

IDE : Integrated Dev. Environment

file a.out → description

objDump -D a.out | more

classmate

Date _____

Page _____

• Basic steps :

- write a file sort.c & save
- saved on a persistent [non-volatile] storage
- The hierarchical arrangement of these files (stored) : File system

Eg of File systems :- NTFS, FAT32, FAT16, EXT2

→ Compiled Task of the compiler, Read binary encoded data.

(ii) Executable file formats :- ELF, COFF

(S)

- On running objDump -D a.out | more ⚡
Details of the executable are displayed.
- with operations & operands.
col 3 col 4.

To communicate with the processor the only way is through ISA: Instruction set architecture

Col. 2: Represents the operation instruction in hexadecimal representation.

→ x86 : ISA (proposed by INTEL)

BSBC, MIPS, ARM, RISC V

open source. ISA
processors can be built.

Intel i5 is a processor such that the ISA corresponding to x86 can be understood with the same.

If to design a processor:-

i) Speed

iii) Power

iii) Cost

iv) Size

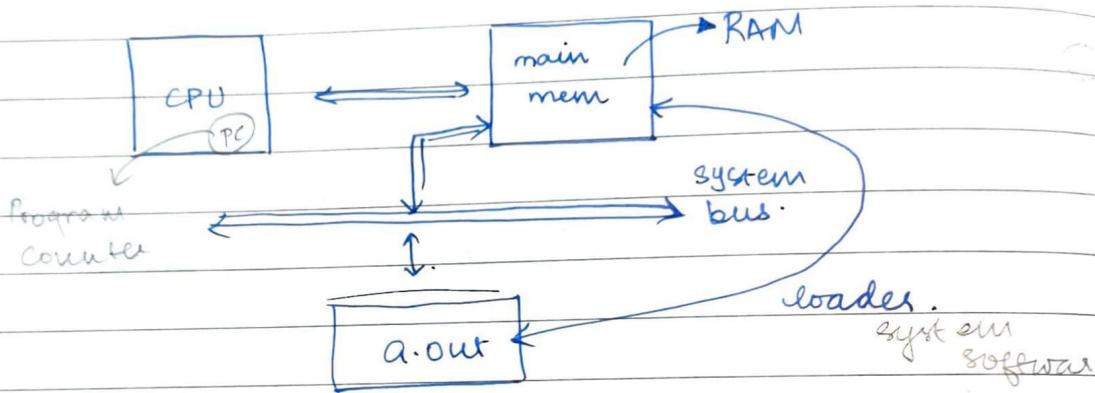
ISA is what the chip
understand.

through.

ISA

- ↳ Variable length encoding (used in CISC) x86
- ↳ Fixed length encoding (ARM)

- Advantage of variable length :-
code length can be adjusted acc. to req.
- Disadvantage :-
Decoding becomes complicated (appropriate machinery & equipment is required).



Main-memory : linear array of bytes. (memory)
It is only byte addressable.

1st col of * it ^{is} assigns the address of the cell where that particular instruction sits.

→ Execution Model.

Step: 1> Fetches instruction (PC)

PC: contains address of next command to be executed.

2> Instruction decoded

3> Execute the instruction.

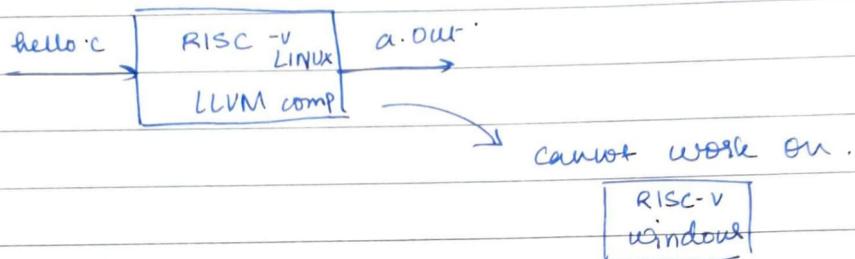
→ ltrace ./a.out .

↳ gives details of all library calls.

Hardware exposes its ISA for the OS .

OS exposes its system call interface for application software .

Library call → User space .
 System call → Kernel space .



Because the system call interface varies .

→ Native compilation :-

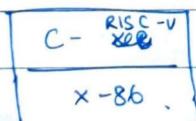
C pro - x86	compilation & run on the same machine
x86 ISA	

→ Cross compilation :-

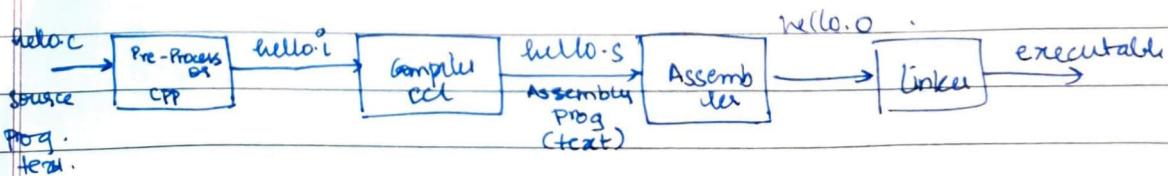
Eg arduino:

Compiled on one machine run on another .

Eg: Developing a software for a sensor node .



• Software interface :



way assembly lang. is an intermediate.

- ① for efficiency & performance.
- ② Most code for an OS can be written in a programming lang. however, sl. of the code will be tied to the processor, (its memory management etc.) ∴ it is processor dependent & you may need to write it in Assembly lang.

Solving a computational problem :-

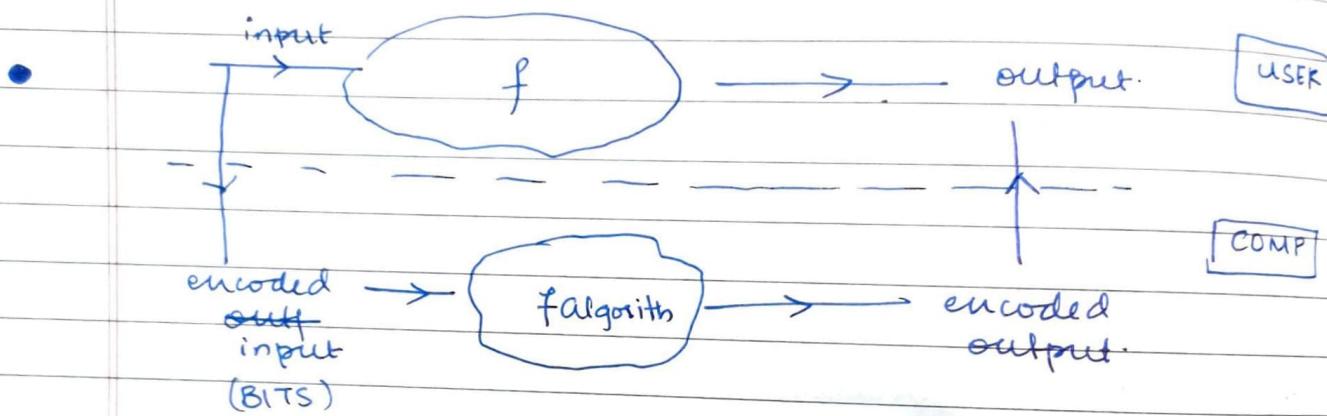
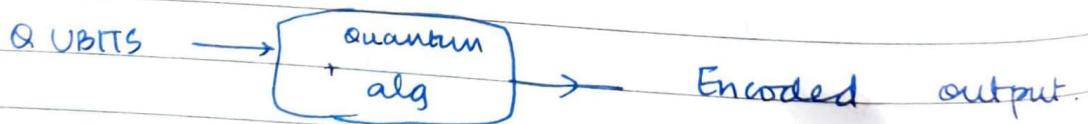
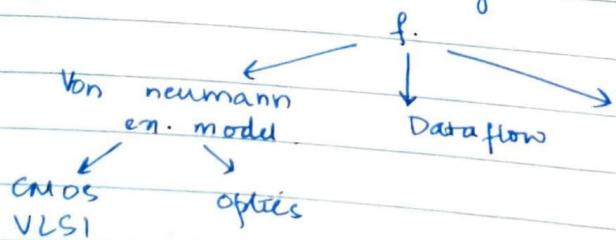


Fig. Execution model.

Quantum execution model.



Physical Realization of ex. models.



otherwise computing is useless

classmate

Date _____

Page _____

Optimizing encoding & decoding processes is essential.

- Objects can be classified as:-

- Primitive
 - Integers.
 - Real nos.
 - Alphabets.
- Composite
 - Rational nos.
 - Complex
 - Matrix
 - Images, video.

Space allocated to different data types varies from machine to machine.

1 bit / 3-bit / 13 bit data types are unheard of because, it is not possible for the processor to represent multiple wide range of nos.

- ~ Unsigned rep

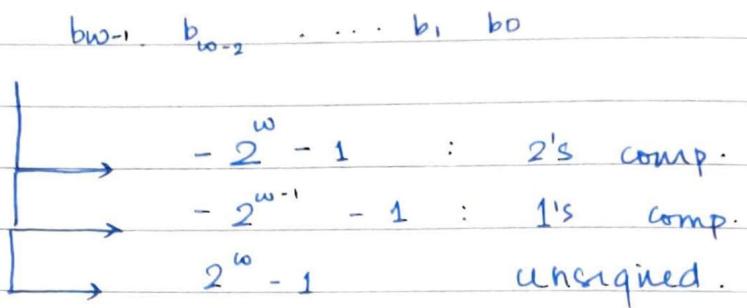
$b_n \ b_{n-1} \dots \ b_2 \ b_1 \ b_0$
represents mag. 2^{b_n} .

- ~ Signed rep.

$b_n \dots b_2 \ b_1 \ b_0$
if 1 → -ve no: } signed magnitude
0 → +ve no. rep.

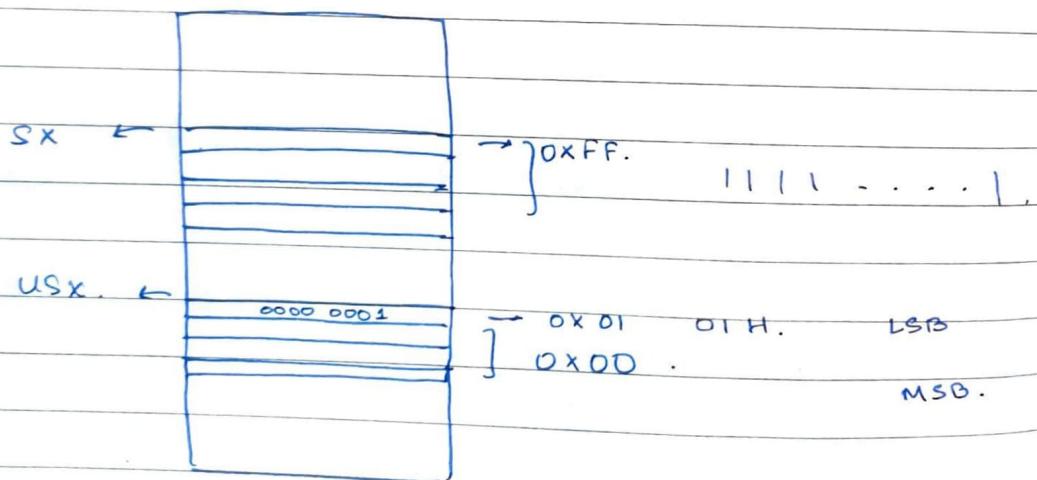
U	SM	1's	2's
0 0 0	0	0	0
0 0 1	1	1	1
0 1 0	2	2	2
0 1 1	3	3	3
1 0 0	4	0	-3
1 0 1	5	-1	-2
1 1 0	6	-2	-1
1 1 1	7	-3	0

for a w bit no:



: Rules of binary arithmetic are the same,
the same ~~*~~ circuit (seq.) is invoked for add
Operating on signed or ~~*~~ unsigned numbers

e.g.: $Sx = -1$.
 $USx = 1$.



Overflow:

Overflow can be determined using a carry bit
that is the final output carry.

Z bit is 1 - when result of arithmetic op
is 0. else it is 0.

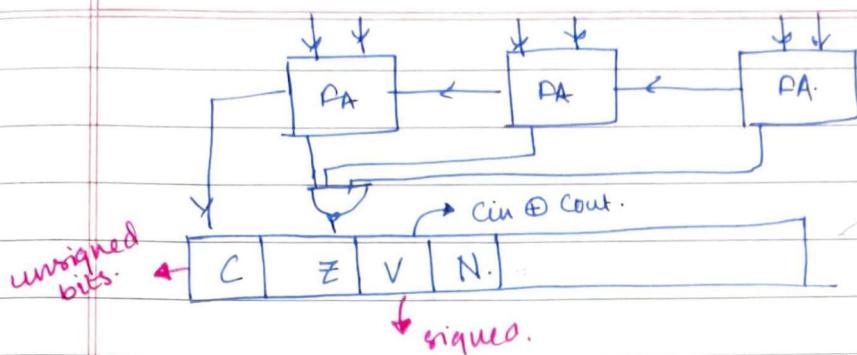
For signed bit nos check $\rightarrow V$.

For unsigned $\rightarrow C$

classmate

Date _____

Page _____



REGISTER Bits
in the process

Case:

$$\begin{array}{r} -1 \\ -3 \\ \hline 110 \\ (1) 001 \end{array}$$

This isn't an overflow
as closure property isn't
violated.

signed
rep

$$\begin{array}{r} 3 & 011 \\ 2 & 010 \\ \hline 5 & 101 \\ (-8) & \text{---} \\ -3 & \end{array}$$

This is a case of overflow.

(Rep)

C - carry out [final].

Z - Zero bit [checks if result is 0].

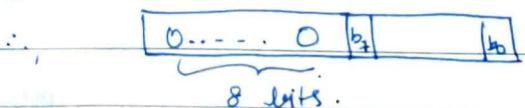
V - Cin \oplus Cout.

→ If the addends are from 2 different domains of bits [Eq: 16 bit no + 8 bit no].



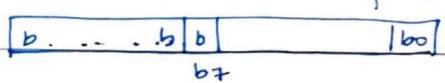
For unsigned nos:

zero-extension solves the problem.

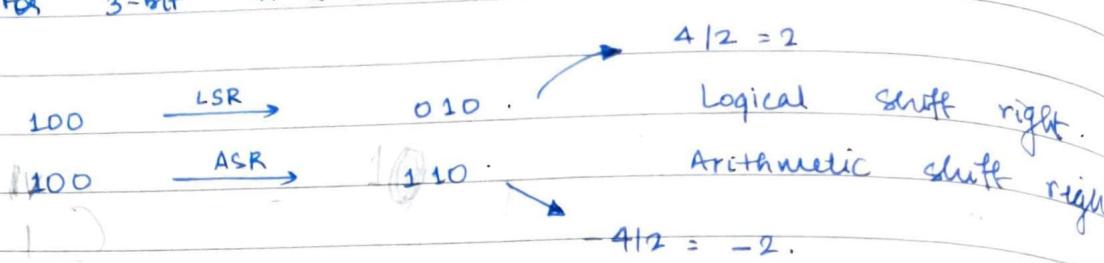


For signed nos:

Signed extension solves the problem.



→ For 3-bit R & L shift



- Using gdb:

→ `p & varname`

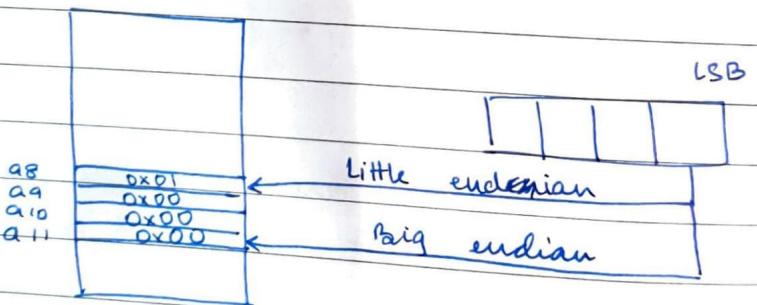
prints address of varname

→ `mvarname 7/4b 0xFFFF...ac`

Shows value stored in the 4 bytes starting from 0xFFFF...ac

→ `p/x 0xFFFF....`

→ `p/x * address . → value in address`



- Programming lang. provides:

1) Primitive data type [ISA]
2) compiler (bridge)

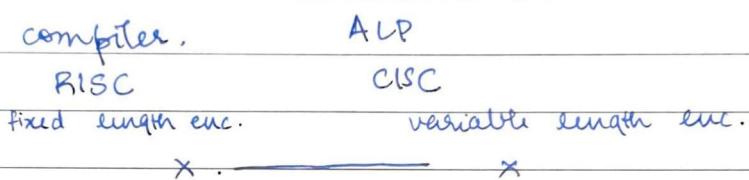
2) Mechanisms to compose new data types
[arrays, structs, unions]

→ Programming lang. must support arithmetic & logical op.

- Should support control flow constructs
 - ↳ if-else
 - ↳ while, for.
- Memory transfer.
- Turing machines.

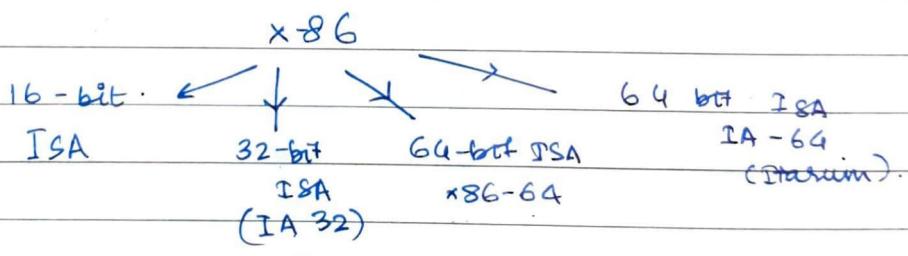
For a lang. to be turing complete :-

- the ISA should have instructions for every line of code. [for every possible computational program]
- ease of programmability



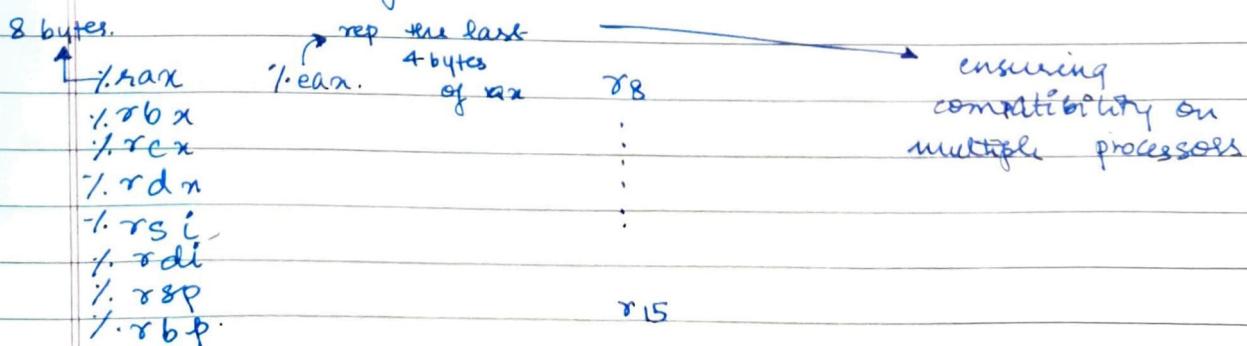
11

• x-86 Instruction set Arch.

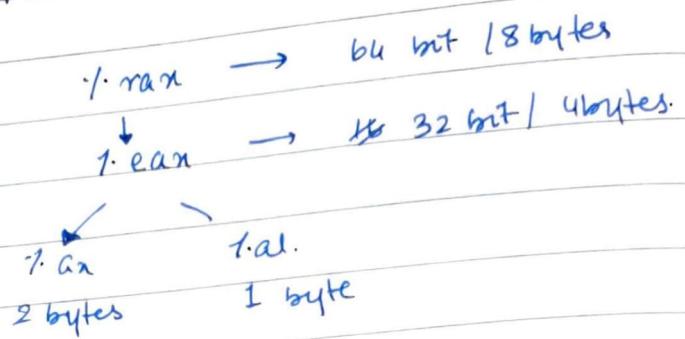


→ x86 - 64.

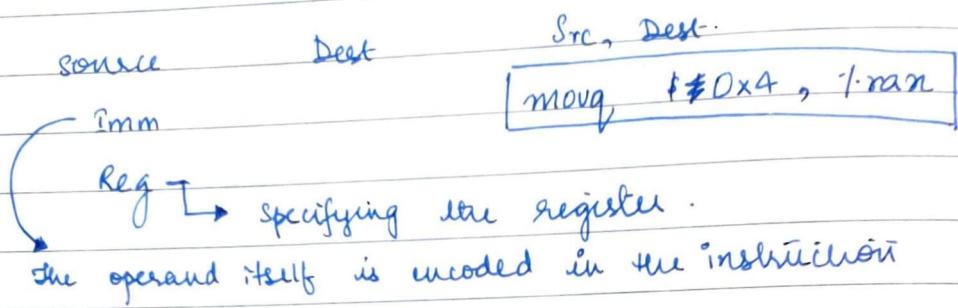
Total 16 registers



Split of .l.ean.



Layout of Instructions:



	Source	Dest	Source, Dest
①	Imm.	{ Reg Mem.	<code>movq \$0x4, %rax</code> <code>movq \$-147, (%rax)</code>
			0x4 moved to rax temp: ↓ 2's complement stored in address stored in %rax.
②	Reg	{ Reg Mem.	<code>movq -1.rax, %rdx</code> <code>movq 1.rax, (%rdx)</code>
			temp=mem ↓ complement=mem
③	Memory	Reg	<code>movq (%rax), %rdx</code>

* Movement from one location in memory to another location in memory is not done
 ∵ not cost-effective.

Ques:

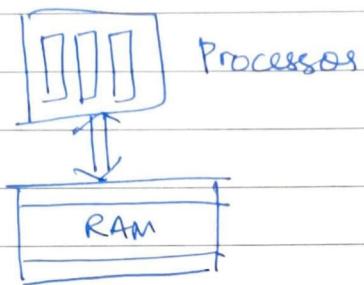
`movq 8(%ebp), %rdn.`

from stack
memory add
in %ebp
offset 8

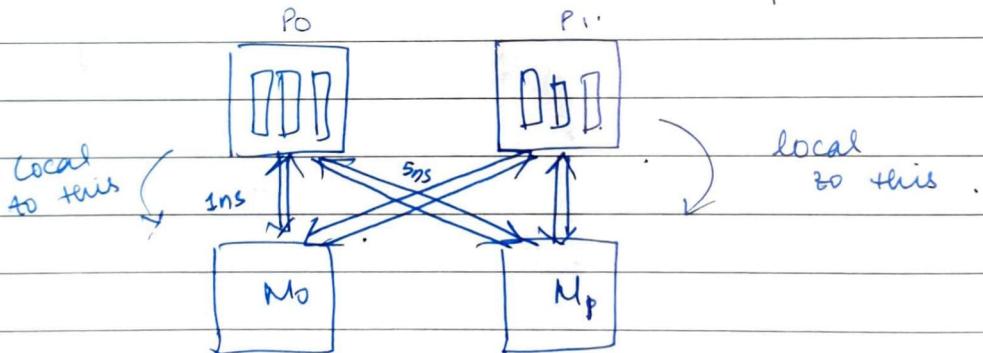
Instruction classification :-

1. Data Transfer.
 2. Arithmetic & logic
 3. Control transfer.
- } - control flow.
 } - data transfer.
 } - logical/arithmetic

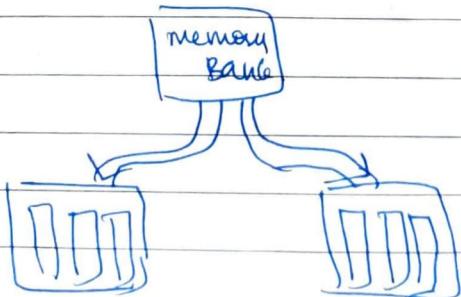
UMA : Uniform memory access (time)



NUMA : Non-uniform memory access.

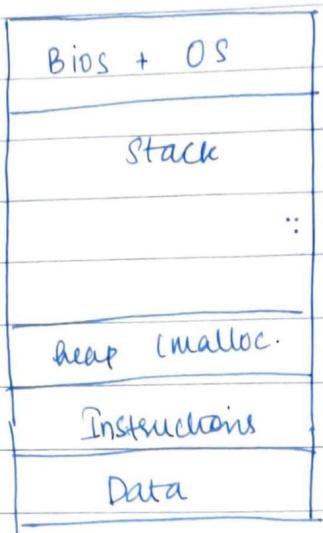


Case :-



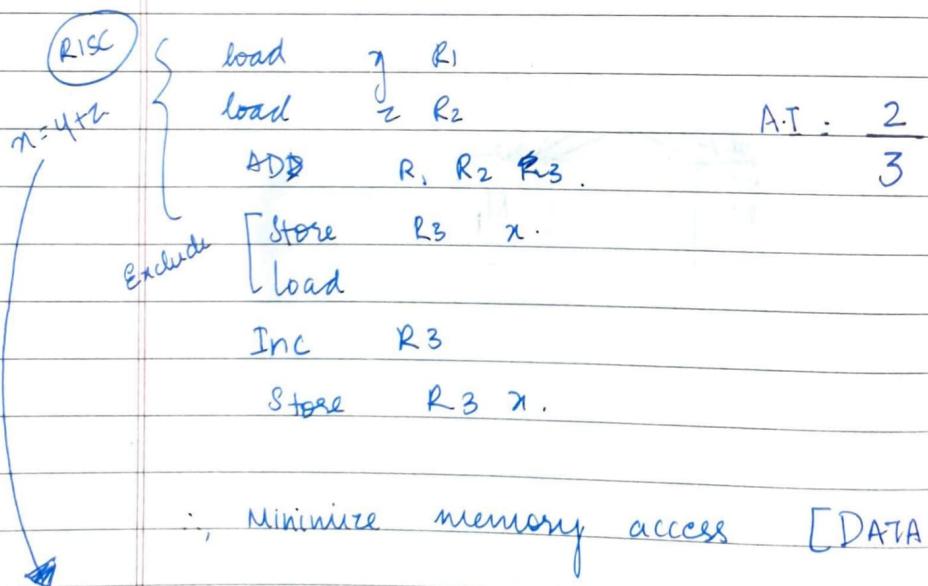
Memory bandwidth reduces since both processors access the same bank (pool) of data

Main memory layout :-



$$\text{Arithmetical intensity} = \frac{\# \text{ computations}}{\# \text{ memory access}}$$

Goal: Maximize A-Intensity so as to keep the processor busy and doing work

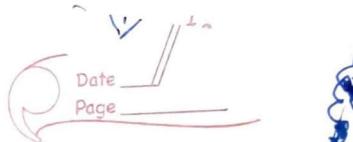


CISC:

LOAD $y \leftarrow R_1$

ADD $z = R_1$

STORE $R_1 \leftarrow x$



→ x86 instruction set.

• Data transfers.

MOV :- used for both LOAD & STORE.

Structure of writing the instructions :-

$D(R_b, R_i, S)$

displacement index reg scale factor
base reg.

(1.rax) → normal.

8 (1.rax) → 8 + (1.rax)

-16 (1.rax, 1.rbx) → -16 + rax + rbx.

32 (1.ram, -1.ran, 4)

Analogous :- $A[x]$.

To access 4th el.

Base 4 × size of index.

⇒ (1.rax → 1.ran, 4).

e.g.: $P = \&x[i];$

$y = *P;$

$P \Leftrightarrow 1.rax$ $i \Leftrightarrow 1.rcx.$

$X \Leftrightarrow 1.rbx.$ $\&y \Leftrightarrow 1.rdx.$

calculates instruction

stores it.

Instr:- load (1.rbx, 1.rcx, 4) 1.rax.

movd (1.rax), 1.rwd

1.rwd, (1.rcx)

Flags are not set by this instruction.

cannot read 2 addresses from memory

similar

Date _____
Page _____

Flags analogy (86 & ARM)

for :	signed	unsigned	
(V)	OV overflow	carry flag	[C]
(N)	SF sign	Zero	[Z]

x86	ARM
ADDQ	ADDS
: changes the flags.	: Arithmetic. ADDS → operates on old flags

→ CMP	a b, a	→ computes	a - b
	unsigned	- - signed	
CF = 1	$\Leftrightarrow a < b$	$a < b$	
ZF = 1	$\Leftrightarrow a == b$	$\Leftrightarrow a == b$	
SF = 1	$\Leftrightarrow a < b$	$\Leftrightarrow \text{Overflow} \quad a < b$	
OF = 1			

OF flag is set if 2's comp
 $(a > 0 \wedge b < 0 \wedge (a-b) < 0) \vee (a < 0 \wedge b > 0 \wedge (a-b) > 0)$

• claim : $a \leq b \Leftrightarrow (SF \wedge OF) \mid ZF$

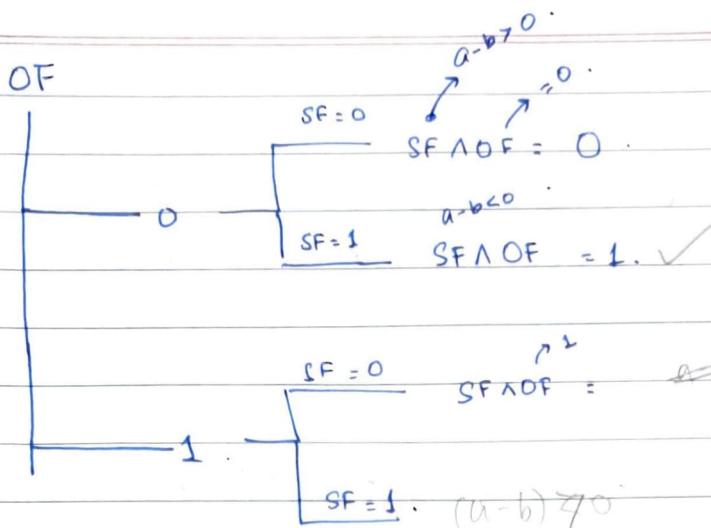
$$a \leq b \quad SF = 1 \quad ZF = 0$$

$$\rightarrow \text{equality} : - \quad a == b \Rightarrow ZF = 1 \\ \therefore ZF = 1 \Rightarrow A == B.$$

$$\rightarrow a < b \quad SF = 1 \quad ZF = 0$$

Computer Science

CLASSMATE
Date _____
Page _____



Set X instruction : sets if satisfied
clears if not satisfied]



Basic Layout:

Data { Data section : .data
N: .nyc 10
sum sum: .long

text [.text
.global main

Squaring:

.text

```
movl %ebx %ecx  
imul %ecx %ecx } function like  
ret
```

run , tui ./a.out , layout , regs run

will enable

run

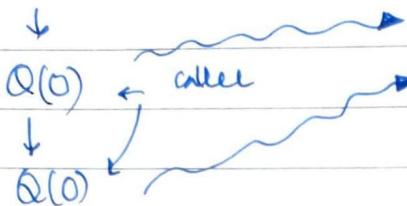
Show reg

classmate
Date _____
Page _____

Control-flow: (Function calls)

local variables cannot be declared in the data section. In case of recursive calls this won't work.

Eg: P() ... caller



On stack
at each
activation

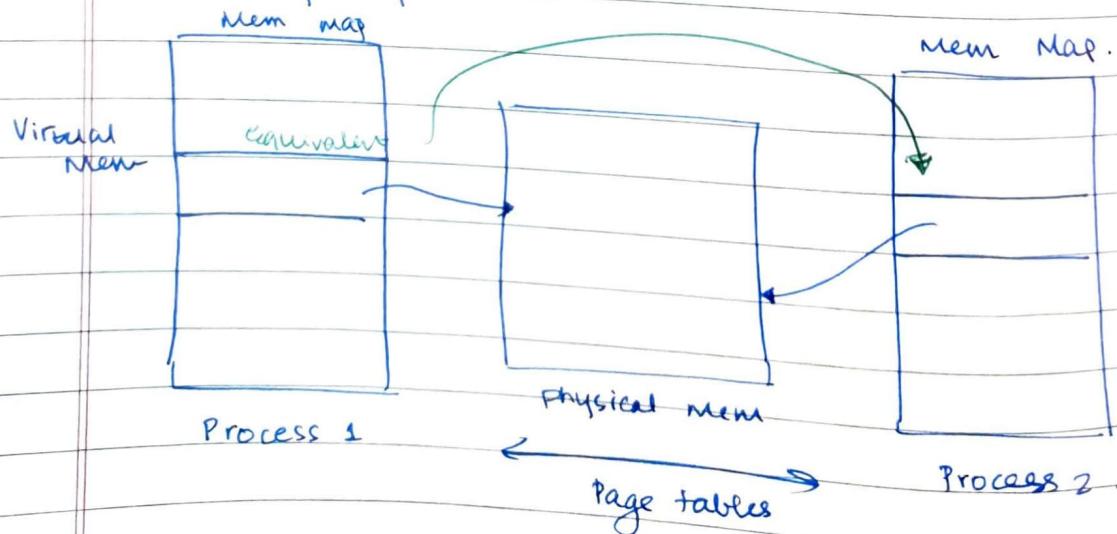
→ Scope:- within which range the variable can be used.
Life:- How long the variable lives.

→ Memory Image Map slide

In case of so many recursive calls, at times the memory exceeds the space that is accessible.

Eq: It may enter a memory-mapped region for shared lib & gives seg-fault ∵ READ-only

→ For multiple processes :-



0xffffffff dead

0x0000000000400485

Hello
classmate

Date _____

Page _____

PROCEDURES & FUNCTIONS:-

→ TPA Primitives

Push Call
Pop Ret

On using multiple abstractions there could be

- performance overheads.
- mitigates the performance.

- Handle control transfer
- Passing parameters

address

~~Global var~~ → always has the same add on repeated runs .. stored in data seq.

Local variables :- Stored in the stack, stack start is randomized.

[This is done to avoid buffer overflow-attack]
→ read.

Passing parameters:

1. calling library

2. x86 → calls square() in file.c.
gcc compiler understands call square
as the parameter is in %rdi

3. conversion of file.c → file.s
accepting arg

* In a jmp instruction
Jumps to a value in memory

Date _____
Page _____

→ Interrupts are asynchronous.

It is necessary to ensure that the return address is at the top of the stack when ret inst. executed.

ABI: Application Binary interface.

Default req for variable parameters:- (6)
%rdi, %rsi, %rdx, %rcx, %r8, %r9

Passing 6+ parameters uses stack so also reduces performance ∵ arithmetic intensity red.

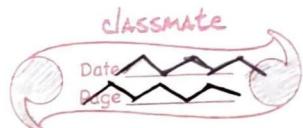
Eg: of a parameter of a struct, accessing multiple values with pointer.

• Stack

- keeps track of the functions running in a program.
- when a function is called, a stack frame is created supporting the function's exec. activation record. This includes local variables & parameters.
 - + info for callee to return to caller.
- Grows towards lower memory.
- esp points to top.
- Housekeeping data → address of prev stack frame + address of the next inst after exit.
- function prologue: saving %esp in %ebp

→ When a function is called within another, the 2 stack frames form something of a linked list.

test a, → performs bitwise AND
Sfr OF & CF flags are set to 0.
SF is set to most sig bit of AND.



• Procedures + Func:

→ Involves passing parameters & return values.
Allocate local add variable & deallocate.

Caller save : %ebp, %esi, %edi ←

Called save :- %eax, %edx, %ecx ←

→ Recursive calls :- Each call is assigned its own space on the stack.

Non-modifying code → Reentrant.

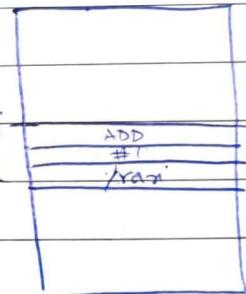
Modifying code → Non-reentrant.

Eg: computing sum of first N numbers.

Modification isn't allowed

"in the linux op. system text section {
", text section is }

READ-ONLY. (security).



→ even code is compiled using -g option,
compiler is less aggressive (removes some
optimisations) unlike when you compile using
-O2, -O3 etc. These are called 'Heisenbug'

→ (Symbol table : Mapping memory & reg. (virtual mem)

Heisenbug is detected when -O flag is utilized but
may not be detected otherwise -

Conditional instructions use PC rel. addresses

classmate

Date _____
Page _____

• When :-

who () {

[push 1.rbx]

who() {

}

[pop 1.rbx]

}

→ Caller save

instead if:- Inside the who() func:-

who() {

push 1.rcx .

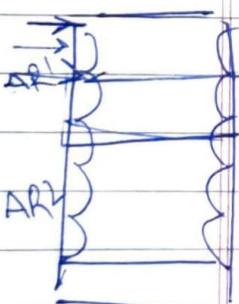
{

pop 1.rcx]

→ callee save

If it is a callee saved reg. the
caller can use it shamelessly
& vice versa -

Callee saved 1.rbx, 1.r12, 1.r13, 1.r14
1.rbp, 1.rsp



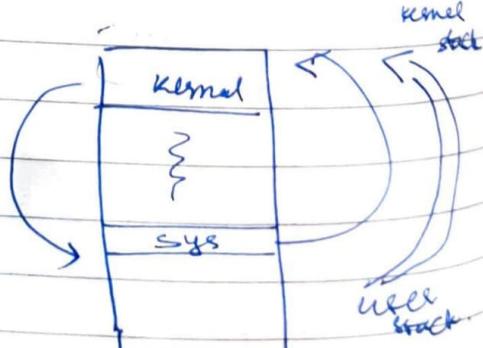
syscall :-

Looks for the syscall no:
in 1.rax

1.rdi

1.rsi → address

1.rdx → no. of bytes.



Switch statements make use of a jump table.

classmate

Date _____

Page _____

→ Encoding any instructions :-

Units	while	8 bytes	8 bits	8 bits
class	function	source dest		
inst code				Imm / Disp.
Eg :-	addq	y.ram rA, rB	y.rsi	

6 0 ↑ 0 6

Encoding: movq \$0xabcd, %rdx

30 F2. abcd

PUSH: A 0 rA F pushes rA

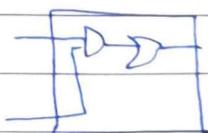
POP B 0 rA F Pops rA



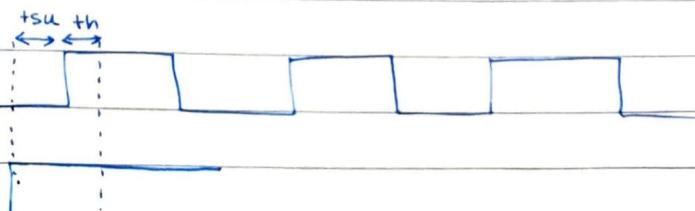
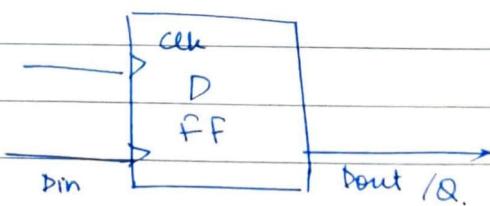
GOAL: Designing a processor.

↓ Digital circuits

Req: Combinational circuits, sequential el
adder/multiplexers



D-flip flop

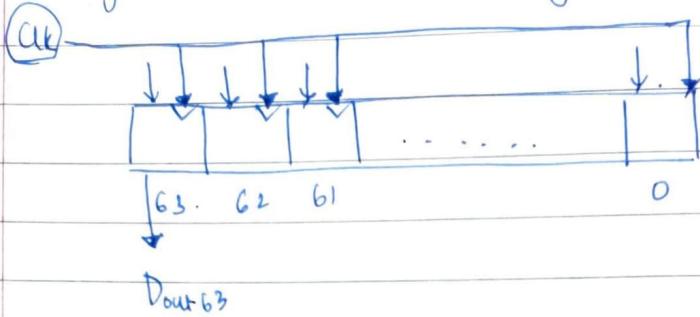


In the interval $[t - tsu, t + th]$ signal must be stable

→ clock - 2-Q time :- delay in obtaining the output
 ∴ output ob. at :- $t + t_{\text{clock}-\text{Q}_2}$.

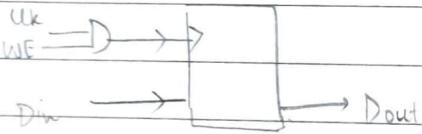
REGISTER FILE:

Q1 → Design a 64 bit register.

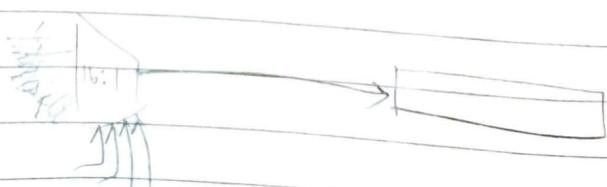
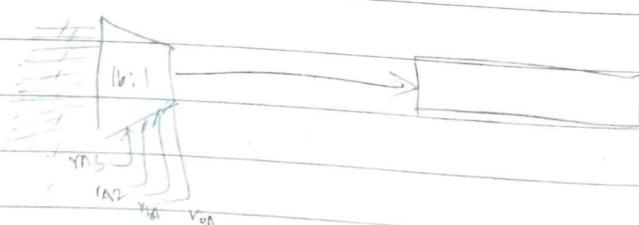
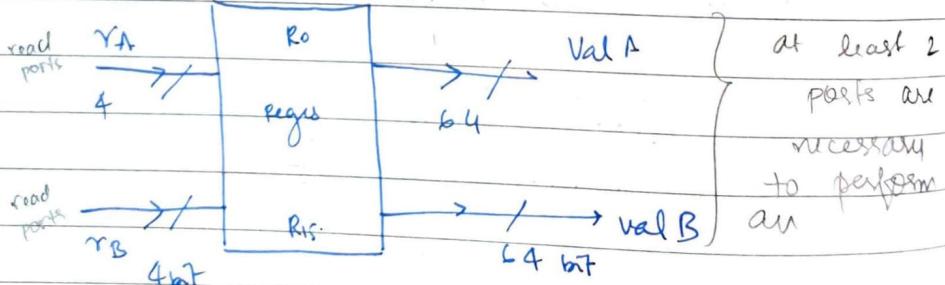


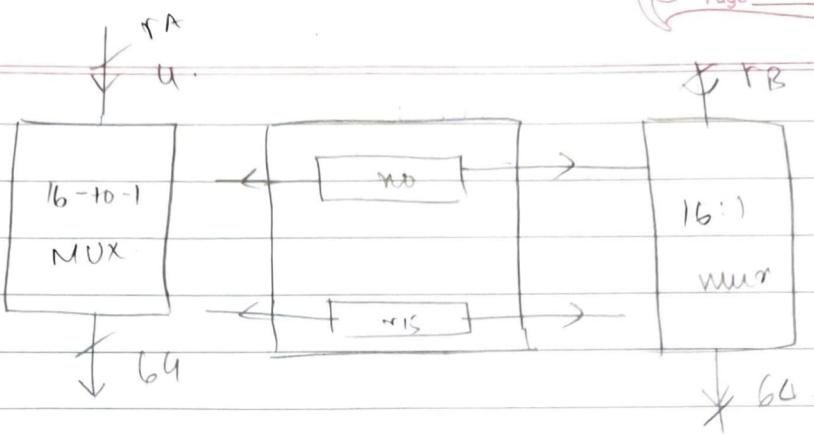
Q2 → Design a WE (Write Enable) D-FF.

Method 1:

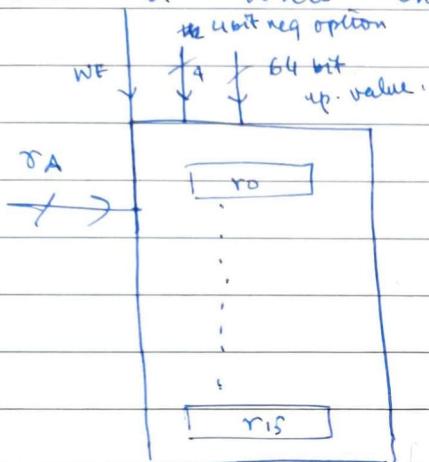


Q3 →

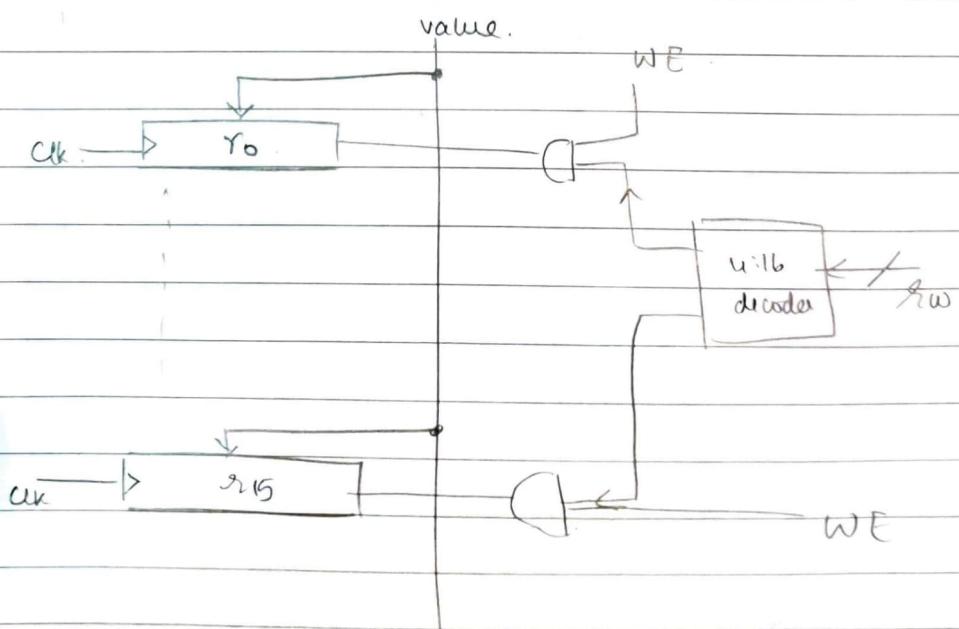




Q4 → Create a write enable in the registers :-



Naive will always update at least one register to avoid this scenario



If the components aren't clock synced, then it may so happen that some wrong register

values are updated.

i.e. At each rising edge \Rightarrow all required values for an update are available.

ALU COMP

Design an ALU to support

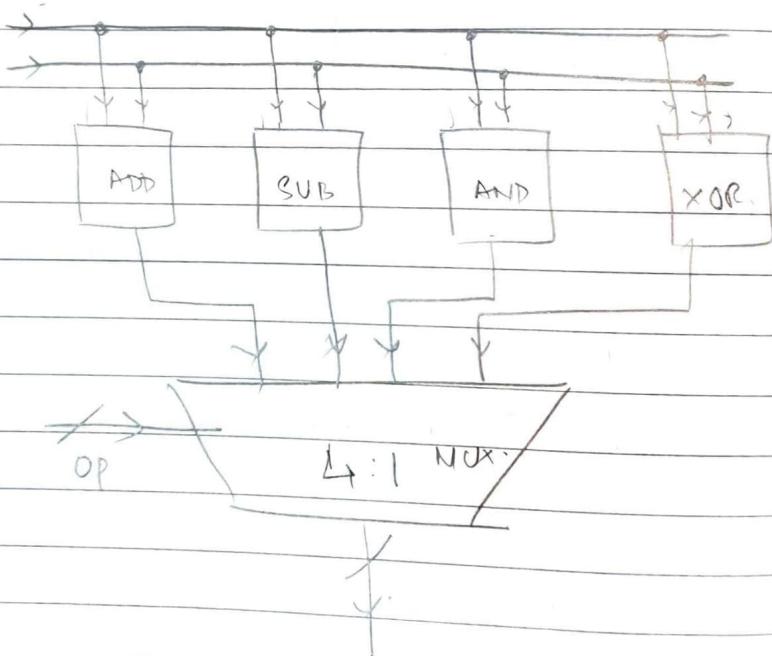
ValA
↓
↓ ValB

ALU

ADD SUB.

C = 0

ValW.

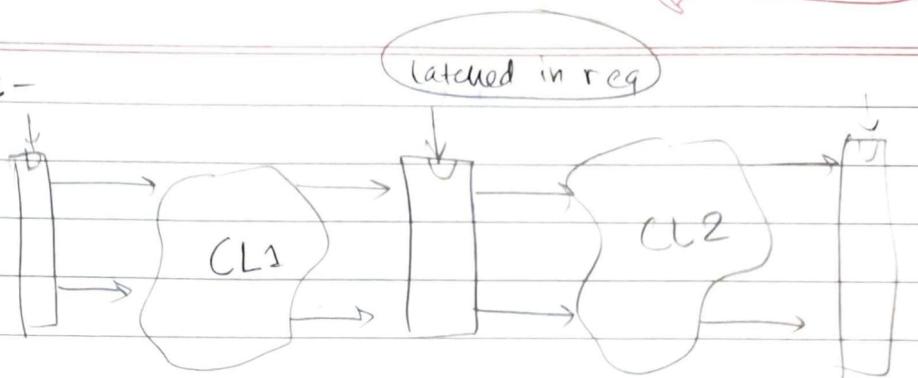


→ Assumption :- $t_{\text{add}} > t_{\text{sub}}$ $\therefore t_{\text{add}} \geq t_{\text{sub}}$ $\Rightarrow t_{\text{add}} \geq \max(t_{\text{add}}, t_{\text{sub}}, t_{\text{and}}, t_{\text{xor}})$.

→ For a series of combinational circuit.

$$C \geq \max(t_{\text{add}}^1, t_{\text{add}}^2)$$

→ Where :-

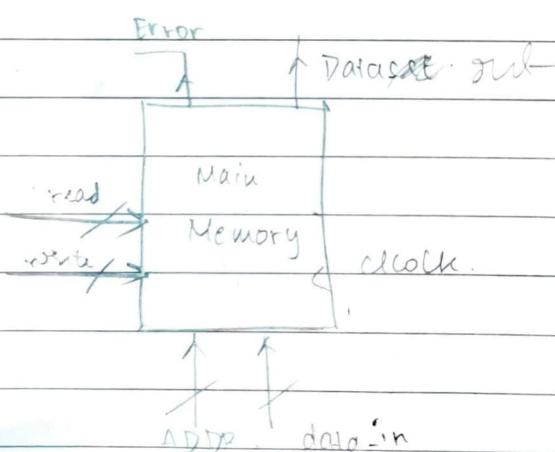


→ When a function is called:-

- Its return address is pushed onto the stack.
- If any caller save reg are used they must be pushed onto the stack.
- Assign additional space on stack for local variables (if necessary).



MEMORY MODULE :-



Assume data is in one memory location
inst. Are in another memory loc

MEMORY HIERARCHY:

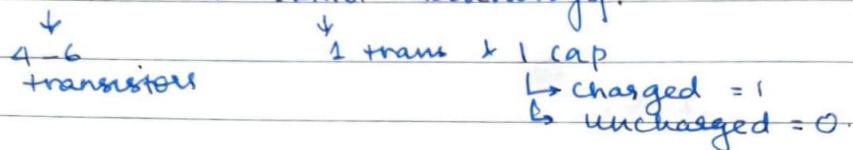
3 important components for any computational devic

- ↳ Processor
- ↳ memory
- ↳ I/O device

Mem. is a linear array of bytes

Physical realisation varies from proc to proc

And each bit of a byte is realised using SRAM or DRAM technology.



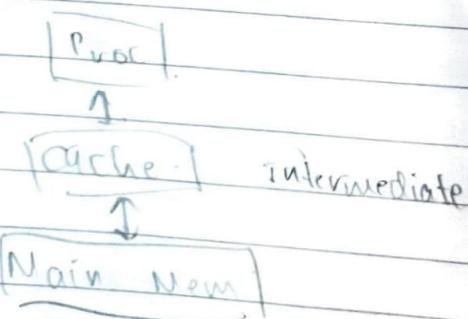
Memory refresh cycles ensure state is maintained in a DRAM Eg: Cap. shouldn't discharge

DRAM & SRAM - volatile

Flash memory → loads the micro loader from the harddisk

Processor \leftrightarrow Memory protocol

- Send add onto add bus
- Enable read control signal
- Latch the value in the register

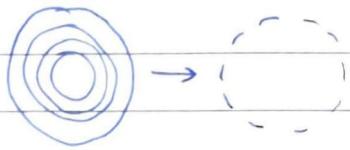


Hard disk - collection of platters. (2 sided)

In each platter there are tracks

In each track there are gaps.

Each sector is a primary storage unit → atomic units.

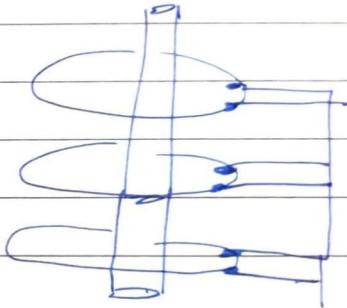


→ Storage density = recording den × track den.

To change any value → Read, Modify, Write

A sector of data ensures that it is either fully written or not written at all. (cannot be partially written).

Every surface has its own R/W handle.



Q: What is the DT time if

R → RPM

k → sectors / track

→ Seek time : latency to move the R/W head

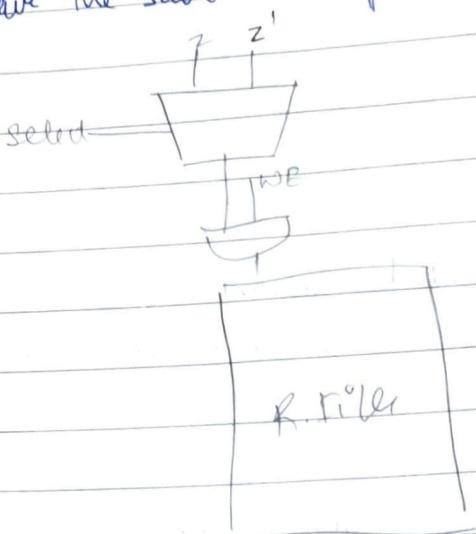
+ → rotational latency

+ → Time for Data transfer

Q: How to use cache to bridge the long b/w processor & memory.

SRAM: bistable - can have only 2 states. unstable
in the remaining states.

With zones, each track set is divided into zones. The tracks in a particular zone all have the same no. of sectors. (most).



- Control transfer
- Jump instruction

$$T_{avg} = T_{transfer} + T_{seek} + T_{avg\ rotational}$$

Two types of starting

→ Internal, external sorting algorithm

→ STEPS TO READ

- ① → Disk drivers / disk controller $(1, 2, 3) +$ param
- ② → Sector to be read detected on disk controller
- ③ → through DMA it accesses the memory
- ④ → Sends an interrupt to the processor.

* Processor does not remain idle till then. It keeps working.

• Memory wall:-

Processor \leftrightarrow Memory mismatch

Instructions & data cannot be supplied at the same rate.

Memory Hierarchy

Combine the 2.

Sram	Expensive, fast, small.
Dram	less exp., slow, large

Data items which will be accessed in the near future, put it in the cache.

what data is required by the user at the moment place that in the Sram. later usage ones place it in the DRAM

Instructions

e.g.: for ($i=0; i<100; i++$) (i)

Temporal locality.
same instruction repeated

Spatial locality
executing inst. that are spatially close.

Accessing A[0], A[1], A[2]
etc.

Stride-1 memory access pattern: accessing array el. in order (difference of indices = 1).

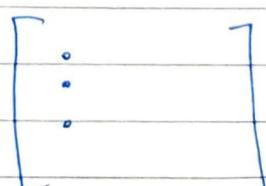
For frequent accesses of variable optimum locations its storage are either cache but preferably a register

Memory column major order
 row major order.

Mem



Stride = 1.



Stride = n.

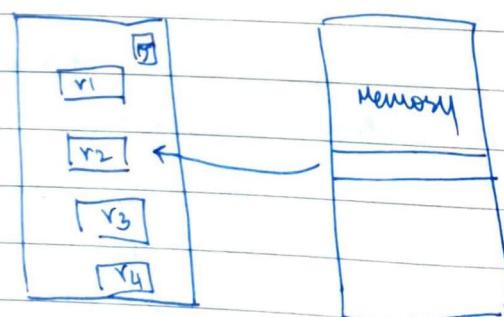
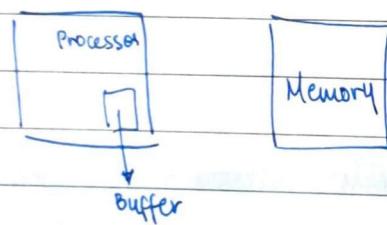


This is however a regular memory access pattern

In case of a linked list, the memory allocated will be random & stride is non uniform. Irregular memory access.

loop perm. trans :- computer tries to get the best stride possible.

→ Fetching data takes time. Significant effort hence maintain a buffer in the processor that holds the data for some interval of time. This is where caching comes into play.



Whenever an address is touched in any instruction, the address, the value & some meta information is cached.

With simple hashing, the register to which the content must be transferred to is decided.

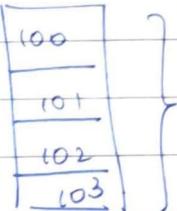
Assume 2^r registers :-

then we will need to consider the last n bits in the LSB of the address.

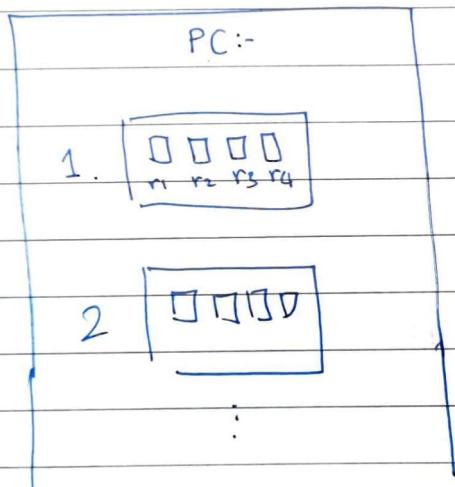
Eviction policy: kick out the already existing data if present.

Block size: no. of bytes we fetch in one-go.

e.g.:

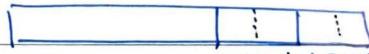


To ensure that the whole block of addresses the processor registers are also blocked



.., by from an address (1 byte) of a block.

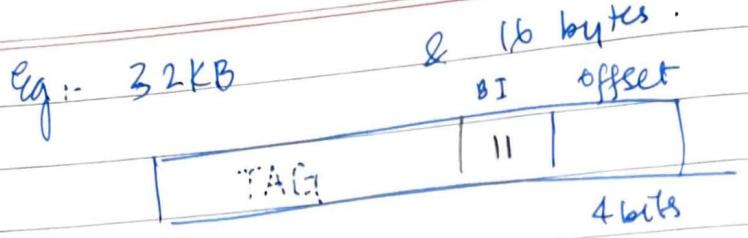
add:



block index offset

depends on no. of registers.

In block addressing, it adjusts the block bounds depending on the fetched address.



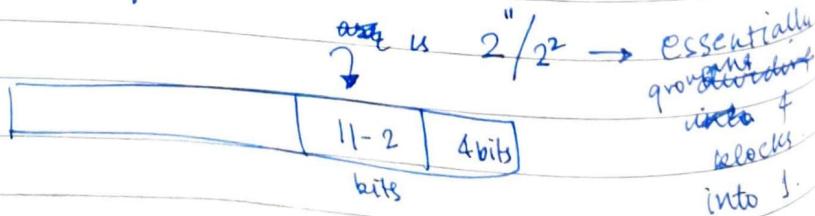
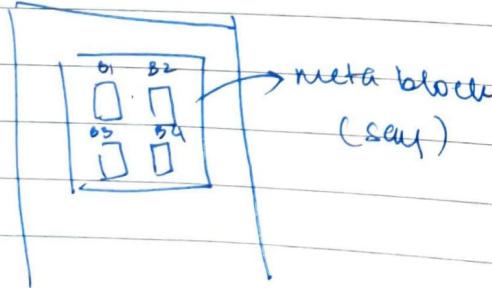
$$\frac{2^5 \times 2^{10}}{2^4} = 2^{11} \quad 11 \text{ bits}$$

If BI & offset bits are same for 2 addresses there could be an error. Hence we first match the TAG bits.

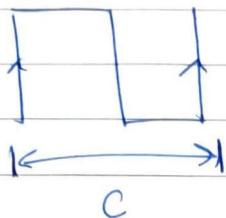
In case of a collision :- linearly probe.
set a max limit (associativity of cache) which is at max - how many blocks you will probe through.

In case all blocks (ass. cache) are occupied then use an eviction policy LRU . kick out the one that is least used by the processor.

Q:- 2MB main mem 32KB cache 16 byte
4 way set associative



- single cycle design

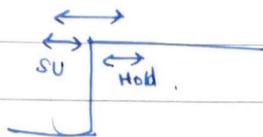


$$C > t_i$$

where $i \in \{ \text{MOV, ADD, SUB} \dots \}$

Disad:- some part of the clock cycle may go waste

Eg:- $t_{add} = \text{fetch time} + \text{ALU delay} + \text{setup time}$



$t_{morrow} = \text{Fetch time} + \text{address comp. delay} + \text{memory access} + \text{setup time}$

- Multicycle design :-

split the instruction into multiple steps

Eg: $t_{add} = t^1_{add} + t^2_{add}$

con. include overheads $\rightarrow t_{morrow} = t^1_{morrow} + t^2_{morrow} + t^3_{morrow}$

} $C > \text{max}$

freq. increases. $\therefore C' = C/k$

Actually the new frequency becomes

$$C' = \frac{C}{k} + \Delta \rightarrow \text{accounting for reg. mov. etc.}$$

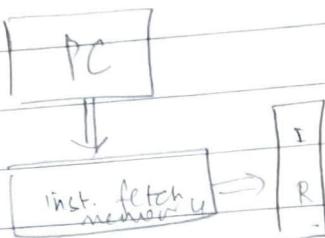
$10^6 C$ vs $10^6 \left(\frac{C}{k} + \Delta \right)$

$10^6 C$ vs $10^6 C + 4 \Delta \times 10^6$:)

this Δ value creeps in due to the set-up time
 \Rightarrow Processor w higher freq \Rightarrow better performance

- For any single instruction only, single cycle is better
- when there are multiple inst. multi-cycle design

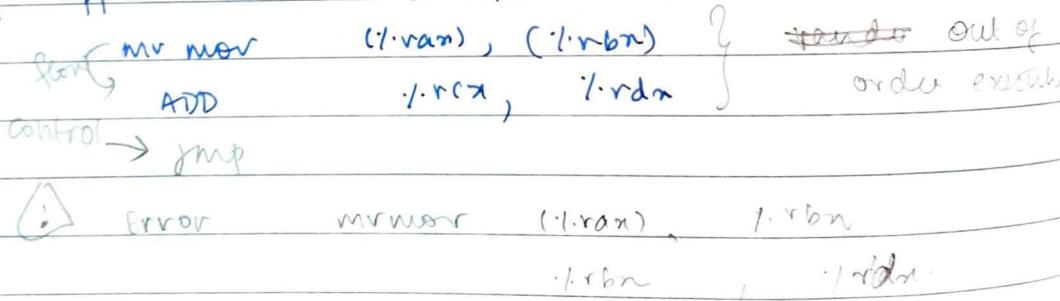
HW:-



• Microarchitecture :-

- Single cycle.
- Multi cycle.
- Pipelined. → embedded system (Inorder)
- Superscalar → desktop servers (out of order)
- VLIW. → embedded.

If approach is order out-of-order:-



A window of 16/ 64 instructions continuously reorders the inst. that must be executed at run time only. (done by hardware)

Static reordering could be done at compilation but requires more chip area & power.

Instruction scheduling → Static
Dynamic

Branch Another issue could be the jump instruction
eg:-

add ...
sub ...

jr table

Branch prediction strategy predicts the probability of picking one of the 2 branches

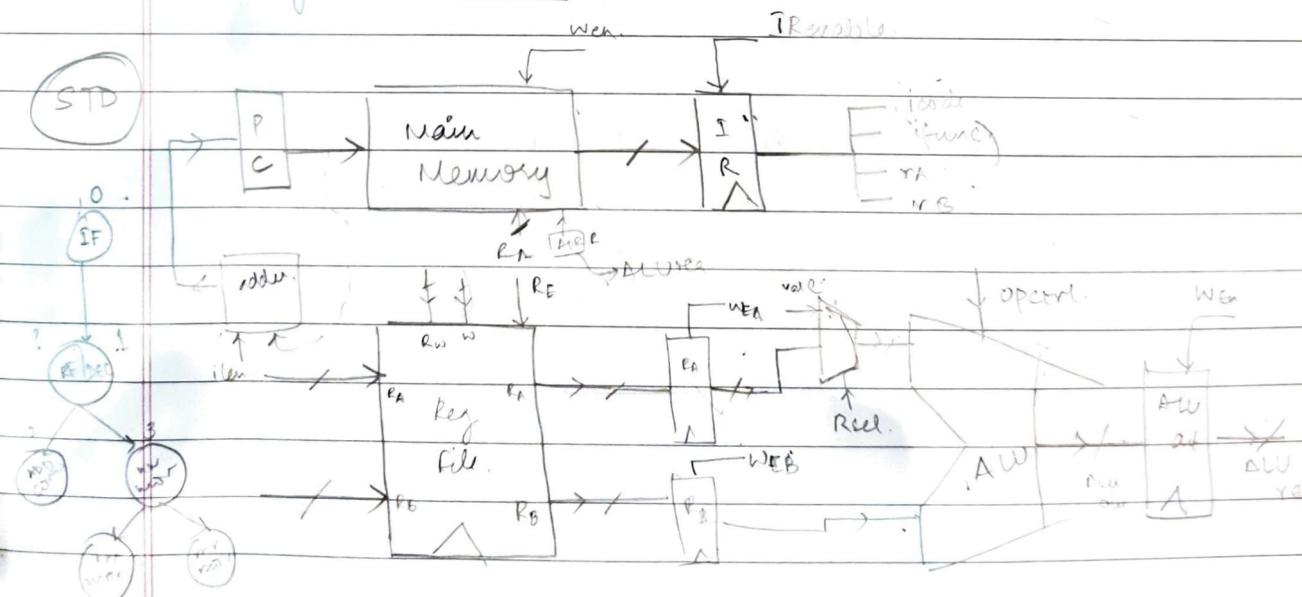
jmp → do not jmp

→ does that what BP recommends.

Speculative execution

If the incorrect branch is taken the processor must restore its state. wipe out the state done

Multicycle Processor :-



DATA PATH DESIGN.

- We cannot update in I.F. cycle if we want know by how much the PC must be inc.
- To optimize the PC inc. we can utilize the ALU itself (in multicycle).

3 - ADD, SUB, XOR,

jmp → 2

415 - MR mov

classmate
Date _____
Page _____

47 rm mov

3 - ir mov

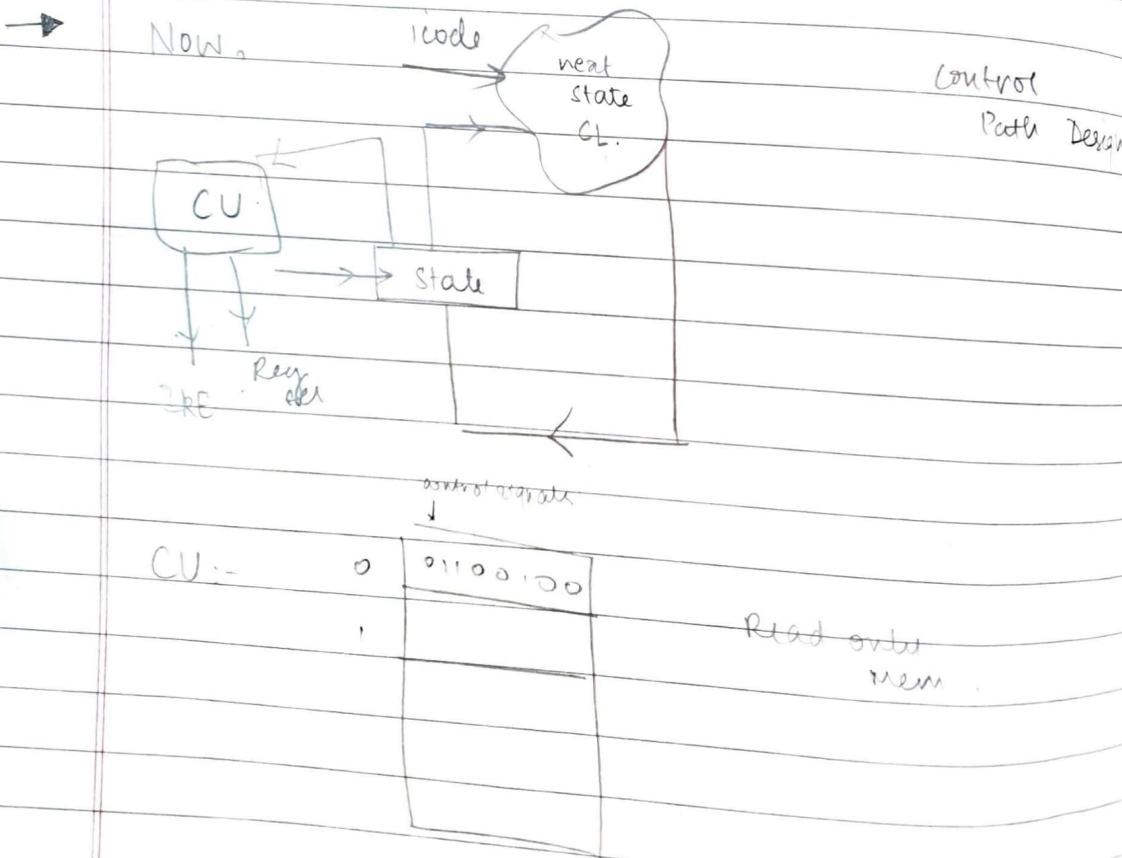
→ jmp instruction can be done in 2 clock cycles.

→ cmov → 3 clock cycles.
Reg to Reg

Single cycle : separate data & instruction
∴ it could lead to a structural hazard

but in multi clock cycle - 1 for IF, 1 for Data:
∴ no problem

Advantage of Multicycle Proc. shorter clock cycle

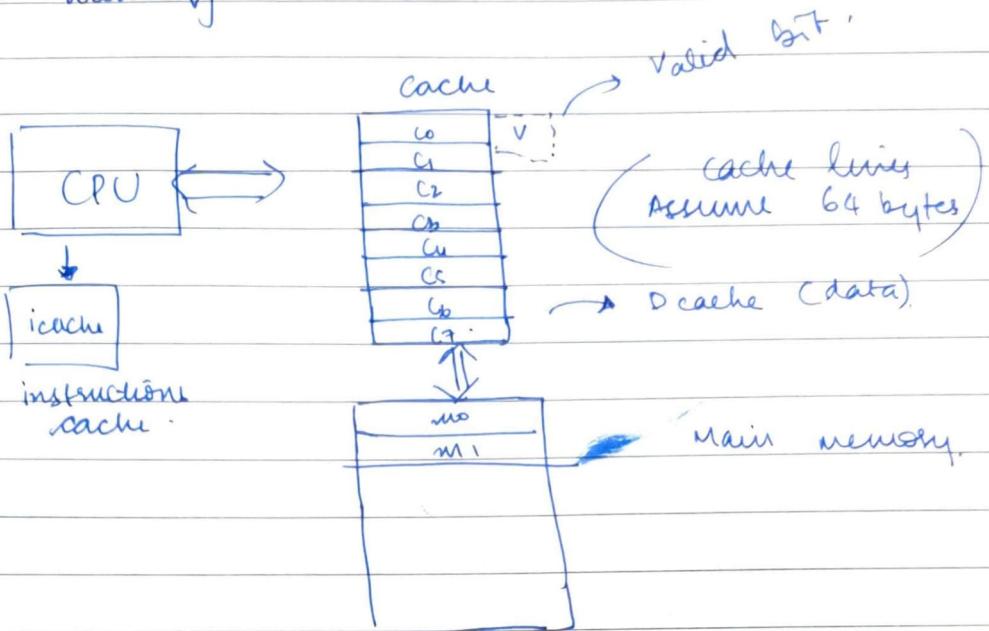


Eviction policy → LRU : least recently used
 CLASSMATE

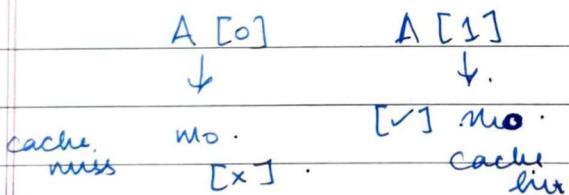
FIFO

Date _____
 Page _____

→ SRAM memory → cache memory.
 Main memory → DRAM.



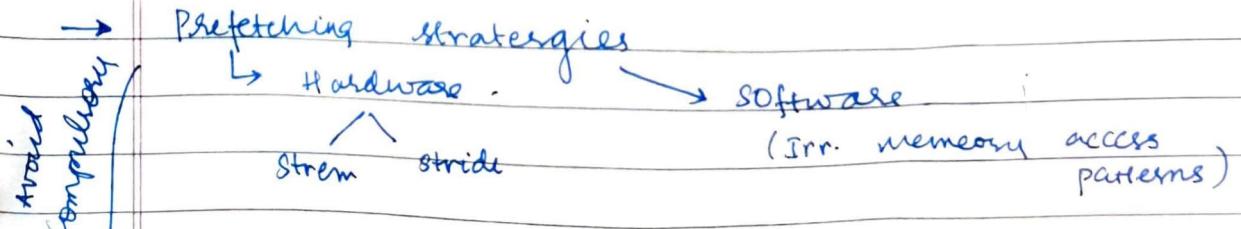
Consider :- for (int i=0; i<128; +i)
 sum += A[i];



$$\text{Cache miss ratio} = \frac{16}{128}$$

"optimal method → farthest use in future".

cache misses incurred - cold misses, compulsory misses



Stream & stride prefetchers fail miserably
case of linked lists.

Code 2:

```

for (j=0 ; j < 32 ; t+j) {
    for (i=0 ; i < 128 ; i++) {
        sum += f(j, A[i]);
    }
}

```

For each iteration of j , the blocks accessed are: m_0, m_1, \dots, m_{15}
 \Rightarrow there are 16 misses.

Capacitive misses :- unintentional. Bootstrapping process of missing the 16 misses when But unfortunately cache size is small & hence there will be a miss.

Here m_0, m_1, \dots, m_{15} \nearrow cache size
 working set

Here this conditions \Rightarrow capacitive misses will arise

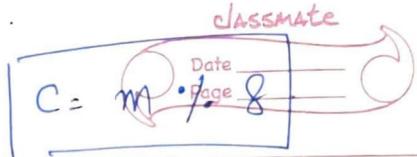
Rectification 1: To reduce cap. misses
 interchange the loops

then the working set size becomes

Direct mapped cache org.

Co : m₀, m₈, m₁₆

Ci : m₀, m₈, m₁₆ etc.



So for

[x]

i = 0

m₀ . . . m₀

i = 1

m₀ . . . m₀

:

i = 7

m₀ . . . m₀

32 times

[x] i = 8

m₁ . . . m₁

32 times

use m₀
8 times



• Rectification 2:

```
for (j=0; j<32; ++j)
```

```
    for (i=0; i<64; ++i)
```

```
        fun (j, A[i]);
```

```
for (j=0; j<32; ++j).
```

```
    for (i=64; i<128; ++i)
```

```
        fun (j, A[i]);
```

Code 3: for (i = 0 ; i < 32, ++i)
sum += A[i] + B[i];

if A[0] = { m₀, m₈ }

B[0] = { m₈ - m₁₆ }.

A[0]	B[0]	A[1]
m ₀	m ₈	A[1]
↓	[x]	[x]
[x]		

Fully associative cache organization

classmate
Date _____
Page _____

Accessing $A[1]$ is also a miss since currently m_8 is present in C_0 .

Therefore all iterations are misses here.
These are called conflict misses.

Resolving this problem:- any main memory block can go to any block FACC. (no conflict misses).

We cannot always have this as best method
 \therefore , searching for which cache block the m_i sits in will be expensive to do it in small clock cycles. Hence we may have to do a fully parallel search. \therefore , chip area & power cons. ↑.

→ 2-way set associative.

S₀ - {C₀, C₁} - m₀, m₄, m₈, M₁₂.

S₁ - {C₂, C₃} - m₁, m₅, m₉, M₁₃

\therefore No. of misses now will be 8.

Base address randomisation of arrays can reduce conflict misses.

Optimisation on one machine may not reflect on another.

Depends on cache size.

Code portability exists, but performance portability does not.

Installing libraries causes a tuning library to keep in which then enables the program to function well.

Atlas/Autotuning

Writes are asynchronous
Reads are synchronous

classmate

Date _____

Page _____

Write miss:

write no-allocate

write allocate

- On missing a memory block while writing, simply write into it don't upgrade.
- Write through cache update cache line as well as the main memory

→ real Write back

Just update the cache block. But on eviction totally rewrite that main memory block. This track is kept with the help of the dirty bit.

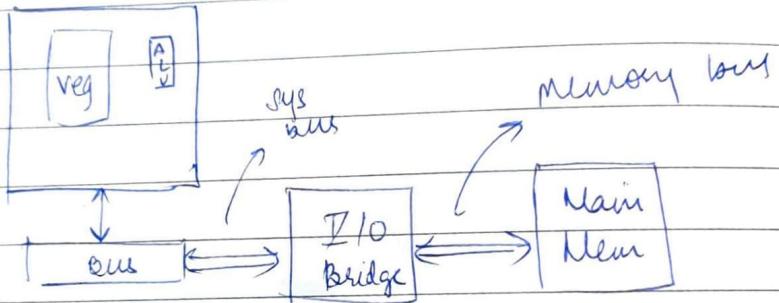
Cache Characteristics:
(E, S, m, B)

$$C = S \times E \times B$$

First identify the S then match tag bits then offset of b block off.

NOTES

- 1). SRAM - static : Exp, Fast, small
 DRAM - dynamic . Slow, large, cheap
- 2) SRAM - bistable , can retain its state indefinitely as long as it is powered.
- 3) SRAM & DRAM - volatile ; ROM, FLASH - non-volatile
- 4) basic CPU structure



5). $T_{access} = T_{seek \ avg} + T_{avg \ rot} + T_{dt.}$

$$T_{avg \ rot} = \frac{1}{2} \times \frac{60}{\text{given}} \times 1000$$

$$T_{avg \ seek} = \frac{60}{\text{given rev}} \times \frac{1}{500} \times 100$$

6) For the various disks are mapped to logical blocks. The firmware creates a tuple of 3 parameters (surface, track, sector) thereby accessing it. (A lookup table / mapping is maintained)

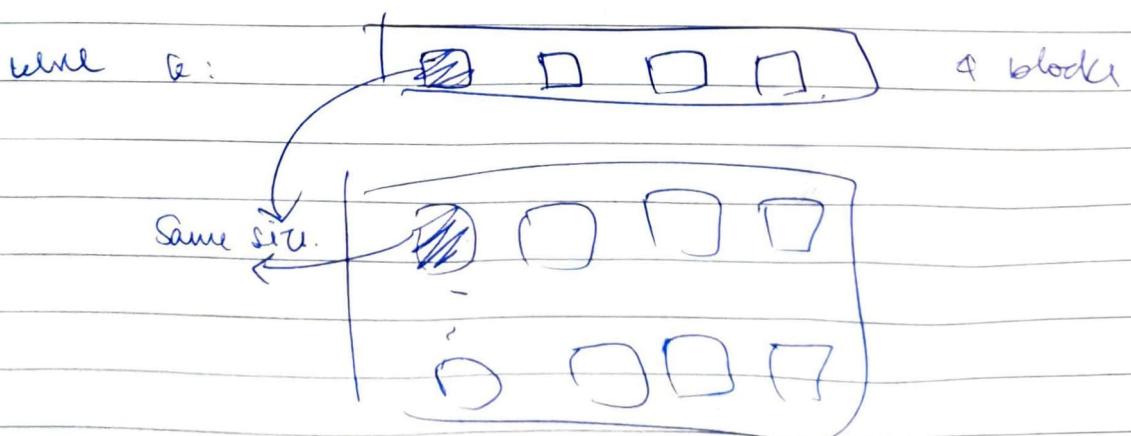
7) System bus & Memory bus are CPU specific

- 8) DMA: A process where a device performs a read/write on its own without the involvement of the CPU.
- 9) Principle of locality: tendency to reference data that is near recently referenced or that was itself recently referenced.
- 10) Smaller the stride \Rightarrow \uparrow spatial locality.
- 11) Memory hierarchy :-



12) Each level in this hierarchy caches ~~less~~ data in levels from lower hierarchy.

13) Cache blocks



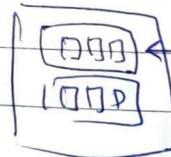
\rightarrow Further, at any point in time the contents of blocks at level k are copies / subsets of k+1

- 14) A transfer unit is a block of unit that is used to move data back & forth b/w the levels $k, k+1$.
- 15) Data already cached at level k & to be read from level $k \Rightarrow$ cache hit
else cache miss
- 16). Evicted block - victim block

- 17) Empty cache \rightarrow cold cache.
- \sim Misses of this kind :- compulsory misses
 - \sim Restrictive misses are :- conflict misses
 - \sim Working set $>$ cache size :- capacity misses

Problem 6.6 :

- 18) Terminology for how cache blocks are mapped
 $S \rightarrow$ set



E - Cache line

b - offset within a set.

- 19) A cache with exactly one cache line per-set is Direct Mapped Cache organization

- 20) Steps involved in accessing / finding a word from cache memory.
- Set selection
 - Line matched
 - Word extraction

- 21) To check if a particular word is present or not the valid bit must be set & tag bits in the cache line match tag bits of word.
- 22) The term thrashing describes when the cache is repeatedly loading & evicting the same set cache block.
- 23) In fully associative ~~cache~~ org. there are no set index bits.
- 24) Write through caches - write no-allocate
- write allocate.
- 25) Cache that holds both instruction & data is called unified cache.