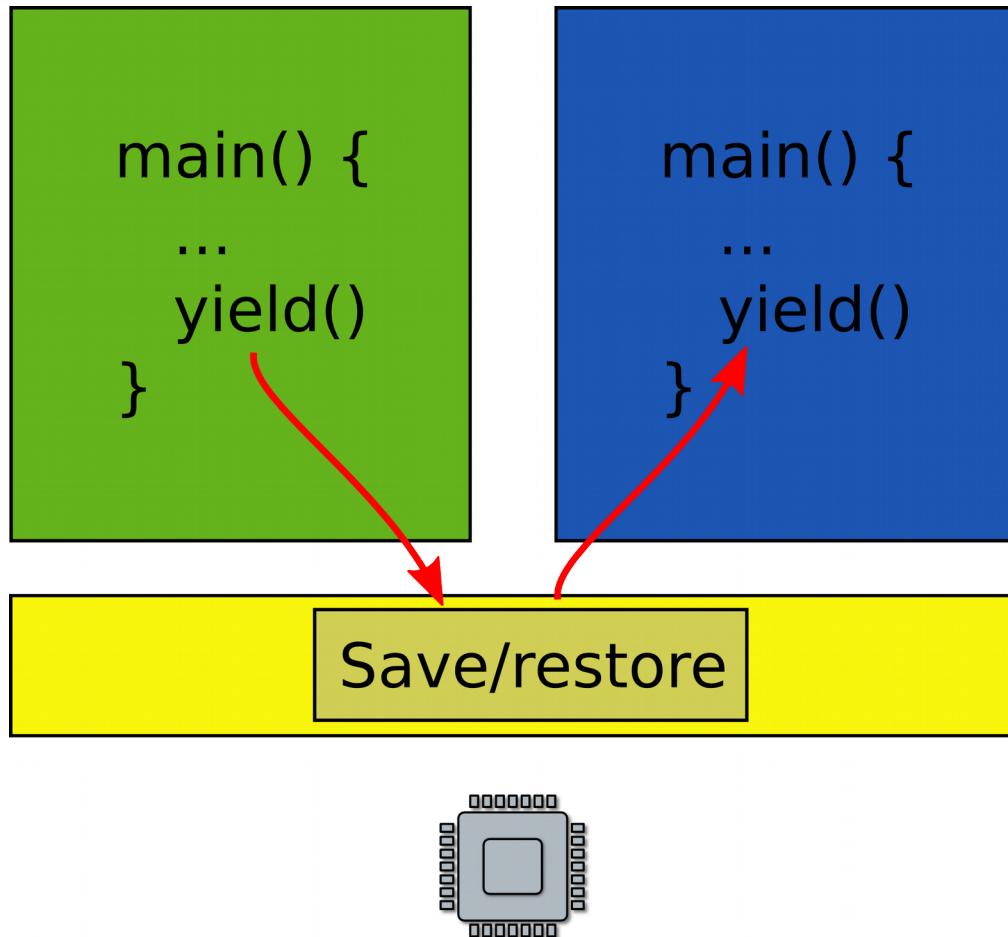


# 143A: Principles of Operating Systems

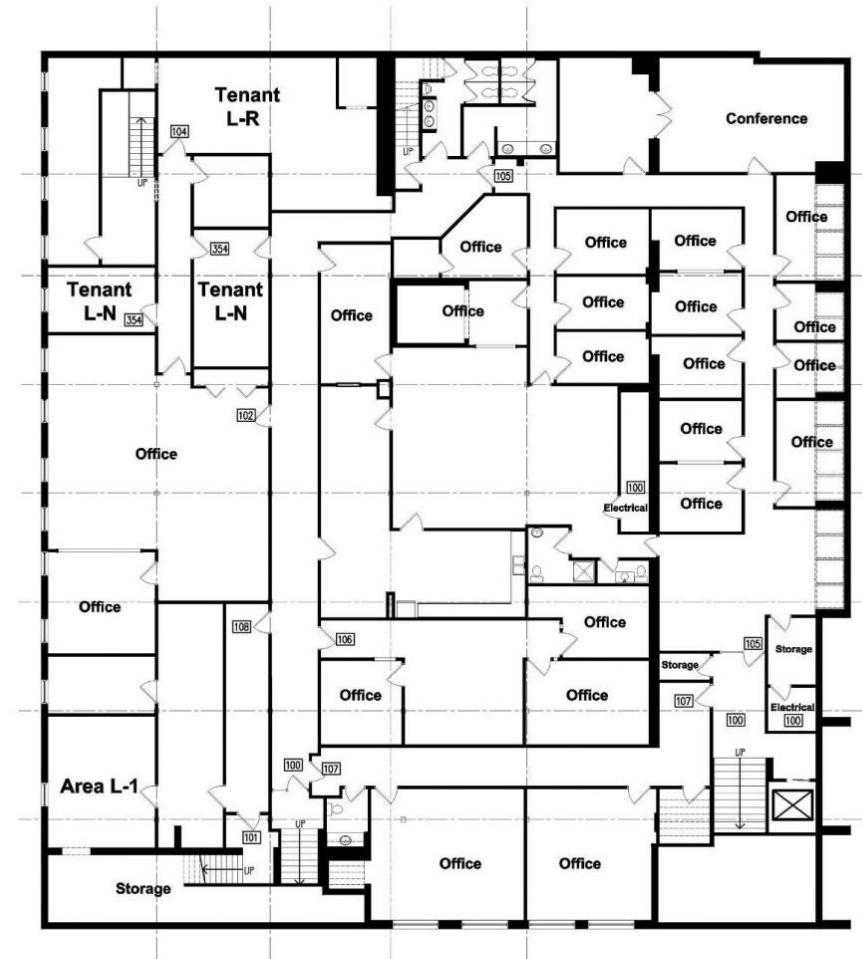
## Lecture 5: Address translation

Anton Burtsev  
October, 2018

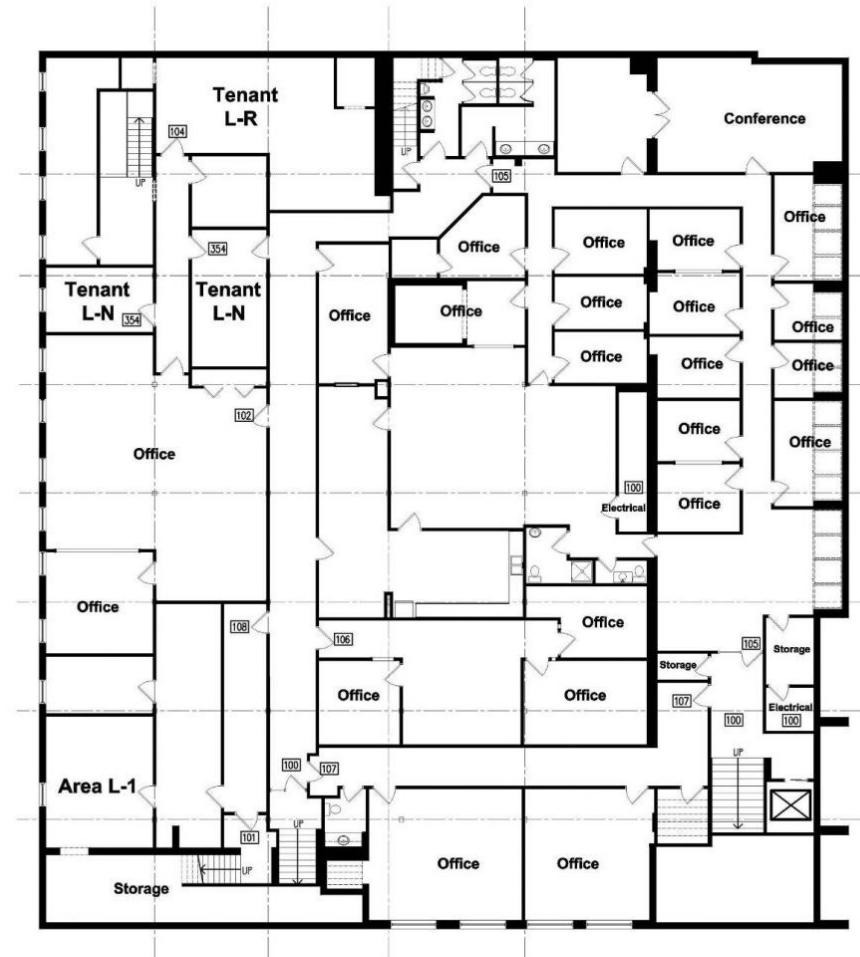
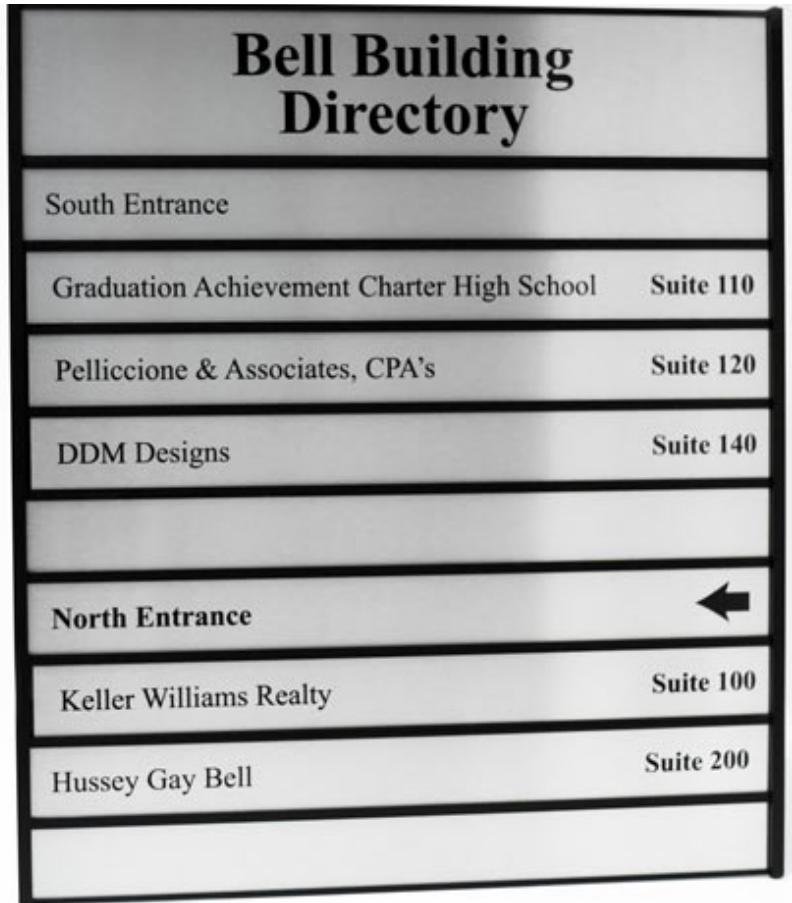
# Two programs one memory



# Or more like renting a set of rooms in an office building



# Or more like renting a set of rooms in an office building

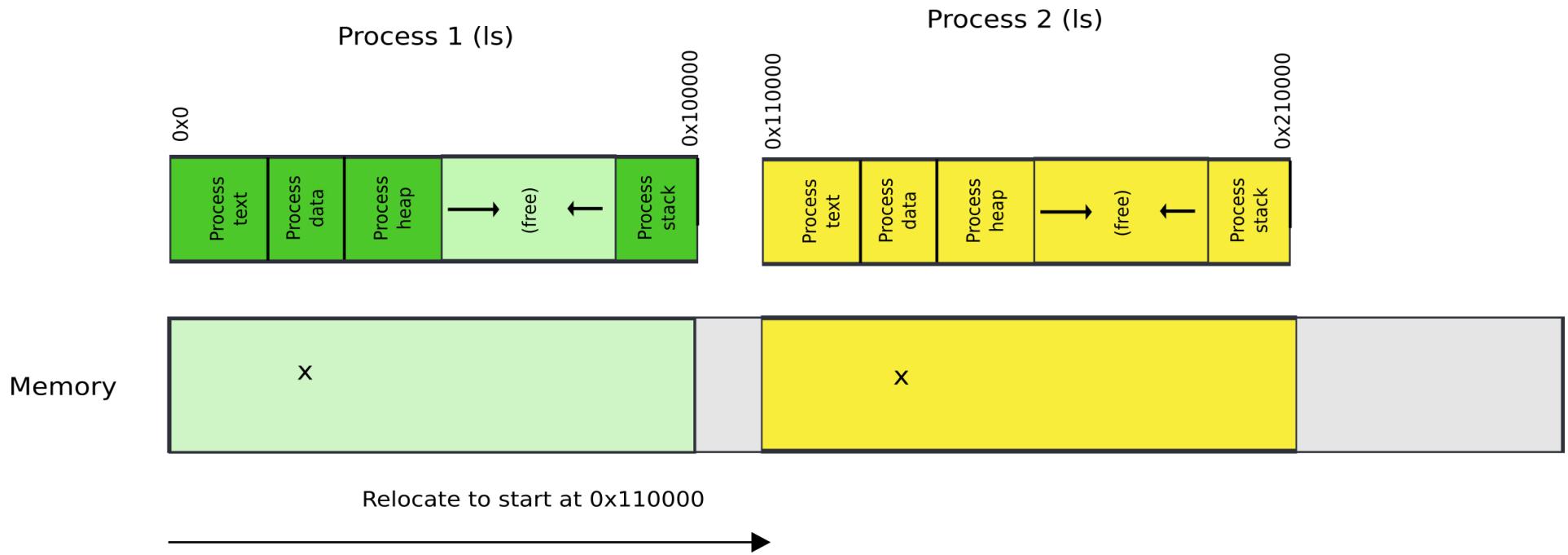


13852 Sq.Ft.  
Lower Level  
1/16" = 1'-0"

# Relocation

- One way to achieve this is to relocate program at different addresses
  - Remember relocation (from linking and loading)

# Relocate binaries to work at different addresses



- One way to achieve this is to relocate program at different addresses
- What is the problem?

- One way to achieve this is to relocate program at different addresses
- What is the problem?
  - No isolation

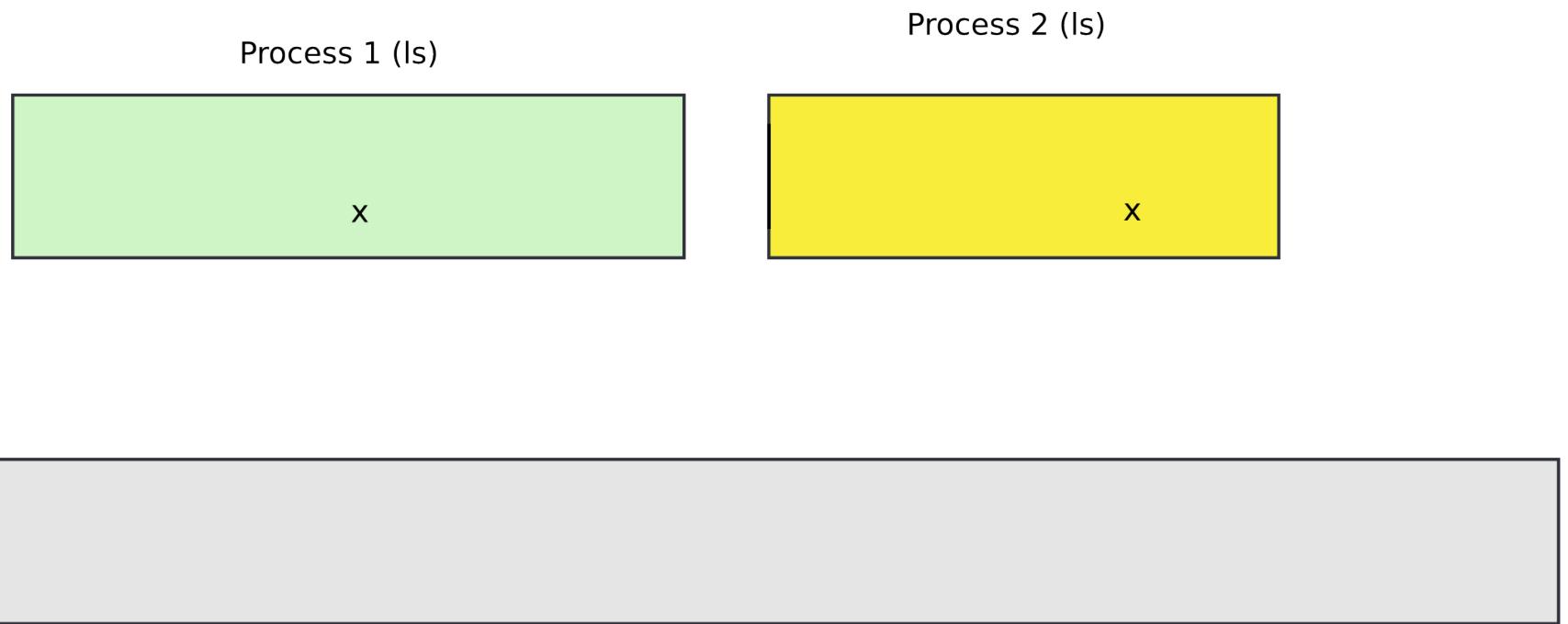
- Another way is to ask for hardware support

This is called segmentation

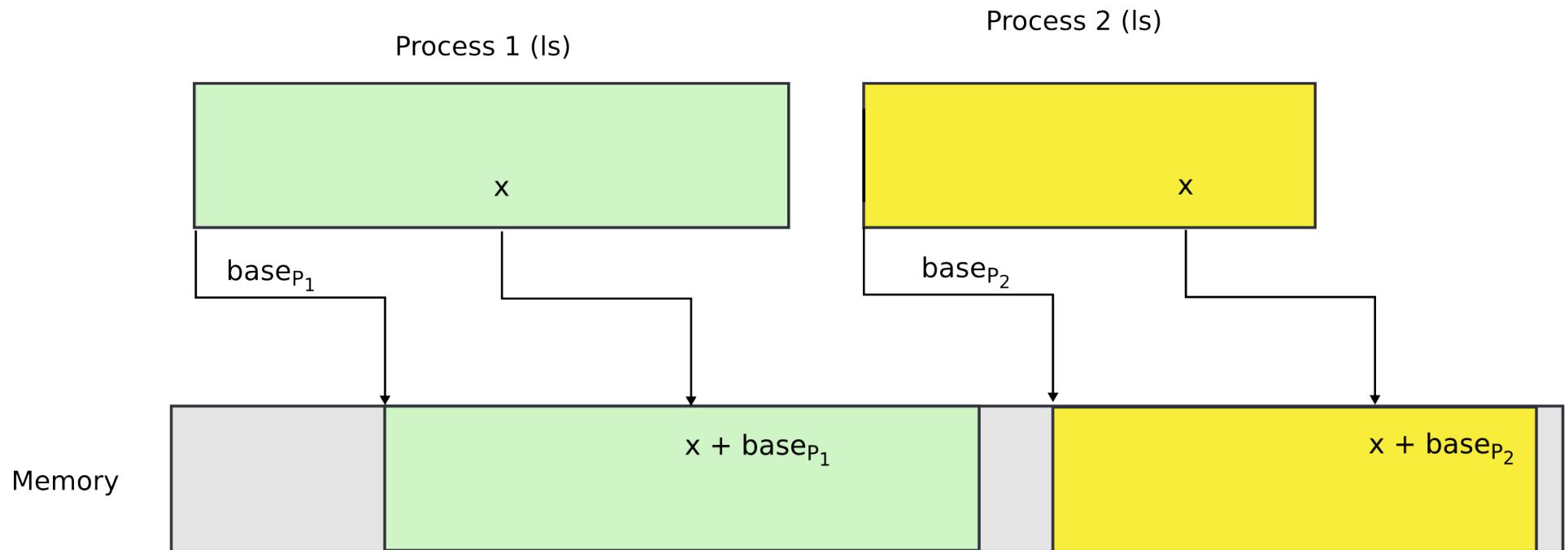
# What are we aiming for?

- Illusion of a private address space
  - Identical copy of an address space in multiple programs
    - Remember `fork()`?
  - Simplifies software architecture
    - One program is not restricted by the memory layout of the others

# Two processes, one memory?



# Two processes, one memory?



- We want hardware to add base value to every address used in the program

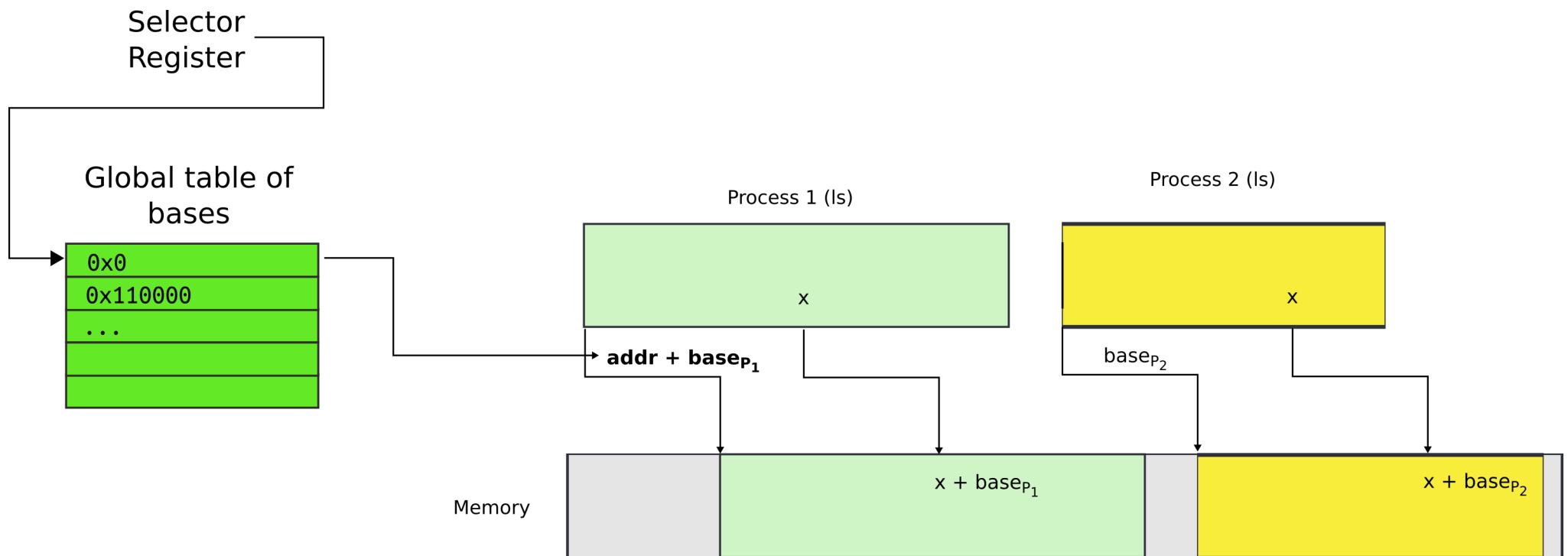
# Seems easy

- One problem
  - Where does this base address come from?

# Seems easy

- One problem
  - Where does this base address come from?
  - Hardware can maintain a table of base addresses
    - One base for each process
  - Dedicate a special register to keep an index into that table

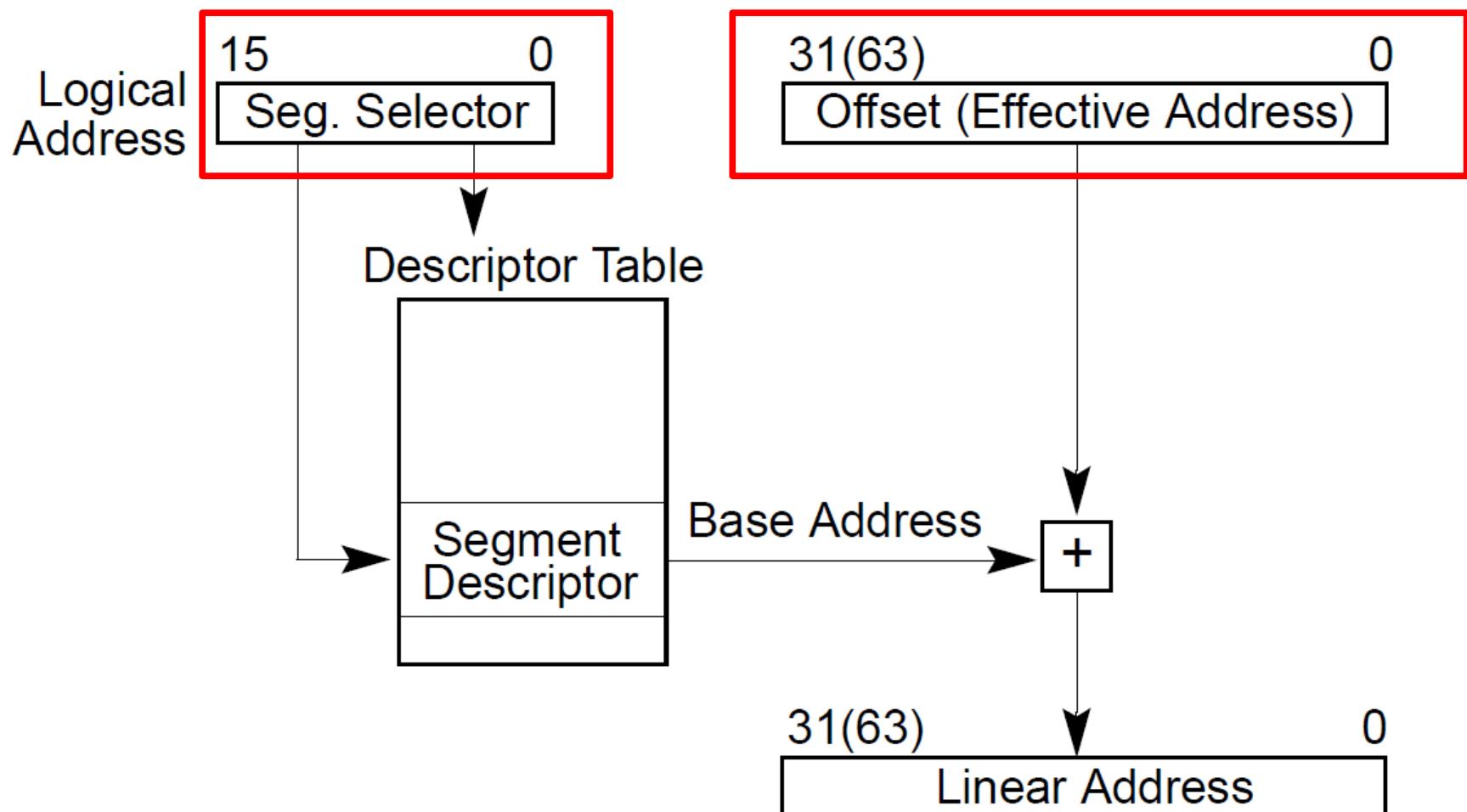
- One problem
  - Where does this base address come from?
  - Hardware can maintain a table of base addresses
    - One base for each process
  - Dedicate a special register to keep an index into that table



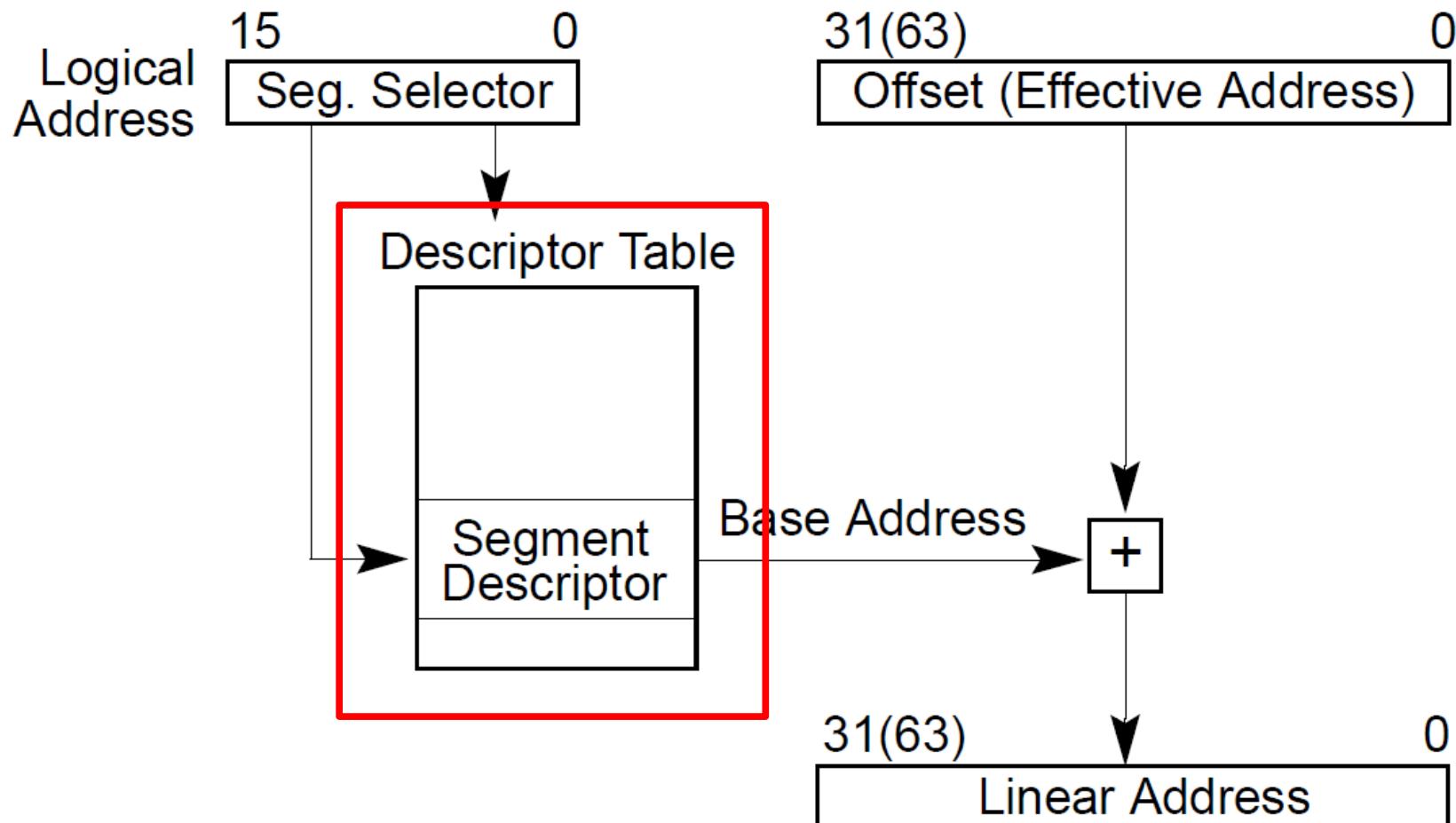
# New addressing mode

All addresses are logical address

- They consist of two parts
    - Segment selector (16 bit) + offset (32 bit)



- Segment selector (16 bit)
  - Is simply an index into an array (Descriptor Table)
  - That holds segment descriptors
    - Base and limit (size) for each segment



# Elements of the descriptor table are segment descriptors

- Base address

- 0 – 4 GB

- Limit (size)

- 0 – 4 GB

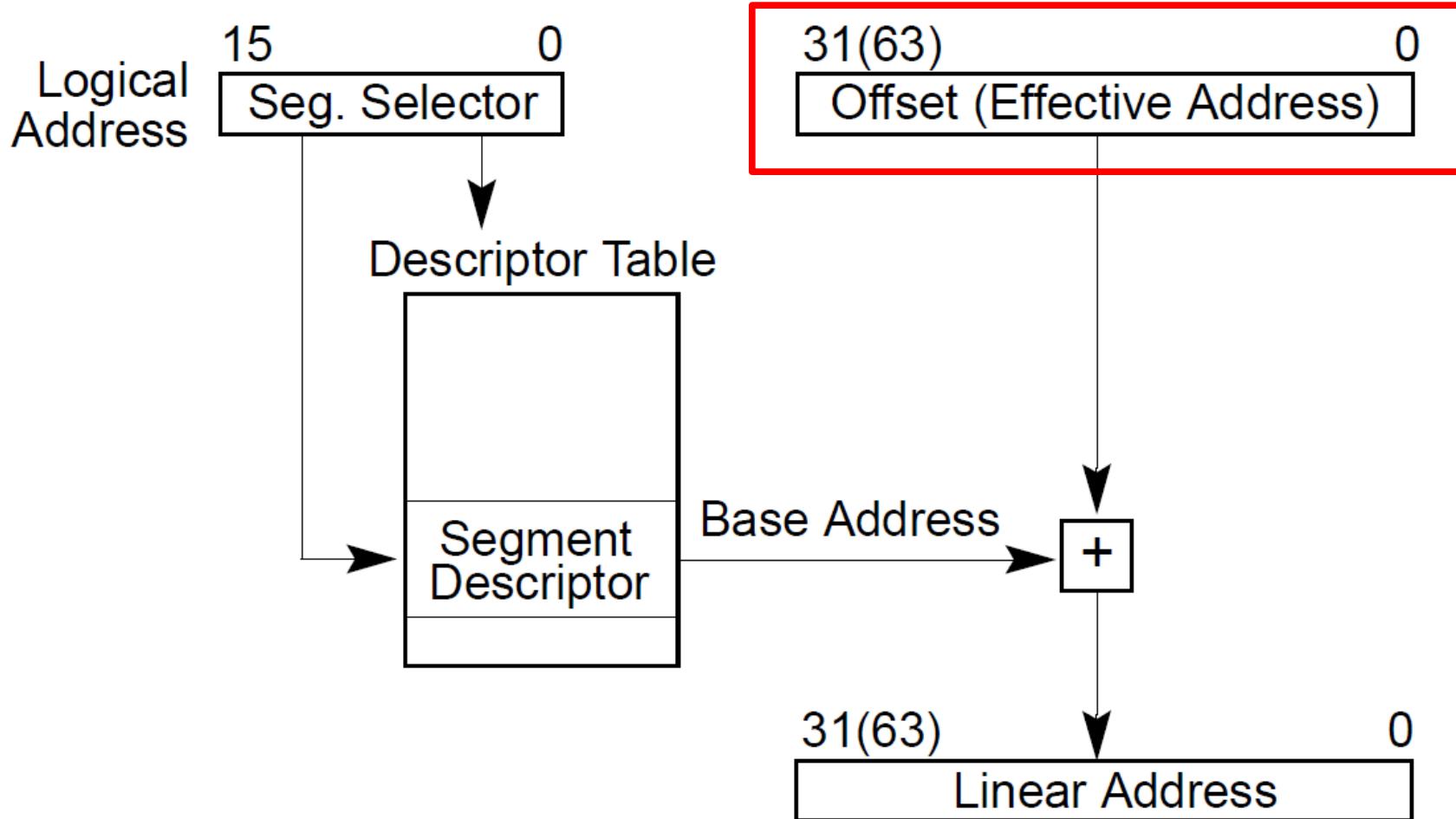
- Access rights

- Executable, readable, writable

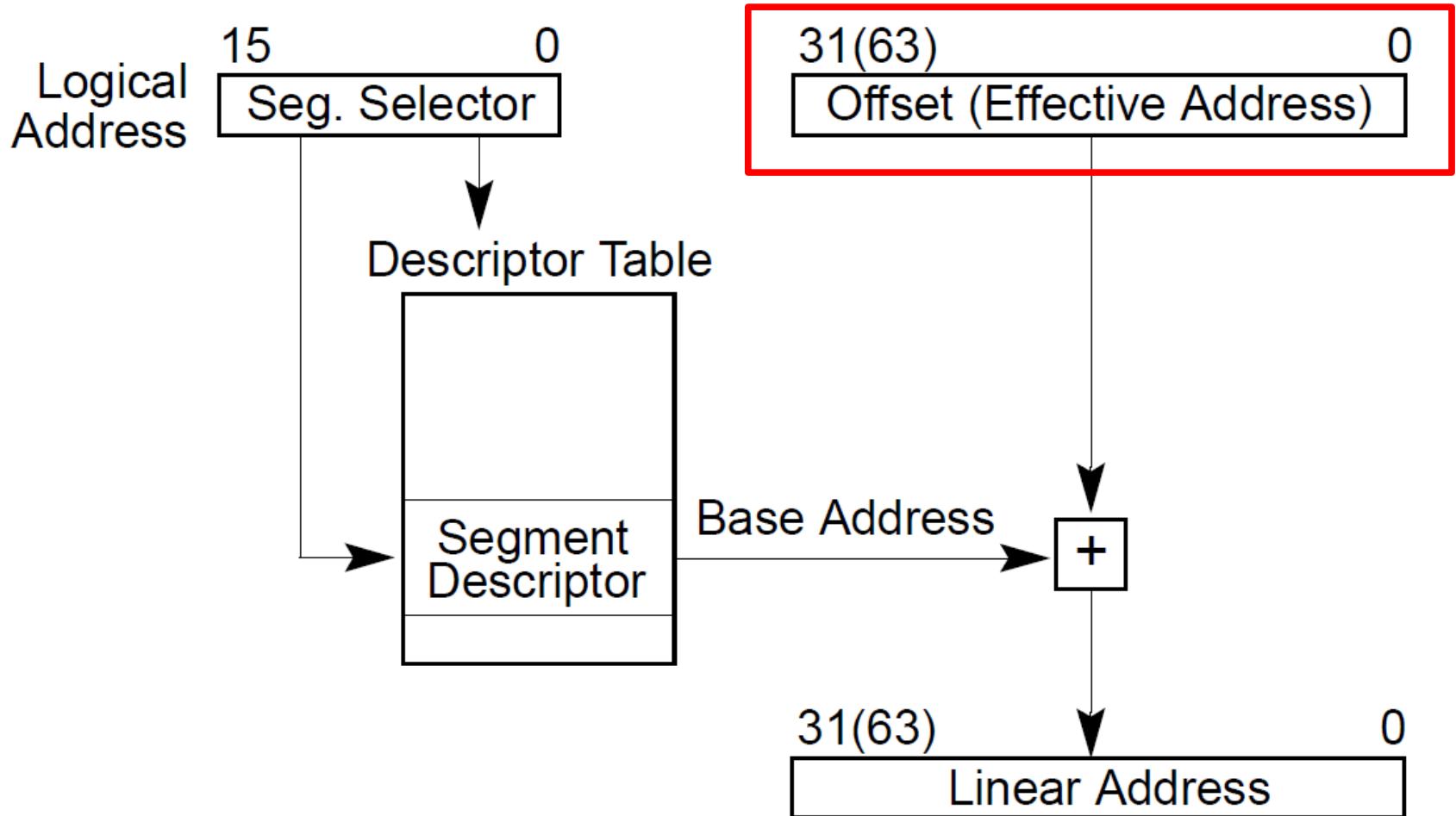
- Privilege level (0 - 3)

Access	Limit
Base Address	

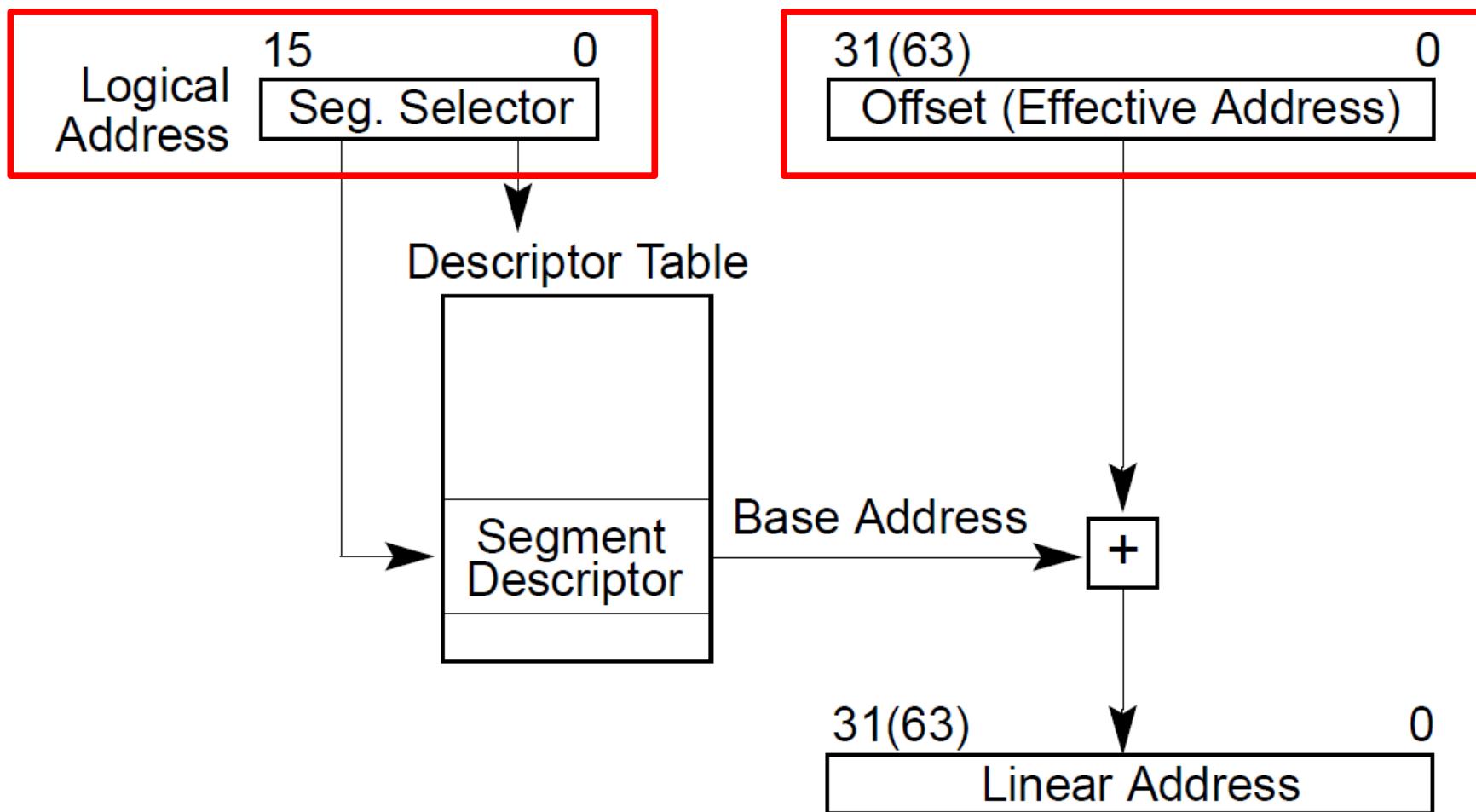
- Offsets into segments (x in our example) or “Effective addresses” are in registers



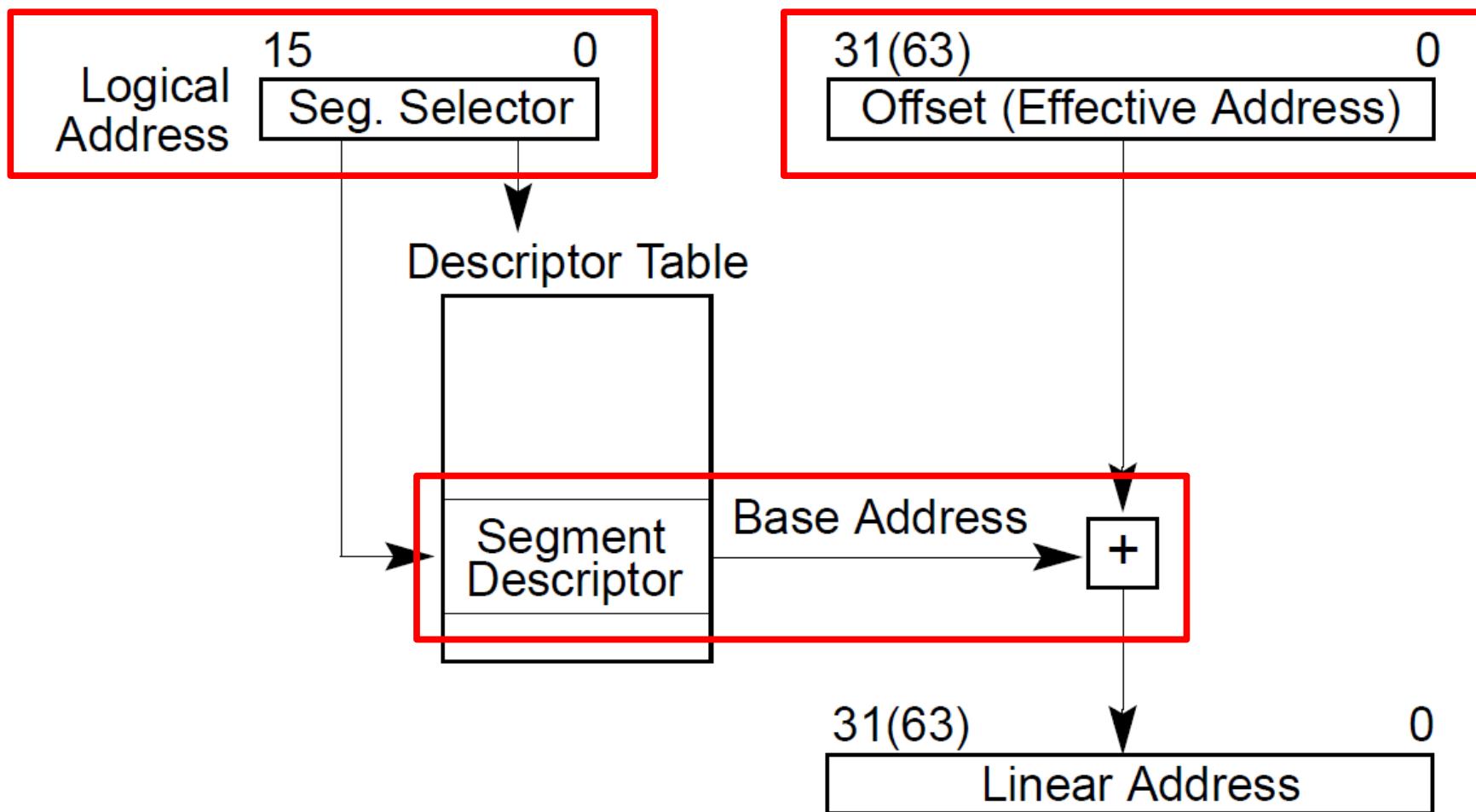
- Logical addresses are translated into physical
  - $\text{Effective address} + \text{DescriptorTable}[selector].Base$



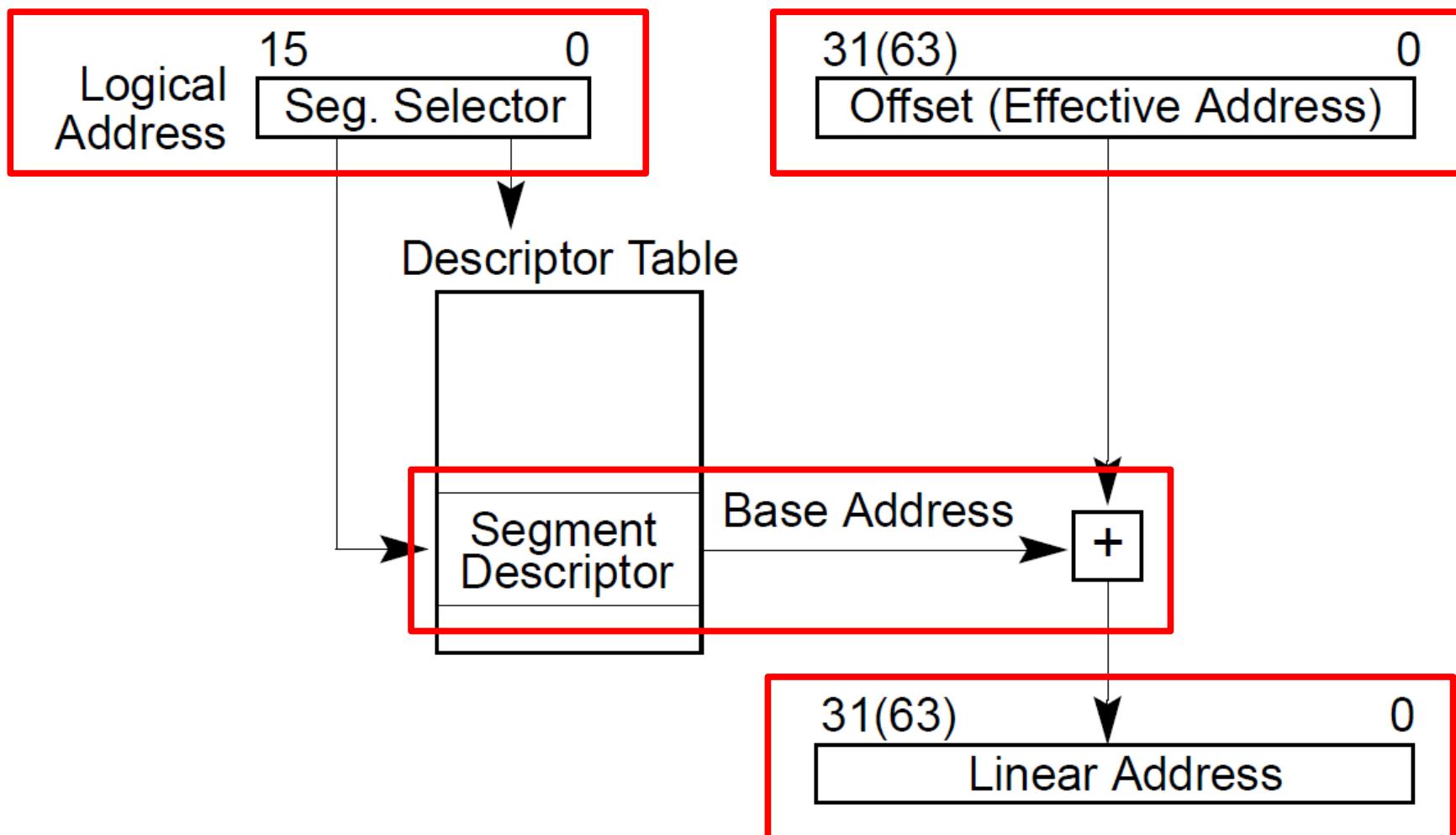
- Logical addresses are translated into physical
  - Effective address + DescriptorTable[selector].Base



- Logical addresses are translated into physical
  - Effective address + DescriptorTable[selector].Base

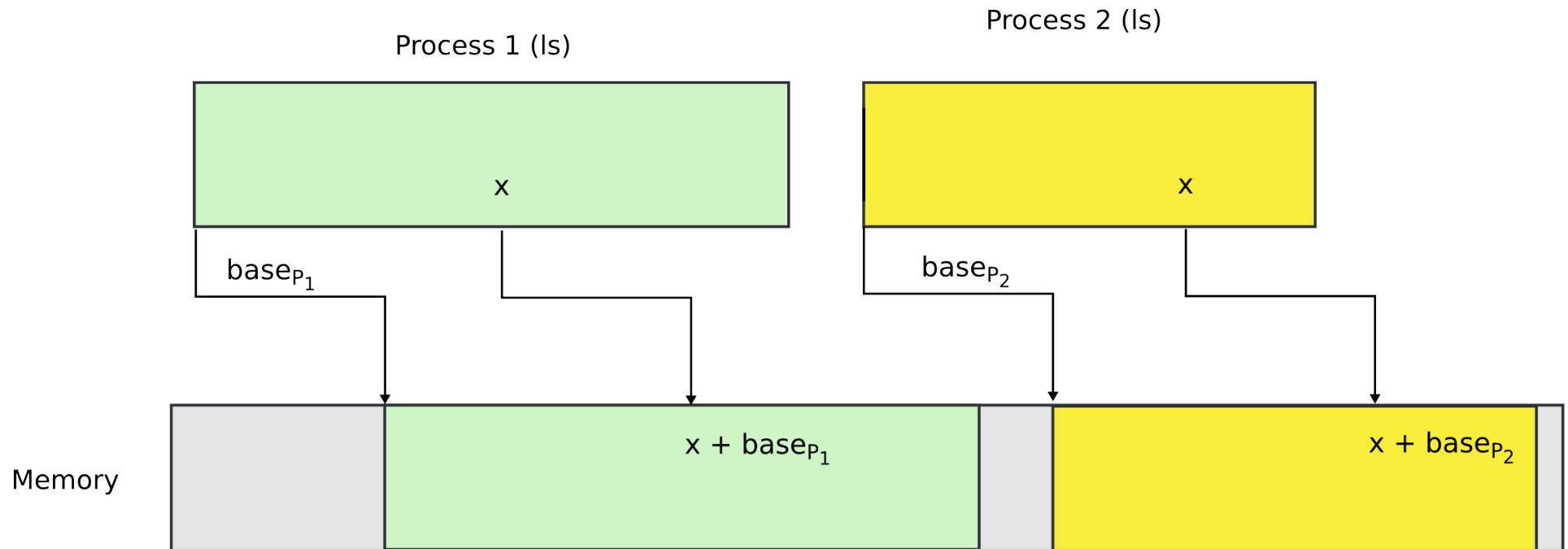


- Logical addresses are translated into physical
  - Effective address + DescriptorTable[selector].Base

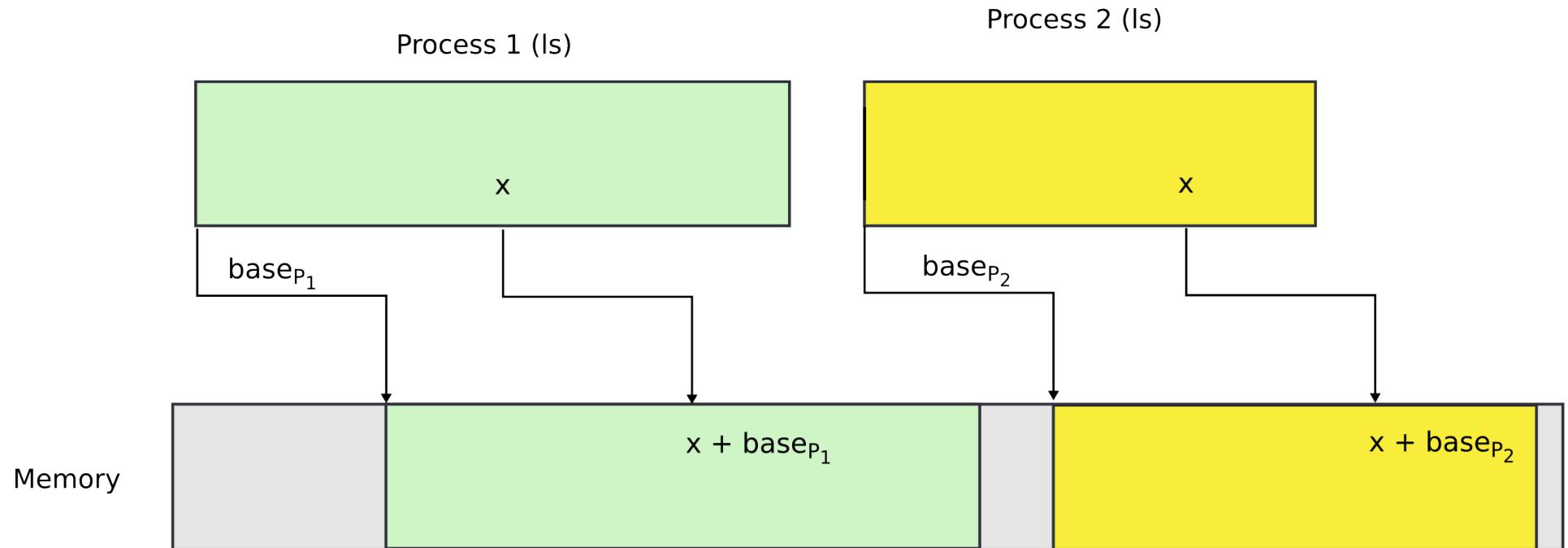


- *Physical address =*  

$$\text{Effective address} + \text{DescriptorTable}[selector].Base$$
- Effective addresses (or offsets) are in registers
- **Selector is in a special register**



- Offsets (effective addresses) are in registers
  - *Effective address + DescriptorTable[selector].Base*
  - **But where is the selector?**



# Segment registers

- Hold 16 bit segment selectors
  - Pointers into a special table
  - Global or local descriptor table
- Segments are associated with one of three types of storage
  - Code
  - Data
  - Stack

# Segmented programming (not real)

```
static int x = 1;           ds:x = 1; // data
int y; // stack           ss:y;      // stack
if (x) {                  if (ds:x) {
    y = 1;                ss:y = 1;
    printf ("Boo");       cs:printf(ds:"Boo");
} else                      } else
    y = 0;                ss:y = 0;
```

# Programming model

- Segments for: code, data, stack, “extra”
  - A program can have up to 6 total segments
  - Segments identified by registers: cs, ds, ss, es, fs, gs
- Prefix all memory accesses with desired segment:
  - `mov eax, ds:0x80` (load offset 0x80 from data into eax)
  - `jmp cs:0xab8` (jump execution to code offset 0xab8)
  - `mov ss:0x40, ecx` (move ecx to stack offset 0x40)

# Programming model, cont.

- This is cumbersome,
- Instead the idea is: infer code, data and stack segments from the instruction type:
  - Control-flow instructions use code segment (jump, call)
  - Stack management (push/pop) uses stack
  - Most loads/stores use data segment
- Extra segments (es, fs, gs) must be used explicitly

# Code segment

- Code
  - CS register
  - EIP is an offset inside the segment stored in CS
- Can only be changed with
  - procedure calls,
  - interrupt handling, or
  - task switching

# Data segment

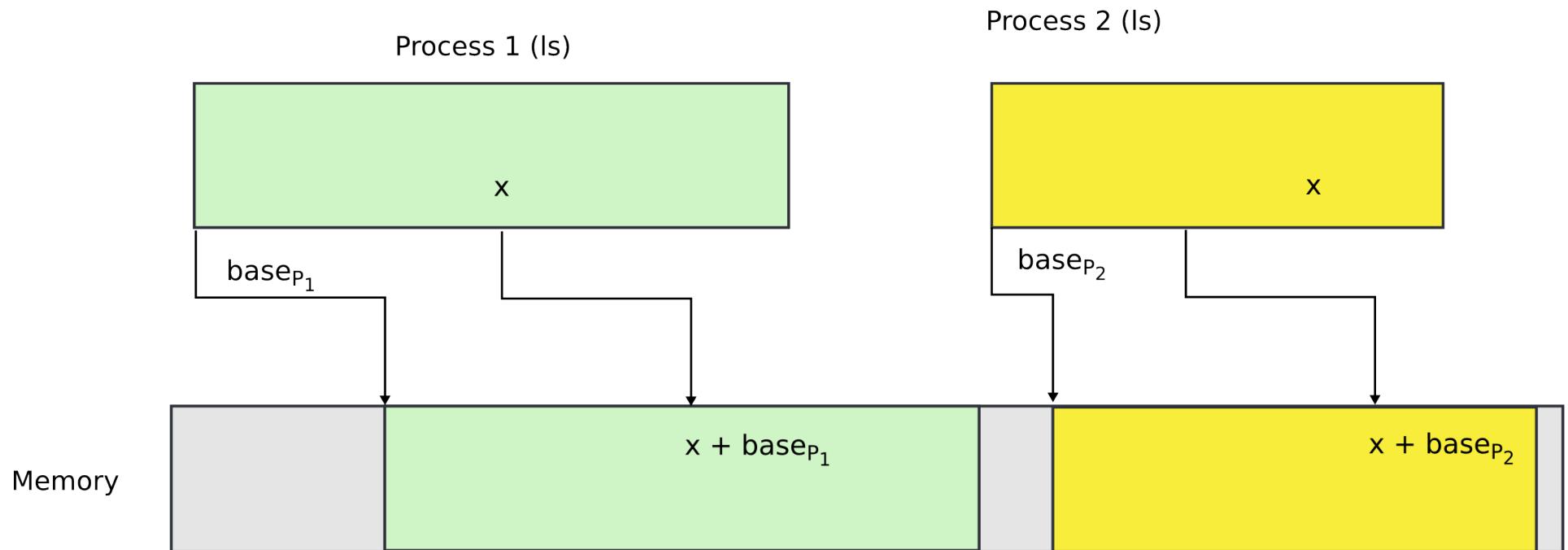
- Data
  - DS, ES, FS, GS
  - 4 possible data segments can be used at the same time

# Stack segment

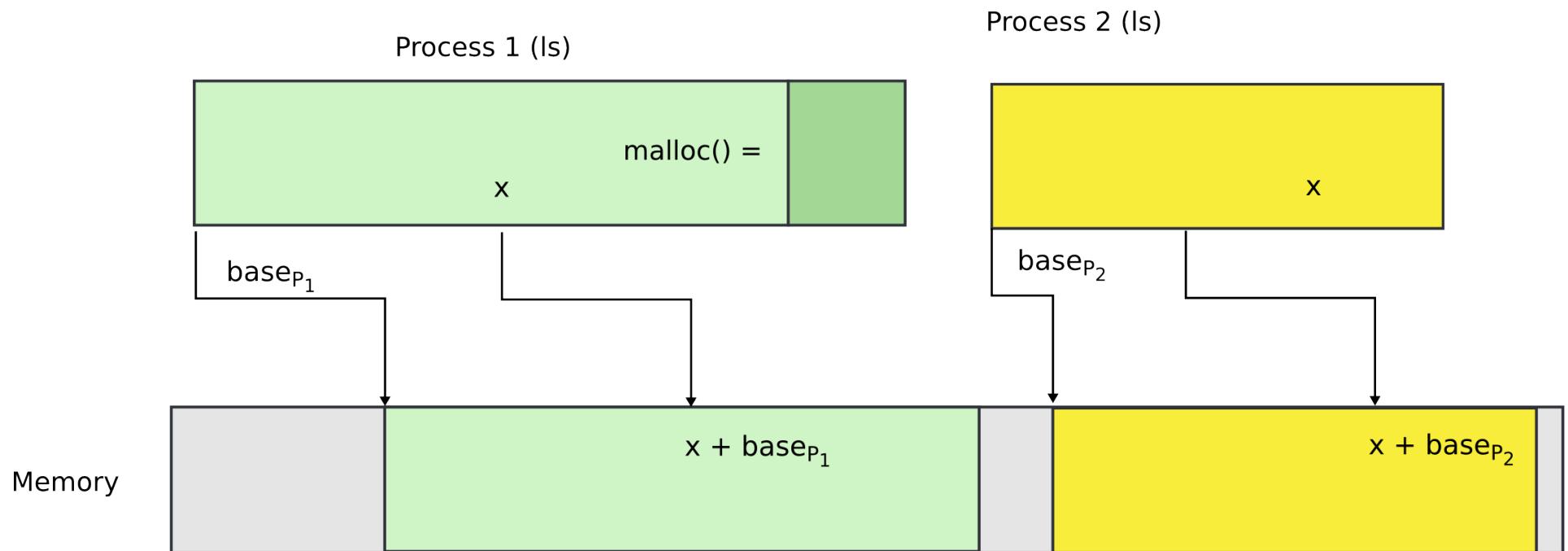
- Stack
  - SS
- Can be loaded explicitly
  - OS can set up multiple stacks
  - Of course, only one is accessible at a time

Segmentation works for isolation, i.e., it does provide programs with illusion of private memory

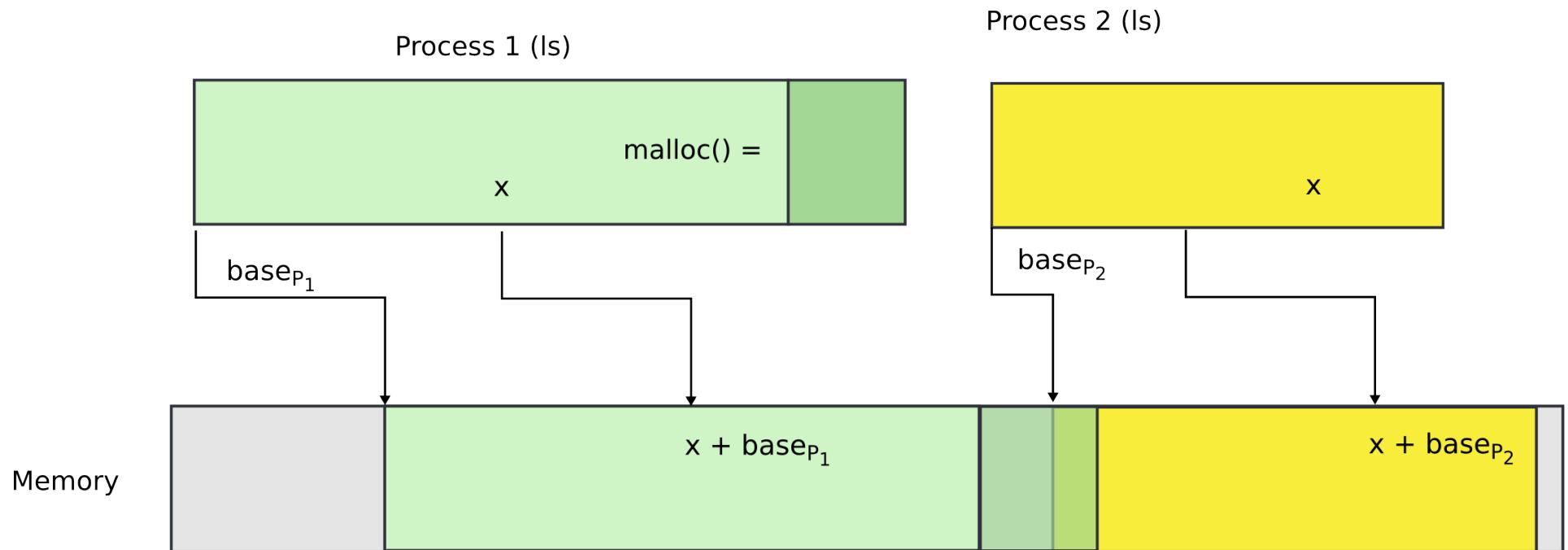
# Segmentation is ok... but



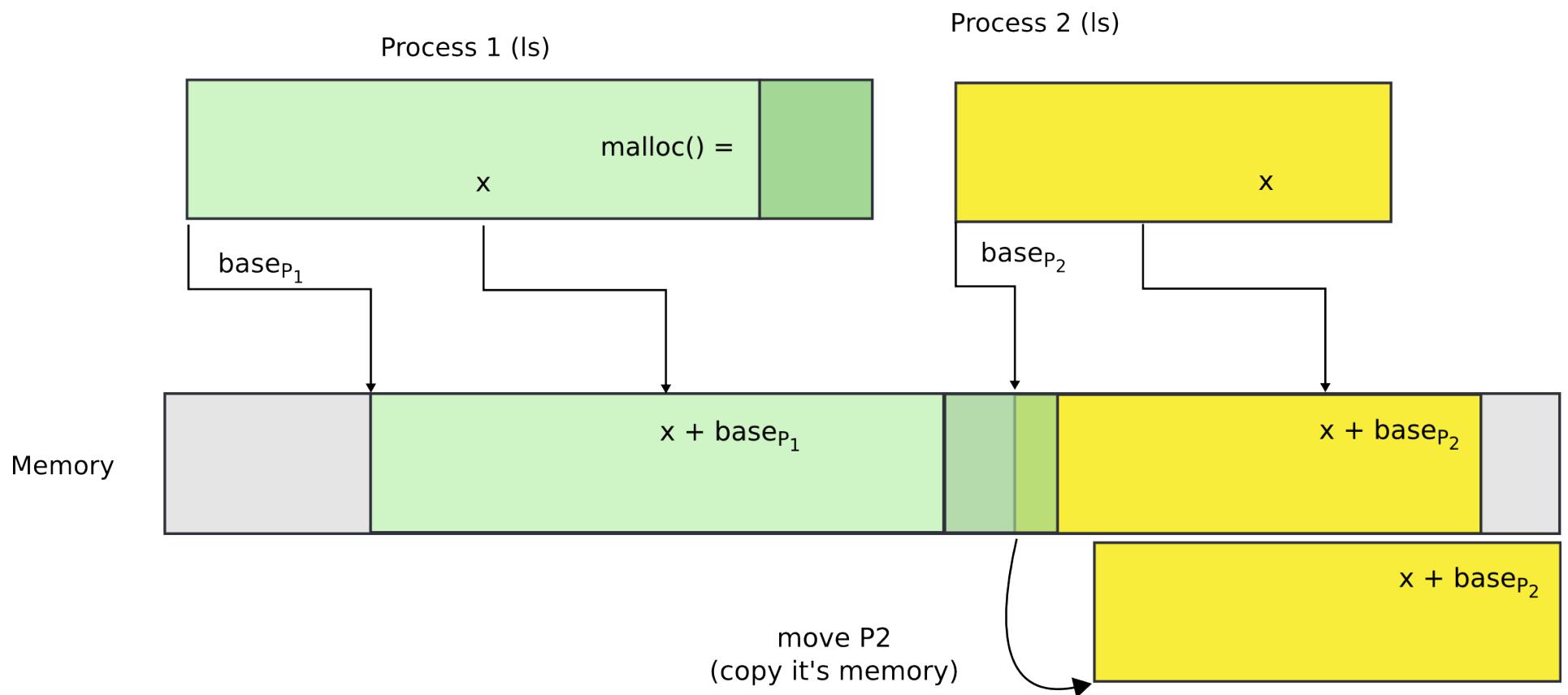
# What if process needs more memory?



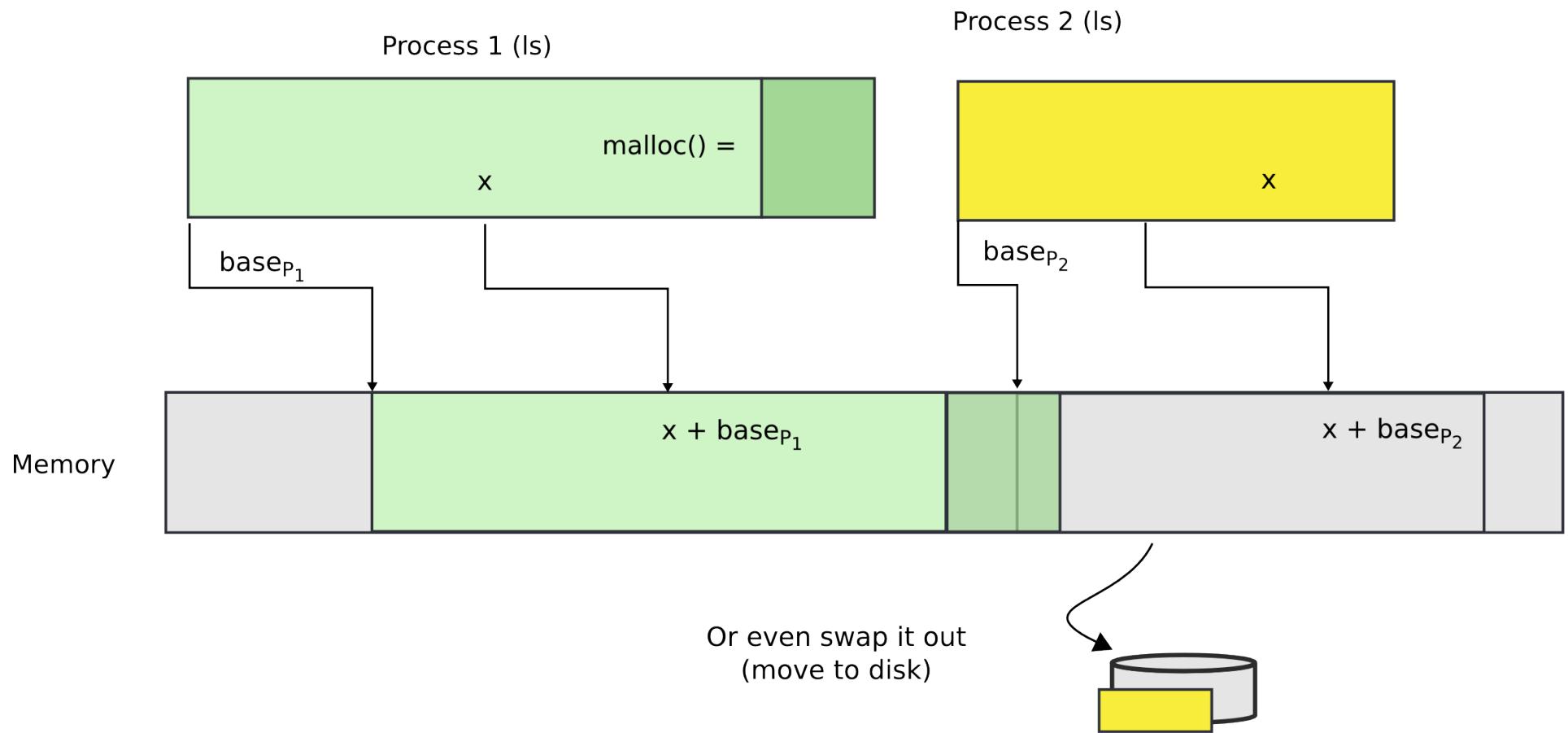
# What if process needs more memory?



# You can move P2 in memory



# Or even swap it out to disk



# Problems with segments

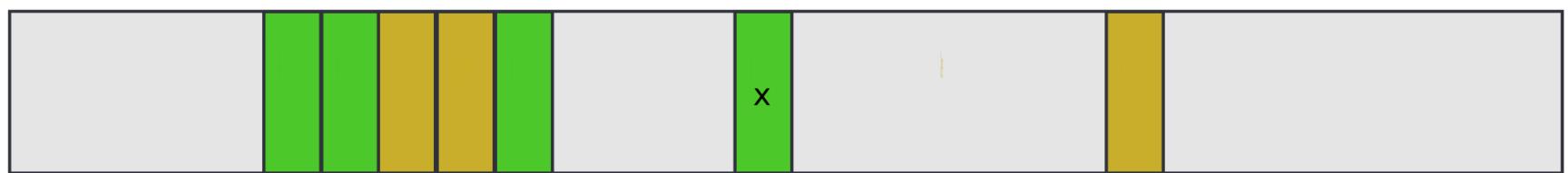
- But it's inefficient
  - Relocating or swapping the entire process takes time
- Memory gets fragmented
  - There might be no space (gap) for the swapped out process to come in
  - Will have to swap out other processes

# Paging

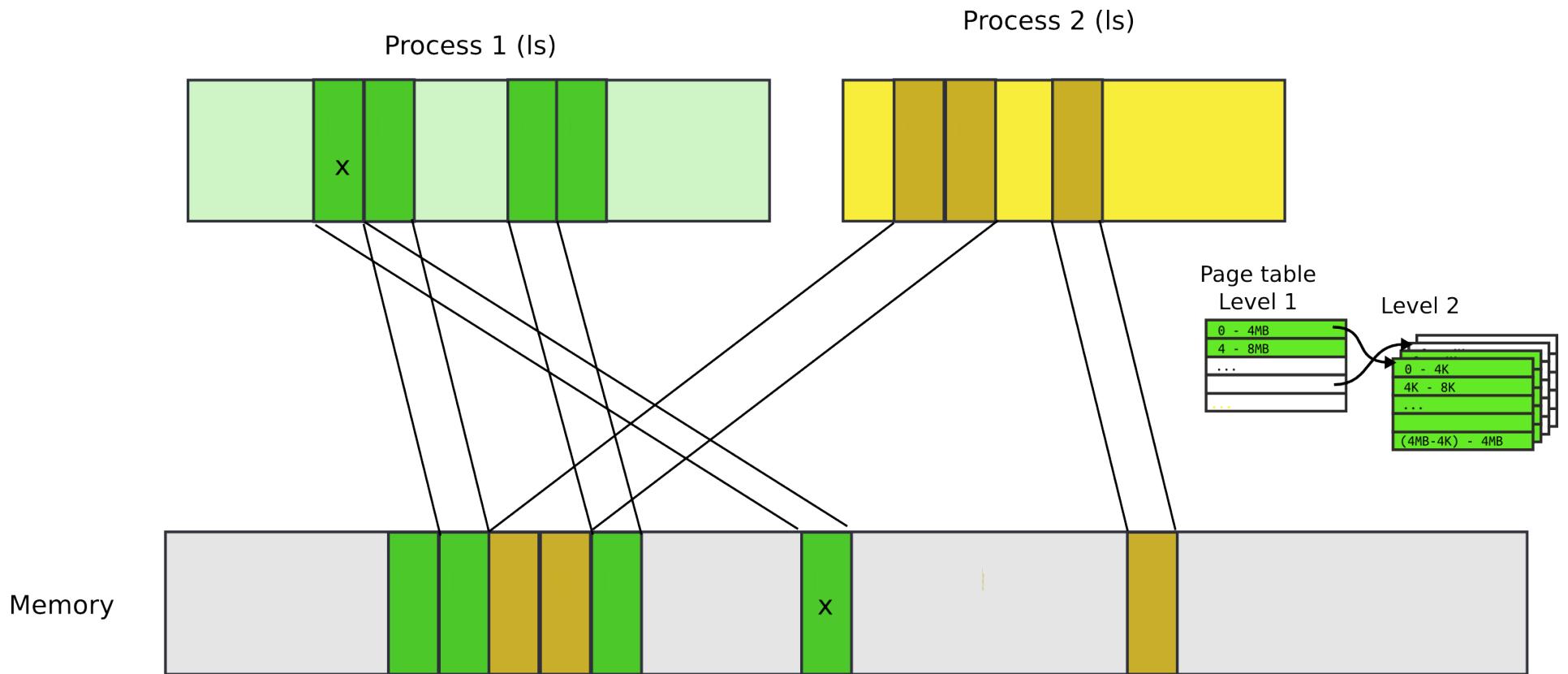
# Pages



Memory



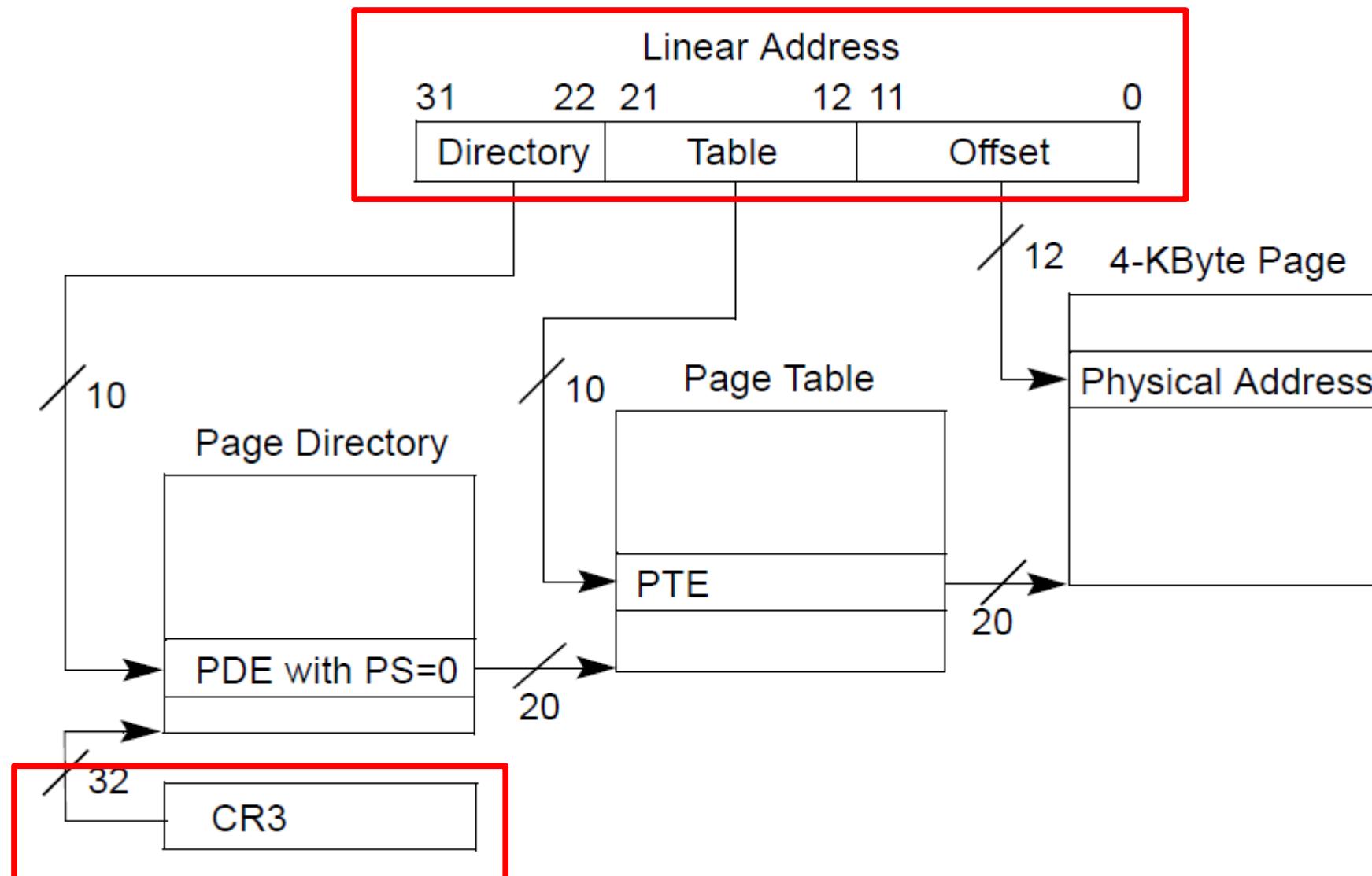
# Pages



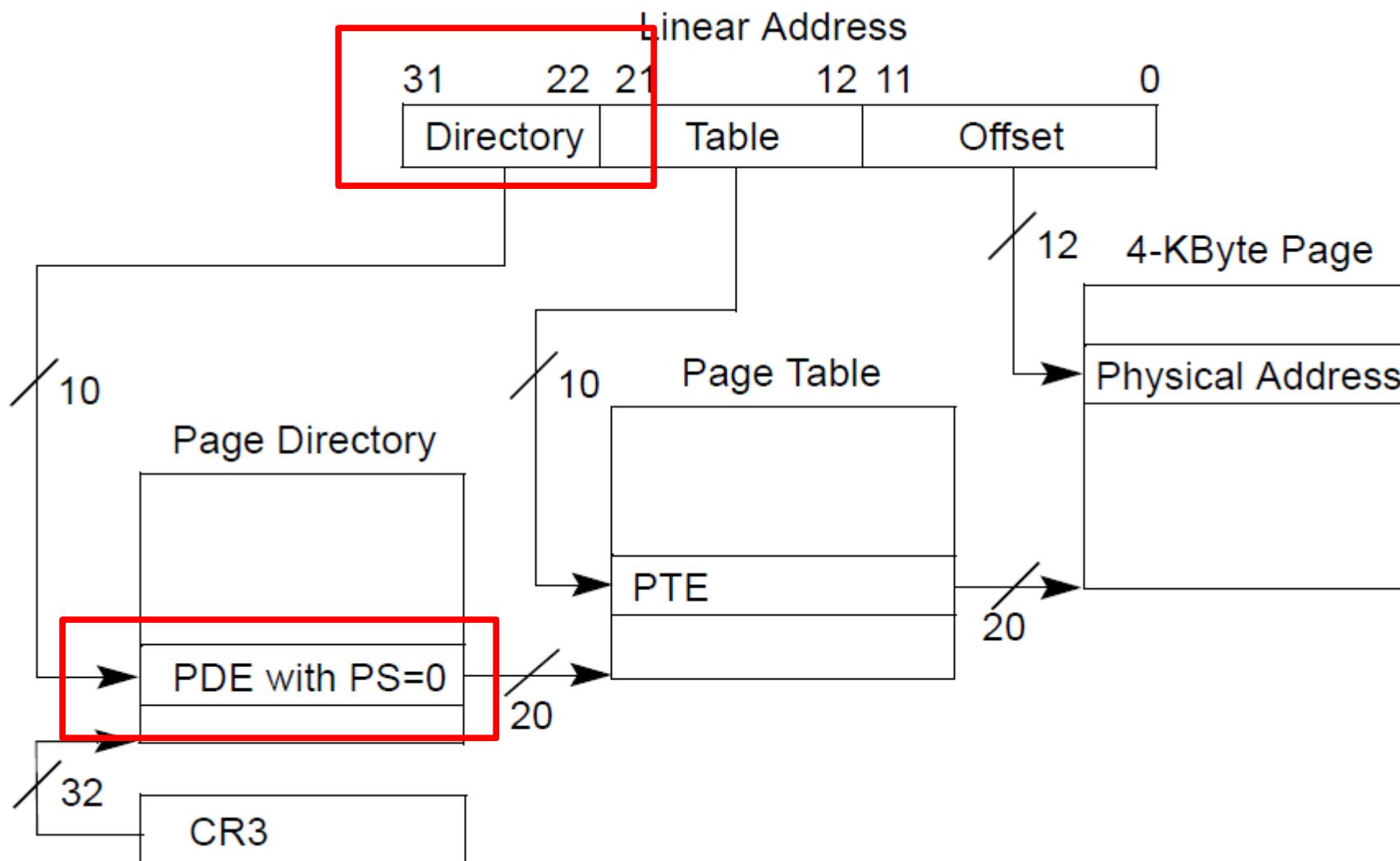
# Paging idea

- Break up memory into 4096-byte chunks called pages
  - Modern hardware supports 2MB, 4MB, and 1GB pages
- Independently control mapping for each page of linear address space
- Compare with segmentation (single base + limit)
  - many more degrees of freedom

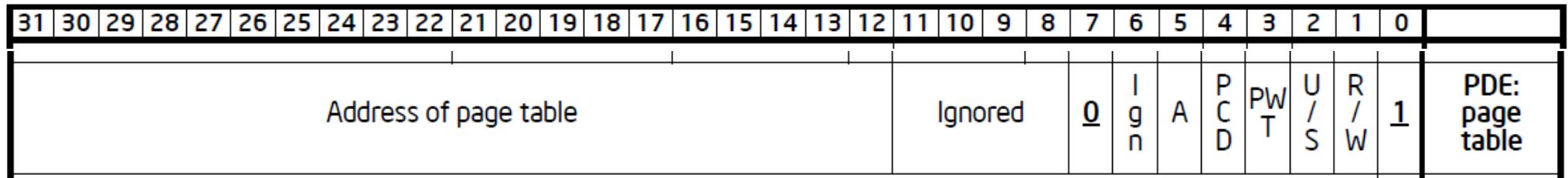
# Page translation



# Page translation

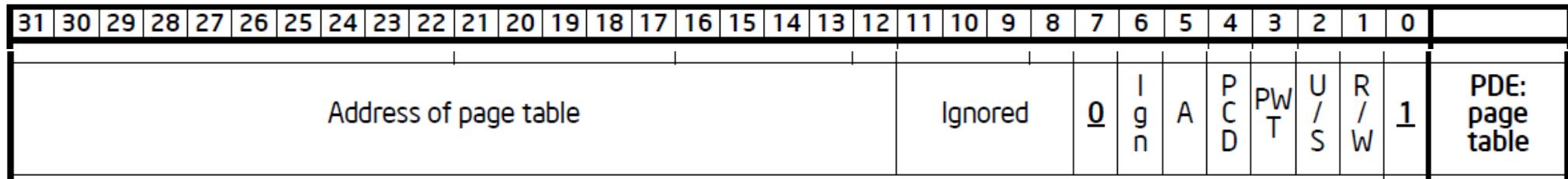


# Page directory entry (PDE)



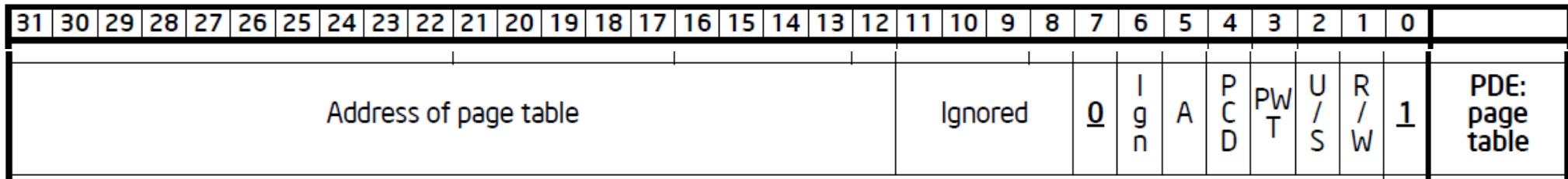
- 20 bit address of the page table

# Page directory entry (PDE)



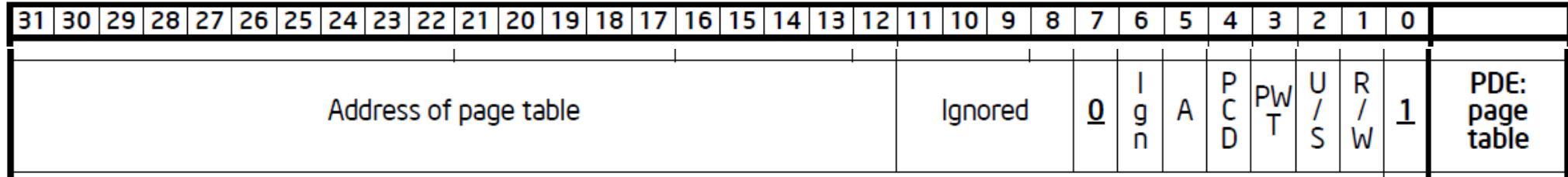
- 20 bit address of the page table
- Wait... 20 bit address, but we need 32 bits

# Page directory entry (PDE)



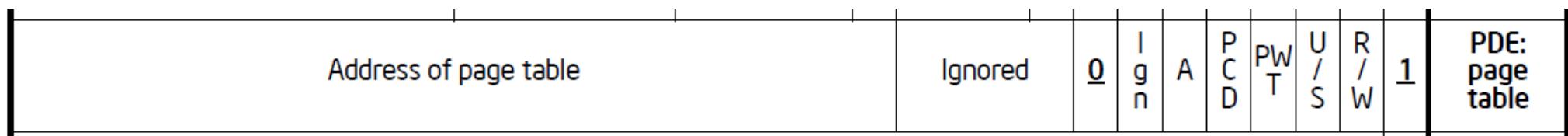
- 20 bit address of the page table
- Wait... 20 bit address, but we need 32 bits
  - Pages 4KB each, we need 1M to cover 4GB
  - Pages start at 4KB (page aligned boundary)

# Page directory entry (PDE)



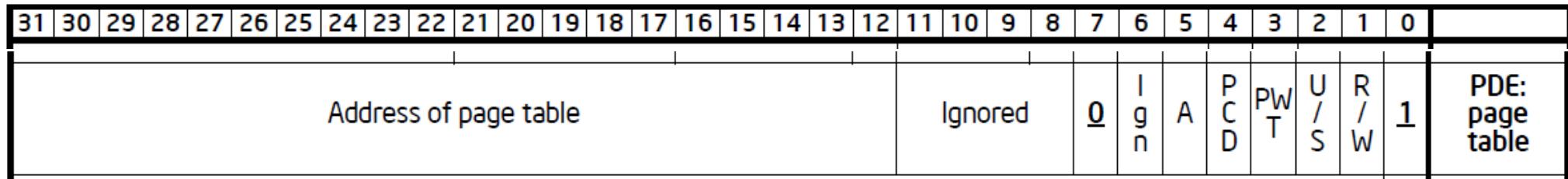
- Bit #1: R/W – writes allowed?
  - But allowed where?

# Page directory entry (PDE)



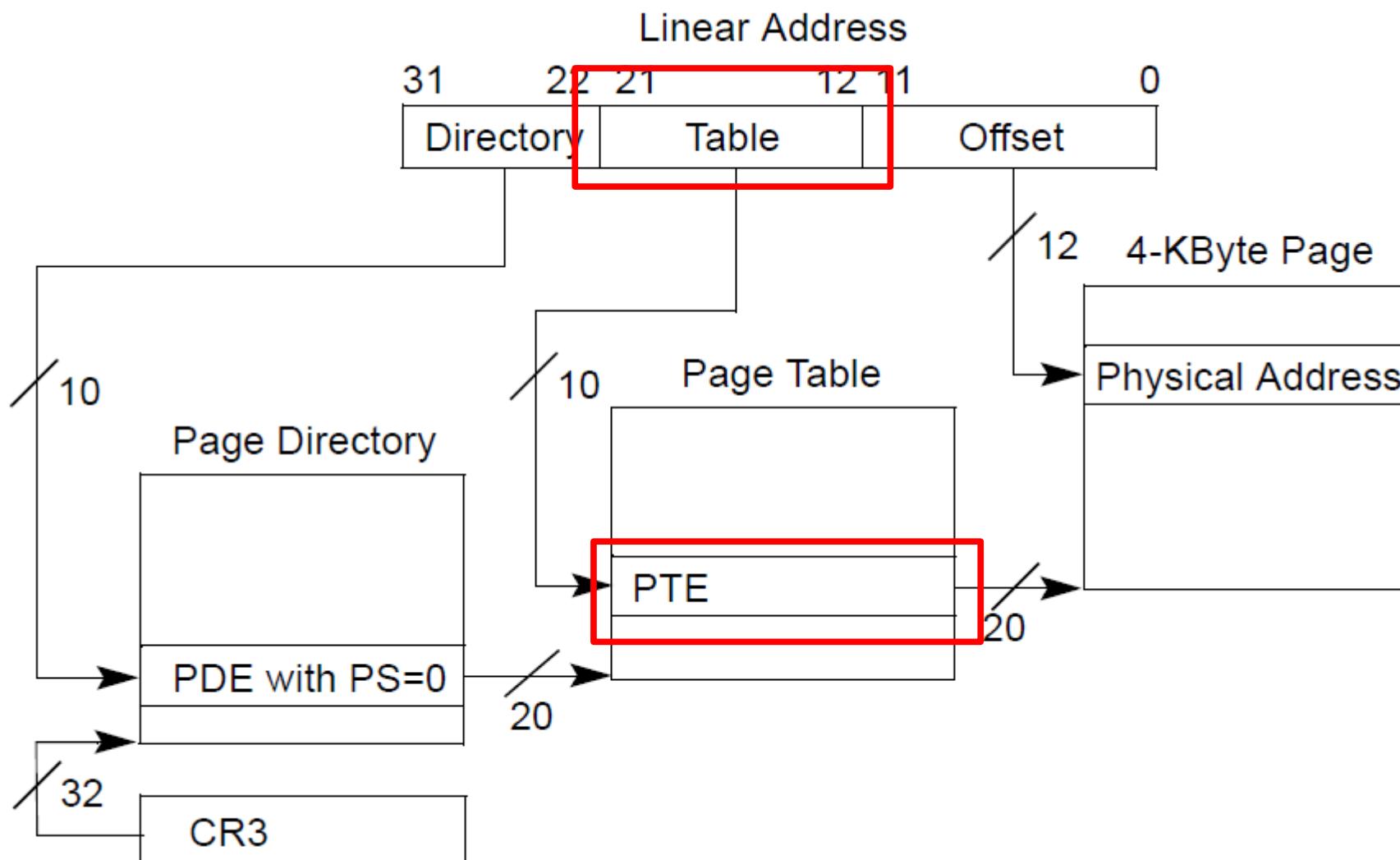
- Bit #1: R/W – writes allowed?
  - But allowed where?
  - One page directory entry controls 1024 Level 2 page tables
    - Each Level 2 maps 4KB page
  - So it's a region of  $4\text{KB} \times 1024 = 4\text{MB}$

# Page directory entry (PDE)

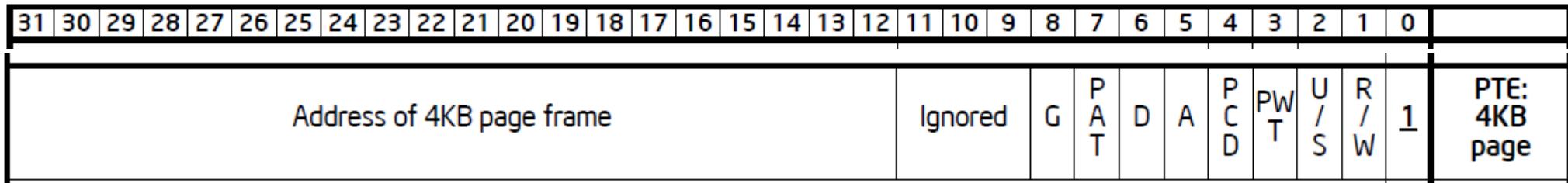


- Bit #2: U/S – user/supervisor
  - If 0 – user-mode access is not allowed
  - Allows protecting kernel memory from user-level applications

# Page translation

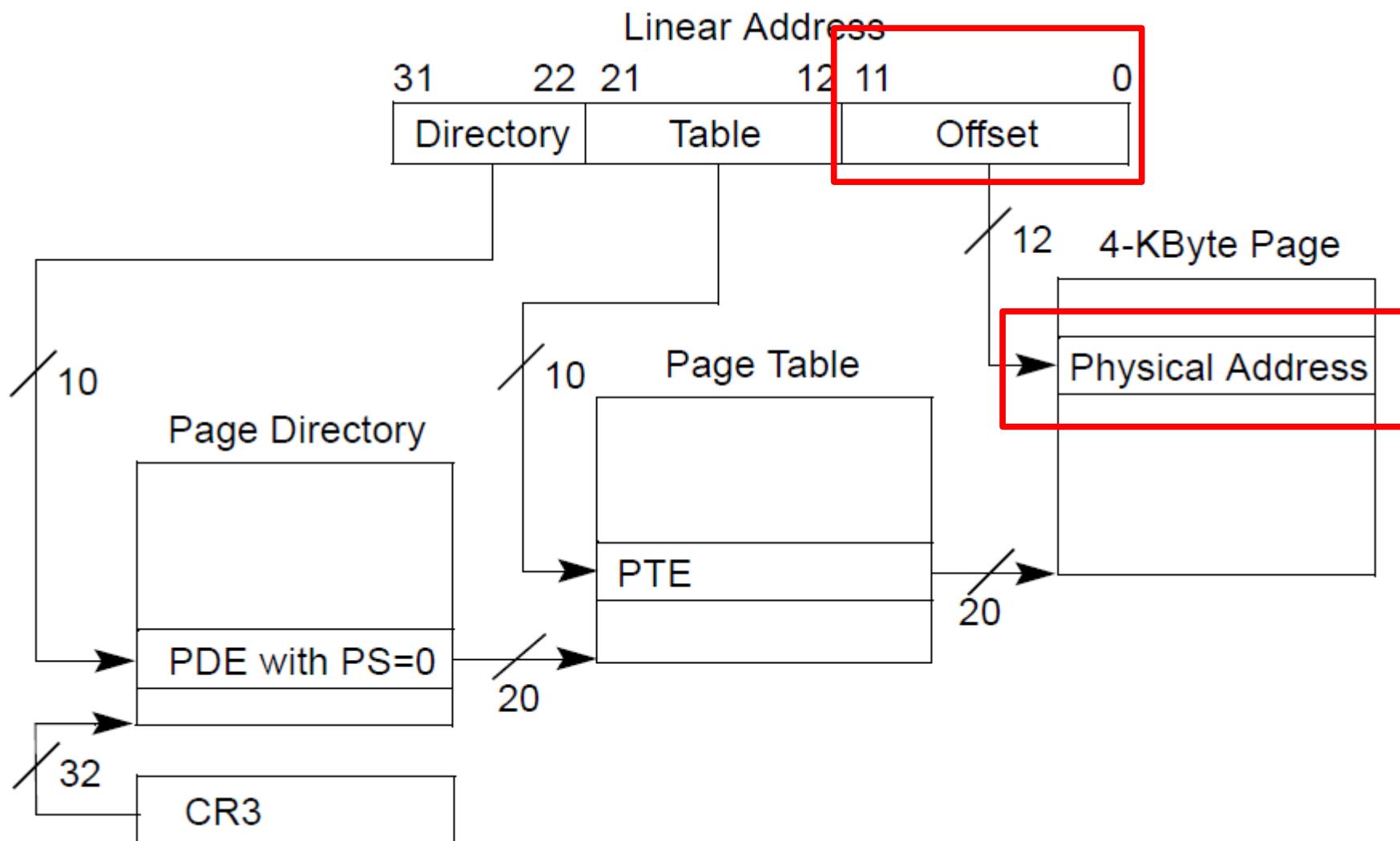


# Page table entry (PTE)



- 20 bit address of the 4KB page
  - Pages 4KB each, we need 1M to cover 4GB
- Bit #1: R/W – writes allowed?
  - To a 4KB page
- Bit #2: U/S – user/supervisor
  - If 0 user-mode access is not allowed
- Bit #5: A – accessed
- Bit #6: D – dirty – software has written to this page

# Page translation



```
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0  
EBX = 20 983 809
```

20 983 809 =  00 0000 0101|00 0000 0011|0000 0000 0001

{  
page number

1M (1,048,575)

Virtual Address  
Space (or Memory)  
of the Process

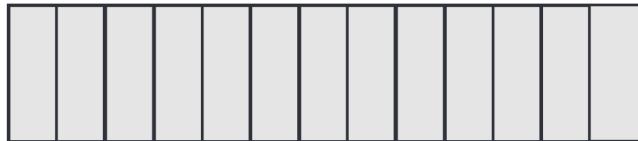


0 1 2

page number = 5123  
or (0b1 0100 0000 0011)

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

Physical  
Memory



```
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0  
EBX = 20 983 809
```

20 983 809 = **00 0000 0101|00 0000 0011|0000 0000 0001**

page number

1M (1,048,575)

Virtual Address  
Space (or Memory)  
of the Process



0 1 2

page number = 5123  
or (0b1 0100 0000 0011)

CR3 = 0

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

Physical  
Memory



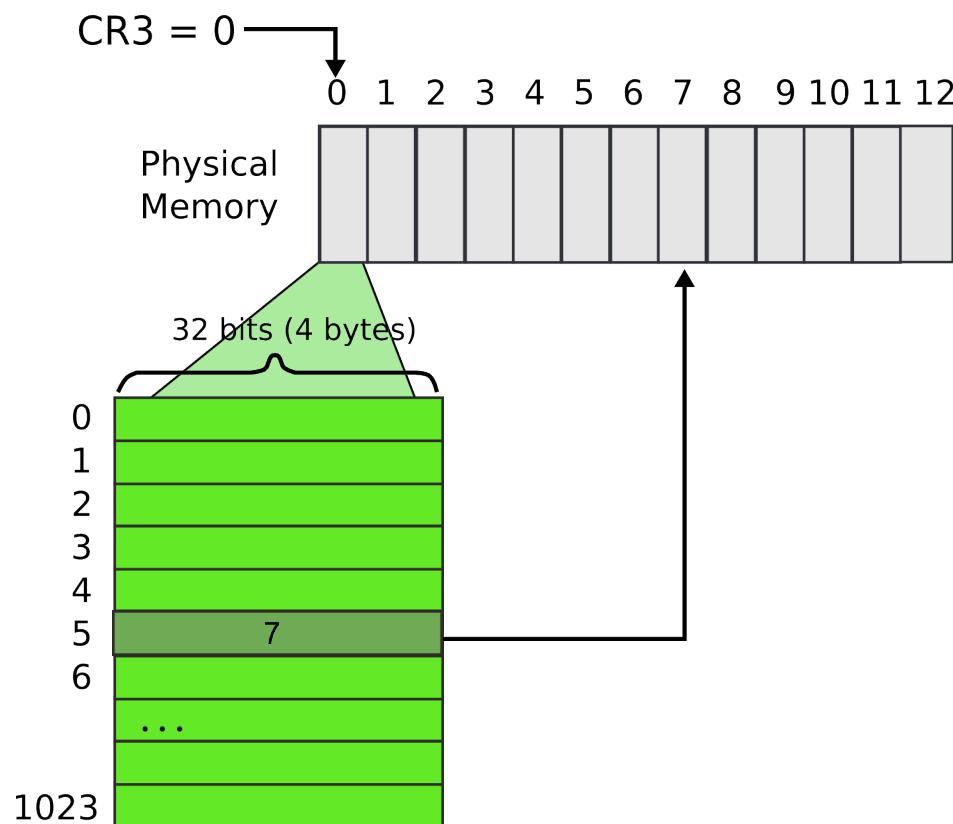
```

mov (%EBX), EAX # mov value from the location pointed by EBX into EAX
EAX = 0
EBX = 20 983 809

```

20 983 809 =  1M (1,048,575)

page number



```

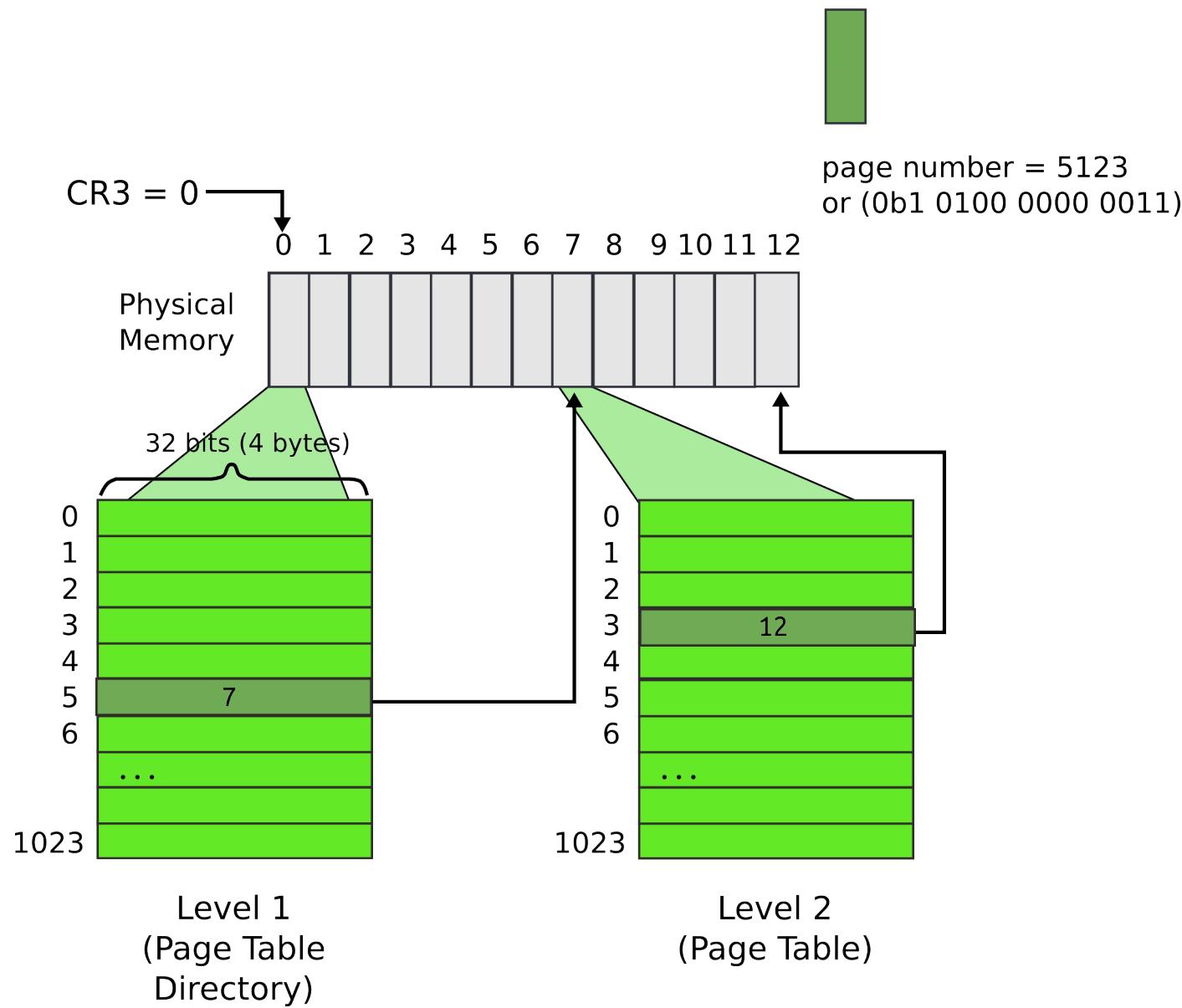
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX
EAX = 0
EBX = 20 983 809

```

$20\ 983\ 809 = \boxed{00\ 0000\ 010} \boxed{00\ 0000\ 0011} \boxed{0000\ 0000\ 0001}$

page number

1M (1,048,575)



```

mov (%EBX), EAX # mov value from the location pointed by EBX into EAX
EAX = 0
EBX = 20 983 809

```

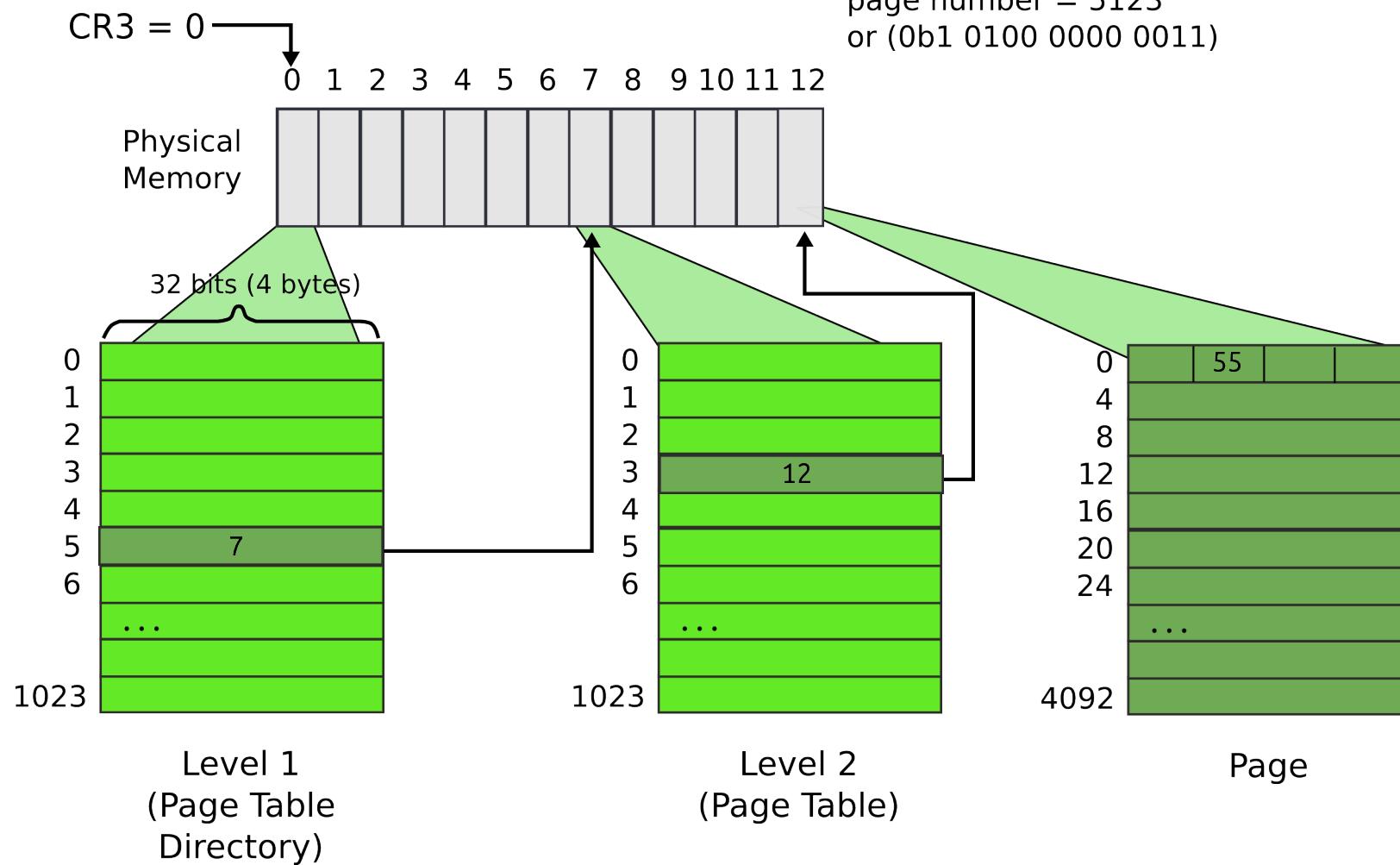
$20\ 983\ 809 = \boxed{00\ 0000\ 0101} \boxed{00\ 0000\ 0011} \boxed{0000\ 0000\ 0001}$

page number

1M (1,048,575)



page number = 5123  
or (0b1 0100 0000 0011)



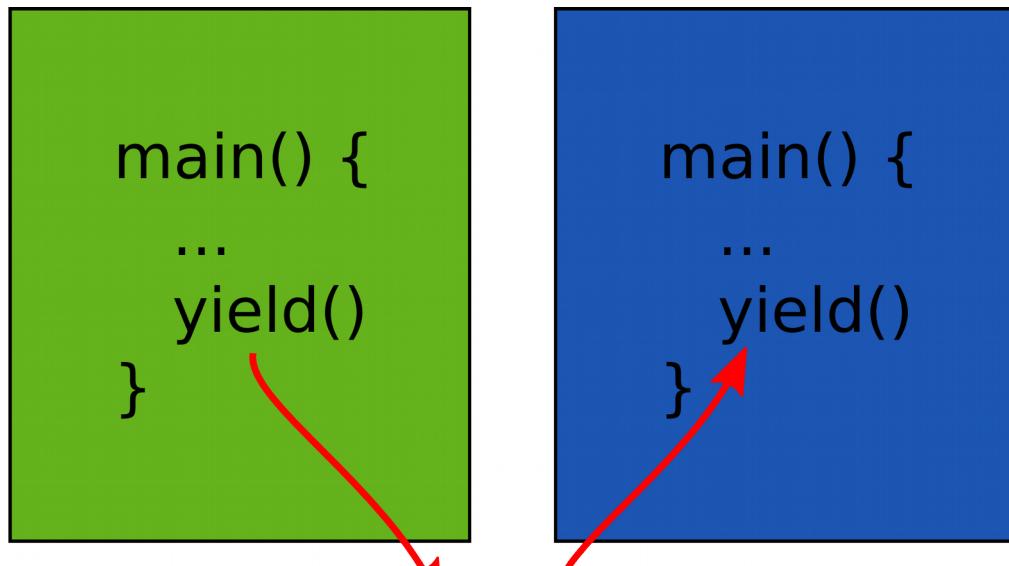
- Result:
  - $EAX = 55$

# But why do we need page tables

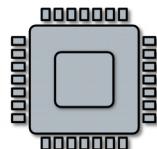
... Instead of arrays?

- Page tables represent sparse address space more efficiently
  - An entire array has to be allocated upfront
  - But if the address space uses a handful of pages
  - Only page tables (Level 1 and 2 need to be allocated to describe translation)
- On a dense address space this benefit goes away
  - I'll assign a homework!

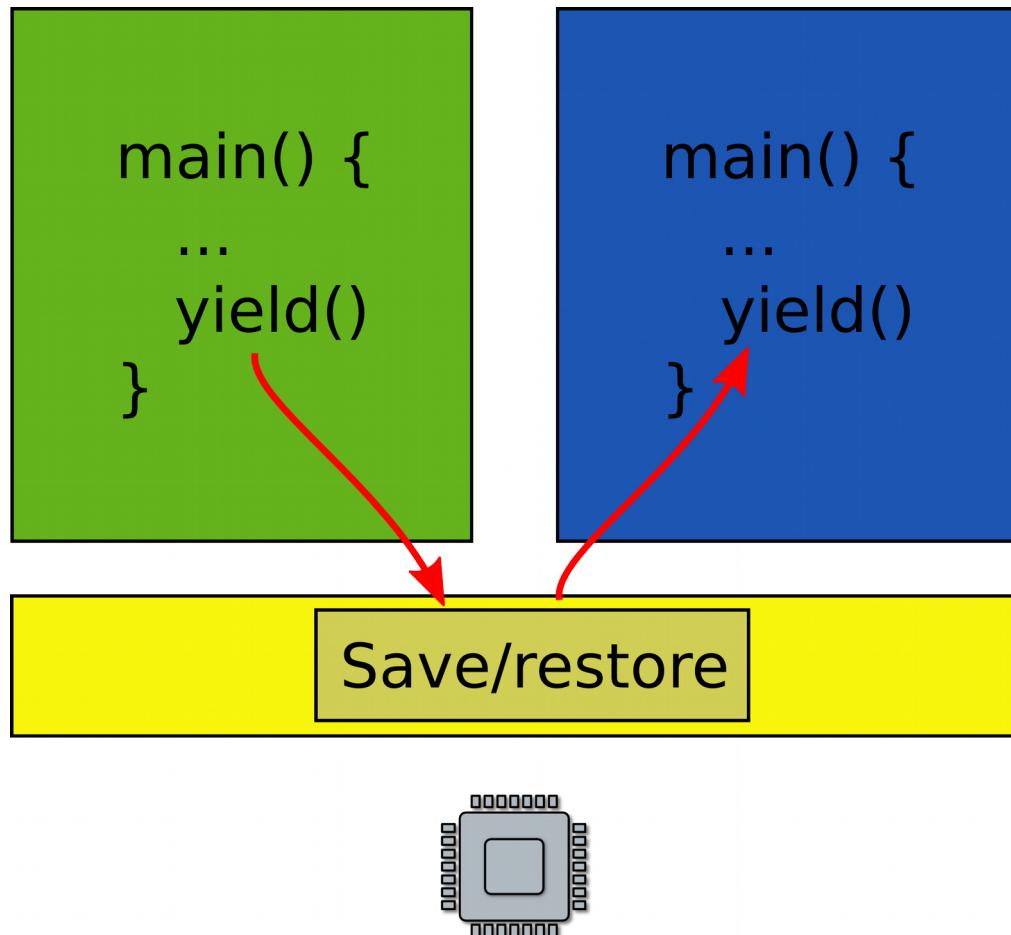
# What about isolation?



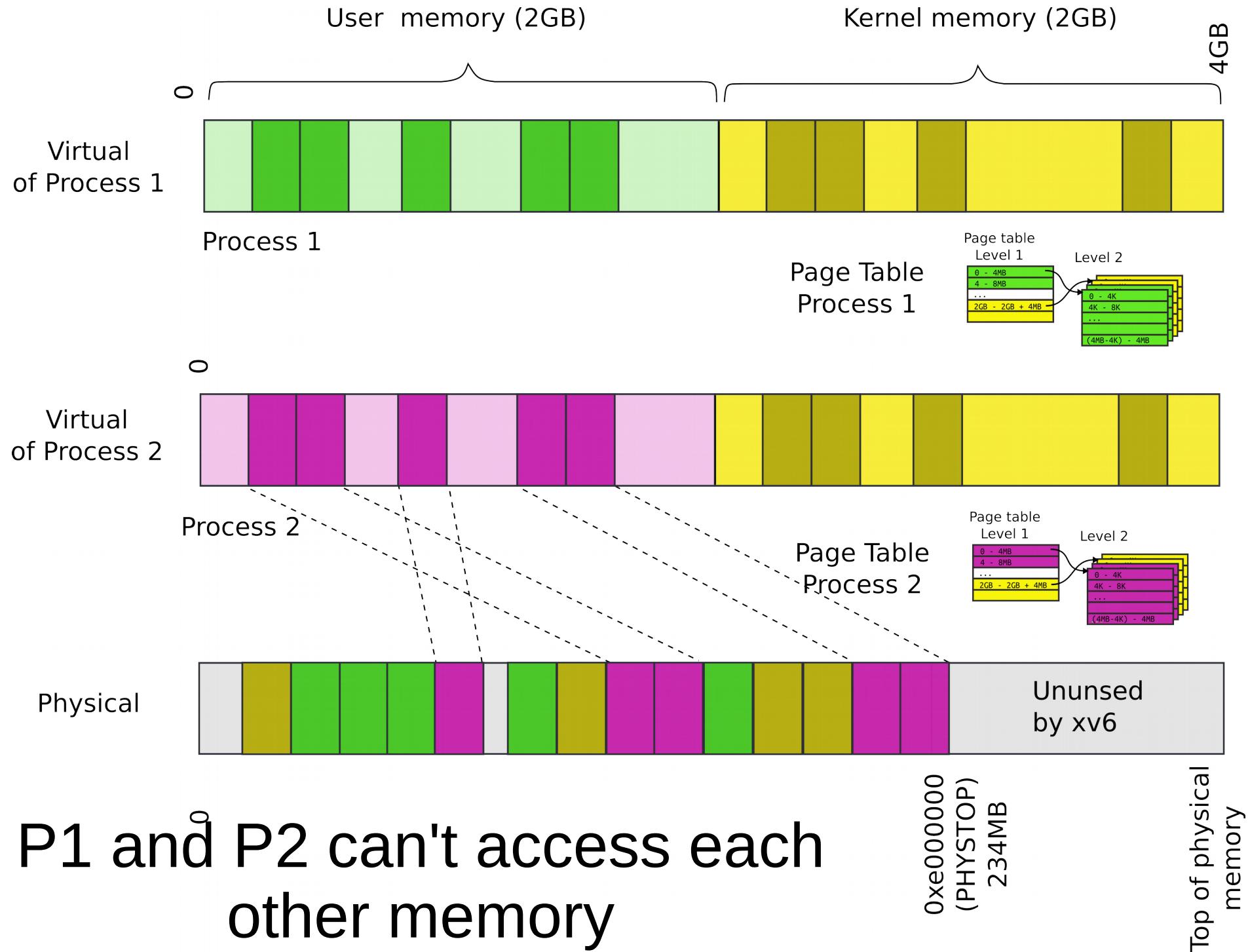
- Two programs,  
one memory?



# What about isolation?



- Two programs, one memory?
- Each process has its own page table
  - OS switches between them

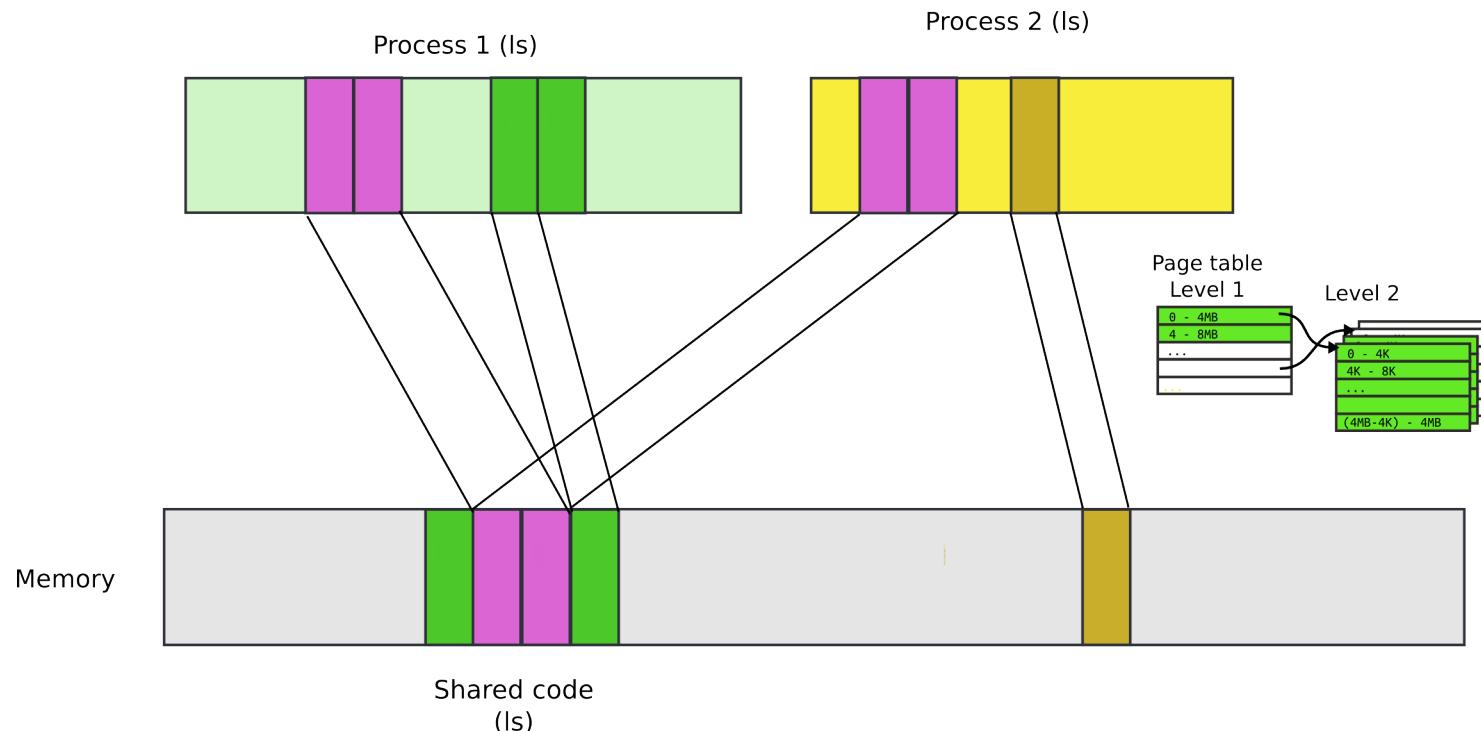


# Compared to segments pages allow ...

- Emulate large virtual address space on a smaller physical memory
  - In our example we had only 12 physical pages
  - But the program can access all 1M pages in its 4GB address space
  - The OS will move other pages to disk

# Compared to segments pages allow ...

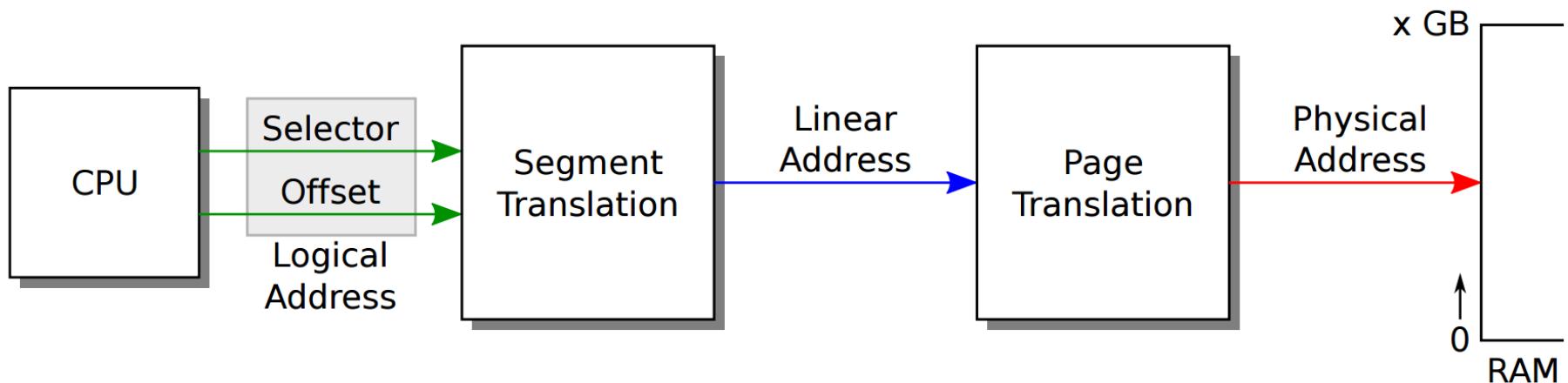
- Share a region of memory across multiple programs
  - Communication (shared buffer of messages)
  - Shared libraries

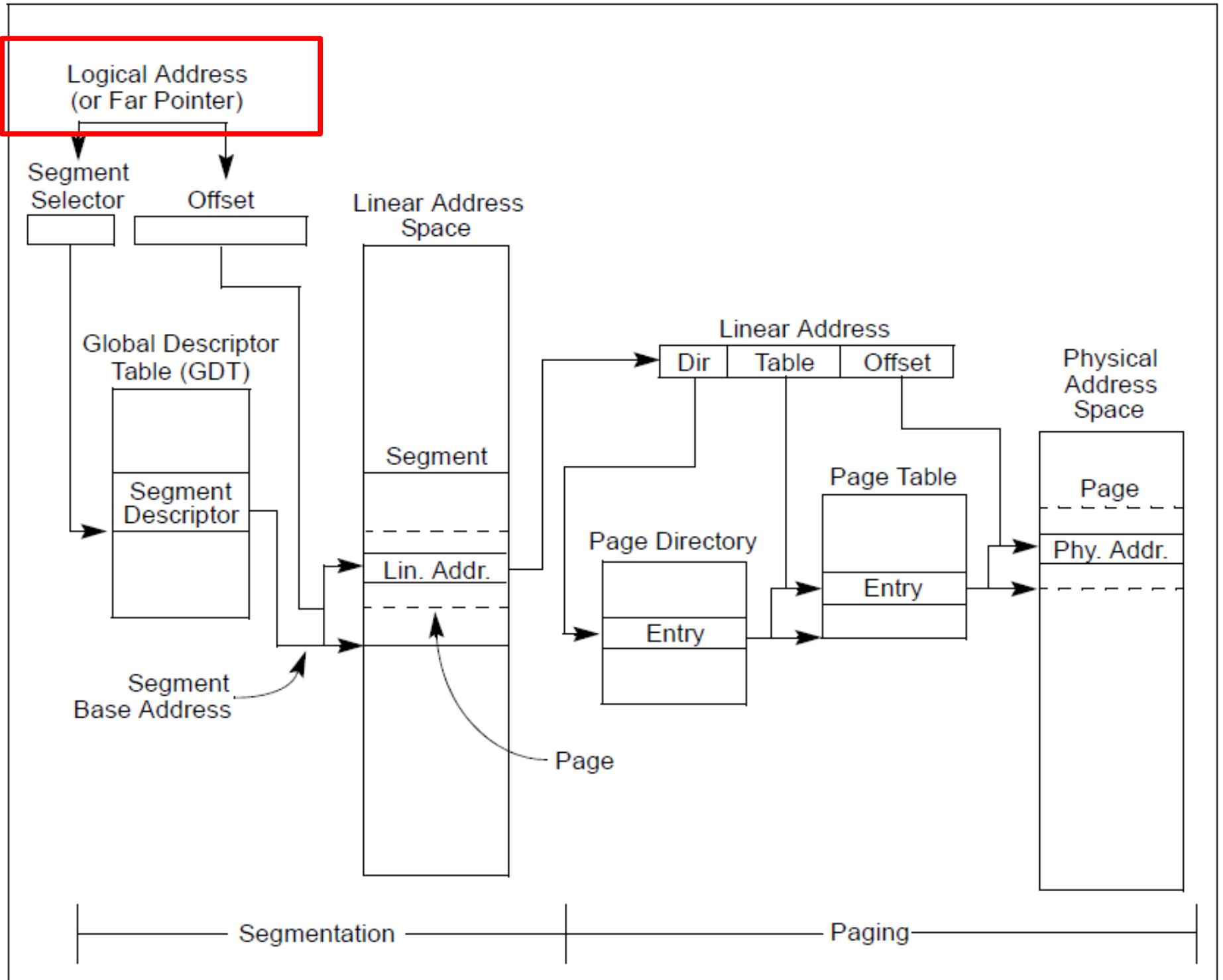


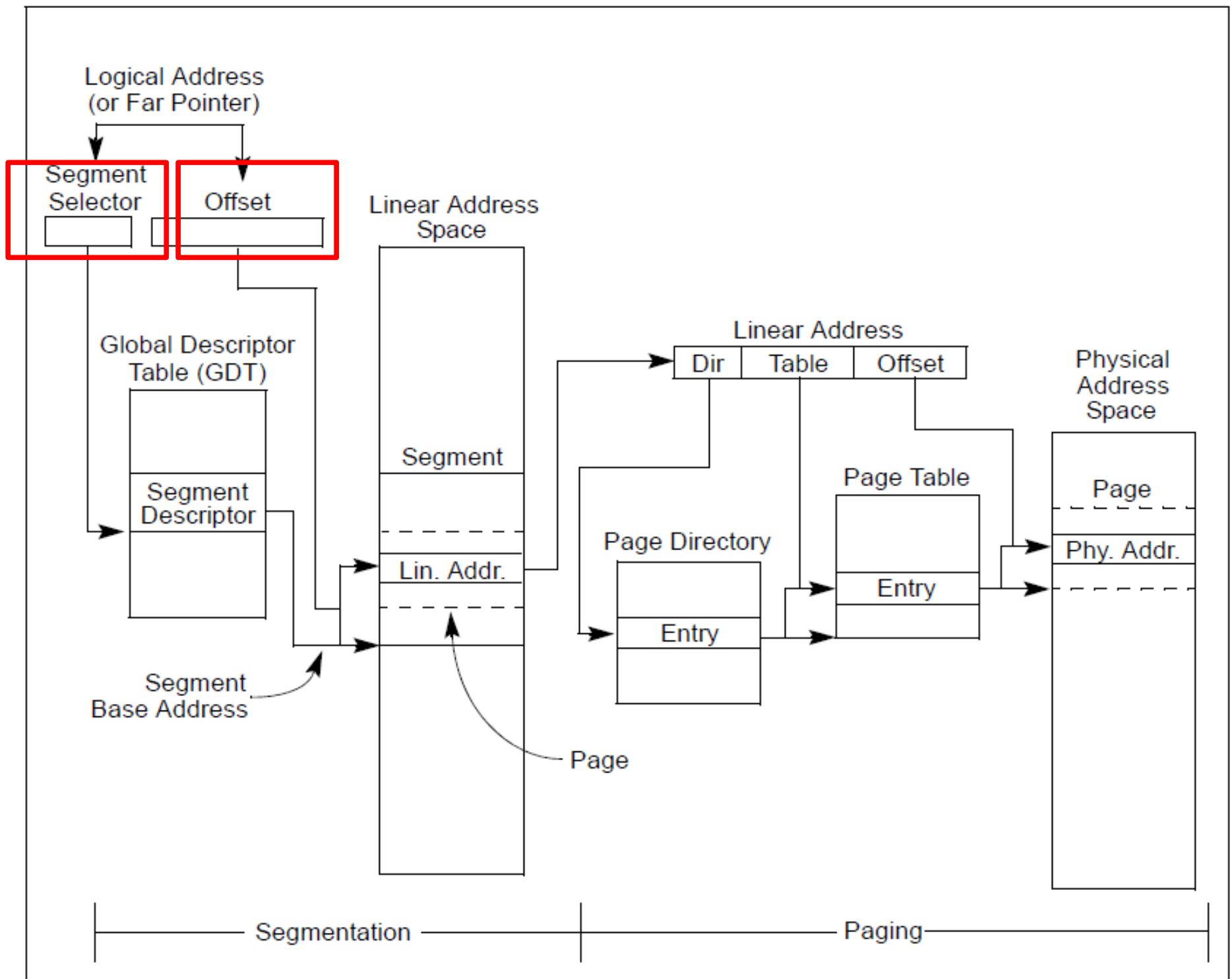
# More paging tricks

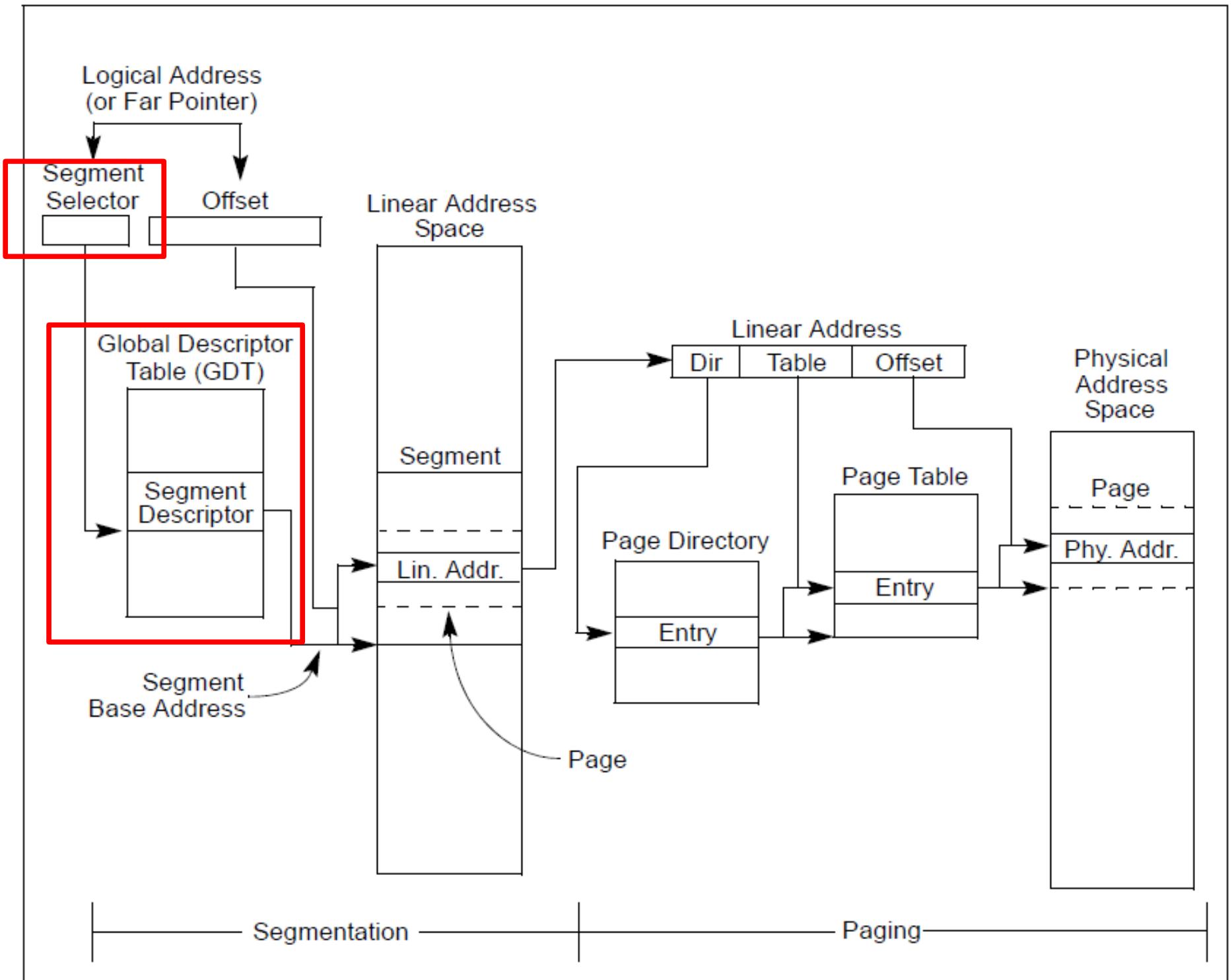
- Protect parts of the program
  - E.g., map code as read-only
    - Disable code modification attacks
    - Remember R/W bit in PTD/PTE entries!
  - E.g., map stack as non-executable
    - Protects from stack smashing attacks
    - Non-executable bit

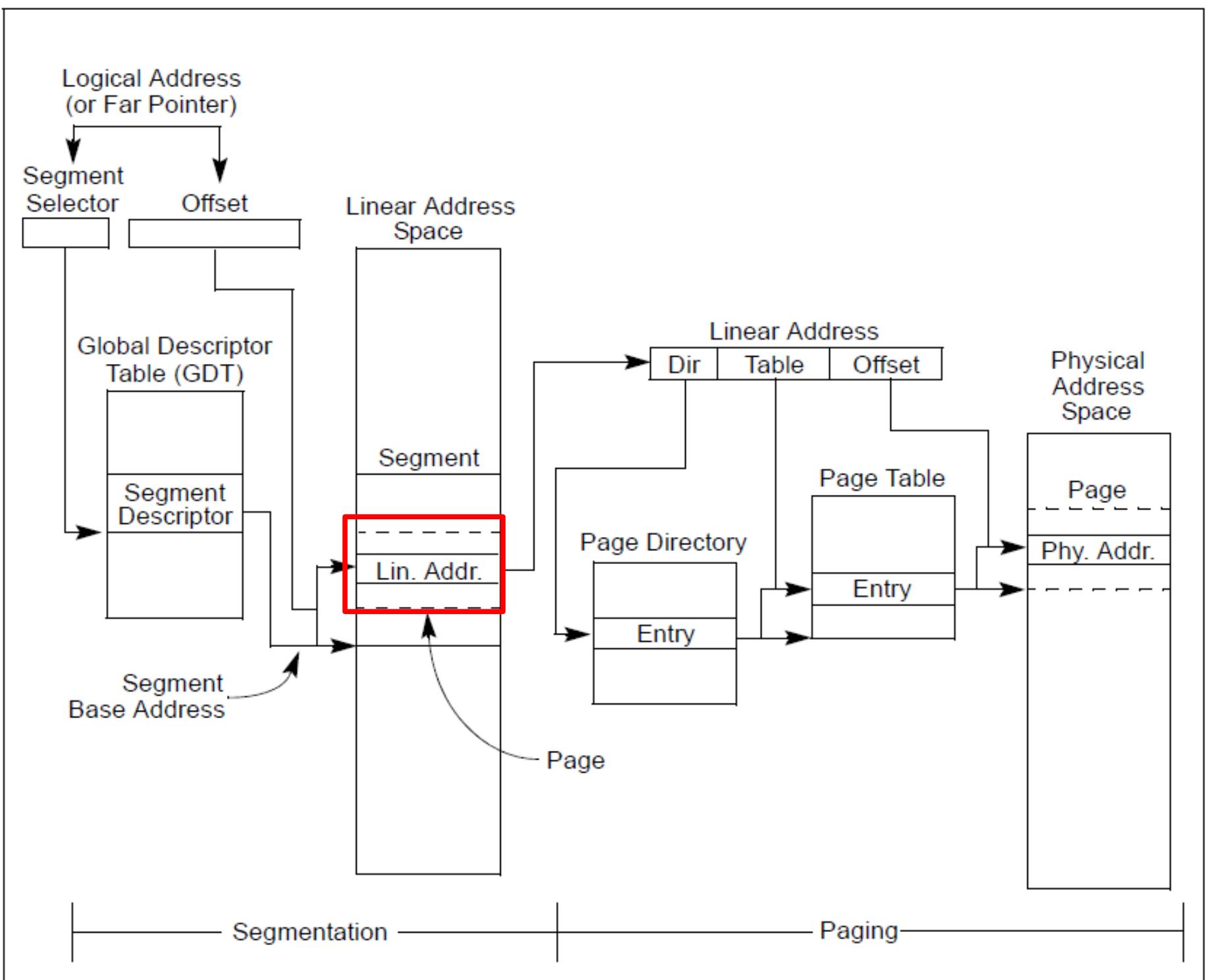
# Recap: complete address translation

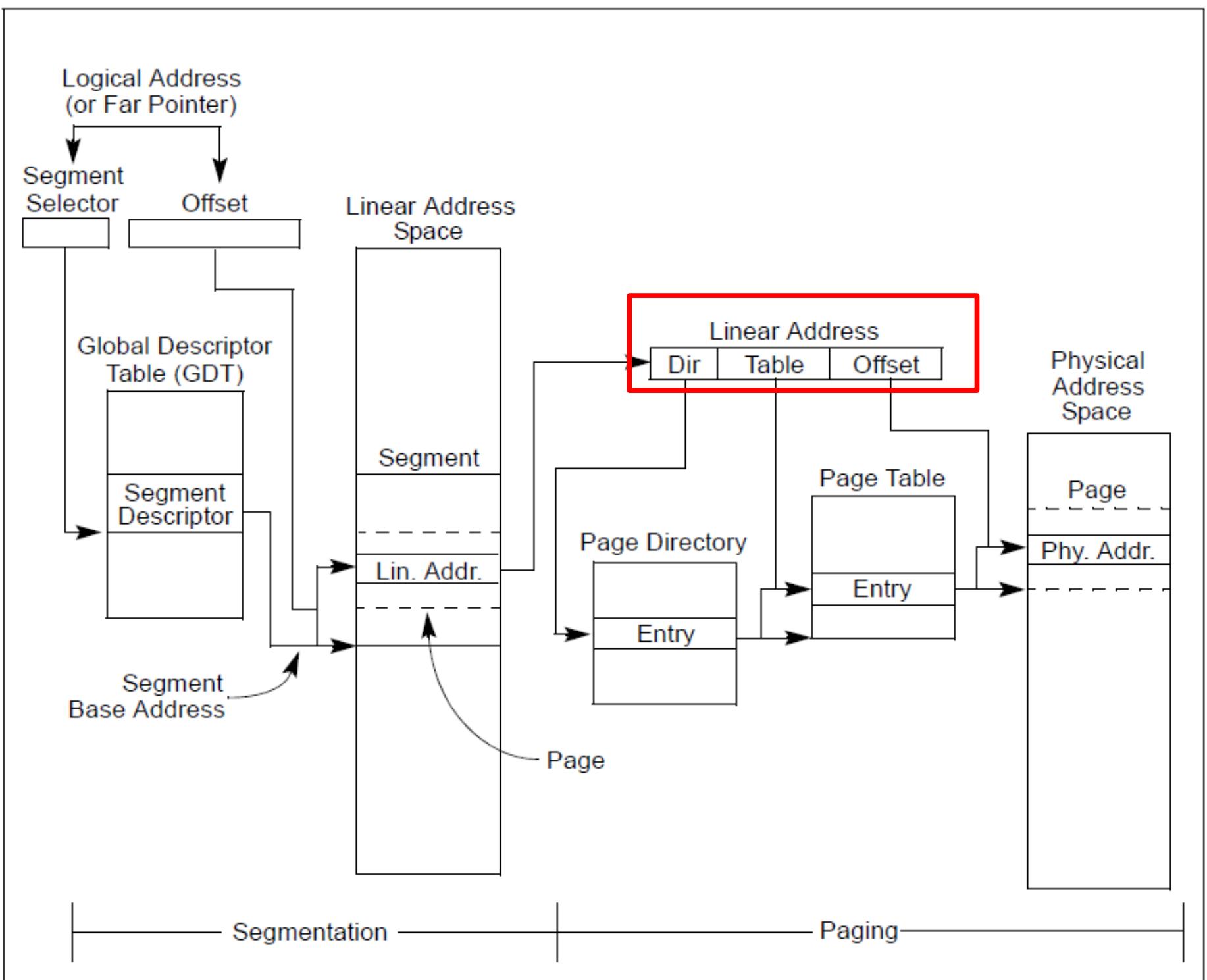


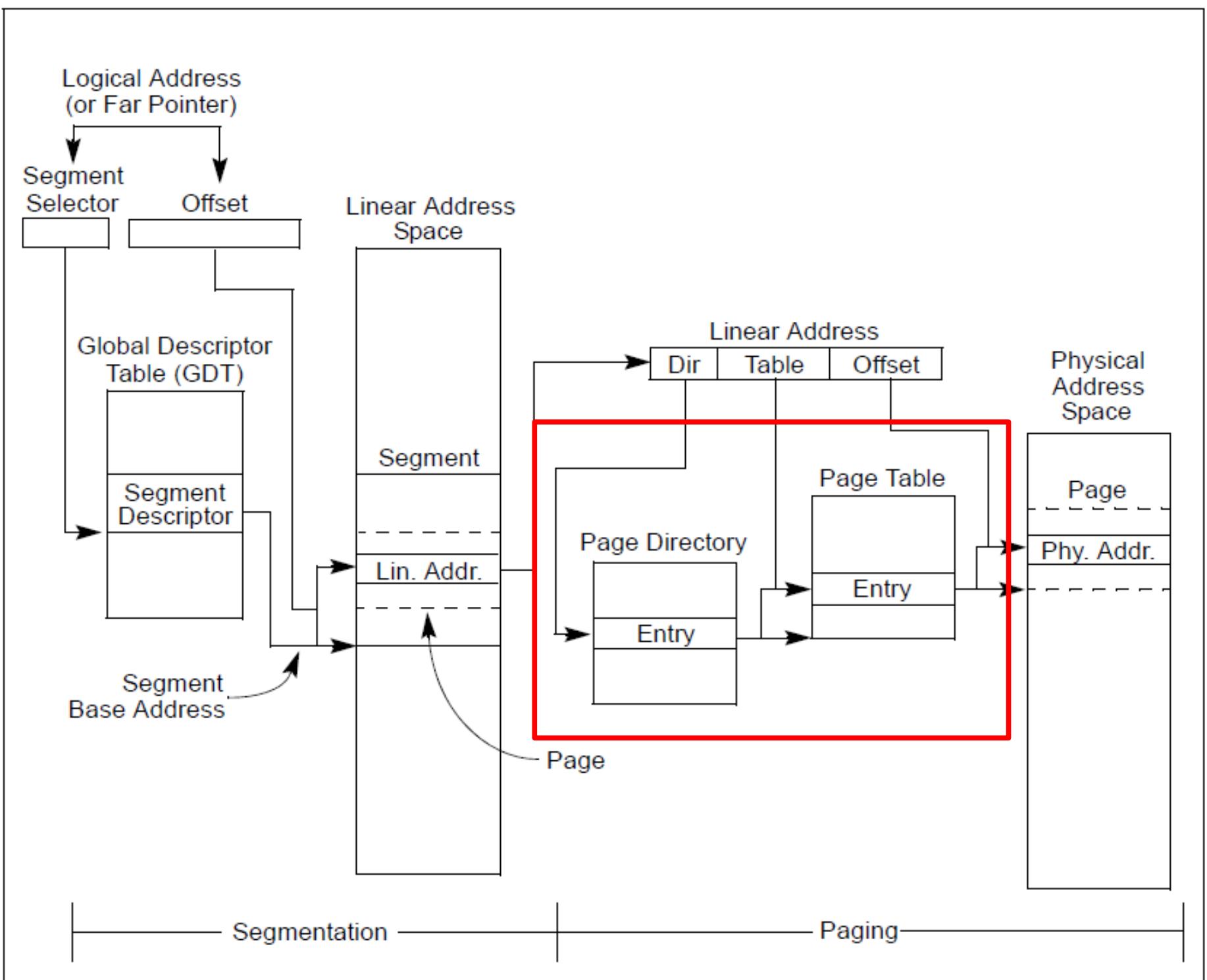


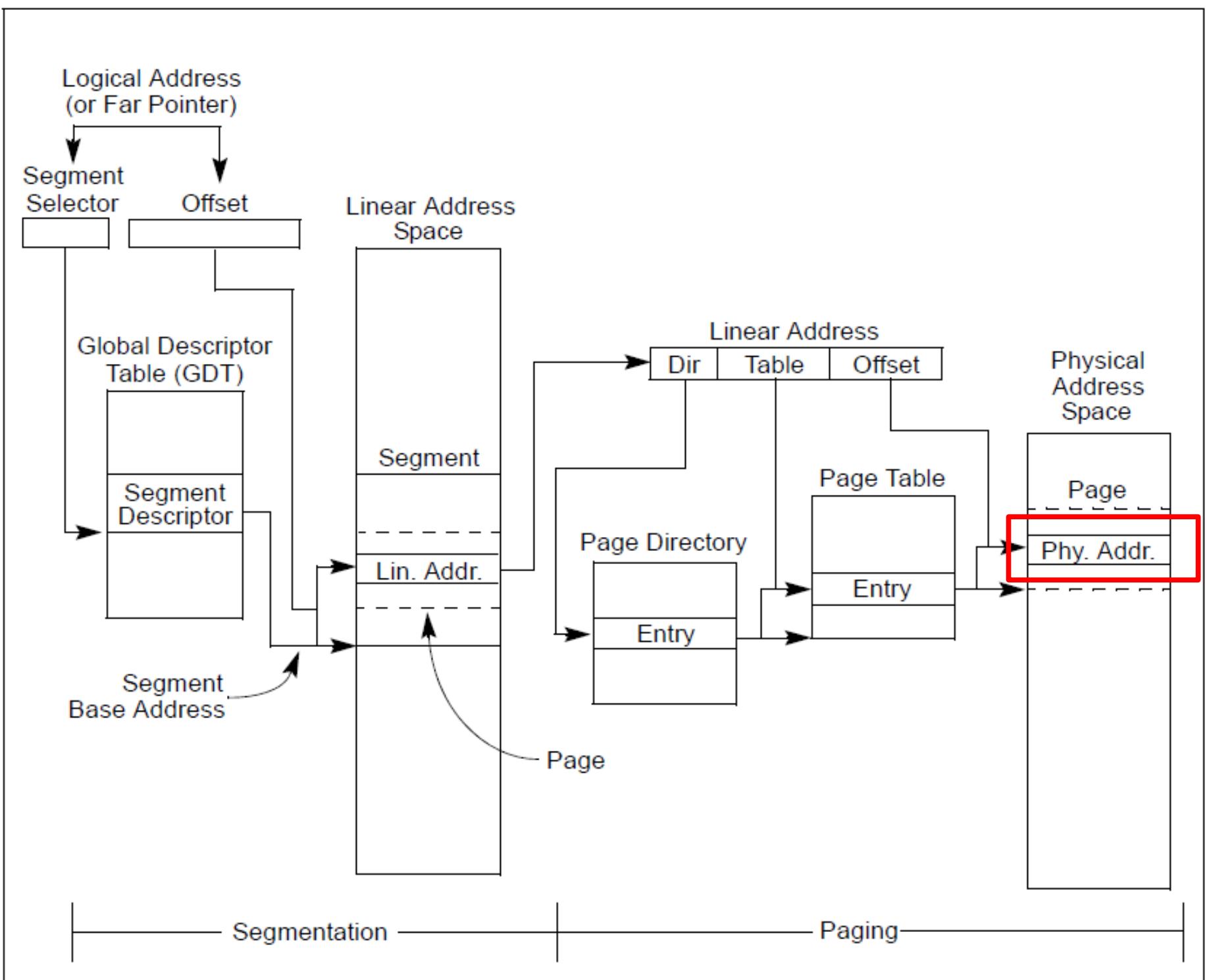


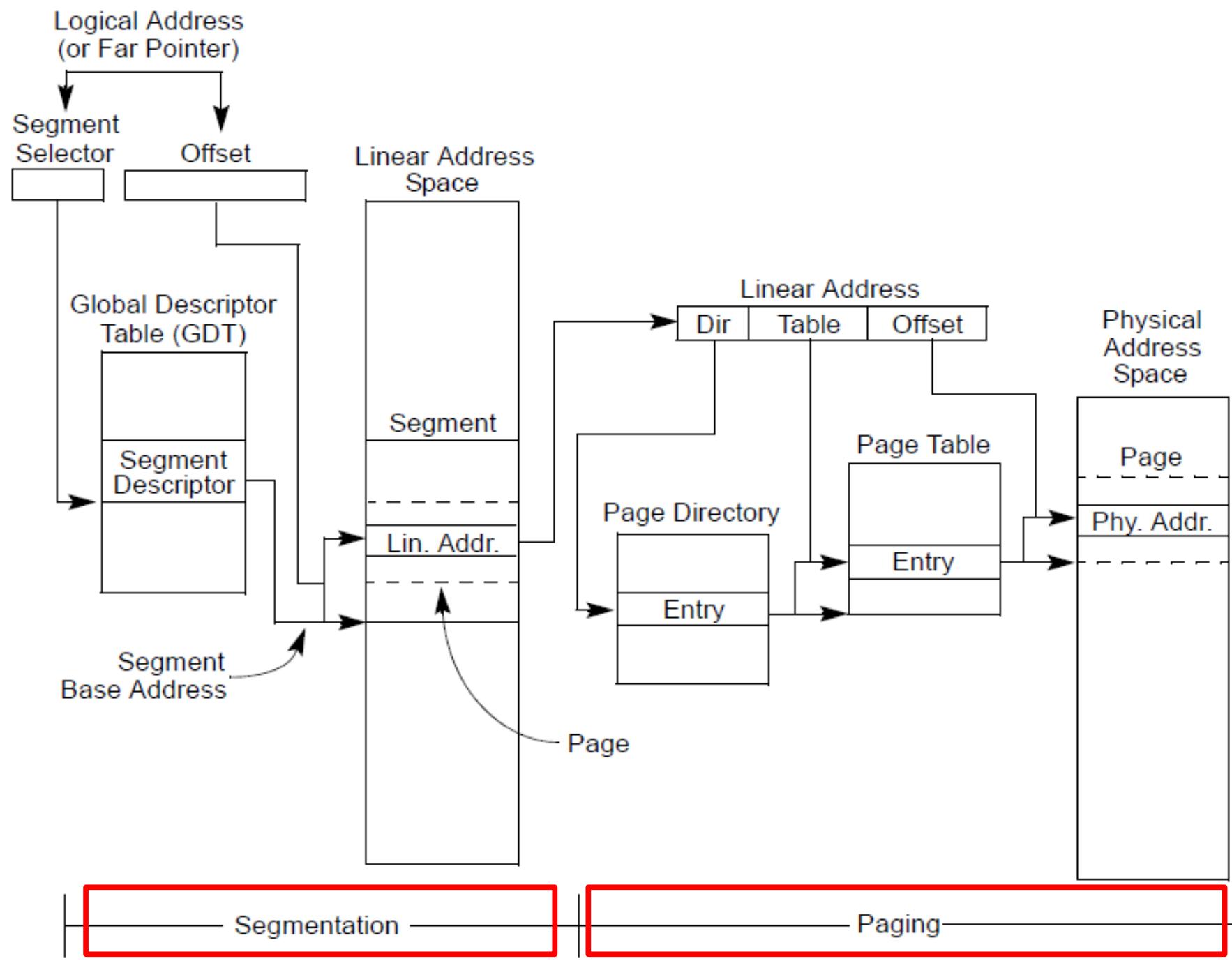












# Why do we need paging?

- Compared to segments pages provide fine-grained control over memory layout
  - No need to relocate/swap the entire segment
    - One page is enough
    -
- You're trading flexibility (granularity) for overhead of data structures required for translation

Questions?