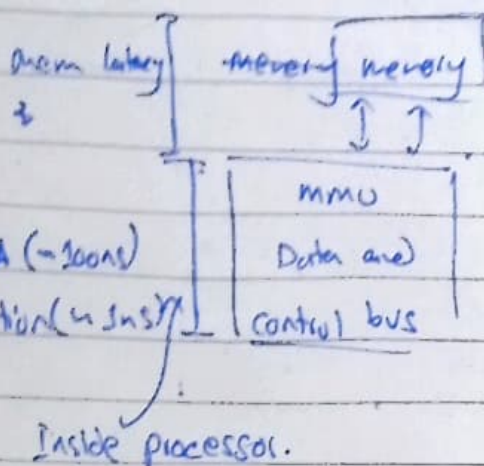K200353 Mohsin Ali Mirza

## ① Implicit Parallelism

Instructions {

RD, R1, [2000H]

RD R2, [3000H]
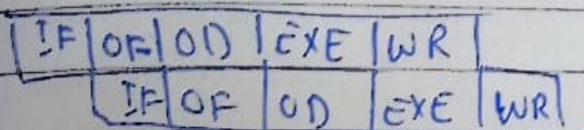
Add R3, R1,10 → Memory MA (~200ns)

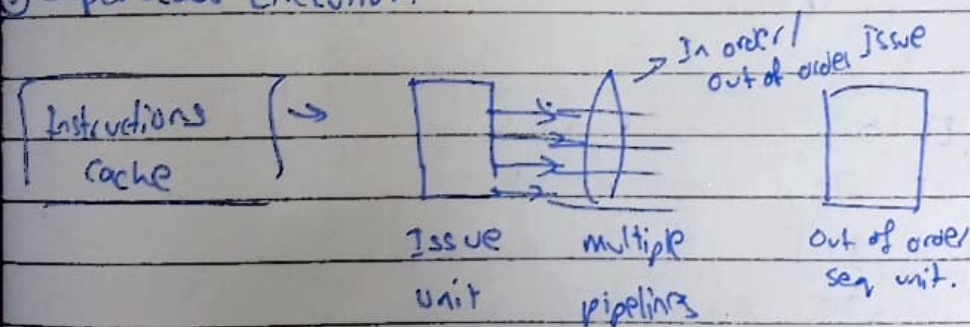Add R3, R2, R3 → Alu Execution (~ Insn)

mem latency    memory   memory
↕                        ↕  ↕

mmu
Data and
control bus

Inside processor.

## ② Pipelining (

| IF | OF | OD | EXE | WR |
| | IF | OF | OD | EXE | WR |

## ③ Super Scalar Execution

Instructions
Cache

→ In order/
out of order    Issue

Issue      multiple     Out of order
unit       pipelines    seq unit.

## ④ VLIW (Very Large Instructions word).

| Ins1 | Ins2 | Ins3 | . . . . |

→ multiple instructions
↳ Block of instructions

# 5) Cache Line

↳ Prefetching. four data is fetched and loaded in Cache in 1 access.



data & control bus.

Memory

cache line

First word :- 10ns

other words
$$\begin{cases} 1\,ns \\ 2\,ns \\ 1\,ns \end{cases}$$

Execution Schedule for code fragments for super scalar pipe

① load R1, @1000
② load R2, @1008
③ add R1, @1004

④ add R2, @100C

⑤ add R1, R2

⑥ Store R1, @2000

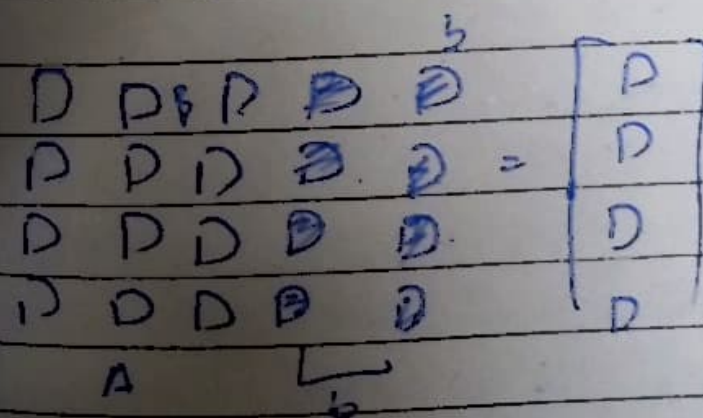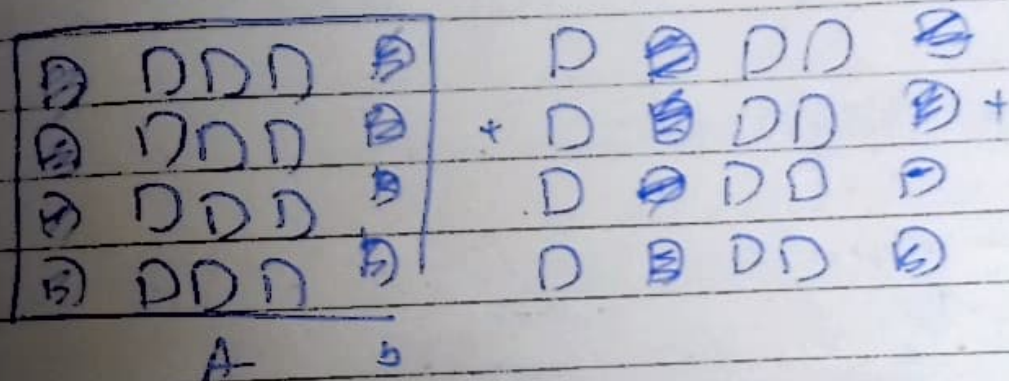| IF | ID | OF | | | | | | load R1, @1000 |
|----|----|----|---|---|---|---|---|---|
| IF | ID | OF | | | | | | load R2, @1008 |
| | IF | ID | OF | E | | | | add R1, @1004 |
| | IF | ID | OF | E | | | | add R2, @100C |
| | | IF | ID | NA | E | | | add R1, R2 |
| | | | IF | ID | NA | WB | | Store R1, @2000 |

# Hardware Utilization trace for schedule

Full issue slots

Empty issue slots

Horizontal waste

Vertical waste

Column Majority data access

$$D \quad D D D \quad B$$

(matrix diagrams with D symbols and shaded boxes)

A     b

$$D \quad D\&R \quad B \quad B = \begin{bmatrix} D \\ D \\ D \\ D \end{bmatrix}$$
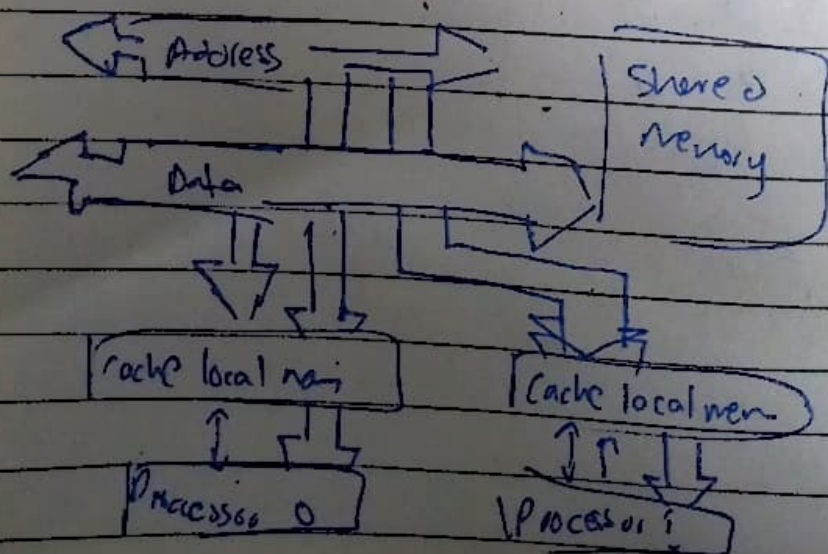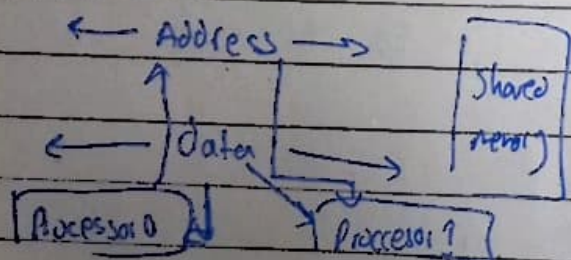
A     b

# Row majority data access?



## Network Topologies: Buses

Cross bars



Switching element.

Ex 2.2  **Effect of memory latency on performance**

Consider a processor operating at 1 GHz (1ns clock) connected to a DRAM with a latency of 100ns(no caches). Assume that processor has 2 mul add units and is capable of executing 4 instructions in each cycle of 1ns. The peak processor = $\frac{4}{1ns}$ = 4GFlops.

Since memory latency is equal to 2 cycles and block size is 1 word, everytime a memory request is made, the processor must wait 100 cycles before it can process the data. dot product of two vector pbran. It is easy to see the ratio of computation is limited to 1 flo FlOP every 100ns time or speed of 10 Mflops.

A very small fraction of the peak processor rating. The example highlights the need for ems effective memory system performance in achieving high computation rate.

Ex 2.3 Impact of cache on memory system performance

Consider 1 GHz processor with 100ns latency DRAM. In this case we introduce a cache of size 32kb with a latency of 1 ns or we use this setup to multiply 2 matrices A & B, of dimensions 32×32. we once again we assume an ideal cache placement strategy in which none of data items are over written by others fetching the two matrices into the cache corresponds to fetching 2k words which takes approx 2μs. Multiplying n*n takes $2n^3$ operations which can be performed at 16k cycle. The total computation is 216μs as peak processor computation of 64/216 = 303 MFlops.

Ex 2.4 Dot product of 2 vectors a memory system with a single cycle cache 100 cycle latency DRAM with the processor operating at 1 GHz. If the block size is 1 word, the processor takes 100 cycle to fetch each word. The dot product performs are mul-add ie 2 Flops therefore the algo performs are Flops every 100 cycles for a peak speed of 20 M Flops as estimated. Now let us consider what happens if the block size is increased to too words i.e a processor can fetch a 4-word cache i.e every 100 cycle, Assuming that vectors are laid out linearly in every: 4 mul-adds can be performed in 0 cycles. This is because a single memory access fetches 4 elements of vector. This corresponds to a flop every 25ns