ARCOS Group

uc3m | Universidad **Carlos III** de Madrid

# Lesson 5 (III)
# Memory hierarchy

Computer Structure

Bachelor in Computer Science and Engineering
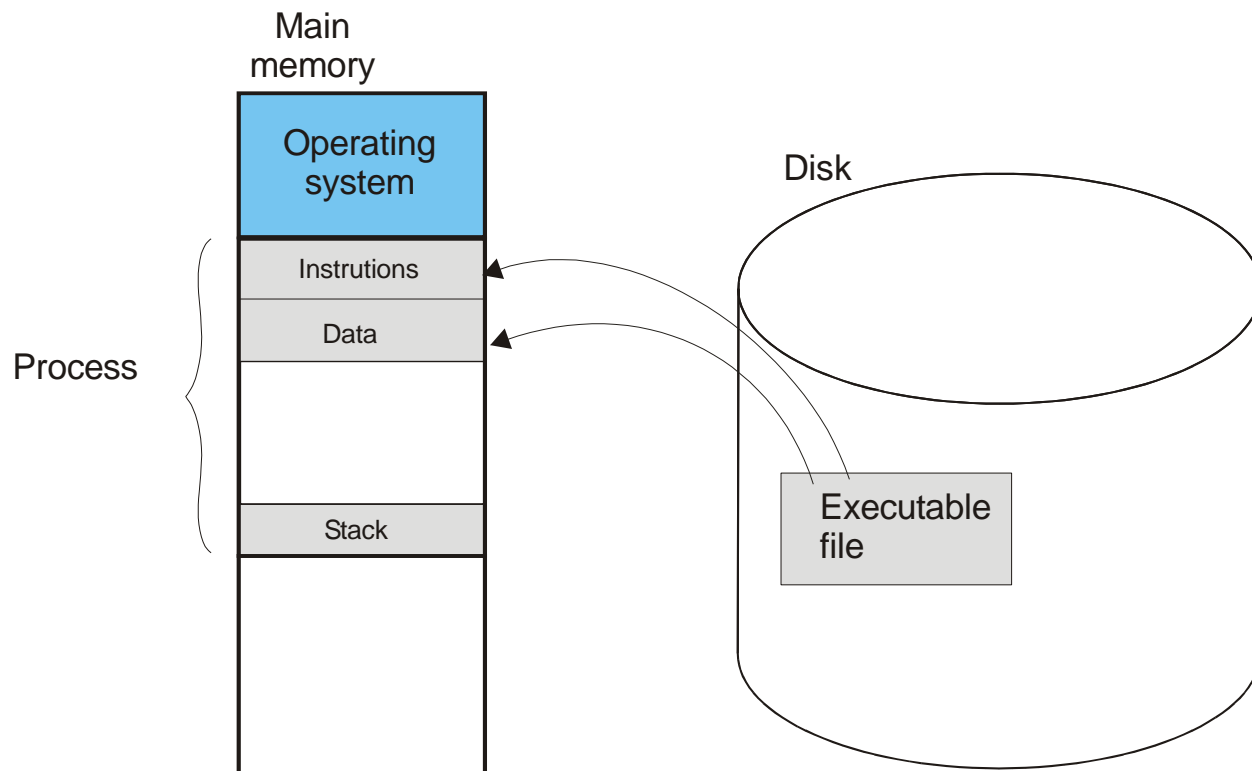
# Contents

1. Types of memories

2. Memory hierarchy
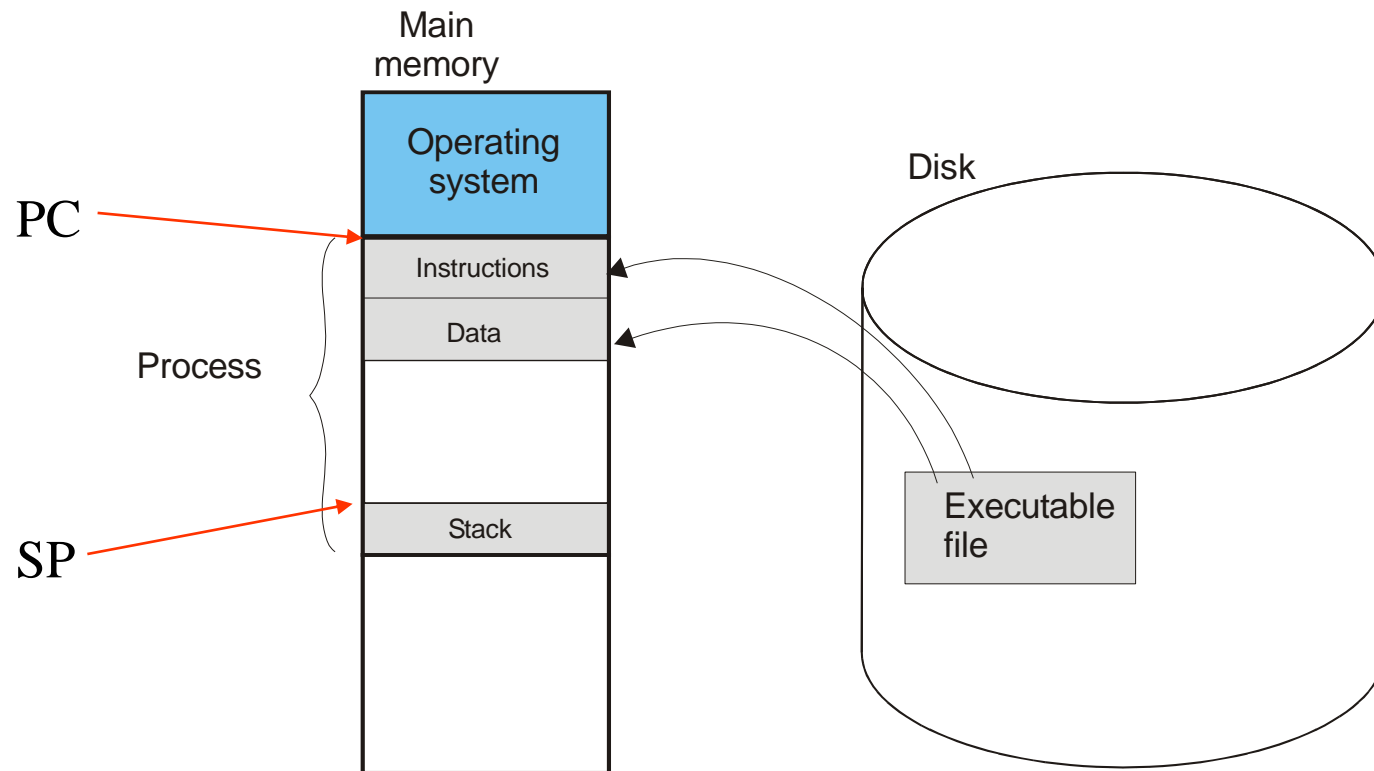
3. Main memory

4. Cache memory

5. Virtual memory

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Systems without virtual memory

▸ In systems without virtual memory, the program is completely loaded in memory before the execution

Main memory

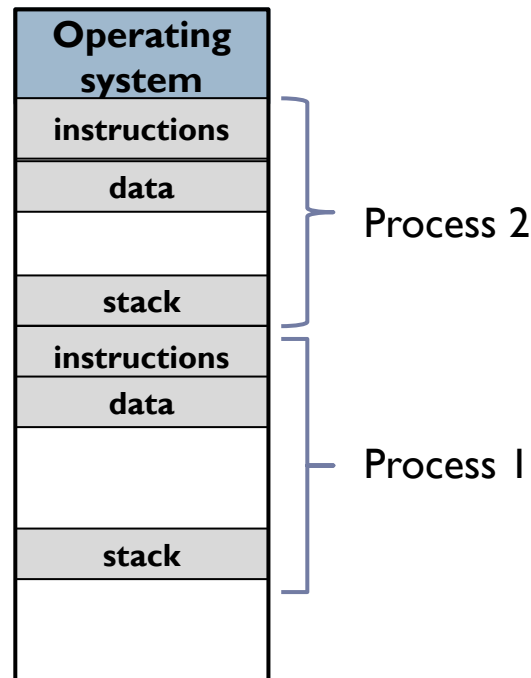| Operating system |
|---|
| Instrutions |
| Data |
| |
| Stack |
| |

Process

Disk

Executable file

# Systems without virtual memory

▶ Registers are initialized

# Multiple programs loaded in memory

Main memory

| |
|:---:|
| **Operating system** |
| **instructions** |
| **data** |
| |
| **stack** |
| **instructions** |
| **data** |
| |
| **stack** |
| |

Process 2

Process 1

# Multiple programs loaded in memory

Main memory

| Operating system |
|:---:|
| instructions |
| data |
| |
| stack |
| instructions |
| data |
| |
| stack |
| |

PC →

SP →

Process 2

Process 1

# Multiple programs loaded in memory

Main memory

PC

SP

| Operating system |
|---|
| instructions |
| data |
| |
| stack |
| instructions |
| data |
| |
| stack |
| |

Process 2

Process 1

# Multiple programs loaded in memory

Main memory

| Operating system |
|:---:|
| **instructions** |
| **data** |
| |
| **stack** |
| **instructions** |
| **data** |
| |
| **stack** |
| |

Process 2

Process 1

PC

SP

# Hypothetical executable file

```
int v[1000];  // global
int i;
for (i=0; i < 1000; i++)
    v[i] = 0;
```

Félix García Carballeira, Alejandro Calderón Mateos

# Hypothetical executable file

```
int v[1000];   // global
int i;
for (i=0; i < 1000; i++)
    v[i] = 0;
```

```
.data:
        v: .space 4000
.text:  li    $t0, 0
        li    $t1, 0
        li    $t2, 1000
loop:   bgt   $t0,  $t2, end
        sw    $0, v($t1)
        addi $t0, $t0, 1
        addi $t1, $t1, 4
        b     loop
end:        …
```

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Hypothetical executable file

```
int v[1000];  // global
int i;
for (i=0; i < 1000; i++)
    v[i] = 0;
```

assembly

```
.data:
        v: .space 4000
.text:  li   $t0, 0
        li   $t1, 0
        li   $t2, 1000
loop:   bgt  $t0,  $t2, end
        sw   $0, v($t1)
        addi $t0, $t0, 1
        addi $t1, $t1, 4
        b    loop
end:        …
```

executable

| | |
|---|---|
| 0 | |
| 4 | |
| | Header |
| 96 | |
| 100 | li   $t0, 0 |
| 104 | li   $t1, 0 |
| 108 | li   $t2, 1000 |
| 112 | bgt  $t0, $t2, 16 |
| 116 | sw   $0,  2000($t1) |
| 120 | addi $t0, $t0, 1 |
| 124 | addi $t1, $t1, 4 |
| 128 | b    -20 |
| 132 | |
| . | . |
| . | . |
| . | . |

Address 2000 is assigned to v
Assumes that program starts in address 0

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Loading the program in memory

▶ The Operating System reserves a contiguous free portion in memory for the entire process image

executable

| | |
|---|---|
| 0 | |
| 4 | Header |
| 96 | |
| 100 | li    $t0, 0 |
| 104 | li    $t1, 0 |
| 108 | li    $t2, 1000 |
| 112 | bgt   $t0, $t2, 20 |
| 116 | sw    $0,  2000($t1) |
| 120 | addi  $t0, $t0, 1 |
| 124 | addi  $t1, $t1, 4 |
| 128 | b     -24 |
| 132 | |

memory

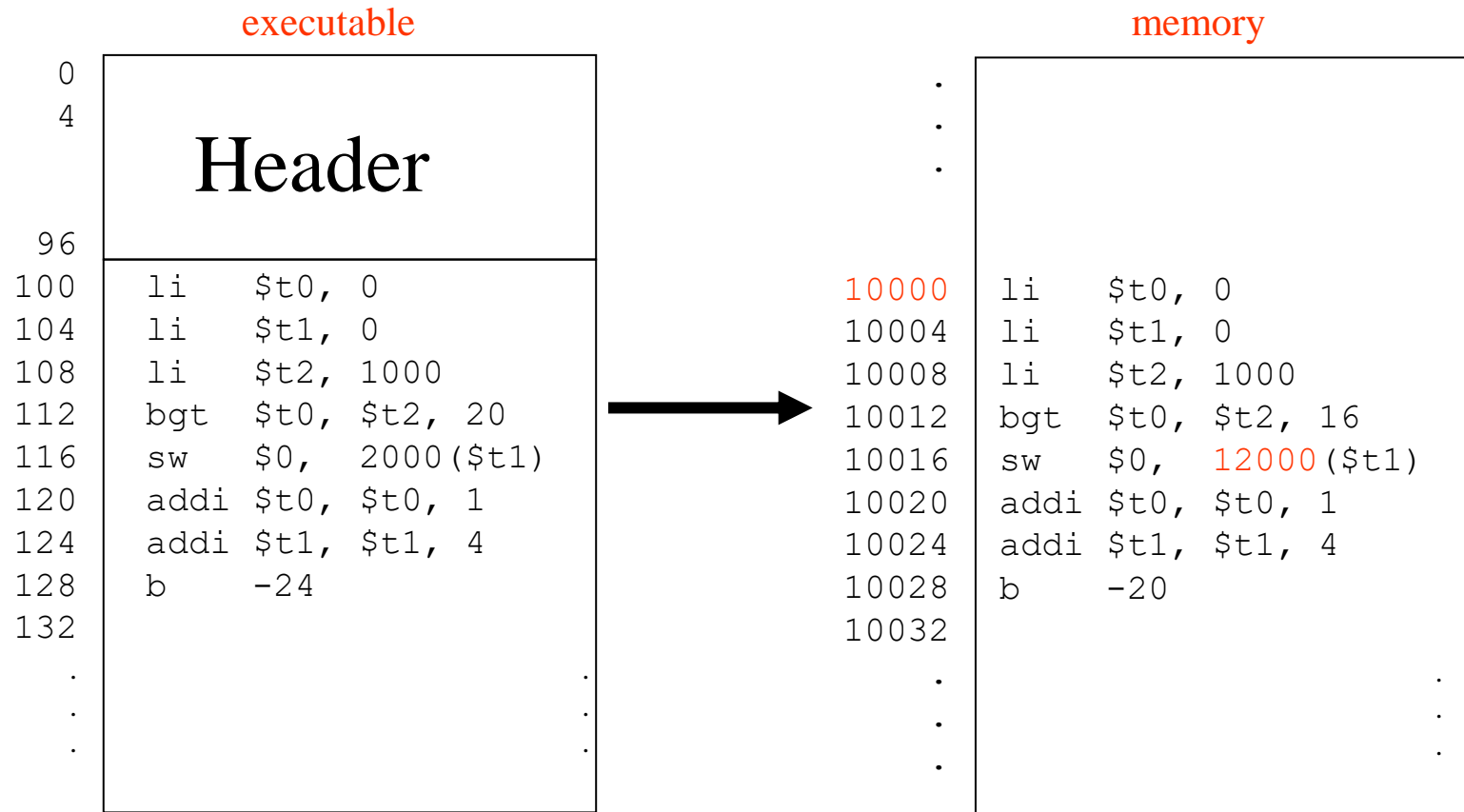| |
|---|
| . |
| . |
| . |
| 10000 |
| 10004 |
| 10008 |
| 10012 |
| 10016 |
| 10020 |
| 10024 |
| 10028 |
| 10032 |
| . |
| . |
| . |

# Loading the program in memory

- In the executable file the address 0 is considered as the init address
  - Logical address
- In memory, the init address is10000
  - Physical address
- Address translation is needed
  - From logical address to physical
- The array in memory is in:
  - The logical address 2000
  - The physical address 2000 + 10000
- This process is called relocation
  - Software relocation
  - Hardware relocation

# Software relocation

▶ Occurs in the loading process

executable

|      |                      |
|------|----------------------|
| 0    |                      |
| 4    | **Header**           |
| 96   |                      |
| 100  | li    $t0, 0         |
| 104  | li    $t1, 0         |
| 108  | li    $t2, 1000      |
| 112  | bgt   $t0, $t2, 20   |
| 116  | sw    $0,  2000($t1) |
| 120  | addi  $t0, $t0, 1    |
| 124  | addi  $t1, $t1, 4    |
| 128  | b     -24            |
| 132  |                      |

memory

|        |                      |
|--------|----------------------|
| 10000  | li    $t0, 0         |
| 10004  | li    $t1, 0         |
| 10008  | li    $t2, 1000      |
| 10012  | bgt   $t0, $t2, 16   |
| 10016  | sw    $0,  12000($t1)|
| 10020  | addi  $t0, $t0, 1    |
| 10024  | addi  $t1, $t1, 4    |
| 10028  | b     -20            |
| 10032  |                      |

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Software relocation

▶ What happens with the instructions loaded in 10012 and 10028 addresses?

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Software relocation

▶ What happens with the instructions loaded in 10012 and 10028 addresses?

executable

```
 0
 4          Header
96
100    li    $t0, 0
104    li    $t1, 0
108    li    $t2, 1000
112    bgt   $t0, $t2, 20
116    sw    $0,  2000($t1)
120    addi  $t0, $t0, 1
124    addi  $t1, $t1, 4
128    b     -24
132
```

memory

```
10000  li    $t0, 0
10004  li    $t1, 0
10008  li    $t2, 1000
10012  bgt   $t0, $t2, 16
10016  sw    $0,  12000($t1)
10020  addi  $t0, $t0, 1
10024  addi  $t1, $t1, 4
10028  b     -20
10032
```

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Problem with memory protection

▸ What happens if the program executes these instructions?

```
li $t0 ,8
sw $t0, ($0)
```
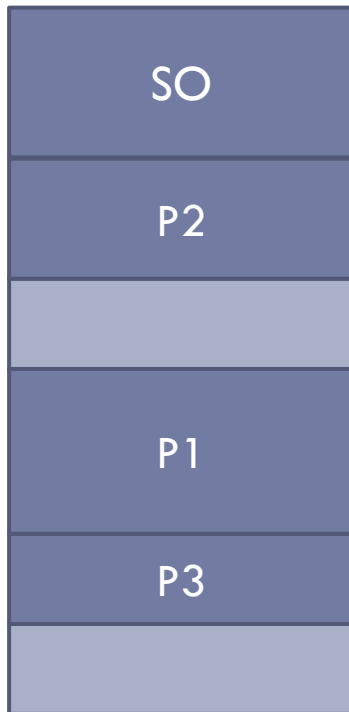
# Problem with memory protection

▸ What happens if the program executes these instructions?

```
li $t0 ,8
sw $t0, ($0)
```

Illegal access to physical address 0 that is not assigned to the program
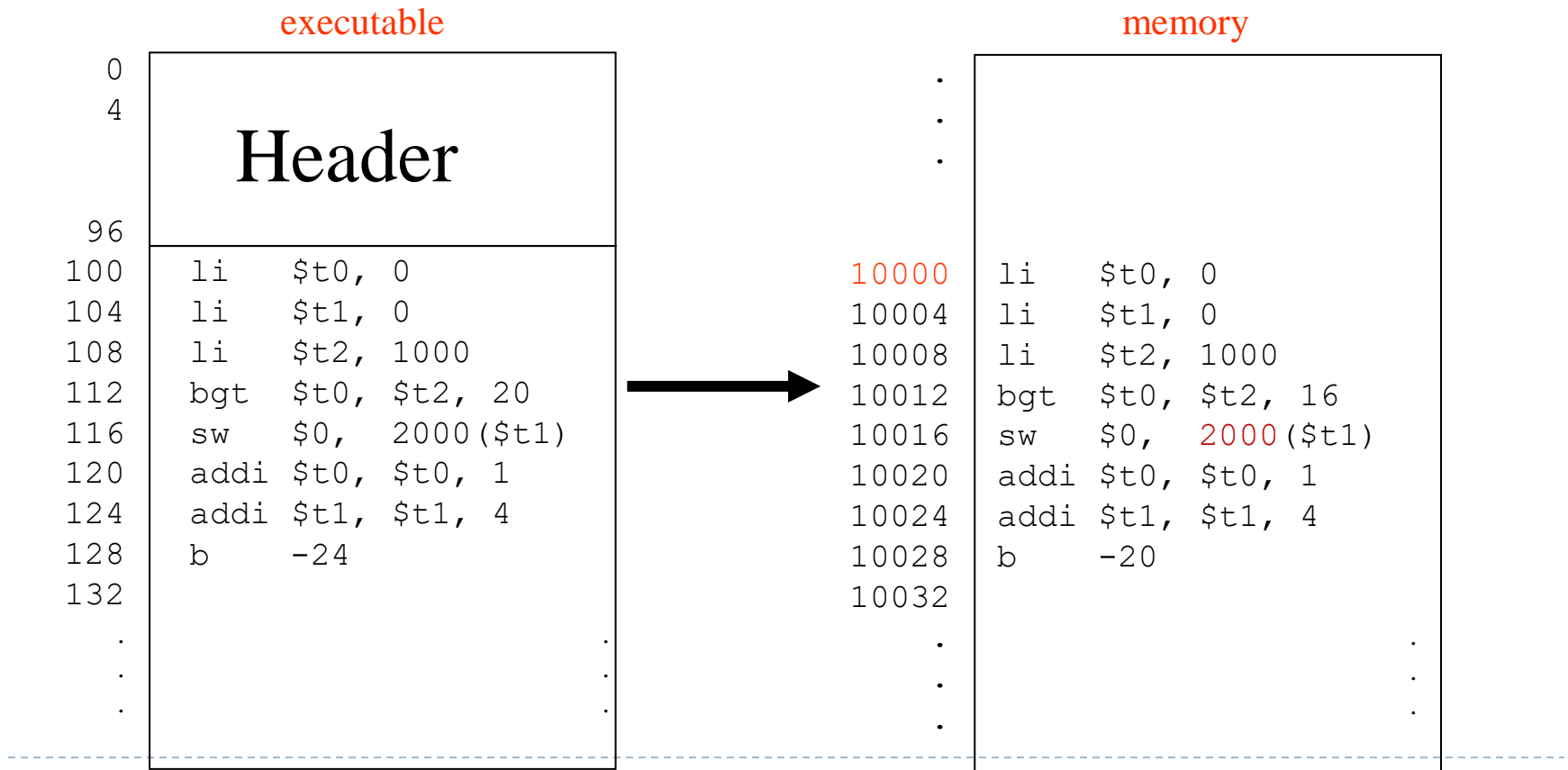
# Multiprogramming

▶ A computer can store several programs in memory

▶ Each program needs an address space in memory

| SO |
|----|
| P2 |
|    |
| P1 |
| P3 |
|    |

We need to ensure that a program does not access to the address space of other program

# Hardware relocation

▸ The translation occurs in the execution

▸ Special HW is needed. Ensure protection



executable

| | |
|---|---|
| 0 | |
| 4 | Header |
| 96 | |
| 100 | li    $t0, 0 |
| 104 | li    $t1, 0 |
| 108 | li    $t2, 1000 |
| 112 | bgt   $t0, $t2, 20 |
| 116 | sw    $0,  2000($t1) |
| 120 | addi  $t0, $t0, 1 |
| 124 | addi  $t1, $t1, 4 |
| 128 | b     -24 |
| 132 | |

memory

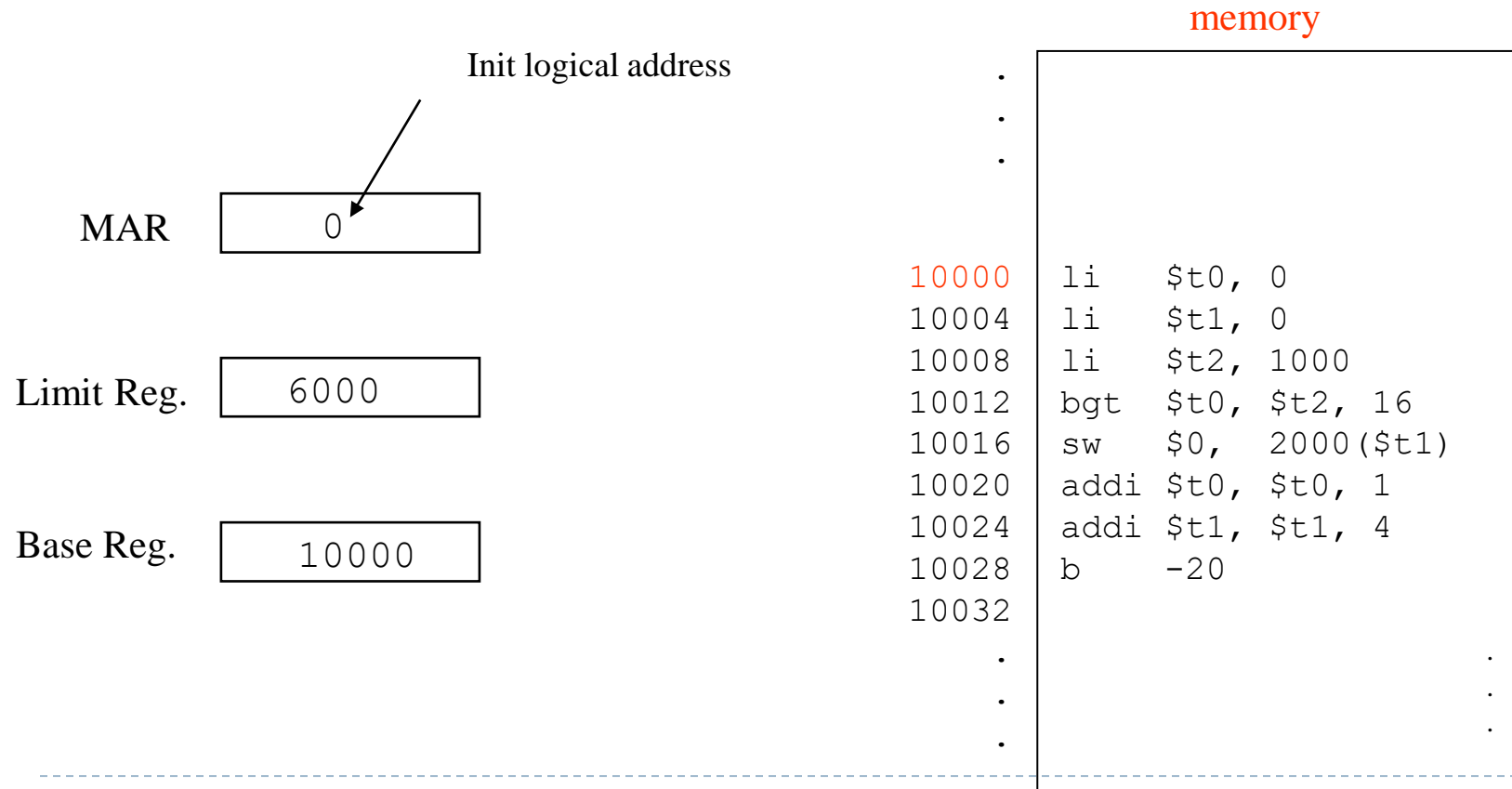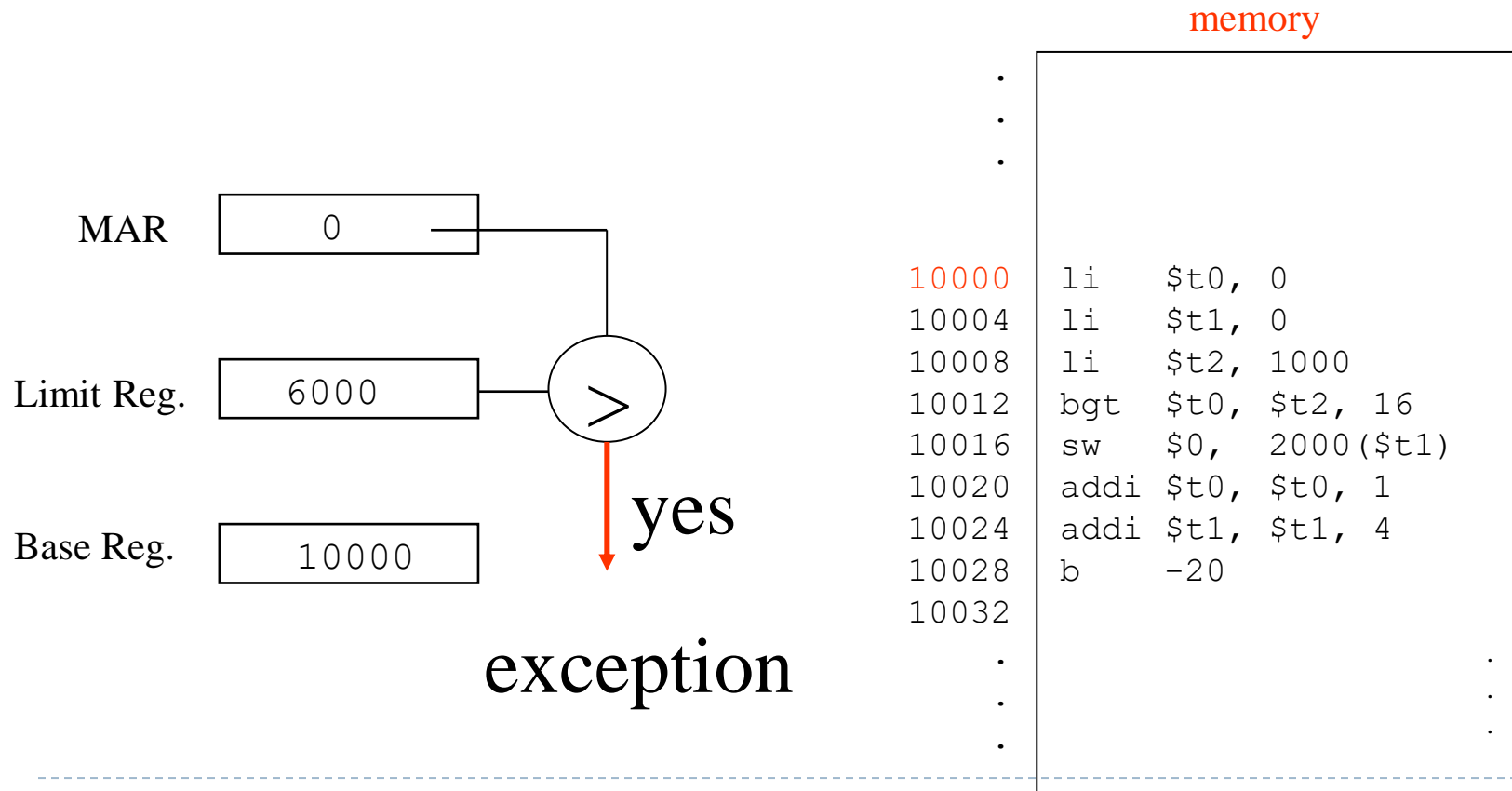| | |
|---|---|
| 10000 | li    $t0, 0 |
| 10004 | li    $t1, 0 |
| 10008 | li    $t2, 1000 |
| 10012 | bgt   $t0, $t2, 16 |
| 10016 | sw    $0,  2000($t1) |
| 10020 | addi  $t0, $t0, 1 |
| 10024 | addi  $t1, $t1, 4 |
| 10028 | b     -20 |
| 10032 | |

# Example of hardware support

▶ Limit register: maximum logical address assigned to the program

▶ Base register: program init address in memory

memory

MAR

Limit Reg.

Base Reg.

```
10000  li    $t0, 0
10004  li    $t1, 0
10008  li    $t2, 1000
10012  bgt   $t0, $t2, 16
10016  sw    $0,  2000($t1)
10020  addi  $t0, $t0, 1
10024  addi  $t1, $t1, 4
10028  b     -20
10032
```
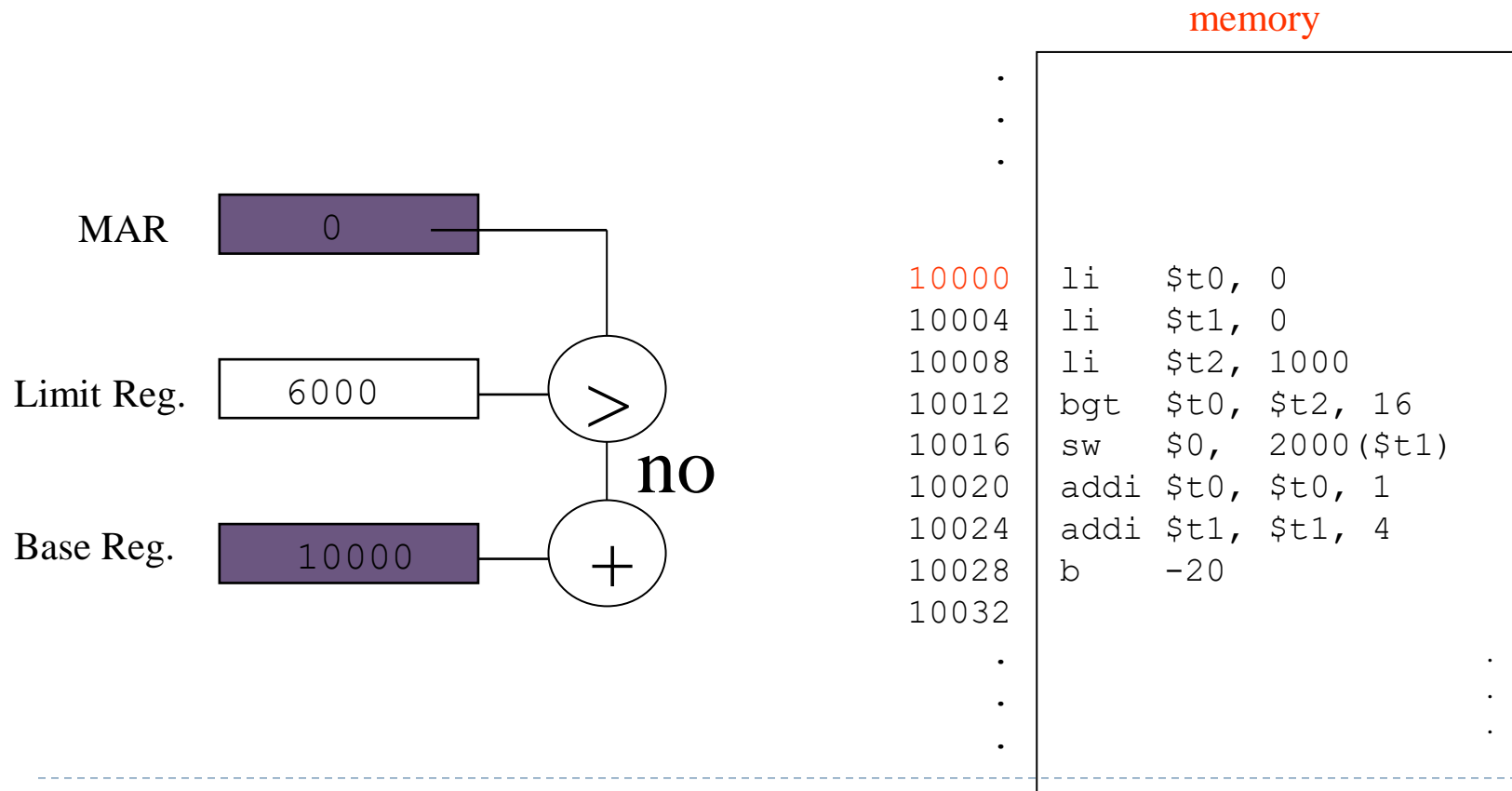
# Example of hardware support

▸ Limit register: maximum logical address assigned to the program
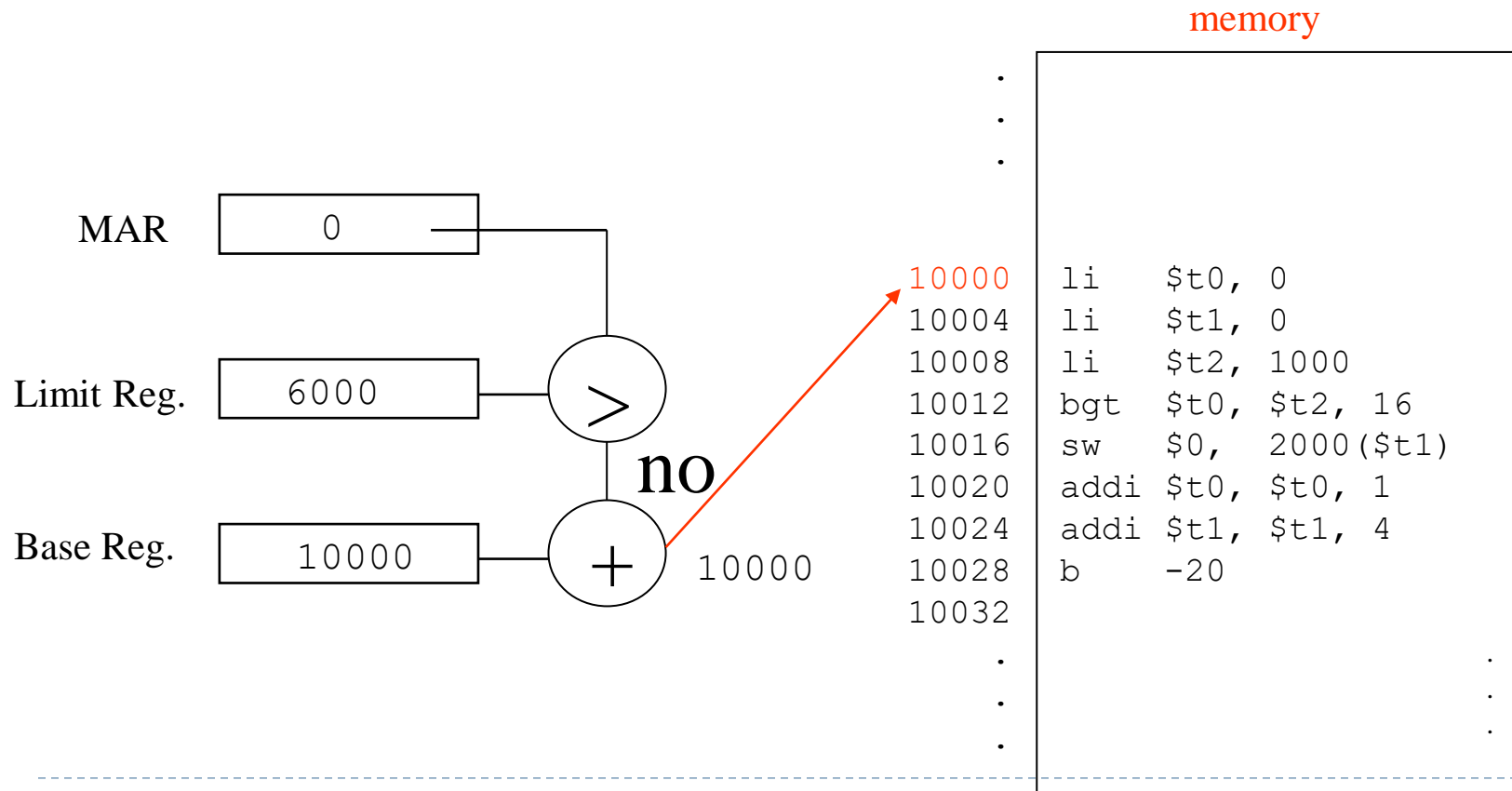
▸ Base register: program init address in memory

Init logical address

MAR    0

Limit Reg.    6000

Base Reg.    10000

memory

```
          .
          .
          .
10000 | li    $t0, 0
10004 | li    $t1, 0
10008 | li    $t2, 1000
10012 | bgt   $t0, $t2, 16
10016 | sw    $0,  2000($t1)
10020 | addi  $t0, $t0, 1
10024 | addi  $t1, $t1, 4
10028 | b     -20
10032 |
          .                    .
          .                    .
          .                    .
```

# Example of hardware support

▸ Limit register: maximum logical address assigned to the program
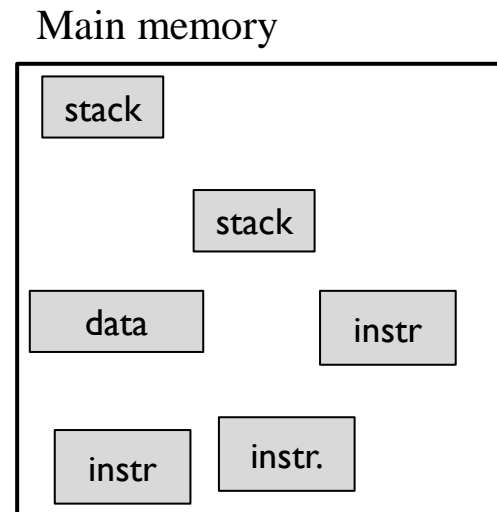
▸ Base register: program init address in memory

memory

MAR | 0

Limit Reg. | 6000

Base Reg. | 10000

> yes

exception

```
10000  li    $t0, 0
10004  li    $t1, 0
10008  li    $t2, 1000
10012  bgt   $t0, $t2, 16
10016  sw    $0,  2000($t1)
10020  addi  $t0, $t0, 1
10024  addi  $t1, $t1, 4
10028  b     -20
10032
```

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Example of hardware support

▶ Limit register: maximum logical address assigned to the program
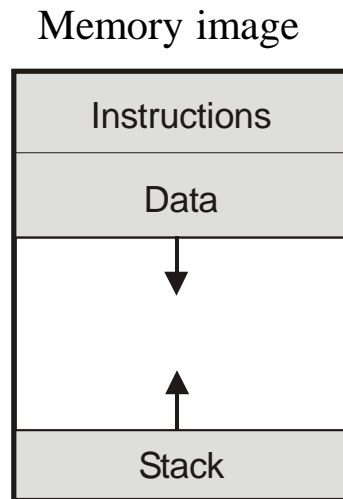
▶ Base register: program init address in memory

memory

```
       .
       .
       .

10000 │ li    $t0, 0
10004 │ li    $t1, 0
10008 │ li    $t2, 1000
10012 │ bgt   $t0, $t2, 16
10016 │ sw    $0,  2000($t1)
10020 │ addi  $t0, $t0, 1
10024 │ addi  $t1, $t1, 4
10028 │ b     -20
10032 │
       .              .
       .              .
       .              .
```

MAR  `0`

Limit Reg.  `6000`

`>`

no

Base Reg.  `10000`

`+`

# Example of hardware support

▶ Limit register: maximum logical address assigned to the program

▶ Base register: program init address in memory

memory

|  | MAR | 0 |
|---|---|---|

| | Limit Reg. | 6000 |

| | Base Reg. | 10000 |

\>

no

\+        10000

```
10000   li    $t0, 0
10004   li    $t1, 0
10008   li    $t2, 1000
10012   bgt   $t0, $t2, 16
10016   sw    $0,  2000($t1)
10020   addi  $t0, $t0, 1
10024   addi  $t1, $t1, 4
10028   b     -20
10032
```

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Systems without virtual memory
## Main problems

▶ **If the process image is bigger than the available memory, the process can not be executed**

▶ **In a 32-bit computer:**

  ▶ What is the theoretical maximum size of a program?

  ▶ What if this size if the memory has 512 MB?

▶ **The number of active programs is reduced**

# Virtual memory

- ▶ It is not needed to load the entire process in memory
- ▶ Only the program portions needed are loaded in memory
- ▶ Main advantages:
  - ▶ We can execute a program bigger than the main memory available
  - ▶ More programs can be active in memory

Memory image

| Instructions |
|:---:|
| Data |
| ↓ |
| ↑ |
| Stack |

Main memory

stack

stack

data

instr

instr

instr.

Félix García Carballeira, Alejandro Calderón Mateos

# Main concepts on virtual memory

Virtual memory uses:
- ❑ Main memory
- ❑ Disk

**Mvirtual memory map**
(addresses generate by the program)

Physical address
(hit)

**Main memory**

**Processor**

Virtual
address

**MMU**

Page fail

OS transfers the
page to memory

**Disk**

**(Swap)**

# Pages virtual memory

▶ Processors generate <span style="color:red">virtual addresses</span>

▶ The virtual address space is divided in equal size blocks called <span style="color:red">pages</span>

▶ Main memory is divided in equal size blocks called <span style="color:red">page frames</span>

▶ The part of the disk that supports the virtual memory is divided in equal size blocks called <span style="color:red">swap pages</span>

# Physical address and virtual address

- ## Virtual address space:
  - Memory addresses that use the processor.
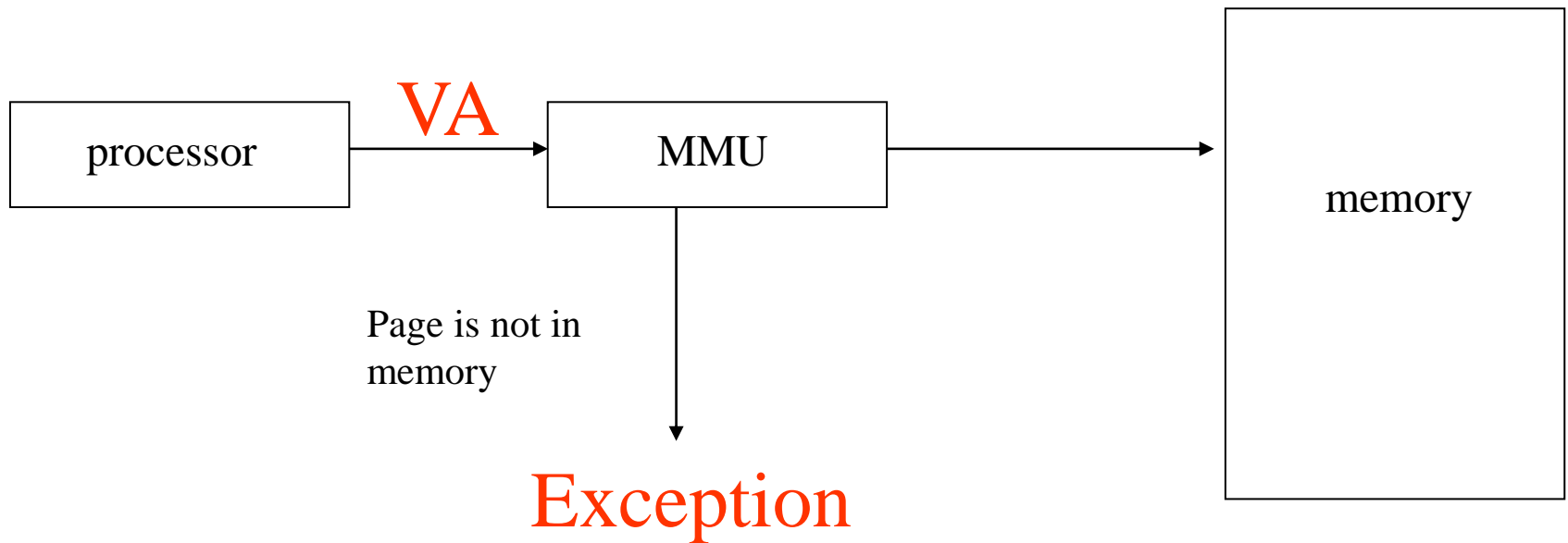- ## Physical address space:
  - Main memory addresses

**Virtual address** **Physical address**

Program → HW translation → Physical memory

# Paged virtual memory

▸ **The memory image of the programs are stored in disk**

# Address translation

processor → **VA** → MMU → **PA** → memory

Page is memory

# Address translation

processor → **VA** → MMU → memory

MMU → Page is not in memory → **Exception**

# Address translation

# Address translation

processor — **VA** → MMU → memory

**Operating systems**

Disk (swap)

# Address translation



processor → **VA** → MMU → **PA** → memory

Page is in memory

# Address translation

processor → VA → MMU → PA → memory

The data is sent to the processor

# Paged virtual address

32 bits

**011010001001010**

**Virtual address**

32 bits

**1111100100000**

**Physical address**

blocks with the same size

Disk (swap)

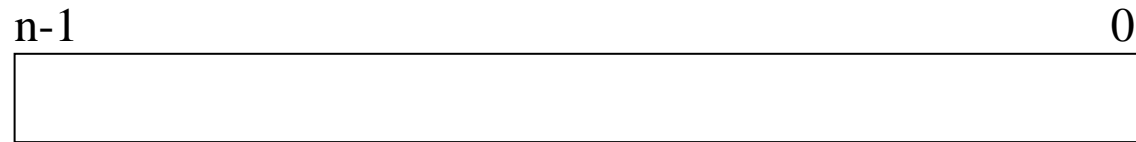# Structure of a virtual address

- ## A n bit computer has:

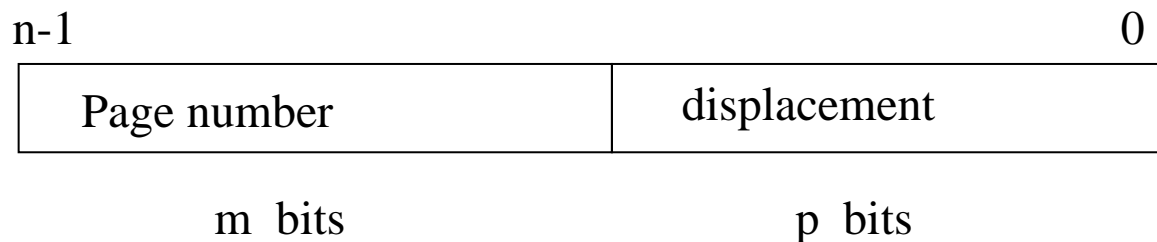  - Addresses of n bits

  n-1                                                              0

  | |
  |---|

  - Can address $2^n$ bytes

# Structure of a virtual address

▶ Memory image consists of pages with the same size(4 KB, 8 KB)

| | |
|---|---|
| n-1 | 0 |
| Page number | displacement |
| m bits | p bits |

▶ n = m + p
▶ Addressable memory: $2^n$ bytes
▶ Page size: $2^p$ bytes
▶ Maximum number of pages: $2^m$

# Example

32 bits

**0 1 1 0 1 0 0 0 1 0 0 1 0 1 0**

**Virtual address**

32 bits

**1 1 1 1 1 0 0 1 0 0 0 0 0**

**Physical address**

Disk (swap)

# Example

20 bits      12 bits

| Page id. | Displacement |
|----------|--------------|

**Virtual address**

20 bits      12 bits

| Frame | displacement |
|-------|--------------|

**Physical address**

With pages of 4 KB

**Frame**

Disk (swap)

# Exercise

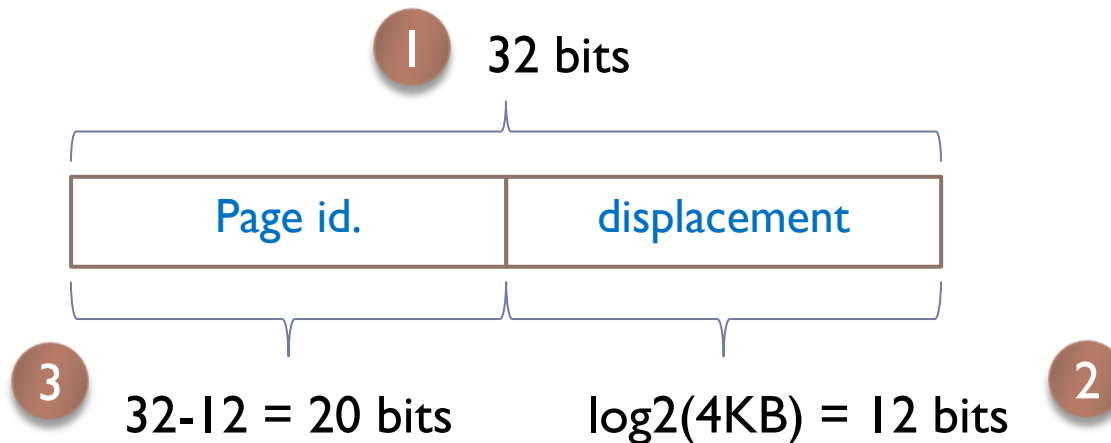▶ **A 32 bit computer has a memory of 512 MB and pages of 4 KB**

▶ **Answer:**

a) Indicate the format of a virtual address and the number of page frames

# Solution

- **Virtual address format:**

① 32 bits

| Page id. | displacement |
|----------|-------------|

③ 32-12 = 20 bits          log2(4KB) = 12 bits ②
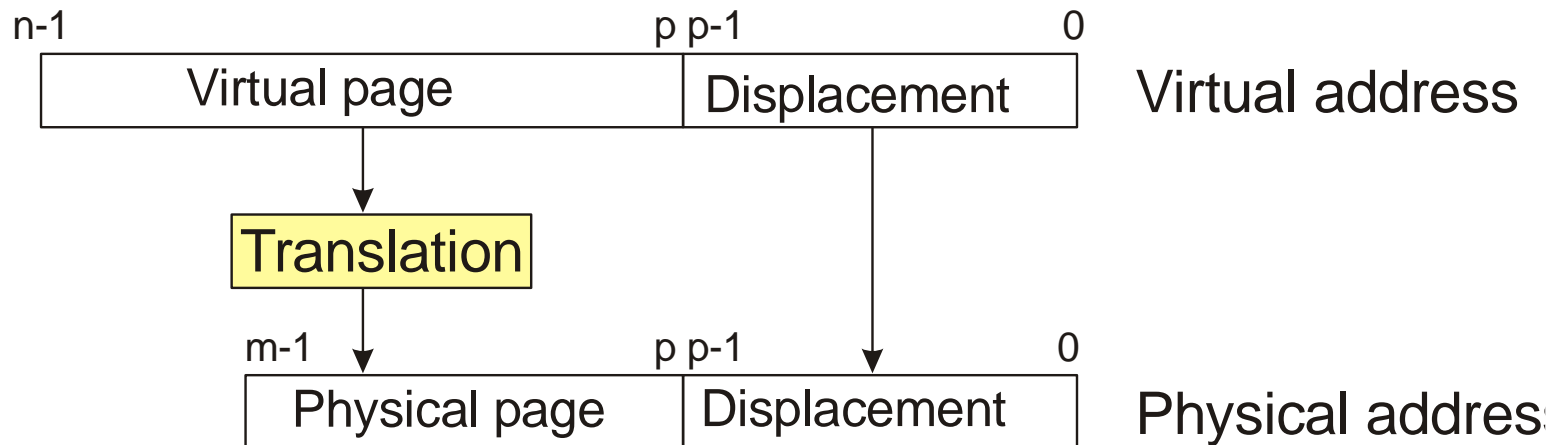
### Number of page frames

$$\frac{\text{Main memory size}}{\text{Page size}} \quad \frac{512 \text{ MB}}{4 \text{ KB}} = \frac{512 * 2^{20}}{4 * 2^{10}} = 128 * 2^{10}$$

# Address translation

# Page table

Page table

Present

ETP | Yes

No

Virtual address space
of a program

Main
memory

Disk

| Page 0 |
| Page 1 |
| Page 2 |
| Page 3 |
| Page 4 |
| Page 5 |

Page
frame

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Page table entry

Modified

Protection

Frame number/
disk block

Present

# Page table structure

- Operating system creates the page table when a program is going to be executed

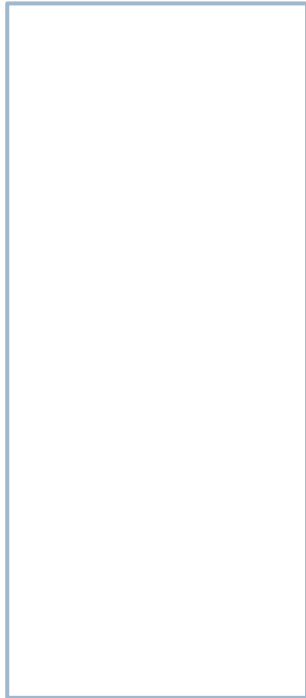- The page table is accessed by the MMU in the translation process

- The page table is modified by the operating system when a page fail occurs

# Example

- Pages of 1 KB
- Process of 8 KB
    - Number of pages: 8
- Size of sections:
    - Instructions: 1.5 KB
    - Data: 1 KB
    - Stack: 0.2 KB

# Example

| | |
|---|---|
| Instr. | Page 0 |
| Instr. | Page 1 |
| | Page 2 |
| Data | Page 3 |
| | Page 4 |
| | Page 5 |
| | Page 6 |
| | Page 7 |
| Stack | |

▸ **Pages of 1 KB**

▸ **Process of 8 KB**

  ▸ Number of pages: 8

▸ **Size of sections:**

  ▸ Instructions: 1.5 KB -> 2 pages

  ▸ Data: 1 KB ->   1 page

  ▸ Stack: 0.2 KB -> 1 page

# Example



| | |
|---|---|
| Instr. | Page 0 |
| Instr. | Page 1 |
| Data | Page 2 |
| | Page 3 |
| | Page 4 |
| | Page 5 |
| | Page 6 |
| | Page 7 |
| Stack | |

▸ Init virtual address (VA): 0

▸ Final virtual address: 8191

▸ Pags. 3, 4, 5 and  6 are not assigned to the program at the beginning

# Example
# Process image initially in disk

| Instr. | Page 0 |
| Instr. | Page 1 |
| Data | Page 2 |
| | Page 3 |
| | Page 4 |
| | Page 5 |
| | Page 6 |
| | Page 7 |
| Stack | |

| 0 | |
| 1 | |
| 2 | 0 |
| 3 | |
| 4 | 1 |
| 5 | 2 |
| 6 | |
| 7 | |
| 8 | 7 |
| 9 | |
| 10 | |
| 11 | |
| 12 | |

## Swap

Pages of the process

# Example
## OS creates the page table

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 0 | 0 | 2 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

All pages in swap at the begining

| | |
|---|---|
| Instr. | Page 0 |
| Instr. | Page 1 |
| Data | Page 2 |
| | Page 3 |
| | Page 4 |
| | Page 5 |
| | Page 6 |
| | Page 7 |
| Stack | |

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | 0 |
| 3 | |
| 4 | 1 |
| 5 | 2 |
| 6 | |
| 7 | |
| 8 | 7 |
| 9 | |
| 10 | |
| 11 | |
| 12 | |

Swap

Pages of the process

# Example
# Access to VA 0

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 0 | 0 | 2 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

Process → MMU

DV= 0

Swap

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | 0 |
| 3 | |
| 4 | 1 |
| 5 | 2 |
| 6 | |
| 7 | |
| 8 | 7 |
| 9 | |
| 10 | |
| 11 | |
| 12 | |

Pages of the process

| Instr. | Pag. 0 |
|---|---|
| Instr. | Pag. 1 |
| | Pag. 2 |
| Data | Pag. 3 |
| | Pag. 4 |
| | Pág. 5 |
| | Pag. 6 |
| | Pag. 7 |
| Stack | |

# Example
## Access to VA 0

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 0 | 0 | 2 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

processor → MMU

VA= 0

VA= 0 | 0 | 0 |
PN    D

| | |
|---|---|
| Instr. | Pag. 0 |
| Instr. | Pag. 1 |
| | Pag. 2 |
| Data | Pag. 3 |
| | Pag. 4 |
| | Pág. 5 |
| | Pag. 6 |
| | Pag. 7 |
| Stack | |

Swap

| 0 | |
|---|---|
| 1 | |
| 2 | 0 |
| 3 | |
| 4 | 1 |
| 5 | 2 |
| 6 | |
| 7 | |
| 8 | 7 |
| 9 | |
| 10 | |
| 11 | |
| 12 | |

Pages of the process

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Example
# Access to VA 0

|   | P | M | frame/swap |
|---|---|---|---|
| 0 | 0 | 0 | 2 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

Page fault
Page 0 is not in memory

processor

VA= 0

MMU

VA= 0 | 0 | 0 |
PN | D

| Instr. | Pag. 0 |
| Instr. | Pag. 1 |
| Data | Pag. 2 |
|  | Pag. 3 |
|  | Pag. 4 |
|  | Pág. 5 |
|  | Pag. 6 |
|  | Pag. 7 |
| Stack |  |

| 0 |  |
| 1 |  |
| 2 | 0 |
| 3 |  |
| 4 | 1 |
| 5 | 2 |
| 6 |  |
| 7 |  |
| 8 | 7 |
| 9 |  |
| 10 |  |
| 11 |  |
| 12 |  |

# Swap

Pages of the process

# Example
# handling the page fault

memory

| P | M | frame/swap |
|---|---|---|
| 0 | 0 | 0 | 2 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

processor → MMU

VA= 0    | 0 | 0 |
         PN    D

Swap

Pages of the process

OS reserves a free page frame in
memory (5) and copies the block 2
in the frame 5

# Example
# handling the page fault

Memory

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | 0 |
| 6 | |
| 7 | |
| 8 | |
| 9 | |
| 10 | |

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

processor → MMU

VA= 0

| 0 | 0 |
|---|---|
| PN | D |

OS updates the page table

Swap

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | 0 |
| 3 | |
| 4 | 1 |
| 5 | 2 |
| 6 | |
| 7 | |
| 8 | 7 |
| 9 | |
| 10 | |
| 11 | |
| 12 | |

Pages of the process

# Example
# Restoring the process

Memory

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

| | Memory |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | 0 |
| 6 | |
| 7 | |
| 8 | |
| 9 | |
| 10 | |

processor → MMU

VA= 0

| VA= 0 | 0 | 0 |
|---|---|---|
| | PN | D |

**VA 0 is generated again**

| | Swap |
|---|---|
| 0 | |
| 1 | |
| 2 | 0 |
| 3 | |
| 4 | 1 |
| 5 | 2 |
| 6 | |
| 7 | |
| 8 | 7 |
| 9 | |
| 10 | |
| 11 | |
| 12 | |

# Swap

Pages of the process

# Example
# Restoring the process

Memory

|   |   |
|---|---|
| 0 |   |
| 1 |   |
| 2 |   |
| 3 |   |
| 4 |   |
| 5 | 0 |
| 6 |   |
| 7 |   |
| 8 |   |
| 9 |   |
| 10 |   |

|   | P | M | frame/swap |
|---|---|---|---|
| 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

processor → MMU

VA= 0

VA= 0

| 0 | 0 |
|---|---|
| PN | D |

**VA 0 is generated again**

Swap

|   |   |
|---|---|
| 0 |   |
| 1 |   |
| 2 | 0 |
| 3 |   |
| 4 | 1 |
| 5 | 2 |
| 6 |   |
| 7 |   |
| 8 | 7 |
| 9 |   |
| 10 |   |
| 11 |   |
| 12 |   |

Pages of the process

# Example
# Restoring the process

Memory

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | 0 |
| 6 | |
| 7 | |
| 8 | |
| 9 | |
| 10 | |

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

processor → MMU

VA= 0

VA= 0

| 0 | 0 |
|---|---|
| PN | D |

PA

| 5 | 0 |
|---|---|
| PN | D |

Page in memory
Obtain the physical address

# Example
# Restoring the process

Memory

| | P | M | frame/swap |
|---|---|---|---|
| 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 5 |
| 3 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 |

| | Memory |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | 0 |
| 6 | |
| 7 | |
| 8 | |
| 9 | |
| 10 | |

processor → MMU

VA= 0

| VA= 0 | 0 | 0 |
|---|---|---|
| | PN | D |

| PA | 5 | 0 |
|---|---|---|
| | PN | D |

Access to memory

# Translation



**Virtual address**

**Physical address**

| Pag. Num. | Offset |
|---|---|

| Frame num. | Offset |
|---|---|

**Register**

Page table pointer

**Page table**

N.º of page

+

Frame #

**Offset**

**Frame**

**Program**

**Translation mechanism**

**Main memory**

# Memory protection

Virtual address space
of process 1

Virtual address 0

| Virtual page 0 |
| |
| |
| Virtual page k |
| |

Page table of
process 1

| | Frame A |
| | |
| | |
| | Frame C |
| | |

Main memory

| |
| Frame for page A |
| |
| Frame for page B |
| |
| Frame for page C |
| |
| Frame for page D |
| |

Different physical
addresses

Virtual address space
of process 2

Virtual address 0

| Virtual page 0 |
| |
| |
| Virtual page j |

Page table of
process 2

| | Frame B |
| | |
| | Frame D |

# Two-level page table

Virtual page

Virtual address

| 1st level | 2nd level | Byte |

Page table
Base register

First level
Page table

Sencond level
Page table

frame

# Inverted page table

Virtual address

| PID | Virtual page | Byte |
| --- | --- | --- |

Inverted page Table

| pid | Página virtual |

i

| i | Byte |

Physical address
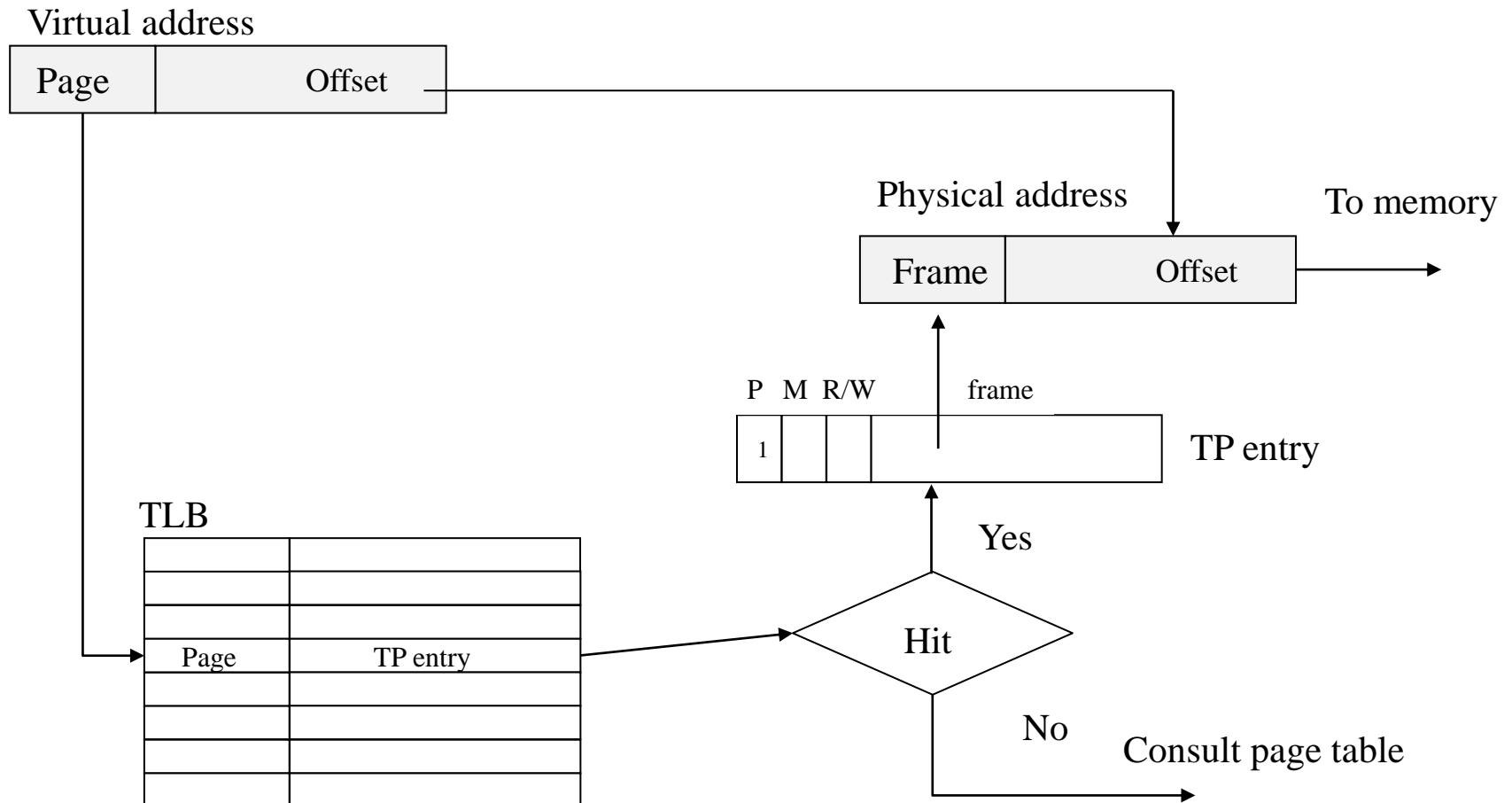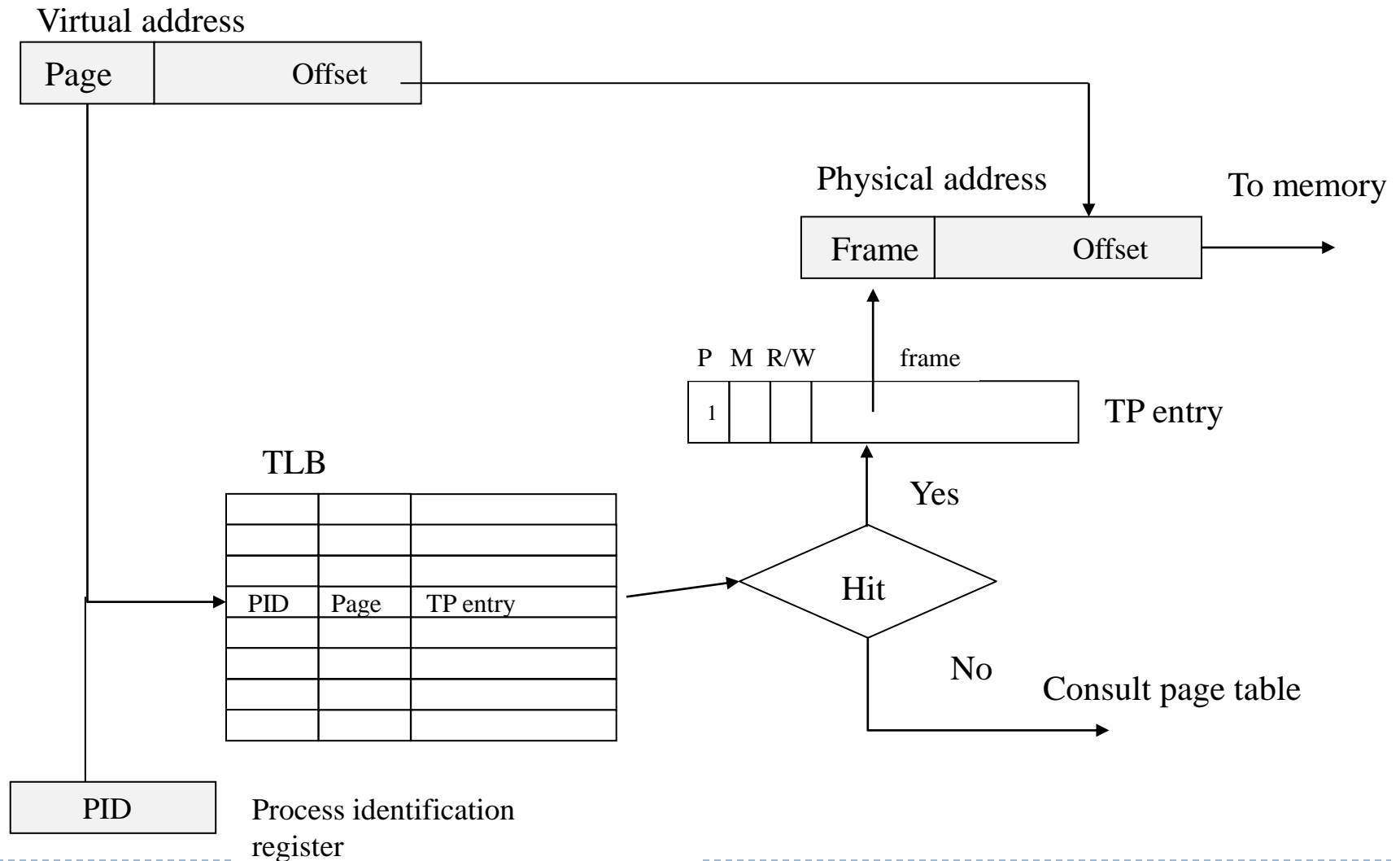
# TLB (Translation Lookaside Buffer)

▶ With virtual memory, two memory accesses are needed for each memory reference:

  ▶ One access to the page table

  ▶ One access to the page in memory

▶ TLB is used to optimize the memory access:

  ▶ Table with reduced access time located in the MMU

  ▶ Each entry has the page number and the corresponding page table entry

  ▶ In case of hit, the page table is not accessed

▶ Two types:

  ▶ TLB with process identification

  ▶ TLB without process identification

# TLB without process identification

Virtual address

| Page | Offset |
|------|--------|

Physical address

To memory

| Frame | Offset |
|-------|--------|

P  M  R/W          frame

| 1 |   |   |          |

TP entry

TLB

|      |          |
|------|----------|
|      |          |
|      |          |
| Page | TP entry |
|      |          |
|      |          |
|      |          |

Hit

Yes

No

Consult page table

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# TLB witht process identification

Virtual address

| Page | Offset |
|------|--------|

Physical address    To memory

| Frame | Offset |
|-------|--------|

P  M  R/W        frame

| 1 |  |  |  |     TP entry

TLB

| | | |
|---|---|---|
| | | |
| | | |
| PID | Page | TP entry |
| | | |
| | | |
| | | |

Yes

Hit

No    Consult page table

| PID |    Process identification register

ARCOS @ UC3M
Félix García Carballeira, Alejandro Calderón Mateos

# Virtual memory and cache memory

Processor

MMU

Virtual address

TLB

Physical address

L1 cache

word

L2 cache

line

line

Main memory

line

page

Disk (swap)

# Read access with cache and virtual memory

ARCOS Group

uc3m | Universidad **Carlos III** de Madrid

# Lesson 5 (III)
# Memory hierarchy

Computer Structure

Bachelor in Computer Science and Engineering