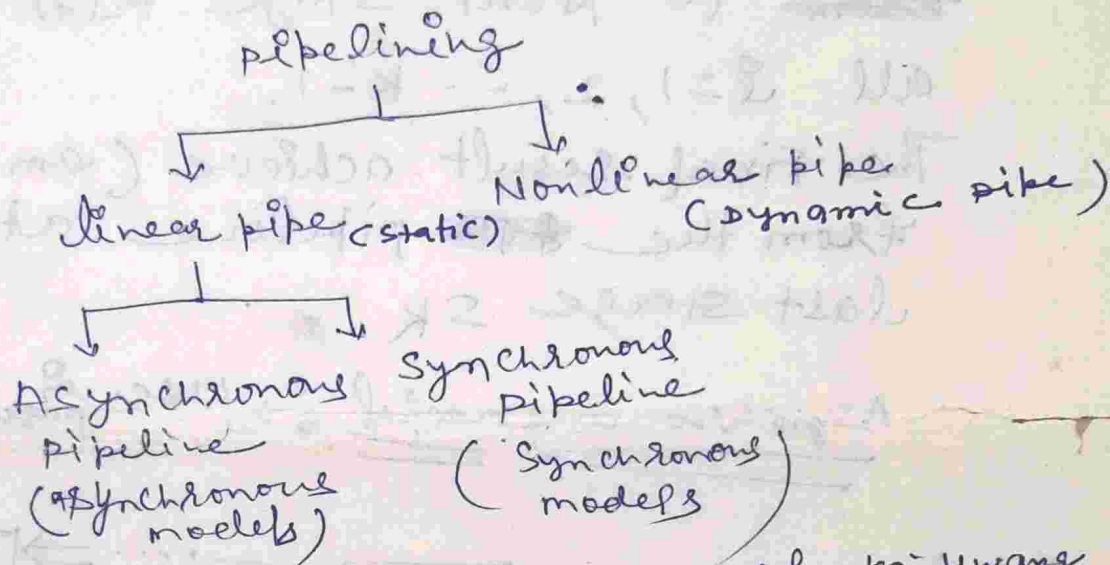


Pipelining \rightarrow cascade of the processing stages which are linearly or nonlinearly (variable function) connected to perform the operation on the stream of data, flowing from one end to other end.



Linear pipeline processors \rightarrow ref. - Kai-Hwang
Ch. Pipelining & Parallel Processing

The Linear pipeline processor is a cascade of processing stages which are linearly connected to perform a fixed function over a stream of data flowing from ~~one~~ one end to other end.

In modern computers linear pipeline are applied for instruction execution: Arithmetic Computation and Memory Access Operations

It is divided into two types. or
Asynchronous & Synchronous models

A linear pipeline processor is constructed with K processing stages.

External Inputs (operands) are fed into the pipeline at the first stage.

The Result of S_i stage is passed ~~from~~ to Next Stage S_{i+1} for all $i = 1, 2, \dots, K-1$.

The Final result achieve (emerges) from the ~~the~~ pipeline at the last stage S_K .

Asynchronous model \rightarrow (use in message passing)

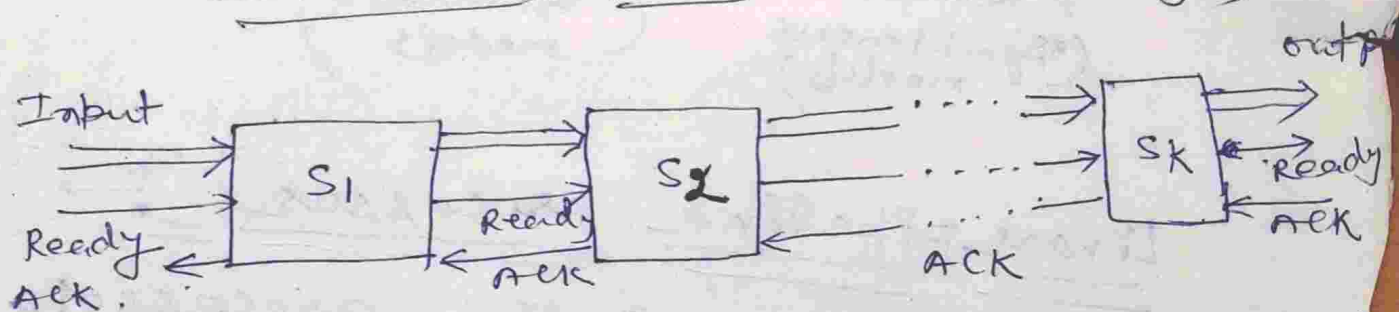


Fig: An Asynchronous pipeline model

The data flow between ~~stage~~ adjacent stages in an Asynchronous pipeline is controlled by a handshaking protocol.

When stage S_i is ready to transmit, it sends a ready signal to stage S_{i+1} . After stage S_{i+1} receives the incoming data, it returns an Acknowledge signal to S_i .

Synchronous Model \Rightarrow

In synchronous pipeline, the latch is use to store data between adjacent stages.

In synchronous pipeline, there is approximately ~~same~~ signal delay on all stages.

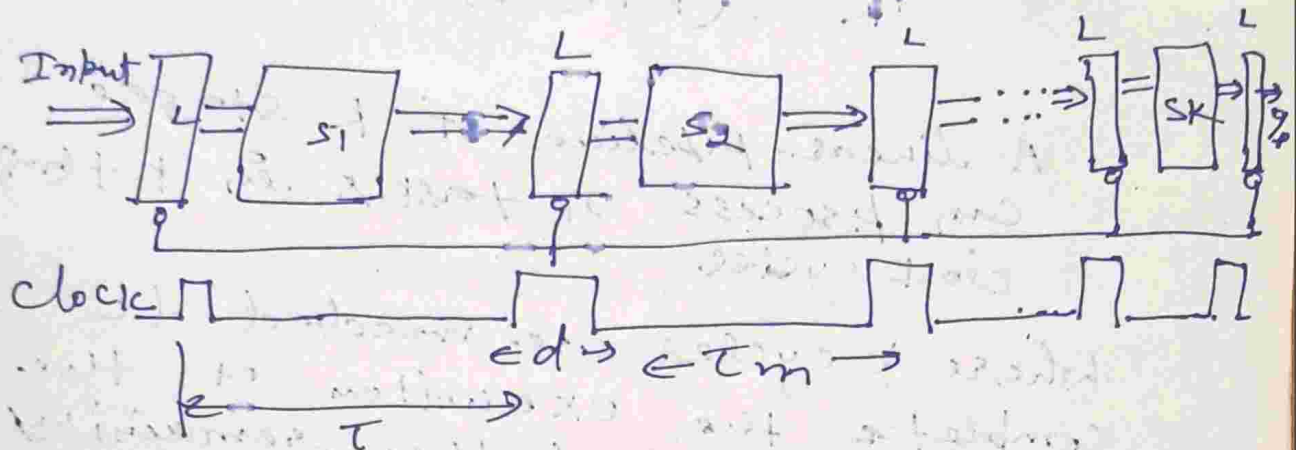


fig: A synchronous pipeline model.

S_i = stage i

L - Latch

T - clock period or clock cycle

t_m - Maximum stage delay

d - delay on Latch.

ACK - Acknowledge signal.

So we can say that all the stages are synchronize with the help of Latch.

$$T = t_m + d = \max_i \{ \tau_i \} + d$$

here - $t_m \gg d$

$$\text{Pipeline Frequency} = f = \frac{1}{T}$$

Speedup, Efficiency and Throughput

$$\text{Speedup} = \frac{\text{Non pipeline } \overset{\text{process.}}{\text{Time}}}{\text{pipeline process Time}}$$

$$\frac{T_1}{T_k} = \frac{n k T}{(k + (n-1)) T} = \frac{n k}{k + (n-1)}$$

A linear pipeline of k stages can process n tasks in $k + (n-1)$ clock cycles.

Where k cycles are needed to complete the execution of the first task and the remaining $(n-1)$ task require $(n-1)$ cycle

$$\text{so } T_k = [k + (n-1)] T$$

& Non pipelined time $T_1 = n k T$

$$\text{so - } S_k = \frac{n k}{k + (n-1)}$$

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

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

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

if $n \rightarrow \infty$

then

$$S_k \rightarrow k$$

The maximum speedup is $SK \rightarrow k$ as $n \rightarrow \infty$.
 This maximum speedup is very difficult to achieve because of data dependences between successive tasks (instructions) interrupts etc.

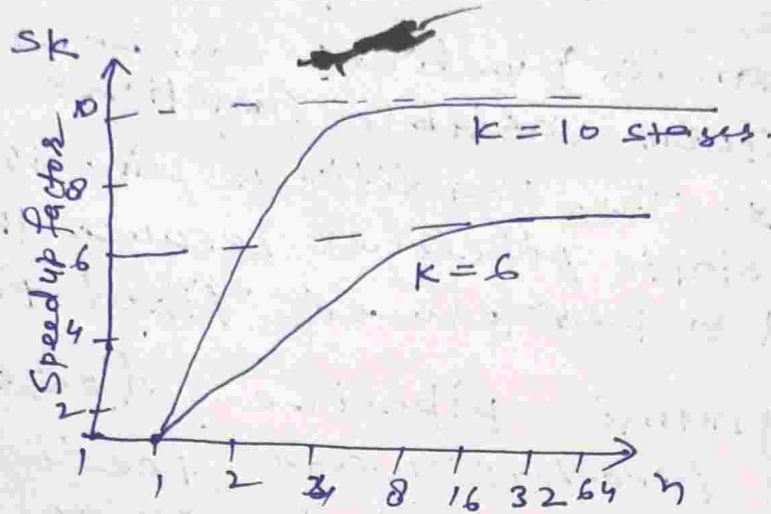


Fig: Speedup factor as a function of number of operations.

Efficiency & Throughput \Rightarrow

The efficiency E_k of a linear k stage pipeline is defined as.

$$E_k = \frac{SK}{k} = \frac{n}{k + (n-1)}$$

the efficiency approaches 1 when $n \rightarrow \infty$, and a lower bound on E_k is $1/k$ when $n=1$. The pipeline throughput H_k is defined as the number of tasks (operations) performed per unit time.

$$H_k = \frac{n}{[k + (n-1)]t} = \frac{nf}{k + (n-1)}$$

$$H_k = \frac{nf}{k + (n-1)}$$

The maximum throughput f occurs when
 $E_k \rightarrow 1$ as $n \rightarrow \infty$.

Nonlinear pipeline Processors \Rightarrow

A dynamic pipeline can be reconfigured to perform variable functions at diff^t times. The traditional linear pipelines are static pipelines because they are used to performed fixed functions.

A dynamic pipeline allows feedforward and feedbackward (feedback) connection. So dynamic ~~is~~ ^{pipeline} is known as nonlinear pipeline.

Reservation and Latency Analysis \Rightarrow

A static pipeline, it is relatively easy to partition a given function into a sequence of linearly ordered subfunction. However function partitioning in a dynamic pipeline become quite involved because the pipeline stages are interconnected with loops in addition to streamline connections.

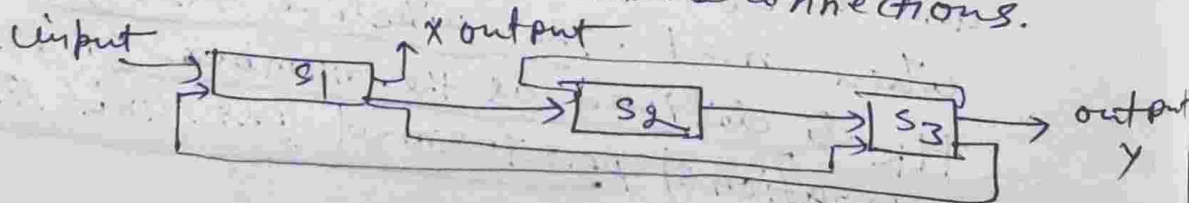


Fig: A three stage pipeline

A multifunctional dynamic pipeline is shown in Fig. This pipeline has 3 stages. Besides the streamline connection from S_1 to S_2 and from S_2 to S_3 .

there is a feed forward connection from s_1 to s_3 and two feedback connections from s_3 to s_2 and s_3 to s_1 .

These feedforward and feedback connections make the scheduling of successive events into the pipeline a nontrivial task. With these tasks the output of pipeline is not necessarily from the last stage. In fact, following different dataflow patterns one can use the same pipeline to evaluate diff'g functions.

	1	2	3	4	5	6	7	8
s_1	X					X		X
s_2		X		X				
s_3			X		X		X	

Fig. reservation table for function x

	1	2	3	4	5	6
s_1	y				y	
s_2			y			
s_3		y		y		y

Fig. reservation table for function y

✓ The reservation table for a dynamic pipeline becomes more interesting because a nonlinear pattern is followed.

✓ Given a pipeline configuration, multiple reservation tables can be generated for the evaluation of diff'g functions.

Latency cycle & latency →

Latency → The number of time units (clock cycles) between two initiations of a pipeline is the latency between them. The latency value must be non negative integer.

Latency cycle → A latency cycle is a latency sequence which repeats some subsequence

Forbidden latency → latency that ~~can~~ cause the collision, which is known as forbidden latency.

INSTRUCTION Pipeline →

A stream of instructions can be executed by a pipeline in an overlapped manner, is known as instruction pipeline.

There is diff^{'s} - diff^{'s} type of instruction pipeline

- 2-stage pipeline
- 5-stage pipeline
- Seven stage pipeline



Fig: — two stage pipe

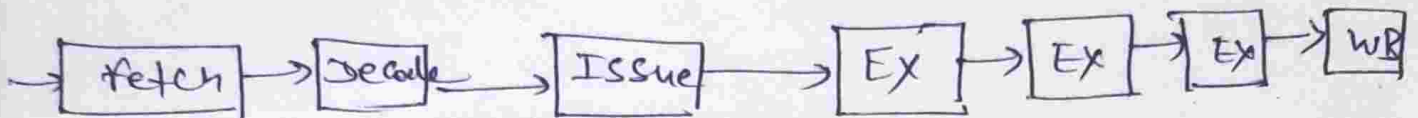


Fig: Seven stage pipeline

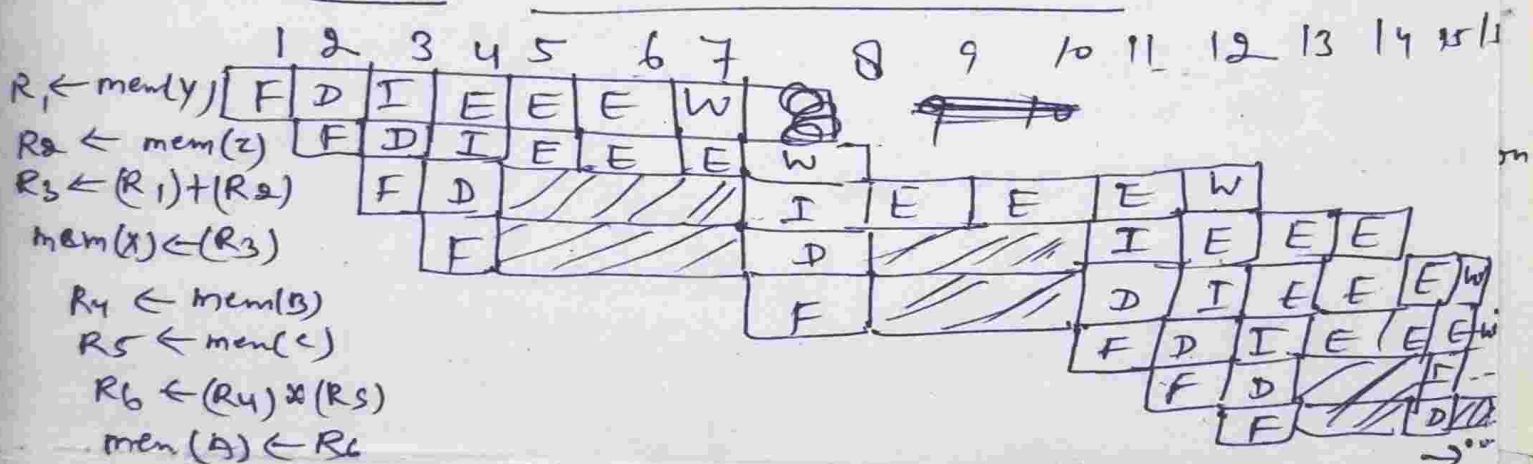


Fig: in order instruction issue

Ref: M. Morris mano / Computer system Arch

Ch. pipeline and vector processing

Page: 312

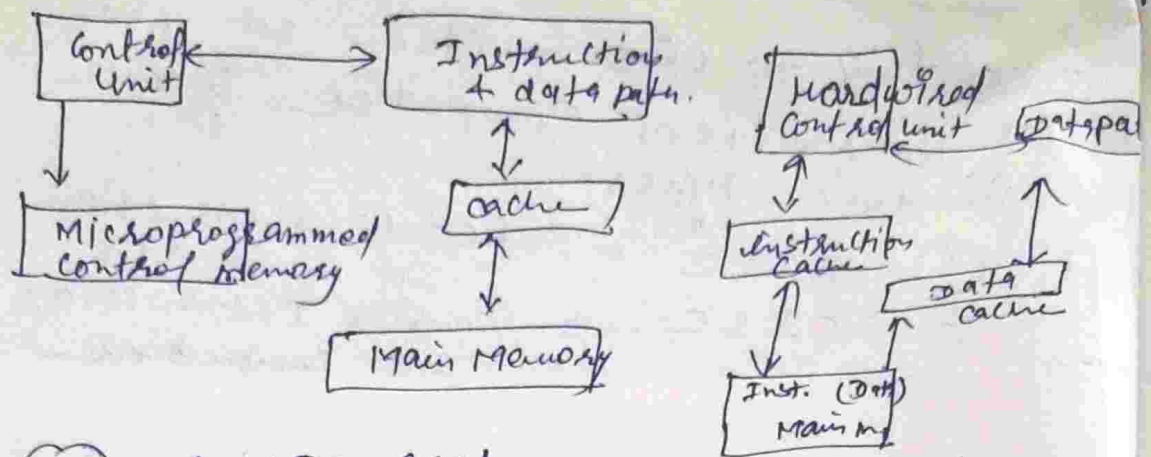
Instruction set Architecture

According to instruction set Arch.
we classify Computer system.

in two types

- ① - RISC → Reduce instruction set computer
- ② - CISC → Complex Instruction set computer.

Architectural Characteristic	CISC	RISC
1 - Instruction set size and instruction set format	Large set of instruction with variable formats (16-64 bits per instruction)	Small set of instruction with fixed (32 bit) format and most Register Based Instruction
2 - Addressing Modes.	12-24	Limited (3 to 5)
3 - General purpose registers & cache design	8-24 GPRs with a unified cache for instructions & data. Recent designs also use split caches	Large number of GPRs (32-192). There is split data & instruction cache.
4 - Requirement	lower end (low speed)	use at higher end (high speed)
5 - CPI	CPI between 2 to 15	Less CPI or $CPI < 1.5$
6 - CPU Control	Microprogramm. Controlled (use ROM)	Hardwired controlled



a. The CISC Arch.
with Microprogrammed
Control & unified cache

b. The RISC Arch.
with Hardwired
Control & split
instruction & data
Cache.

Fig: Distinctions between typical RISC & CISC processor Arch.

CISC → CISC stands for Complex instruction set computers.

A CISC Computer has 120 to 150 instructions, which is variable length.

- It use ~~small~~ less GPRs (8 to 24)
- There is no split cache.
- They ~~has~~ have ~~low~~ High CPI

RISC → RISC stands for Reduce instruction set computers.

A RISC Computer has less than

100 instruction. with fixed length (32 bit)

- It use large number of GPRs (32 to 192)
- There is split cache
- They have less CPI

Ref: Kai Hwang / RISC, CISC

$CPI < 1.5$

Arithmetic Pipeline Design \Rightarrow

pipeline Arithmetic units are usually found in very high speed computers. They are used to implement floating point operations, multiplication of fixed point numbers.

ex- The ~~inputs~~ input of floating point shown in below, is Normalised floating point binary number.

$$\begin{aligned} X &= A \times 2^a \\ Y &= \frac{B}{2} \times 2^b \end{aligned} \quad \begin{array}{l} \text{exponents} \\ \text{mantissa.} \end{array}$$

[A and B are two fractions that represent the mantissa and a and b are the exponent.]

the floating point addition & subtraction can be performed in four segments.

- ① - compare the exponents.
- ② - Align the mantissa.
Add or subtract the mantissa.
- ③ - normalise the result.
- ④ -

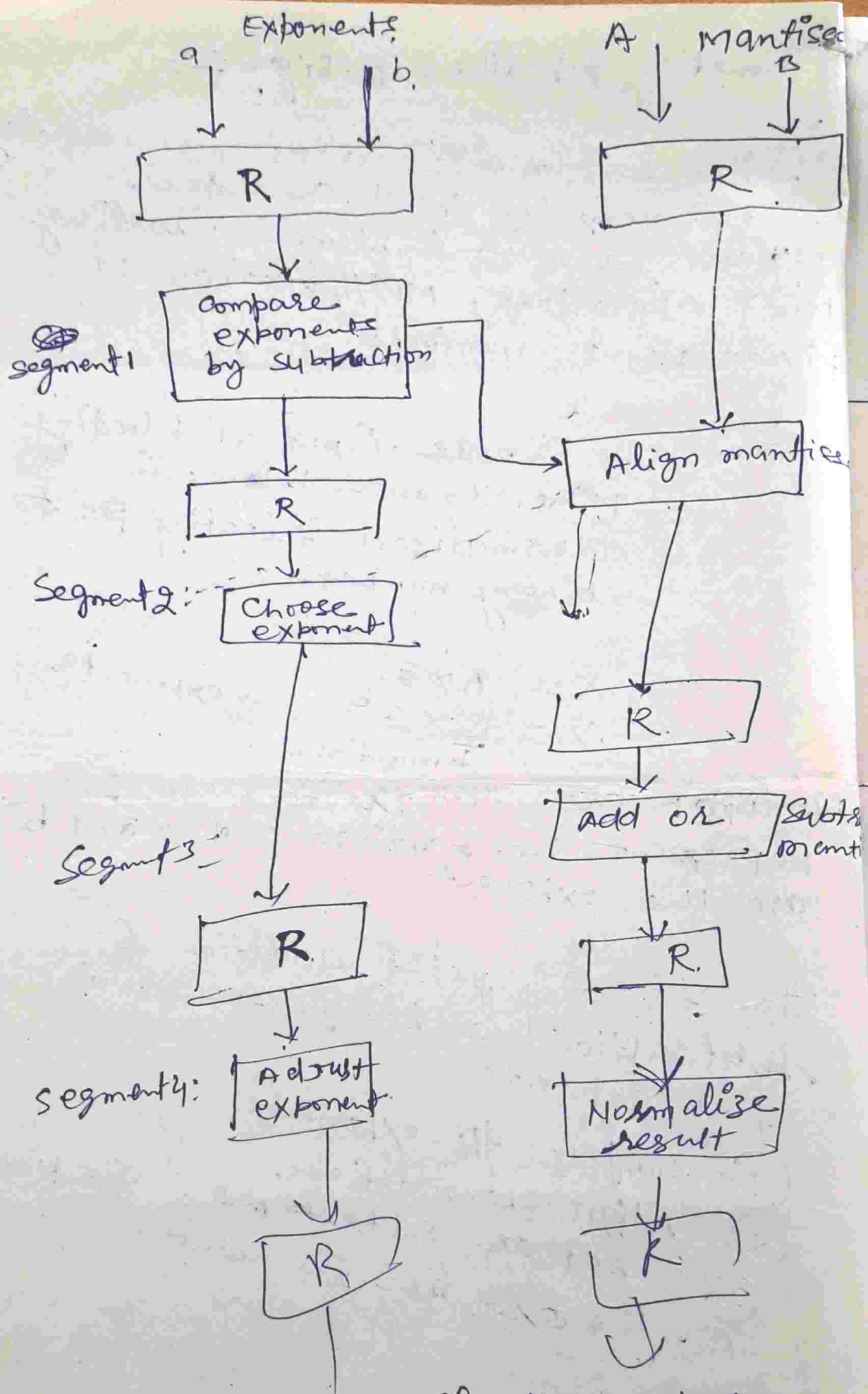


Fig. → pipeline for floating point addition & subtraction

$$X = \underbrace{0.9504}_{\text{mantissa}} \times 10^3$$

$$Y = 0.8200 \times 10^2$$

→ two exponents are subtracted in the first segment to obtain

$$3 - 2 = 1$$

The large exponent. 3 is chosen as the exponent of the result (segment 2)

The next segment shift the mantissa of Y to the right to obtain

$$X = 0.9504 \times 10^3$$

$$Y = 0.0820 \times 10^3$$

Align mantissa

one right shift

This align the two mantissa under the same exponents.

The addition of the two mantissas in segment 3 produces the sum.

$$Z = 1.0324 \times 10^3$$

The sum is adjusted by normalizing the result so that it has a fraction with a nonzero first digit.

This is done by shifting the mantissa once to the right and increment the exponent by one to obtain the normalized sum

$$Z = 0.10324 \times 10^4$$

Ref: Computer system Arch-

M. Morris mano page- 303