# Operating Systems

## Lecture 3
## Physical Memory Management

IIS & CS
Tsinghua University

# Review

- Dual Mode Operation
- What is an Interrupt/Exception/System Call?
- The difference of Interrupt/Exception/System Call
- X86 related
  - How to build IDT
  - The hardware processing when INT happens
  - The software processing when INT happens
  - The system call processing (non-privilege(user)  mode /privilege(supervisor) mode)
  - The different stacks in different privilege mode

# Review: Dual-mode operation

- Why do we have "user mode" and "kernel mode"?

- Problem: Would you trust any users to … read and write memory, manage resource, access I/O, …?

- Solution: dual mode operation

  ◆ CPU has a "mode" when it is executing an instruction

  ◆ "User Mode": can only perform a restricted set of operation (applications)
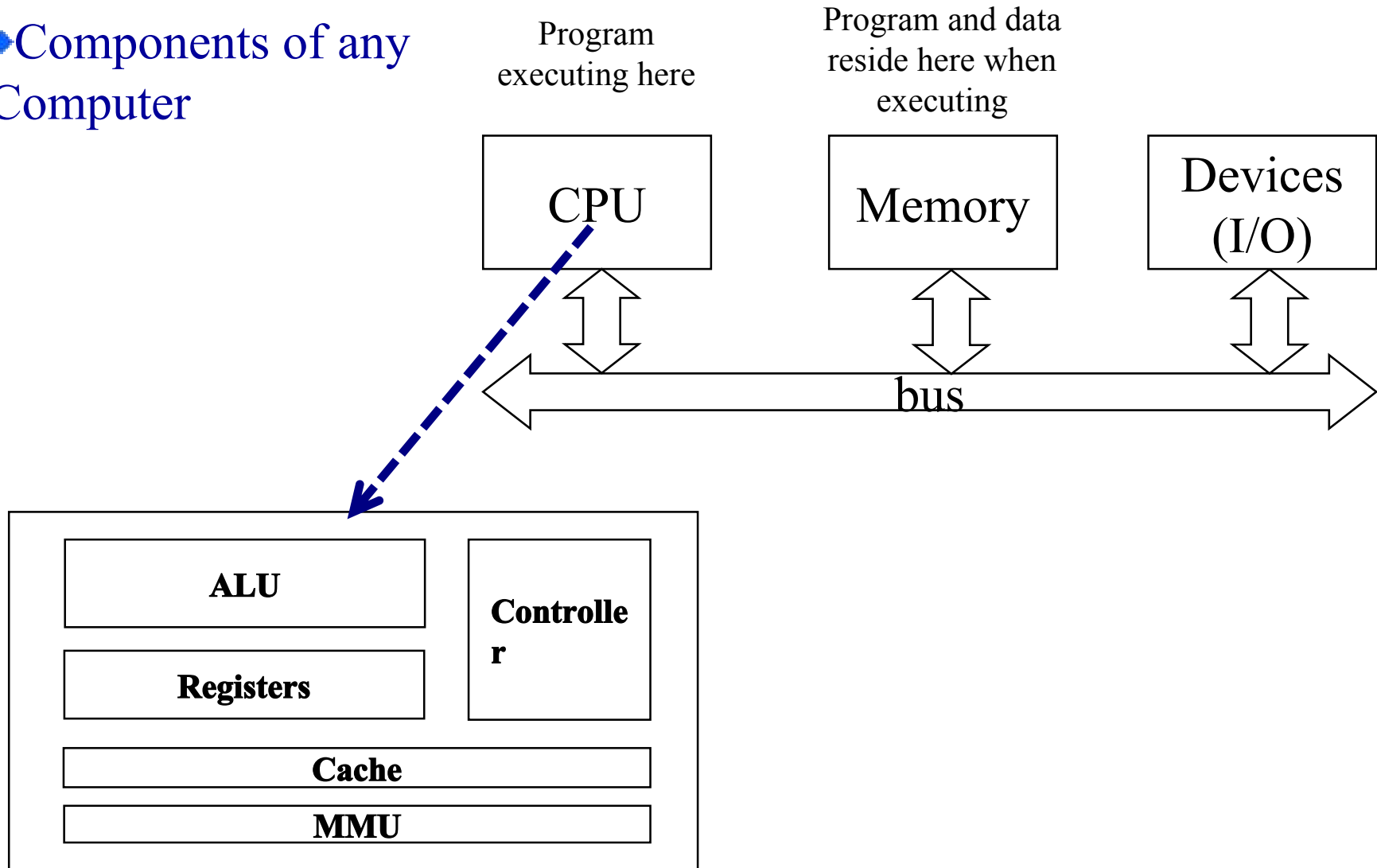
  ◆ "Kernel Mode": can do anything (OS kernel)

# From "User Mode" to "Kernel Mode"

- ## Interrupt: hardware device requests OS service

  - ◆ CPU interrupts current execution and jumps to interrupt handler, and returns when done

  - ◆ None of this is visible to user program

- ## Exceptions: user program acts illegally

  - ◆ CPU executes exception handlers

  - ◆ May cause abnormal execution flow (such as terminated)

- ## System calls: user program requests OS service

  - ◆ User program execute a trap instruction

  - ◆ OS identifies the type of service and parameters, and executes the requested service

  - ◆ OS returns to user program when done

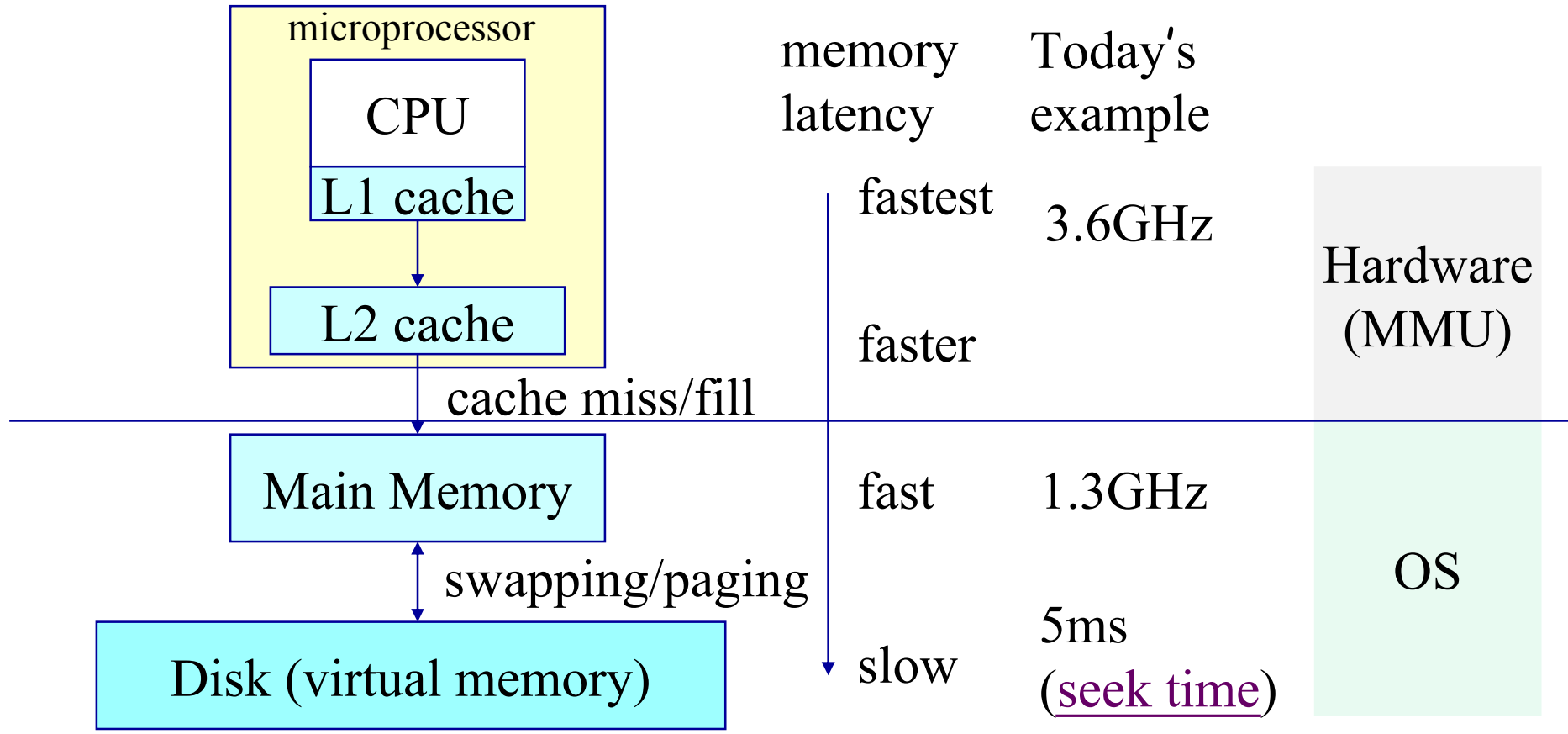  - ◆ This appears as a function call to the user program

- ● Computer Arch/Memory Hierarchy
- ● Address Space & Address Generation
- ● Contiguous Memory Allocation
  - ◆ Dynamic Allocation of Partitions
- ● Non-Contiguous Memory Allocation
  - ◆ Segmentation
  - ◆ Paging
  - ◆ Page Table
    - · Translation Look-aside Buffer (TLB)
    - · Multi-Level Page Table
    - · Inverted Page Table
  - ◆ Paged Segmentation Model
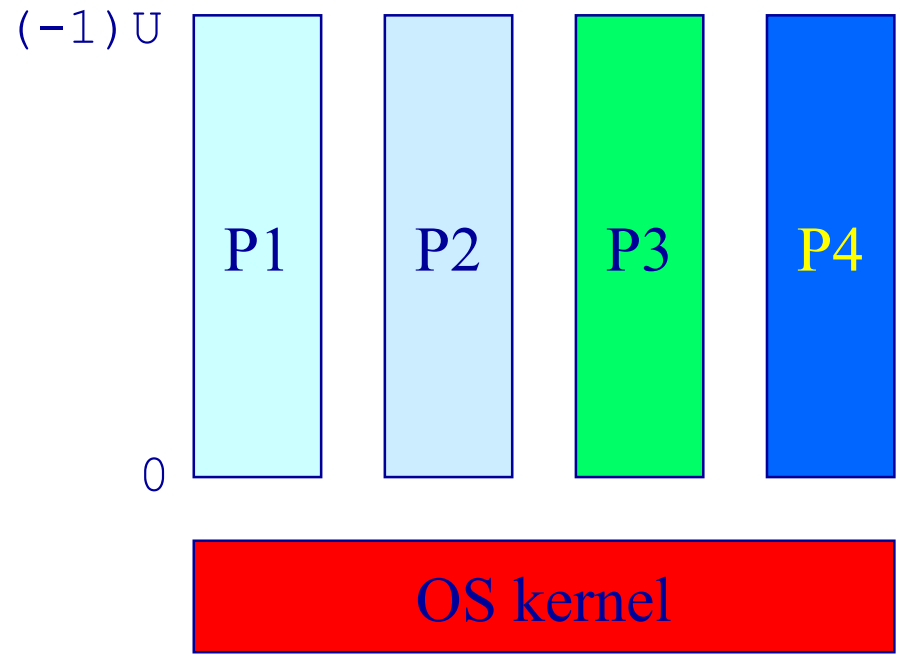
# Brief Introduction to Computer Architecture

● Components of any Computer

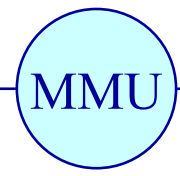Program executing here

Program and data reside here when executing

| CPU | Memory | Devices (I/O) |
|---|---|---|

bus

ALU

Registers

Controller

Cache

MMU

Intel® 64 and IA-32 Architectures Software Developer Manuals

# Memory Hierarchy

| microprocessor | | memory latency | Today's example | |
|---|---|---|---|---|
| **CPU** | | | | |
| L1 cache | fastest | | 3.6GHz | **Hardware (MMU)** |
| L2 cache | faster | | | |
| cache miss/fill | | | | |
| Main Memory | fast | | 1.3GHz | **OS** |
| swapping/paging | | | | |
| Disk (virtual memory) | slow | | 5ms (seek time) | |

# Modern Memory Management Paradigm

(-1)U

P1    P2    P3    P4

**OS kernel**

*Logical (virtual) space*

MMU

*Physical space*

- Abstraction
  - ➤ Logical address space
- Protection
  - ➤ Isolation
- Programming models
  - ➤ Shared memory

**Main memory**

**Disk (virtual memory)**

# OS Memory Models

- **Different ways to manage memory in an OS**
  - Program relocation
  - Segmentation
  - Paging
  - Virtual memory
  - Mostly (e.g., Linux): demand paging virtual memory

- **Implementation highly hardware dependent**
  - Must know memory architecture
  - MMU (Memory Management Unit): hardware components responsible for handling memory accesses requested by the CPU

- ● Computer Arch/Memory Hierarchy
- ● Address Space & Address Generation
- ● Contiguous Memory Allocation
  - ◆ Dynamic Allocation of Partitions
- ● Non-Contiguous Memory Allocation
  - ◆ Segmentation
  - ◆ Paging
  - ◆ Page Table
    - · Translation Look-aside Buffer (TLB)
    - · Multi-Level Page Table
    - · Inverted Page Table
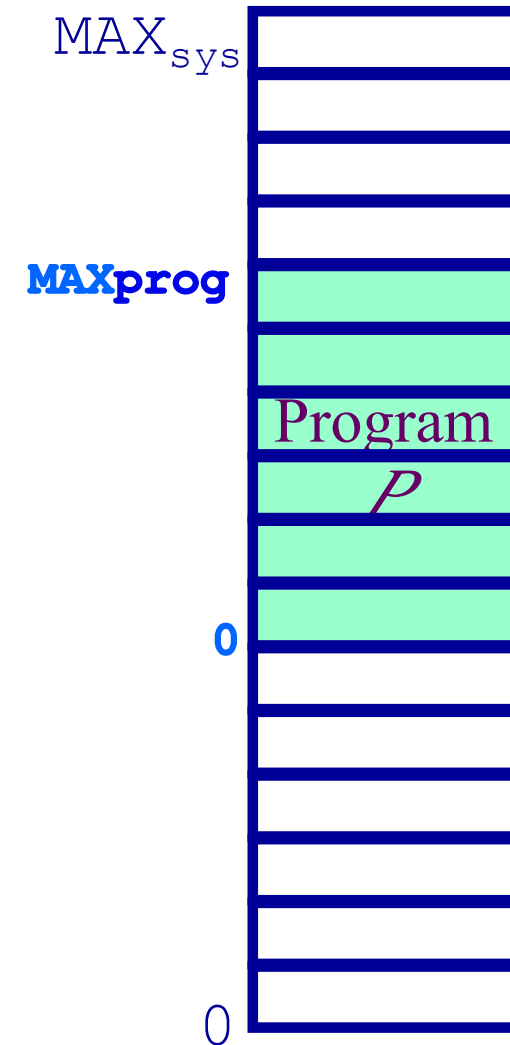  - ◆ Paged Segmentation Model

# Address Space & Address Generation

## address space

♦ *Physical address space* — The address space supported by the hardware

➤ Starting at address $0$, going to address **MAX$_{sys}$**

♦ *Logical address space* — A process's view of its own memory

➤ Starting at address $0$, going to address **MAX$_{prog}$**
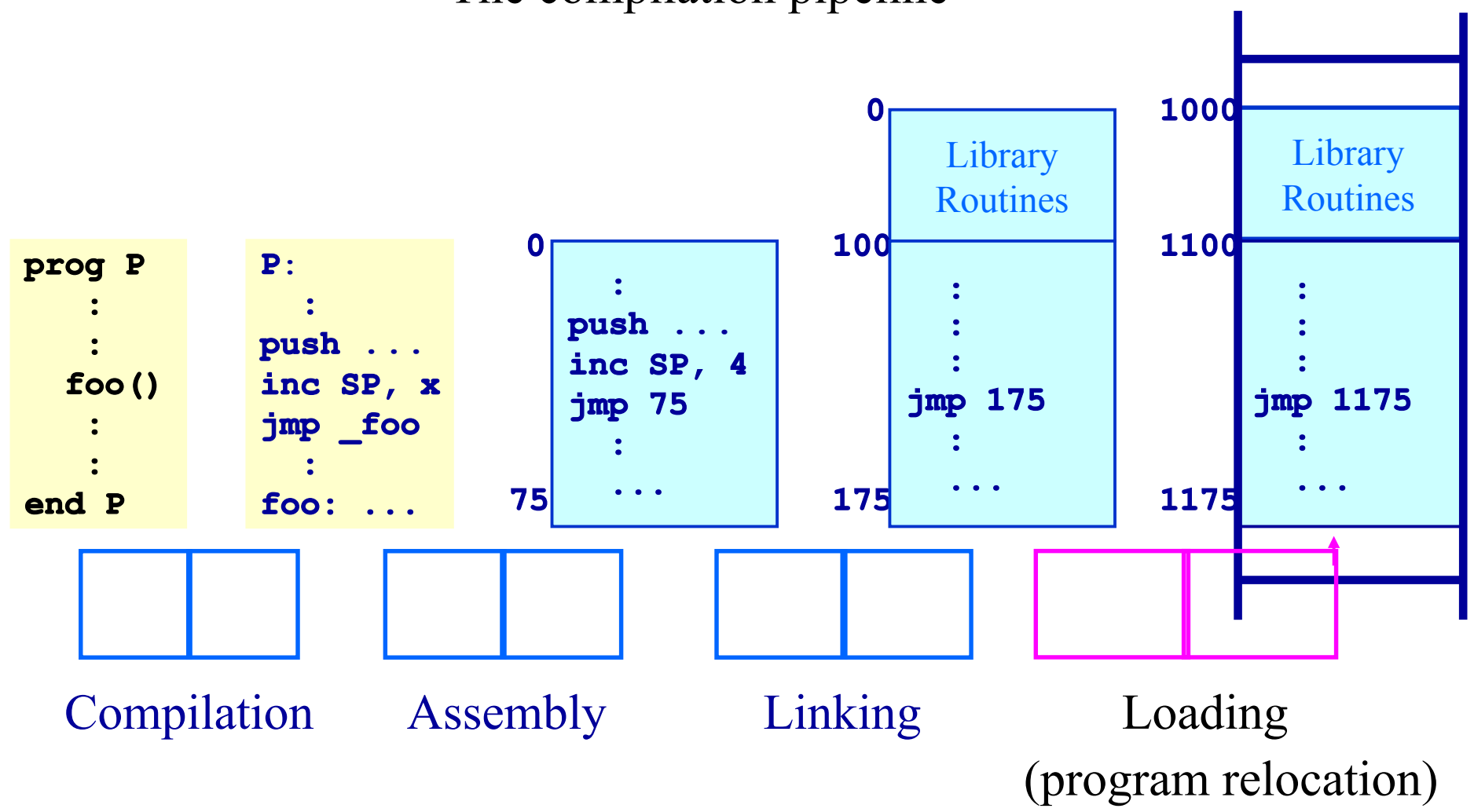
But where do addresses come from?

**movl %eax, $0xfffa620e**

MAX$_{sys}$

**MAXprog**

Program
$P$

0

0

# Address Space & Address Generation

## Address Generation

- The compilation pipeline

```
prog P
   :
   :
   foo()
   :
   :
end P
```

```
P:
   :
   push ...
   inc SP, x
   jmp _foo
   :
foo: ...
```

```
0
   :
   push ...
   inc SP, 4
   jmp 75
   :
75  ...
```

```
0
    Library
    Routines
100
   :
   :
   :
   jmp 175
   :
175  ...
```

```
1000
    Library
    Routines
1100
   :
   :
   :
   jmp 1175
   :
1175  ...
```

Compilation     Assembly     Linking     Loading
(program relocation)

# Address Space & Address Generation

## Address Generation Time

- Compile time

  ➢ If memory location known a priori

  ➢ Must recompile code if starting location changes

- Load time

  ➢ Compiler must generate *relocatable* code if memory location is not known at compile time

  ➢ Absolute addresses generated at load time

- Execution time

  ➢ The process can be moved during its execution

  ➢ Need hardware support for address translation

- Computer Arch/Memory Hierarchy

- Address Space & Address Generation

- Contiguous Memory Allocation

  ◆ Dynamic Allocation of Partitions

- Non-Contiguous Memory Allocation

  ◆ Segmentation

  ◆ Paging

  ◆ Page Table

    · Translation Look-aside Buffer (TLB)

    · Multi-Level Page Table
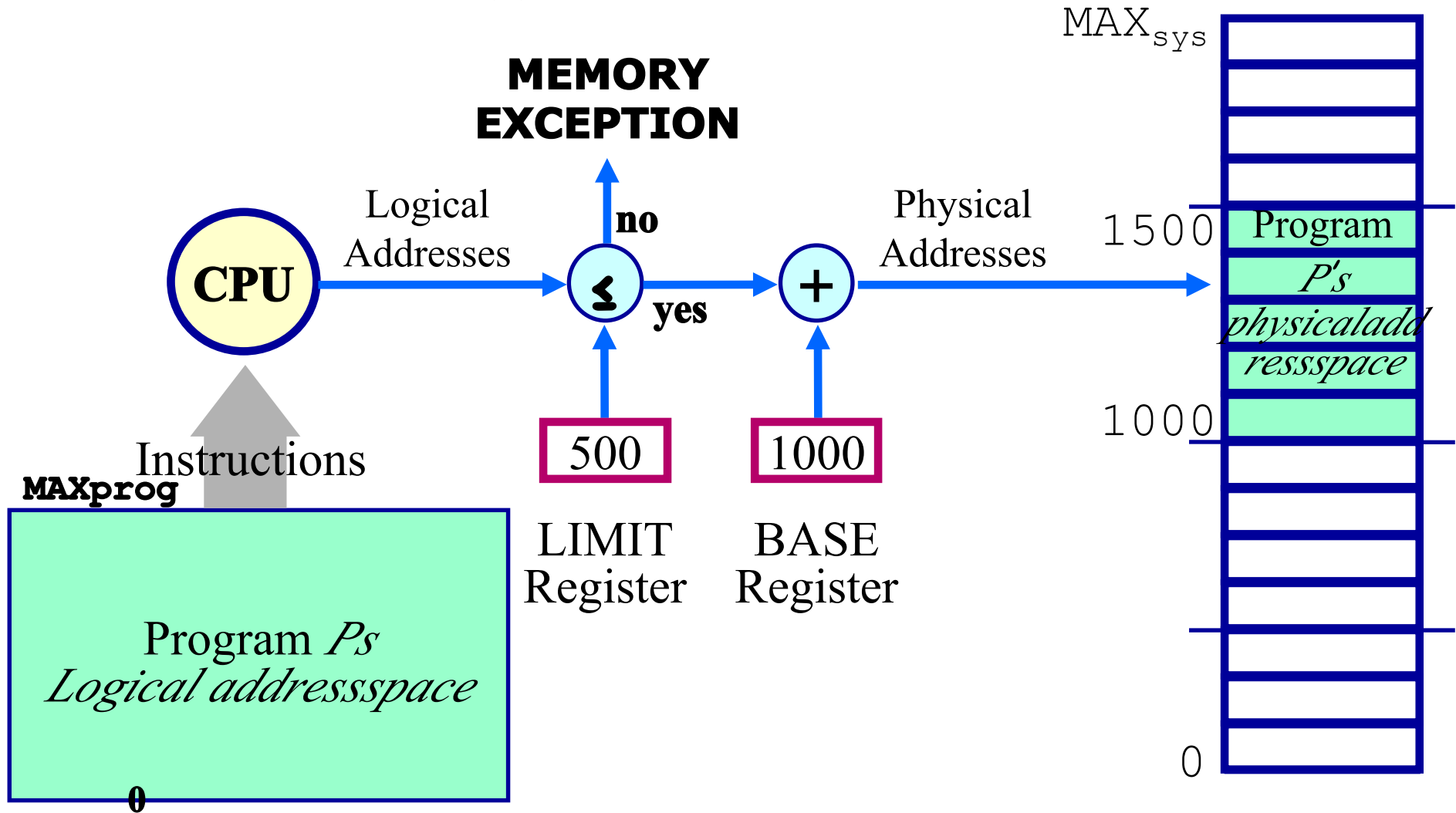
    · Inverted Page Table

  ◆ Paged Segmentation Model

## Program Relocation

- Relocate logical addresses to physical at run time
  - ➢ While we are relocating, also bounds check addresses for safety.

- Require hardware support (MMU)

- Basic component
  - ➢ Address translation with two registers: BASE and LIMIT
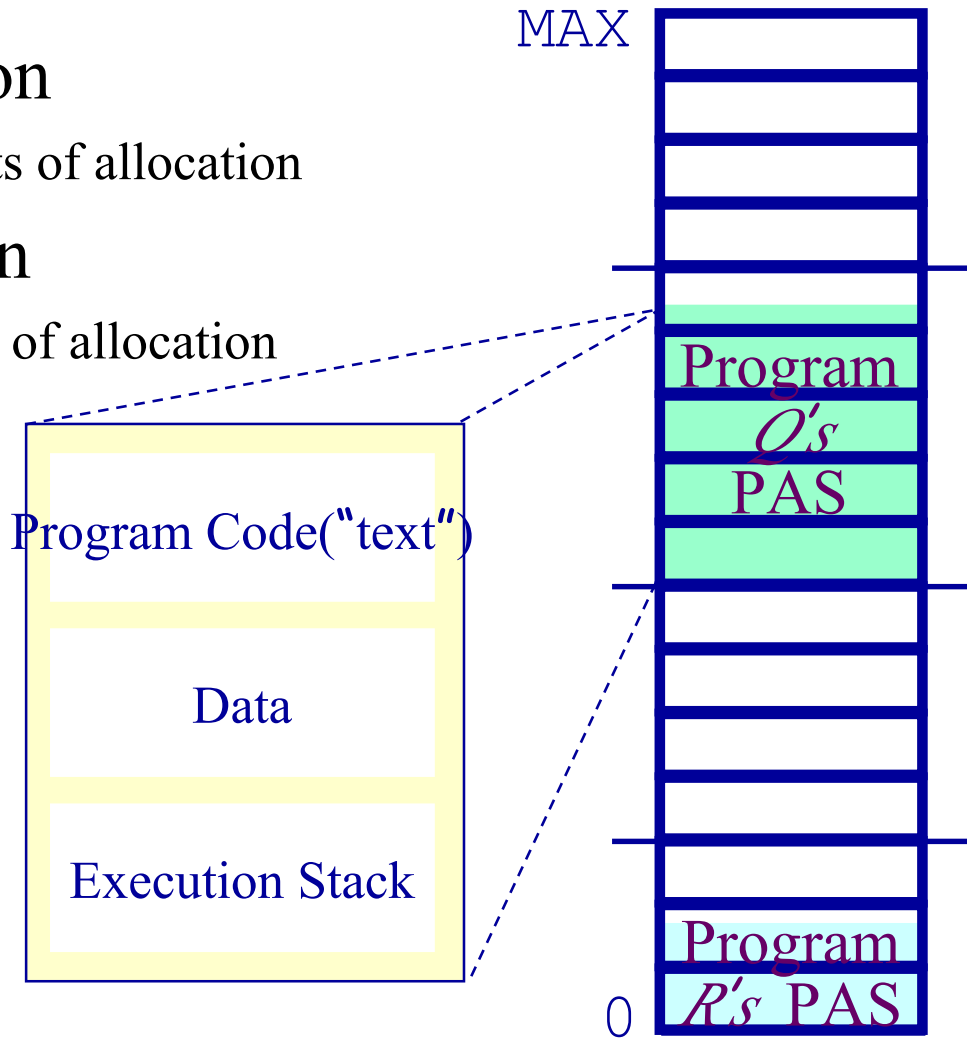
# Contiguous Memory Allocation

## Address Translation

$MAX_{sys}$

**MEMORY EXCEPTION**

CPU

Logical Addresses

no

≤

yes

+

Physical Addresses

1500 — Program

$P's$ *physicaladd ressspace*

1000

500

1000

LIMIT Register

BASE Register

Instructions

**MAXprog**

Program *Ps Logical addressspace*

0

0

# Contiguous Memory Allocation

## The Fragmentation Problem

- Free memory cannot be utilized

- External fragmentation

  ➢ Unused memory between units of allocation

- Internal fragmentation

  ➢ Unused memory within a unit of allocation

MAX

Program $Q's$ PAS

Program Code("text")

Data

Execution Stack

Program $R's$ PAS

0

# Contiguous Memory Allocation

## Dynamic Allocation of Partitions

◆ Simple memory management approach:

➢ Allocate a partition when a process is admitted into the system

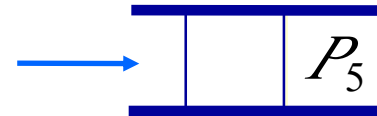➢ Allocate a contiguous memory partition to the process

OS keeps track of...
Full-blocks
Empty-blocks ("holes")

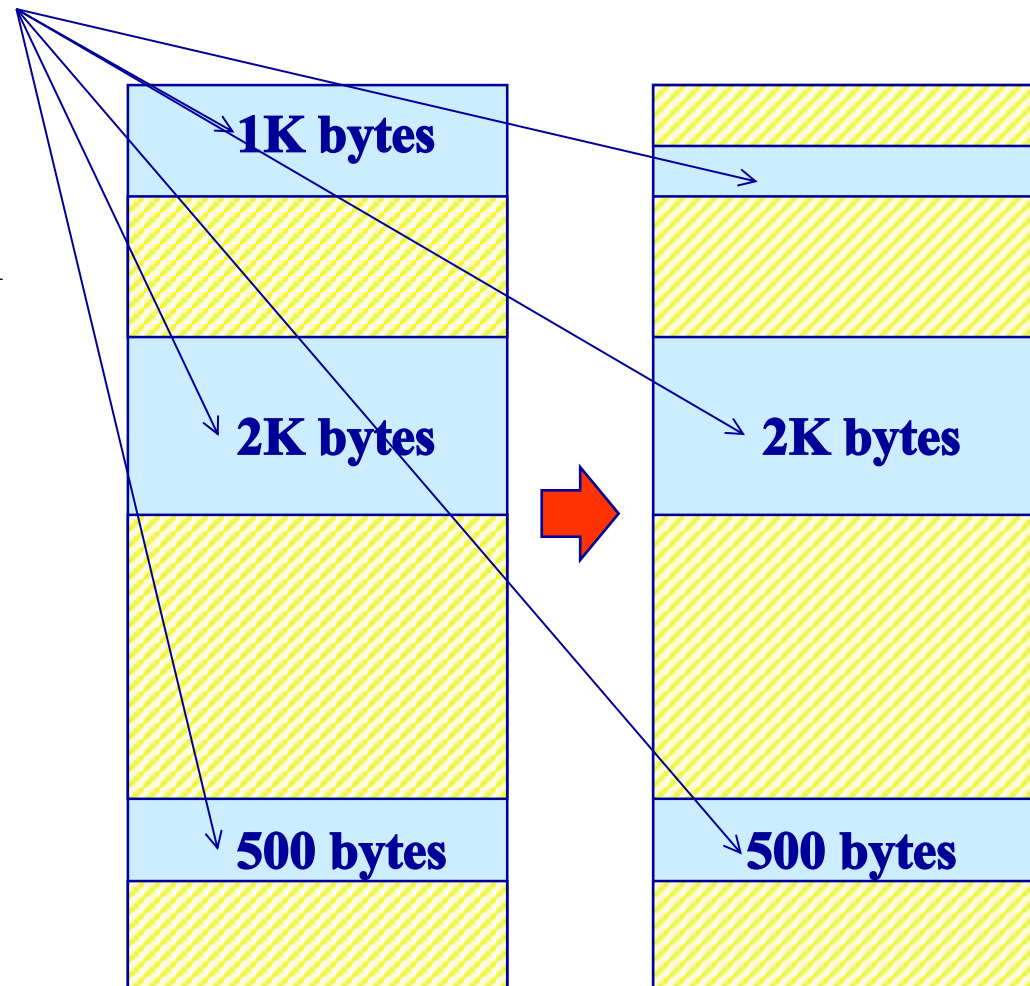Allocation strategies
First-fit
Best-fit
Worst-fit

MAX

| Program |
| $P1$ |
| |
| |
| Program |
| $P2$ |
| |
| |
| |
| Program |
| $P3$ |
| |
| |
| Program |
| $P4$ |
| |

$P_5$

0

# Contiguous Memory Allocation

## First Fit Allocation

FreeBlock

To allocate *n* bytes, use the *first* available free block such that the block size is larger than *n*.

To allocate 400 bytes, we use the 1st free block available



1K bytes

2K bytes

500 bytes

2K bytes

500 bytes

# Contiguous Memory Allocation

## Rationale & Implementation

- Simplicity of implementation
- Requires:
  - Free block list sorted by address
  - Allocation requires a search for a suitable partition
  - De-allocation requires a check to see if the freed partition could be merged with adjacent free partitions (if any)

### Advantages

- Simple
- Tends to produce larger free blocks toward the end of the address space
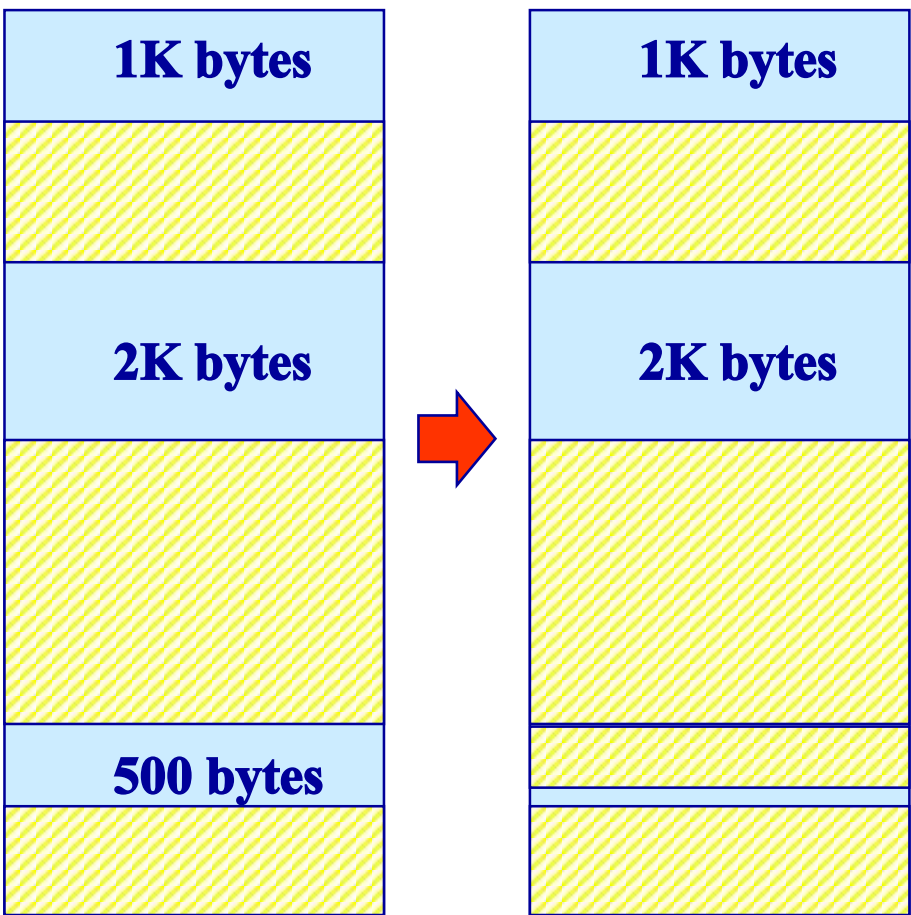
### Disadvantages

- External fragmentation
- Uncertainty

# Contiguous Memory Allocation

## Best Fit Allocation

To allocate $n$ bytes, use the *smallest* available free block such that the block size is larger than $n$.

To allocate 400 bytes, we use the 3rd free block available (smallest)



| 1K bytes |
| 2K bytes |
| 500 bytes |

| 1K bytes |
| 2K bytes |

# Contiguous Memory Allocation

## Rationale & Implementation

- To avoid fragmenting big free blocks
- To minimize the size of external fragments produced
- Requires:
  - Free block list sorted by size
  - Allocation requires search for a suitable partition
  - De-allocation requires search + merge with adjacent free partitions, if any

| Advantages | Disadvantages |
|---|---|
| Works well when most allocations are of small size | External fragmentation |
| Relatively simple | Slow de-allocation |
| | Tends to produce many useless tiny fragments (not really great) |

# Contiguous Memory Allocation

## Worst Fit Allocation

To allocate *n* bytes, use the *largest* available free block such that the block size is larger than *n*.

To allocate 400 bytes, we use the 2nd free block available (largest)

| 1K bytes |
|:---:|
| |
| 2K bytes |
| |
| 500 bytes |
| |

→

| 1K bytes |
|:---:|
| |
| |
| |
| |
| |

# Contiguous Memory Allocation

## Rationale & Implementation

◆ To avoid having too many tiny fragments

◆ Requires:

➢ Free block list sorted by size

➢ Allocation is fast (get the largest partition)

➢ De-allocation requires merge with adjacent free partitions, if any, and then adjusting the free block list

### Advantages

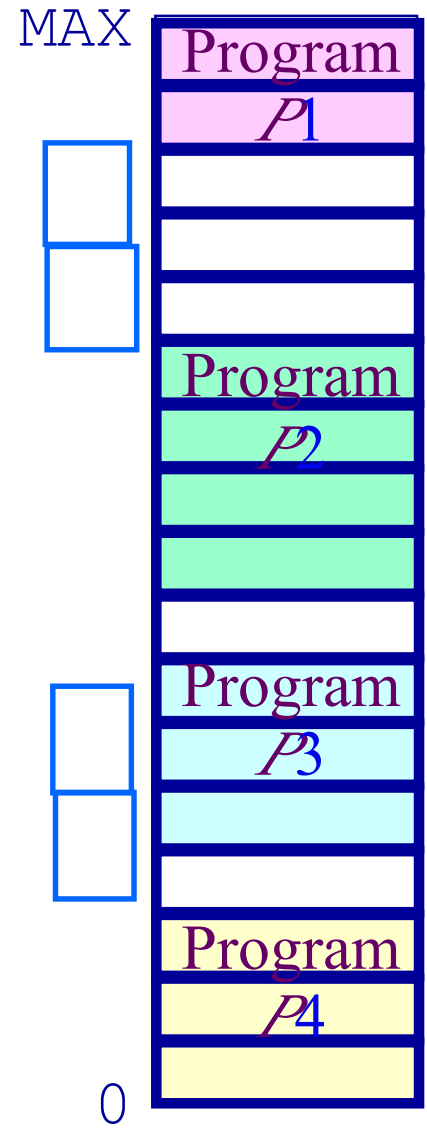❧ Works best if allocations are of medium sizes

### Disadvantages

❧ Slow de-allocation

❧ External fragmentation

❧ Tends to break large free blocks such that large partitions cannot be allocated

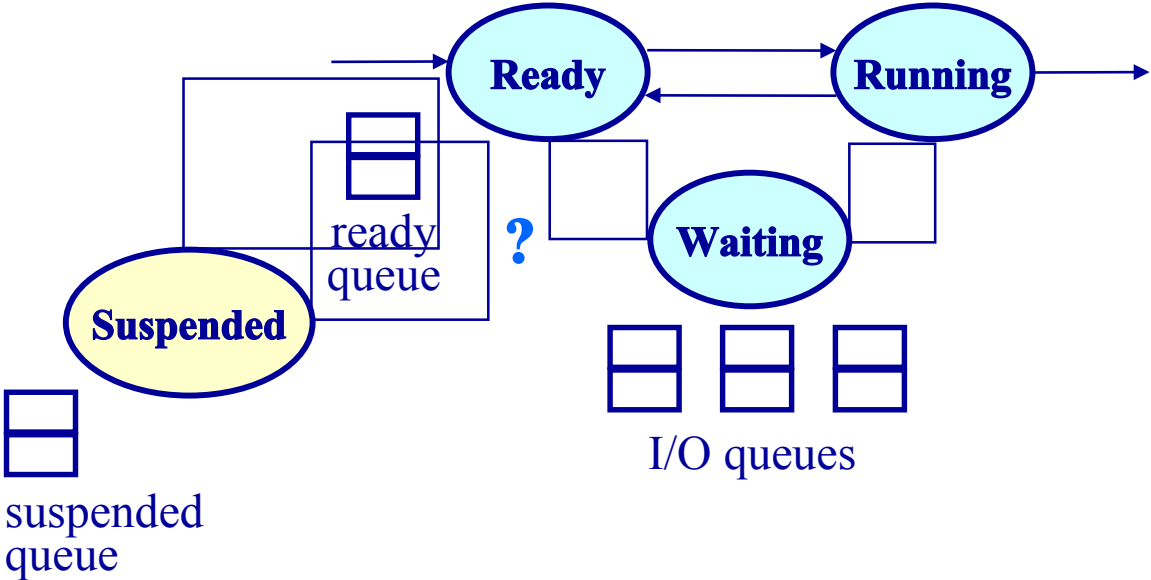# Contiguous Memory Allocation

## De-fragmentation by Compaction

- Relocate programs to coalesce holes
- Require all programs to be dynamically relocatable
- Issues
  - When to relocate?
  - Overhead

MAX

| Program |
| P1 |

| Program |
| P2 |

| Program |
| P3 |

| Program |
| P4 |

0

# Contiguous Memory Allocation

## De-fragmentation by Swapping
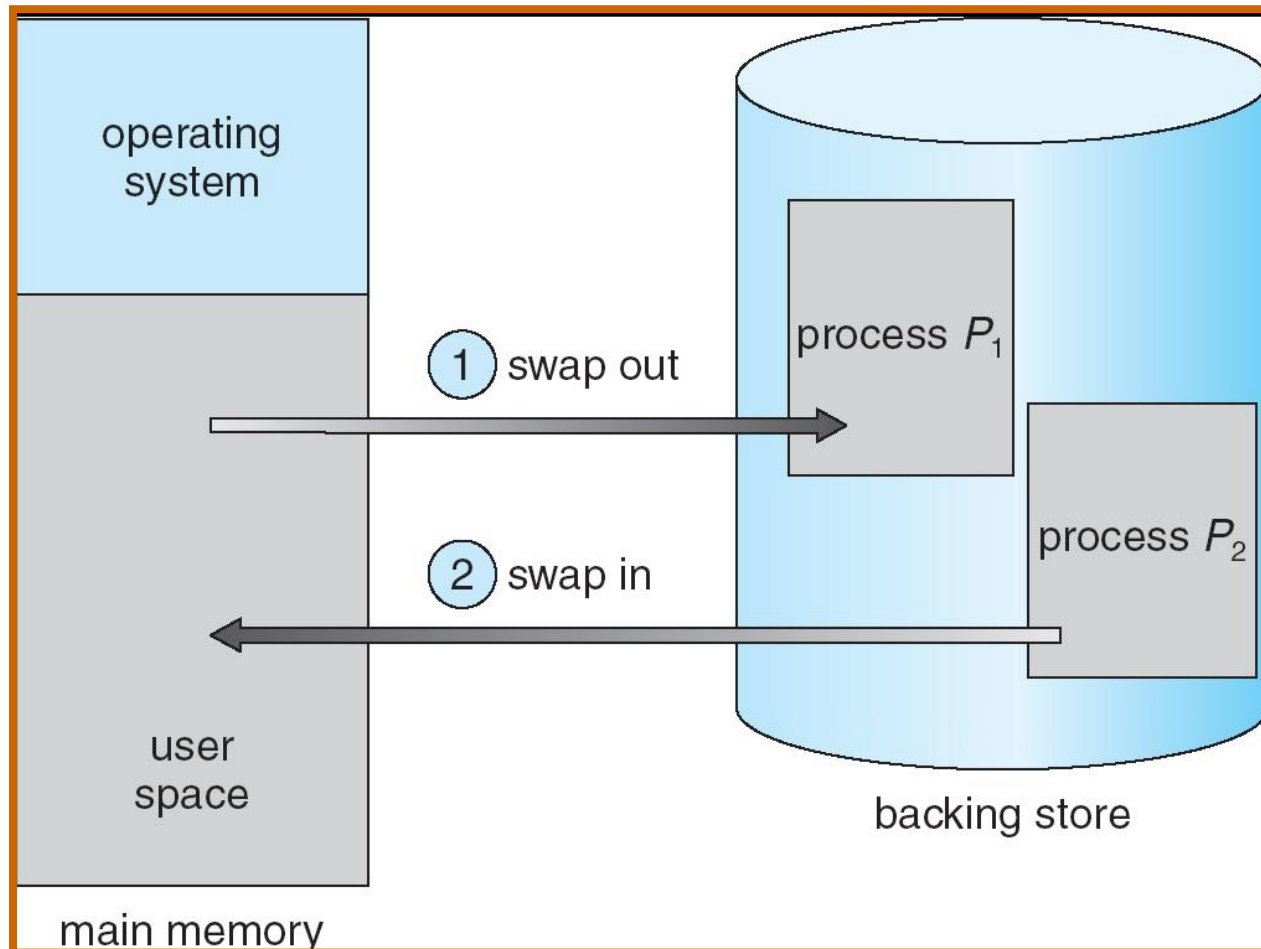
◆ Preempt processes & reclaim their memory



◆ Issue: which process(es) to swap?

# Schematic View of Swapping



operating
system

(1) swap out

process $P_1$

(2) swap in

process $P_2$

user
space

backing store

main memory

- Computer Arch/Memory Hierarchy
- Address Space & Address Generation
- Contiguous Memory Allocation
  - ◆ Dynamic Allocation of Partitions
- Non-Contiguous Memory Allocation
  - ◆ Segmentation
  - ◆ Paging
  - ◆ Page Table
    - Translation Look-aside Buffer (TLB)
    - Multi-Level Page Table
    - Inverted Page Table
  - ◆ Paged Segmentation Model

# Non-contiguous Allocation : Segmentation

- ## Previously,
  - ➢ Physical memory allocated to a process is contiguous
  - ➢ Poor memory utilization
  - ➢ Suffers from external fragmentation

- ## Noncontiguous allocation
  - ➢ Physical address space of a process is noncontiguous
  - ➢ Better memory utilization and management
  - ➢ Allow sharing of common blocks (code, data, library, etc.)
  - ➢ Support dynamic loading and dynamic linking

- ## Two schemes: segmentation and paging

# Non-contiguous Allocation : Segmentation

## Dynamic Loading

- Routine is not loaded until it is called

- Better memory-space utilization; unused routine is never loaded

- Useful when large amounts of code are needed to handle infrequently occurring cases

- Most OS allows user programs to do dynamic loading of components (relocatable object code)
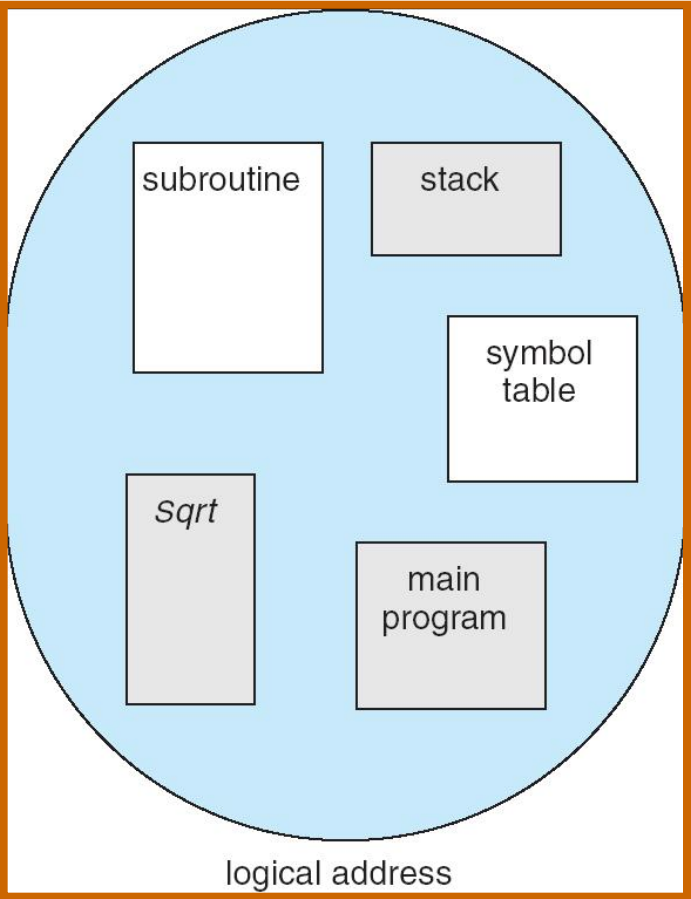
- Some OS supports loadable kernel modules

# Non-contiguous Allocation : Segmentation

## Dynamic Linking

- Linking postponed until execution time
  - ➢ Small piece of code, *stub*, used to locate the appropriate memory-resident library routine
  - ➢ Stub replaces itself with the address of the routine, and executes the routine
  - ➢ Operating system needed to check if routine is in processes' memory address
- Dynamic linking is particularly useful for libraries
  - ➢ Better known as *shared libraries*
- Dynamic linking in ucore

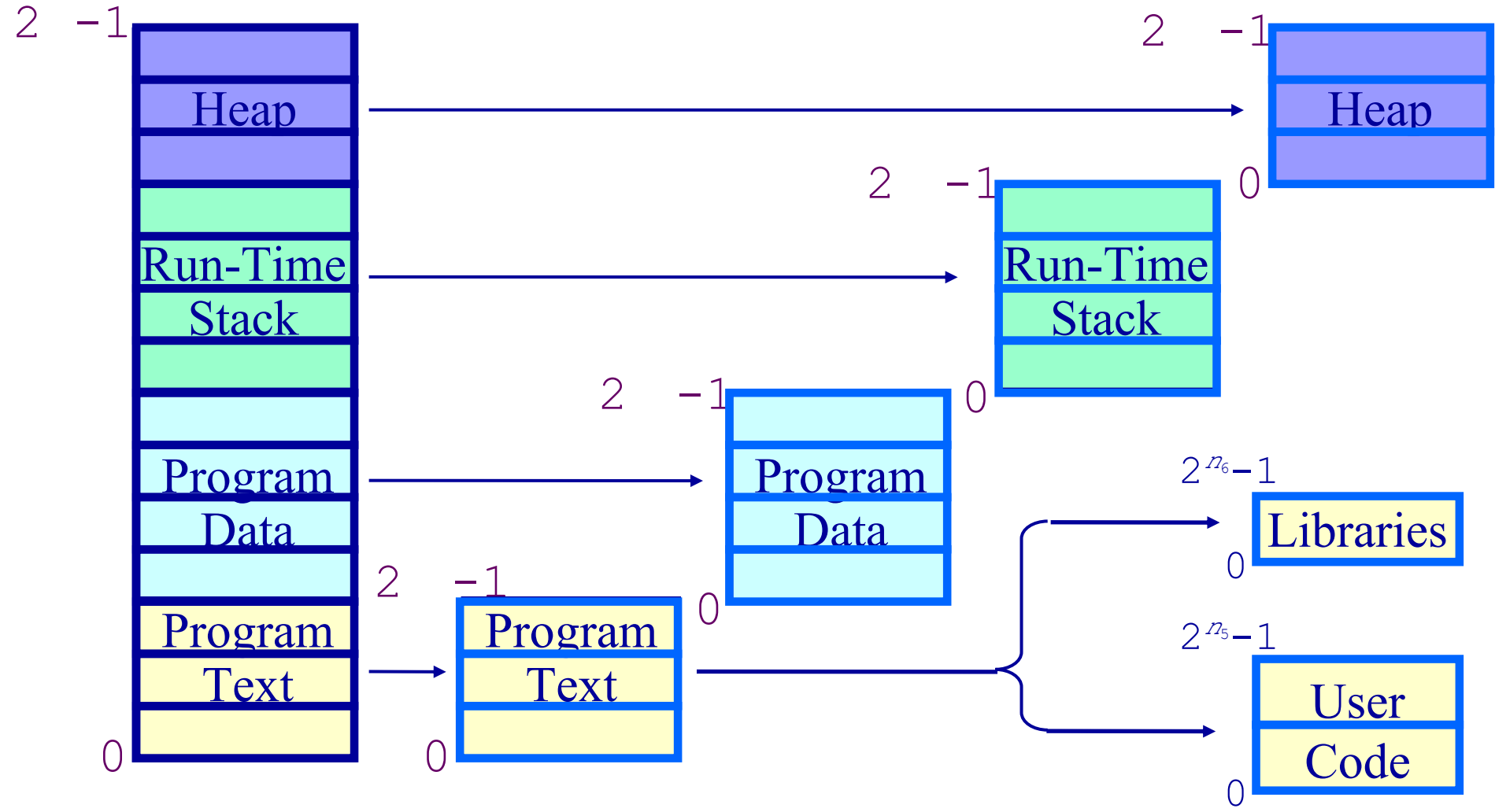# Non-contiguous Allocation : Segmentation

## Segmentation



- A program is a collection of segments, such as
  - Main program
  - Subroutines
  - Stack
  - Symbols
  - Data
  - Common libraries
  - Common blocks

- Purpose: enable finer grain isolation and sharing
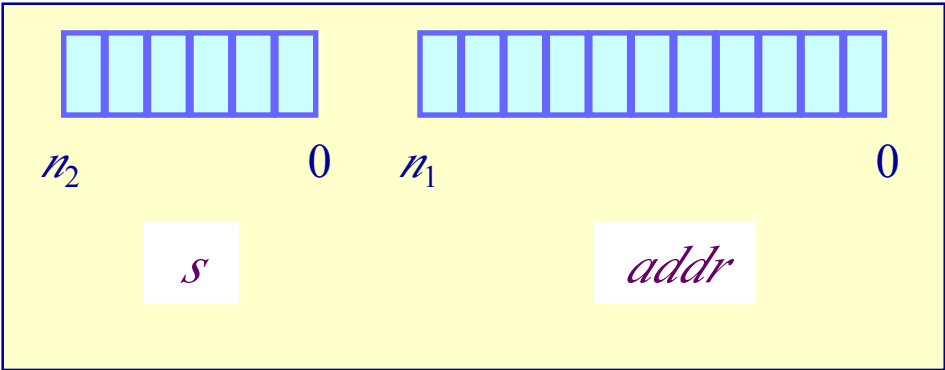
# Non-contiguous Allocation : Segmentation

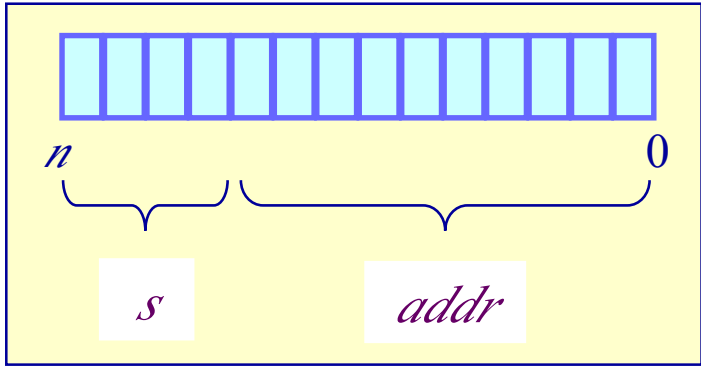## Separating into Multiple Address Spaces

$2-1$

| Heap |
| --- |
| |
| Run-Time Stack |
| |
| |
| Program Data |
| |
| Program Text |
| |

$0$

$2-1$ Heap → $2-1$ Heap $0$

$2-1$ Run-Time Stack → Run-Time Stack $0$

$2-1$ Program Data → Program Data $0$

$2-1$ Program Text → Program Text $0$

$2^{n_6}-1$ Libraries $0$

$2^{n_5}-1$ User Code $0$

# Segmentation Schemes

◆New concept: A segment — a memory "object"

➢ A logical address space

◆A process now addresses objects —a pair ($s$, $addr$)

➢ $s$ — segment number

➢ $addr$ — an offset within an object



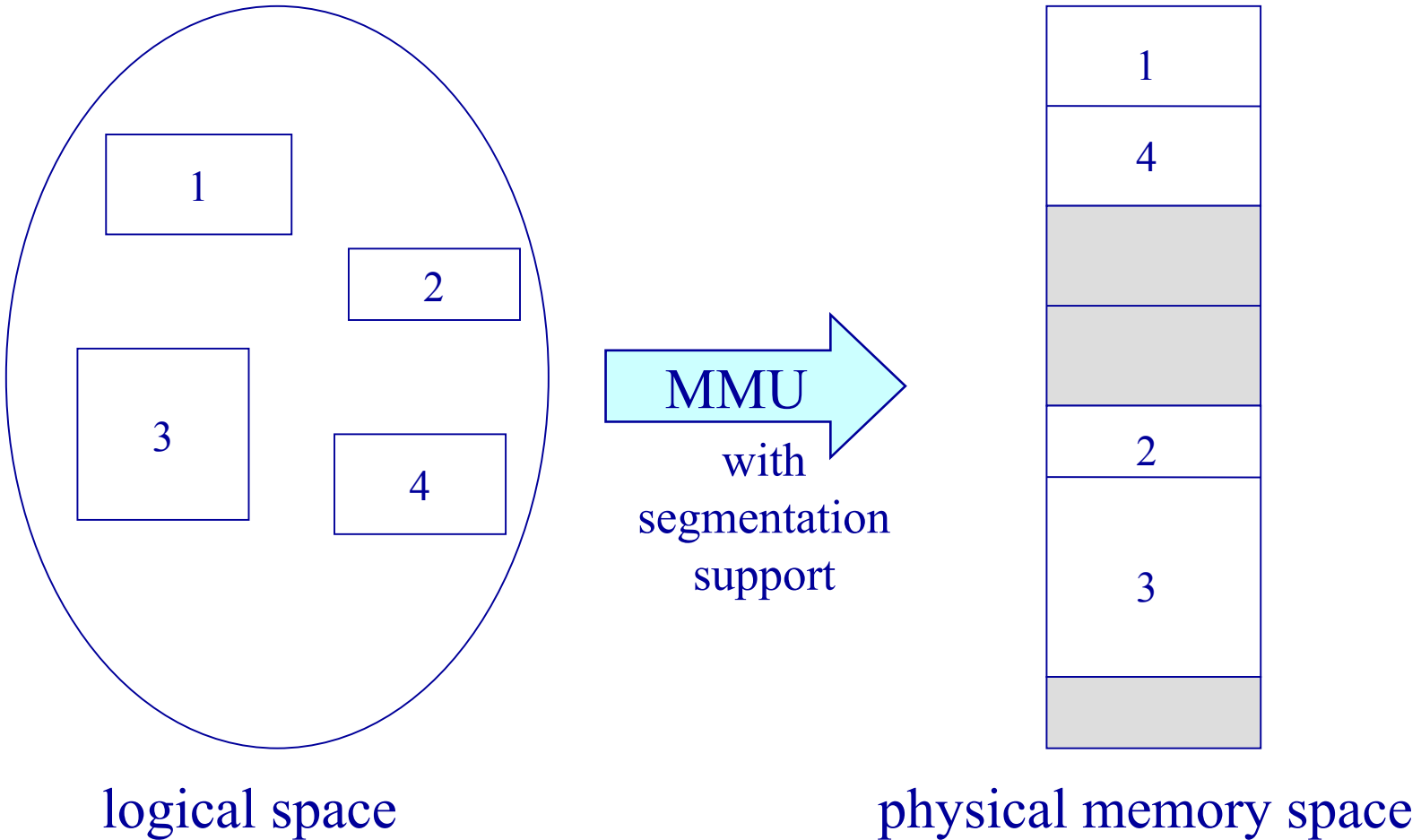| $n_2$ | $0$ | $n_1$ | $0$ |
|---|---|---|---|
| $s$ | | $addr$ | |

Segment + Address register scheme

$n$ ... $0$

$s$     $addr$

Single address scheme

# Non-contiguous Allocation: Segmentation

## Logical View of Segmentation



logical space

physical memory space

# Segmentation Hardware Architecture

- Segment table (in cache)
- STBR: Segment Table Base Register

- Computer Arch/Memory Hierarchy

- Address Space & Address Generation

- Contiguous Memory Allocation

    ◆ Dynamic Allocation of Partitions

🔴 ● Non-Contiguous Memory Allocation

    ◆ Segmentation

    ◆ Paging

    ◆ Page Table

        ‧ Translation Look-aside Buffer (TLB)
        ‧ Multi-Level Page Table
        ‧ Inverted Page Table

    ◆ Paged Segmentation Model

# Non-contiguous Allocation : Paging

- Divide physical memory into fixed-sized frames
  - ➢ Size is power of 2, e.g., 512, 4096, 8192
- Divide logical address space into same size pages
- To run a program of size n pages, find n free frames and load program
- Set up a page table to translate logical to physical addresses (pages to frames)
- Frame/page: basic units of memory allocation
  - ◆ OS keep track of all free frames
  - ◆ Same-sized frame eliminates external fragmentation
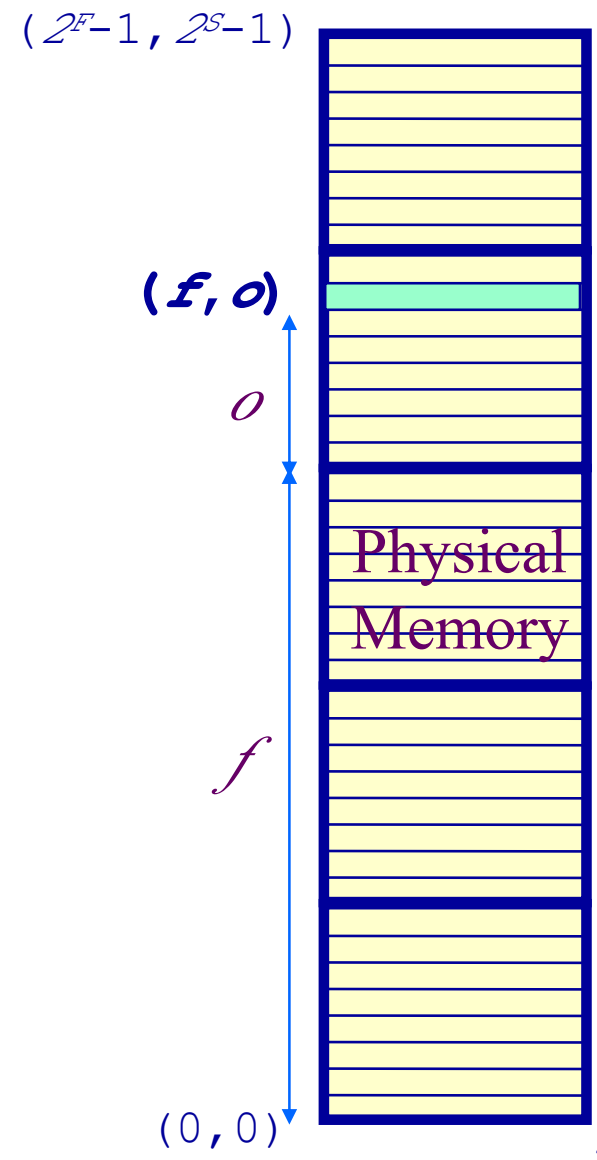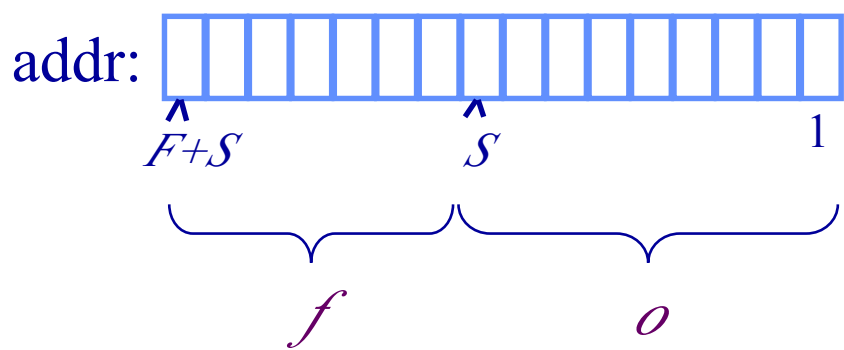
# Non-contiguous Allocation : Paging

## Frames

◆Physical memory partitioned into equal sized *frames*

A memory address is a pair ($f$, $o$)

$f$ — frame number (total $2^F$ frames)
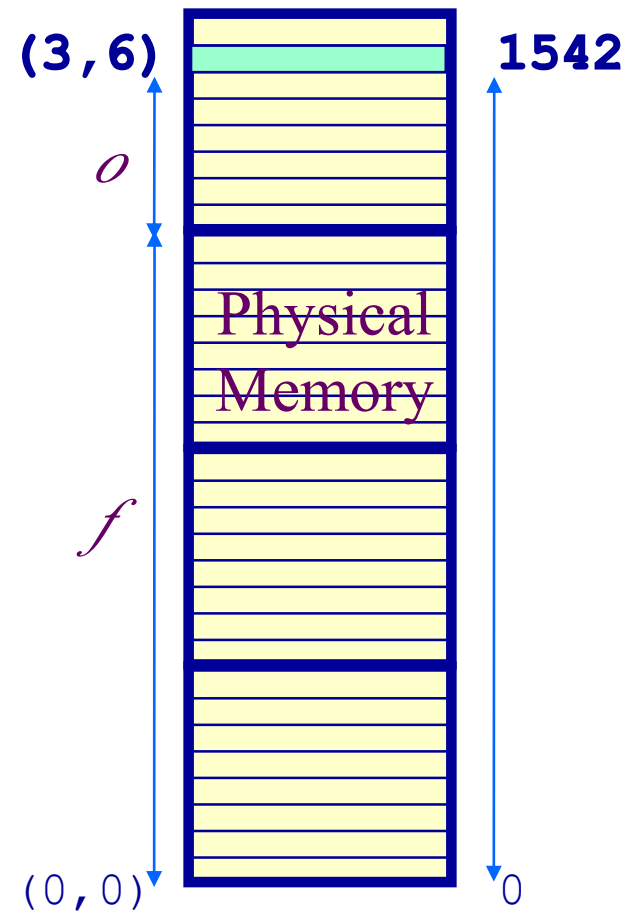$o$ — frame offset ($2^S$ bytes/frames)
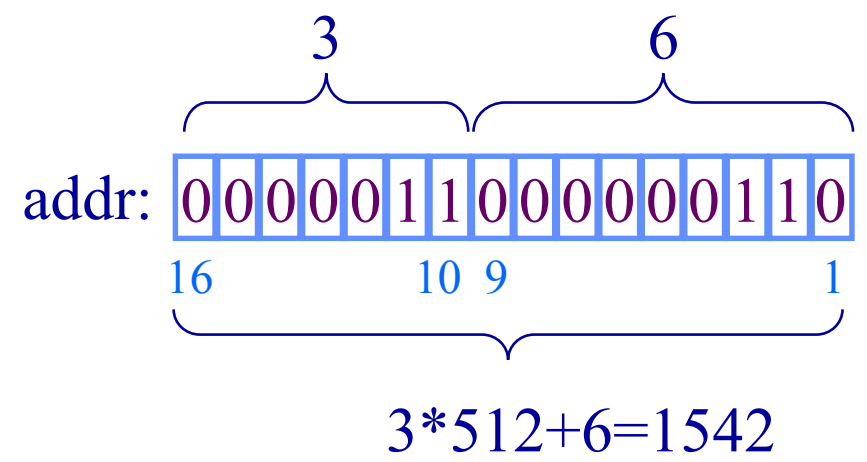Physical address = $2^S \cdot f + o$

addr:

$\overset{\wedge}{F+S} \qquad \overset{\wedge}{S} \qquad 1$

$\underbrace{\qquad}_{f} \qquad \underbrace{\qquad}_{o}$

$(2^F-1, 2^S-1)$

($f$, $o$)

$o$

Physical
Memory

$f$

$(0, 0)$

# Non-contiguous Allocation : Paging

## Frame Example

◆ Example: A 16-bit address space with 9-bit (512 byte) page frames

> ➢ Addressing location (3, 6) = 1542

**(3,6)**                                        **1542**

$o$

```
        3              6
addr: 0 0 0 0 0 1 1 0 0 0 0 0 0 1 1 0
     16        10 9              1
```

Physical Memory

$f$

3*512+6=1542

**(0,0)**                                        **0**

# Non-contiguous Allocation : Paging

## Pages

$2^n-1 =$
$(2^P-1, \ 2^S-1)$

◆ A process's logical address space is partitioned into equal sized *pages*
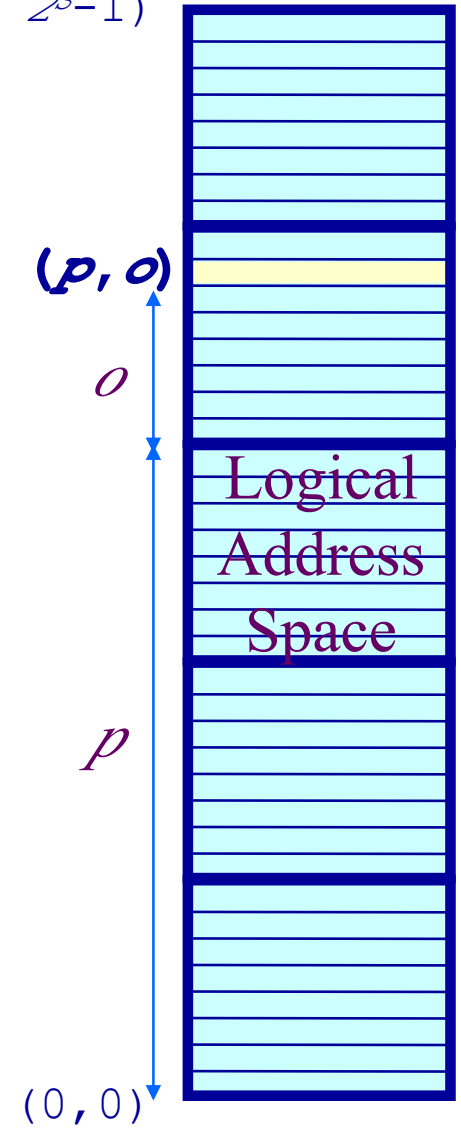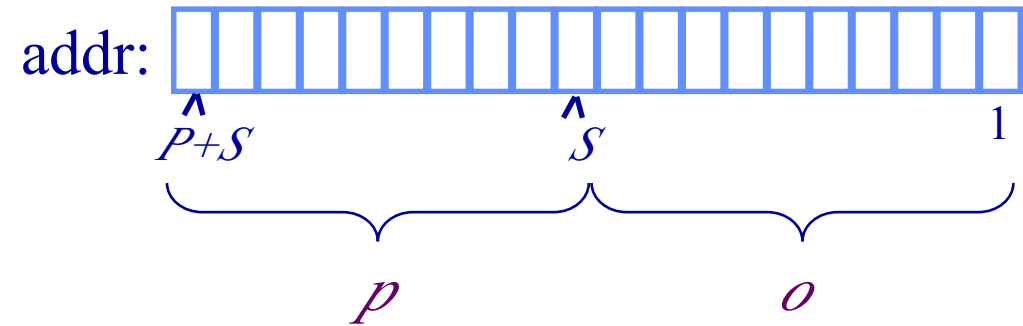
➤ |*page*| = |*frame*|

A logical address is a pair ($p$, $o$)

$p$     — page number ($2^P$ pages)

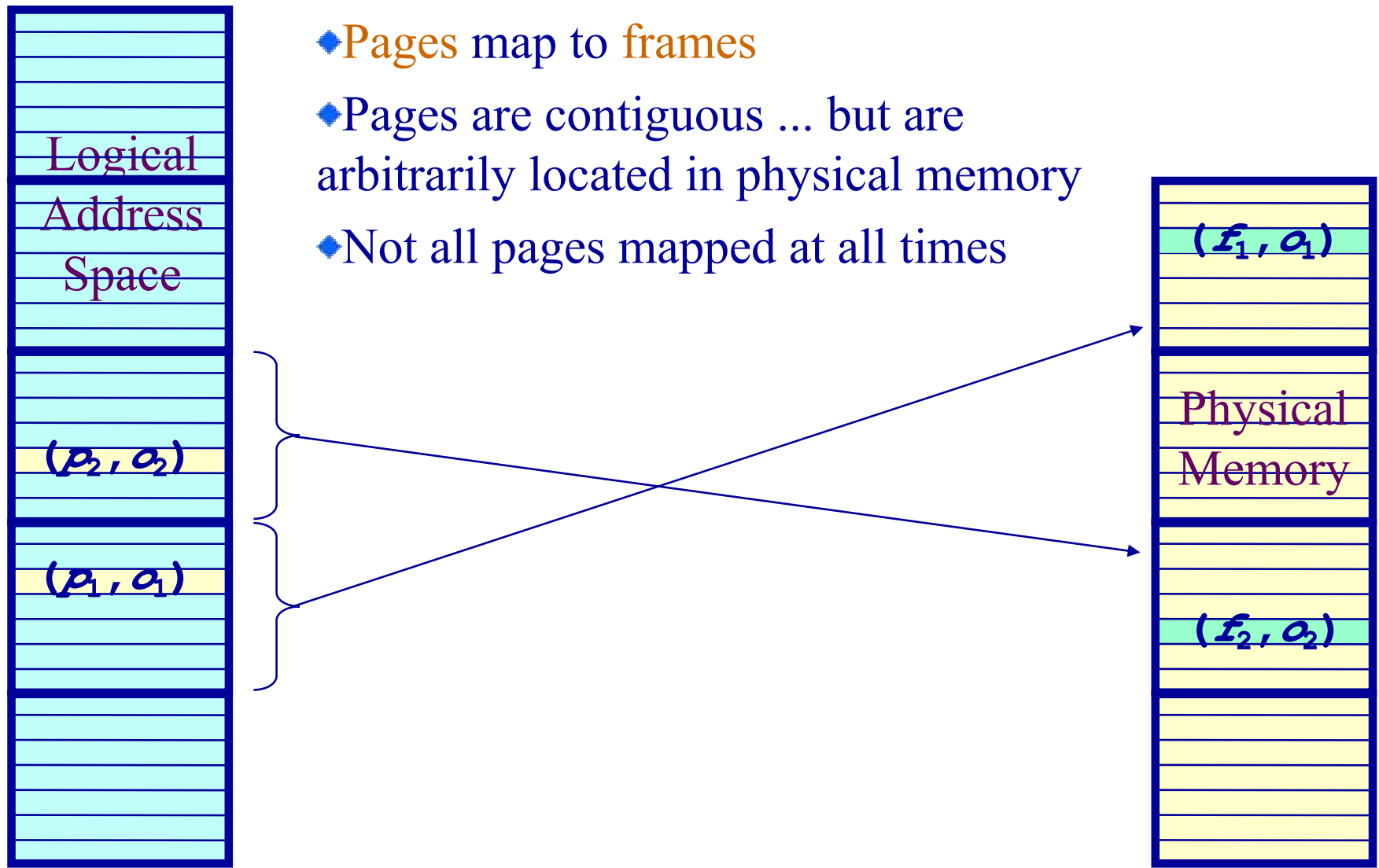$o$     — page offset ($2^S$ bytes/pages)

Virtual address = $2^S$ $p + o$

($p, o$)

$o$

Logical
Address
Space

$p$

(0,0)

addr:

$P+S$      $S$      1

$p$      $o$

# Non-contiguous Allocation : Paging

## Paging Model

Logical Address Space

$(p_2, o_2)$

$(p_1, o_1)$
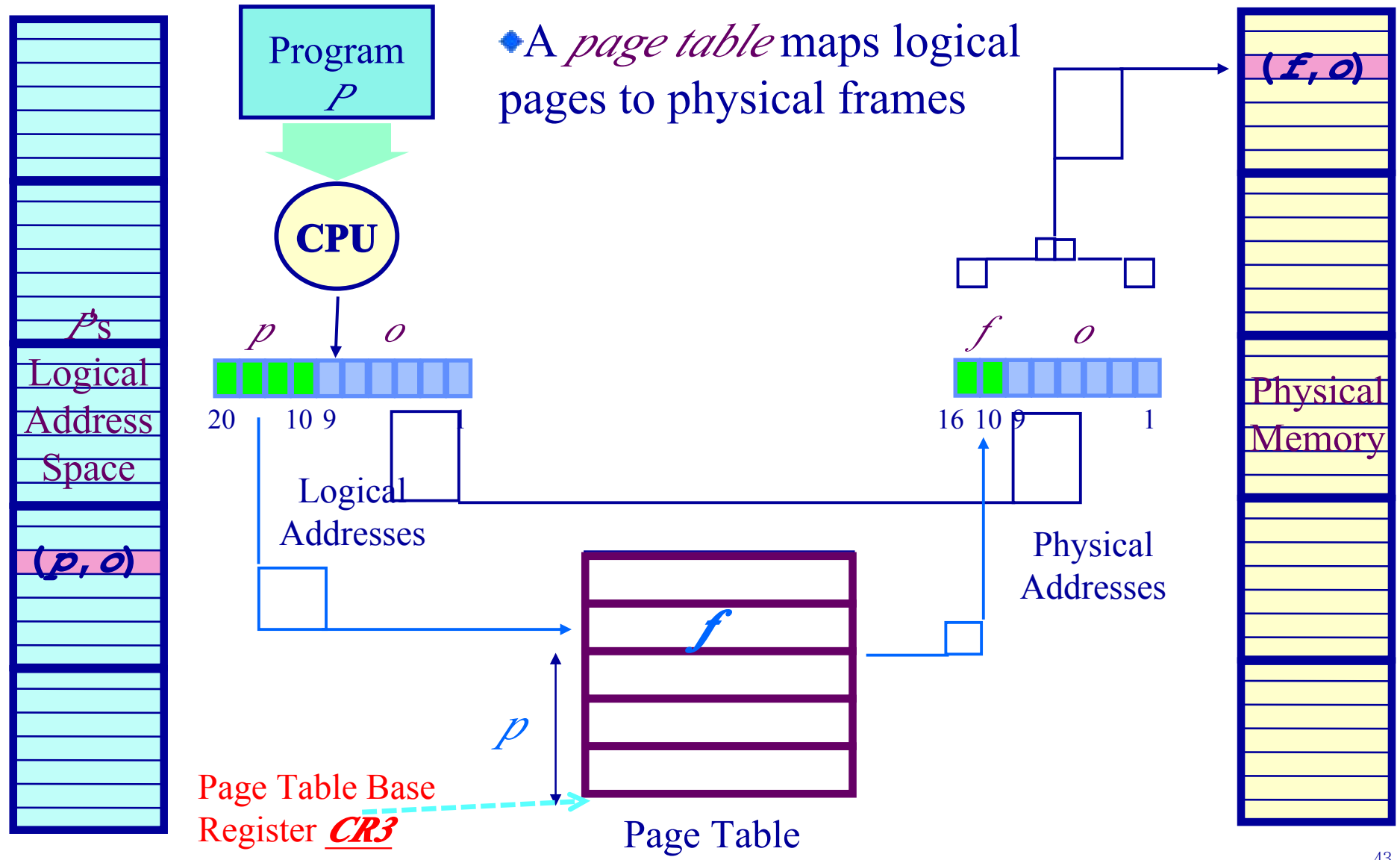
- Pages map to frames

- Pages are contiguous ... but are arbitrarily located in physical memory

- Not all pages mapped at all times

$(f_1, o_1)$

Physical Memory

$(f_2, o_2)$

# Non-contiguous Allocation : Paging

## Paging Hardware Architecture

Program $P$

CPU

A *page table* maps logical pages to physical frames

$(f, o)$

$P$s Logical Address Space

$(p, o)$

$p \quad o$

20 10 9 1

Logical Addresses

$f \quad o$

16 10 9 1

Physical Addresses

Physical Memory

$f$

$p$

Page Table Base Register **CR3**

Page Table

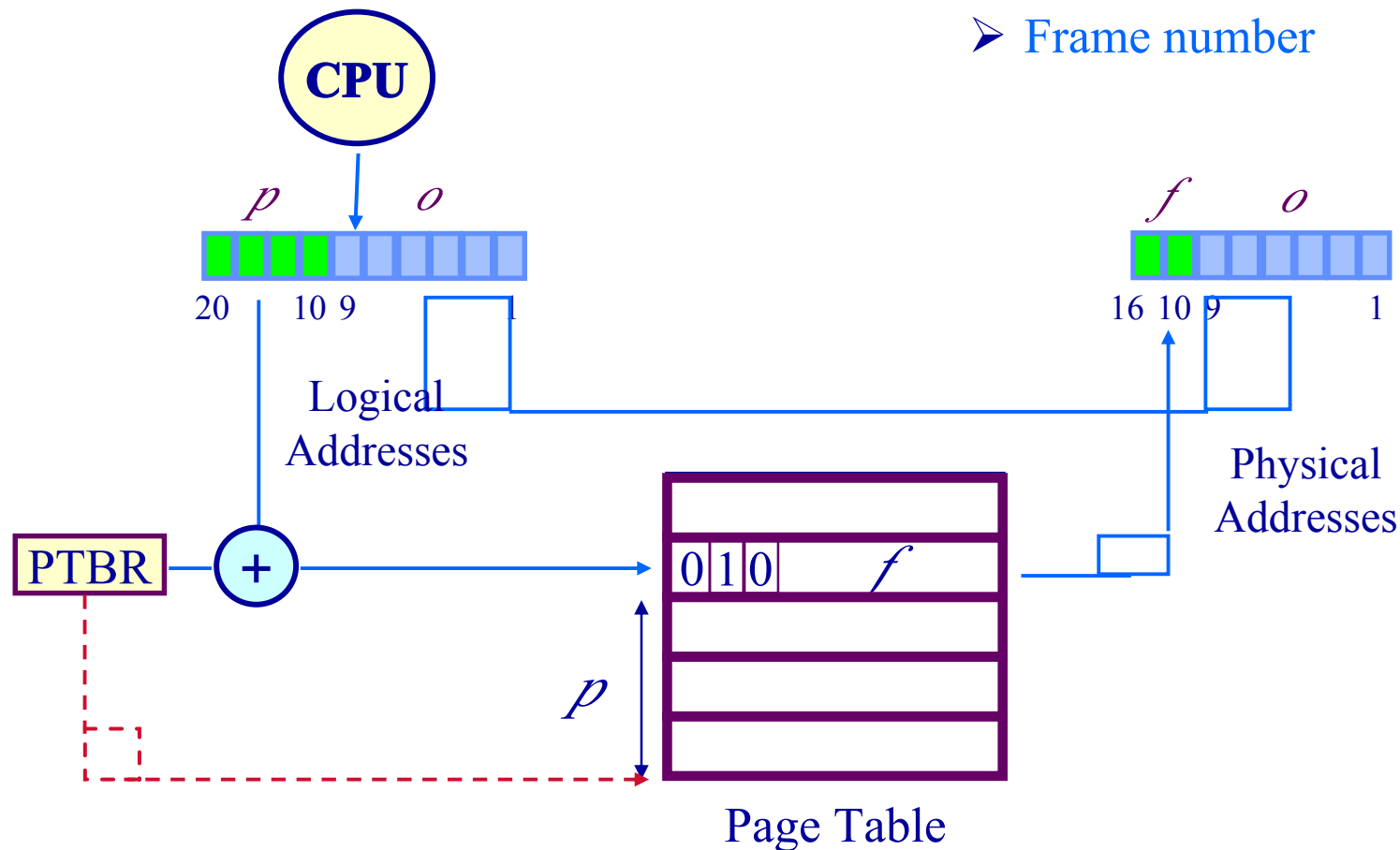- Computer Arch/Memory Hierarchy
- Address Space & Address Generation
- Contiguous Memory Allocation
  - ◆ Dynamic Allocation of Partitions
- Non-Contiguous Memory Allocation
  - ◆ Segmentation
  - ◆ Paging
  - ◆ Page Table
    - Translation Look-aside Buffer (TLB)
    - Multi-Level Page Table
    - Inverted Page Table
  - ◆ Paged Segmentation Model
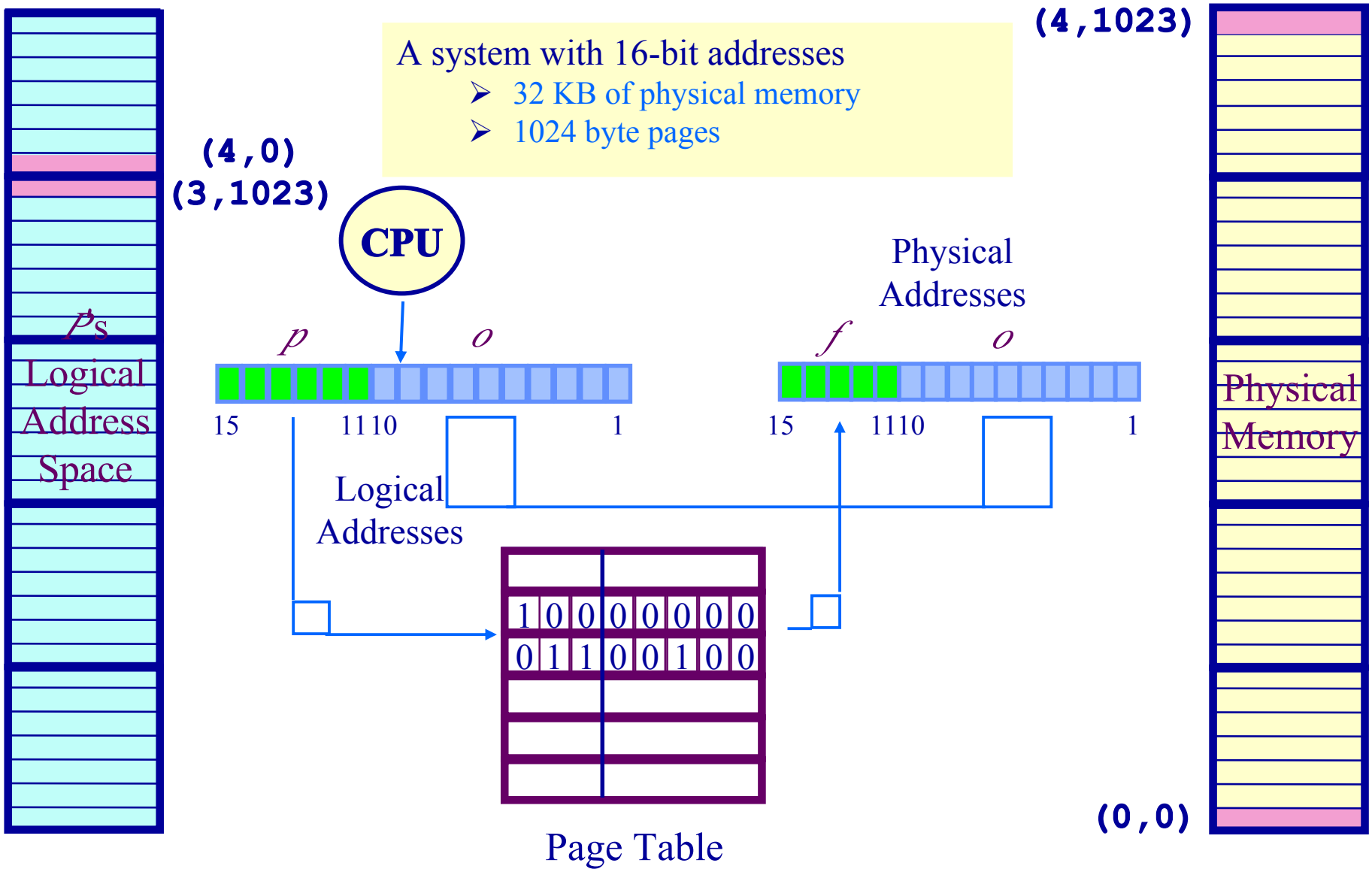
# Non-contiguous Allocation : Page Table

## Page Table Structure

- **One table per process**
  - Part of process's state
  - PTBR: Page Table Base Register

- Contents:
  - Flags — dirty bit, resident bit, clock/reference bit
  - Frame number



Page Table

## Example Address Translation

A system with 16-bit addresses
- 32 KB of physical memory
- 1024 byte pages

**(4,0)**
**(3,1023)**

**CPU**

$Ps$

Logical Address Space

$p$      $o$

15    1110          1

Logical Addresses

Physical Addresses

$f$      $o$

15    1110          1

**(4,1023)**

Physical Memory

| | |
|---|---|
| 1 0 0 | 0 0 0 0 0 |
| 0 1 1 | 0 0 1 0 0 |
| | |
| | |
| | |

Page Table

**(0,0)**
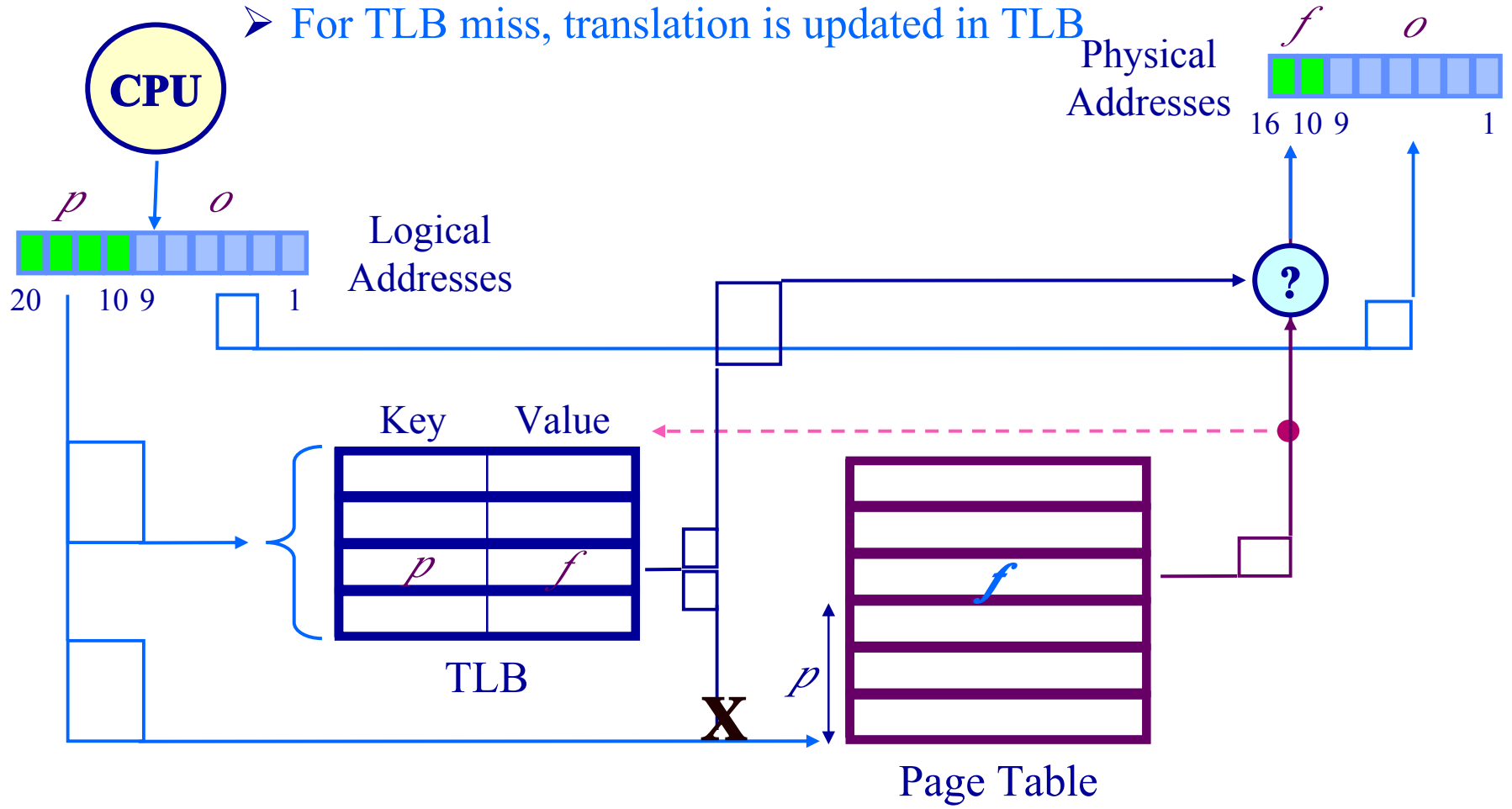
# Paging Performance Issue

◆ Problem — Requires 2 memory references!

> ➢ One access to get the page table entry

> ➢ One access to get the data

◆ Page table can be very large

> ➢ For a machine with 64-bit addresses and 1024 byte pages, what is the size of a page table?

◆ What to do? Hint: most computing problems are solved by some form of…

> ➢ Caching

> ➢ Indirection

- Computer Arch/Memory Hierarchy

- Address Space & Address Generation

- Contiguous Memory Allocation

  ◆ Dynamic Allocation of Partitions

- Non-Contiguous Memory Allocation

  ◆ Segmentation

  ◆ Paging

  ◆ Page Table
    - Translation Look-aside Buffer (TLB)
    - Multi-Level Page Table
    - Inverted Page Table

  ◆ Paged Segmentation Model

# Non-contiguous Allocation : Page Table

## Translation Look-aside Buffer (TLB)

- Cache recently accessed page-to-frame translations
  - ➢ TLB implemented in associative memory for fast access
  - ➢ For TLB hit, physical page number obtained in 1 cycle
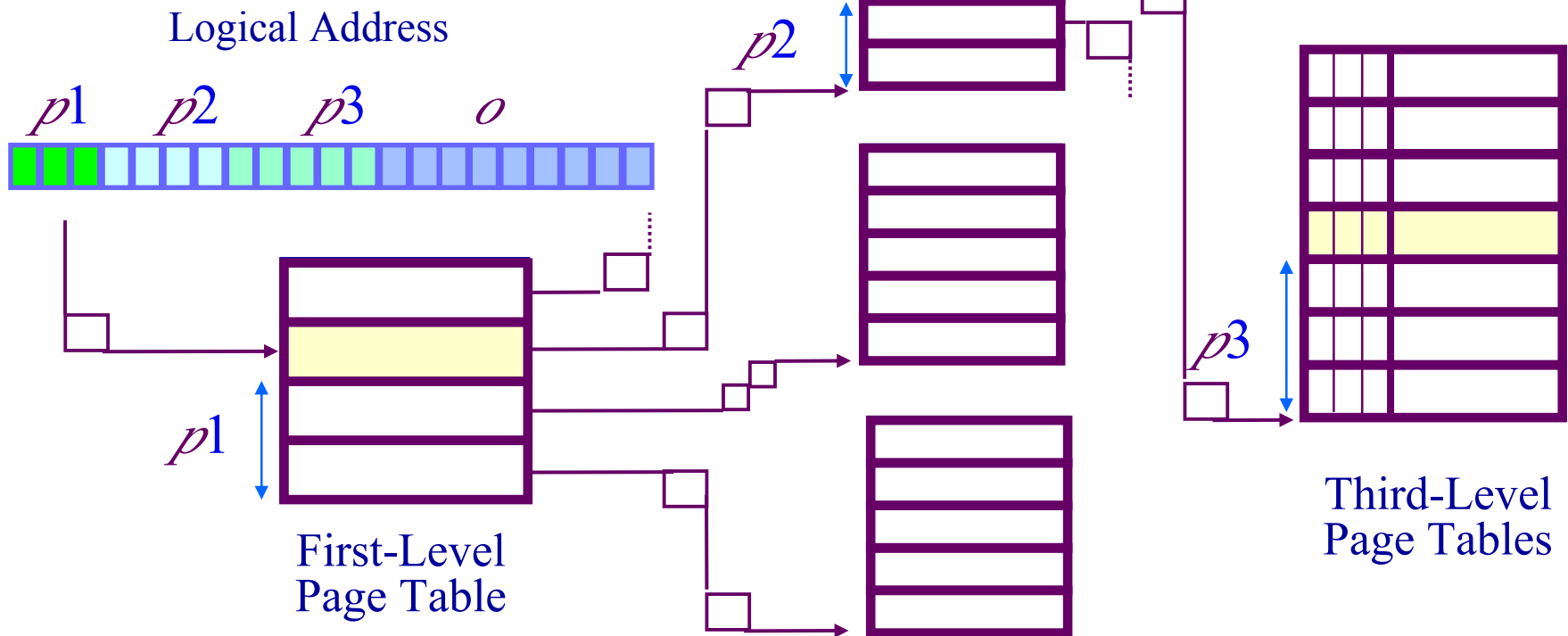  - ➢ For TLB miss, translation is updated in TLB



CPU

Logical Addresses

Physical Addresses

Key    Value

TLB

Page Table

## Multi-level Paging
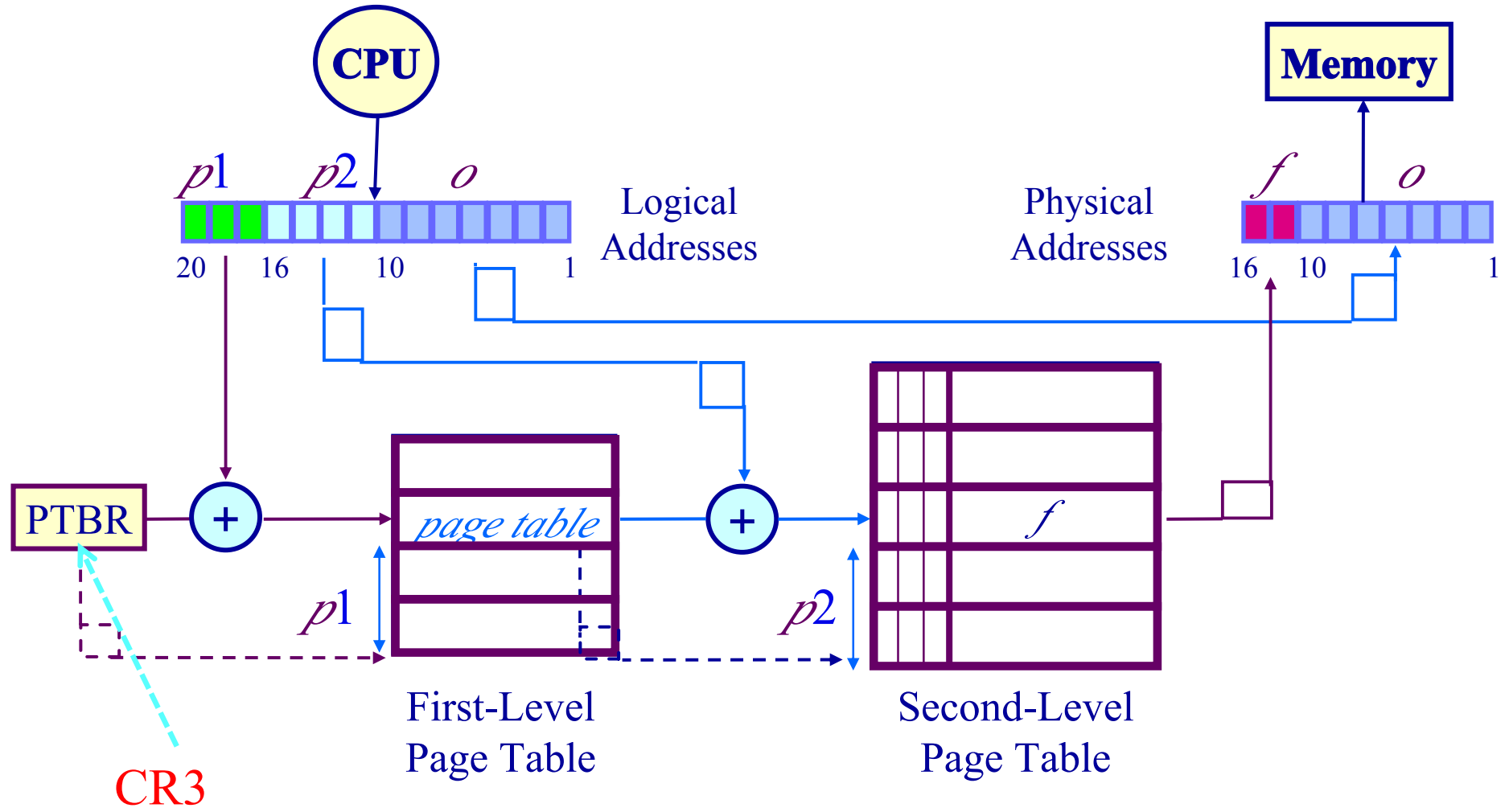
◆Add additional levels of indirection to the page table by sub-dividing page number into k parts

➢ Create a "tree" of page tables

Logical Address

$p1$   $p2$   $p3$   $o$

Second-Level Page Tables

$p2$

First-Level Page Table

$p1$

$p3$

Third-Level Page Tables

# Non-contiguous Allocation : Page Table

## Example: Two-level Paging

# Non-contiguous Allocation : Page Table

## The Problem of Large Address Spaces

◆With large address spaces (64-bits) forward mapped page tables become cumbersome.

➢ E.g. 5 levels of tables.

◆Instead of making tables proportional to size of logical address space, make them proportional to the size of physical address space.

➢ Logical (virtual) address space is growing faster than physical.

## Using Page Registers (aka Inverted Page Tables)

- Each frame is associated with a register containing
  - Residence bit: whether or not the frame is occupied
  - Occupier: page number of the page occupying frame
  - Protection bits

- Page registers: an example
  - Physical memory size: 16 MB
  - Page size: 4096 bytes
  - Number of frames: 4096
  - Space used for page registers (assuming 8 bytes/register): 32 Kbytes
  - Percentage overhead introduced by page registers: 0.2%
  - Size of virtual memory: irrelevant

## Page Registers Tradeoffs

◆ Advantages:

➢ Size of translation table occupies a very small fraction of physical memory

➢ Size of translation table is independent of logical address space size

◆ Disadvantages:

➢ We have reverse of the information that we need….

➢ How do we perform translation ?

➢ Search the translation table for the desired page number

## Searching for a Page in Inverted Page Tables

◆If the number of frames is small, the page registers can be placed in an associative memory

◆Logical page number looked up in associative memory

➢ Hit: frame number is extracted

➢ Miss: results in page fault

◆Limitations:

➢ Large associative memories are expensive

⚶ Difficult to make large and accessible in a single cycle.
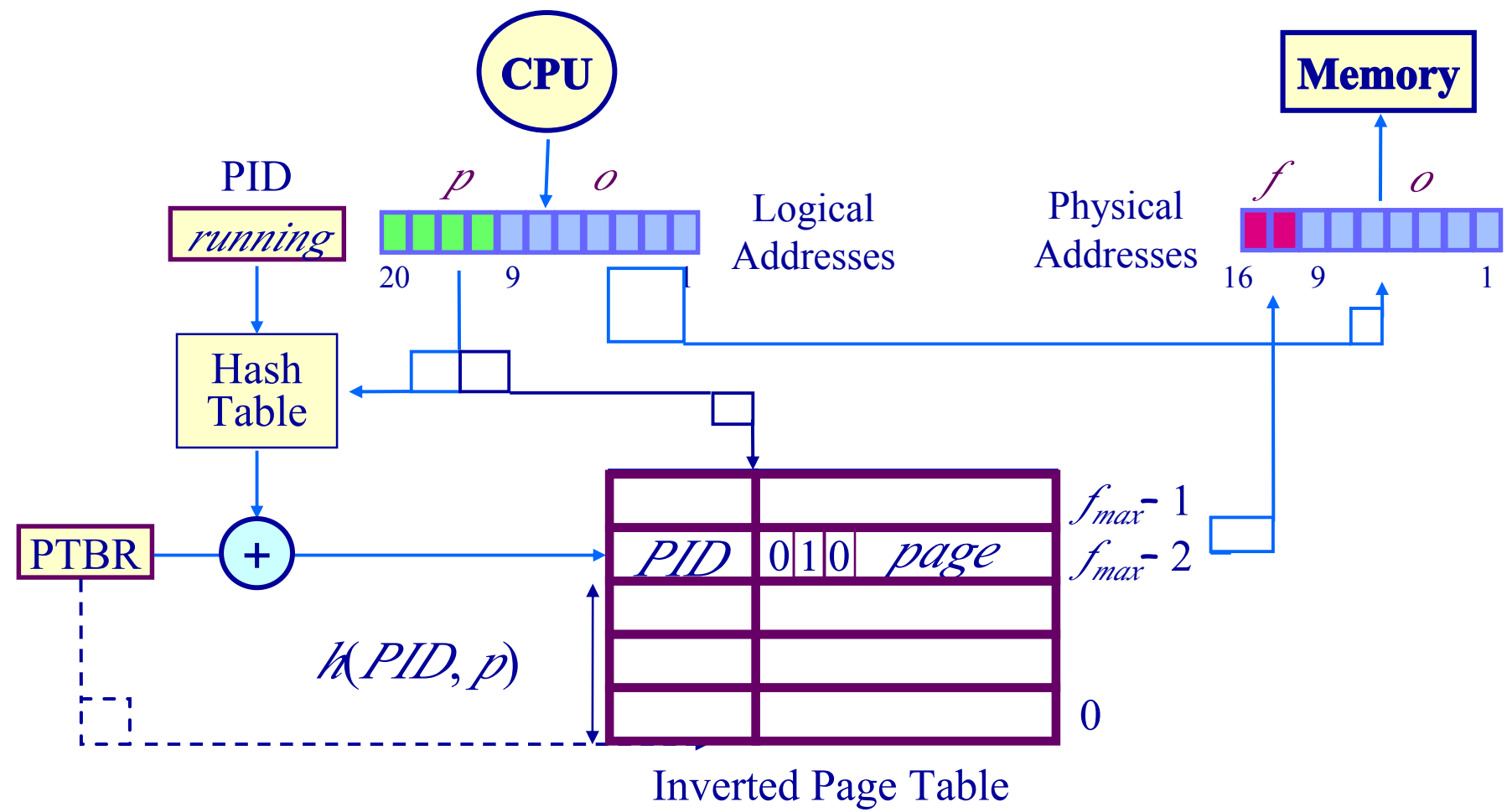
⚶ They consume a lot of power

# Non-contiguous Allocation : Page Table

## Hashing Large Inverted Page Tables

Hash page numbers to find corresponding frame numbers in a "frame" table with one entry per frame

Page i is placed in slot f(i) where f is an agreed-upon hash function

To lookup page i, perform the following:

- Compute f(i) and use it as an index into the table of page registers
- Extract the corresponding page register
- Check if the register tag contains i, if so, we have a hit
- Otherwise, we have a miss

## Hashed Inverted Page Table Architecture



CPU

Memory

PID

*running*

$p$    $o$

20    9

Logical
Addresses

Physical
Addresses

$f$    $o$

16    9    1

Hash
Table

PTBR

$+$

$h(PID, p)$

$PID$  $0$ $1$ $0$  $page$
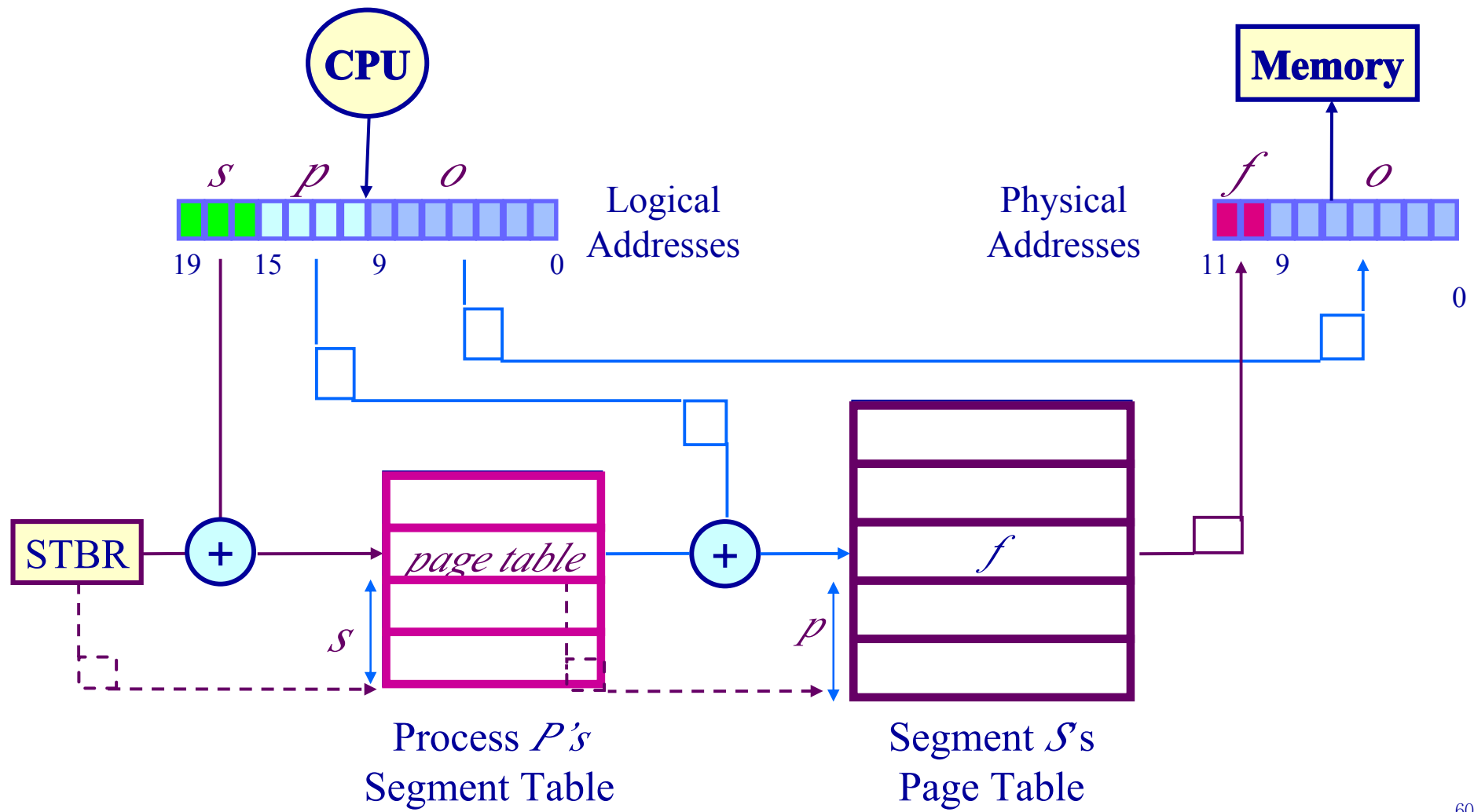
$f_{max} - 1$

$f_{max} - 2$

0

Inverted Page Table

● Computer Arch/Memory Hierarchy

● Address Space & Address Generation

● Contiguous Memory Allocation

◆ Dynamic Allocation of Partitions

● Non-Contiguous Memory Allocation

◆ Segmentation

◆ Paging

◆ Page Table

· Translation Look-aside Buffer (TLB)

· Multi-Level Page Table

· Inverted Page Table

● ◆ Paged Segmentation Model

# Paged Segmentation Model

- Segmentation has advantages for protection, paging has advantages for memory utilization and optimizing transfer to backing store.
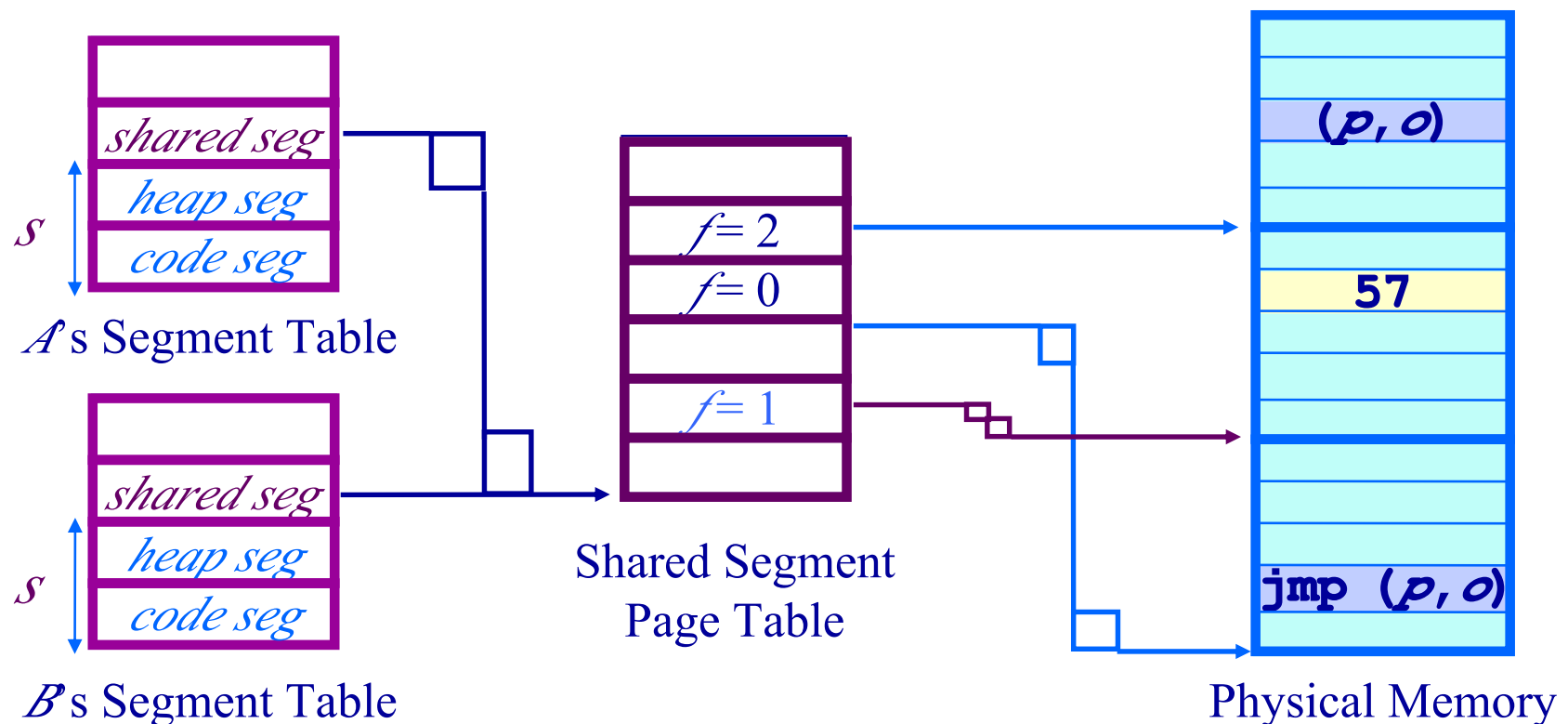
- Can we combine segmentation and paging?

# Paged Segmentation Hardware Architecture

◆ Add an additional level of indirection to page table

# Sharing in Paged-Segmented Systems

- If segments are paged then page tables are automatically shared
  - Processes need only agree on a number for the shared segment



$A$'s Segment Table

$B$'s Segment Table

Shared Segment
Page Table

Physical Memory

# This week's Work

Lab1 should be finished!