

## X86 Processor

→ 8085, 8086 : initial registers.

386, 786, 5086, 7086 → X86 Generation (X86 Assembly).

X86 Processor → 32 bits (4 bytes can be processed in parallel)

→ 64 bits (8 bytes " " " )

→ X86 (32/64 bits) Assembly.

If source data is 256 bits → 64 bits will take 4 cycles to complete  
→ 32 bits will take 8 cycles to complete.

1. CPU fetch instruction from memory (instruction queue).

- FDE cycle.
2. CPU decodes instruction by looking at binary bit pattern
  3. If operands are involved - operands are fetched from memory + registers

4. Instruction is Executed.

5. If O/P → operand was part of instruction, the CPU stores the result of its execution in the operand

## Modes of Operation

↳ Protected : all instructions are available (memory protection + multitasking)

↳ Real-address : implements programming environment of Intel 8086 processor.

↳ System-management : provides OS management, system security, diagnostics etc.

↳ Virtual 8086 - hybrid of protected - each program has its own 8086 computer

## Registers (high Speed)

1. General Purpose Register - hold data for operation [arithmetic, logical and could store address]

2. Address Registers - holds address of data/instructions

× EIP : points to next instruction .

3. Status Registers - keeps current status of processor

AX - 16 bits      AH    AL      2 bits      2 bits      EAX - Extended AX (32 bits)      RAX - (64 bits)

Note: \*downward compatibility ; 64 bits can be used process 32 bits not the other way around.

## General Purpose Register (available for data manipulation).

1. EAX (Extended Accumulator Register) - used in arithmetic, logical & control

2. EBX (Extended Base Register) - serve as address register

3. ECX (Extended Counter Register) - serves as loop counter

4. EDX (Extended Data Register) - used for I/O & used with AX for \* & ÷

→ increase RAM capacity.  
 → done to implement backward compatibility  
 Segments. → efficient execution

Von-Neuman: data & address are in same memory divided

## 1. Code

## 2. Data

segment

3. Stack: holds local variables + function parameters.

how to access →  
 - using memory addresses which are stored in registers  
 - stores addresses of next instructions.



memory segment - block of consecutive memory bytes. Each segment is identified by segment number.

offset: where next address is (distance from beginning to next instruction).

To keep track of the various program segments, segment registers are used.

EES - Extended Extra Segment  
Used for accessing second data segment

## Pointer and Index Registers

ESP, EBP, ESI, EDI - point to memory locations (has offset addresses).

A memory location may be specified by providing a segment number & offset value  
Segment: offset

1. ESP: used in conjunction with ESS for accessing the stack segment.

2. EBP: access data on stack and other segments.

3. ESI: used to point to memory locations in the data segments addressed by EDI.

4. EDI: used to access memory location addressed by EES.

The program's code, data & segment are located into different memory Segment ie code segment data segment, stack segment

## Flag Registers

- indicates status of microprocessor - sets individual bits called FLAGS

0 = reset/clear  
1 = set.

- Status Registers: stores info whether an instruction should be executed or not.

↳ Control and System Flags: control CPU operation

- 8 general purpose

↳ Status Flags: status of arithmetic + logical operations.

- 6 segment registers

- EIP

## → Basic Language Elements (Assembly Language)

### main PROC

move eax, 5 ; move 5 to EAX register

add eax, 6 ; add 6 to EAX Register.

INVOKE ExitProcess, 0;

end program

### main ENDP

\* Operands: a value that is

### → Mnemonic Description

MOV assign (move) one value to another

ADD add 2 values

SUB subtract 2 values.

MUL multiply 2 values

JMP jump to new location

LOL

;Add Two.asm - adds two 32 bit integers

.386 //represents 32 bit architecture.

.model flat, stdcall //telling how much memory is reserved for data & code

.stack 4096 //memory assigned to stack.

ExitProcess PROTO dwExitCode:DWORD //acts as return 0; stops all threads of processes.

.code //code segment

main PROC

MOV eax, 5 ;

ADD eax, 6 ;

INVOKE ExitProcess, 0

main ENDP

END main //Program ends //

→ Little Endian Order

- refers to order in which Intel stores integers in memory.

- multi-byte integers are stored in reverse order when most significant byte is stored at MSB and smallest is stored at LSB.

→ Big Endian Order

- most significant byte is stored at lowest address.

→ Directives vs Instructions

↳ tells assembler what to do.

↳ to tell CPU what to do.

→ Procedure

[NAME] PROC

[NAME] ENDP

LABEL (optional), Mnemonic, Operands → Instruction Format.

Integer Literals:

Encoded real=r

{+/-} digits [ radix ] binint = b/y octal = 9/o hexadecimal = h... decimal = t/d

Identifier → identifies variable.

MOV destination, source.

can't IP, EIP, CS

otherwise  
Assembler gives  
error

not immediate to segments!

must

destination = source  
size size

**data**  
 bVal BYTE 100  
 bVal2 BYTE ?  
 wVal WORD 2  
 dVal DWORD 5

possible

register, register  
 register, value  
 value, register  
 register, memory

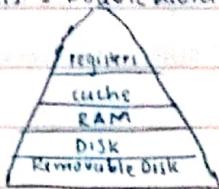
1, 0 = bit

4 bits = Nibble

8 bits = Byte

16 bits = Word

32 bits = DoubleWord



Memory hierarchy

• code

mov ds, 45

mov es, wVal

mov eip, dVal

mov 25, bVal

mov bVal2, bVal

// Immediate value can't because its segment

// size mismatch (32bit + 16bit)

// Invalid src transfer (doubleword)

// 25 + not address

// bytes are memory so no data transfer

as memory to memory transfer is not possible ∵ transfer to register → register to memory.

### Direct Memory Operations

- DMO is a

Little Endian Order VS Big Endian Order

eg (12345678h)

0000 78

0000 12

0001 56

0001 34

0002 34

0002 56

0003 12

0003 78

### MOV Instruction

MOV destination, source

no more than one memory operand is allowed.  
 ↪ mostly CPU has control.

CS, EIP, IP cannot be destination

SWORD - Signed WORD (16-bits)

SDWORD - Signed double

DB - define byte (1 Byte)

DW - define word (2 Bytes)

DD - define doubleword (4 bytes)

DQ - define quad. (8 bytes)

GANT - define 10 bytes

LLOL

Operand	Description
reg8	8-bit general purpose
reg16	16-bit general purpose
reg32	32-bit general purpose
reg	any general register
seg	16-bit segment CS, DS, SS
imm	8, 16, 32 bit immediate value

### Multiple Initializers.

label refers to the offset of the first initializer

list BYTE 10, 20, 30, 40.

list BYTE 10, 32, 41h, 00100010b

list BYTE 0Ah, 20h, 'A', 22h

to ensure that assembler doesn't think this is a label

offset = 0000

### .data

Array1 BYTE 10h, 20h, ?, 40h

BYTE 50h, 60h, 70h, 80h

WORD 81h  
BYTE 82h, 83h

list 3 BYTE ?, 32, 41h

### .code

MOV AL, Array1 + 7 ; AL = 80h.

### Defining strings

String ends with null byte (containing 0) called null-terminated string.

greeting1 BYTE "Good afternoon", 0 → shows it's null terminated.

greeting2 BYTE 'Good night', 0 → 15 bytes

each character uses a byte of storage.

rule that byte must be separated by commas doesn't apply on strings

DUP operator.

↳ allocates storage for multiple data items, using integer expression as a counter.

X BYTE 20 DUP(0) ; 20 bytes, all equal to zero.

Y BYTE 20 DUP(?) ; 20 bytes, uninitialized.

array S BYTE 4 DUP ("STACK"). ; 20 bytes STACKSTACKSTACKSTACK

Z BYTE 7 DUP (2 DUP(17,19)) ; 17 19 17 19

Var4 BYTE 10 DUP(0),20 ; 10 000 20

∴ DUP directive is used to duplicate or replicate data values in an array or a block of memory.

17 19 17 19 17 19 17 19 17 19 17 19 17 19  
17 19 17 19 17 19 17 19 17 19 17 19

→ signed word.

Defining WORD and SWORD.

assigns 16 bit.

The legacy DW allows us to store signed num without explicit definition.

Declaring Uninitialized Data.

DATA? - used to declare an uninitialized data segment

- executable memory is used less - assigned at compile time.

Symbolic Constants.

EQU - integer values

TEXT EQU - string

Current Location Counter : symbol \$ is used which returns offset of current location

selfptr DWORD \$

declares selfptr and initializes it with variable's offset value.

Calculating sizes of Array and String.

List BYTE 10,20,30,40,

listsize ~~DB~~ = \$ - list → this is equivalent to list[0].

listsize must follow immediately after list

offset\\_arr3 : (d - arr3) / 4

so we get in word.

\$ current offset address.

= used for assignment/updation of variables

equ integer can't be updated but string can be updated  
↳ constant symbol

eg exp equ 2\*5

exp . equ 21 // can't redefine.

exp = 20

eg txt equ <"ABC"> // can redefine

txt equ <"abc"> as in angular brackets

txt2 equ <10\*10> ; text

eg symb equ exp  
symbol  
mv equ mov // instead of mov mv can be used.

axmov equ <mv al, v1>

axmove equ <mov, al, v1>

list db exp, txt, txt2, symb.

axmov equ <mv, bl> N1

GAMA

LOL

## NOTES.

Directives tell compiler what to do while Instructions tell CPU what to do

Procedure [Name] PROC

Procedure

[Name] ENDP

Character literal 'A1', 'N'

    ↑ Directive

myVar DW D 26

String literal 'ABC'

Instruction ←

MOV eax, myVar

MOV destination, source

EIP, CS, IP cannot be destination

• 386 directive that identifies the program as 32 bit program (can access 32 bit address) register

• model flat, stdcall selects the program memory model and identifies the calling convention

stdcall keyword tells the assembler how to manage the runtime stack

when procedures are called

calling convention: how subroutines receive parameters from their caller and how they return a result.

• Stack 4096 sets aside 4096 bytes of storage for runtime stack.

## STACK

- holds passed parameters

- holds address of the code that called function. The CPU uses this address to return when the function call finishes, back to the spot where the function was called.

- holds local variables

• Code marks beginning of the code area of program

Multiple Initializers → list BYTE 10h, 32, 00101111b, 'A1', ?, 38

Defining Strings: ends with null character ex g1 BYTE "Good", 0 (byte)

each letter uses a byte so total 5 bytes used.

.Dup - duplicates values.

Symbolic constants (symbol definition) created by associating an identifier

name = expression [ can be ~~not~~ redefined ] ; name is symbolic constant

1. Carry Flag - subtract smaller value from larger.

when addition results in

2. Sign : MSB = 1 sign flag on

3. Overflow : Overflow XOR of <sup>carry in</sup> last bit & last carry out

4. Auxiliary : 3<sup>rd</sup> bit to 4<sup>th</sup> bit value shift.

5. Parity : even number of 1s in lower byte.

## NOTES.

1 /

### Overflow

#### → Unsigned Overflow

when there is carry out of the MSB

correct answer larger than biggest unsigned number

On subtraction unsigned overflow occurs when there is borrow into

MSB. Correct answer is smaller than 0

#### → Signed Overflow

- addition of numbers with same sign, signed overflow occurs when sum has different sign.

- subtraction of numbers with different signs is like adding numbers of same sign.  $A - (-B) = A + B$        $-A - (+B) = -A + (-B)$

Signed overflow occurs when result has different signs than expected.

- OVERFLOW IS IMPOSSIBLE WITH DIFFERENT SIGNS

## SHIFT INSTRUCTIONS

for shift instructions bits shifted out are lost.

for rotate instruction bits shifted out are moved in from other end.

opcode destination, 1 [shift/rotate by ]

opcode destination, N [shifts/rotates by N ]

#### → SHL (multiples by 2)

MSB shifted to CF ; 0 is moved into rightmost [LSB]

Effects on flags

SF, PF, ZF reflect the result.

AF : undefined

CF : last bit out of MSB

OF : 1 if result changes sign on last shift

e.g. 10001010, 3

1. CF 1 00010100 OV=1 SF=0 ZF=0 AF=NULL PF=1

2. CF 0 00101000 OV=0 SF=0 ZF=0 " " " " " " PF=1

3. CF 0 01010000 OV=0 SF=0 ZF=0 " " " " " " PF=1

## SAL Instructions

used for signed (emphasizes on extraneous nature)

1000 0000, 2

CF=1 0000 0000

CF=0 0000 0000

CF=0 OV=0 SF=0 ZF=0

→ SHR Instruction. (Divides by 2).

0 shifted into MSB and LSB moved to CF

SHR destination, N

→ SAR Instruction.

MSB retains original value (rest works like SHR).

SAR destination, N.

for odd numbers halves +  
rounds down

→ Signed and Unsigned Division

for unsigned use SHR

for signed use SAR

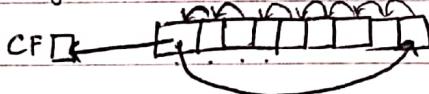
e.g. To divide by 4 use 2 shifts or  $2^2 = 4$ .

## ROTATE INSTRUCTIONS.

→ Rotate Left (ROL)

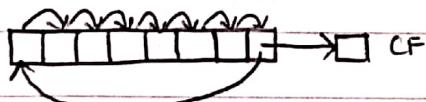
MSB is shifted to LSB (rightmost)

CF also gets bit shifted out of MSB.



→ Rotate Right (ROR)

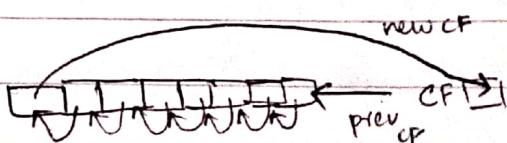
LSB shifted to MSB and LSB also goes in CF



→ Rotate Carry Left

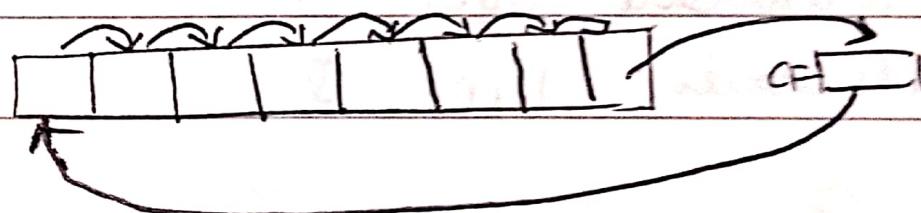
shifts bits to left · MSB shifted to CF

previous CF value is moved to LSB



→ Rotate carry Right

bits rotated to right



Effects ON Flags .

SF, PF, ZF reflect result .

AF undefined

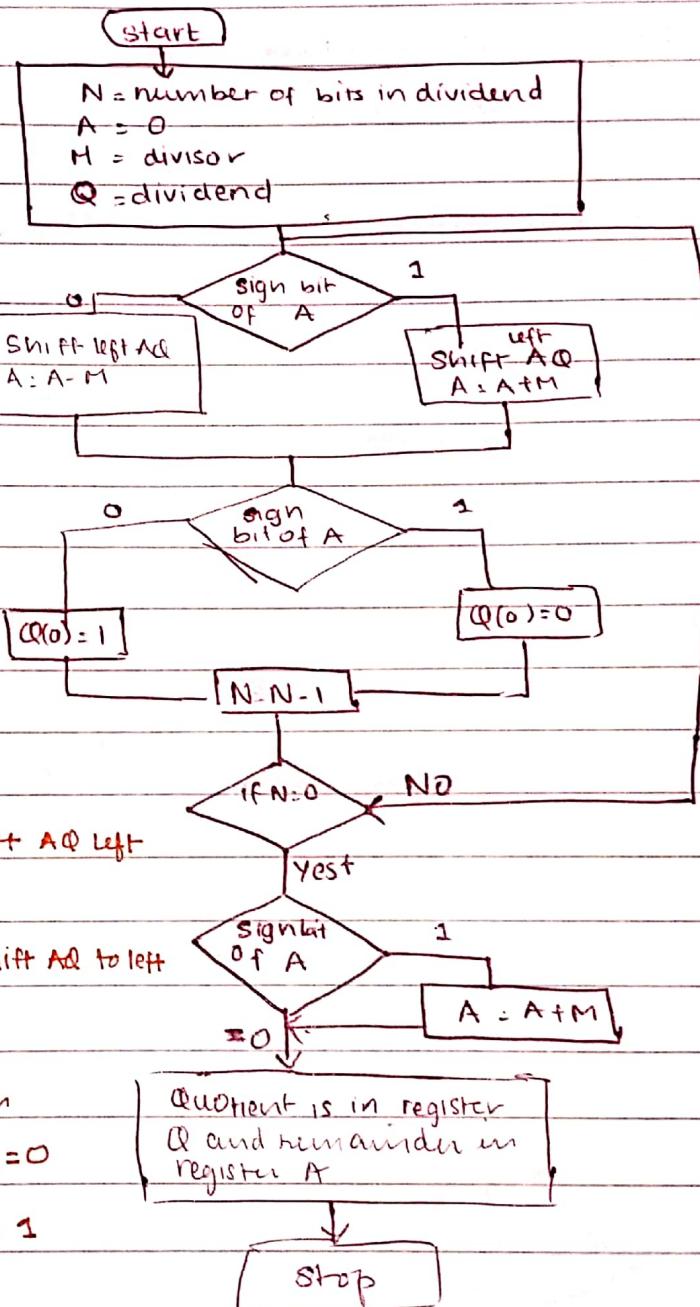
BF > 1 if result on changes of sign on last rotation

CF > last bit shifted out

### NON-RESTORING ALGORITHM (Division)

Convert all algorithms to assembly

+ get description.



1. Initialize values

2. Check sign bit of A

3. If sign bit of A = 1 then shift A<sub>Q</sub> left

and perform  $A + M = A$

If sign bit of A = 0 then shift A<sub>Q</sub> to left

and perform  $A = A - M$

4. Check sign bit of A again

5. if sign bit is 1 then  $Q_0 = 0$

If sign bit is 0 then  $Q_0 = 1$

[ $Q_0$  is LSB of Q]

6. N is decremented

7. If  $N \neq 0$  loop from step again

else

check sign bit of A if 1 then  $A = A + M$

else

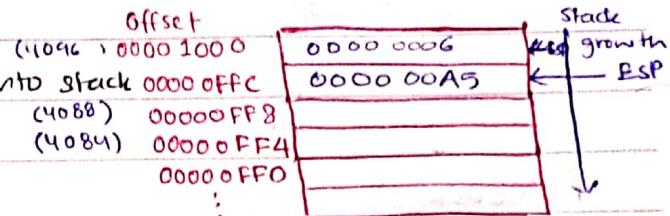
8. A contains remainder

Q contains Quotient

## Stack Operations

- LIFO (last in first out).
- ESP (Extended Stack Pointer) is used to manage memory
- ESP holds 32 bit offset into some location on stack
- Use instructions such as CALL, RET, PUSH, POP to manipulate ESP

- ESP always points to last value pushed onto stack



## PUSH Operation

- (32 bit) decrements stack pointer (ESP) by 4 (32 bits) and copies value into the location in stack pointed by ESP.  
↳ if 16 bit then decrements by 2
- Stack grows downward in memory.
- Each time new value is added location of ESP points to newly added value

## POP Operation

removes value from stack

The value pointed by ESP is removed and ESP increase by 4 (32 bits)

if 16 bits then increments  
by 2

\* when calling a subroutine, you pass input values all arguments by pushing them on stack

\* provides temporary storage for local variables inside subroutines

## PUSH / POP Instructions

pushes reg/mem value onto stack

PUSH reg/mem16

// decrements ESP by 2

PUSH reg/mem32

// decrements ESP by 4

PUSH imm32

// decrements ESP by 4

POP reg/mem16

// increments ESP by 2

POP reg/mem32

// increments ESP by 4

↳ copies popped value in reg/mem

### PUSHFD / POPFD Instructions

- PUSHFD is used to push 32-bit EFLAGS register onto stack
- POPFD is used to pop/reload EFLAGS from stack
- MOV instruction cannot be used to copy the flags to a variable, so PUSHFD is the way to save value of flags
- This allows us to restore value of flags from stack.

`pushfd ; save current value of flags  
// code`

`popfd ; restores value of flags`

\* Ensure if pushfd is used then your program executes popfd also and not skip it

\* This is however ERROR-PRONE, so there is another less error-prone way

push EFLAGS on stack and pop immediately and store in a variable.

o data  
`saveflags DWORD ?`  
 o code  
`pushfd` ; this stores in saveflags variable  
`pop saveflags`  
`push saveflags` ; copies value of saveflags to EFLAGS.  
`popfd`

### PUSHAD, PUSHA, POPAD and POPA

1. **PUSHAD**: pushes all 32-bit general purpose registers onto stack. [EAX, ECX, EDX, EBX, ESP (value of ESP before execution of PUSHAD), EBP, ESI, EDI]

The stack stores in the above mentioned order ↑

2. **POPAD**: pops registers off the stack in reverse order EDI; ESI .... EAX

3. **PUSHA**: pushes 16-bit general purpose registers (AX, CX, DX, BX, SP, BP, SI, DI)

4. **POPA**: pops 16-bit general purpose registers in reverse order

\* for 16-bit programming mode use PUSHA; POPA

\* for 32-bit programming mode use PUSHAD; POPAD

- If you create a procedure which uses general purpose register you may use pusha/pushad and popa/popad.

`SUB proc`

`pushad`

`:`

`popad`

**GAMA** ref  
**BNDP**

\* procedures returning results in one/more registers should not use ↙

LOL



Stack Frame is created by following steps:

1. Passed arguments are pushed on stack
  2. The subroutine is called, causing the subroutine return address to be pushed on the stack. (pushed address (العنوان المنشئ) ) .
  3. As the subroutine begins to execute, EBP is pushed on the stack
  4. EBP is set equal to ESP
- From this point on, ESP is used as reference for all of subroutine

Passing by Value is when value of register is sent

Passing by Reference is when offset is passed.

+ why index starts at 0 from C++

#### o code

Main PROC

PUSH 5

PUSH 6

CALL AddTwo

ADD ESP, 8 [write this statement]

MAIN ENDP

AddTwo PROC

PUSH EBP

MOV EBP, ESP

MOV EAX, [EBP + 8]

ADD EAX, [EBP + 12]

POP EBP

RET 8 [or write this statement]

AddTwo ENDP

When a local variable is passed through procedure always include

```
push EBP
mov EBP,ESP
```

Two ways to clear stack.

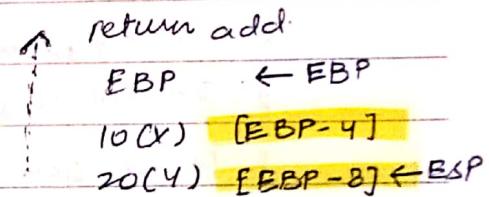
- POP elements
- or just decrement ESP
- RET8 means release 8 bytes after function return.

## Creating Local Variables.

MYSub PROC

```
PUSH EBP
MOV EBP,ESP
MOV DWORD PTR [EBP-4], 10
MOV DWORD PTR [EBP-8], 20
MOV ESP, EBP
POP EBP
RET
```

MYSub ENDP



- Parameters can be passed on stack as well as in registers

### Stack frames

- Passed arguments are pushed on stack
- The subroutine called causing the return address to be pushed on stack.
- As the subroutine begins to execute, EBP is pushed on the stack
- EBP is set equal to ESP
- If any local variables were created ESP is decremented to reserve space for variables on the stack
- If any registers need to be saved they are pushed on stack.

### Enter and leave instructions

- ENTER performs 3 ops:

1. Push EBP on stack

2. Set EBP == ESP

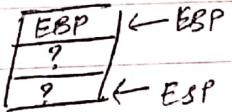
3. Reserve space for local variables ie sub esp, numbytes

ENTER numbytes nesting level

                ↑      ↑  
                immediate values

e.g. mySub PROC

ENTER 80



leave terminates the stack frame procedure.

written before ret instruction

### LOCAL directive

Saves space for local variable on stack

declaring must happen after MySub Proc instruction.

LOCAL var[10] : BYTE

Saves space for one/more variables

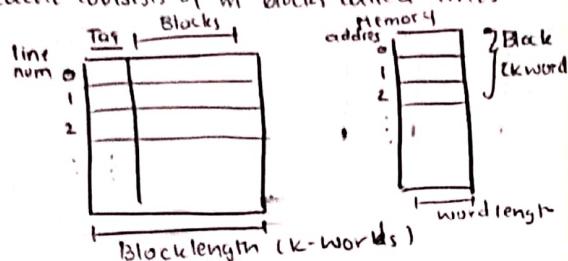
## Cache

There are 3 cache levels L1, L2 and L3

L2 is slower and larger than L1

L3 is slower and larger than L2

Cache consists of m blocks called lines.



## Mapping Function

- algorithm for mapping main memory into cache lines
- A means of knowing which main memory occupies a cache line
- o direct      o Associative      o Set Associative

## EXAMPLE

- Cache can hold 64 KBytes
  - data is transferred b/w M.M and cache in blocks of 4 bytes each
  - Main memory consists of 16 Mbytes with each byte directly addressable by 24 bit address.
- Thus for mapping purpose we can consider main memory to consist of 4 M blocks of 4 bytes each.

## Direct Mapping

maps each block of main memory into one possible cache line.

$$l = j \bmod m$$

$l$  = cache line number

$j$  = main memory block number

$m$  = number of lines in the cache.

Address length =  $(s + w)$  bits

Addressable units no =  $2^{s+w}$  words/bytes

Block size = Line size =  $2^w$  words/bytes.

No. blocks in main memory =  $\frac{2^{s+w}}{2^w} = 2^s$

Number of lines in cache =  $m = 2^r$

Size of cache =  $2^{r+w}$

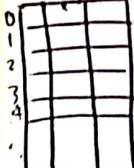
Size of tag =  $(s - r)$  bits

- Cache is divided into fixed-size blocks containing multiple words of data

- Principles of temporal and spatial locality tell us that recently accessed data and data close to it, are likely reused in future

- when accessing data if it is already present in cache we call it cache hit, else cache miss

V Tag Data



- data is stored in blocks

- blocks contain multiple words of data

- each block is selected by index

- valid bit indicates if data is valid or not

- Tag is remaining part of address

## Terminologies

Capacity(C): Total number of bytes that can be stored in cache.

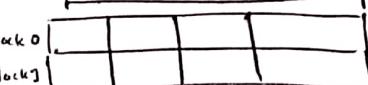


block/cache line

$$\text{eg } b = 32 \text{ bits} = 4B$$

$$C = 64 \text{ bits} = 8B$$

cache line



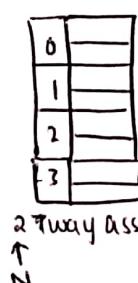
$$\text{Offset bits } \log_2 b = \log_2 4 = 2$$

• Block size(b): bytes of blocks in cache

$$B = C/b \quad \text{in this eg } \frac{8}{4} = 2$$

• Degree of associativity (N): number of blocks in a set

• Number of sets =  $S = B/N$   
4-way associative



1011  
1101  
1110  
1111

Page 1

eg  $b = 4B$

$$C = 64B$$

$N = 1$  (Direct)

$$B = C/b = 64/4 = 16$$

$$S = B/N = 16/1 = 16$$

$$\log_2(S) = \text{Set bits} = 4$$

$$\log_2(b) = \text{Offset bits} = 2$$

sets	00	01	10	11
0000				
0001				
0010				
0011				
0100				
0101				
...				
1111				

highlighted part  
will all add up  
to capacity  
ie 64.

eg  $b = 32B$

$$C = 128B$$

$N = 1$  (Direct)

$$B = 128/32 = 4$$

$$S = B/N = 4/1 = 4$$

$$\log_2 S = \text{Set bits} = \log_2(4) = 2$$

$$\log_2 b = \text{Offset bits} = \log_2(32) = 5$$

sets	00000	00001	00010	00011	00100	...	10000
00							
01							
10							
11							

adds up to  
128

eg  $b = 4B$

$$C = 32B$$

$N = 4$  (4-way associative)

$$B = C/b = 32/4 = 8 \quad (8 \times 4)$$

$$S = B/N = 8/4 = 2$$

$$\text{offset bits} = \log_2(4) = 2$$

$$\text{set bits} = \log_2(2) = 1$$

sets(S)	00	01	10	11
0				
1				

each block  
yields up to  
C ie 32

eg  $b = 8B$

full Associative cache.

single set with four ways

$$C = 8B * 4 = 32B$$

$$B = C/b = 32/8 = 4$$

$$S = B/N = 4/4 = 1$$

set	00	01	011	100	101	110	111
0							

adds up to  
32.

**Tag bits:** we add tags to the cache which lets us supply the address bits to support let us distinguish b/w different memory locations that map to same cache block

**Valid bit:** indicates whether data is valid or not (1).

1. at start all N bits are 0

2. when data is loaded into a particular cache block, the corresponding valid bit is set to 1.

### Example 1.

$$b = 4B$$

$$C = 32B$$

$$N = 4$$

$$B = C/b = 32/4 = 8B$$

$$S = B/N = 8$$

$$\text{Set bits } \log_2(8) = 3$$

$$\text{Offset bits } \log_2(4) = 2$$

Tag bits = Address - Set bits - Offset

sets	V	Tag	00	01	10	11
000						
001						
010						
011						
100						
101	1	000000	00	00	00	11
110						
111						

32 block capacity.

0X014

0000 0001 0100  
Tag      Set      Offset

data at 014 in MM is

0000 0014

LRU (Least Recently Used) replacement

In 2-way associative cache, 1-bit is required to keep track of recently used data block.

- When block is accessed LRU is set to 0 and increment all other values accordingly
- When data is evicted from a block, replace the block with max value of LRU bit

### Example 2.

$$b = 2B$$

$$C = 16B$$

$$N = 2 \text{ way}$$

$$B = C/b = 16/2 = 8B$$

$$S = B/N = 8/2 = 4B$$

$$\text{Set bits} = 2$$

$$\text{Offset bits} = 1$$

$$\text{Address bits} = 3$$

$$\text{Tag bits} = 3 - 2 - 1 = 5 \text{ bits}$$

sets	V	LRU	Tags	0	T
00					
01					
10					
11	0	0	1111	FE	FF

e.g. 0xFF

1111 1111  
tag      set      offset

### Example 3.

$$b = 1B$$

Full Associate cache.

Singlset with 4ways.

Address bits 8bits

$$C = 1 \times 4 = 4$$

$$B = C/b = 4/1 = 4$$

$$\text{Set} = 4/1 = 2$$

$$\text{Set bits} = 0$$

$$\text{Offset bits} = 0$$

$$\text{Tag} = 8 - 0 - 0$$

	V	LRU	Tag	Data
0				
1				
2				
3				

## Type of Cache Miss

- Compulsory Miss: The first request to cache block is called compulsory miss, because the block must be read from memory regardless of cache design.  
↳ at first cache is empty so ↳

- Capacity Miss: occurs when the cache is too small to hold all concurrently used data
- Conflict Miss: caused when several addresses map to the same set and evict blocks that are still needed

## Cache Parameters.

1. Increasing cache capacity reduces capacity + conflict miss but doesn't effect compulsory miss.
2. Increasing block size could reduce compulsory miss (due to spatial locality) But as block size increases, the number of sets in fixed-size cache decreases increases probability of conflicts.
3. N-way associative cache reduces conflicts by providing N blocks in each set where data mapping to that set might be found Each memory address still maps to a specific set but it can map to any one of N-blocks in the set.  
\* Associative caches are usually slower and expensive

## Locality

1. Temporal Locality: if data is used recently and is likely to be used again Such data is kept in higher levels of memory
2. Spatial Locality: if data used nearby data may be used eg array values. bring nearby data of the data currently accessed in higher levels of memory

$$\text{Miss Rate} = \frac{\text{Misses}}{\text{Total memory accesses}} = 1 - \text{Hit Rate}$$

$$\text{Hit Rate} = \frac{\text{Hits}}{\text{Total memory accesses}}, 1 - \text{Miss Rate}$$

## Cache Write Policies.

### - Write Through

- when cache memory is updated simultaneously main memory is also update.
- at any time data in cache and main memory is same.

### - Write Back

- during write operation only cache location is updated
- when update occurs in cache this is marked by a flag also known as modified/dirty bit
- when word is replaced from cache based on value of dirty bit, the data is written to main memory.

### - Write Around (aka write no allocate).

- 'the write operation directly goes to main memory without affect the cache

### - Write Allocate.

- load newly written data to cache and if that data is needed soon it will be available in cache.

