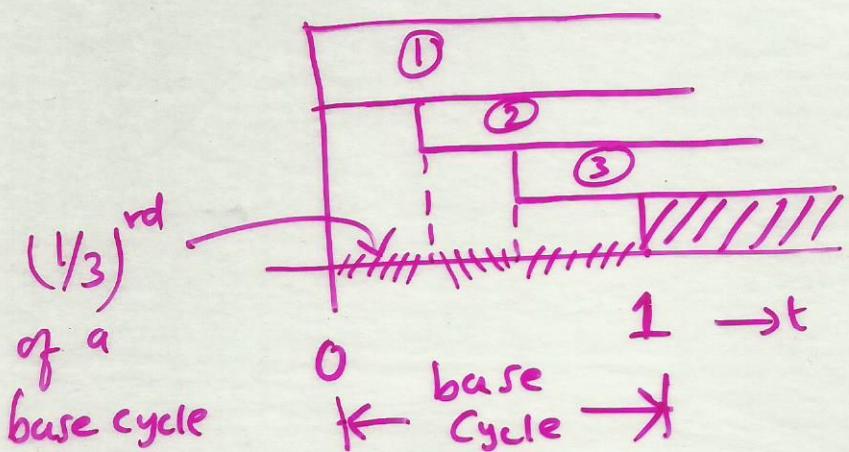
$$T(1,\underline{n}) = k + \frac{1}{n}(n-1) \quad (\text{base cycles})$$

degree n : pipeline cycle time = $(\frac{1}{n})$ of the base cycle.

Thus, $n=3$

(57)

→ an instruction is issued per $(\frac{1}{3})^{rd}$ of a base cycle



Speed-up $S(1, n) = \frac{T(1, 1)}{T(1, n)}$

$$S(1, n) = \frac{\frac{k+N-1}{n}}{k + \frac{1}{n}(n-1)} = \frac{n(k+N-1)}{nk + N - 1}$$

$$\Rightarrow S(1, n) \rightarrow n, \text{ as } n \rightarrow \infty$$

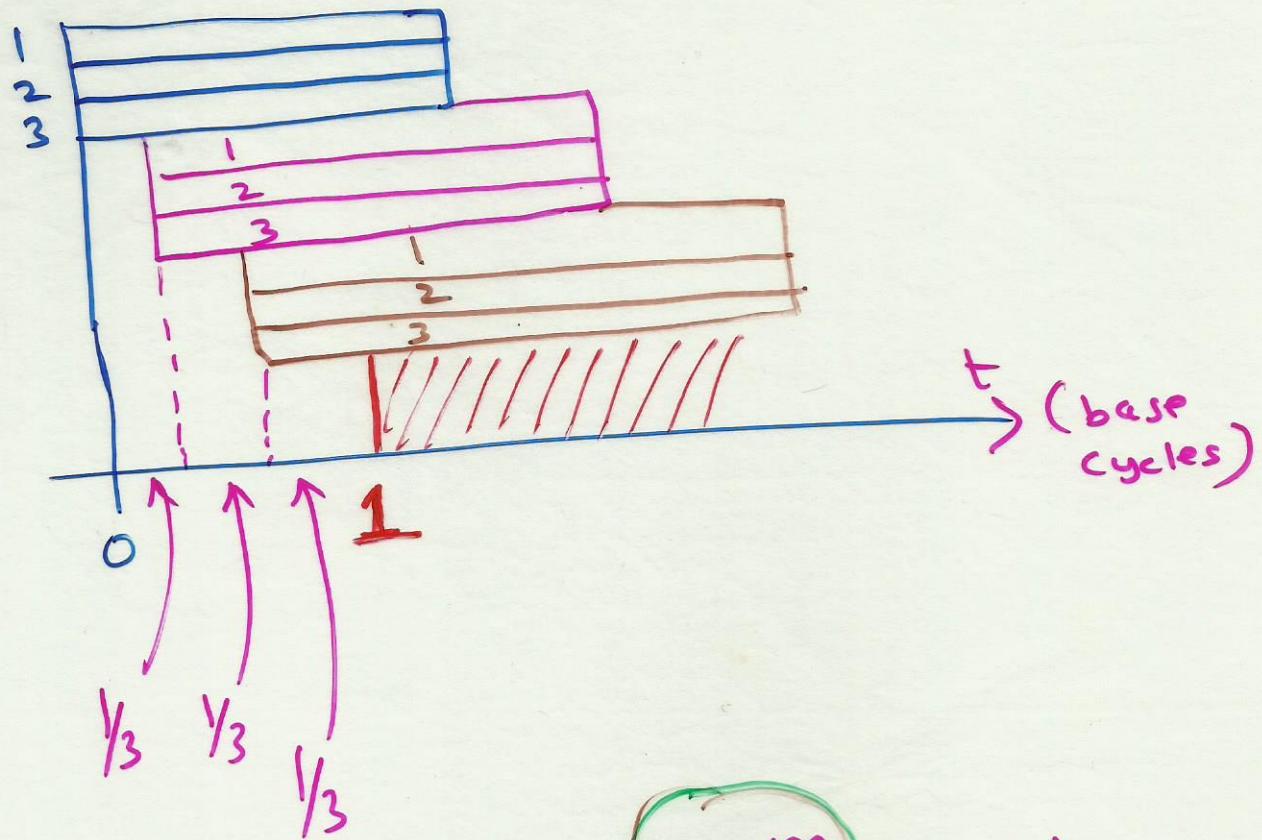
Superpipelined Superscalar Design

(58)

→ degree (m, n)

m-instructions per cycle

pipeline cycle (Y_n)
of the base cycle.



$$T(m,n) = k + \frac{N-m}{m \cdot n} \quad (\text{base cycles})$$

$$S(m,n) = \frac{T(1,1)}{T(m,n)} = \frac{m \cdot n (k+N-1)}{m \cdot n \cdot k + N-m}$$

$\Rightarrow S(m,n) \rightarrow \underline{\underline{m \cdot n}}$ as $N \rightarrow \infty$, as expected!

Performance Metrics

Other flavors of "speed-up"

I : problem instance ; P : # of processors

Ω : parallel program ; n : size of I

① Relative Speed-up

Time to solve I using program Ω and 1 processor

Time to solve I using program Ω and P processors

Comments : Purely depends on the characteristics of the instance I being solved as well as the size of P

② Real Speed-up

Time to solve I using best serial program and 1 proc.

Time to solve I using program Ω and P processors

Comments : The fastest algorithm might not be known and no single algorithm might be fastest in all instances for some applications

③ Absolute Speed up

Time to solve I using best serial program and
1 fastest processor

Time to solve I using program Q and P processors

Comments: we can also use the best known sequential algorithm most often used in practice

④ Asymptotic Speed up

Asymptotic complexity of the best serial program

Asymptotic complexity of Q using as many processors as needed.

Comments: for problems such as sorting, where the asymptotic complexity is not uniquely characterized by the instance size n , we use the worst-case complexity

⑤ Asymptotic relative Speed-up

Asymptotic complexity of θ using 1 processor

Asymptotic complexity of θ using as many processors
as needed

Comments : - None -

Ref: IEEE Parallel & Distributed Technology

"Performance Metrics: keeping the focus on
Runtime" pp. 43-56 , Spring 1996,
by Sartaj Sahni & Venkat Jhanvantri