

Chapter 3: Memory management

Operatieng System

@ptit.edu.vn



Posts and Telecommunications
Institute of Technology
Faculty of Information Technology 1



IT1 Faculty
PTIT

August 15, 2023

Mục lục

Address and related issues

Logical address and physical address

Some ways to organize the program

Loading during runtime

Dynamic linking and shared libraries

Memory management requirements

Reissue

Protect

Share

Logical structure and physical structure

Memory division

Content of chapter 3

1. Address and related issues
2. Some ways to organize the program
3. Memory management requirements
4. Memory chaptering
5. Memory paging
6. Memory segment
7. Virtual memory

Memory management

Bộ nhớ *is the second most important resource after the CPU in a computer system, consisting of bytes and addressed memory words.*

- ▶ Memory is where the process and data for the process are stored
- ▶ Memory management affects the processing speed and computing power of the entire system
- ▶ Tasks related to memory management include:
 - Manage free memory
 - Allocate and release memory for processes
 - Prevent unauthorized access to memory areas
 - Address mapping between logical addresses and physical addresses

Address and related issues

- ▶ Computer memory is addressed in bytes or words

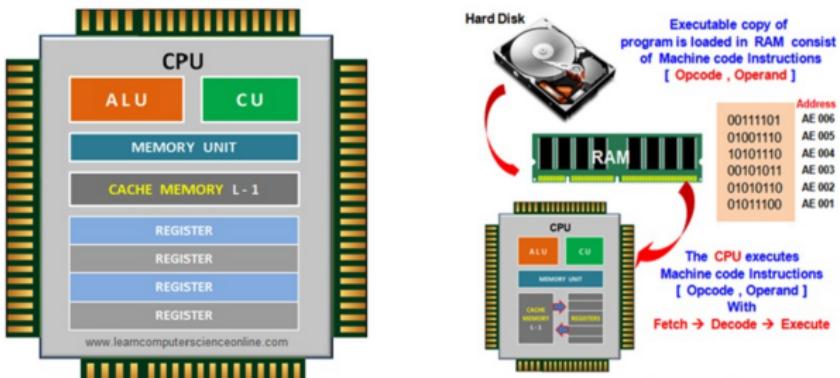
Dữ liệu hoặc lệnh	Địa chỉ byte (theo Hexa)	Dữ liệu hoặc lệnh	Địa chỉ word (theo Hexa)
byte (8-bit)	0x0000 0000	word (32-bit)	0x0000 0000
byte	0x0000 0001	word	0x0000 0004
byte	0x0000 0002	word	0x0000 0008
byte	0x0000 0003	word	0x0000 000C
byte	0x0000 0004	word	0x0000 0010
byte	0x0000 0005	word	0x0000 0014
byte	0x0000 0006	word	0x0000 0018
byte	0x0000 0007	.	
.		.	
.		.	
.			
byte	0xFFFF FFFF	word	0xFFFF FFF4
byte	0xFFFF FFFC	word	0xFFFF FFF8
byte	0xFFFF FFFD	word	0xFFFF FFFC
byte	0xFFFF FFFE		
byte	0xFFFF FFFF		

2^{32} bytes 2^{30} words

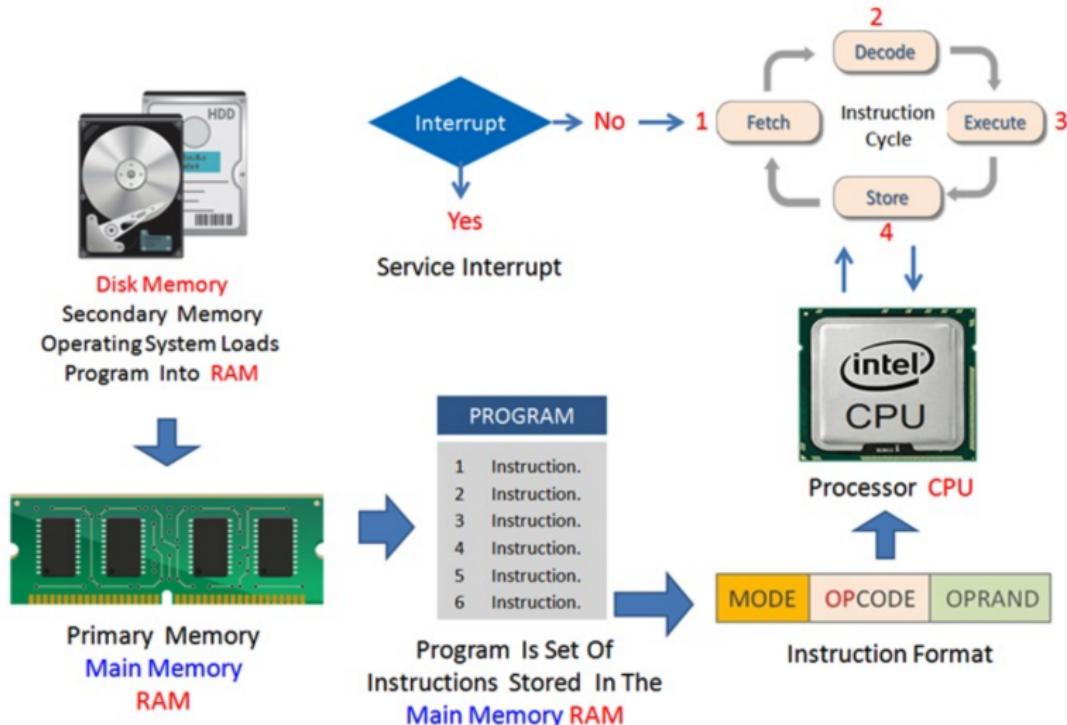
Computer Architecture

Address and related issues (cont.)

- ▶ The CPU process reads instructions from memory one by one and executes these instructions
- ▶ Data and instructions are both localised
- ▶ The CPU will use addresses to identify specific instructions and data



Address and related issues (cont.)



Address and related issues (cont.)

- ▶ Programs use addresses as names (variables, functions)
- ▶ When translating, the translation program maps names to relative addresses from the beginning of the obj file
- ▶ The link program maps addresses to relative addresses from the beginning of the program

The screenshot shows the Godbolt compiler explorer interface. On the left, there is a C++ source code editor with the following content:

```

1 // Type your code here, or load an example.
2 void square();
3 {
4     int A[2];
5     int B[3] = {1, 4, 2};
6     int C[3] = {2, 3, 4, 5, 6};
7     A[0] = 1;
8     A[2] = 3;
9     if(A[0] > C[4])
10    {
11        B[2] = C[4];
12        C[5] = B[2] + 1;
13        A[1] = B[0] + C[3];
14    }

```

On the right, there is an assembly output window titled "x86-64 gcc 12.2 (C++, Editor #1, Compiler #1)". The assembly code generated by the compiler is as follows:

```

1 square():
2     push    rbp
3     mov     rbp, rbp
4     mov     DWORD PTR [rbp-20], 1
5     mov     DWORD PTR [rbp-16], 4
6     mov     DWORD PTR [rbp-12], 2
7     mov     DWORD PTR [rbp-48], 2
8     mov     DWORD PTR [rbp-44], 3
9     mov     DWORD PTR [rbp-40], 4
10    mov    DWORD PTR [rbp-36], 5
11    mov    DWORD PTR [rbp-32], 6
12    mov    DWORD PTR [rbp-8], 1
13    mov    DWORD PTR [rbp+8], 3
14    mov    es, DWORD PTR [rbp-8]
15    mov    es, DWORD PTR [rbp-32]
16    cmp    es, es
17    jle    .L3
18    mov    es, DWORD PTR [rbp-32]
19    mov    DWORD PTR [rbp-12], es
20    mov    es, DWORD PTR [rbp-12]
21    add    es, 1
22    mov    DWORD PTR [rbp-28], es
23    mov    eds, DWORD PTR [rbp-20]
24    mov    es, DWORD PTR [rbp-36]
25    add    es, eds
26    mov    DWORD PTR [rbp-4], es
27 .L3:
28    pop    rbp
29    ret
30

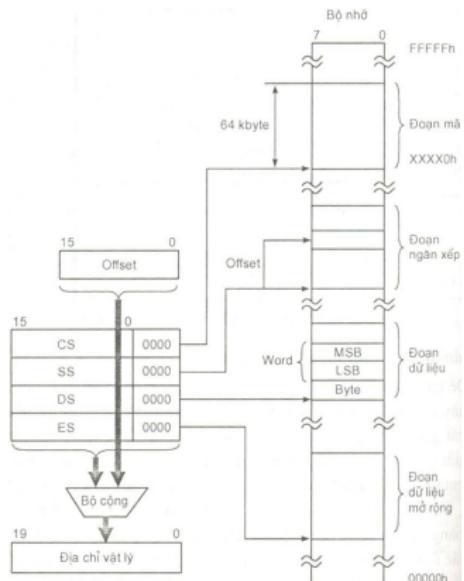
```

Godbolt

Address and related issues

Logical address and physical address

VXL 8086 uses 20 bits to encode 1MB addresses from 00000h - FFFFFh



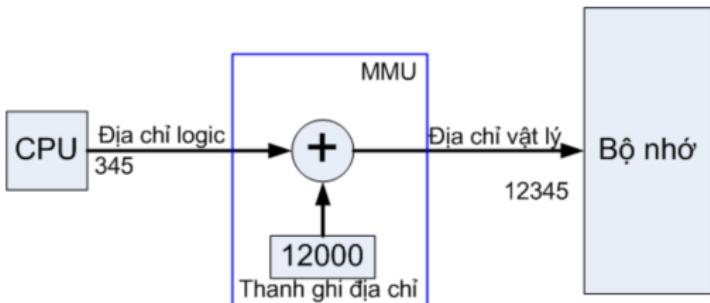
20-bit physical address of a memory cell = 16-bit segment address shifted 4 bits left (multiplied by 16) + 16-bit instruction address.

► Logical address

- Assignment to instructions and data does not depend on the specific location of the process in memory. Spawned in process (CPU launched)
- CTUD only cares about logical addresses and physical addresses are determined by the Operating System
- When executing a program, the CPU "sees" and uses this logical address to point to different parts of instructions and data.
- To access memory, logical addresses need to be mapped to physical addresses and converted to physical addresses by the memory management unit (MMU) when accessing objects in the program.

Address and related issues (cont.)

Logical address and physical address



For example: JMP 010A: Jump to memory location 010Ah at the same code (CS - Code Segment)

If CS = 1555h, will go to position: $1555h * 10h + 010Ah = 1560Ah$

► Physical address:

- Is the exact address in computer memory
- Memory circuits are used to access programs and data

Address and related issues (cont.)

Logical address and physical address



- ▶ Logical addresses are converted to physical addresses using the address mapping block. (MMU=Memory Mapping Unit)

Some ways to organize the program

Important issues in program organization and memory management:

- ▶ Reduce the space the program takes up on disk and memory
- ▶ Use memory space effectively
- ▶ Techniques for organizing and using memory effectively:
 - Loading during runtime
 - Dynamic linking and shared libraries

The entire program is normally loaded into memory for execution.

For large programs, during one session, some parts of the program may not be used. These functions will take up unnecessary memory and increase the initial program load time.

Solution:

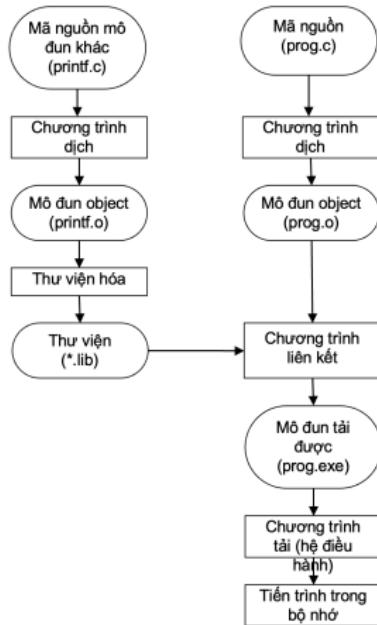
- ▶ A function that has not been called has not been loaded into memory
- ▶ The main program is loaded into memory and run
- ▶ When there is a function call:
 - The program will check if the function is loaded.
 - If not, the program will proceed to load then map the function addresses into the program's common space and change the address table to record those mappings.
- ▶ The programmer is responsible for checking and loading functions.
The operating system only provides library functions for the module

Some ways to organize the program (cont.)

Loading during runtime



- ▶ Programs must go through a translation and linking process before becoming downloadable and executable.



Some ways to organize the program

Dynamic linking and shared libraries



During static linking, functions and libraries are always linked into the program code

$\text{Program size} = \text{Size of newly compiled program} + \text{Size of library functions}$

- ▶ Functions will appear repeatedly in programs
- ▶ Waste of space both on disk and internal memory

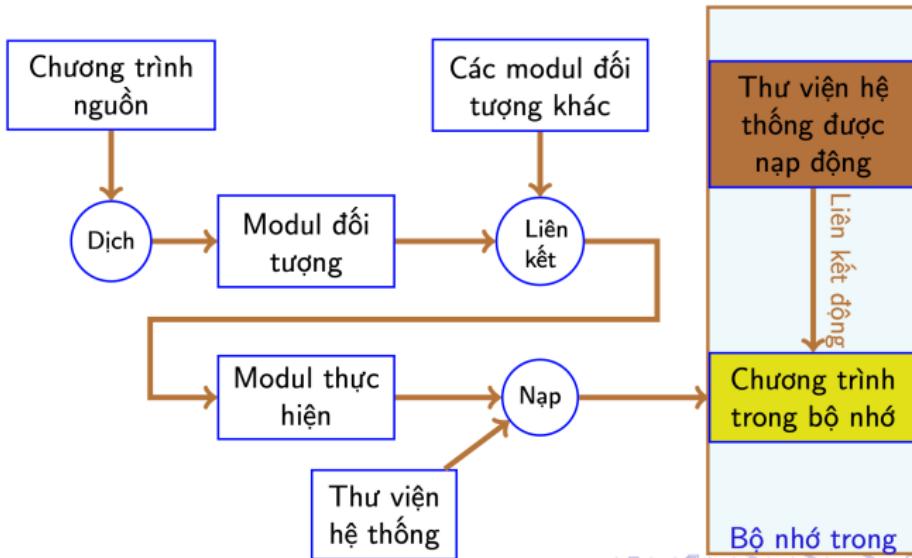
Solution use **dynamic linking technique** : *During the linking phase, do not connect library functions into the program but only insert information about that library function (stub).*

Some ways to organize the program (cont.)

Dynamic linking and shared libraries



- ▶ Programs must go through a translation and linking process before becoming downloadable and executable.



Some ways to organize the program (cont.)

Dynamic linking and shared libraries



► Library modules are linked in the implementation:

- The process does not keep copies of library modules, but the process keeps a small piece of code containing information about the library module
- When the small code is called, it checks whether the corresponding module is already in memory. If not, then download the rest and use it.
- The next time you need to use it, the library module will be run directly
- Each library module has only one copy contained in MEM
- Need support from OS

- ▶ Memory management:
 - **Memory allocation:** for processes
 - **Reposition,** maps instruction and data addresses from the process into memory
 - **Protect the process:** Some memory for one process is not used by another process
 - **Share:** Processes can share memory spaces with each other

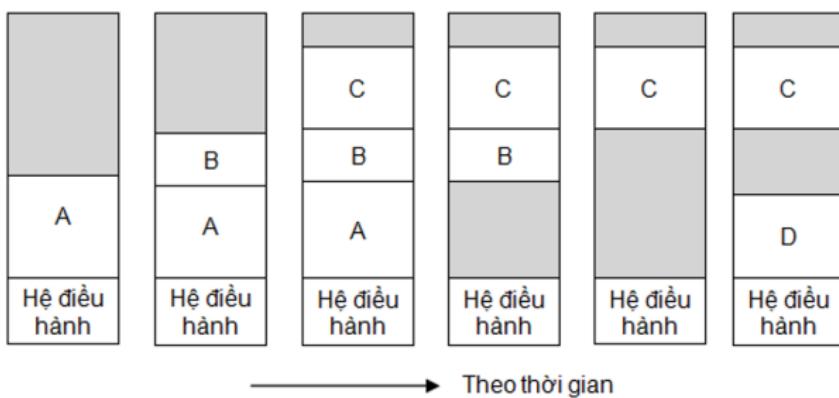
Requirements related to memory management include:

- ▶ **Reallocation** to processes
- ▶ **Protection:** any memory for any process cannot be used by another process
- ▶ **Shared:** processes can share memory spaces when interacting with each other
- ▶ **Structure:** logical structure and physical structure

Memory management requirements

Reissue

- ▶ It is necessary to be able to swap processes in and out of MEM to maximize processor utilization
- ▶ You cannot request that a process be moved back into MEM, it must be in the same place it was used before being moved out



- ▶ Each process must be protected from unwanted references from other processes to its reserved memory
- ▶ When a memory area is used by this process, it does not allow another process to reference that memory area.
- ▶ All memory references of a process must be checked at runtime
- ▶ The operating system cannot predict every MEM reference, so the microprocessor hardware takes care of it

- ▶ Multiple processes need and are allowed to access the same memory area
- ▶ Collaborating processes need to share access to a data structure
- ▶ Allow access to shared areas

Logical structure:

- ▶ MEM is structured in a linear way consisting of bytes, while the program is organized into modules
- ▶ Must respond to:
 - Modules can be written and interpreted independently
 - The level of protection may vary
 - Modules can be shared between processes

Physical structure:

- ▶ 2-level physical structure:
 - Main memory: fast; High cost, low capacity
 - Secondary memory: large capacity, allowing to save programs and data for a long time

Logical structure and physical structure (cont.)

Share



- ▶ The system is responsible for converting information between two levels

Memory chaptering

- ▶ To execute the process, the Operating System needs to allocate the necessary memory space to the process.
- ▶ Memory allocation and management is an important function of the operating system
- ▶ The simplest allocation technique is for each process to be allocated a contiguous area of memory
- ▶ The operating system divides memory into continuous parts called partitions, each process will be provided with a chapter to contain its instructions and data.
- ▶ The process of dividing memory into chapters is called *memory chaptering*.
- ▶ Depending on the choice of chapter location and size, fixed and dynamic chaptering can be distinguished

Memory chaptering (cont.)

Memory management strategies:

- ▶ Fixed chapter division
- ▶ Dynamic chapter division
- ▶ Pagination
- ▶ Segment
- ▶ Combination of segmentation - paging

Chapter 3 Memory management

- ▶ Address and related issues
- ▶ Some ways to organize the program
- ▶ Memory management requirements

Chapter 3 Memory management

- ▶ Memory chaptering
- ▶ Memory paging
- ▶ Memory segmentation
- ▶ Virtual memory