

Chapter 8: Memory Management

Prof. Li-Pin Chang
CS@NYCU



Chapter 8: Memory Management

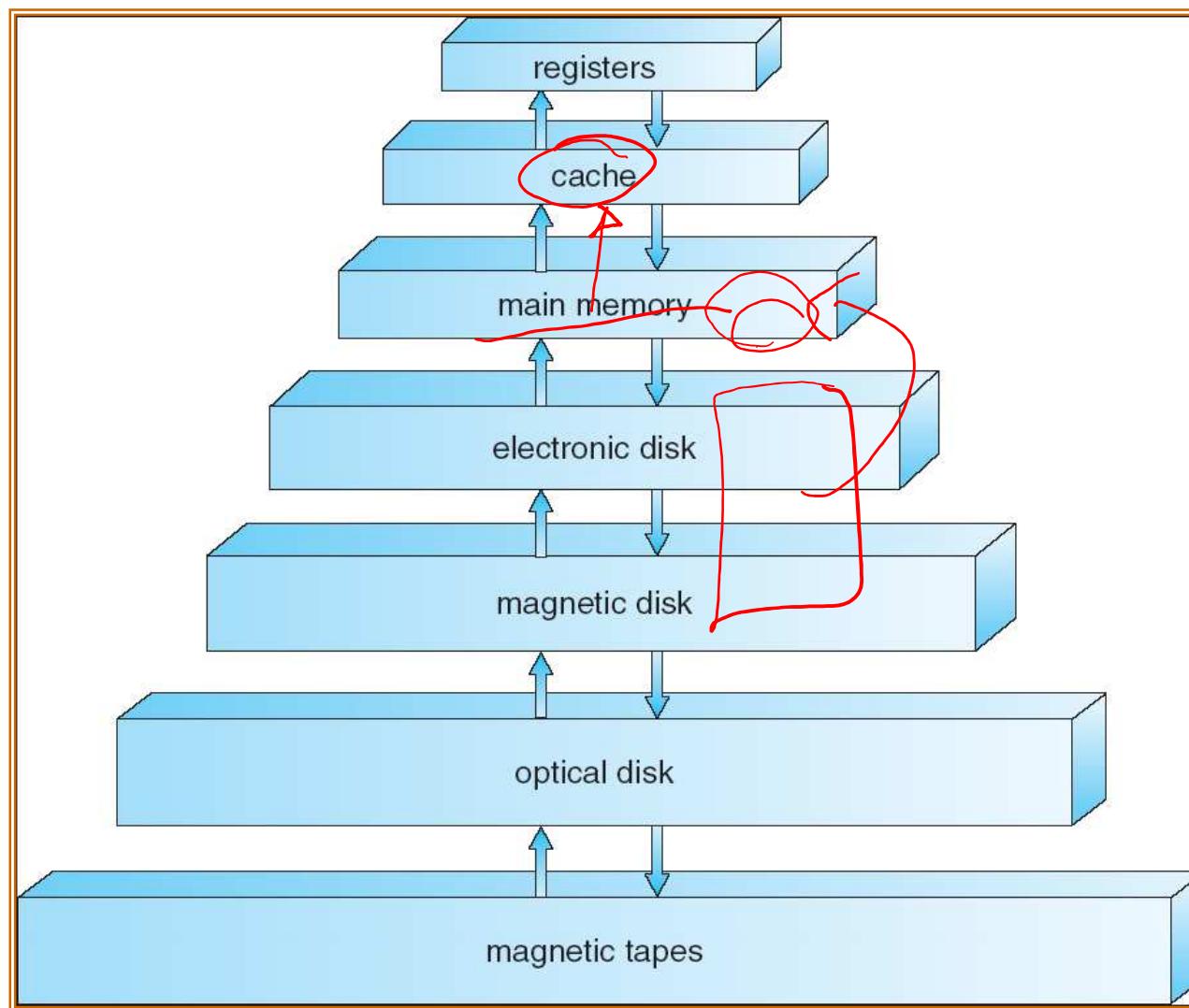
- Memory hierarchy and caching
- Address Binding
- Contiguous Memory Allocation
- Paging
- Structure of the Page Table
- Swapping
- Segmentation
- Example: The Intel Pentium

Memory Hierarchy

Memory Hierarchy

- Main memory – the only large storage media that the CPU can access directly
 - RAM (random access)
- Secondary storage – extension of main memory that provides large nonvolatile storage capacity
 - No random access, Magnetic disks – rigid metal or glass platters covered with magnetic recording material
 - Disk surface is logically divided into tracks, which are subdivided into sectors
 - The disk controller determines the logical interaction between the device and the computer

Storage Structure (Memory hierarchy)



Storage Structure (Memory hierarchy)

- Storage systems organized in hierarchy
 - Speed
 - Cost
 - Volatility
- Caching – copying information into faster storage system; main memory can be viewed as a last cache for secondary storage

Caching

- Important principle, performed at many levels in a computer (in hardware, operating system, software)
- Information in use copied from slower to faster storage temporarily
- Faster storage (cache) checked first to determine if information is there
 - If it is, information used directly from the cache (fast)
 - If not, data copied to cache and used there
- Cache smaller than storage being cached
 - Cache management important design problem
 - Cache size and replacement policy
- What we need?
 - An efficient lookup service
 - A replacement policy that minimizes cache misses

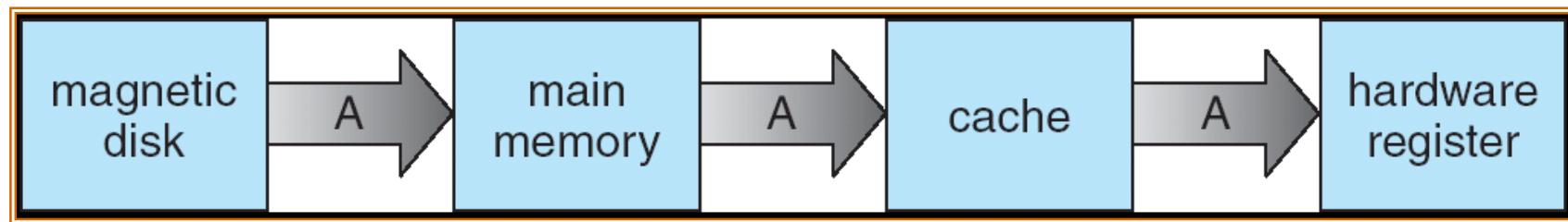
Performance of Various Levels of Storage

- Moving data among storage levels can be explicit or implicit

feeling

Level	1	2	3	4	5
Name	registers	cache	main memory	solid state disk	magnetic disk
Typical size	< 1 KB	< 16MB	< 64GB	< 1 TB	< 10 TB
Implementation technology	custom memory with multiple ports CMOS	on-chip or off-chip CMOS SRAM	CMOS SRAM	flash memory	magnetic disk
Access time (ns)	0.25 - 0.5	0.5 - 25	100ns	25,000 - 50,000	5,000,000
Bandwidth (MB/sec)	20,000 - 100,000	5,000 - 10,000	1,000 - 5,000	500	20 - 150
Managed by	compiler	hardware	operating system	operating system	operating system
Backed by	cache	main memory	disk	disk	disk or tape

Migration of Integer A from Disk to Register



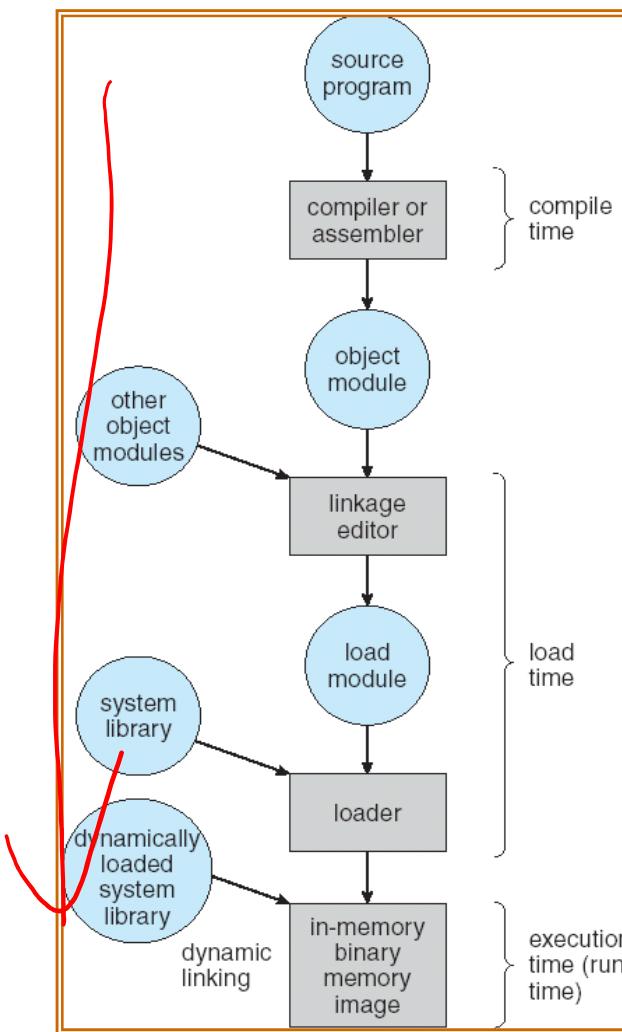
Address Binding

Address Binding

- Assigning memory addresses to instructions and data
- Program must be brought into memory and placed within a process for it to be run
- Input queue – collection of processes on the disk that are waiting to be brought into memory to run the program
- User programs go through several steps before running

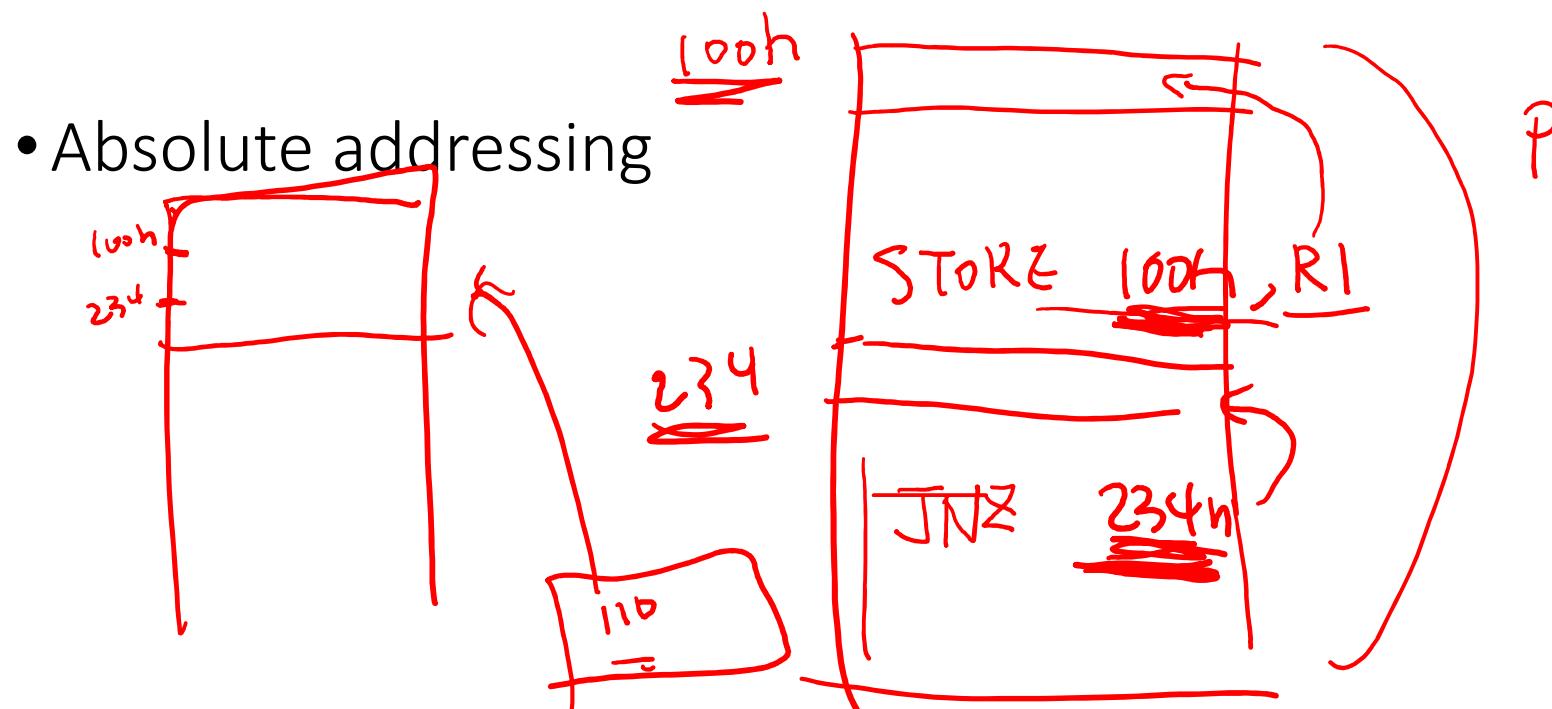
Binding of Instructions and Data to Memory

- Address binding of instructions and data to memory addresses can happen at three different stages
 - Compile time
 - Load time
 - Execution time



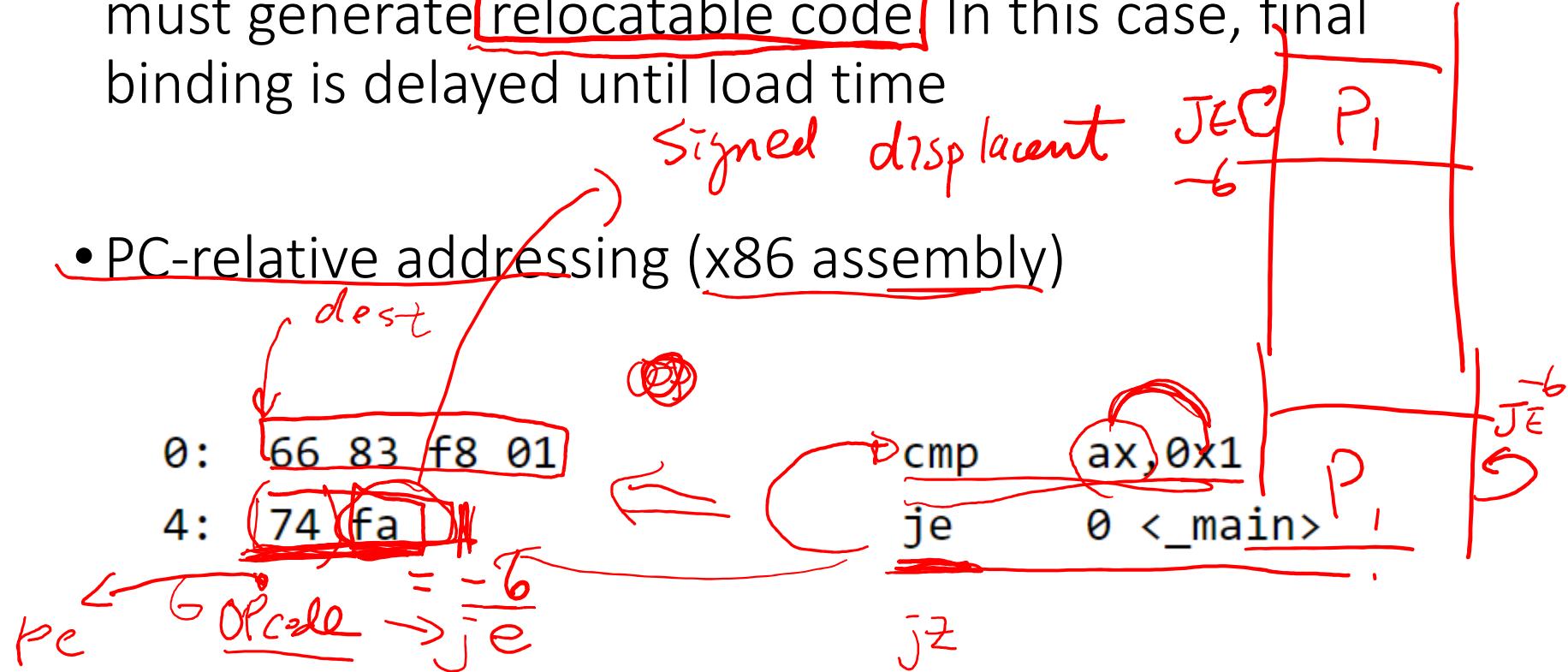
Binding of Instructions and Data to Memory

- Compile time: If memory location known a priori, absolute code can be generated; must recompile code if the starting location changes



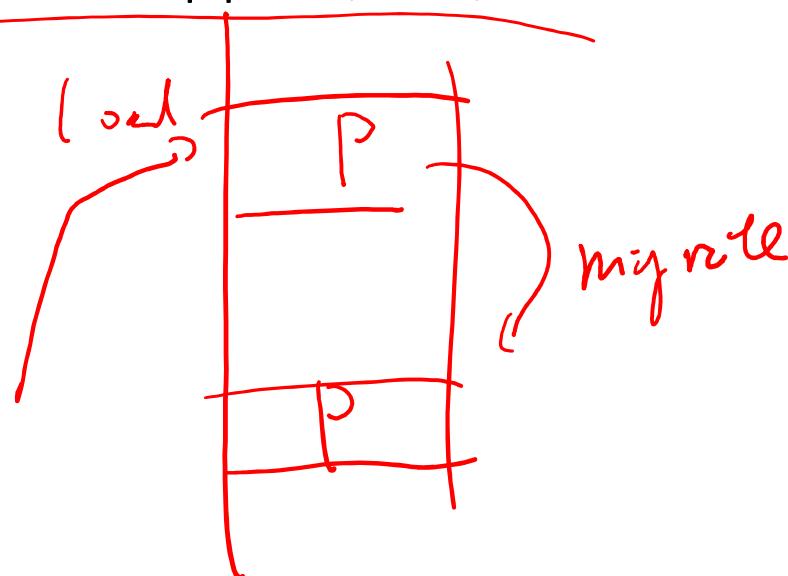
Binding of Instructions and Data to Memory

- Load time: If it is not known at compile time where the process will reside in memory, then the compiler must generate relocatable code. In this case, final binding is delayed until load time



Binding of Instructions and Data to Memory

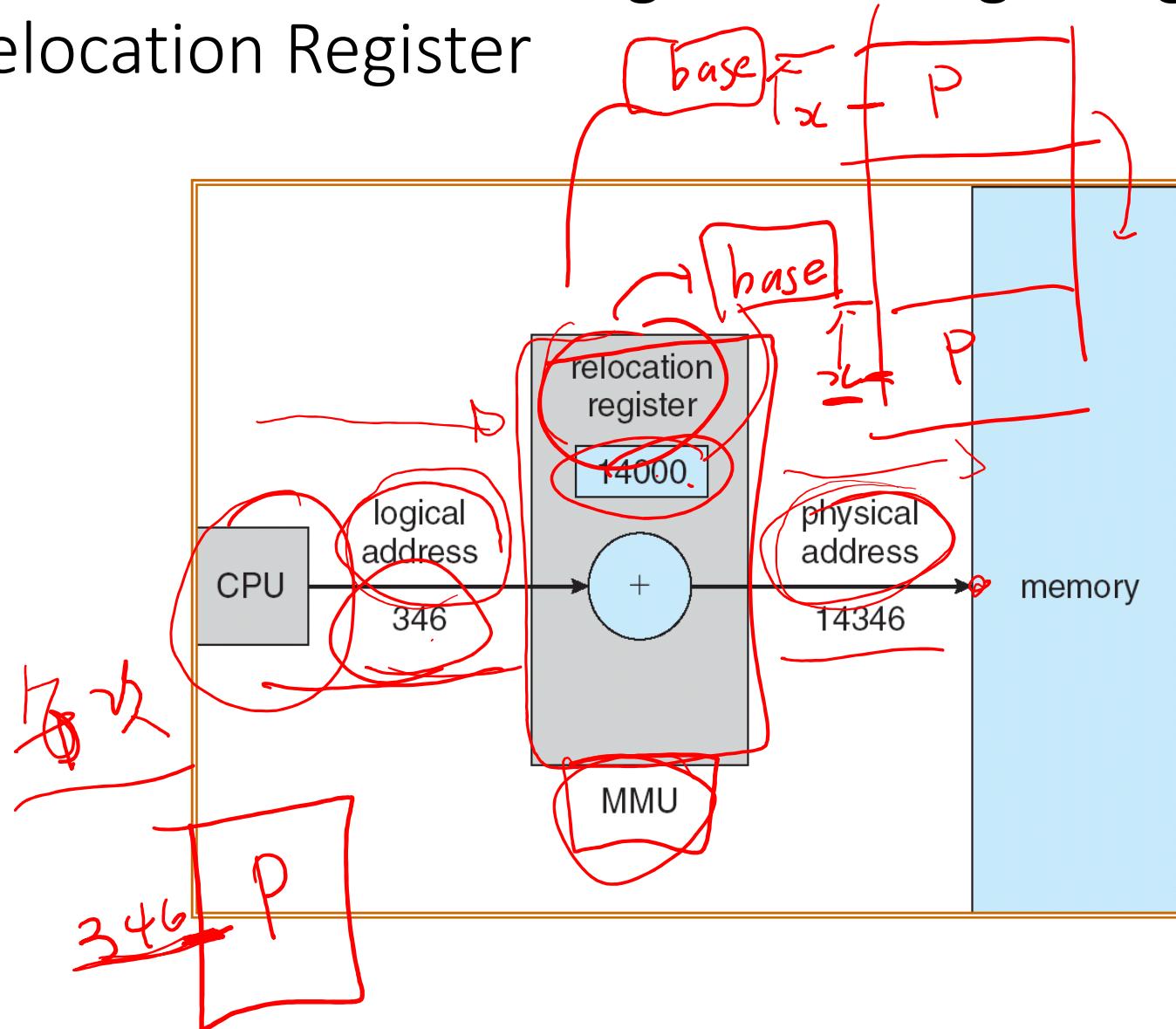
- Execution time: Binding delayed until run time if the process can be moved during its execution from one memory segment to another
- Requiring hardware support, i.e., MMU



Execution Time Binding

- The concept of a logical address space that is bound to a separate physical address space is central to proper memory management
 - **Logical address** – generated by the CPU; also referred to as **virtual address**
 - **Physical address** – address seen by the memory unit
- The user program deals with logical addresses; it never sees the real physical addresses
- Need hardware support to translate **logical addresses** into **physical addresses** during runtime
 - Relocation registers (base/limit)

Execution Time Binding/Relocating using a Relocation Register



Memory-Management Unit (MMU)

- A hardware component in the CPU that translates logical addresses into physical addresses

• Relocation

• Paging *虚地址*

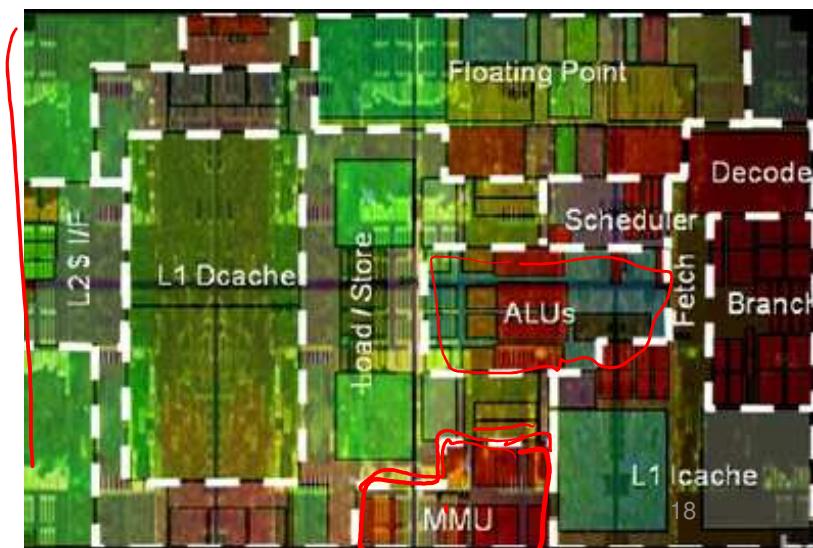
- Address mapping

- Virtual memory

• Segmentation *实地址*

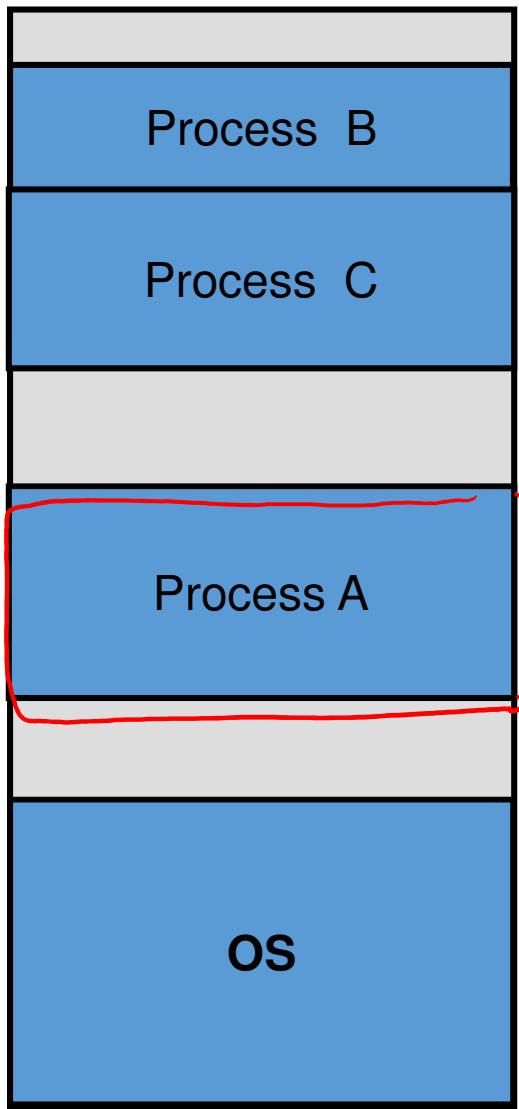
- Memory protection

[Project Denver 64-bit CPU core](#)



CONTIGOUS MEMORY ALLOCATION

physical mem



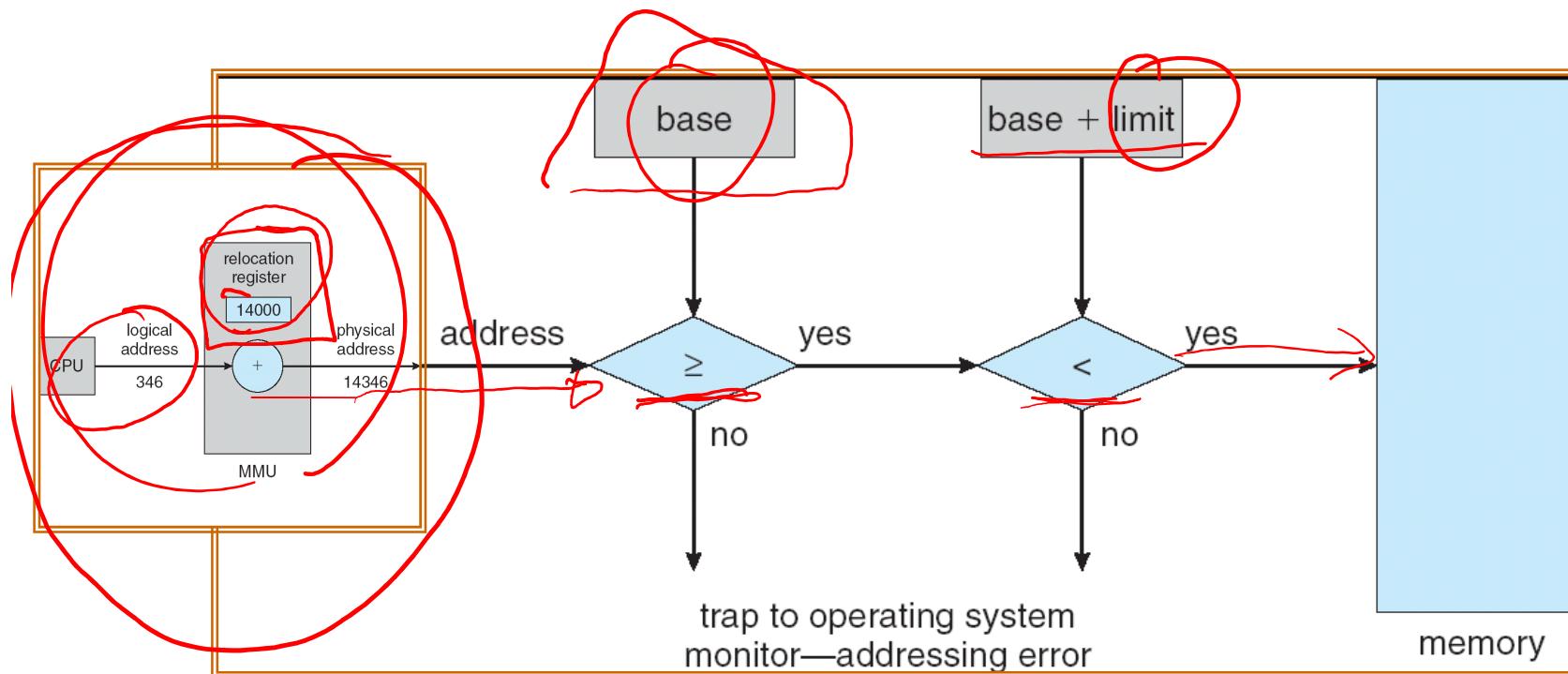
- Processes are allocated to contiguous memory space
- Processes can be loaded into any contiguous and sufficiently large memory space
- As processes arrive and leave, free space may be fragmented into many pieces

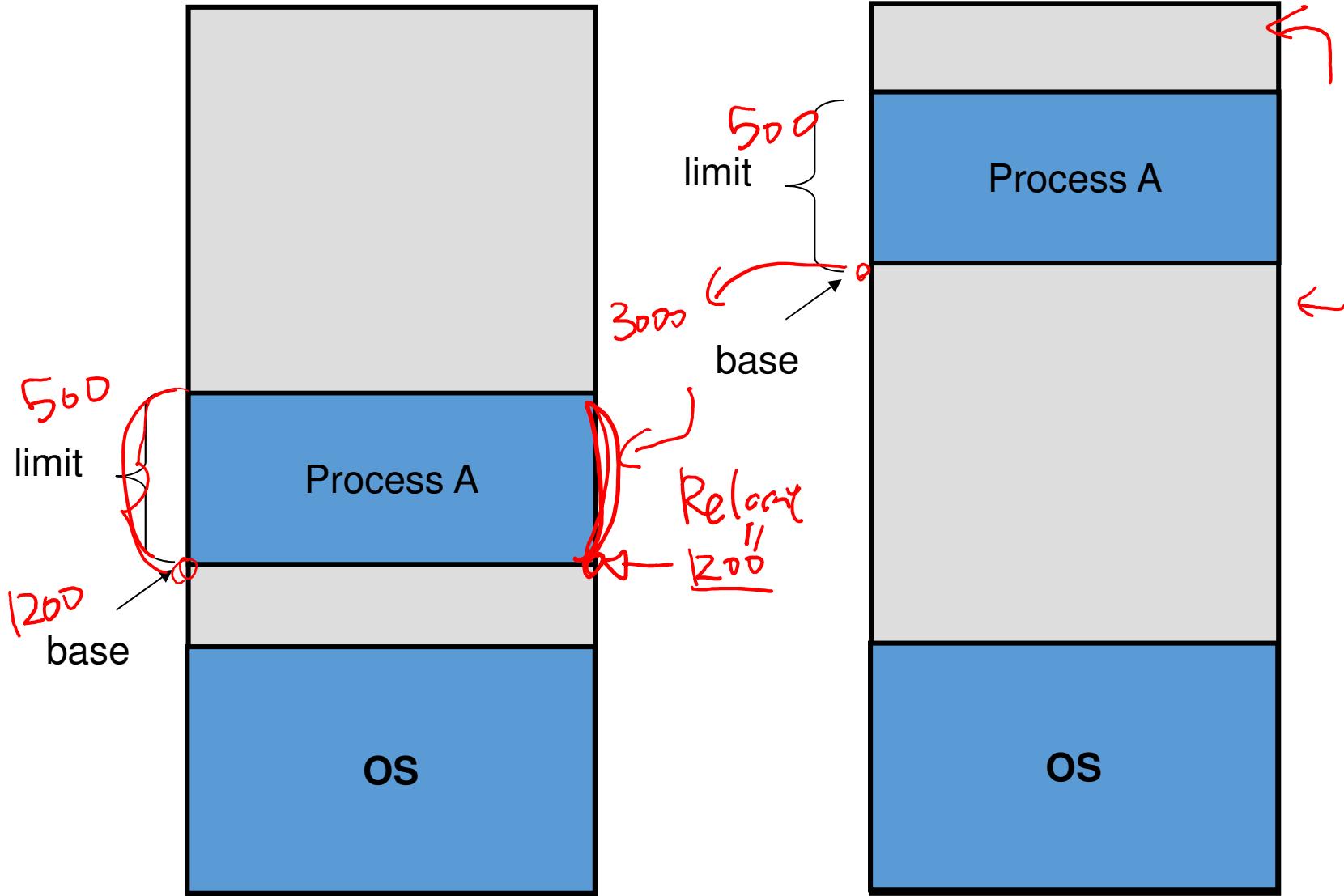
Relocation Ry

Contiguous Allocation

- Main memory usually into two partitions:
 - Resident ~~operating system~~, usually held in low memory with interrupt vector
 - User ~~processes~~ then held in high memory
- Single-partition allocation
 - Relocation-register scheme used to protect user processes from each other, and from changing operating-system code and data
 - Relocation register contains value of smallest physical address; limit register contains range of logical addresses – each logical address must be less than the limit register

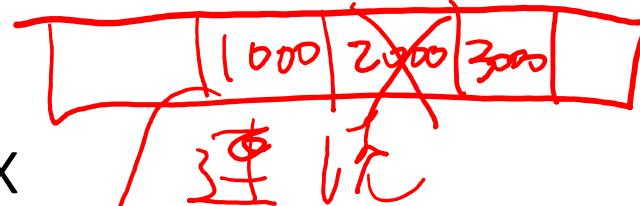
Memory Mapping and Protection





Heap Management in Linux

- Also a problem of contiguous memory allocation
- Heap is managed by `malloc()`
 - Part of the data segment
- `malloc()` finds free spaces for applications
- If heap is full, call `brk()/sbrk()` to increase the size of data segment

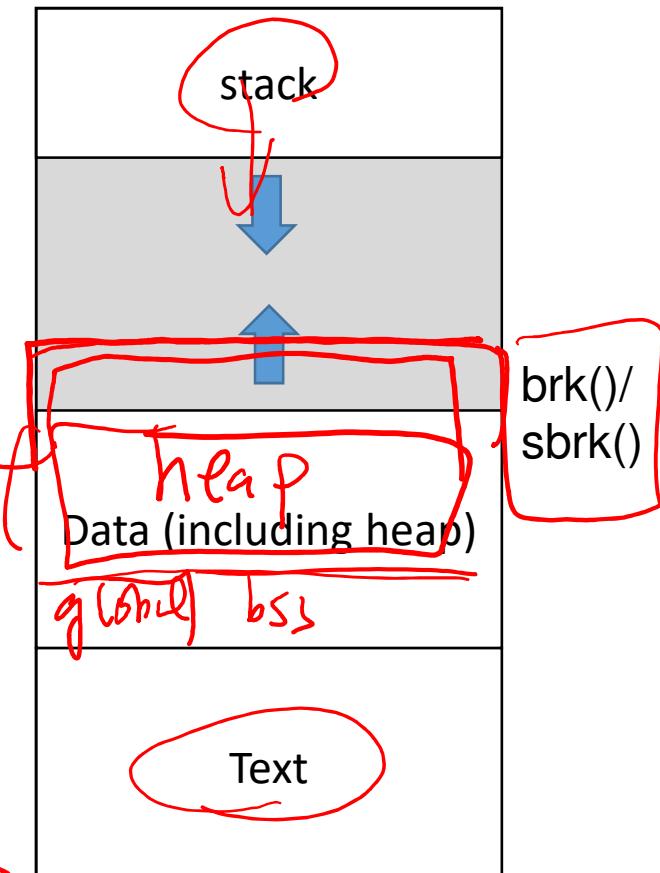


MAX
*P

*

Virtual

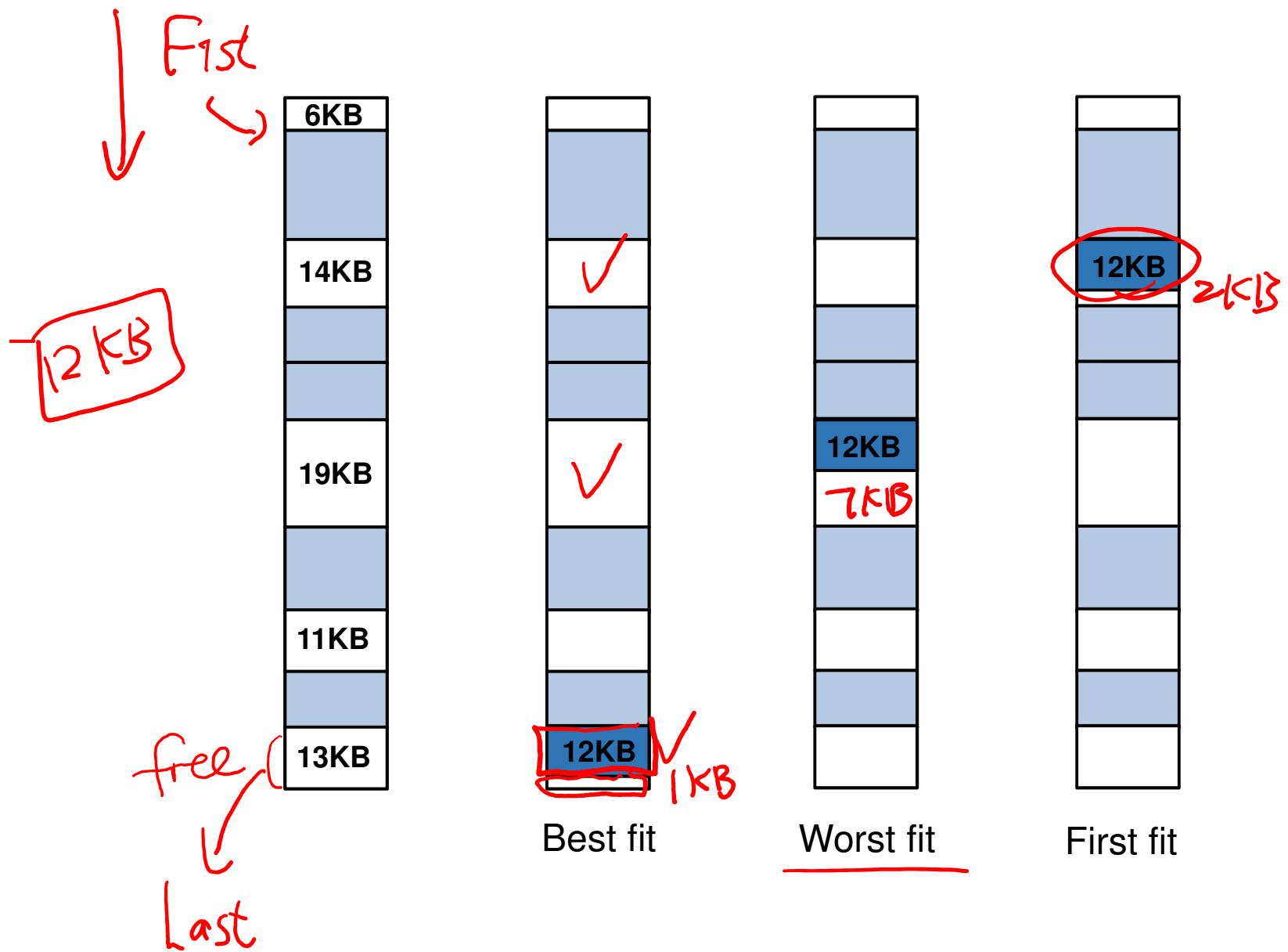
9



Memory Allocation

How to satisfy a request of size n from a list of free
holes? *Policy*

- **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.



12 8

	20			108	
A20	20	15		93	
A15	20	15	10	83	
A10	20	15	10	25	58
A25	20	15	10	25	58
D20	20	15	10	25	58
D10	20	15	10	25	58
A8	8	12	15	10	25
A30	8	12	15	10	25
D15	8	37		25	30
A15	8	15	22	25	30
					28

P D

First fit

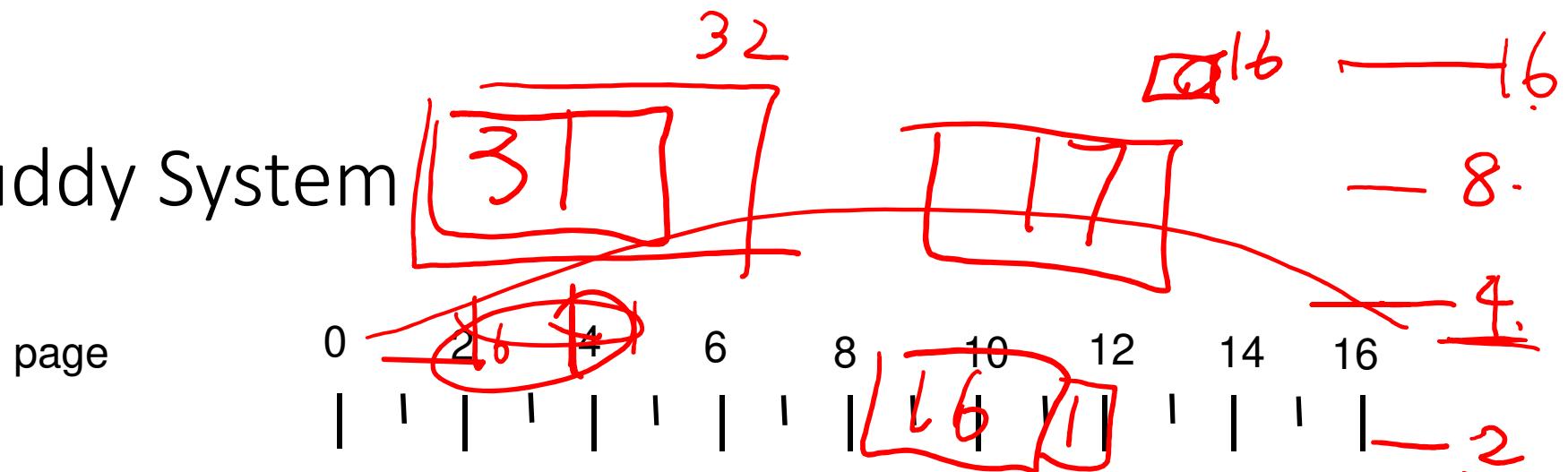
D10	20	15	10	25		58
A8	20	15	8 2	25		58
A30	20	15	8 2	25	30	28
D15	35		8 2	25	30	28
A15	35		8 2	25	30	15 13

Best fit

D10	20	15	10	25		58
A8	20	15	10	25	8	50
A30	20	15	10	25	8	30 20
D15	45		25	8	30	20
A15	15	30	25	8	30	20

Worst fit

Buddy System



Initialization

requestA (2)

requestB (1)

requestC (2)

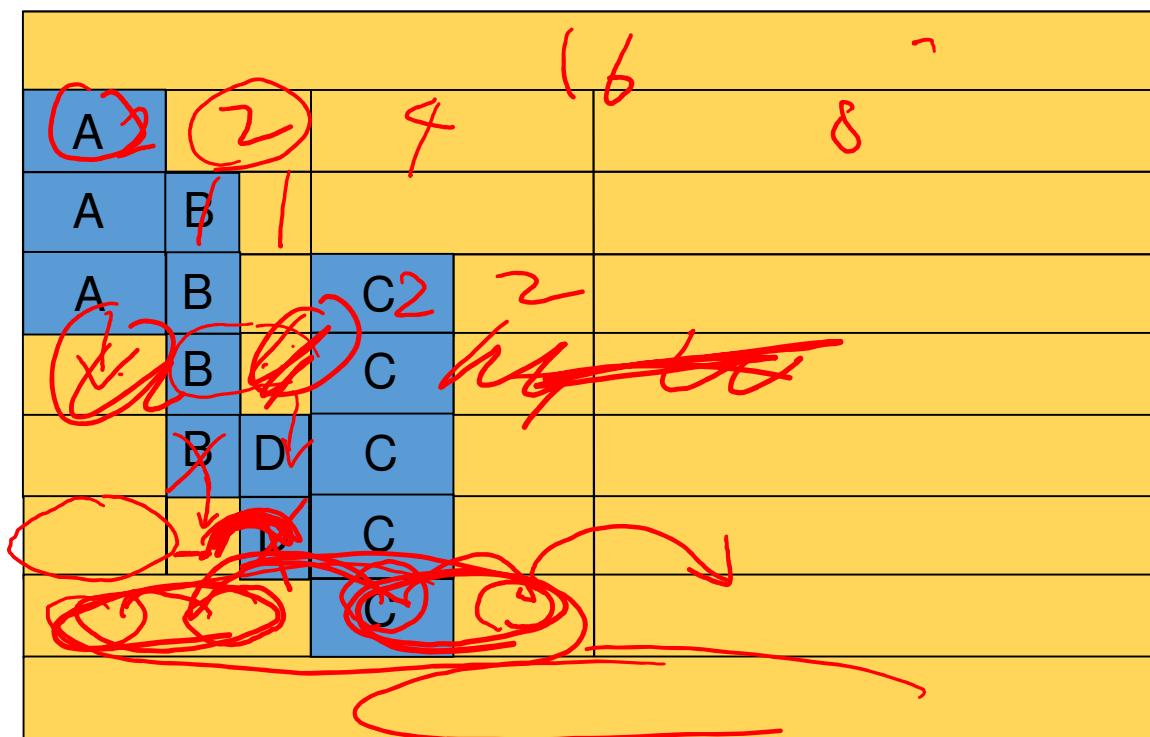
Free A*

requestD (1)

Free B

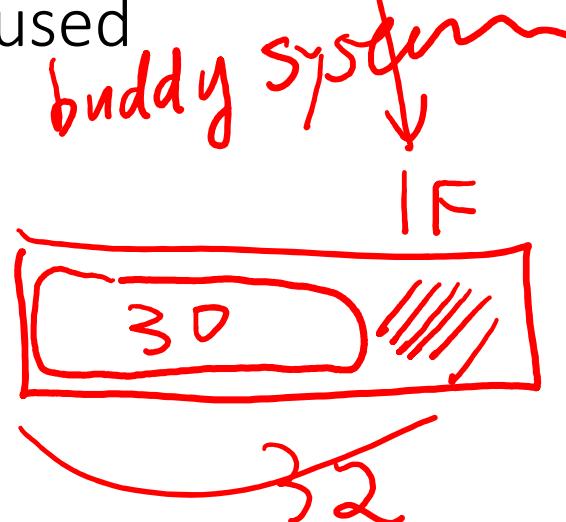
Free D

Free C



Fragmentation

- External Fragmentation – total memory space exists to satisfy a request, but it is not contiguous
- Internal Fragmentation – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used



Comparison

Lower is better

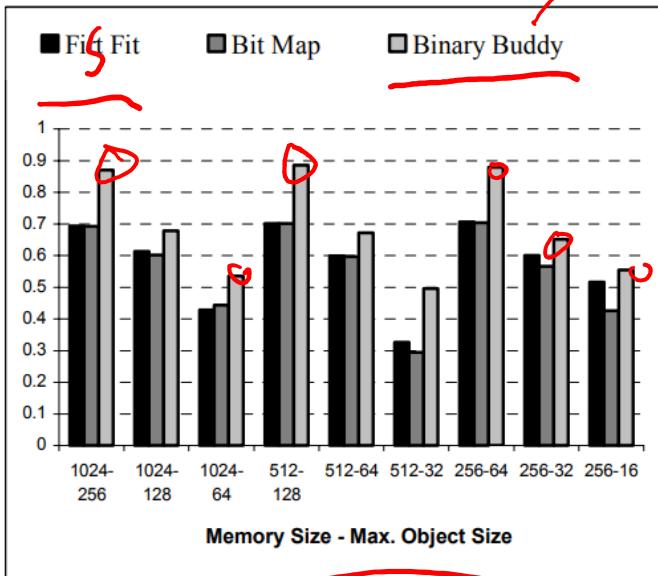


Figure 8: Total fragmentation values of techniques

External frag

Lower is better

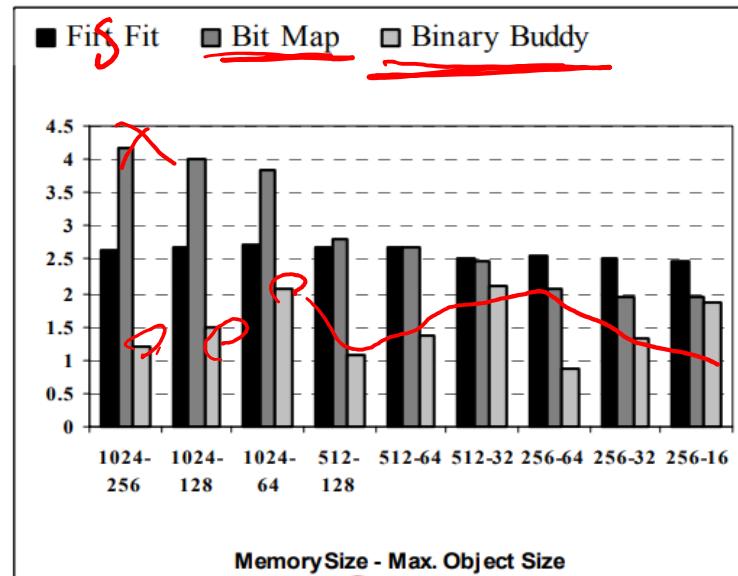


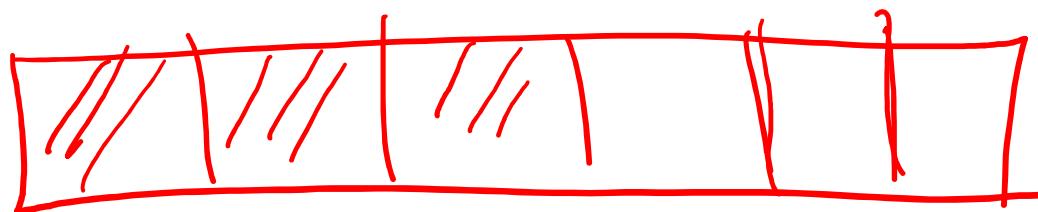
Figure 9: Allocation duration of techniques (micro second)

- BF, FF, and Buddy System are practical choices
- There is no optimal solution to the contiguous memory allocation problem as the problem has been proven intractable (NP-hard)

Slab Allocation

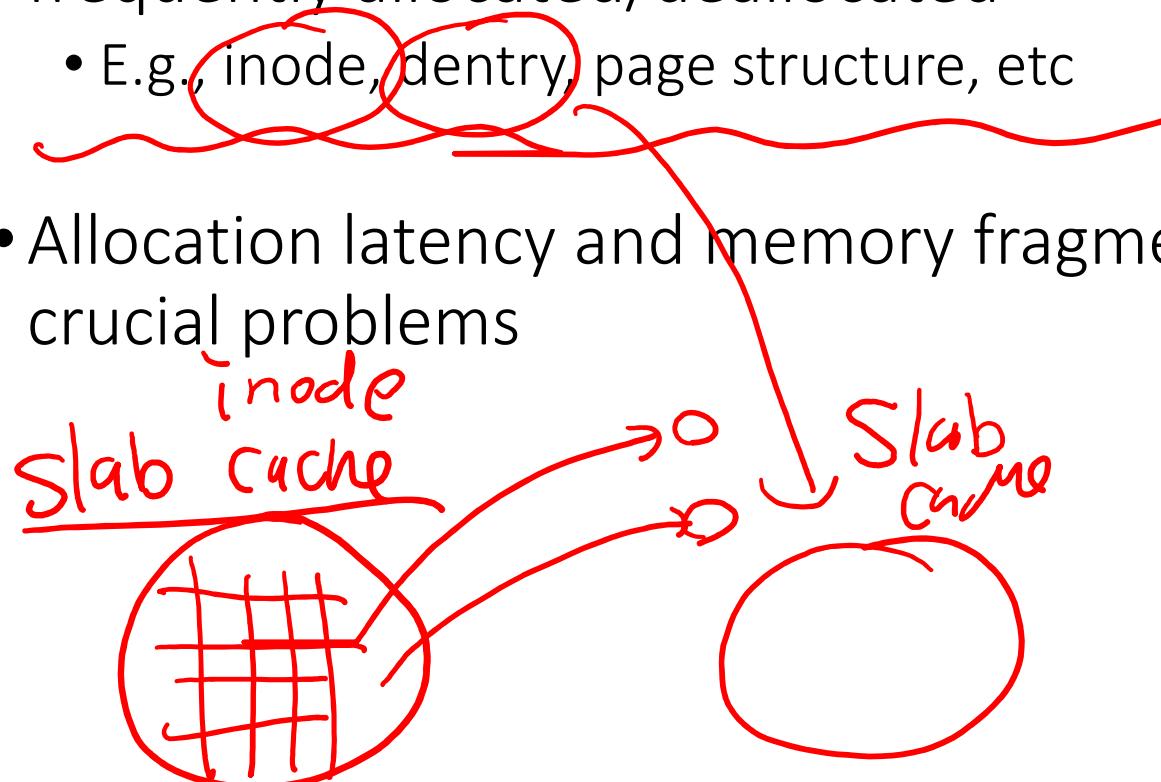
5K 7K 2K

- The root cause of external fragmentation is that allocation sizes are not uniform
- Idea: allocate objects of the same size from the same memory pool
- Solution: slab allocation

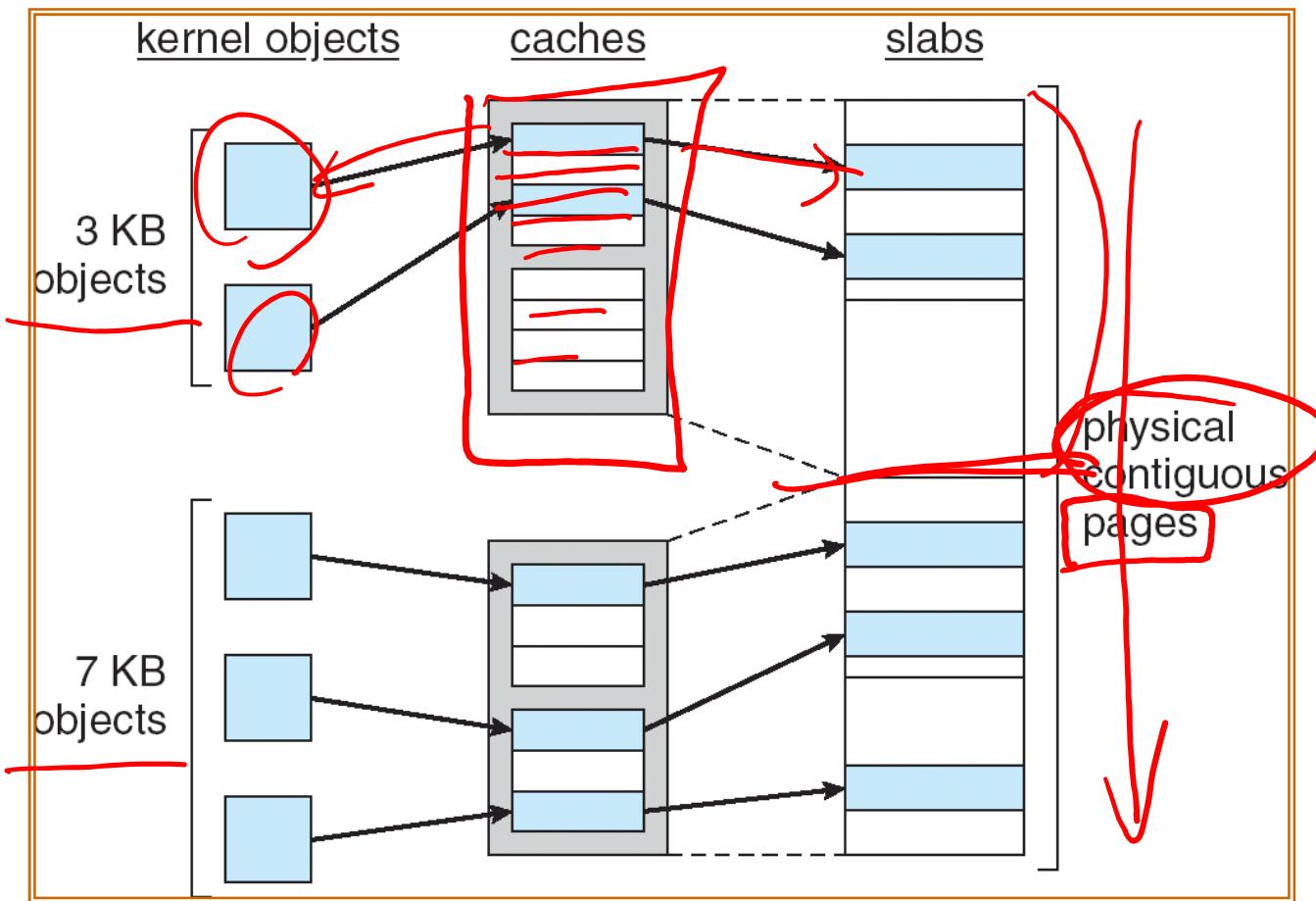


Linux Kernel Memory Allocation

- Kernel objects (of different sizes) are small, frequently allocated/deallocated
 - E.g., inode, dentry, page structure, etc
- Allocation latency and memory fragmentation are crucial problems



Linux kernel slab cache



Benefits:

- No fragmentation
- Quick memory allocation
- Objects get physically contiguous memory

[cache[obj][obj][obj]]
[page][page][page]

Case Study

- malloc() of glibc 2.28
 - A variation of Best Fit
- Linux kernel provides slab allocator and buddy system
 - Frequent allocation/deallocation of objects of the same size, use slab cache
 - Allocation/deallocation of objects of various sizes, use buddy system

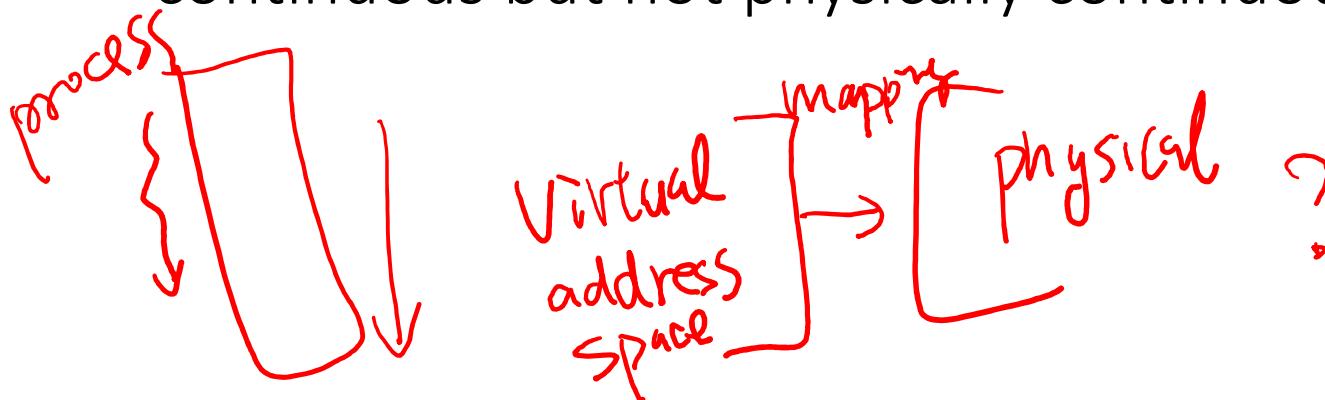
PAGING

External Fragmentation

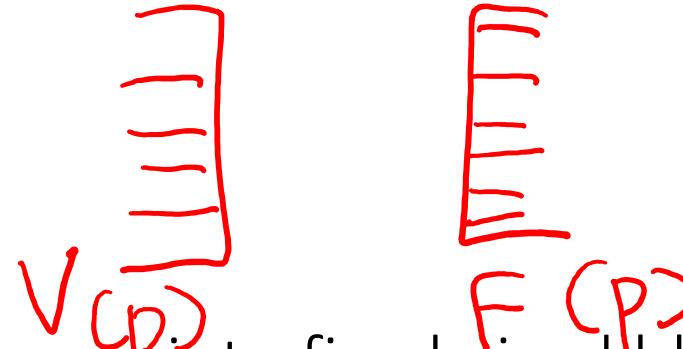


- Compaction

- Migrate allocated memory chunks together to make free space contiguous
- Will relocate programs— need execution time binding
- Occasionally slows down the system
- A better approach: What if memory can be logically continuous but not physically continuous?



Paging

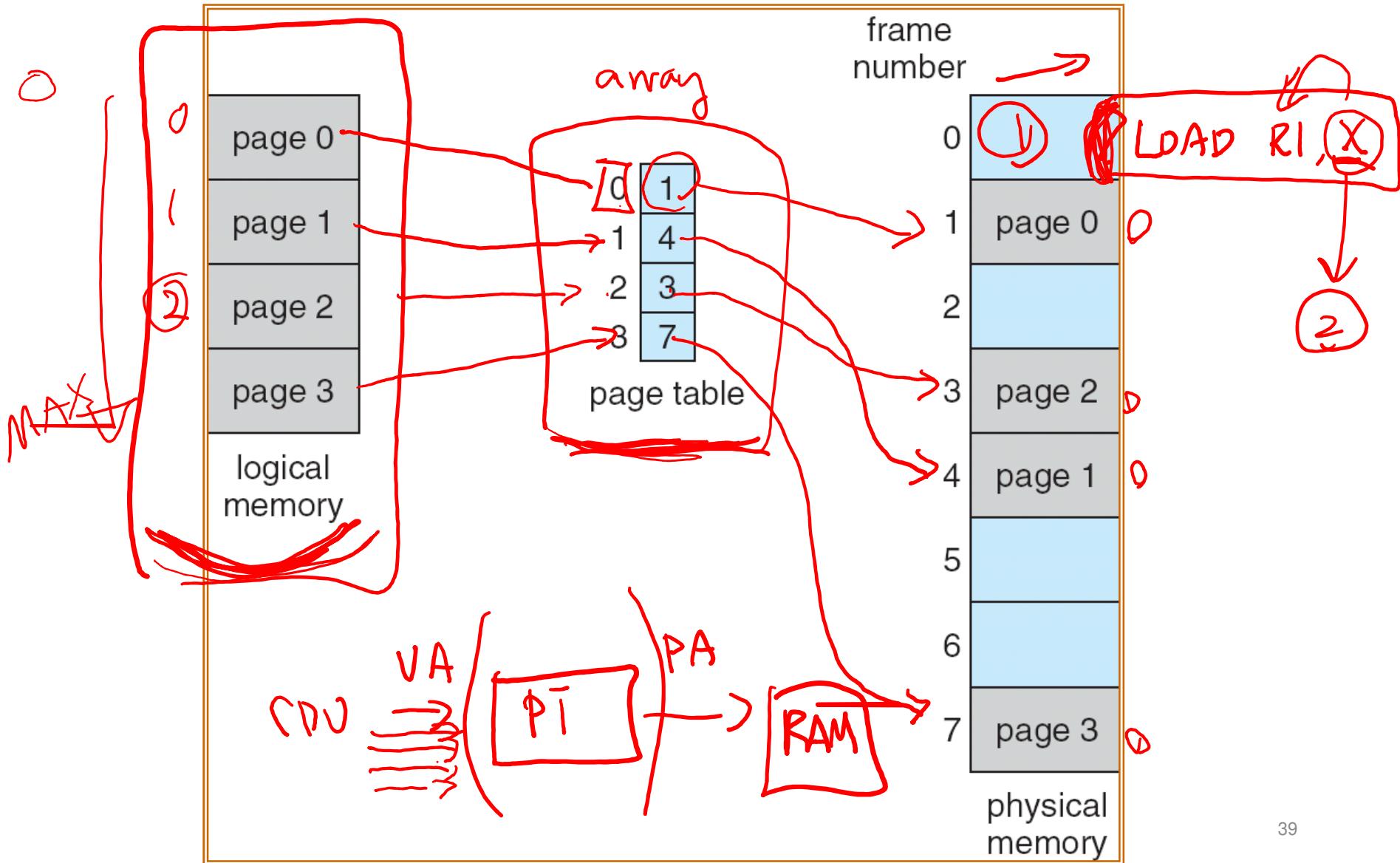


- Divide physical memory into fixed-sized blocks called frames (size is power of 2, between 512 bytes and 8192 bytes)
 $P\text{ size} = F\text{ size}$
- Divide logical memory into blocks of same size called pages.
 $4KB$ 4MB
- Keep track of all free frames
- To run a program of size n pages, need to find n free frames and load program
- Set up a page table to translate logical to physical addresses

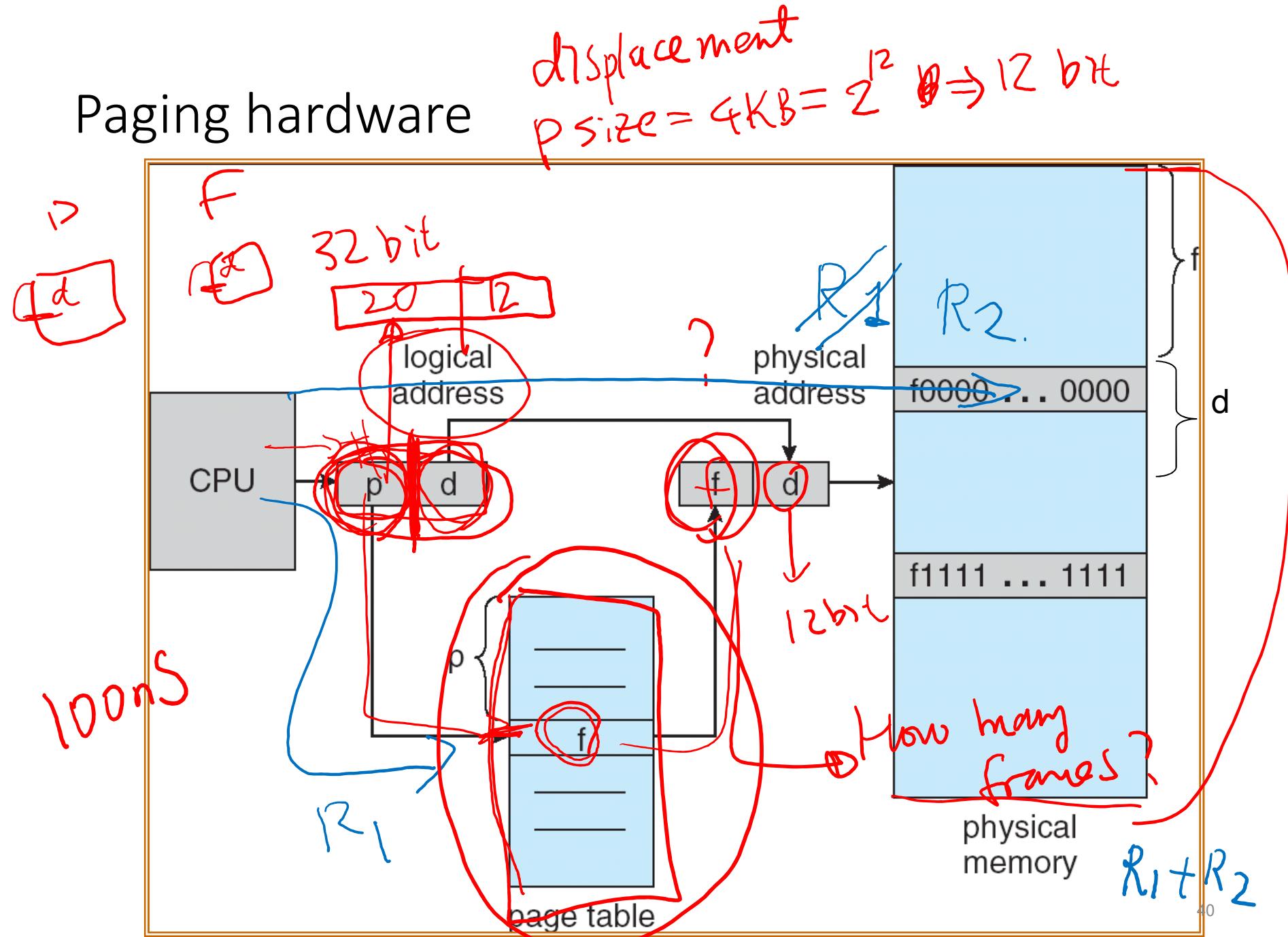
P [2] d

F [3] d

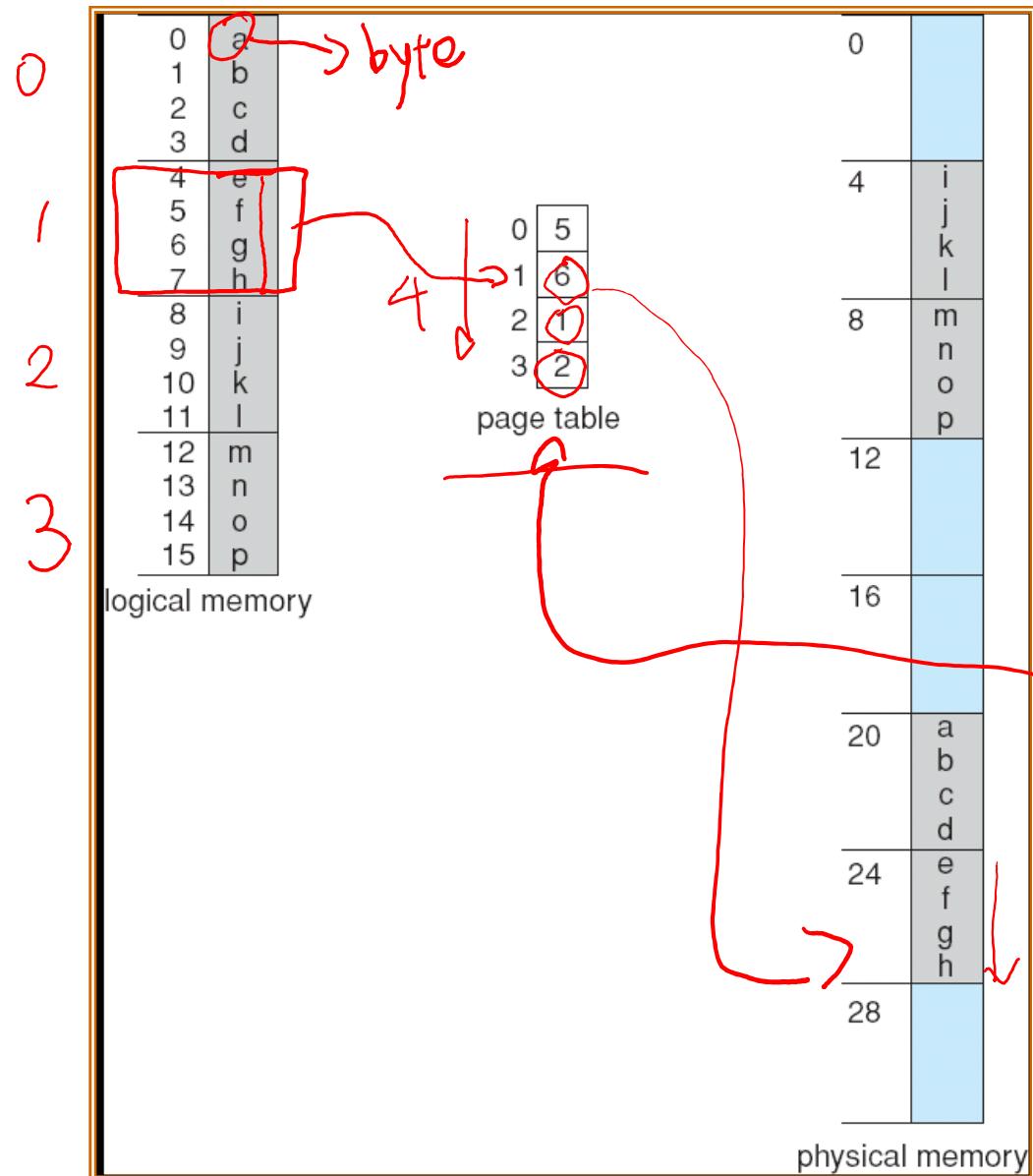
Paging Model of logical and physical memory



Paging hardware



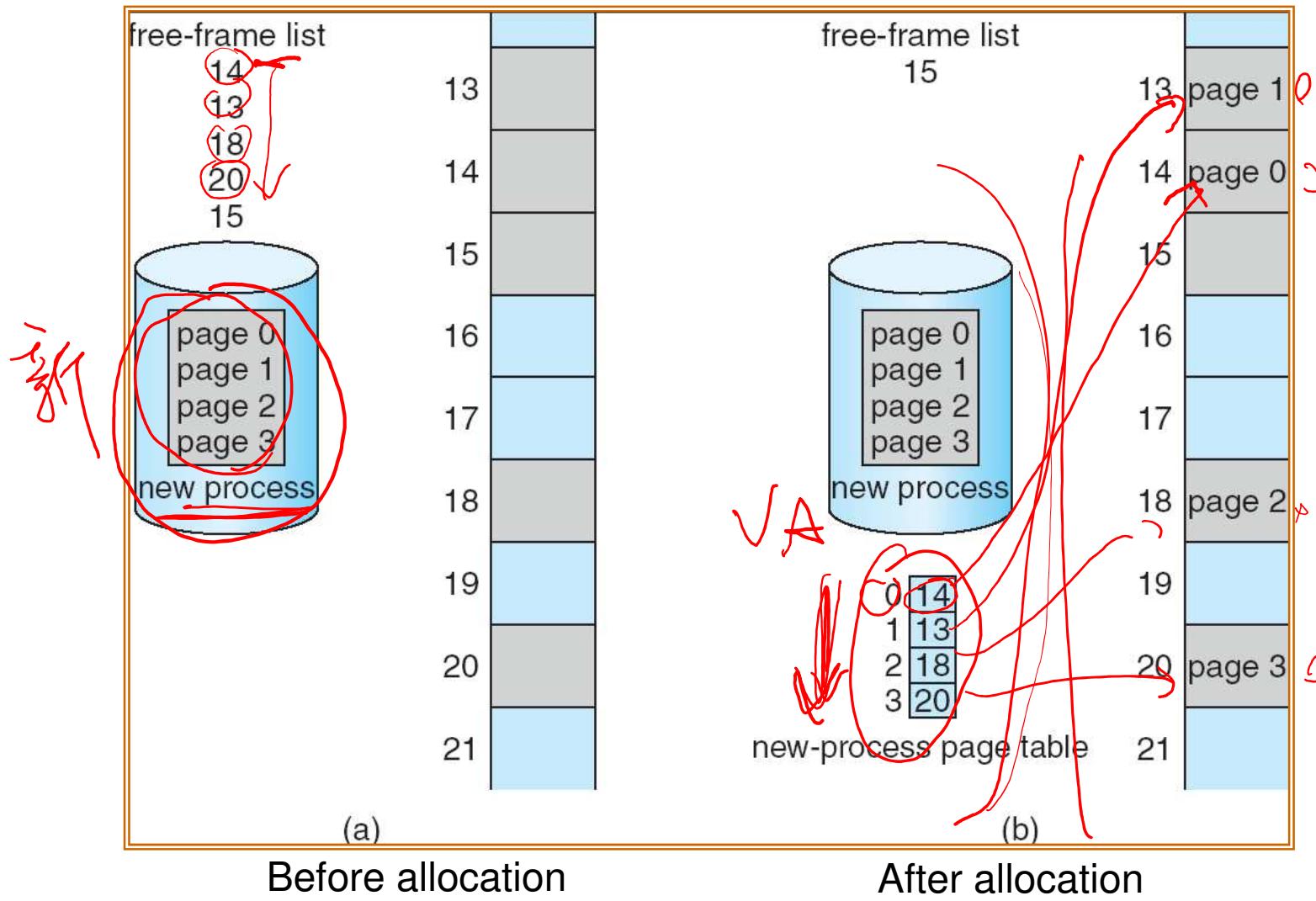
Paging Example



Q: how many bits are required to represent

1. The page number 2
2. The frame number 3
3. The displacement 2
4. A logical address 4
5. A physical address 5
6. A page-table entry? 3

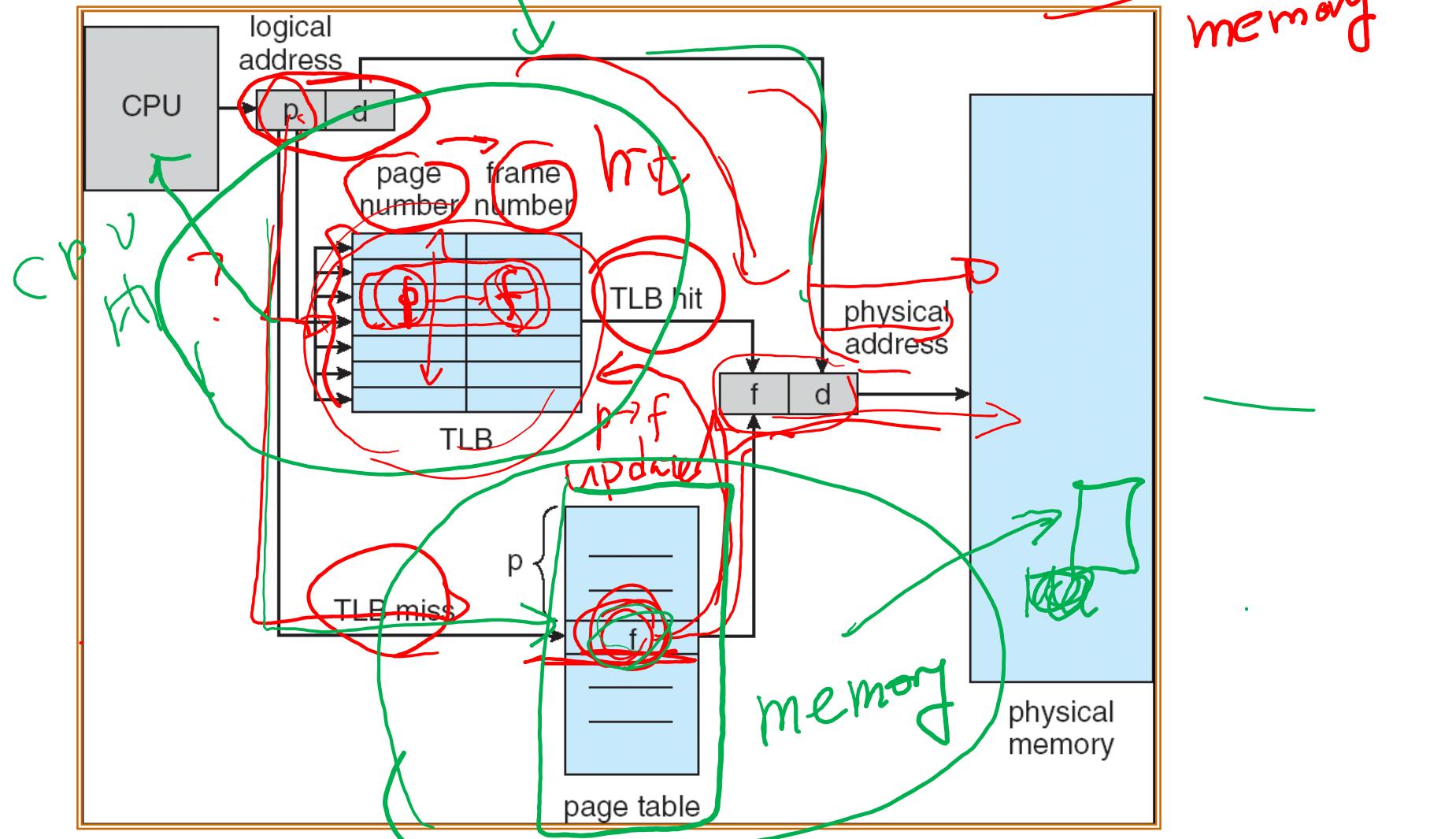
Free Frames



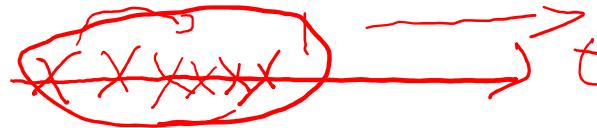
Implementation of Page Table

- Page table is kept in main memory (!) ~~process context~~
- A table per process
 - Page-table base register (PTBR) points to the page table
 - Page-table length register (PRLR) indicates size of the 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) Cache

Paging Hardware With TLB



All the steps are handled by HW



~ 128 entries

TLB hit ratio vs. Locality of Reference

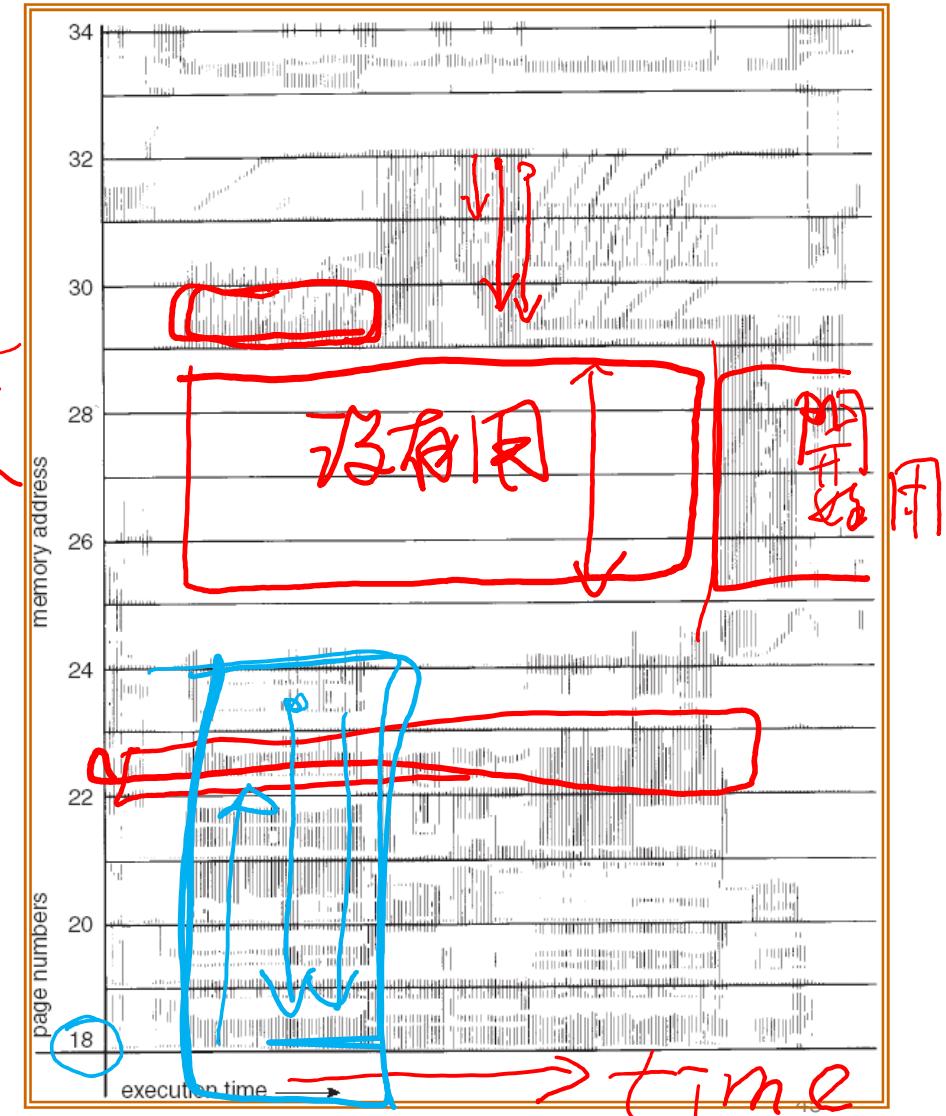
Locality

P

- TLB is small, usually holds 64~1024 entries
 - A replacement policy should be used.
 - E.g., random policy or LRU
 - Page references have **temporal** localities and **spatial** localities
- Important entries can be wired down (nailed)
 - E.g., kernel code

t
o →
P P

page
S. locality
min



Effective Access Time



TLB

- Associative Lookup = ϵ time unit
- Assume memory cycle time is $1 \mu\text{sec}$
- 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 $0 \sim 1$ \rightarrow hit
- Effective Access Time (EAT)

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

\rightarrow $2 \rightarrow 1$ $\underline{1 + \epsilon \sim 1}$

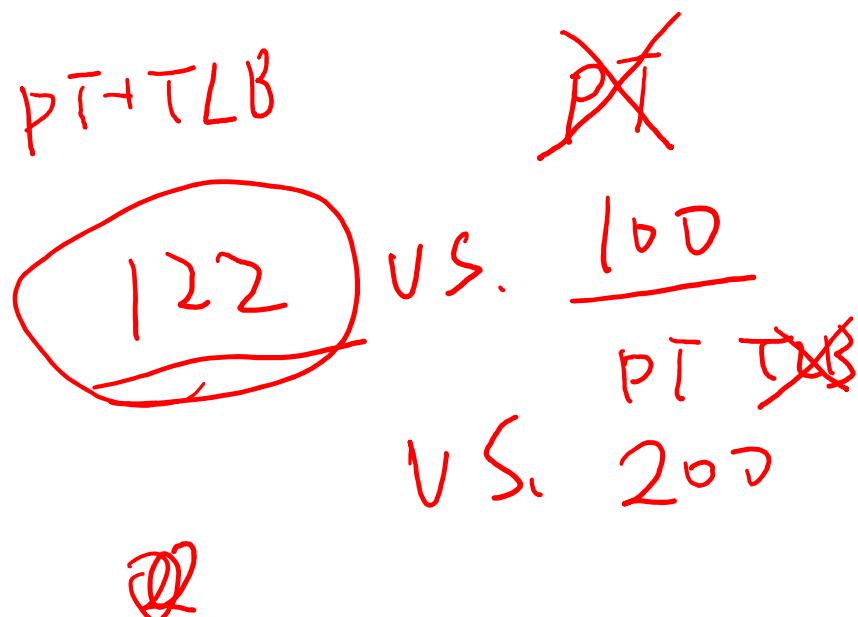
Let the TLB hit ratio be 98%

20ns to lookup the ~~TLB~~
100 ns to access memory

$$\text{TLB hit: } 20 + 100 = 120 \text{ ns}$$

$$\text{TLB miss: } 20 + 100 \text{ (page table)} + 100 \text{ (target address)} = 220 \text{ ns}$$

$$EAT = 0.98 * 120 + 0.02 * 220 = 122 \text{ ns}$$



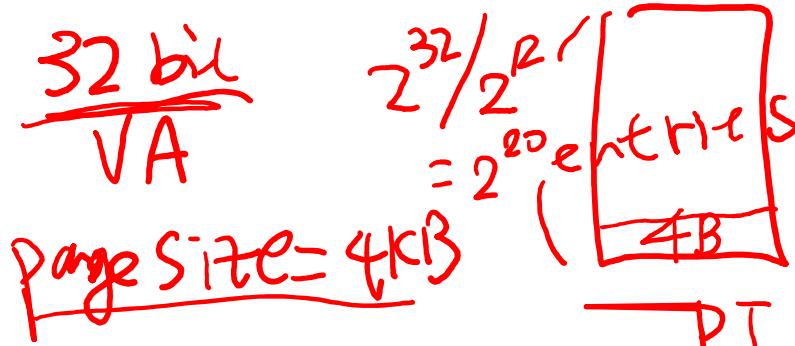
STRUCTURE OF THE PAGE

TABLE

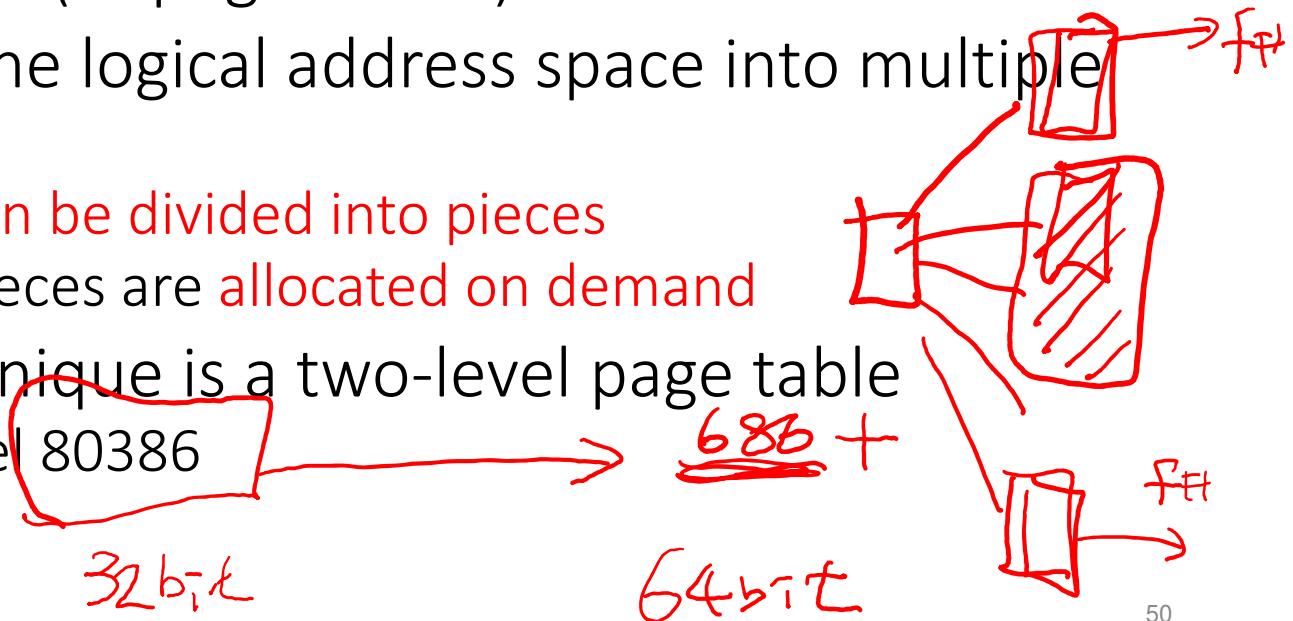
Structure of Page Table

- Hierarchical Paging
- Hashed Page Tables
- Inverted Page Tables

Multi-level Hierarchical Page Tables

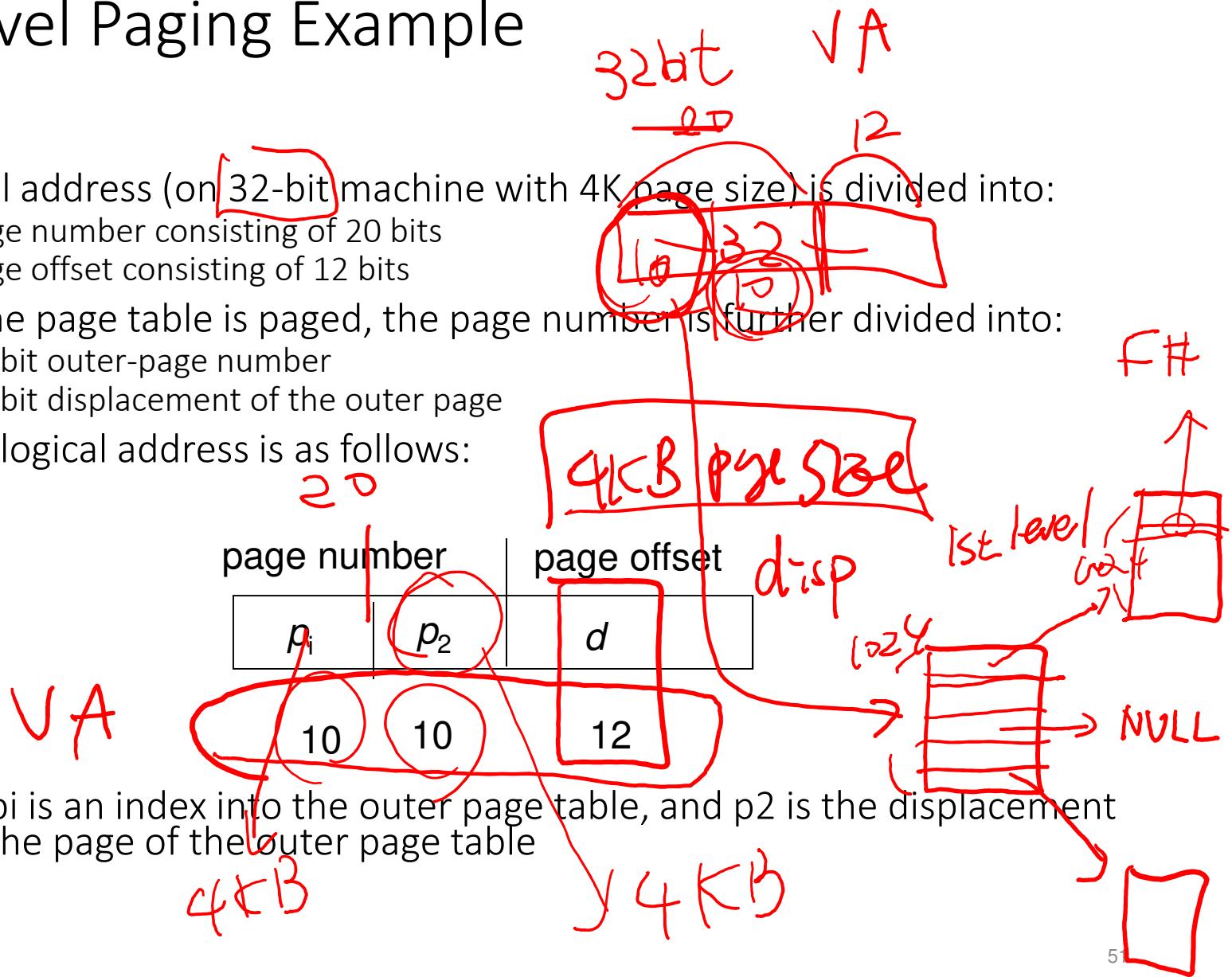


- A per-process page table could be very large and sparse
- Allocating **Large** and **contiguous** page tables for processes is inconvenient and may under-utilize memory space (of page tables)
$$2^{\text{20}} \times 4\text{KB} = 4\text{MB}$$
- Breaking up the logical address space into multiple page tables
 - Page table **can be divided into pieces**
 - Page table pieces are **allocated on demand**
- A simple technique is a two-level page table
 - Example: Intel 80386

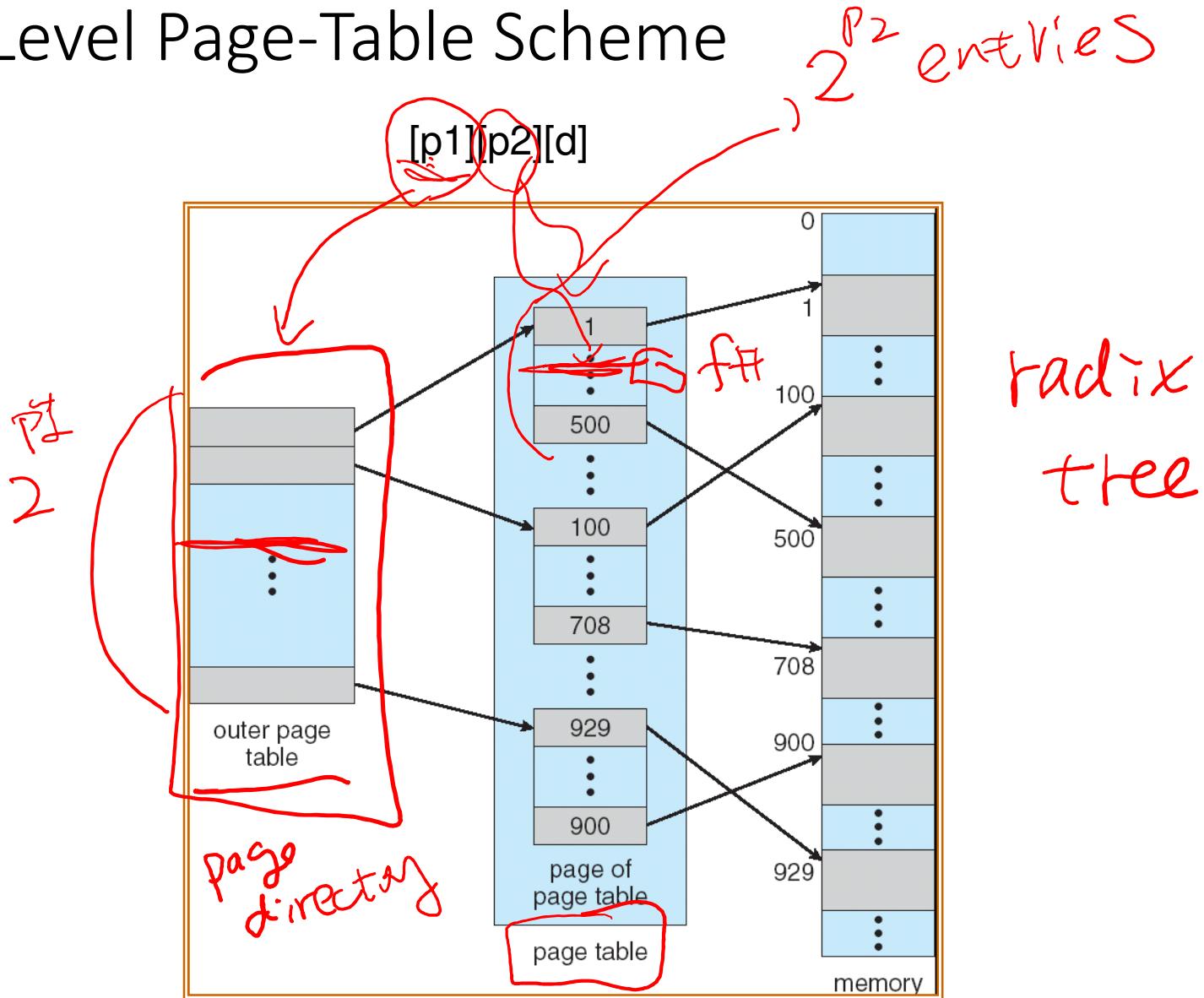


Two-Level Paging Example

- 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
- Since the page table is paged, the page number is further divided into:
 - a 10-bit outer-page number
 - a 10-bit displacement of the outer page
- Thus, a logical address is as follows:



Two-Level Page-Table Scheme

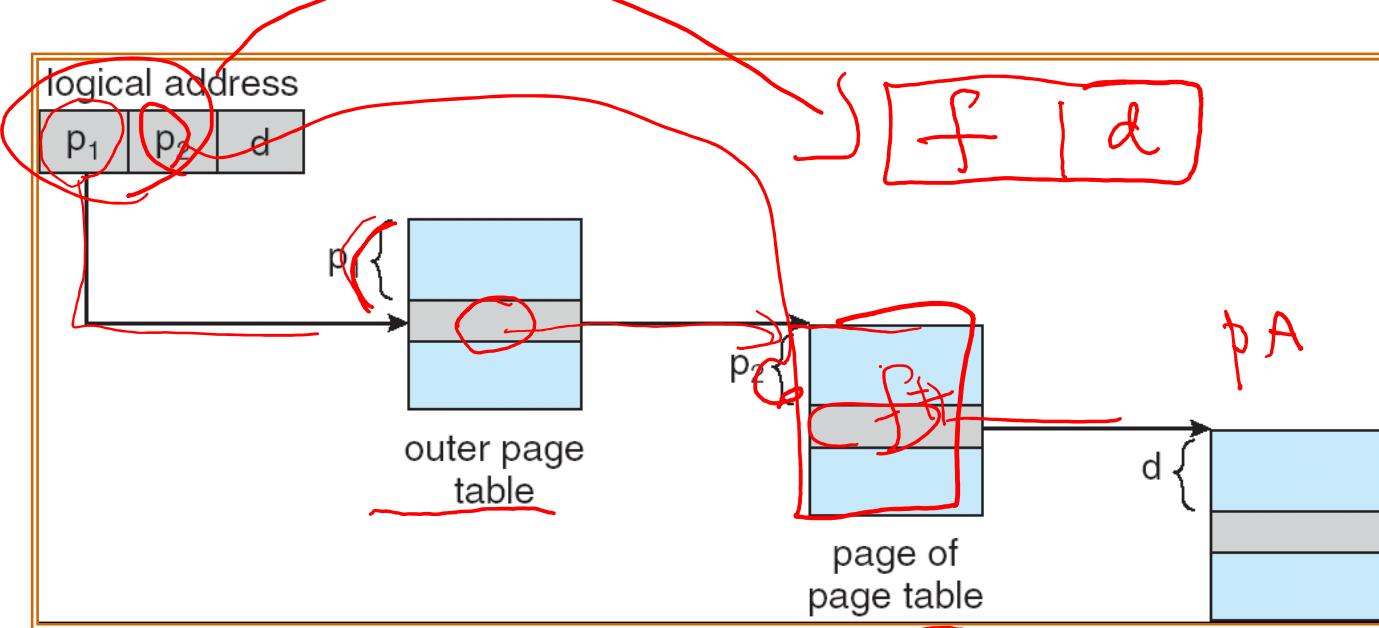


Address-Translation Scheme

IA 64

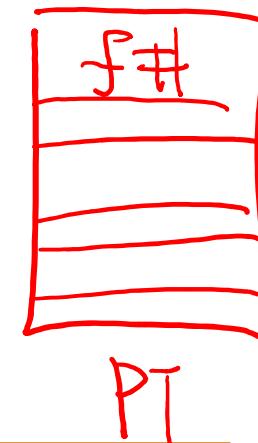
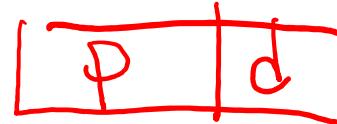
3-level

- Address-translation scheme for a two-level 32-bit paging architecture

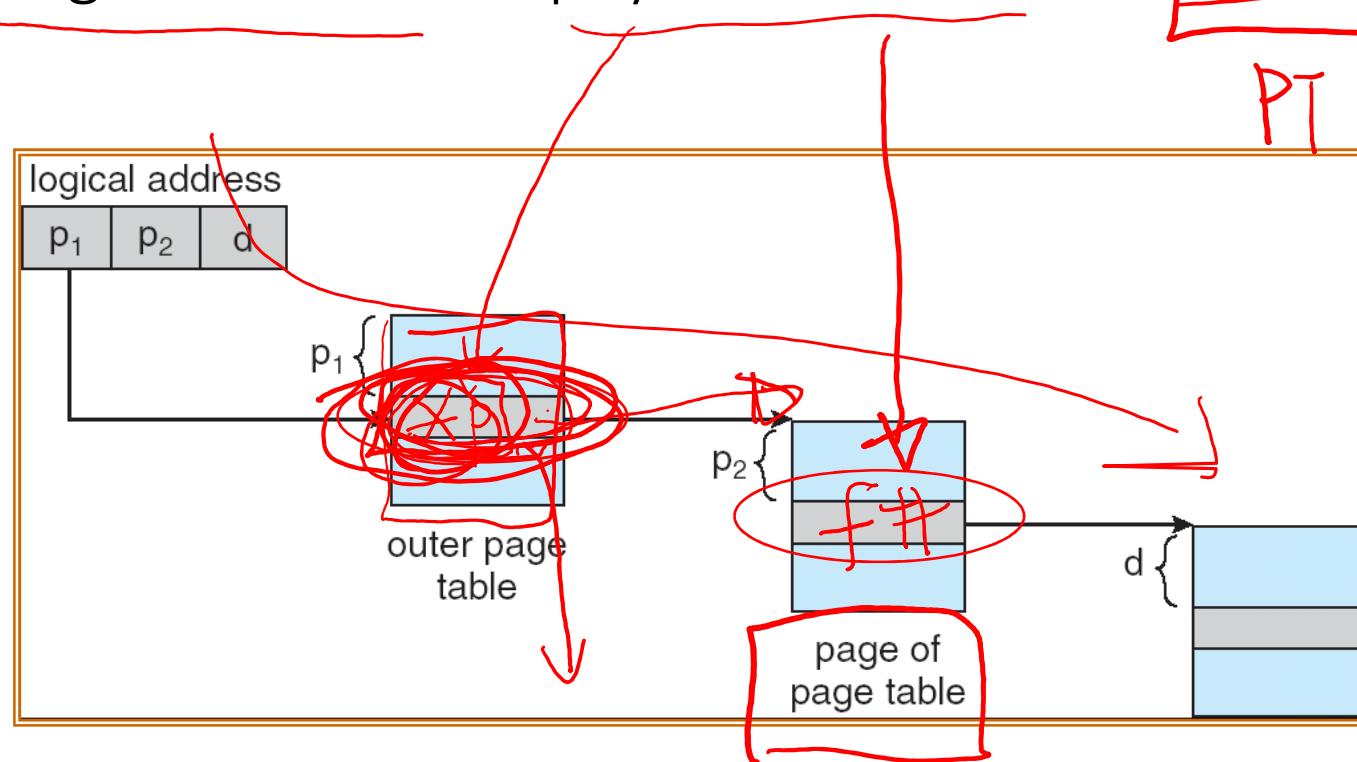


Pros: page tables need not to be contiguous, and need not all present in memory
Cons: multiple memory accesses on TLB miss. 7-level paging in UltraSparc 64

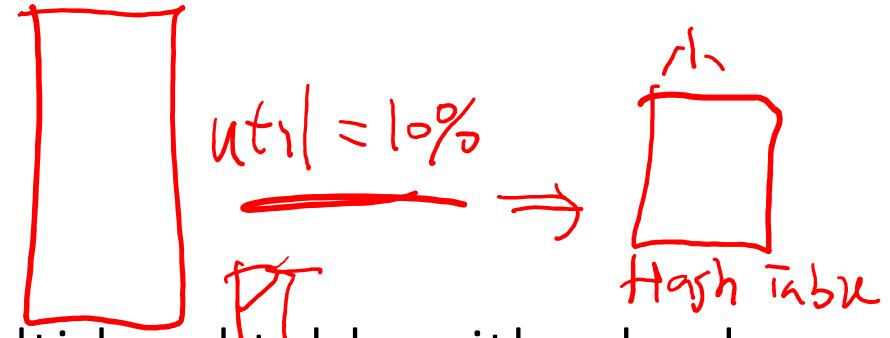
Think about it...



- Logical address or physical address?

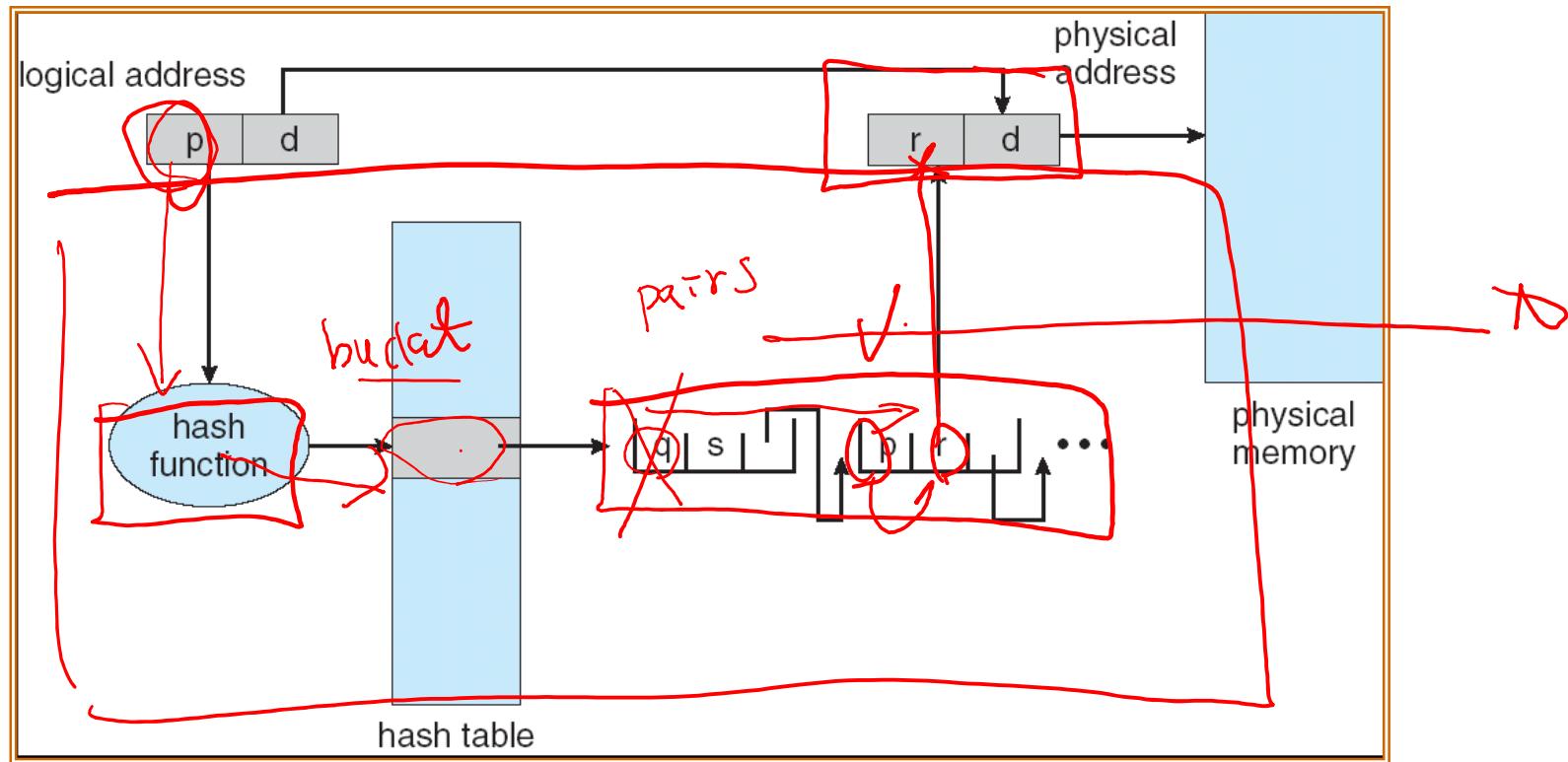


Hashed Page Tables



- Replace the radix-based multi-level table with a hash table
- Common in address spaces > 32 bits
 - The logical address footprint of a process << the logical address space
- The page number is hashed into a page table. This page table contains a chain of elements hashing to the same location
- 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



How many extra memory references are needed on TLB miss
is related to the quality of the hash function

Hashed Page Table

字节 313

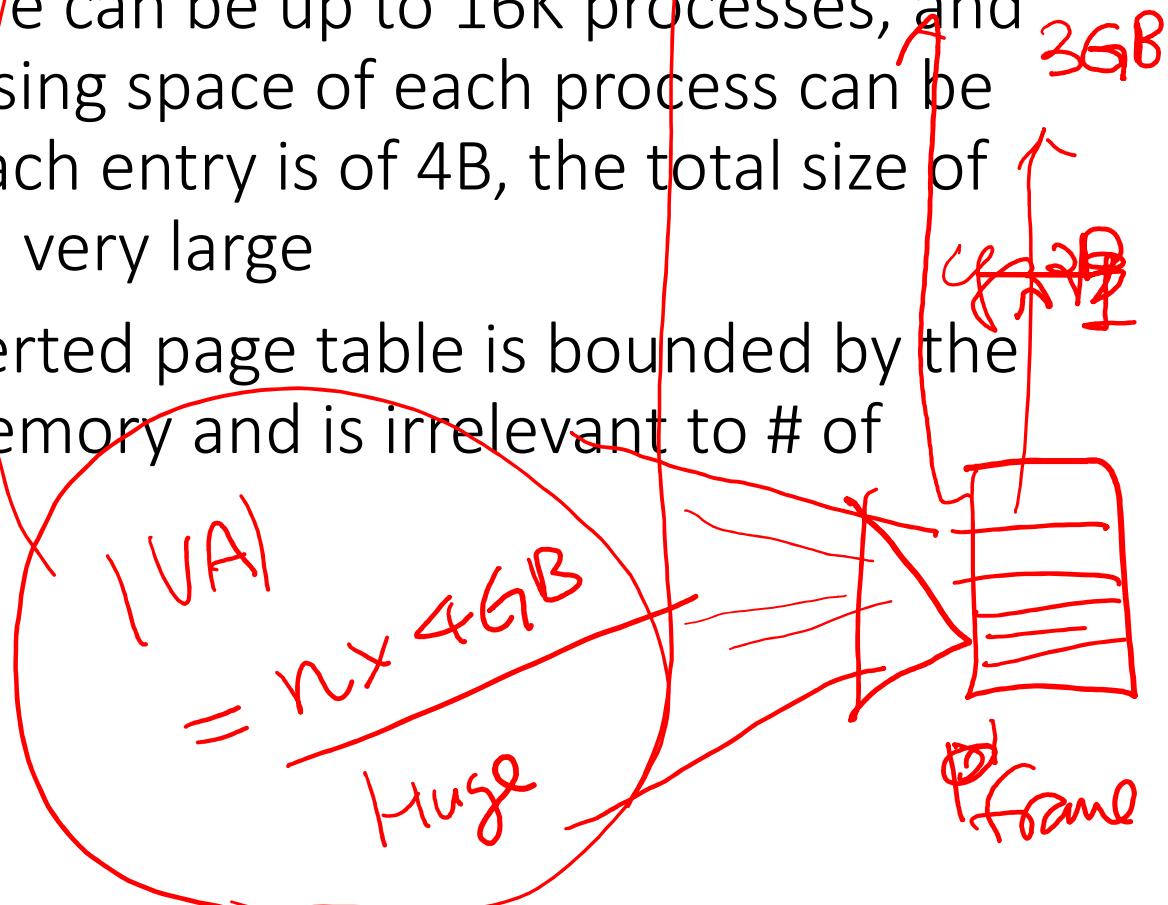
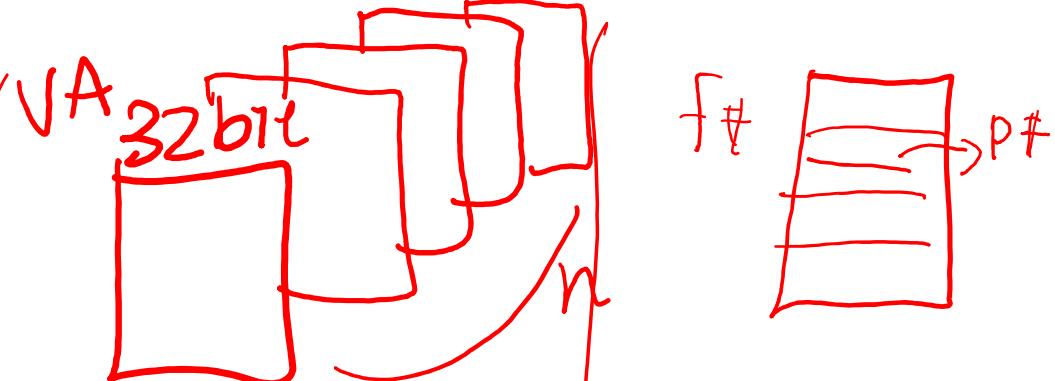
- (Forward) page tables ($p \rightarrow f$) are commonly implemented using multi-level tables, not hash tables
- Inverted page tables are, however, commonly implemented using hash tables

POW6R 603
Co4

IBM

Inverted Page Table

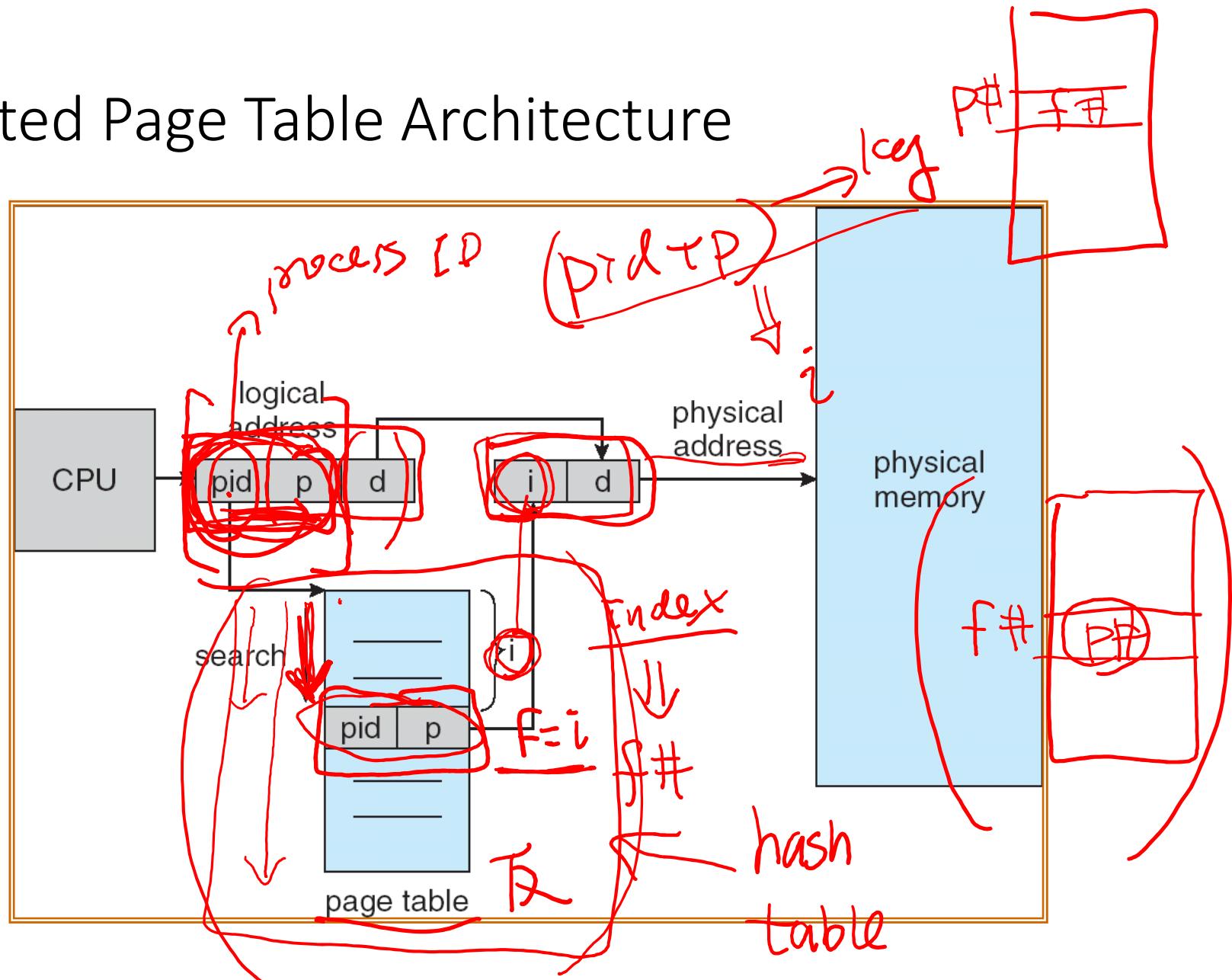
- Consider that there can be up to 16K processes, and the logical addressing space of each process can be of 2²⁰ pages. If each entry is of 4B, the total size of page tables is: very large
- The size of an inverted page table is bounded by the size of physical memory and is irrelevant to # of processes



Inverted Page Table

- One entry for each real page (frame) of memory
 - Entry index number is the frame number
- Each table entry stores a page number and the process ID owns that page
- One global page table shared among all processes
- On memory reference, find the table entry that stores the current process ID and the referenced page number
 - Use hash table to limit the search to one — or at most a few — page-table entries
- Example: PowerPC 603

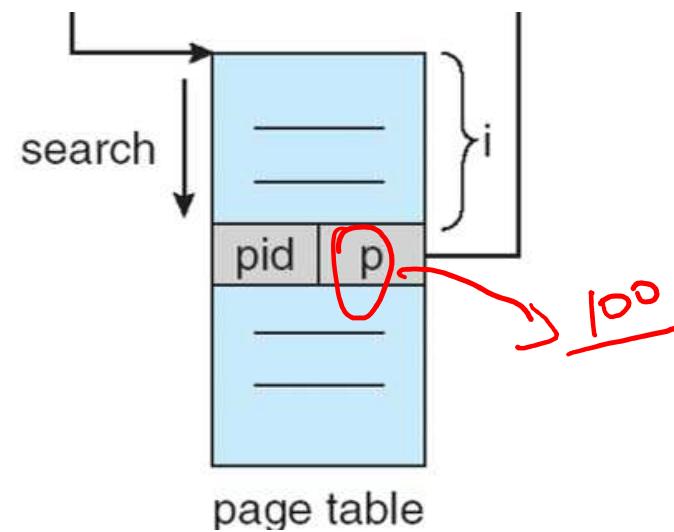
Inverted Page Table Architecture



Usually implemented as a hash table for fast search

Think about it...

- Why pid is necessary in inverted page tables?



Design Considerations of Paging

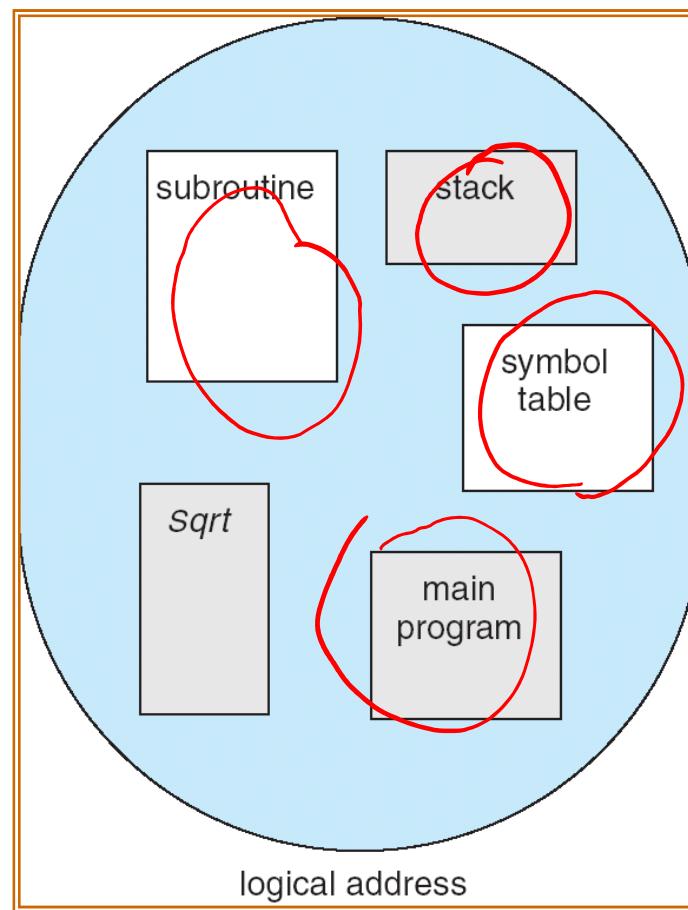
- Increased memory references
 - Solution: Use TLB
- Space requirement of the page table (in main memory)
- Multilevel page tables
 - Break the entire table into small pieces, allocate on demand
- Hash page tables
 - Small hash table with overflow handling (linked lists)
- Inverted page tables
 - Table size is bounded by the physical memory size and is independent of the total number of active processes

SEGMENTATION

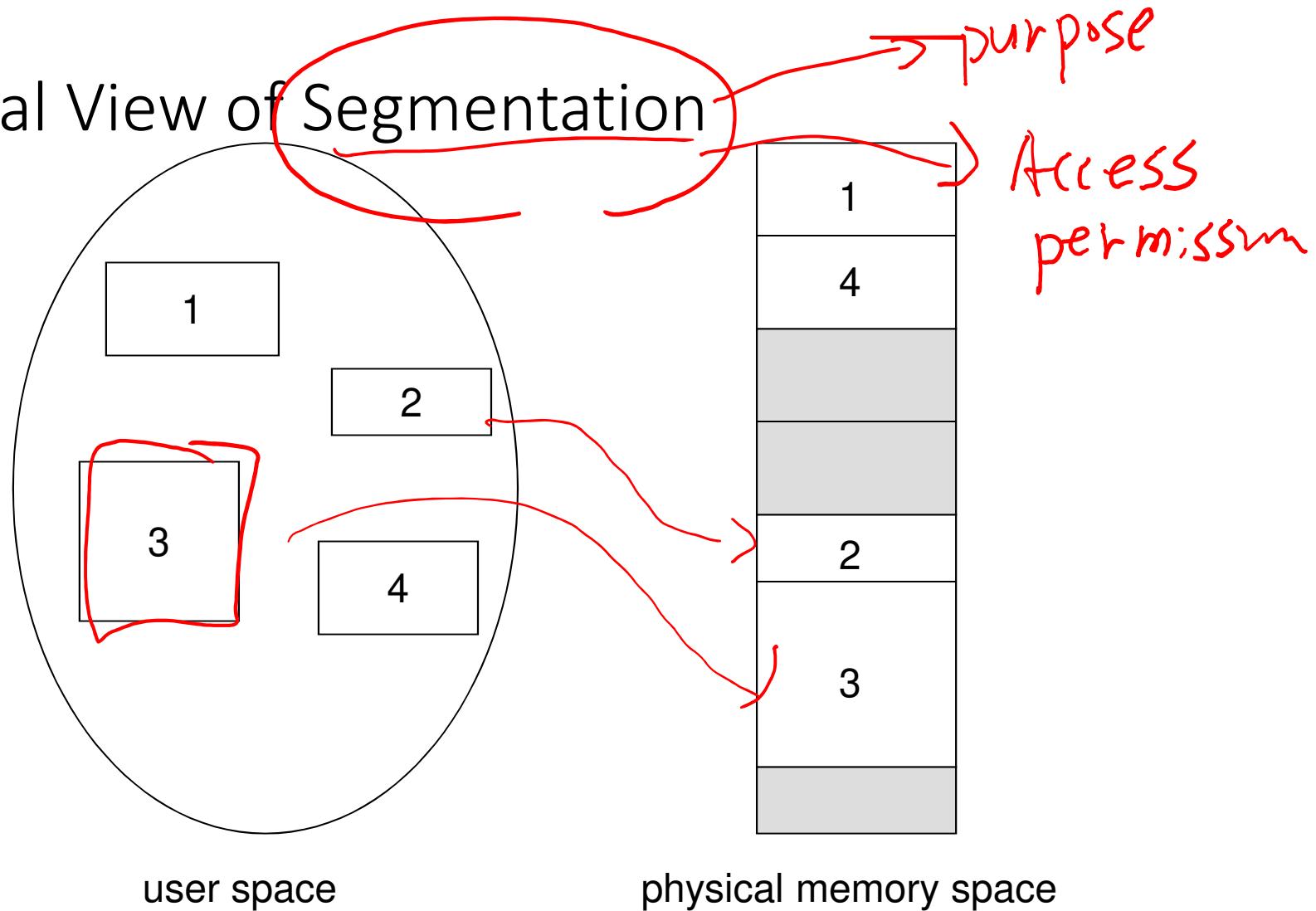
Segmentation

- 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,
 - method,
 - object,
 - local variables, global variables,
 - common block,
 - stack,
 - symbol table, arrays

User's View of a Program



Logical View of Segmentation

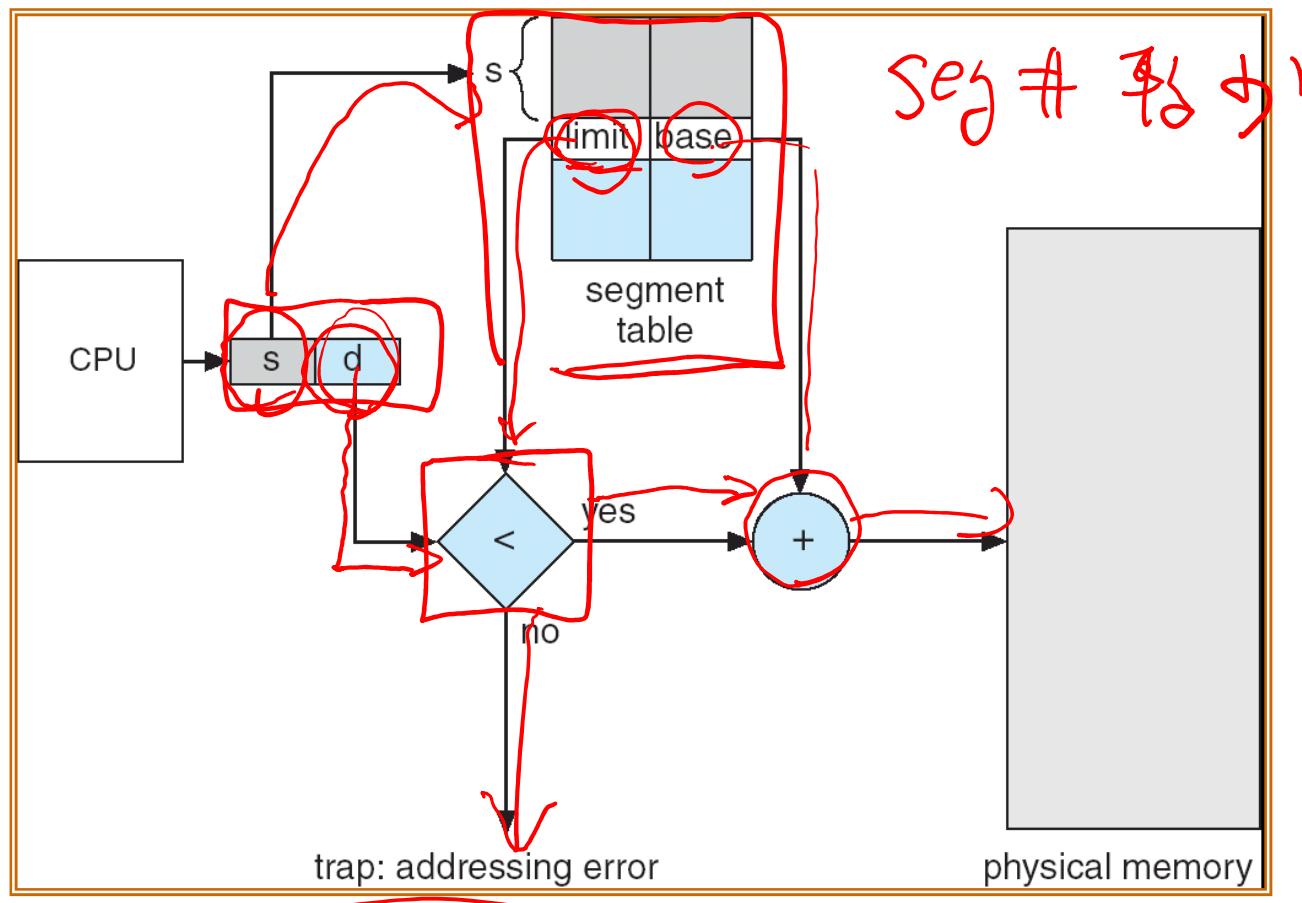


User programs do not know where in physical memory a segment is placed.

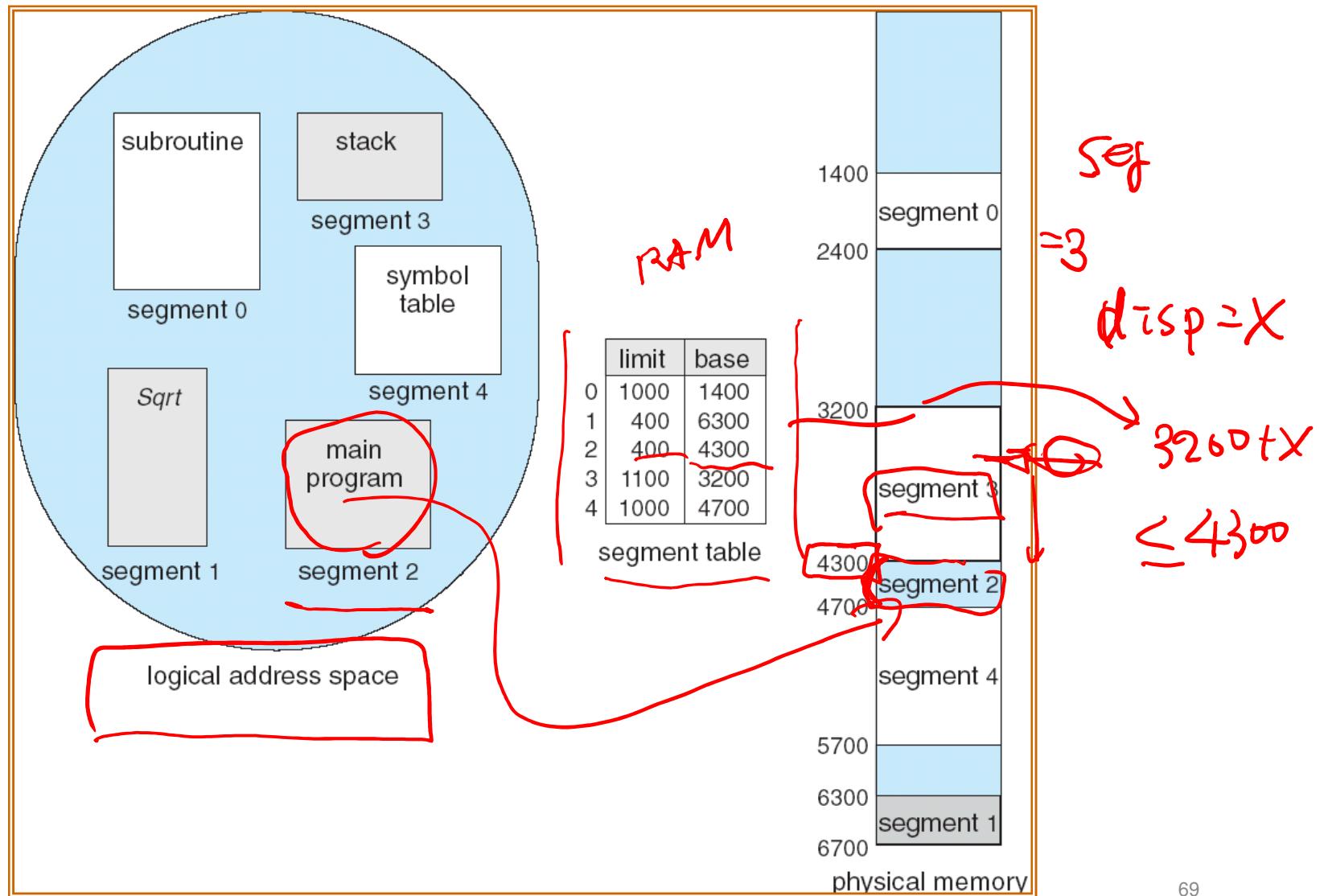
Hardware

- 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 hardware (MMU)



Example of Segmentation



Segment Protection

Code RV
text → X ✓
 W X

- Segments storing execution code normally do not allow modification
 - Allow read, execute but not write
 - Avoid self-modifying malicious code, such as polymorphic viruses
- Stack segments normally do not allow execution
 - Allow read, write but not execute
 - Avoid malicious code injection

funcm stack frame args local var ret
R, W, X
✓ ✓ — X

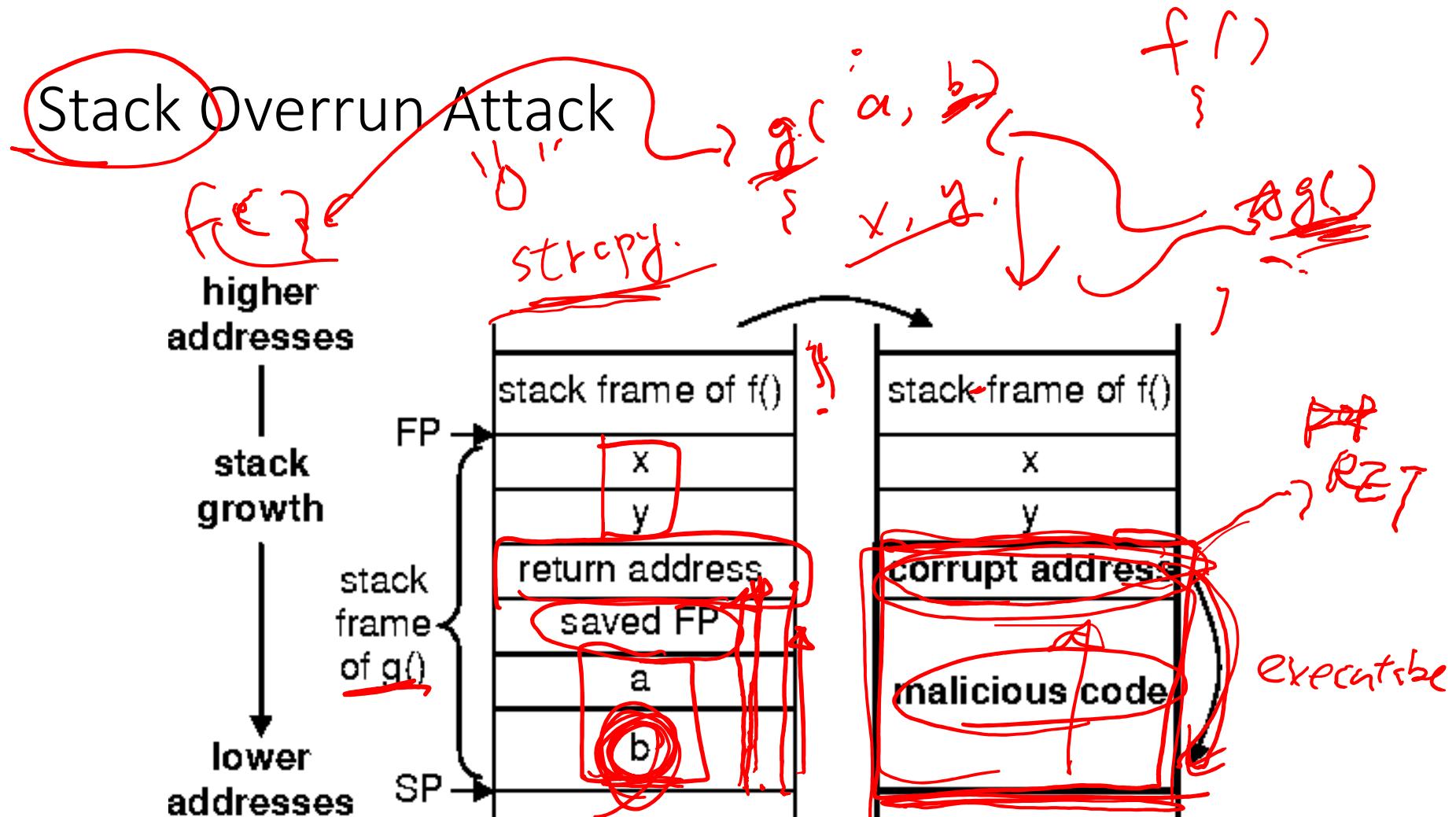


Figure 4: Buffer overrun exploit

- What is the result of executing the program below?

```
#include<stdio.h>
#include<stdlib.h>

int main()
{
    char *p = "hello world";
    p[0] = 'H';
}
```

Annotations on the code:

- The string "hello world" is highlighted with a red box and labeled literal.
- The assignment `p[0] = 'H';` is crossed out with a large red X.
- To the right of the original assignment, there is a red circle with the letters RO and a red square with the letters Sq.
- Below the crossed-out assignment, the word constant is written.
- At the bottom right, there is a red W above a red X.

Hardware Support for Segmentation

- Provided by MMU
- Segment Limit
 - Segment base register
 - Segment limit register
- Segment Protection
 - With each entry in segment table associate:
 - **Permissions** of read, write, and/or execute
 - **Access privileges** of user mode or kernel mode

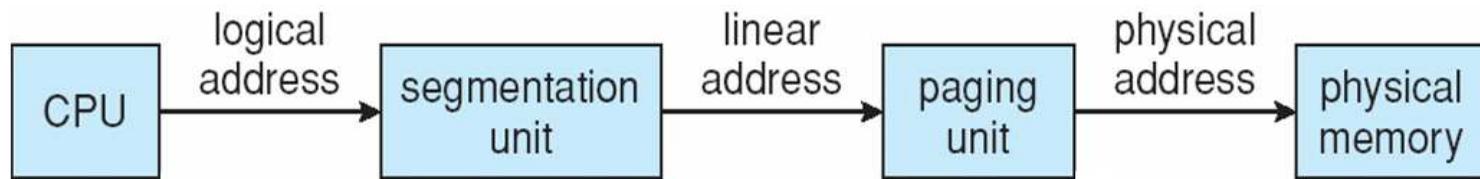
Paging and Segmentation

- Comparison
- Paging solves the external *fragmentation* problem at the process level
- Segmentation provides memory protection according to the purposes (r,w,x) of memory regions
- Combination
- Segmented paging (more common)
- Paged segmentation

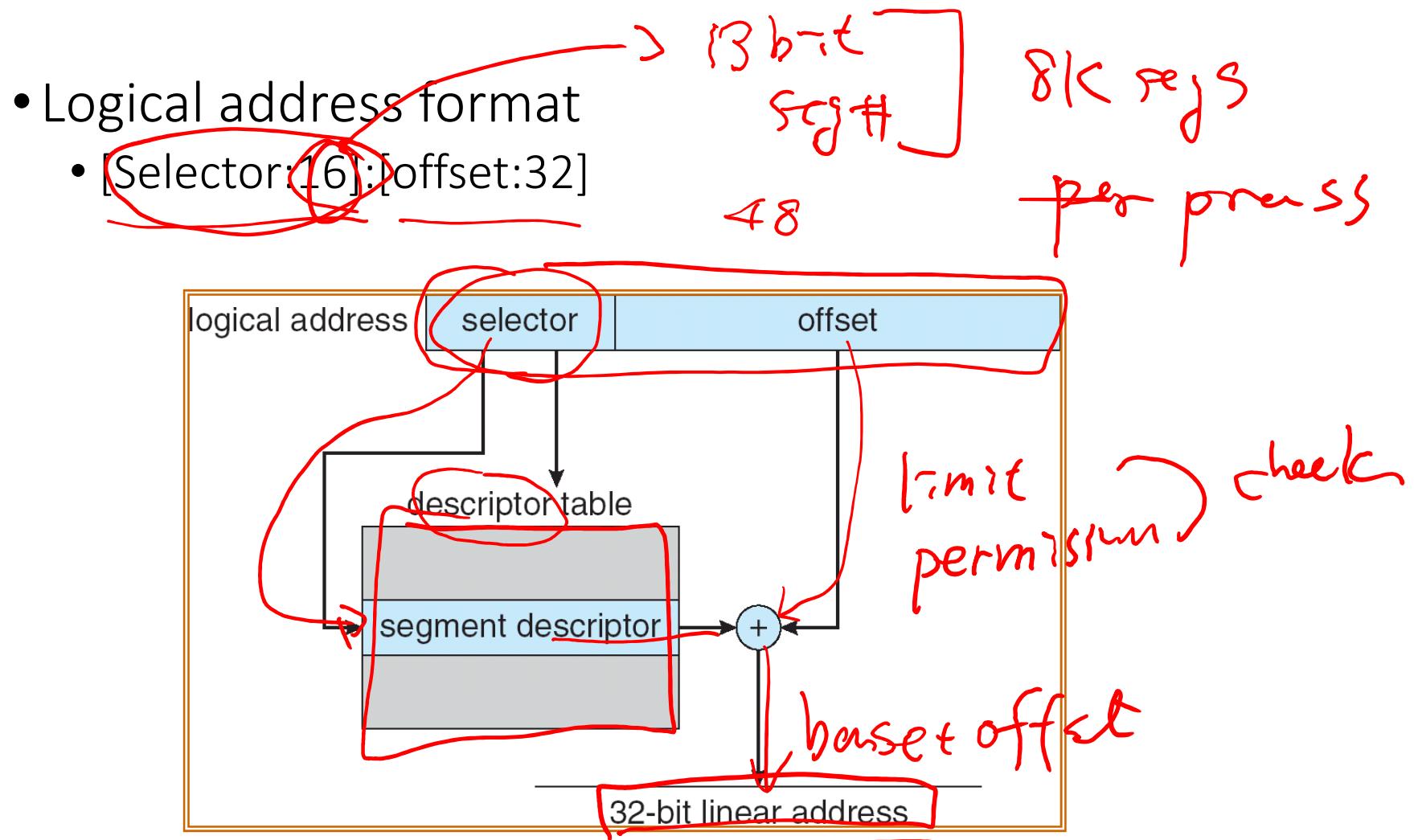
EXAMPLE: THE INTEL 80386+

Segmented Paging

- Intel 80386 *32 bit*
 - Segmentation
 - 2-level paging

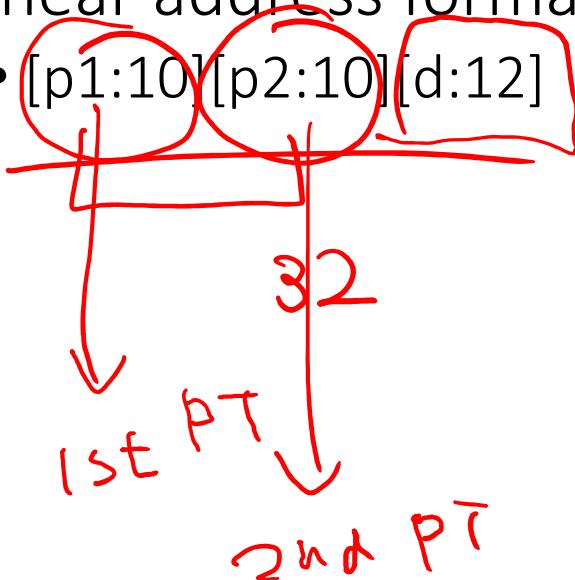


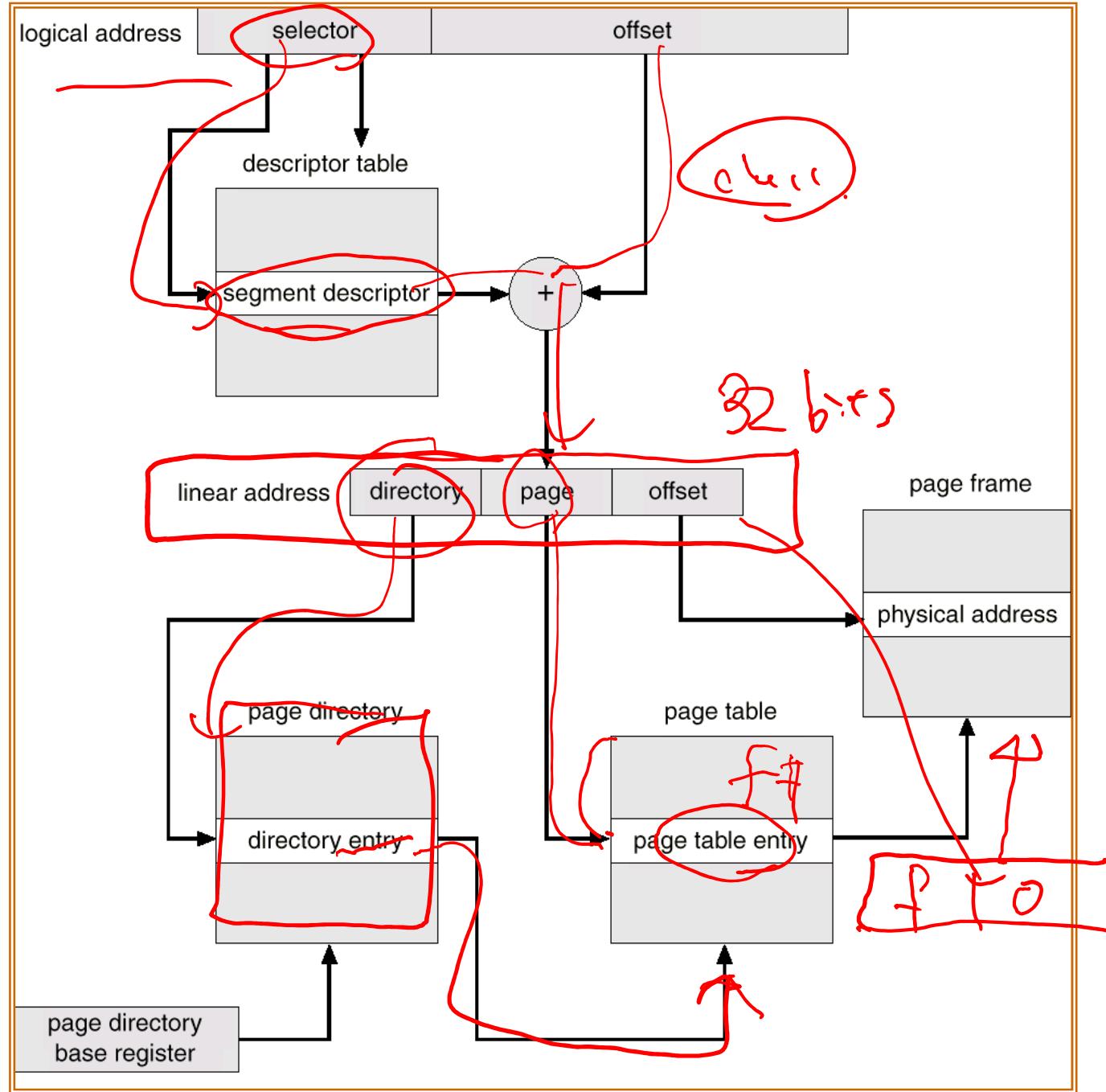
Logical Address to Linear Address



Linear Address to Physical Address

- Once logical address is translated into linear address, it is handled by paging
- Linear address format
 - $[p_1:10][p_2:10][d:12]$





Linux Segmentation on Intel 80386

- Segmentation ✓ and Paging ✓ are somewhat redundant
 - RW protection → supported by both
 - Privilege → supported by both
 - Executable → segmentation only
- RISC architectures often have limited support for segmentation
- Therefore, Linux use segmentation only when required by x86 architecture

Linux Segmentation on Intel 80386

- Uses minimal segmentation to keep memory management implementation more portable
- Uses 6 global segments: (80386 has many)
 - Kernel code
 - Kernel data
 - User code
 - User data
 - The default LDT (usually not used)
 - Per-core
 - Task-state (TSS, used to switch from user mode to kernel mode)
- Uses 2 protection levels: (80386 has 4)
 - Kernel mode
 - User mode

priv

End of Chapter 8