



北京交通大学

Operating System

Main Memory

Di Zhang

School of Software Engineering





北京交通大学

Course Review

What are the four requirements for Deadlock?

- **Mutual exclusion**

- Only one process at a time can use a resource.

- **Hold and wait**

- Process is **holding** at least one resource and is **waiting** to acquire additional resources held by other processes

- **No preemption**

- Resources are released only voluntarily by the process holding the resource, after process is finished with it

- **Circular wait**

- There exists a set $\{T_1, \dots, T_n\}$ of waiting processes
 - T_1 is waiting for a resource that is held by T_2
 - T_2 is waiting for a resource that is held by T_3
 - ...
 - T_n is waiting for a resource that is held by T_1

What are the methods for handling deadlocks?

- Ignore the problem and **pretend** that deadlocks never occur in the system
 - Used by most operating systems, including UNIX and Windows
- Ensure that the system will **never** enter a deadlock state
 - deadlock prevention
 - deadlock avoidance
- Allow the system to enter a deadlock state and then recover
 - deadlock detection
 - deadlock recovery

How the Banker's algorithm works? (safty)

1. Let **Work** and **Finish** be vectors of length **m** and **n**, respectively. Initialize:
 - (a) **Work** = **Available**
 - (b) **Finish** [**i**] = **false** for **i** = 0, 1, ..., **n** - 1.
2. Find an **i** such that both:
 - (a) **Finish** [**i**] = **false**
 - (b) **Need**_{**i**} ≤ **Work**If no such **i** exists, go to step 4.
3. **Work** = **Work** + **Allocation**_{**i**}
Finish[**i**] = **true**
go to step 2.
4. If **Finish** [**i**] == **true** for all **i**, then the system is in a **safe** state.

How the Banker's algorithm works? (request)

- **Request_i** = request vector for process **P_i**.
- If **Request_i[j] = k** then process **P_i** wants **k** instances of resource type **R_j**.
 1. If **Request_i ≤ Need_i** go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim.
 2. If **Request_i ≤ Available**, go to step 3. Otherwise **P_i** must wait, since resources are not available.
 3. Pretend to allocate requested resources to **P_i** by modifying the state as follows:
 - Available = Available – Request_i**;
 - Allocation_i = Allocation_i + Request_i**;
 - Need_i = Need_i – Request_i**;
- If **safe** \Rightarrow the resources are allocated to **P_i**.
- If **unsafe** \Rightarrow **P_i** must wait, and the old resource-allocation state is restored

How to detect deadlocks? (detection algorithm)

- 1. Let *Work* and *Finish* be vectors of length *m* and *n*, respectively. Initialize:
 - (a) *Work* = *Available*
 - (b) For $i = 1, 2, \dots, n$, if $Allocation_i \neq 0$, then *Finish*[*i*] = *false*; otherwise, *Finish*[*i*] = *true*.
- 2. Find an index *i* such that both:
 - (a) *Finish*[*i*] == *false*
 - (b) $Request_i \leq Work$
 - If no such *i* exists, go to step 4.
- 3. *Work* = *Work* + *Allocation*_{*i*}
Finish[*i*] = *true*
go to step 2.
- 4. If *Finish*[*i*] == *false*, for some $i, 1 \leq i \leq n$, then the system is in **deadlock** state. Moreover, if *Finish*[*i*] == *false*, then P_i is deadlocked.

What to do when detect deadlocks?

- Terminate process, force it to give up resources
- Preempt resources without killing off process
- Roll back actions of deadlocked processes
- Many operating systems use other options



Main Memory

Outline

- Background
- Contiguous Memory Allocation
- Paging
- Structure of the Page Table
- Segmentation

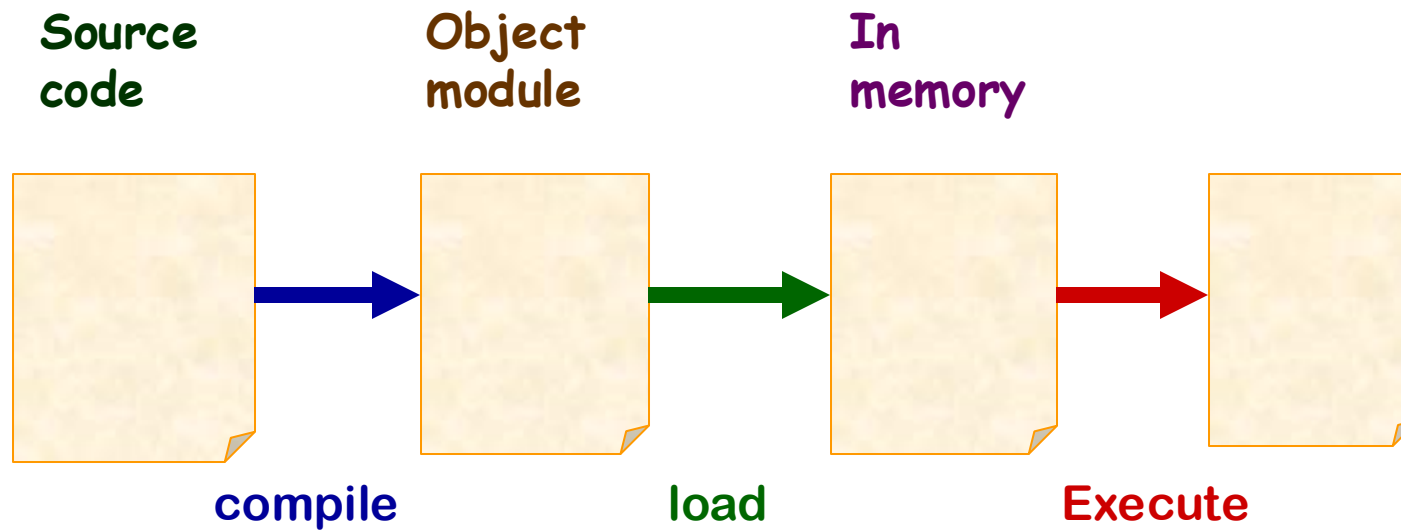


Background

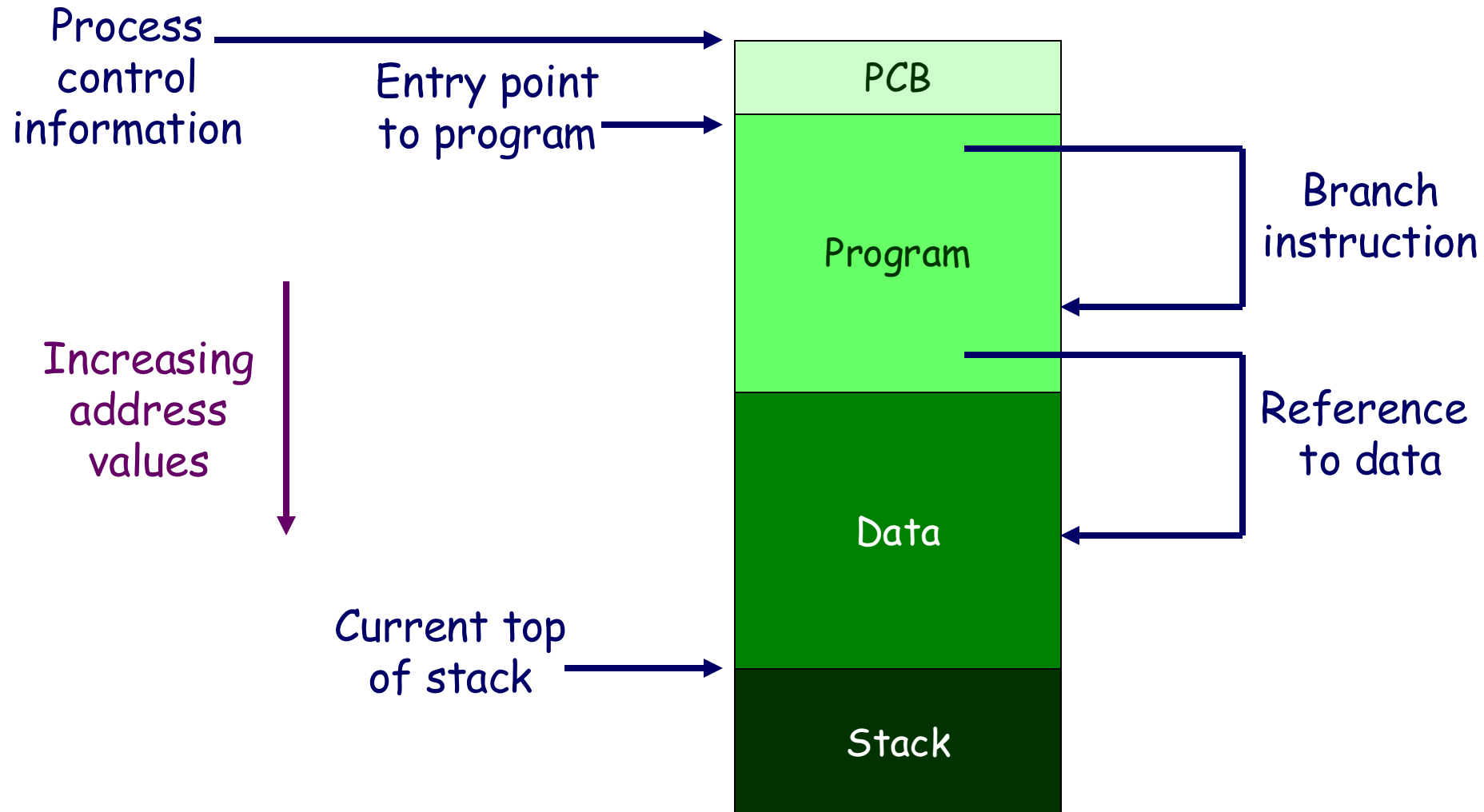
Recall

- Program must be brought (from disk) into memory and placed within a process for it to be run
 - All **data** in memory before and after processing
 - All **instructions** in memory in order to execute

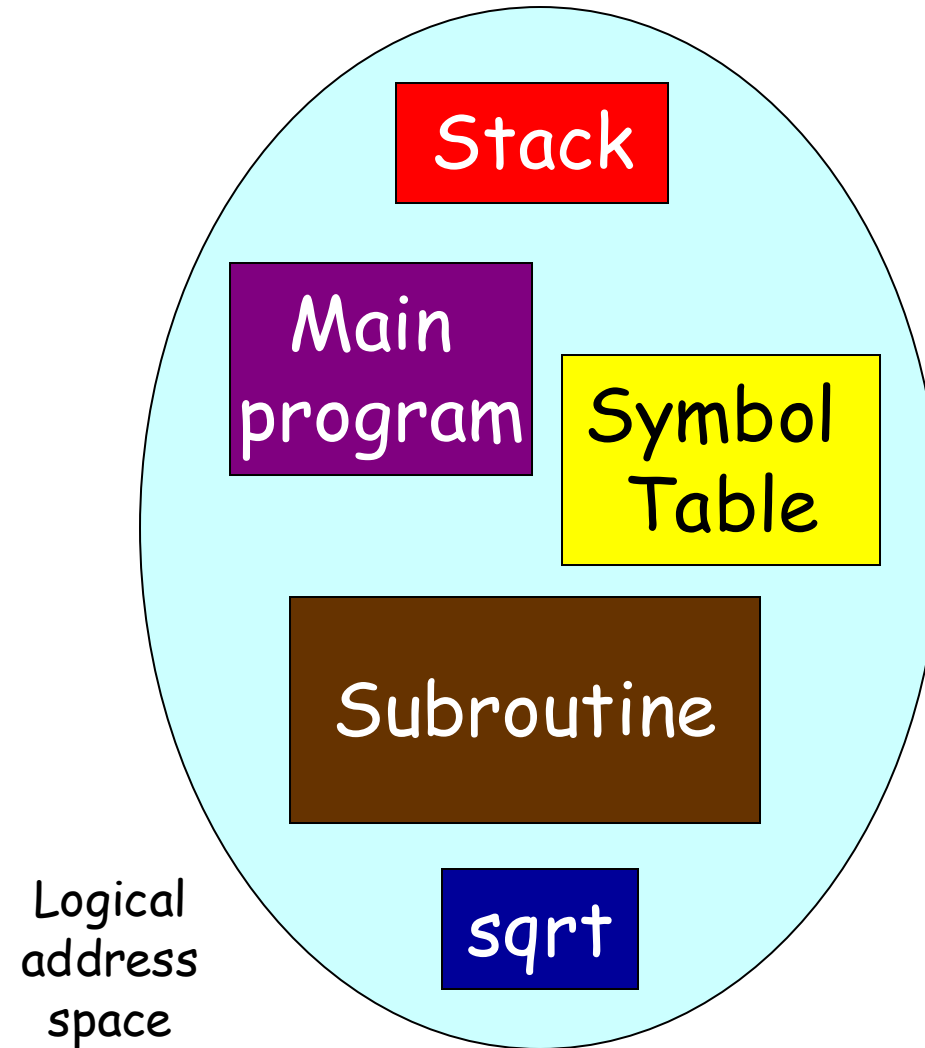
Background



Process in Memory



User's View of a Program

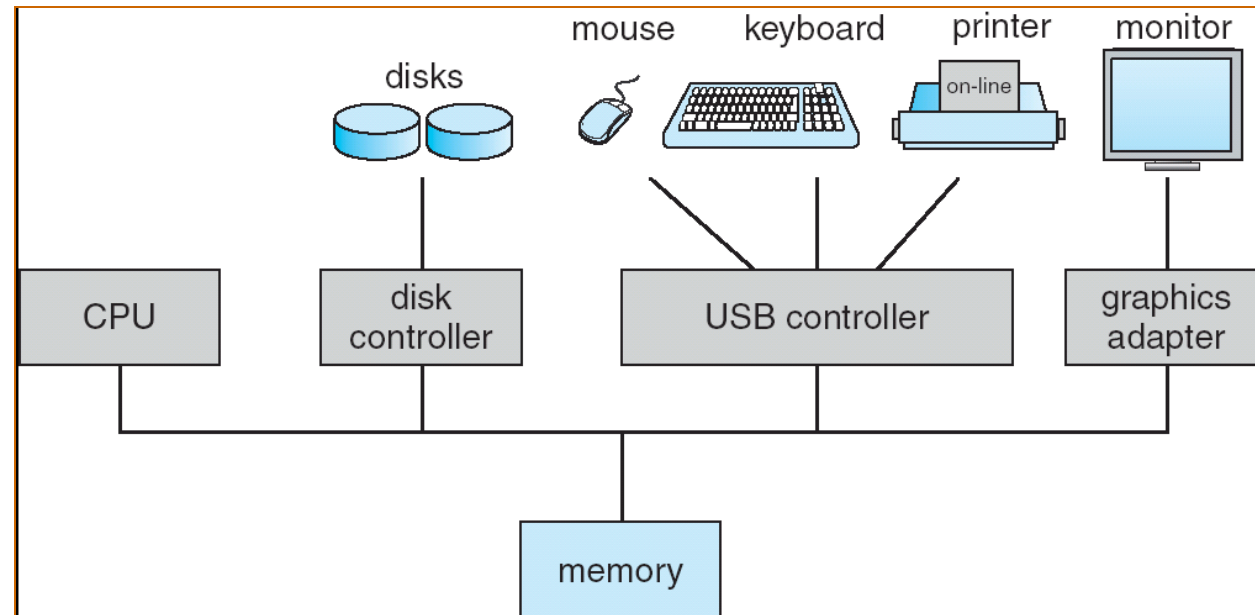


Recall

- Program must be brought (from disk) into memory and placed within a process for it to be run
 - All data in memory before and after processing
 - All instructions in memory in order to execute
- Main memory and registers are the **only** storage which CPU can access directly

Computer System Organization

- Computer-system operation
 - One or more CPUs, device controllers connect through common bus providing access to shared memory
 - Concurrent execution of CPUs and devices competing for memory cycles

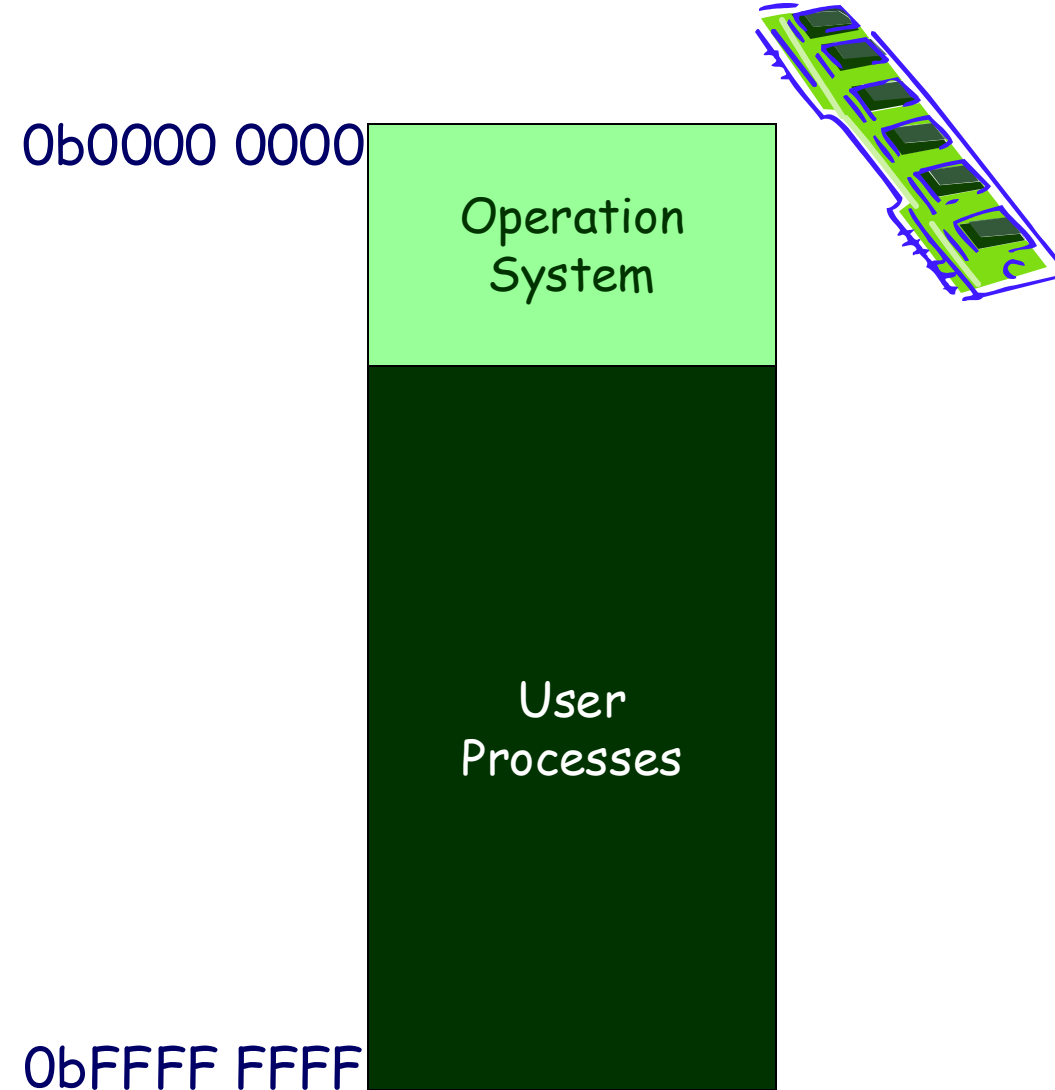


Main memory

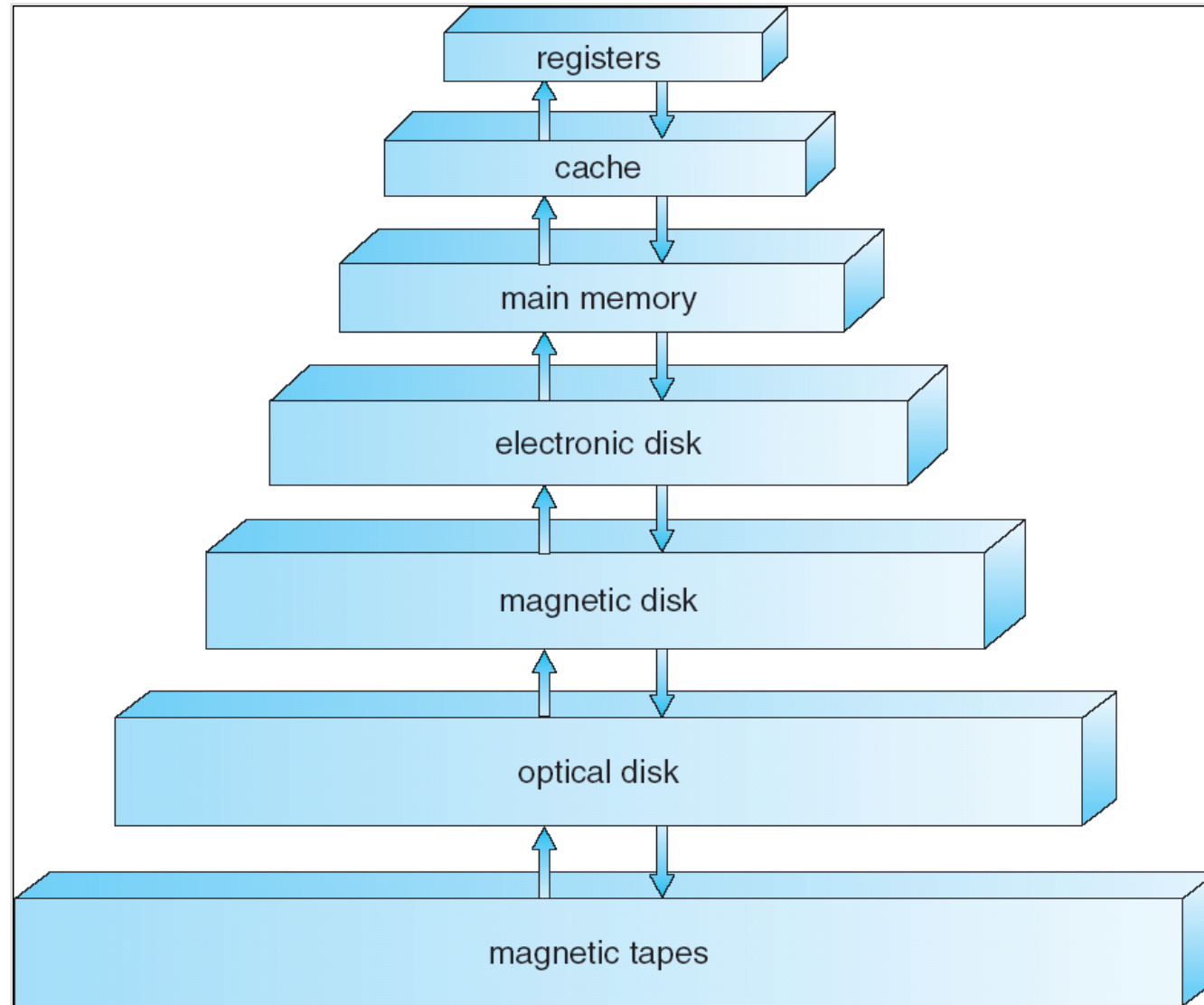


Main memory

- Address Space:
 - All the addresses and state a process can touch
- Main memory usually divided into two partitions:
 - Resident operating system, usually held in low memory
 - User processes then held in high memory

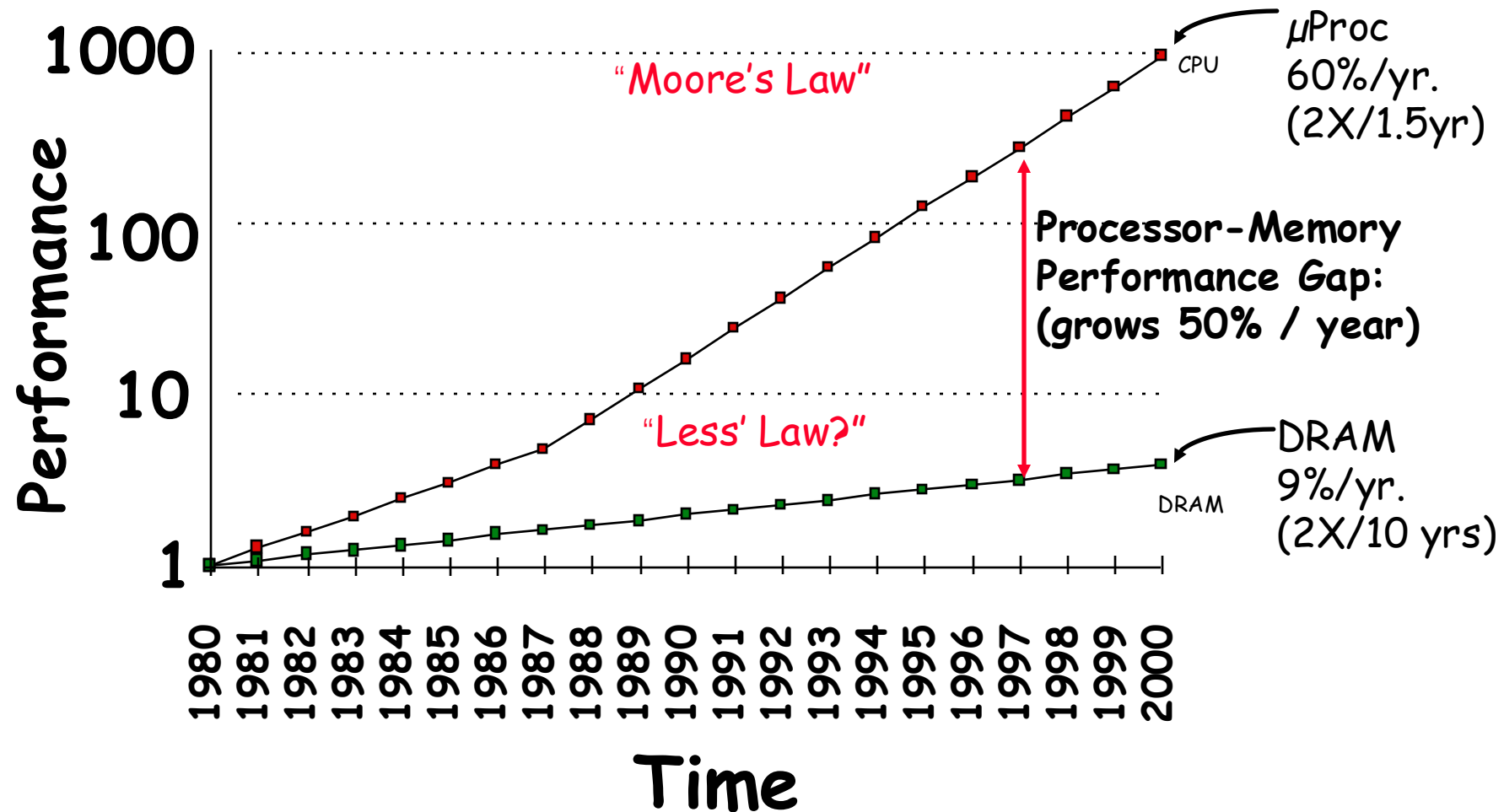


Storage-Device Hierarchy



Memory

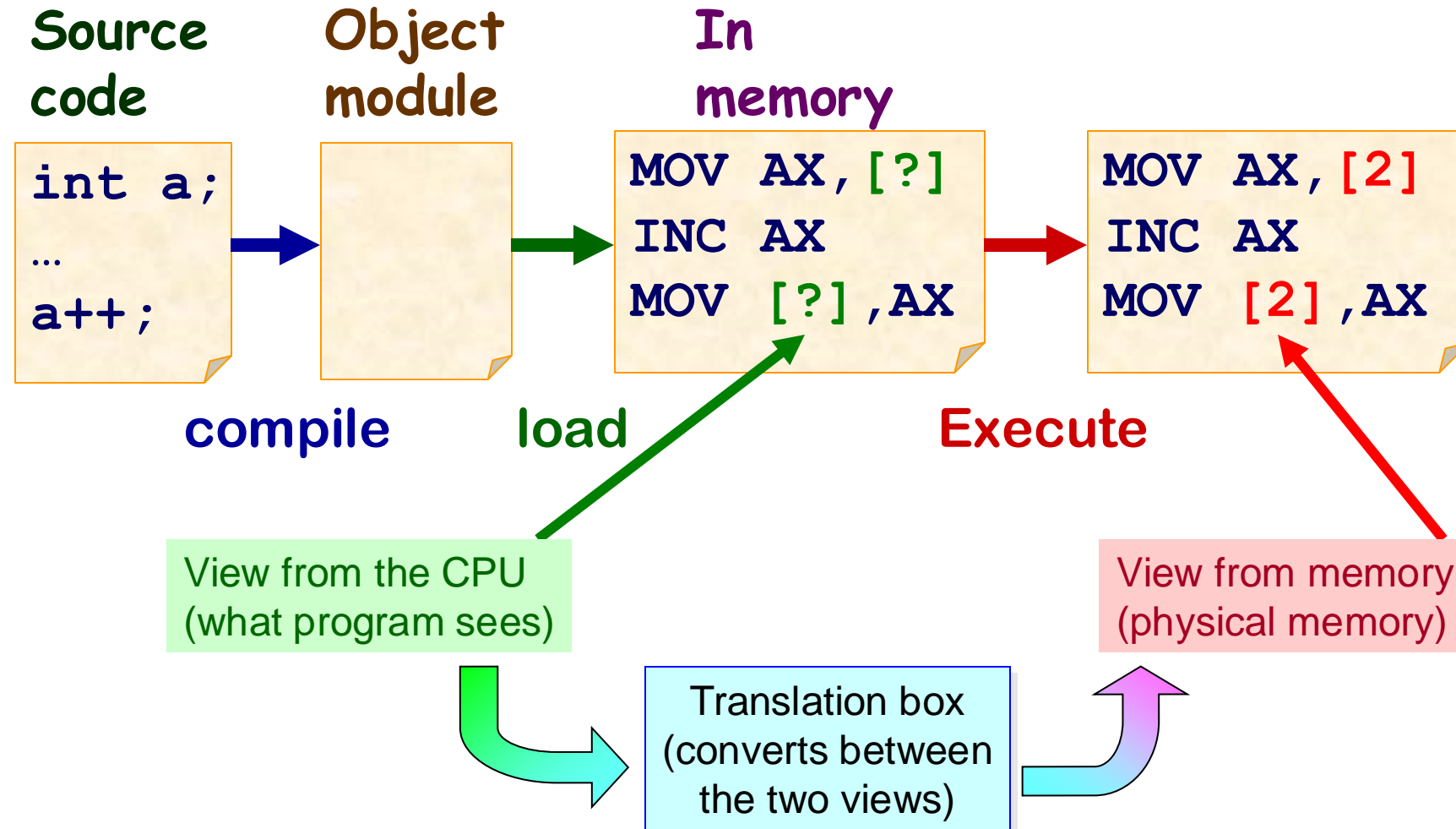
Processor-DRAM Memory Gap (latency)



Memory Management

- Memory management activities
 - Keeping **track** of which parts of memory are currently being used and by whom
 - Deciding which processes (or parts thereof) and data to **move into and out** of memory
 - **Allocating** and **deallocating** memory space as needed

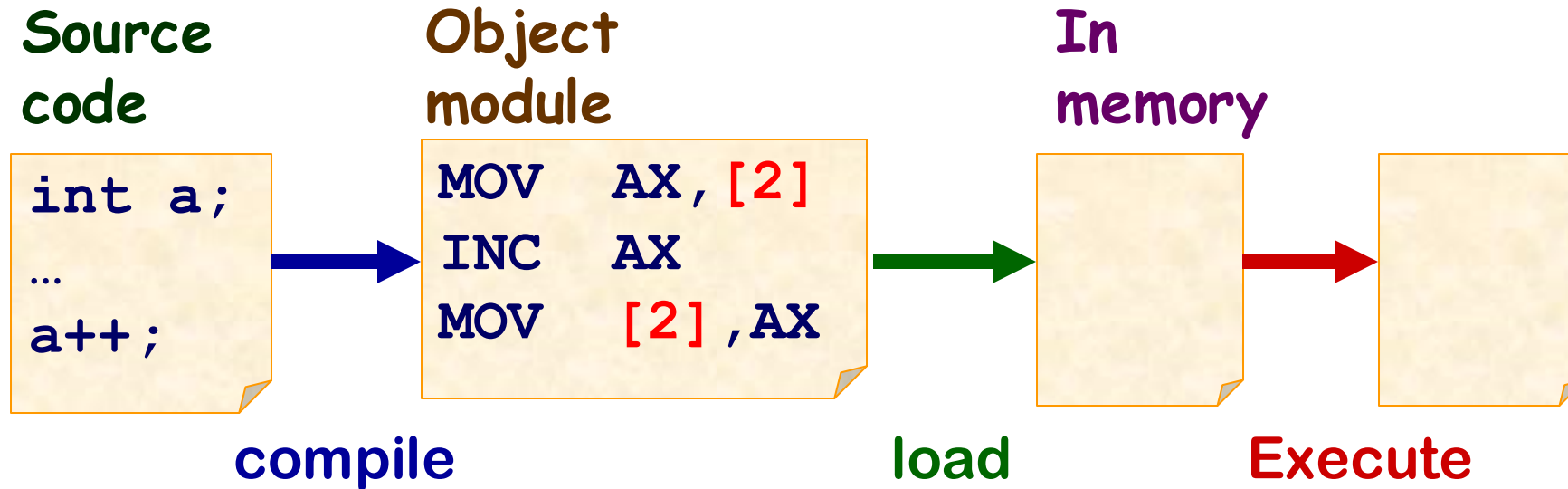
Two Views of Memory



Logical vs. Physical Address Space

- **Logical address** – generated by the CPU
- **Physical address** – address seen by the memory unit
- Logical and physical addresses are **the same in compile-time and load-time address-binding schemes**
- Logical and physical addresses **differ in execution-time address-binding scheme**

Address binding scheme: in compile time



- Address binding of instructions and data to memory addresses happen at:
 - **1. Compile time:** If memory location known a priori, **absolute code** can be generated; must recompile code if starting location changes.

Address binding scheme: in compile time

- Example

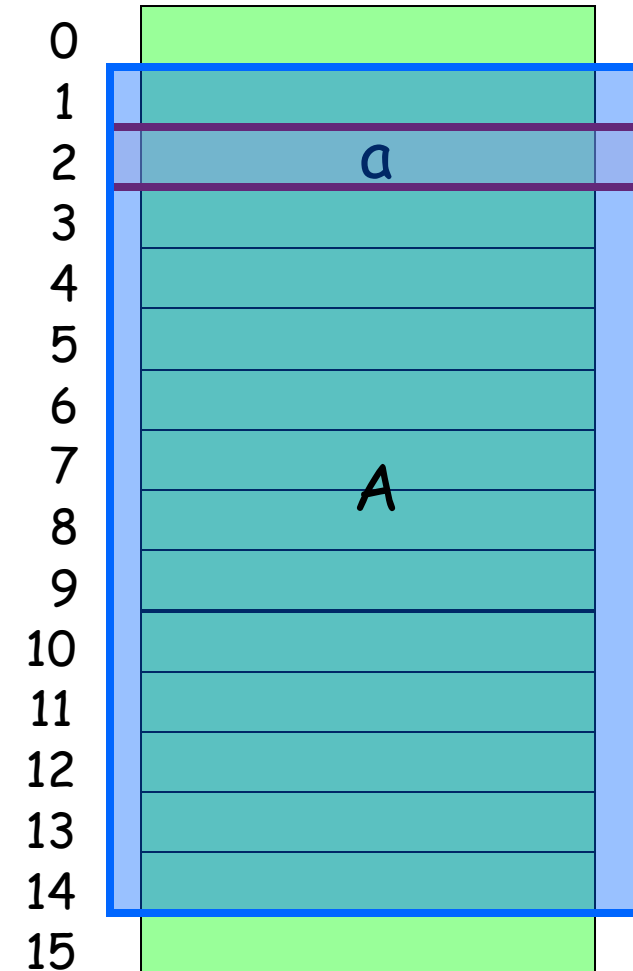
```
int a;  
...  
a++;
```

```
MOV  AX, a  
INC  AX  
MOV  a, AX
```

Where is a

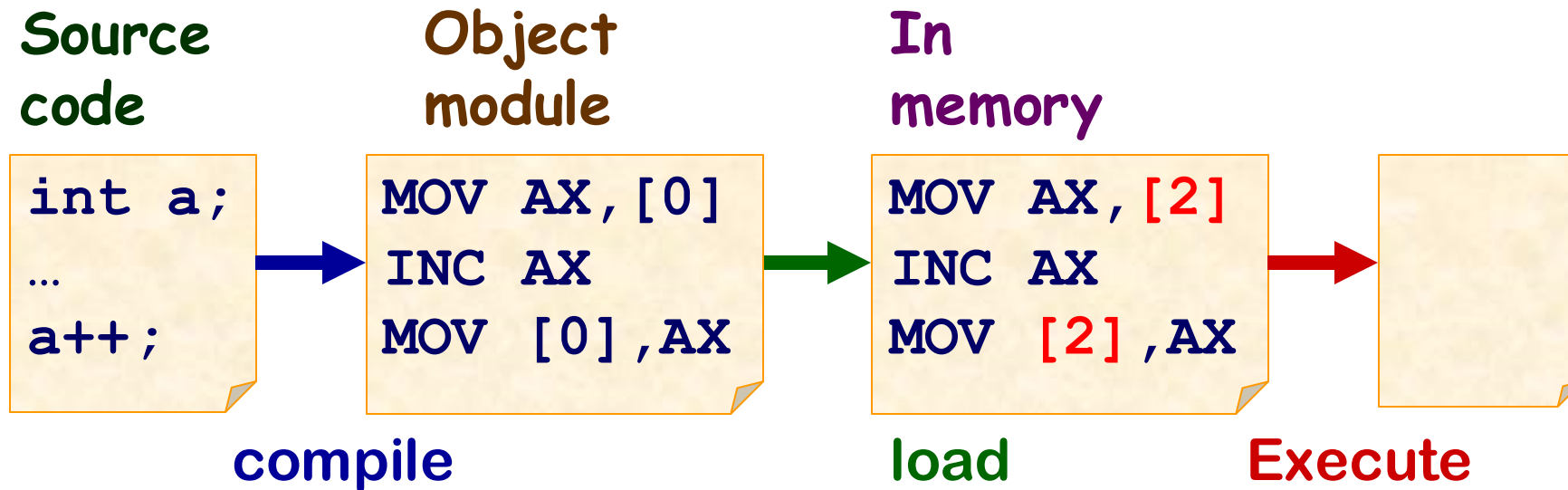
```
MOV  AX, [2]  
INC  AX  
MOV  [2], AX
```

```
Huge.exe  
...  
int A[14];  
...
```



Address binding scheme: in load time

- Address binding of instructions and data to memory addresses happen at
 - **2. Load time**: must generate **relocatable code** if memory location is not known at compile time



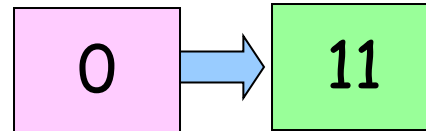
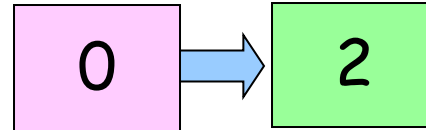
Address binding scheme: in load time

- Example

```
int a;  
...  
a++;
```

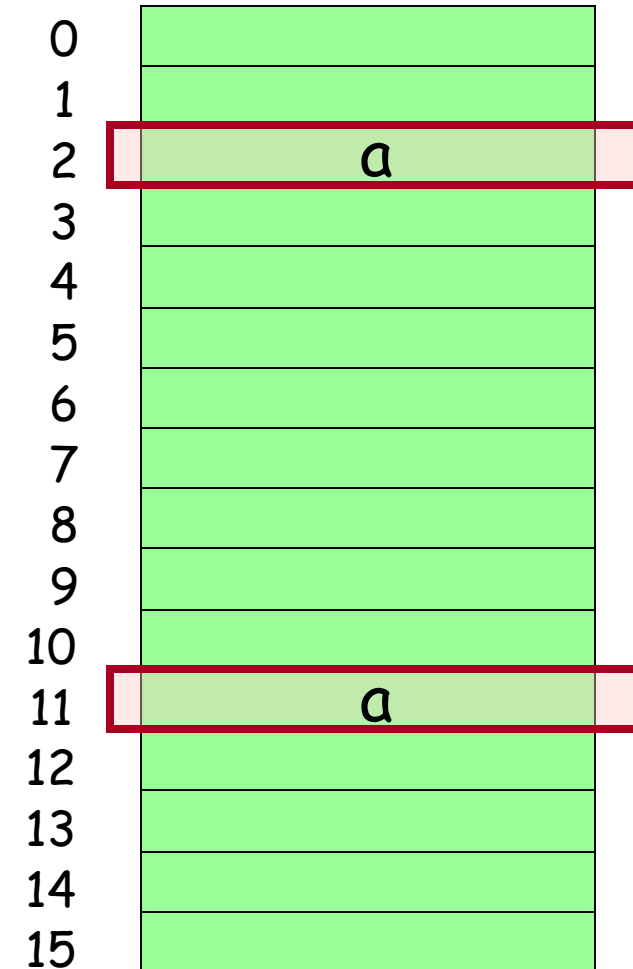
```
MOV  AX, a  
INC  AX  
MOV  a, AX
```

```
MOV  AX, [0]  
INC  AX  
MOV  [0], AX
```



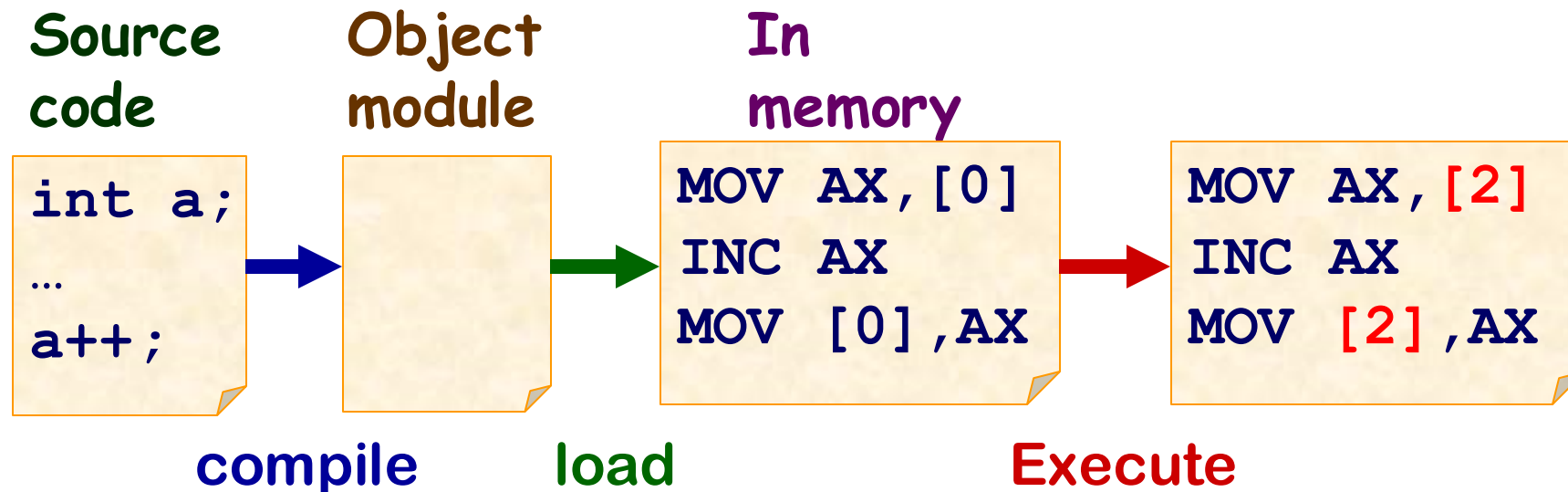
```
MOV  AX, [2]  
INC  AX  
MOV  [2], AX
```

```
MOV  AX, [11]  
INC  AX  
MOV  [11], AX
```



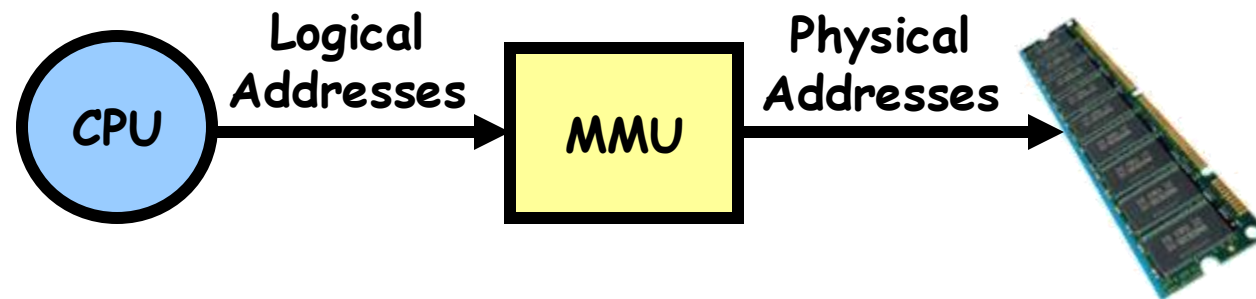
Dynamic Relocation

- Address binding of instructions and data to memory addresses happen at
 - **3. Execution time:** Binding delayed until run time if the process can be moved during its execution from one memory segment to another. Need hardware support for address maps (e.g., base and limit registers)



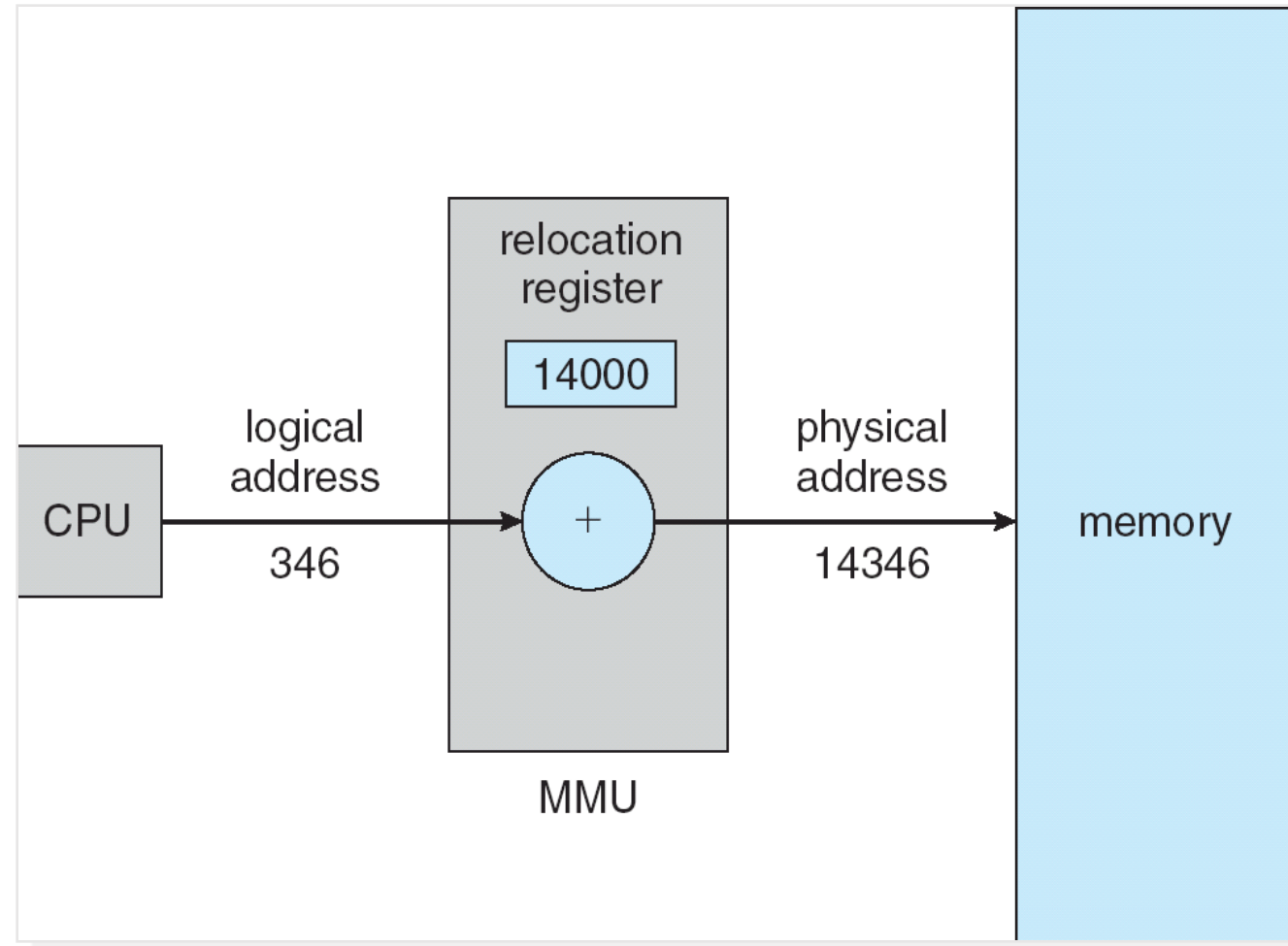
Memory-Management Unit (MMU)

- Hardware device that maps logical address to physical address
- In MMU scheme, the value in the relocation register is added to every address generated by a user process at the time it is sent to memory
- The user program deals with **logical** addresses; it never sees the **real** physical addresses



Memory-Management Unit (MMU)

- Dynamic relocation using a relocation register



Memory Management Requirement

- Protection
 - Processes should not be able to reference memory locations in another process without permission
- Sharing
 - Allow several processes to access the same portion of memory
 - Better to allow each process access to the same copy of the program rather than have their own separate copy

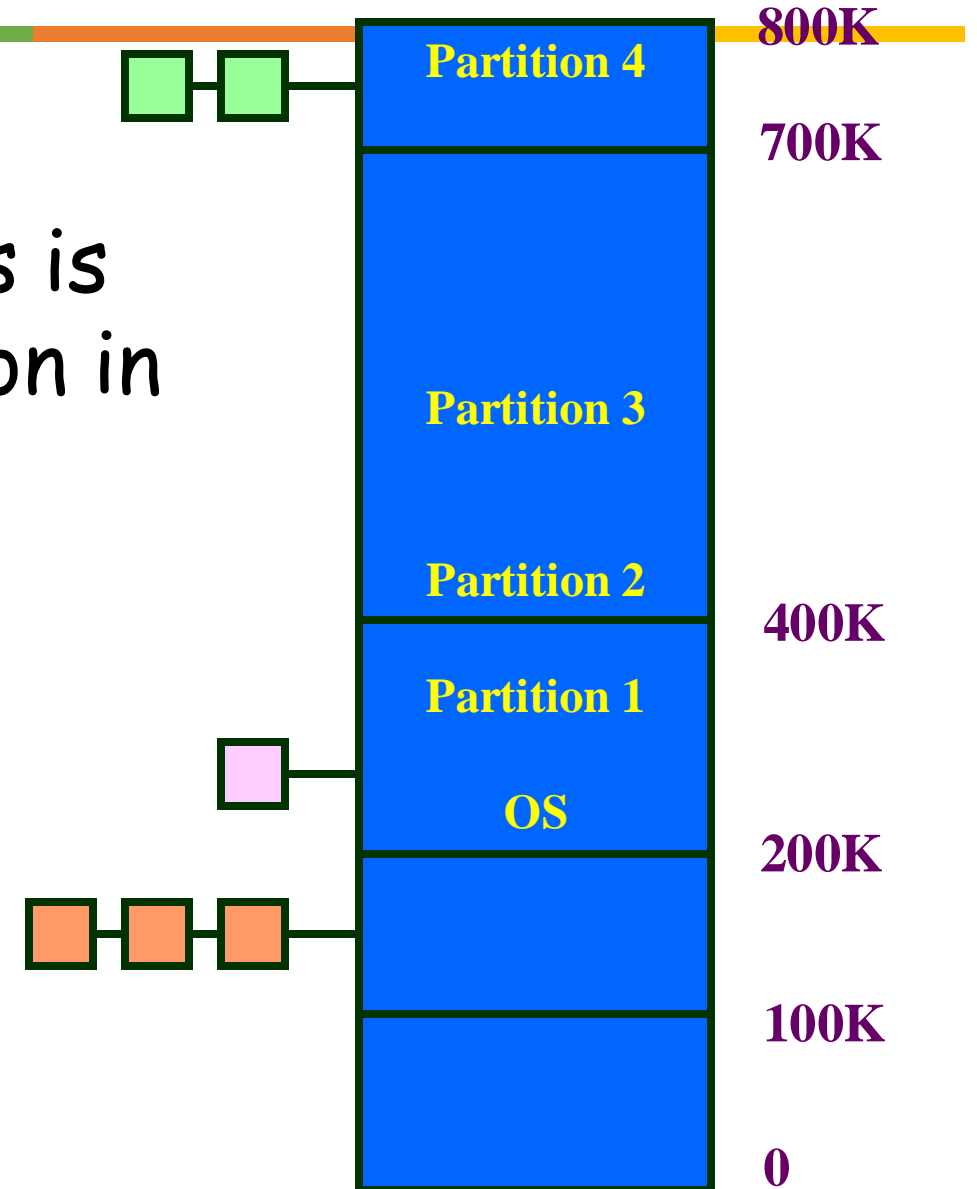


Contiguous Memory Allocation

Fixed Partition

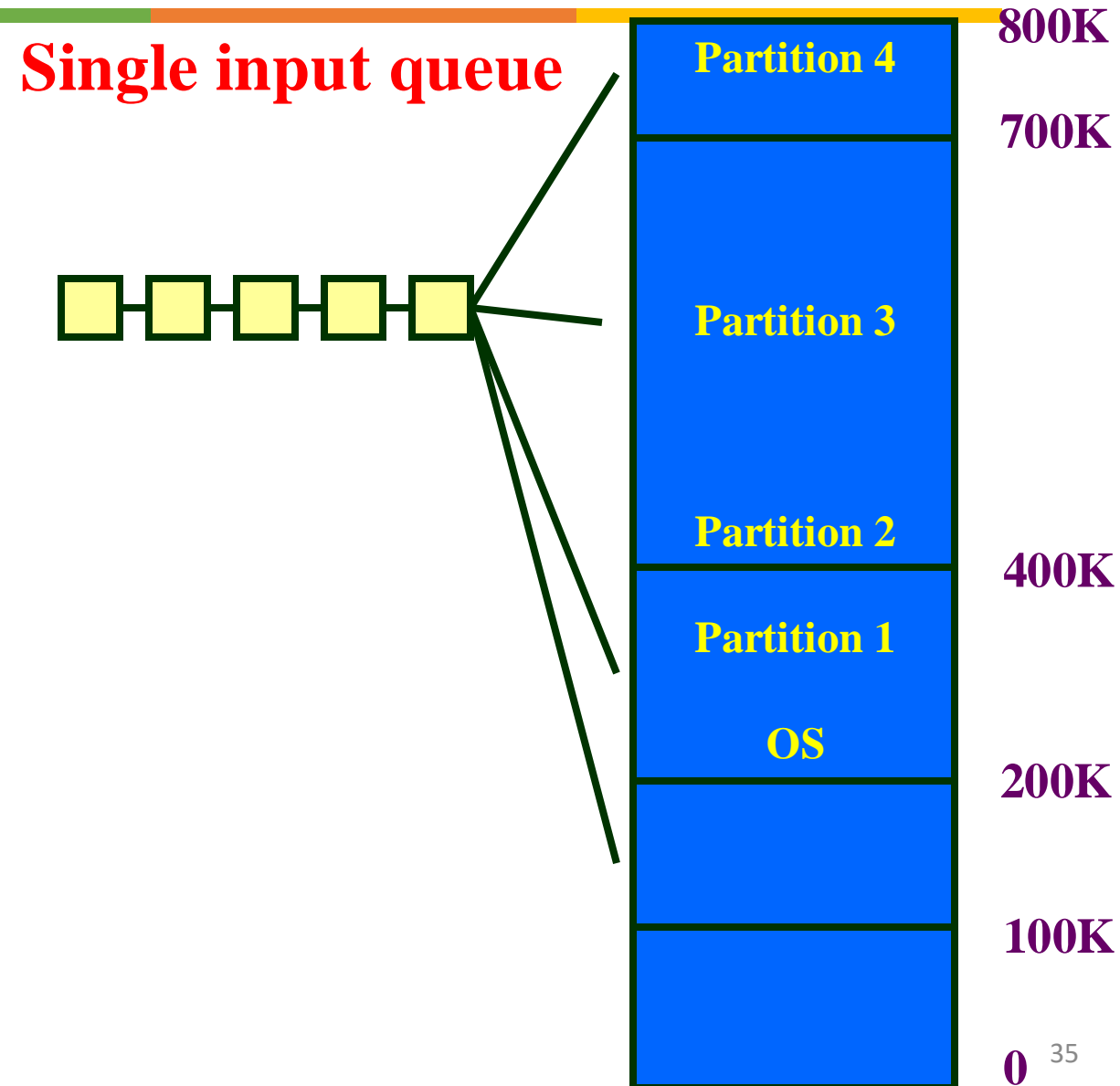
Multiple input queues

- One job each partition
- In multiple queues, each process is assigned to the smallest partition in which it fits and minimizes the internal fragmentation problem.



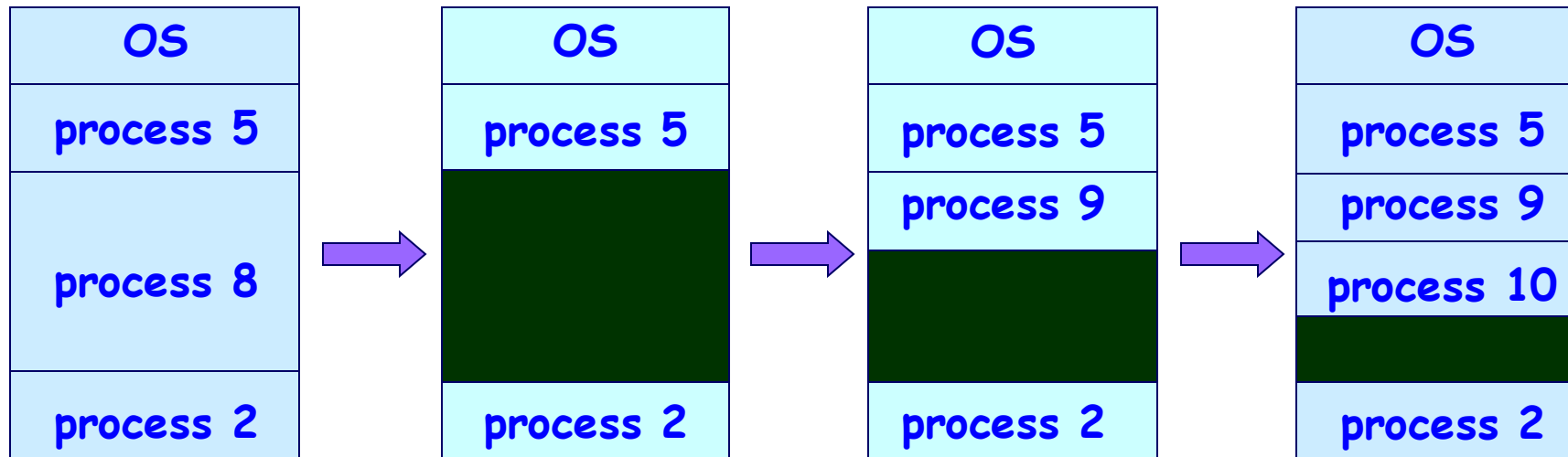
Fixed Partition

- In single queue, the process is assigned to the smallest available partition and the level of multiprogramming is increased.
- Main memory use is inefficient. Any program, no matter how small, occupies an entire partition.

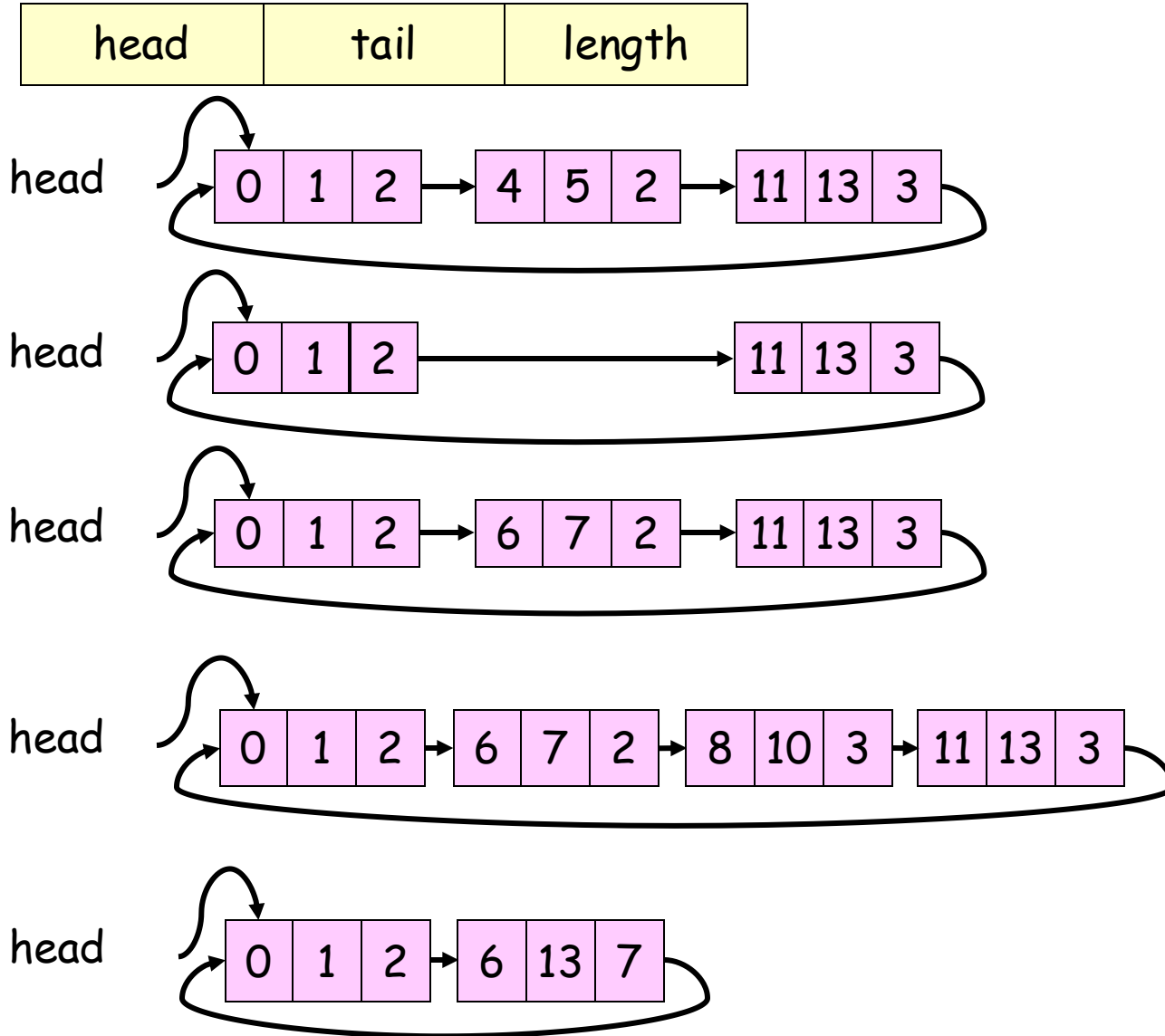


Dynamically Partition

- Multiple-partition allocation
 - **Hole** - block of available memory; holes of various size are scattered throughout memory
 - When a process arrives, it is allocated memory from a hole large enough to accommodate it
 - Operating system maintains information about:
 - a) allocated partitions; b) free partitions (hole)



Dynamic Storage-Allocation Problem



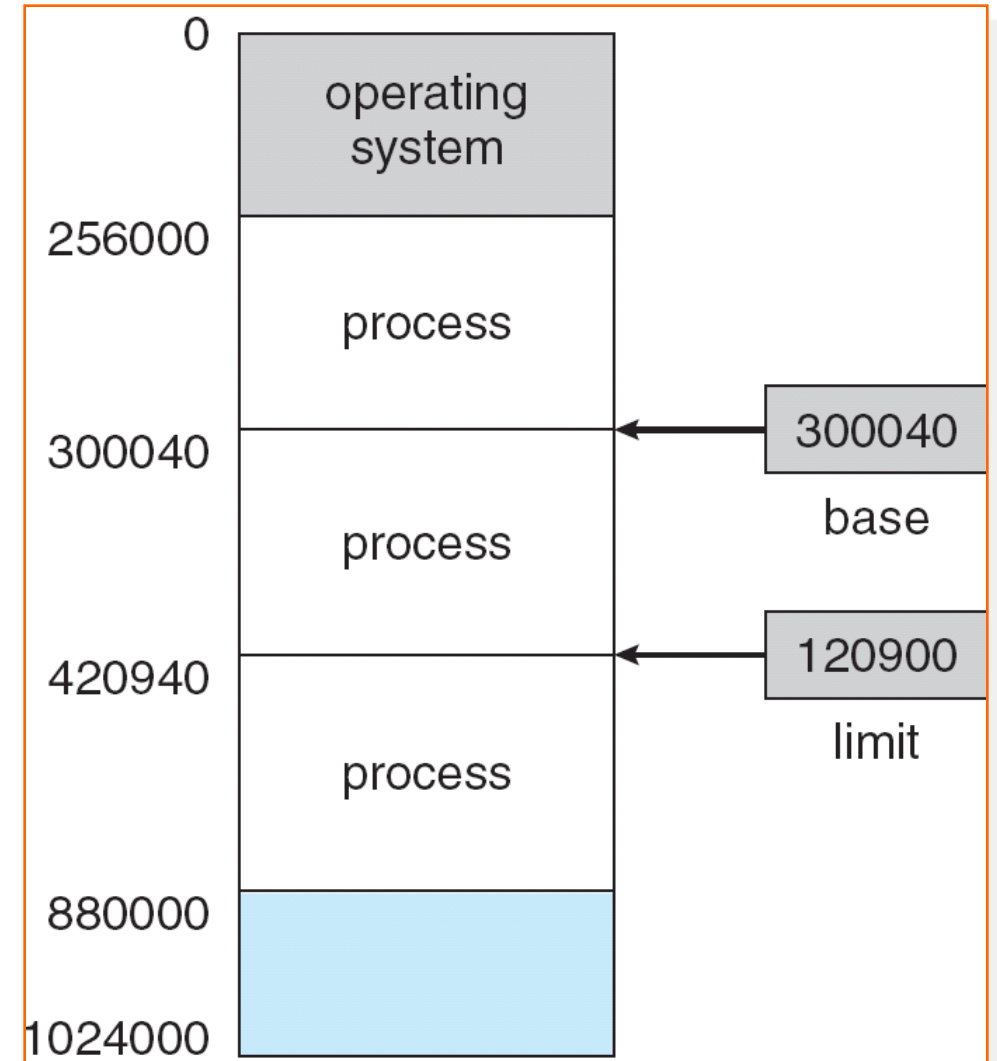
0	
1	
2	Process 1
3	
4	Process 5
5	
6	Process 2
7	
8	Process 3
9	
10	
11	
12	
13	
14	Process 4
15	

Contiguous Memory Allocation

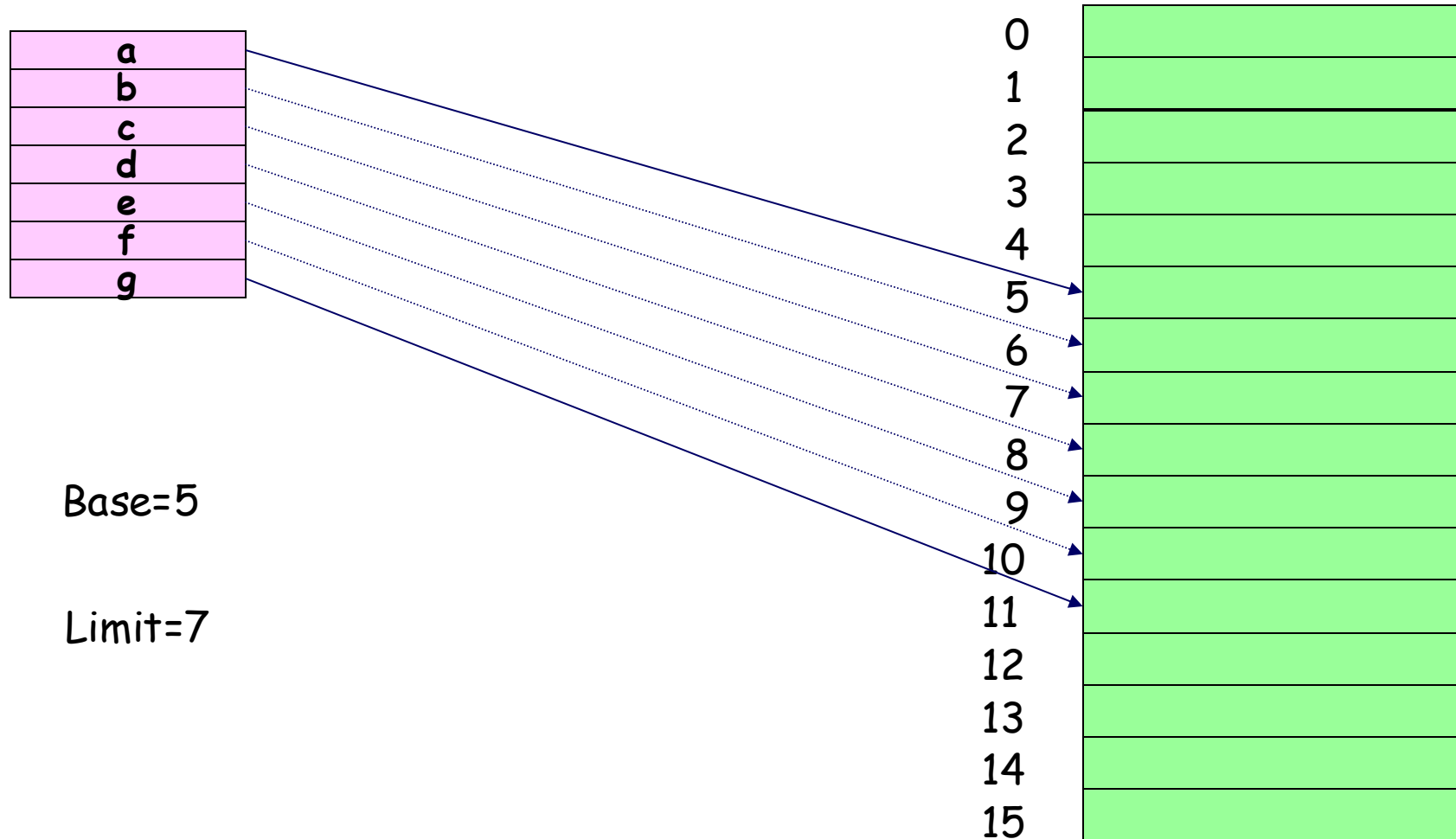
- Relocation registers used to protect user processes from each other, and from changing operating-system code and data
 - **Base** register contains value of **smallest** physical address
 - **Limit** register contains **range** of logical addresses - each logical address must be **less** than the limit register
 - MMU maps logical address **dynamically**

Base and Limit Registers

- A pair of **base** and **limit** registers define the logical address space

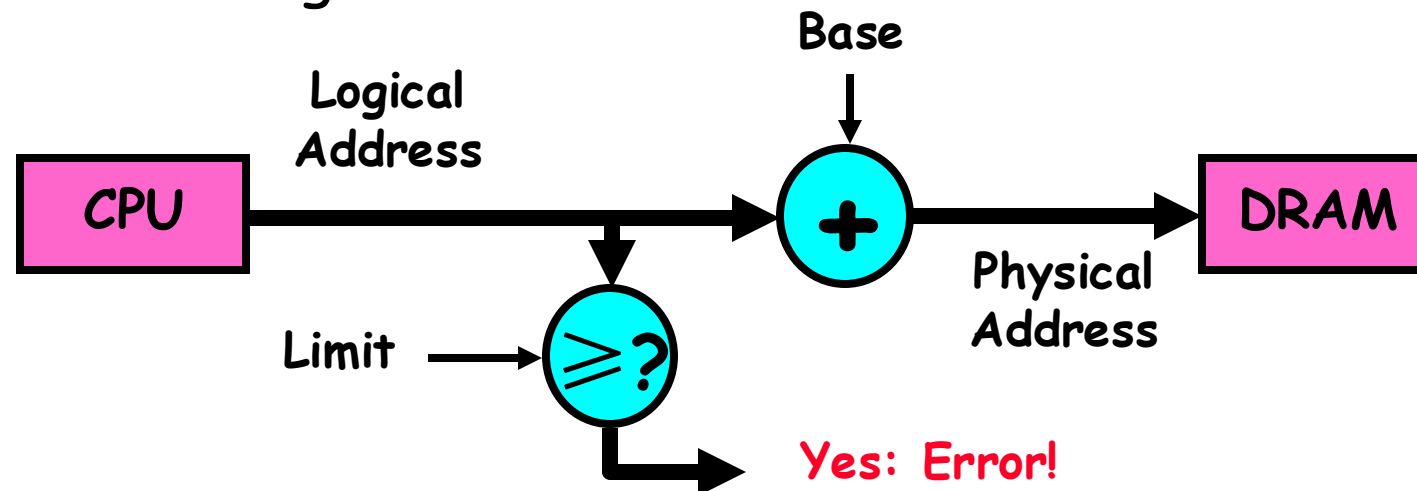


Base and Limit Registers



Base and Limit Registers

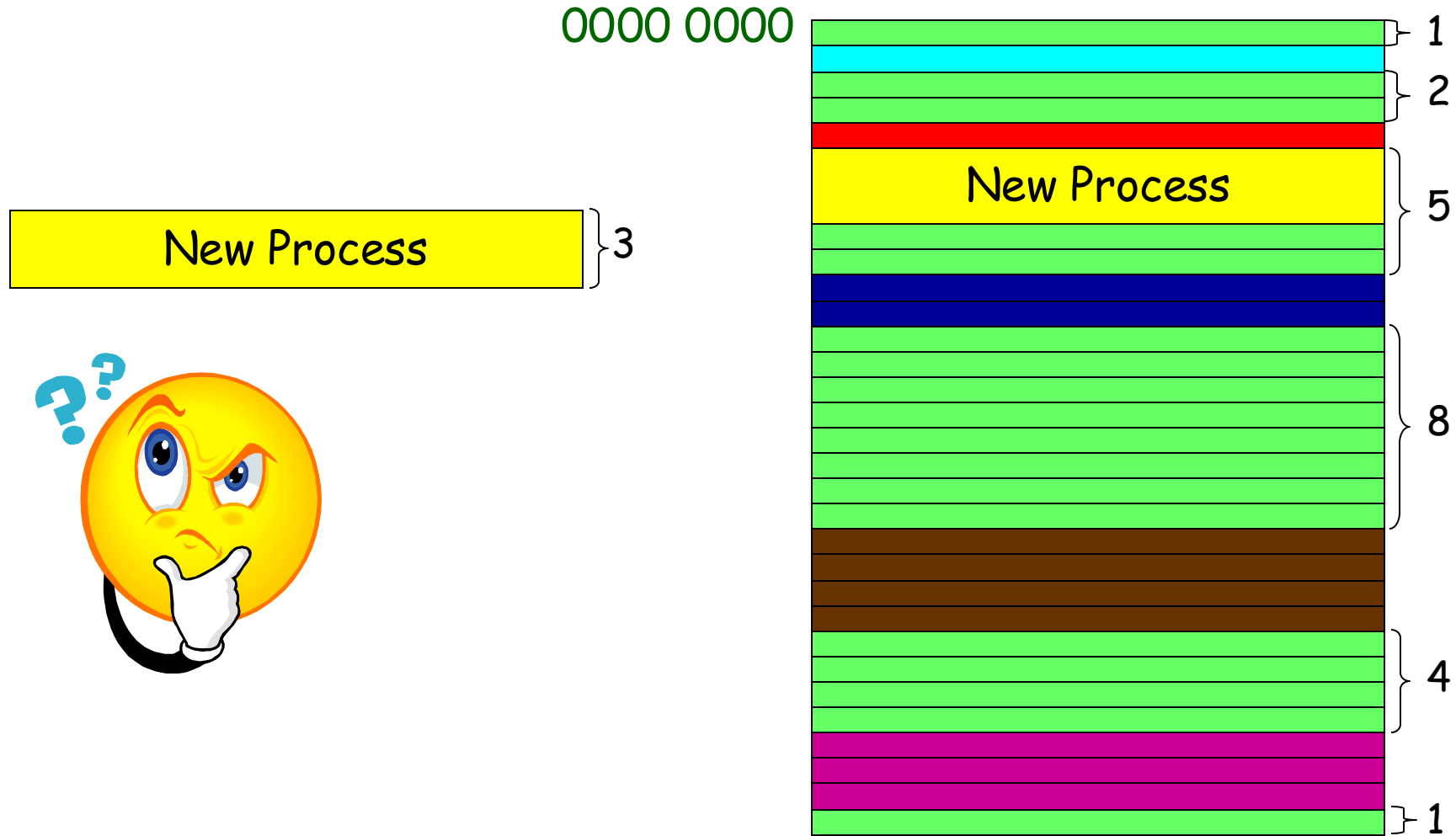
- Can use base & bounds/limit for dynamic address translation (Simple form of "segmentation"):
- This gives program the illusion that it is running on its own dedicated machine, with memory starting at 0
 - Program gets continuous region of memory
 - Addresses within program do not have to be relocated when program placed in different region of DRAM



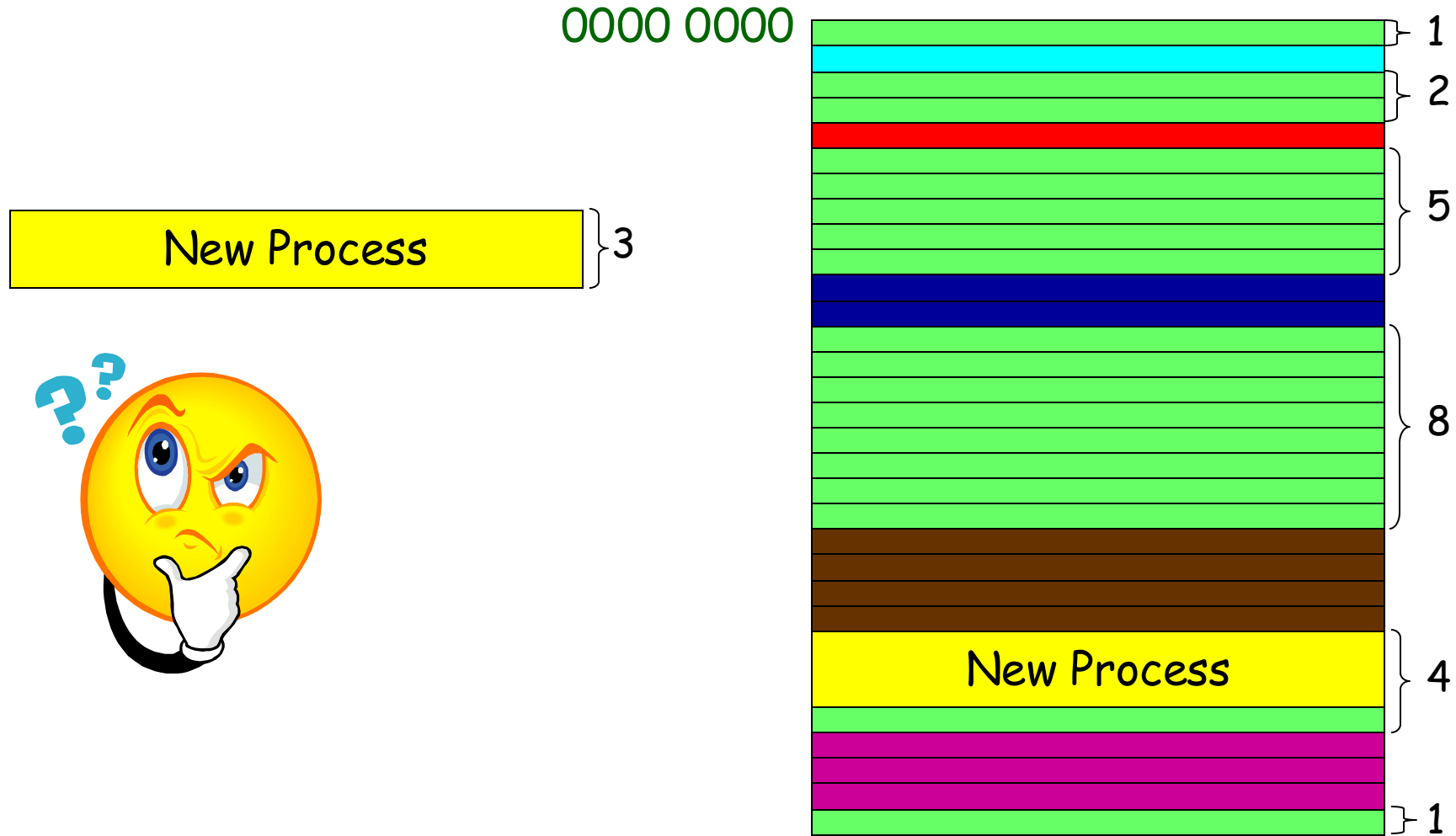
Dynamic Storage-Allocation Problem

- How to satisfy a request of size n from a list of free holes
 - **First-fit**: Allocate the **first** hole that is big enough
 - **Best-fit**: Allocate the **smallest** hole that is big enough; must search entire list, unless ordered by size
 - Produces the smallest leftover hole
 - **Worst-fit**: Allocate the largest hole; must also search entire list
 - Produces the largest leftover hole
 - **Next-fit**: Scans memory from the location of the last placement
- First-fit and best-fit better than worst-fit in terms of speed and storage utilization

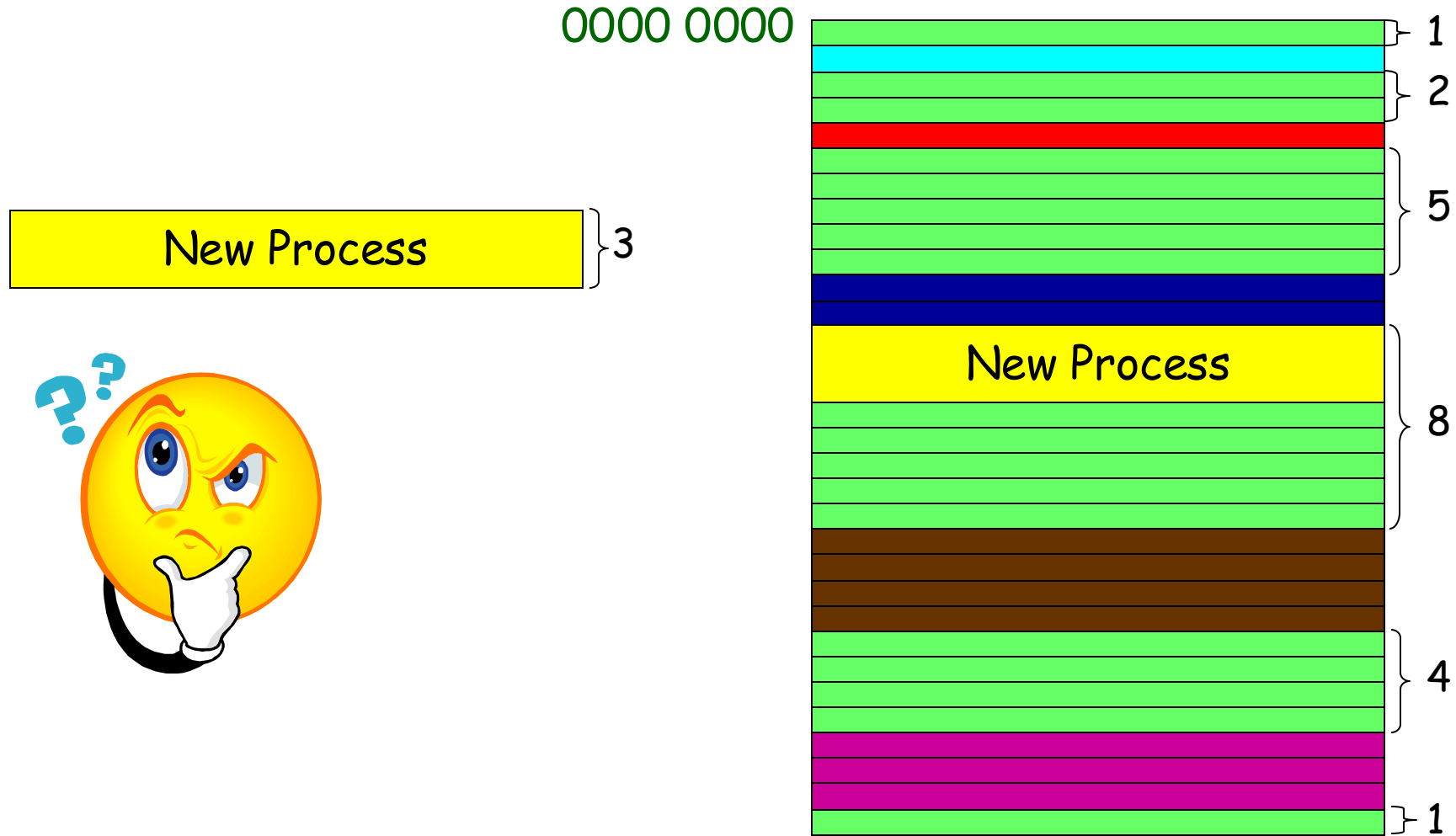
First-fit



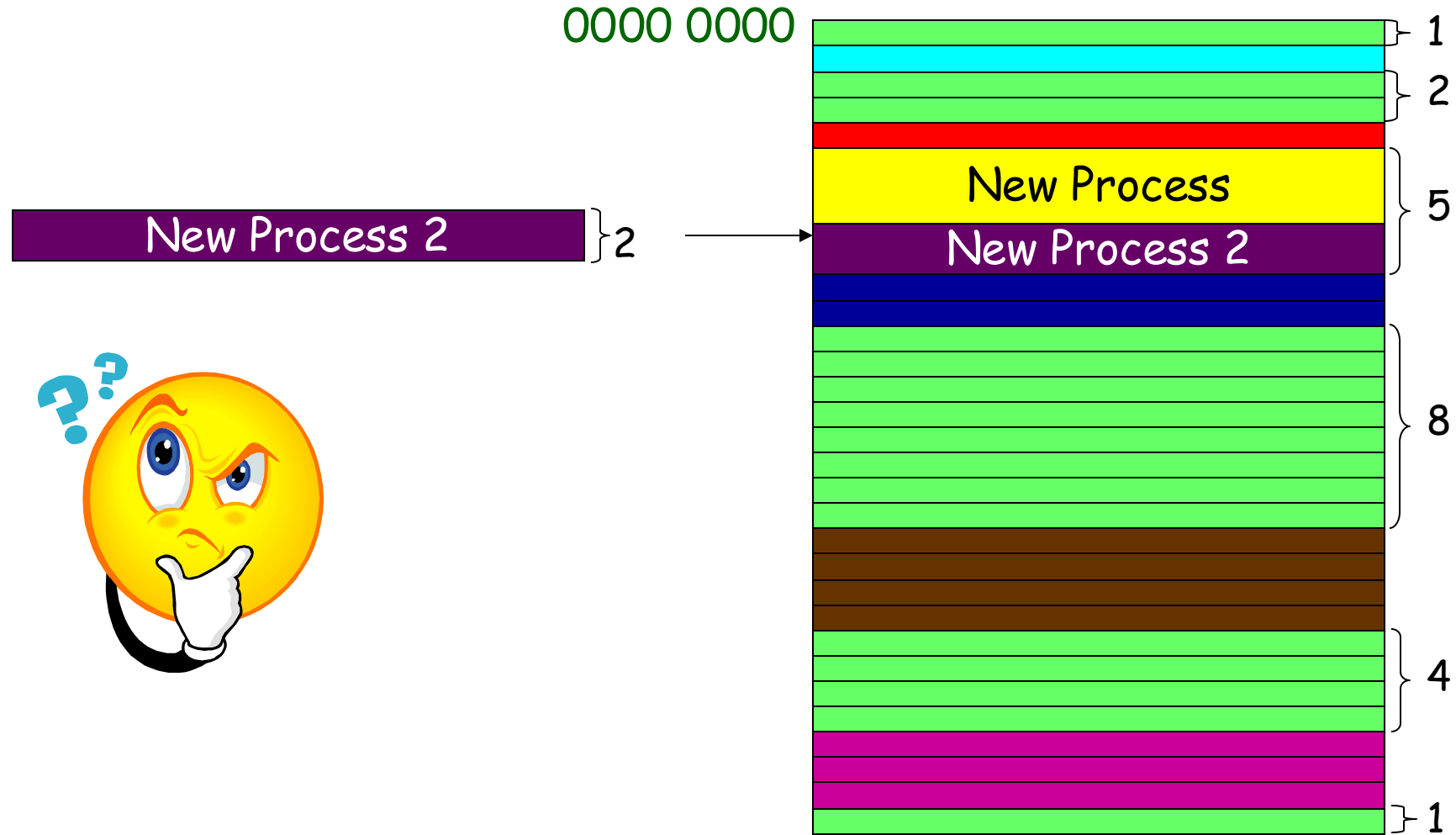
Best-fit



Worst-fit



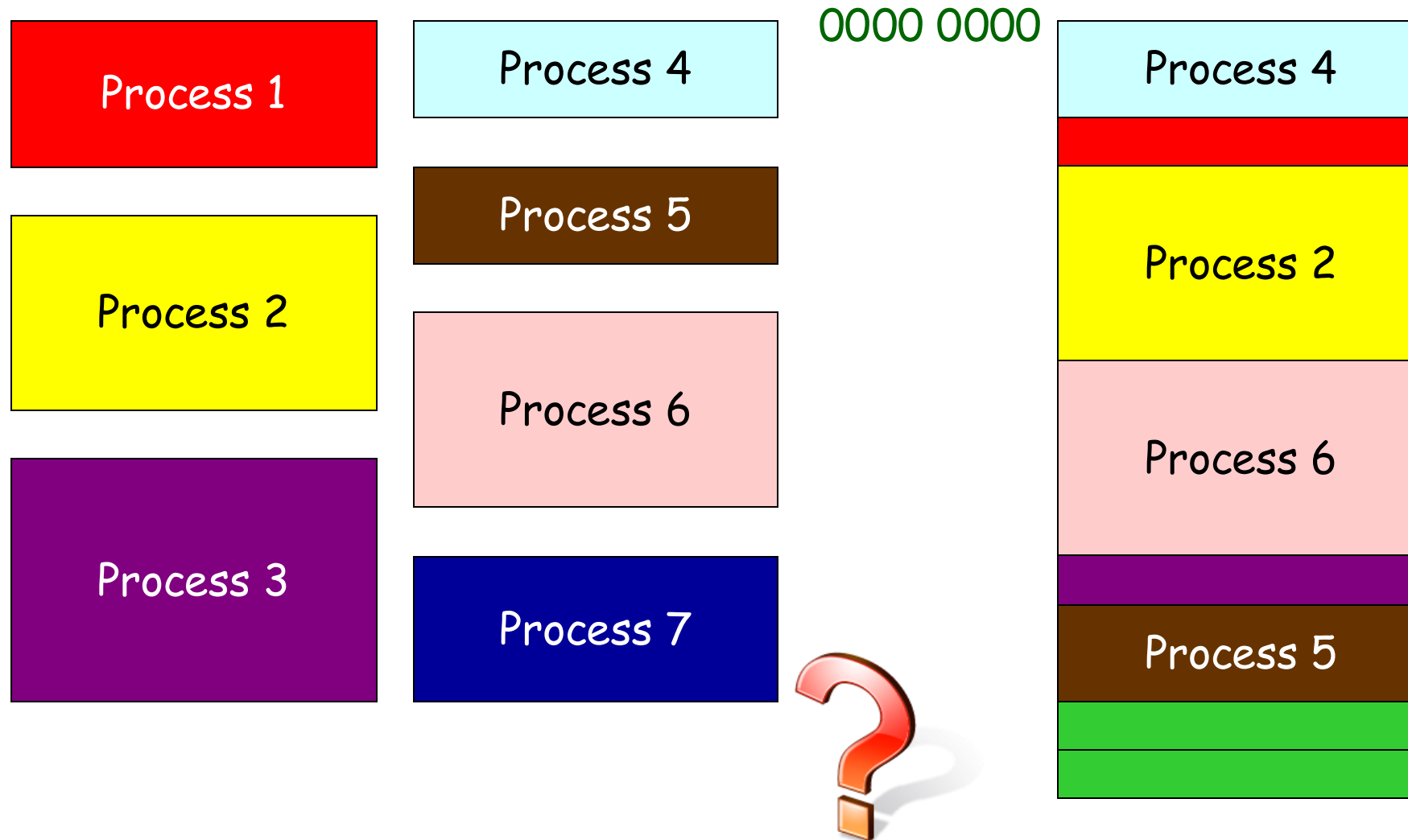
Next-fit



Exercise

- Given memory partitions, which are 100K, 450K, 250K, 300K and 600K
- There are four processes in order: 212K, 417K, 112K, and 426K.
- In order to place processes in the memory, there are three algorithms: the first-fit algorithm, the best-fit algorithm, and the next-fit algorithm. Of the three algorithms, which one makes the best use of memory space?

Fragmentation



Fragmentation

- **External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous
 - External fragmentation is the phenomenon in which free storage becomes divided into many small pieces over time.
 - The term "**external**" refers to the fact that the unusable storage is outside the allocated regions.

Fragmentation

- Consider a multiple-partition allocation scheme with a hole of 18,464 bytes. Suppose that the next process requested 18,462 bytes. If we allocate exactly the requested block, we are left with a hole of 2 bytes.
- The overhead to keep track of this hole will be substantially larger than the hole itself.

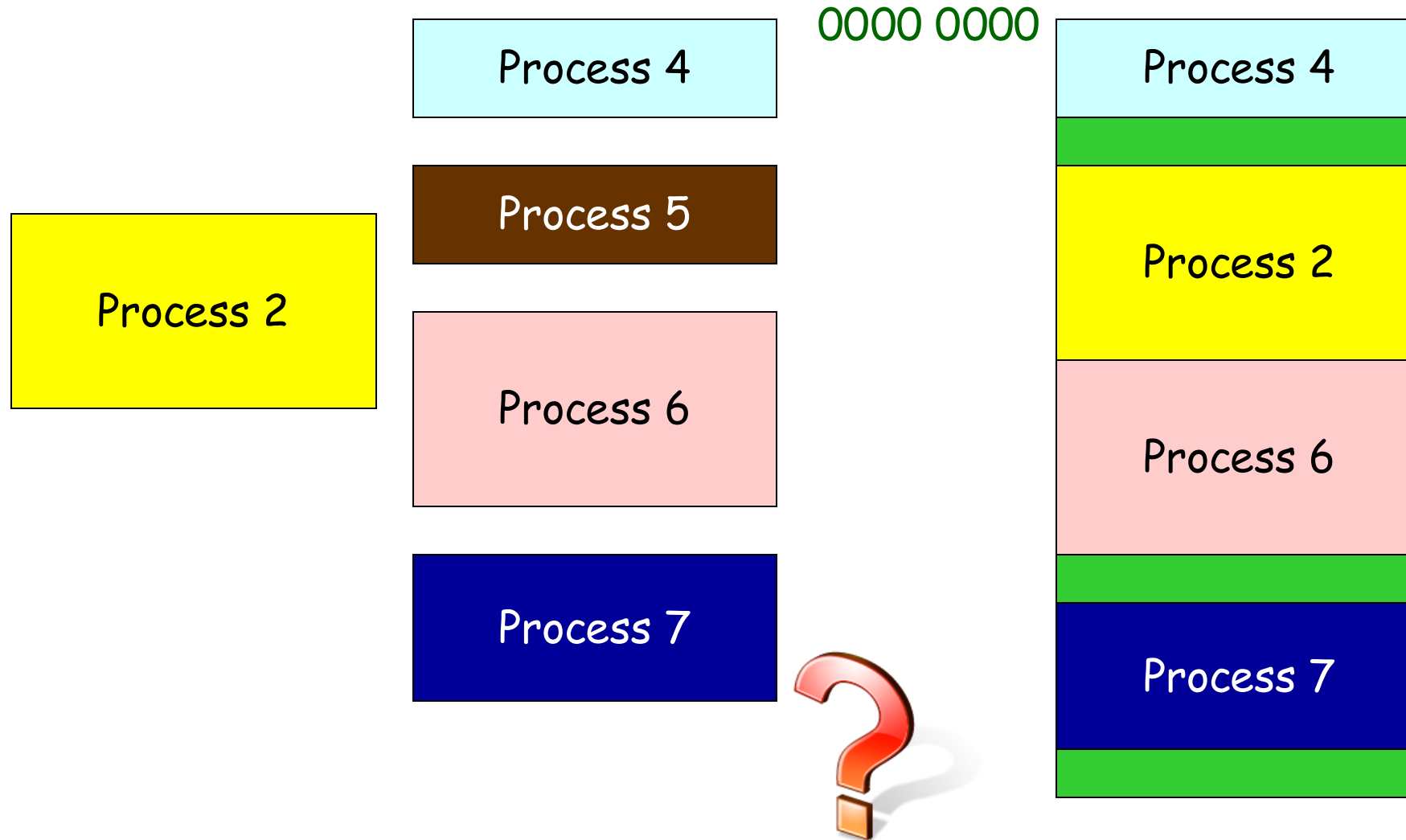
Fragmentation

- The general approach to avoiding this problem is to **break the physical memory into fixed-sized blocks** and allocate memory in units based on **block size**.
- With this approach, the memory allocated to a process may be slightly larger than the requested memory.
- The difference between these two approaches is **internal fragmentation** - memory that is internal to a partition but is not being used.

Fragmentation

- Reduce external fragmentation by **compaction**
 - Shuffle memory contents to place all free memory together in one large block
 - However, it is somewhat expensive.
 - Compaction is possible only if relocation is dynamic, and is done at execution time

Compaction



Contiguous Memory Allocation

- Base and Limit Pros: Simple, fast
- Fragmentation problem
 - Not every process is in the same size
 - Over time, memory space becomes fragmented
- Missing support for sparse address space
 - Would like to have multiple chunks/program
 - E.g.: Code, Data, Stack
- Hard to do inter-process sharing
 - Want to share code segments when possible
 - Want to share memory between processes

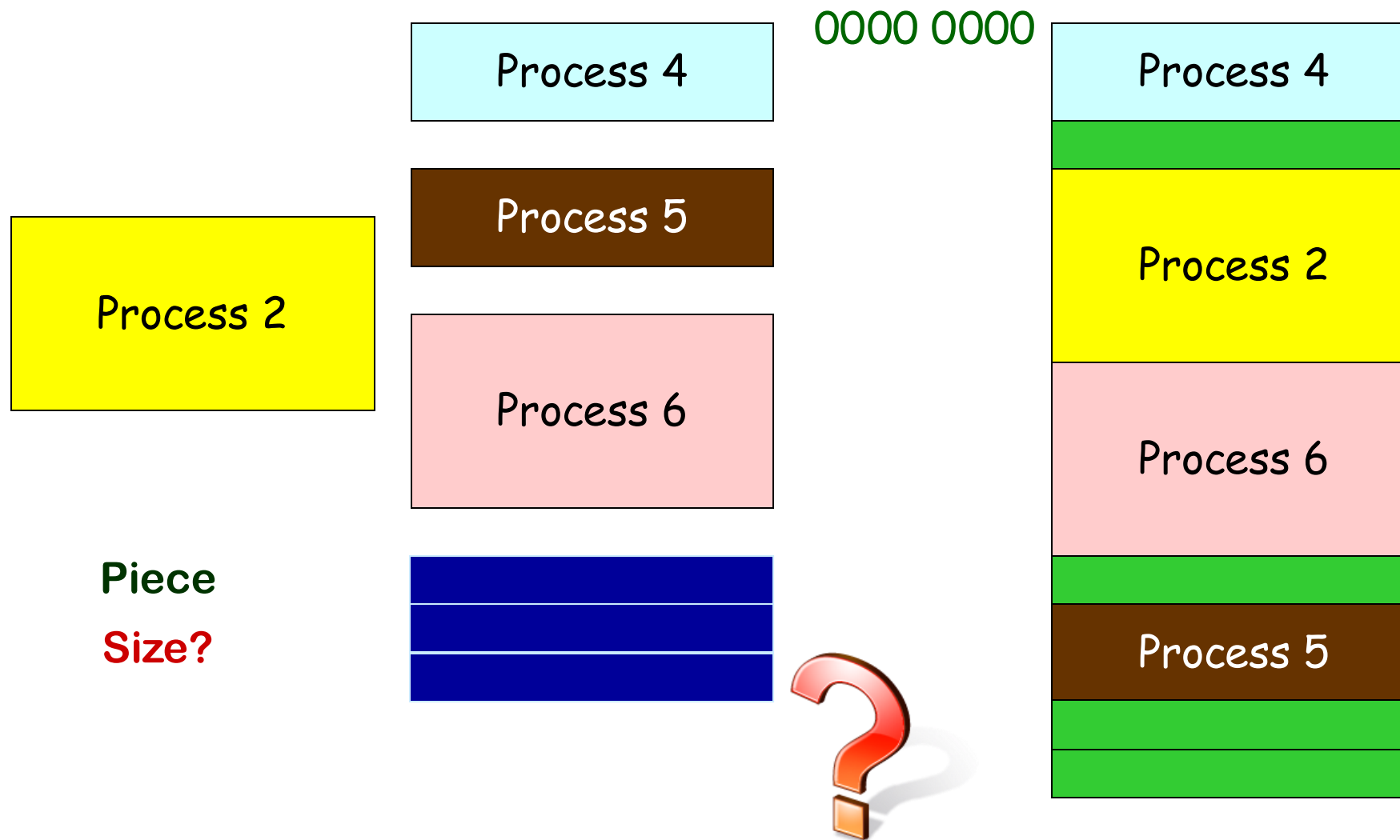
Contiguous Memory Allocation

- Is there a scheme **without** External Fragmentation?
 - Of course!



Paging

Paging



Paging

- Physical address space of a process can be **noncontiguous**; process is allocated physical memory whenever the latter is available
- Divide physical memory into fixed-sized blocks called **frames**
 - Size is power of 2
 - Typically have small pages (1K-16K)
- Divide logical memory into blocks of same size called **pages**

How to allocate?

- 32-byte memory and 4-byte pages

Page0	a
	b
	c
	d
Page1	e
	f
	g
	h
Page2	i
	j
	k
	l
Page3	m
	n
	o
	p

Logical Memory

i	
j	
k	
l	
e	
f	
g	
h	
a	
b	
c	
d	
m	
n	
o	
p	

Physical Memory

How to allocate?

- 32-byte memory and 4-byte pages

Page0	a
	b
	c
	d
Page1	e
	f
	g
	h
Page2	i
	j
	k
	l
Page3	m
	n
	o
	p

Logical Memory

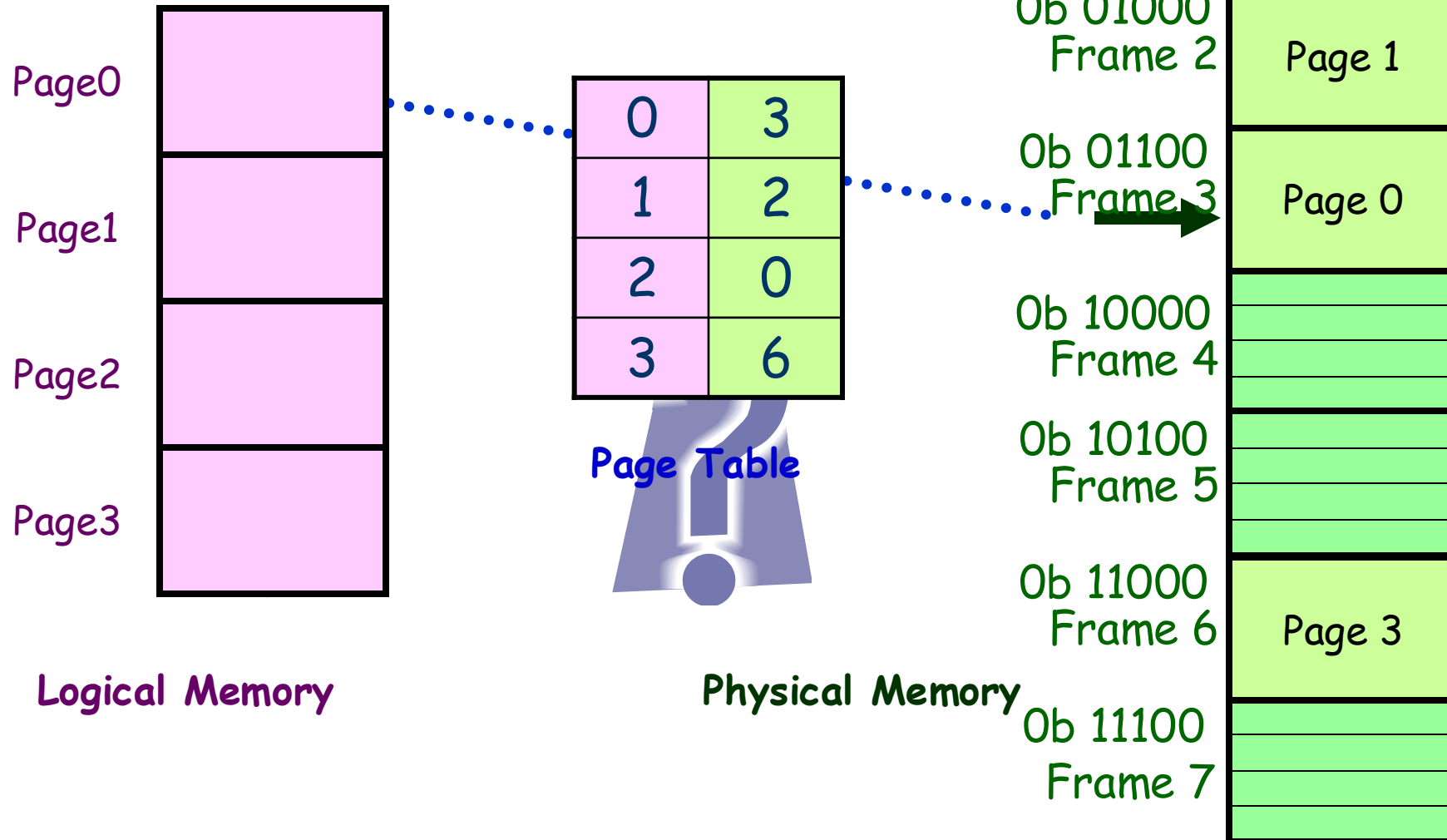


Physical Memory

0b 00000	i
	j
	k
	l
0b 00100	
0b 01000	e
	f
	g
	h
0b 01100	a
	b
	c
	d
0b 10000	
0b 10100	
0b 11000	m
	n
	o
	p
0b 11100	

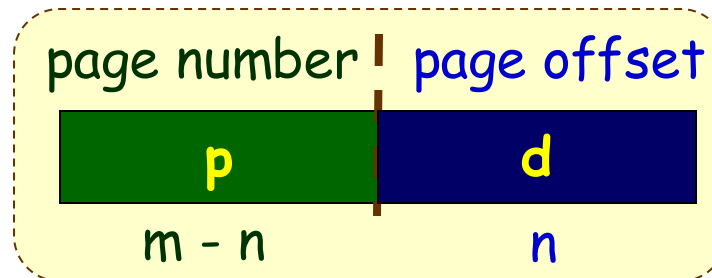
How to allocate?

- 32-byte memory and 4-byte pages



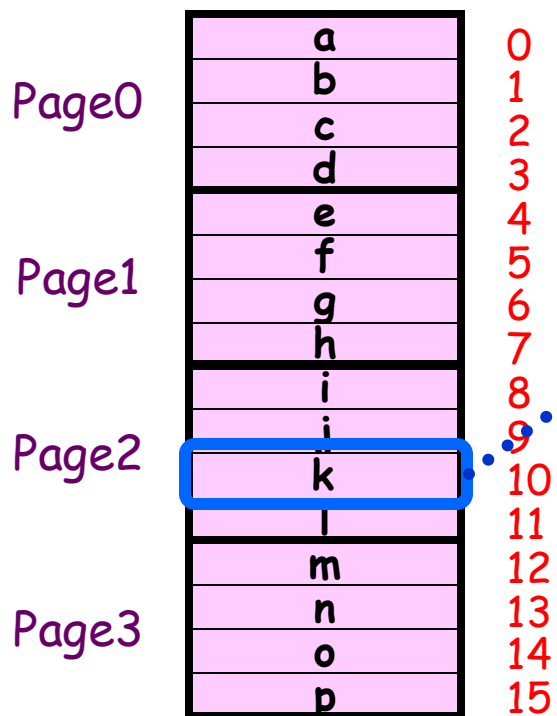
Address Translation Scheme

- Address generated by CPU is divided into:
 - **Page number (p)** – used as an index into a page table which contains base address of each page in physical memory
 - **Page offset (d)** – combined with base address to define the physical memory address that is sent to the memory unit
- For given logical address space 2^m and page size 2^n



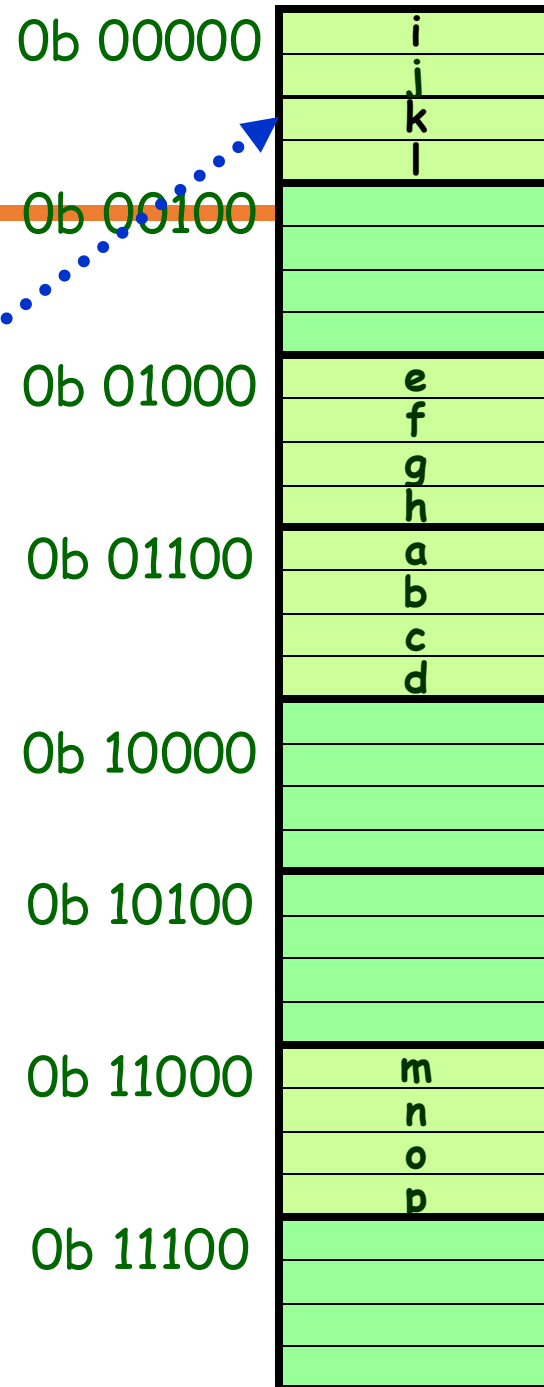
How to allocate?

- 32-byte memory and 4-byte pages

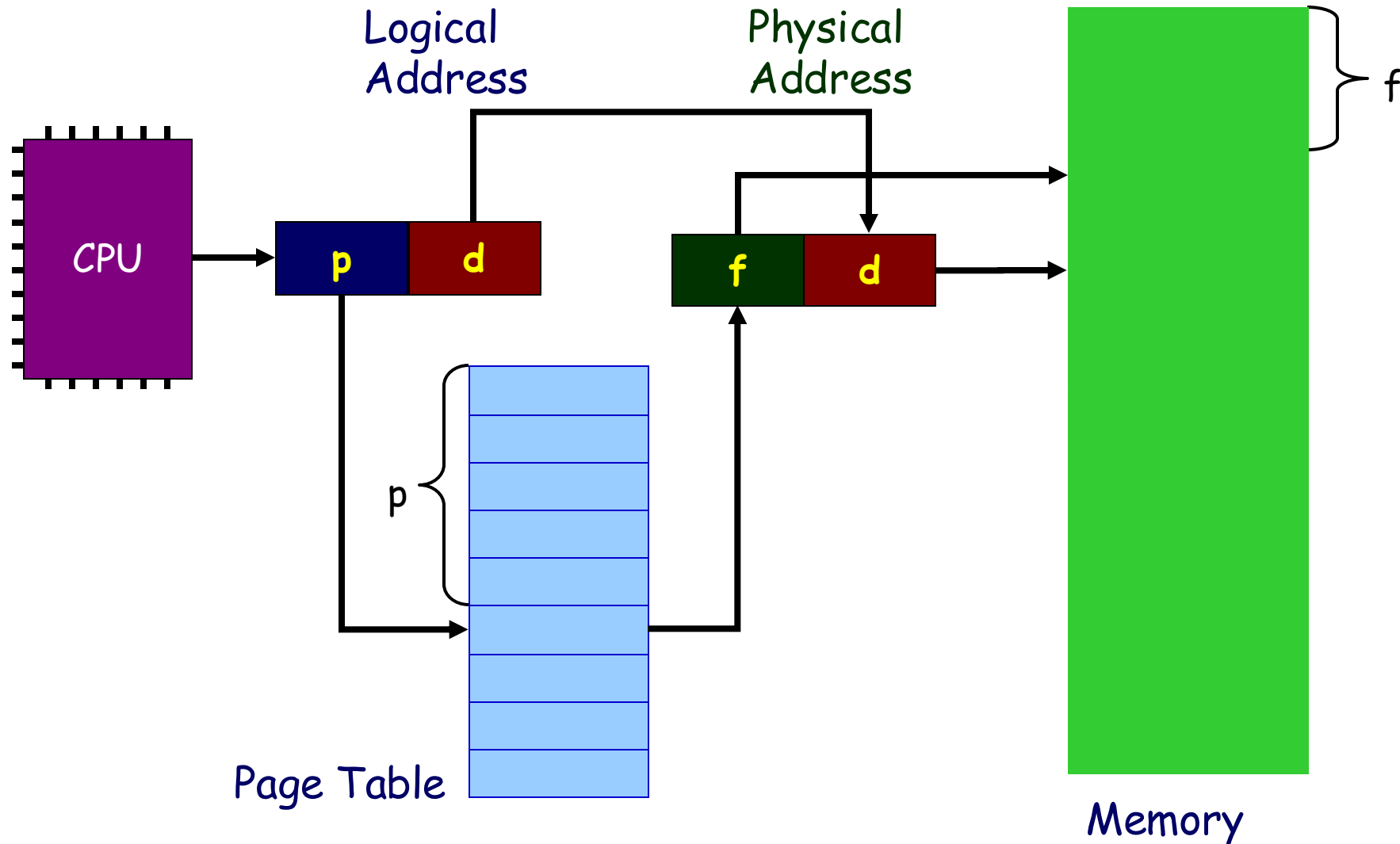


0	3
1	2
2	0
3	6

Physical Memory



Paging Hardware



Paging

- Keep track of all free frames
- To run a program of size n pages, need to find n free frames and load program
- Set up a page table to translate logical to physical addresses
- No external fragmentation
- May has internal fragmentation

Free Frames



Logical
Memory

0	3
1	2
2	0
3	6

Page
Table



Frame 0

Page 2

Frame 1

Frame 2

Page 1

Frame 3

Page 0

Frame 4

Frame 5

Frame 6

Page 3

Frame 7

Physical
Memory

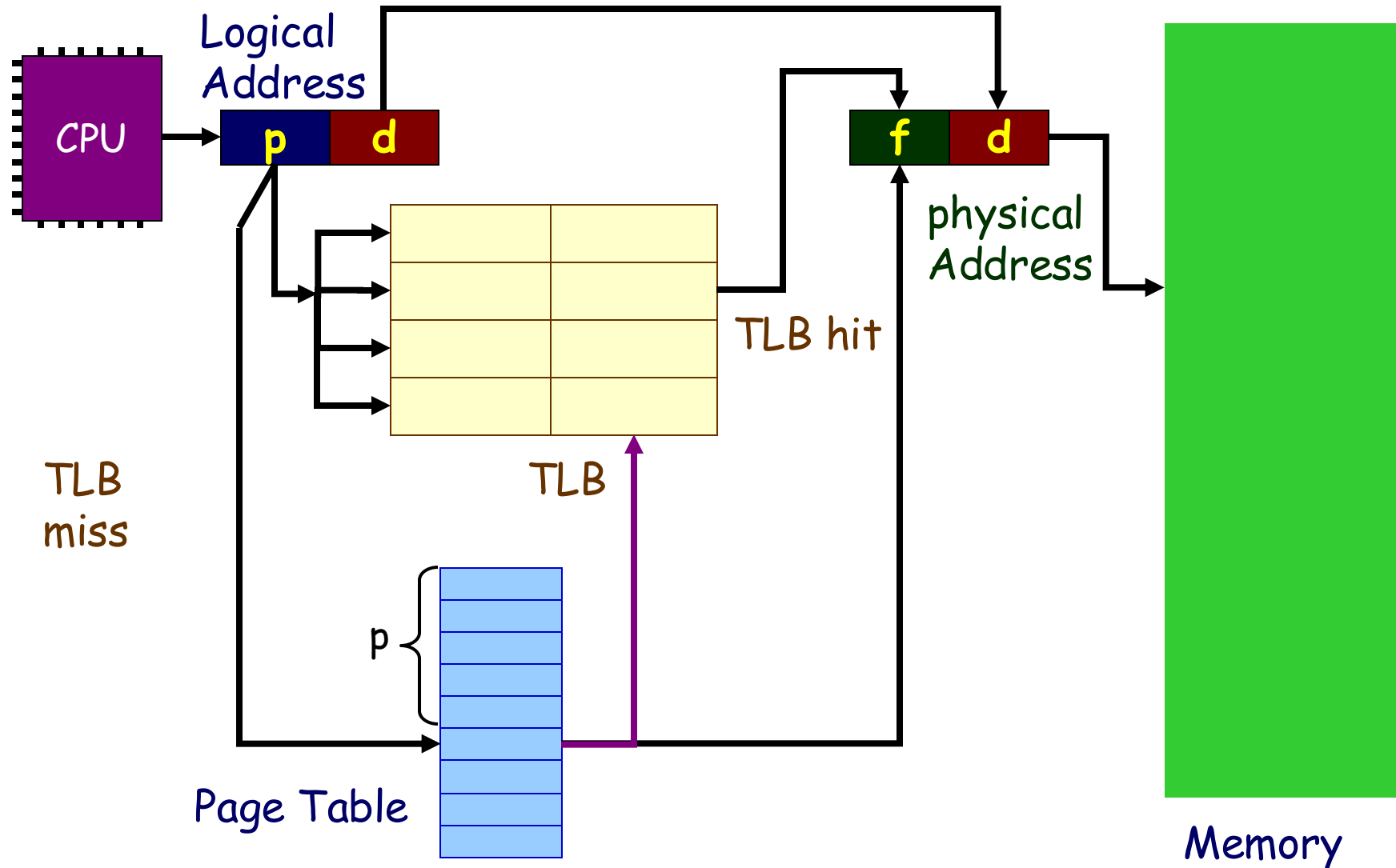
Implementation of Page Table

- The page table is **contiguous** in main memory
- **Page-table base register (PTBR)** points to the page table
- **Page-table length register (PRLR)** indicates size of the page table

Implementation of Page Table

- In this scheme every data/instruction access requires **two** memory accesses.
 - One for the page table, and
 - One for the data/instruction.
- The **two memory access problem** can be solved by the use of a special fast-lookup hardware cache called **associative memory or translation look-aside buffers (TLBs)**

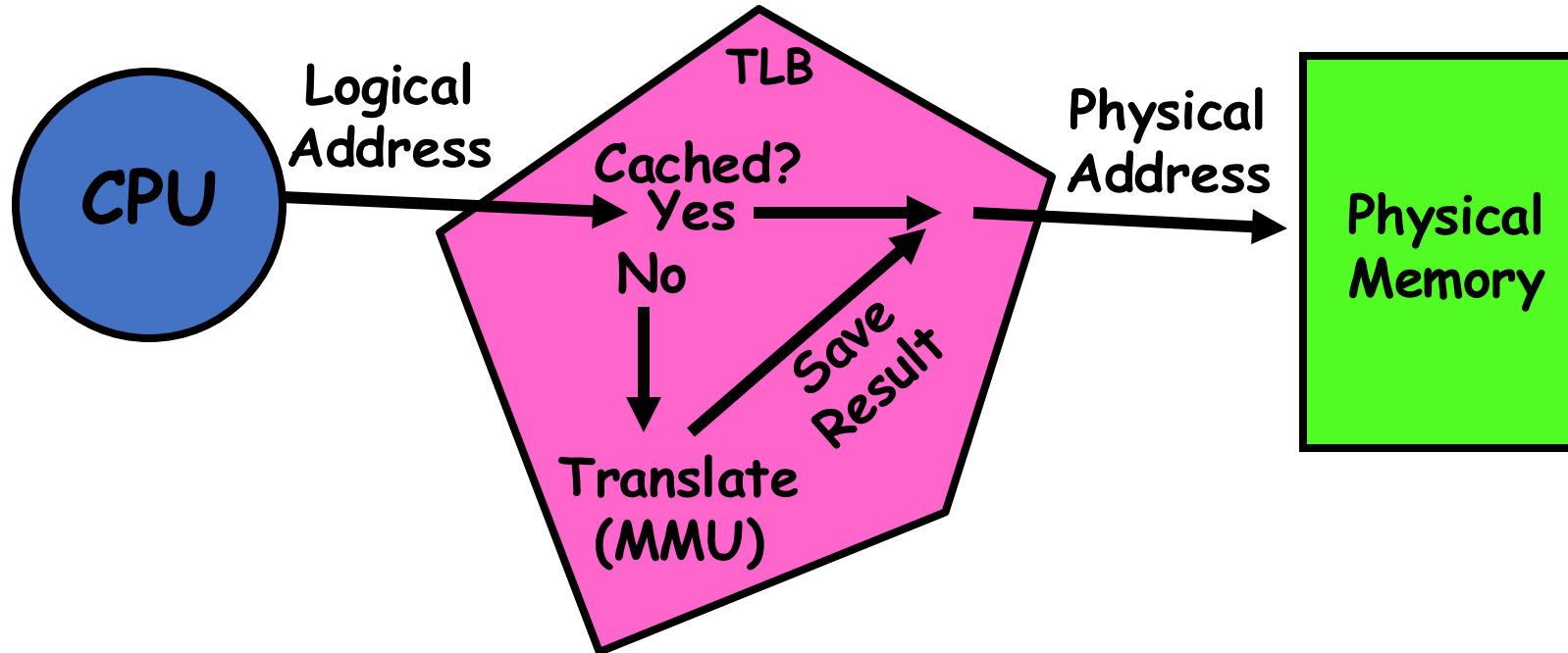
Paging Hardware With TLB



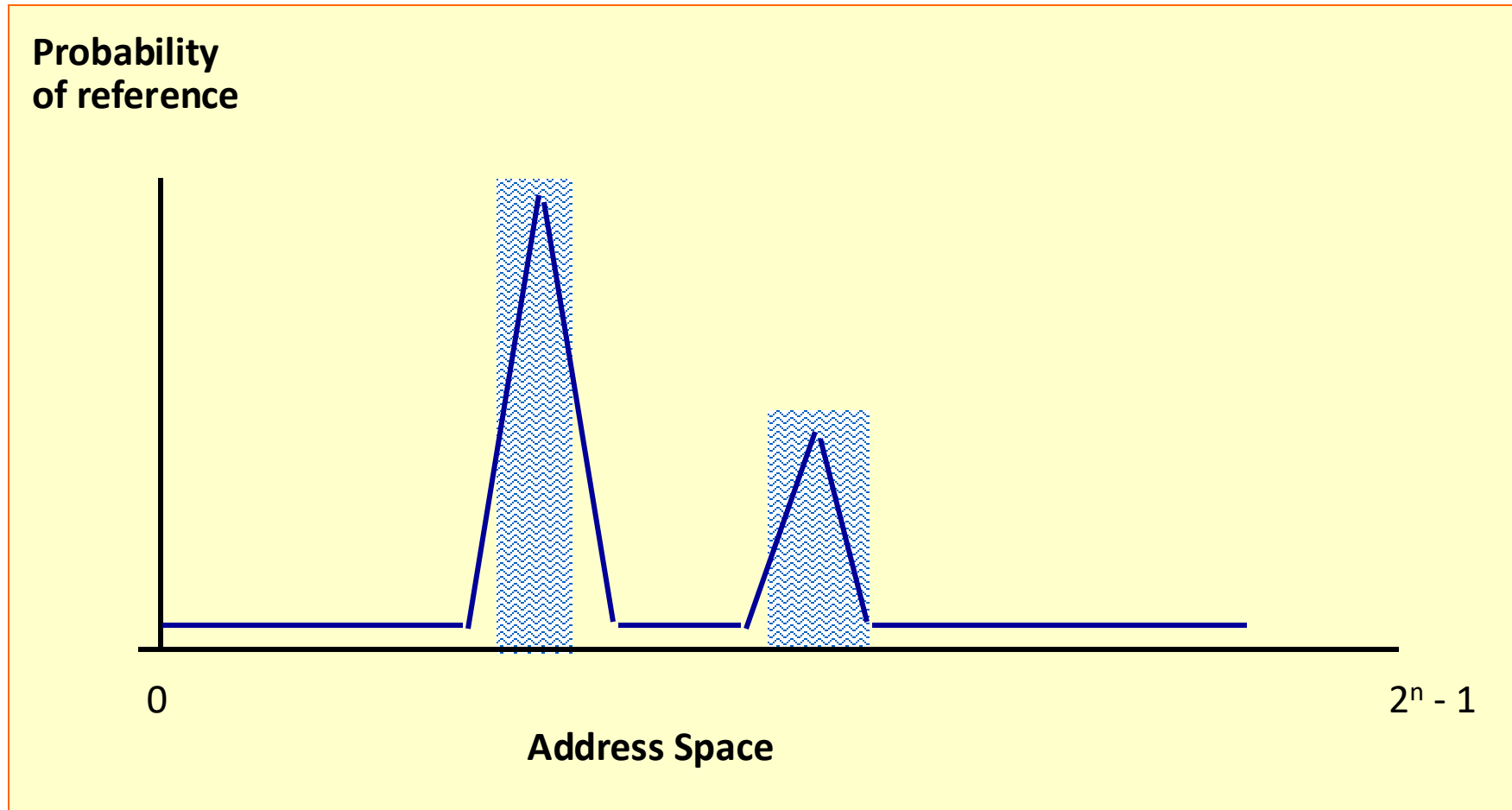
Paging Hardware With TLB

- If the page number is not in the TLB (known as a TLB miss), a memory reference to the page table must be made. When the frame number is obtained, we can use it to access memory.
- In addition, we add the page number and frame number to the TLB, so that they will be found quickly on the next reference.
- If the TLB is already full of entries, the operating system must select one for replacement.

Paging Hardware With TLB



Why Does TLB Help? Locality!



Effective Access Time

- Associative Lookup = ε time unit
- Assume memory cycle time is 1 time unit
- Hit ratio – percentage of times that a page number is found in the associative registers; ratio related to number of associative registers

- Hit ratio = α

- **Effective Access Time (EAT)**

$$\begin{aligned} \text{EAT} &= (1 + \varepsilon) \alpha + (2 + \varepsilon)(1 - \alpha) \\ &= 2 + \varepsilon - \alpha \end{aligned}$$

Effective Access Time

$$\varepsilon = 0.1$$

$$EAT = 2 + \varepsilon - \alpha$$

α	EAT
5%	2.05
10%	2
50%	1.6
80%	1.3
?	1.12

Paging Hardware With TLB

- Needs to be really fast
 - Critical path of memory access
- Needs to have very few conflicts!
 - With TLB, the Miss Time extremely high!
- How big does TLB actually have to be?
 - Usually small: 128-512 entries

Memory Protection

- User cannot modify the page table mapping
- Memory protection implemented by associating protection bit with each frame
- **Valid-invalid bit** attached to each entry in the page table:
 - “**valid**” indicates that the associated page is in the process’ logical address space, and is thus a legal page
 - “**invalid**” indicates that the page is not in the process’ logical address space

Paging

- Valid (v) or Invalid (i) Bit In A Page Table

Page 0
Page 1
Page 2
Page 3
Page 4
Page 5
Page 6

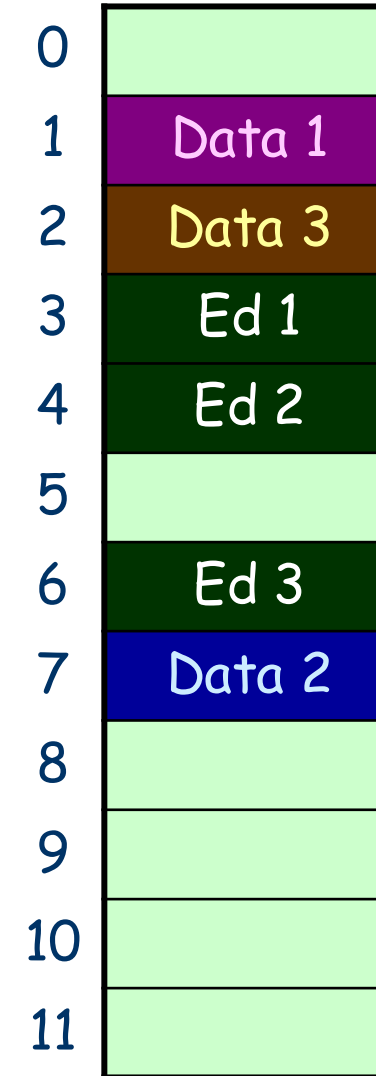
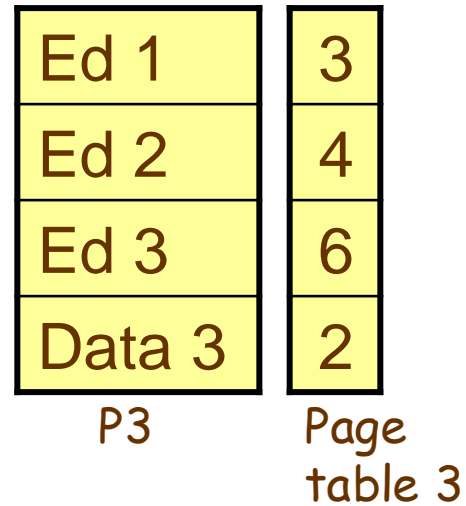
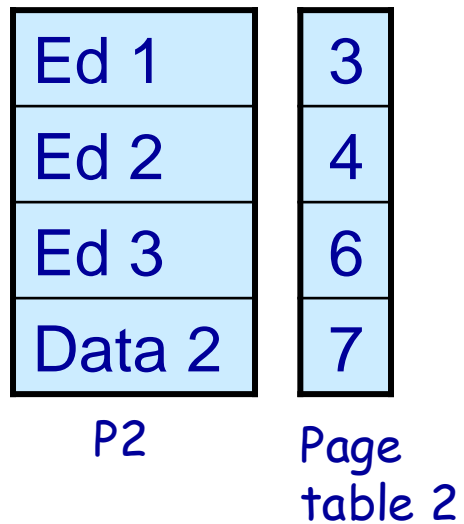
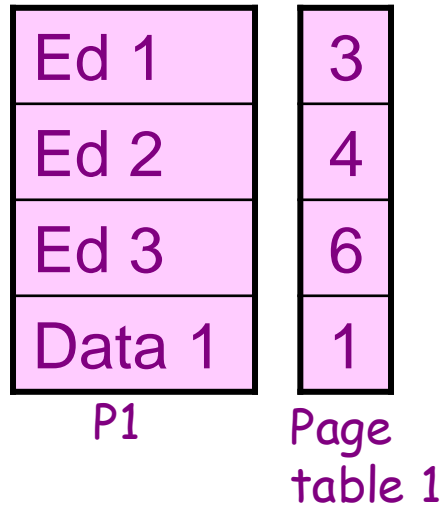
Frame number		Valid-invalid bit
0	2	v
1	3	v
2	4	v
3	7	v
4	8	v
5	9	v
6	0	i
7	0	i

0	
1	
2	Page 0
3	Page 1
4	Page 2
5	
6	
7	Page 3
8	Page 4
9	Page 5
	...
	...

Shared Pages

- Shared code
 - One copy of read-only code shared among processes (i.e., text editors).
- Private code and data
 - Each process keeps a separate copy of the code and data
 - The pages for the private code and data can appear anywhere in the logical address space

Shared Pages Example





Structure of the Page Table

Structure of the Page Table

- A logical address (on 32-bit machine with 4K page size) is divided into:
 - a page number consisting of 20 bits
 - a page offset consisting of 12 bits
- How many pages?
 - $2^{32}/2^{12}=2^{20}$
- What is the size of the page table, if each item of the page table takes 4B?
 - 4MB

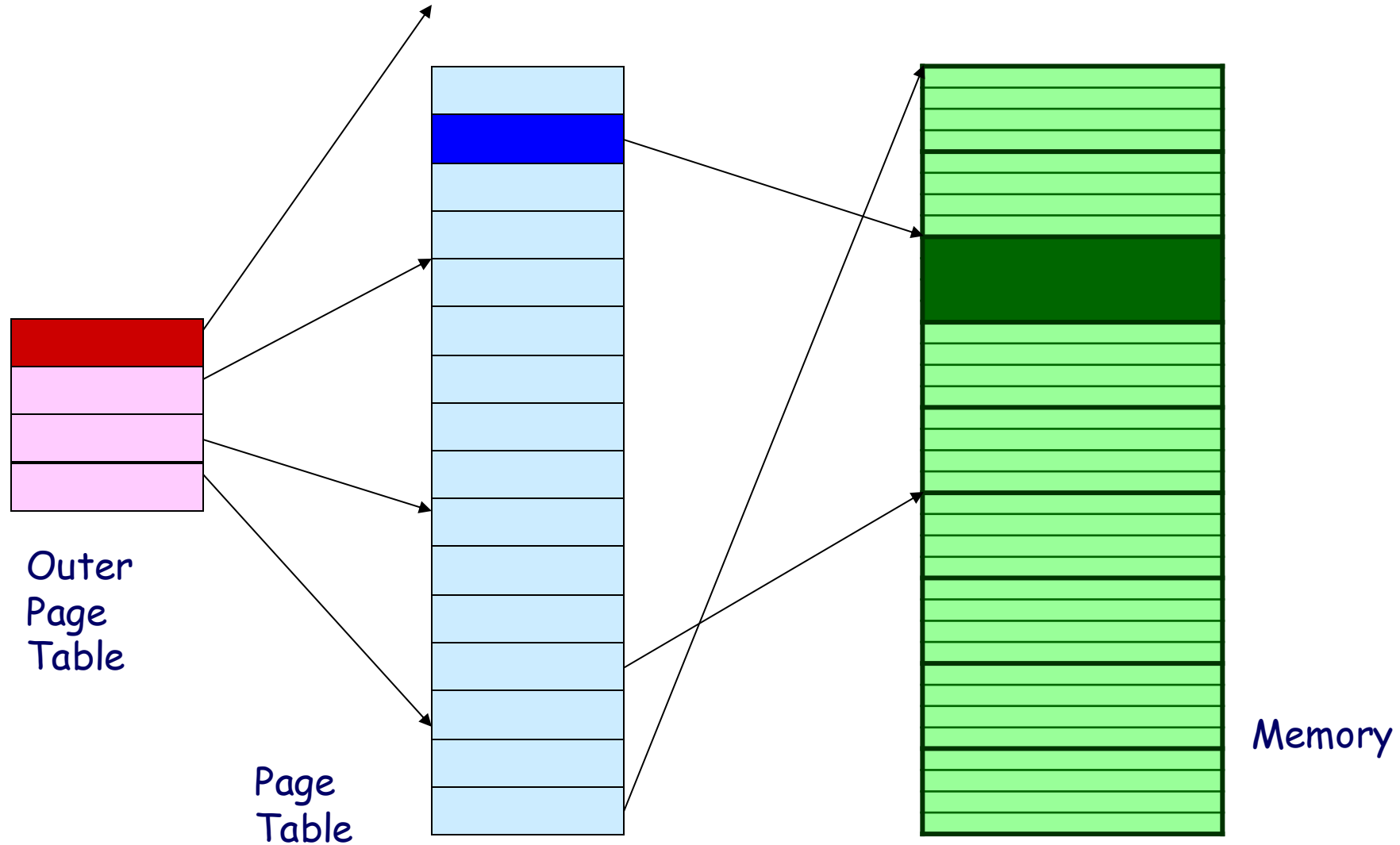
Structure of the Page Table

- Hierarchical Page Tables
- Hashed Page Tables
- Inverted Page Tables

Hierarchical Page Tables

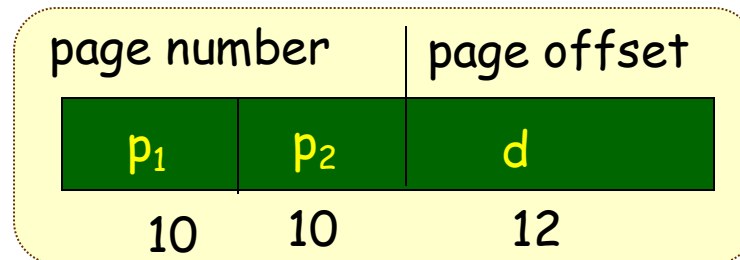
- Break up the logical address space into multiple page tables
- A simple technique is a two-level page table

Two-Level Page-Table Scheme



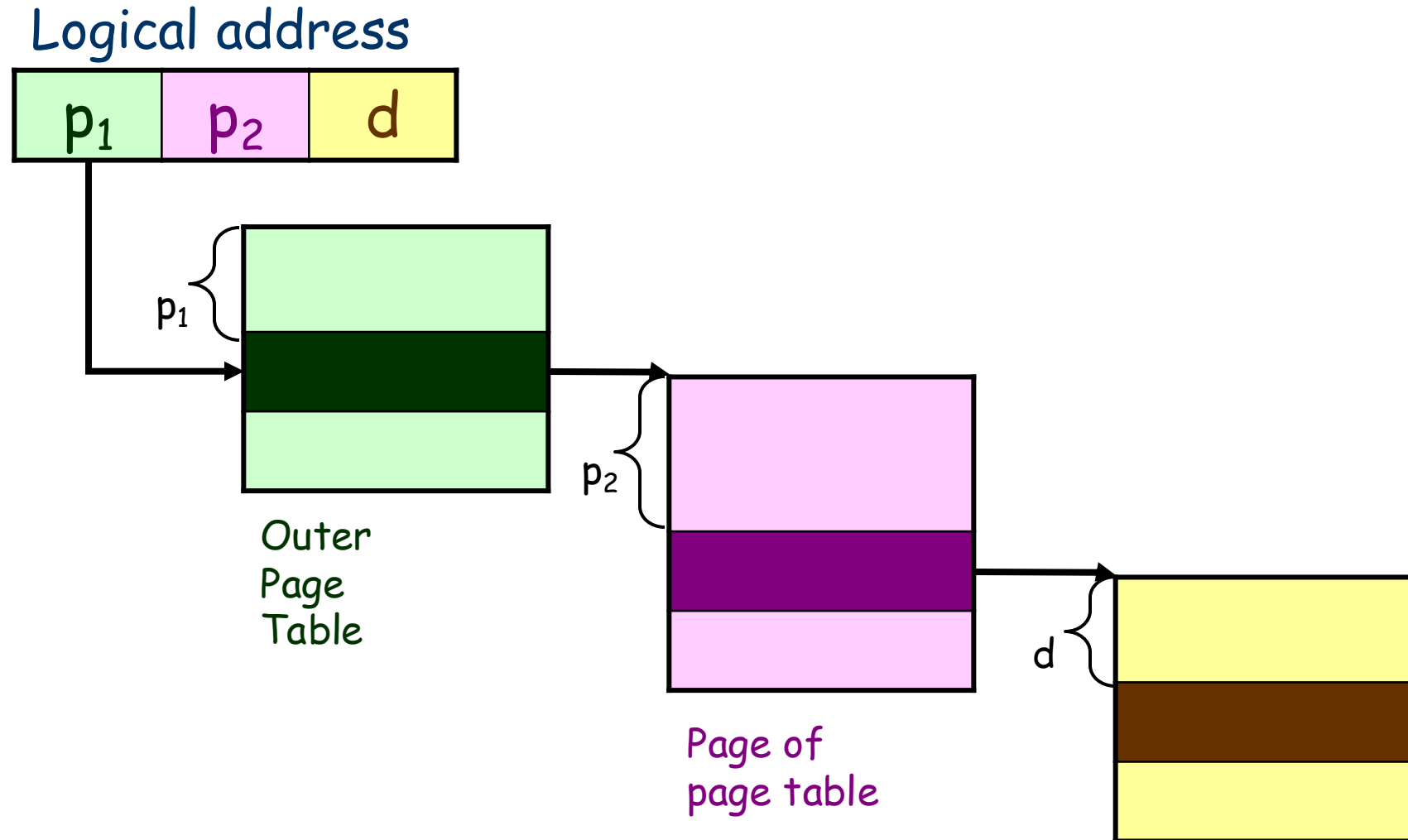
Two-Level Paging Example

- A logical address (on 32-bit machine with 4K page) is divided into:
 - a page number consisting of 20 bits
 - a page offset consisting of 12 bits
- Since the page table is paged, the page number is further divided into:
 - a 10-bit page number
 - a 10-bit page offset
- Thus, a logical address is as follows:



where p_1 is an index into the outer page table, and p_2 is the displacement within the page of the outer page table

Address-Translation Scheme



Three-level Paging Scheme

outer page	inner page	offset
p_1	p_2	d
42	10	12

2nd outer page	outer page	inner page	offset
p_1	p_2	p_3	d
32	10	10	12

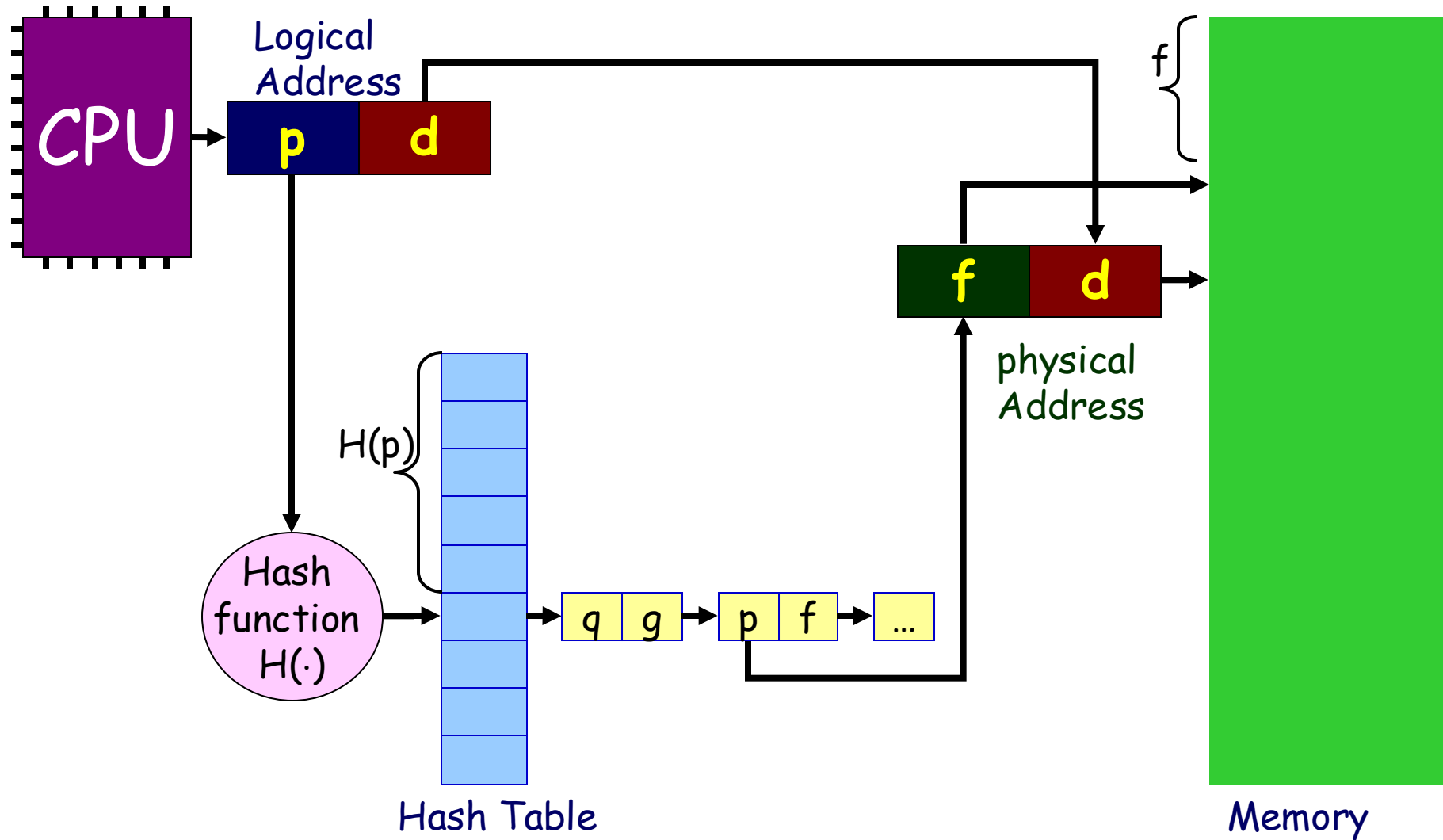
Hierarchical Paging Analysis

- Pros:
 - Easy memory allocation
 - Easy sharing
- Cons:
 - Page tables need to be **contiguous**
 - Two (or more, if >2 levels) lookups per reference

Hashed Page Tables

- In common address spaces > 32 bits
- The logical page number is **hashed** into a page table
 - This page table contains a chain of elements hashing to the same location
- Logical page numbers are compared in this chain searching for a match
 - If a match is found, the corresponding physical frame is extracted

Hashed Page Table



Hashed Page Tables

- Size of page table is at least as large as amount of **logical memory** allocated to processes
- Physical memory may be much less
 - Much of process space may be out on disk or not in use

Inverted Page Table

- One entry for each real page of memory
- Entry consists of the logical address of the page stored in that real memory location, with information about the process that owns that page
- Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs

Inverted Page Table

P1

0	3	v
1	0	v
2	7	v
3	0	i

P2

0	4	v
1	1	v
2	0	i
3	0	i

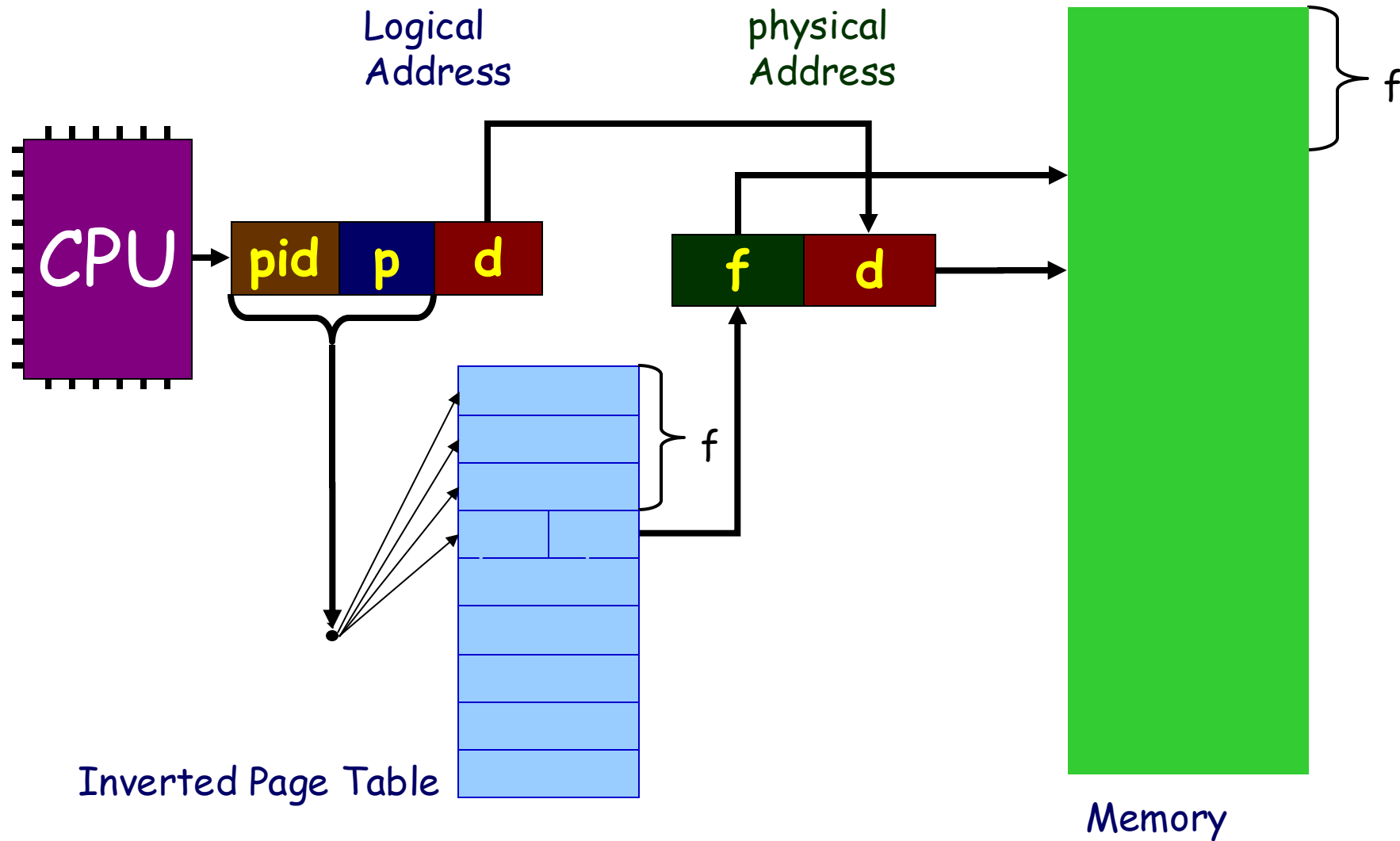
P3

0	5	v
1	2	v
2	6	v
3	0	i
4	0	i

0	P1	Page1
1	P2	Page1
2	P3	Page1
3	P1	Page0
4	P2	Page0
5	P3	Page0
6	P3	Page2
7	P1	Page2

Inverted Page Table

Inverted Page Table



Inverted Page Tables

- Systems that use inverted page tables have difficulty implementing shared memory.



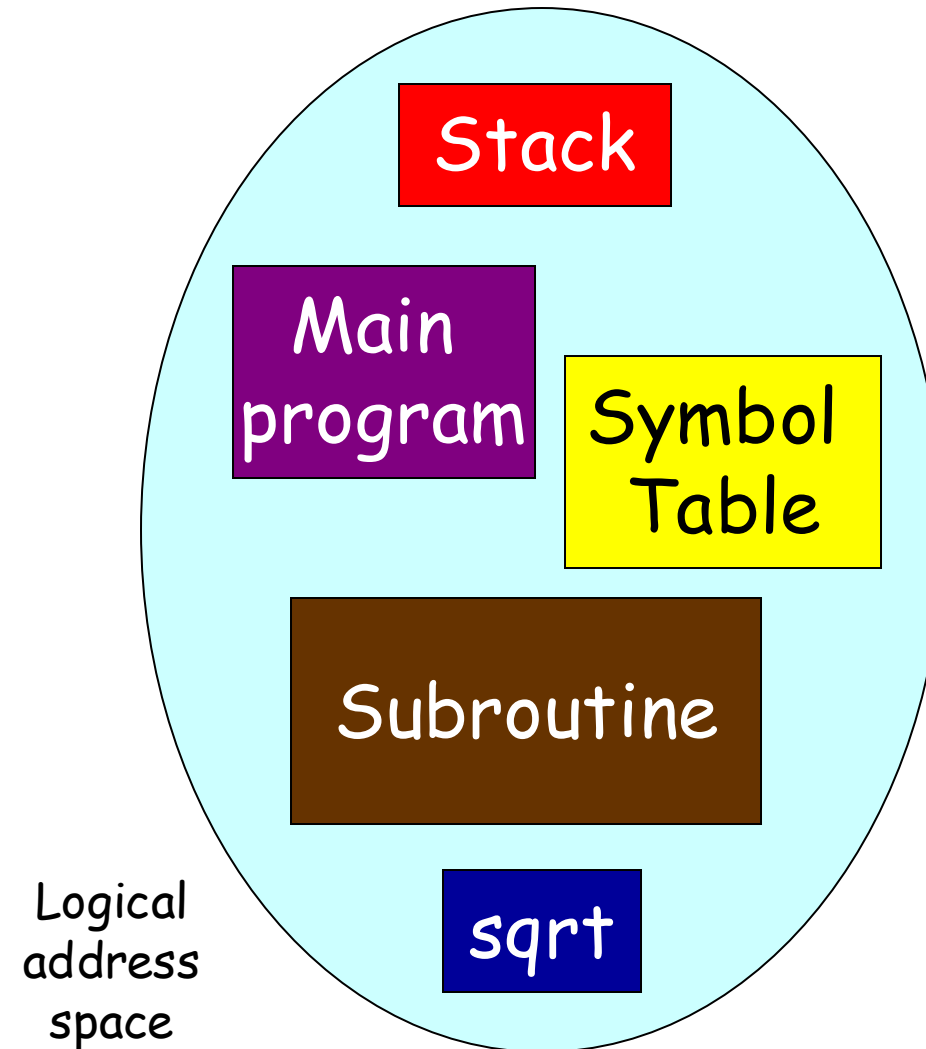
北京交通大学

Segmentation

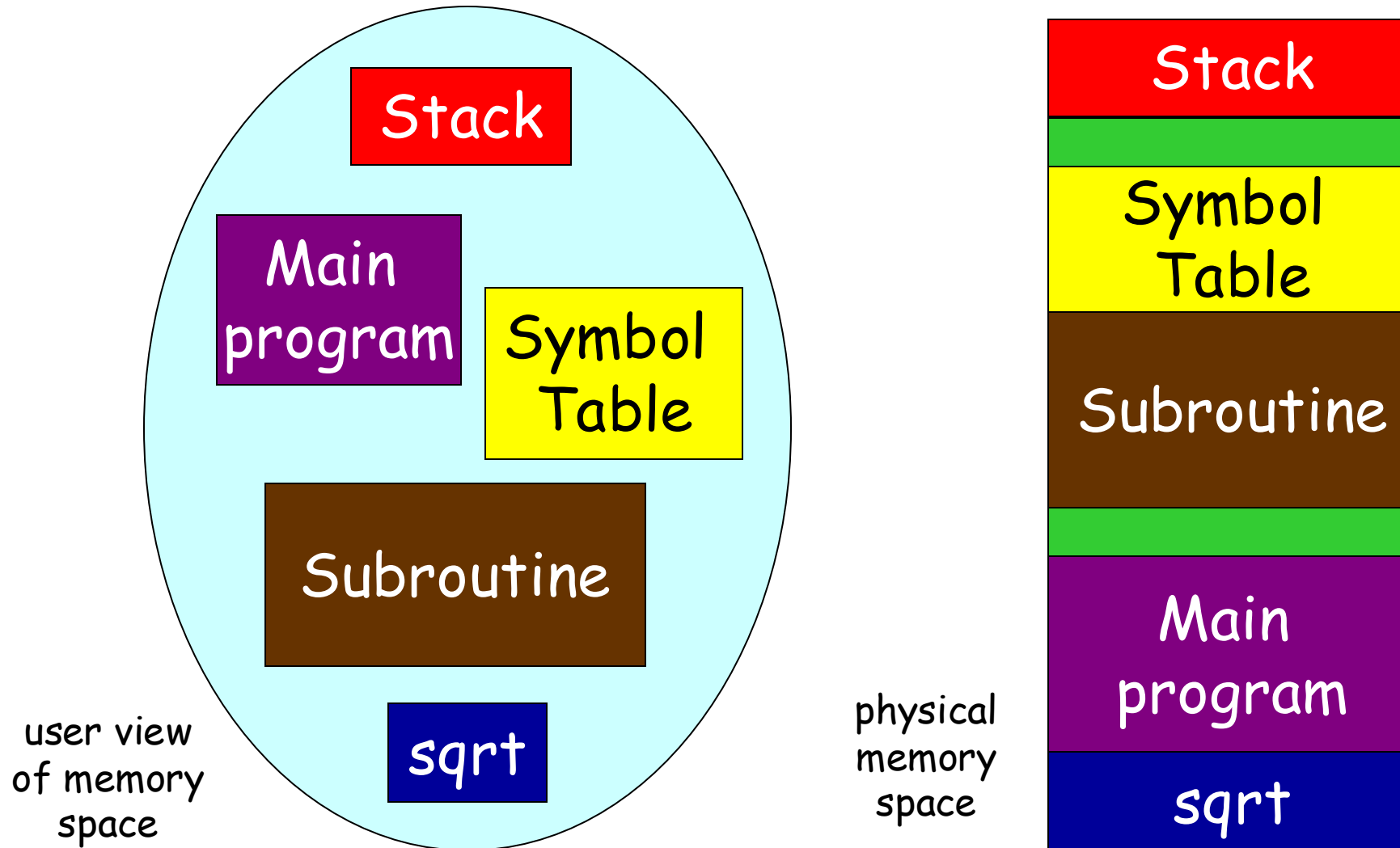
Segmentation

- Memory-management scheme that supports user view of memory
- A program is a collection of segments. A segment is a logical unit such as:
 - main program,
 - procedure,
 - function,
 - local variables, global variables,
 - stack,
 -

User's View of a Program

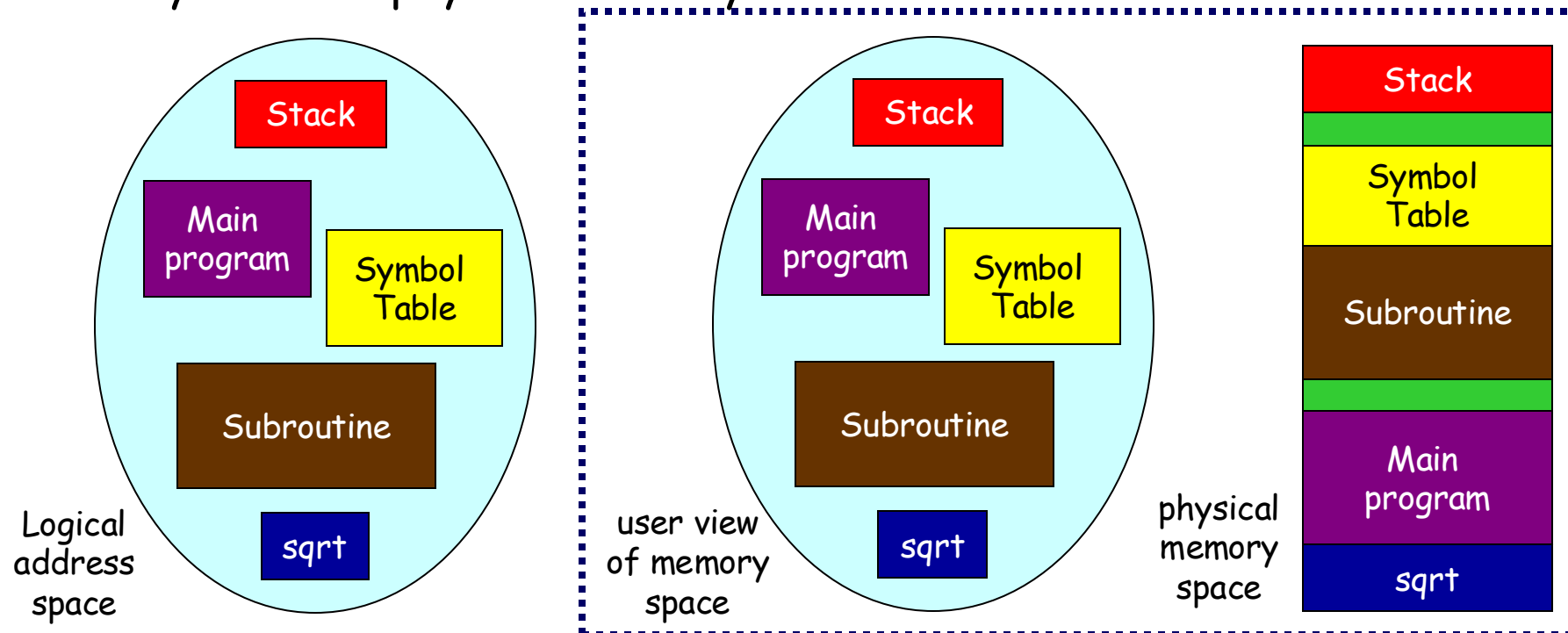


Logical View of Segmentation



More Flexible Segmentation

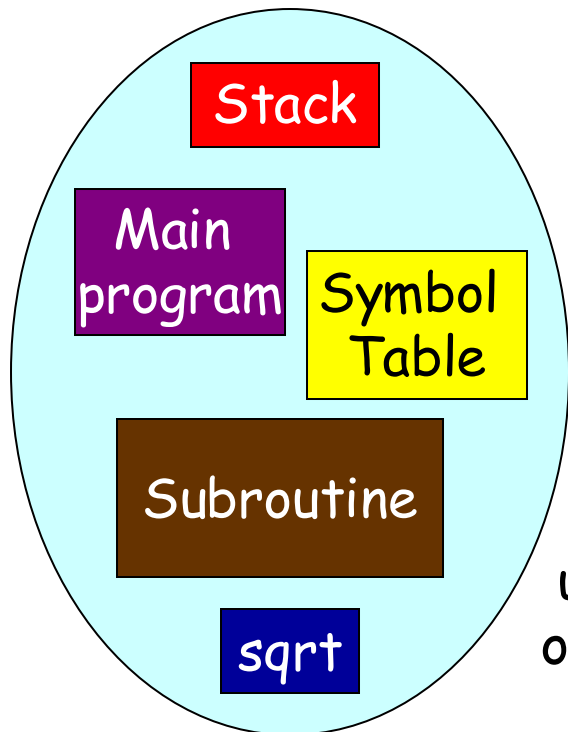
- Logical View: multiple separate segments
 - Typical: Code, Data, Stack
 - Others: memory sharing, etc
- Each segment is given region of contiguous memory
 - Has a base and limit
 - Can reside anywhere in physical memory



Segmentation Architecture

- Logical address consists of a two tuple :
 $\langle \text{segment-number}, \text{offset} \rangle$,
- Segment table – maps two-dimensional physical addresses; each table entry has:
 - base – contains the starting physical address where the segments reside in memory
 - limit – specifies the length of the segment
- Segment-table base register (STBR) points to the segment table's location in memory
- Segment-table length register (STLR) indicates number of segments used by a program;
 segment number s is legal if $s < \text{STLR}$

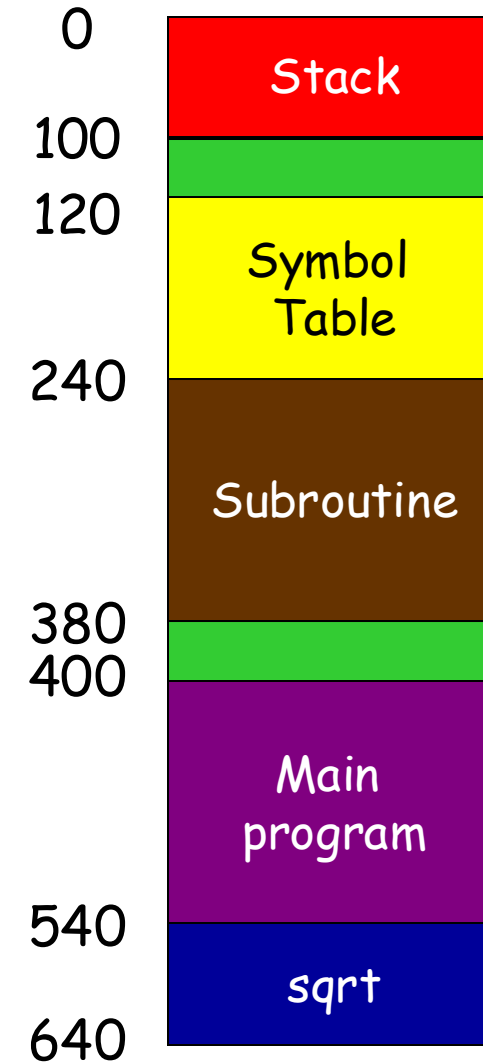
Segmentation Architecture



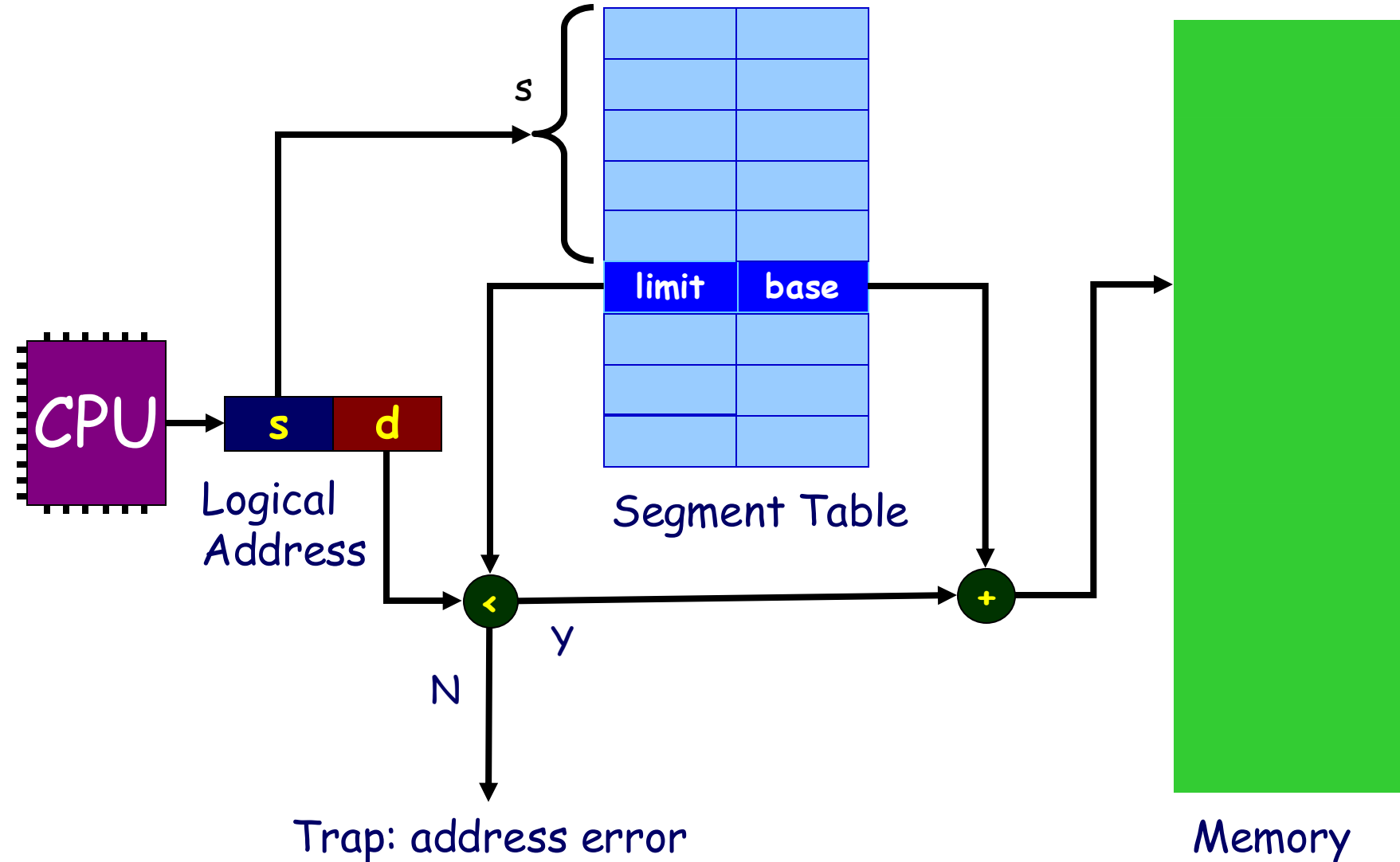
user view
of memory
space

	limit	base
0	100	0
1	120	120
2	140	240
3	140	400
4	100	540

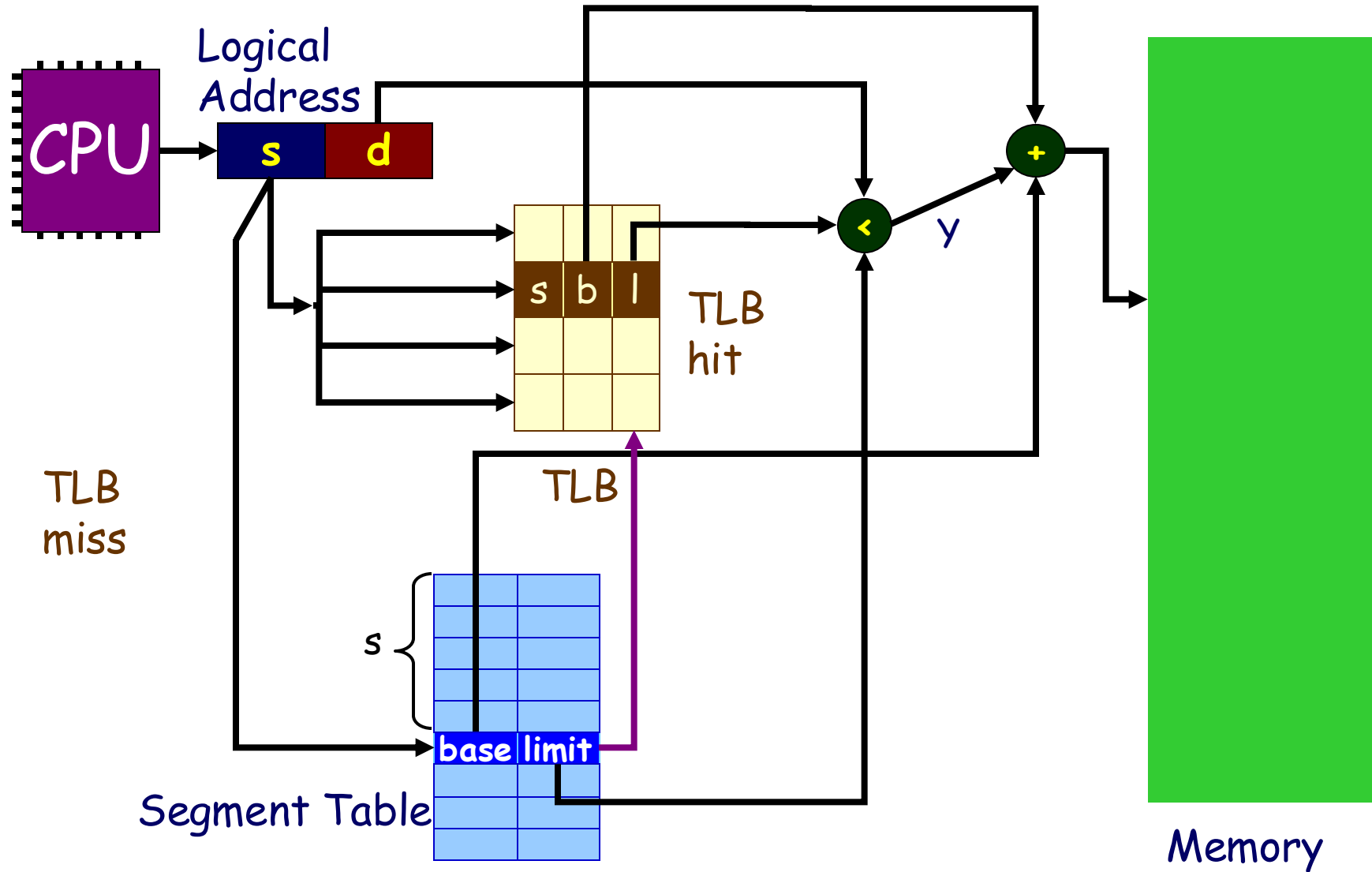
physical
memory
space



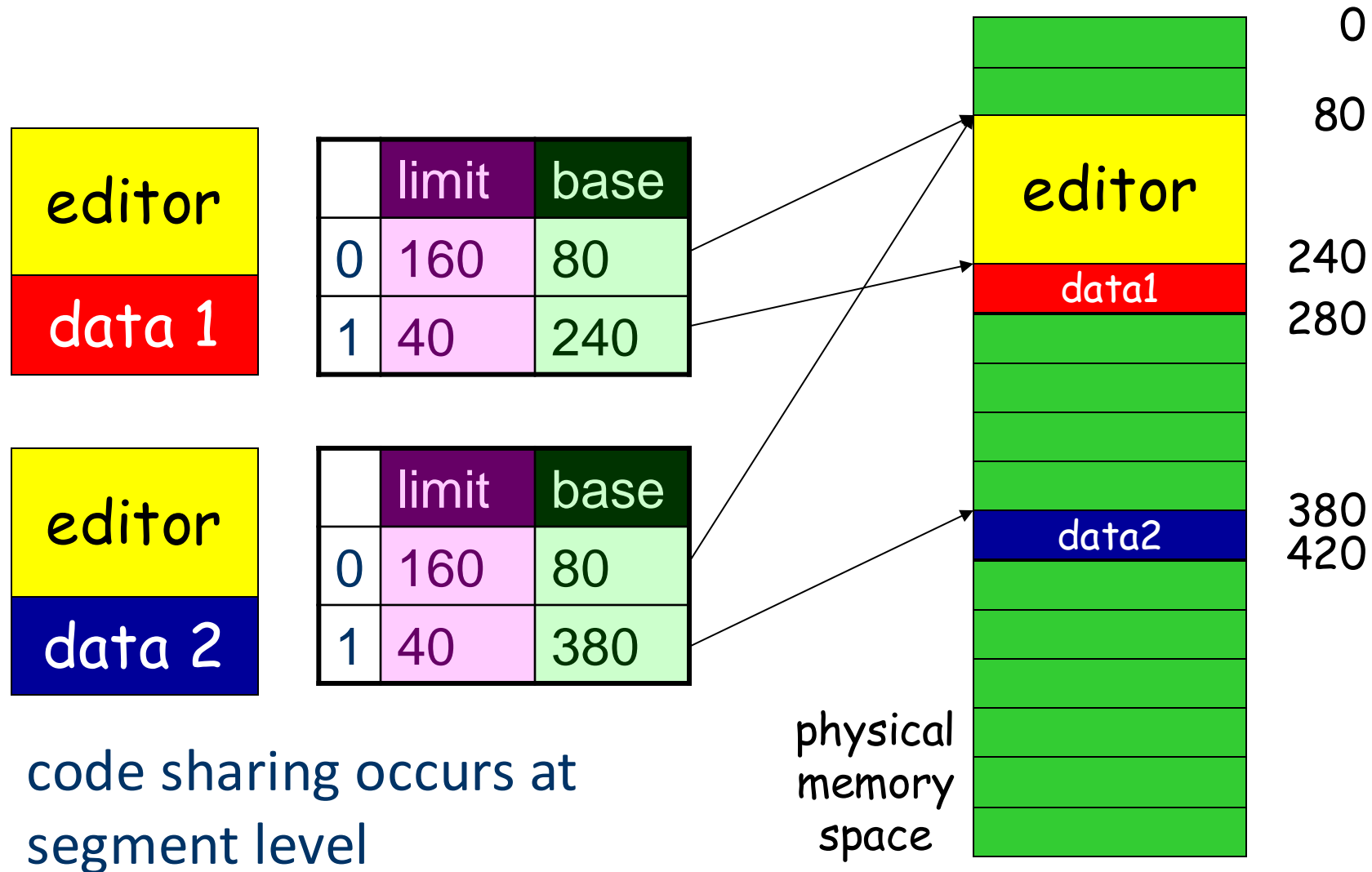
Segmentation Hardware



Segmentation Hardware With TLB



Shared Segments Example



Segmentation Architecture

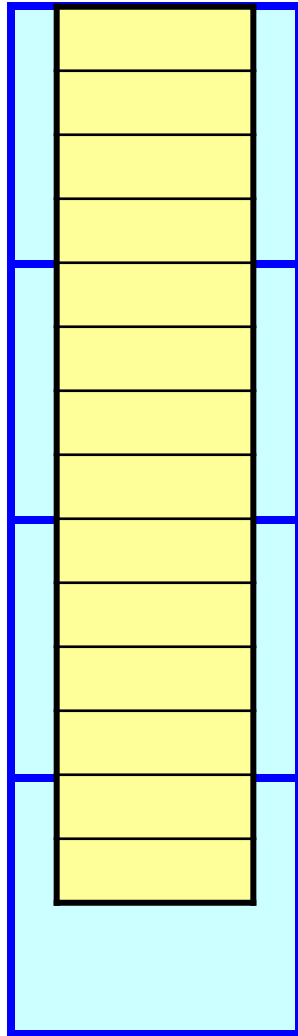
- Protection
 - With each entry in segment table associate:
 - validation bit = 0 illegal segment
 - read/write/execute privileges

Segmentation with Paging

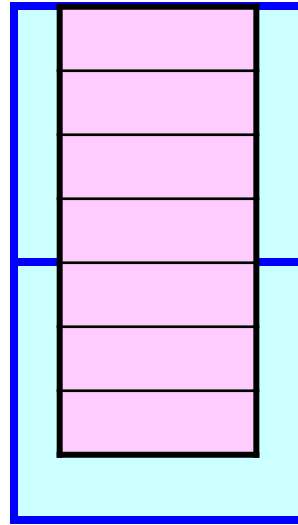
- Divide each segment into blocks of same size.



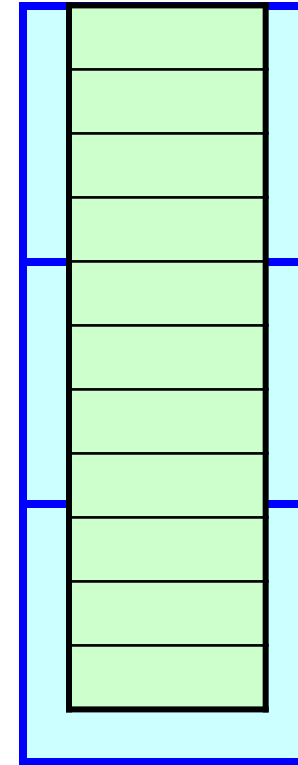
Segmentation with Paging



Segment1

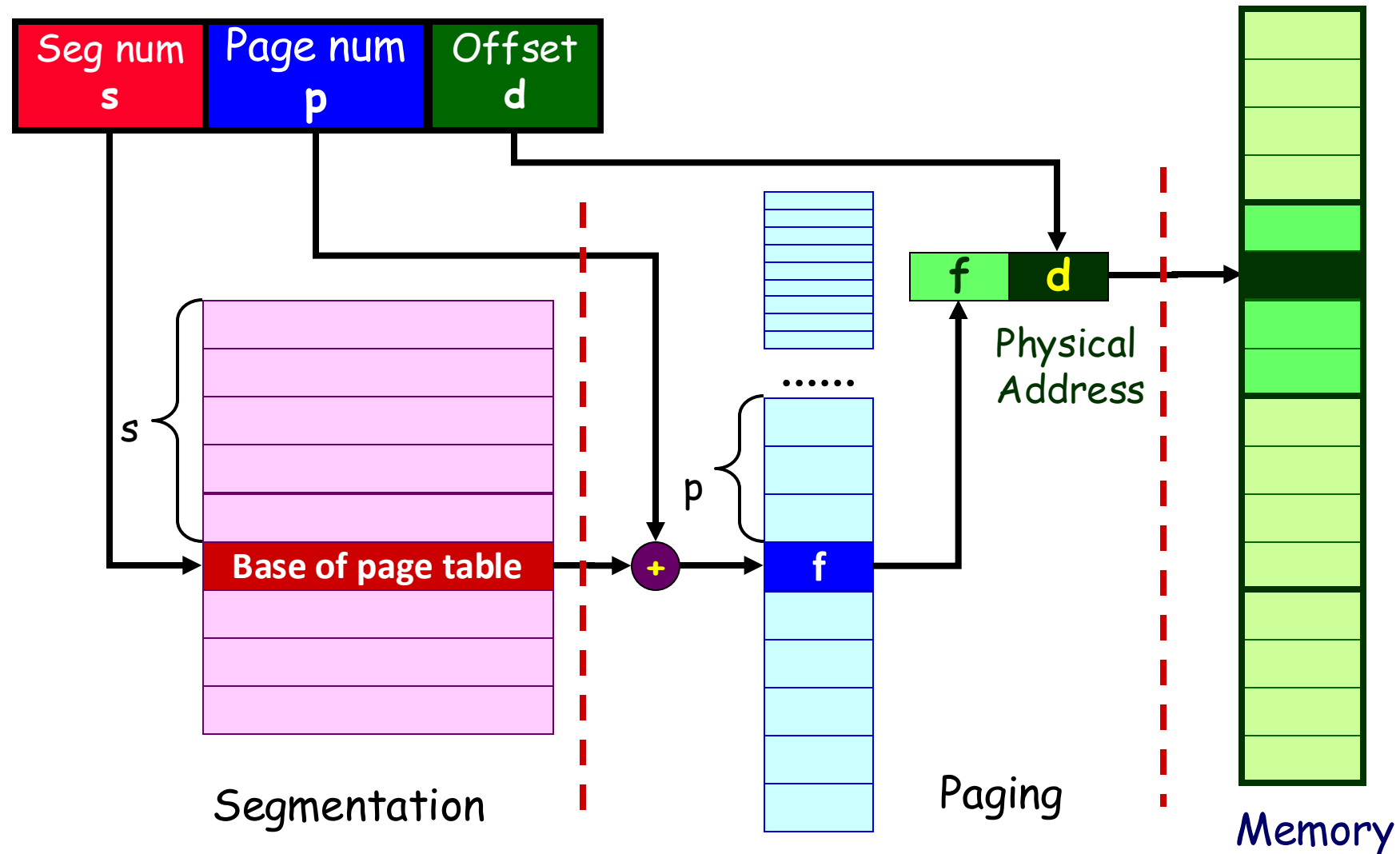


Segment2



Segment3

Segmentation with Paging



Summary

- Memory is a resource that must be shared
 - Translation: Change Logical Addresses into Physical Addresses
 - Protection: Prevent unauthorized Sharing of resources
- Simple contiguous memory allocation
 - Base+limit registers restrict memory accessible to user
 - Can be used to translate as well
- Full translation of addresses through Memory Management Unit (MMU)

Summary

- Page Tables

- Memory divided into fixed-sized chunks of memory
- Logical page number from logical address mapped through page table to physical page number
- Offset of logical address same as physical address
- Large page tables can be placed into logical memory

- Multi-Level Tables

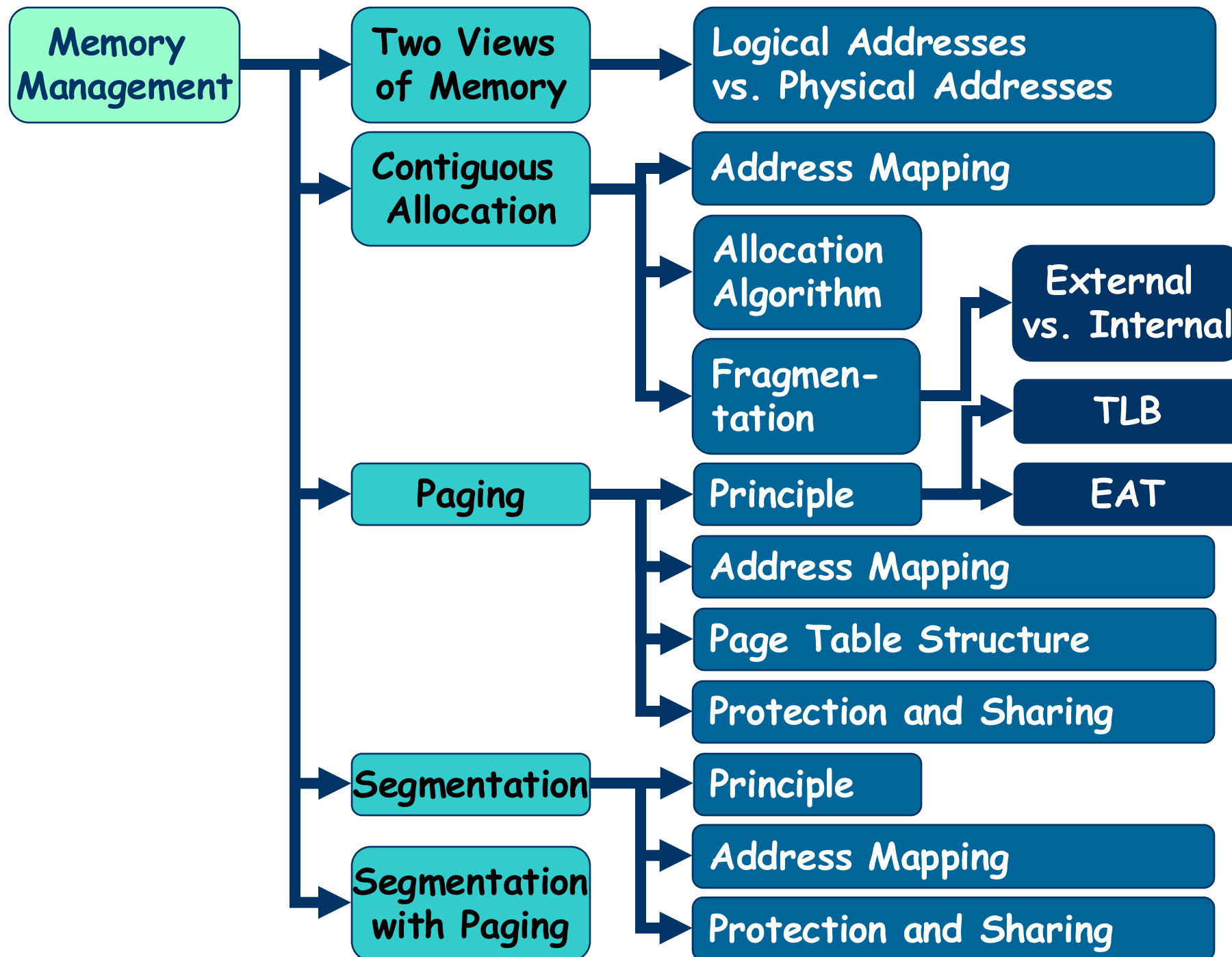
- Logical address mapped to series of tables
- Permit sparse population of address space

- Hashed page table

- Inverted page table

- Segment Mapping

- Each segment contains base and limit information
 - Offset (rest of address) adjusted by adding base



Summary

Logical Addresses vs. Physical Addresses

	Address Mapping	Protection	Sharing	Fragmentation	
				External	Internal
Contiguous Allocation					
Paging					
Segmentation					
Segmentation with Paging					



北京交通大学

Homework

Homework

- Exercises
 - EX-32 (Page 494):
 - 9.13, 9.15



北京交通大学

Thank you!

Q & A

