

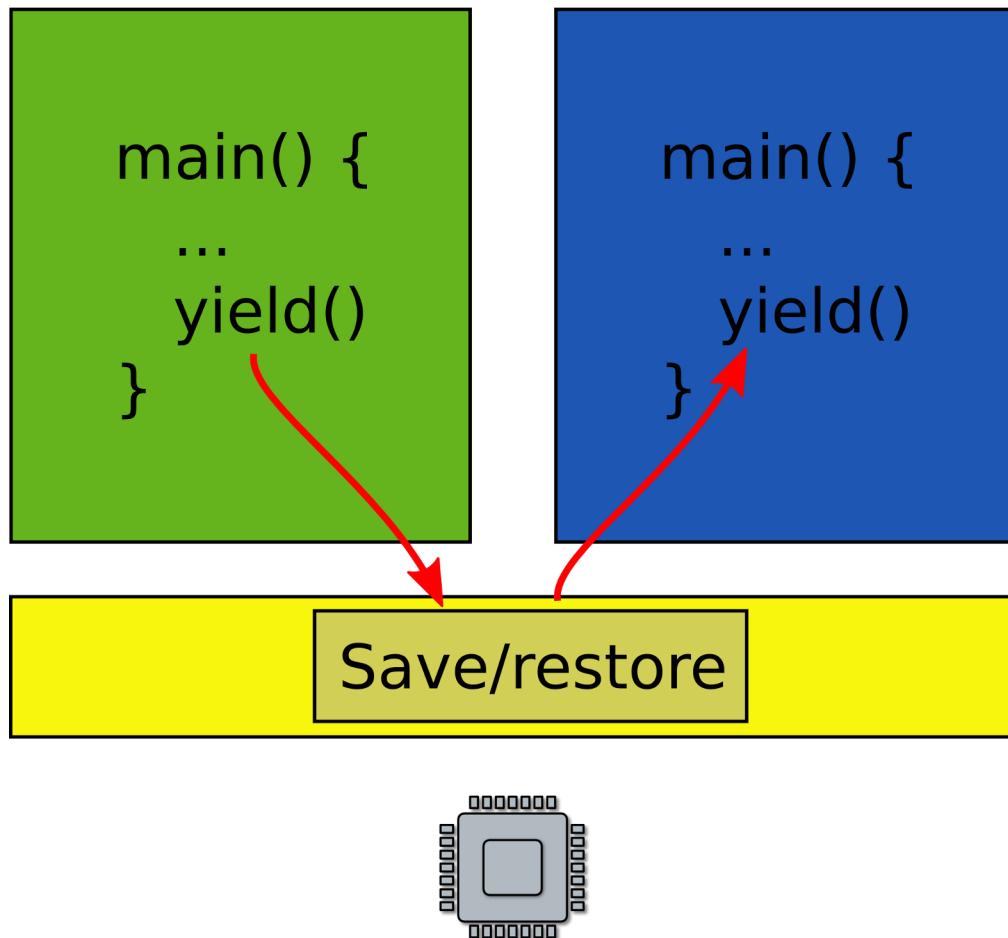
# cs5460/6460: Operating Systems

## Address translation (Segmentation and Paging)

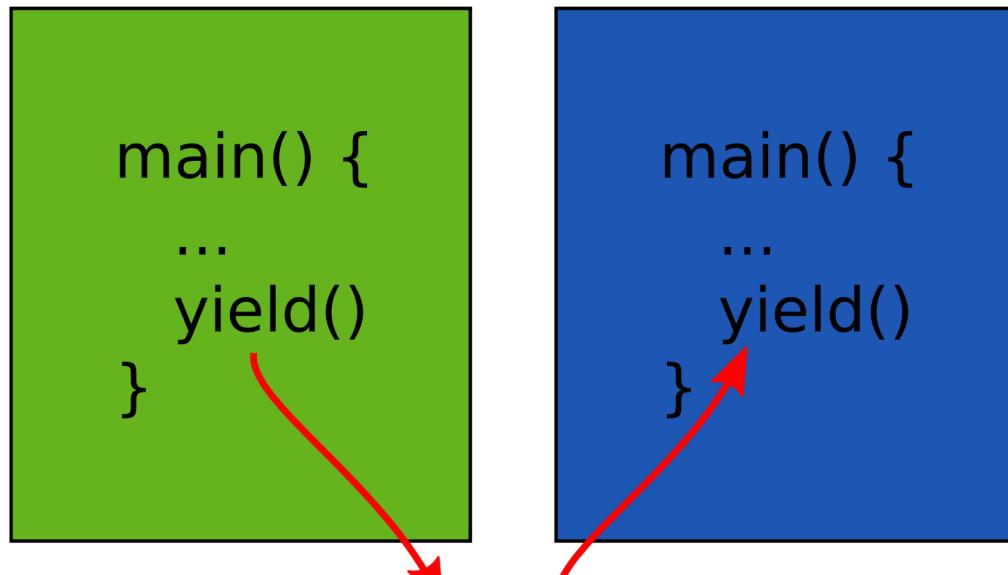
Anton Burtsev

February, 2024

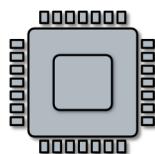
# Two programs one memory



# Two programs one memory



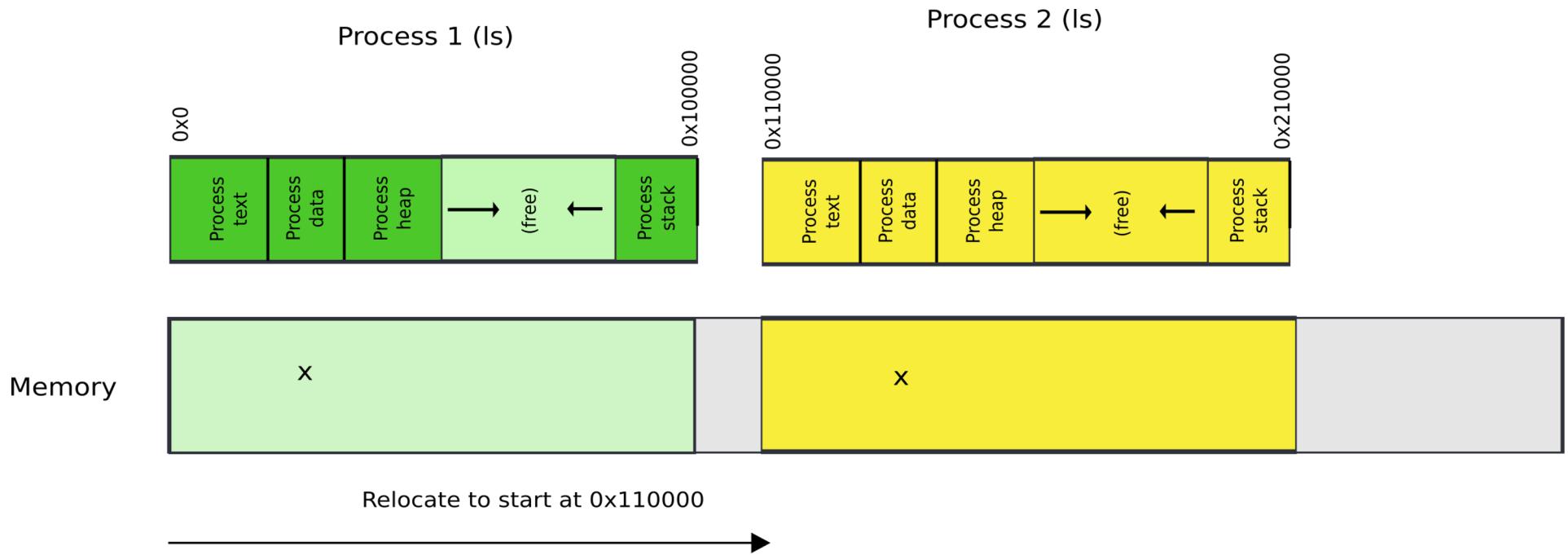
- How can we do this?



# Relocation

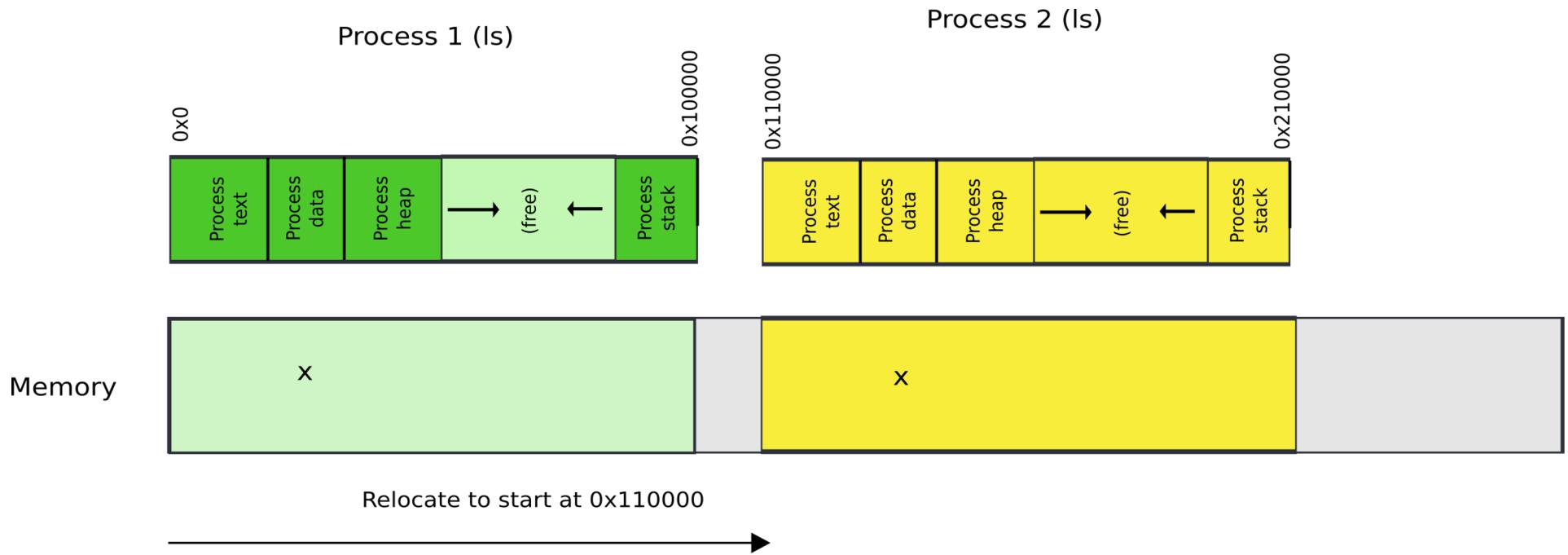
- One way to achieve this is to relocate program at different addresses
- Remember relocation (from linking and loading)

# Relocate binaries to work at different addresses



- One way to achieve this is to relocate program at different addresses
  - Remember relocation (from linking and loading)
  - **This works! But not ideal**
  - What is the problem?

# Relocate binaries to work at different addresses



- What is the problem?

# Problem: isolation

- How can we enforce isolation?

# Problem: isolation

- How can we enforce isolation?
- Isolation can be enforced in **software**
- Software Fault Isolation (SFI)
  - Google NaCl (Chrome Sandbox)
  - WASM (Web Assembly, another sandbox standard)

# Actually, how?

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

```
int main(int ac, char **av)
{
    int a = 5, b = 6;
    return a + b;
}
```

```
00000000 <main>:
```

0:	55	push	ebp
1:	89 e5	mov	ebp,esp
3:	83 ec 10	sub	esp,0x10
6:	c7 45 f8 05 00 00 00	mov	DWORD PTR [ebp-0x8],0x5
d:	c7 45 fc 06 00 00 00	mov	DWORD PTR [ebp-0x4],0x6
14:	8b 45 fc	mov	eax,DWORD PTR [ebp-0x4]
17:	8b 55 f8	mov	edx,DWORD PTR [ebp-0x8]
1a:	01 d0	add	eax,edx
1c:	c9	leave	
1d:	c3	ret	

# PollEv.com/aburtsev

```
#include <stdio.h>

int main(int ac, char **av)
{
    int a = 5, b = 6;
    return a + b;
}
```

00000000 <main>:		
0:	55	push ebp
1:	89 e5	mov ebp,esp
3:	83 ec 10	sub esp,0x10
6:	c7 45 f8 05 00 00 00	mov DWORD PTR [ebp-0x8],0x5
d:	c7 45 fc 06 00 00 00	mov DWORD PTR [ebp-0x4],0x6
14:	8b 45 fc	mov eax,DWORD PTR [ebp-0x4]
17:	8b 55 f8	mov edx,DWORD PTR [ebp-0x8]
1a:	01 d0	add eax,edx
1c:	c9	leave
1d:	c3	ret

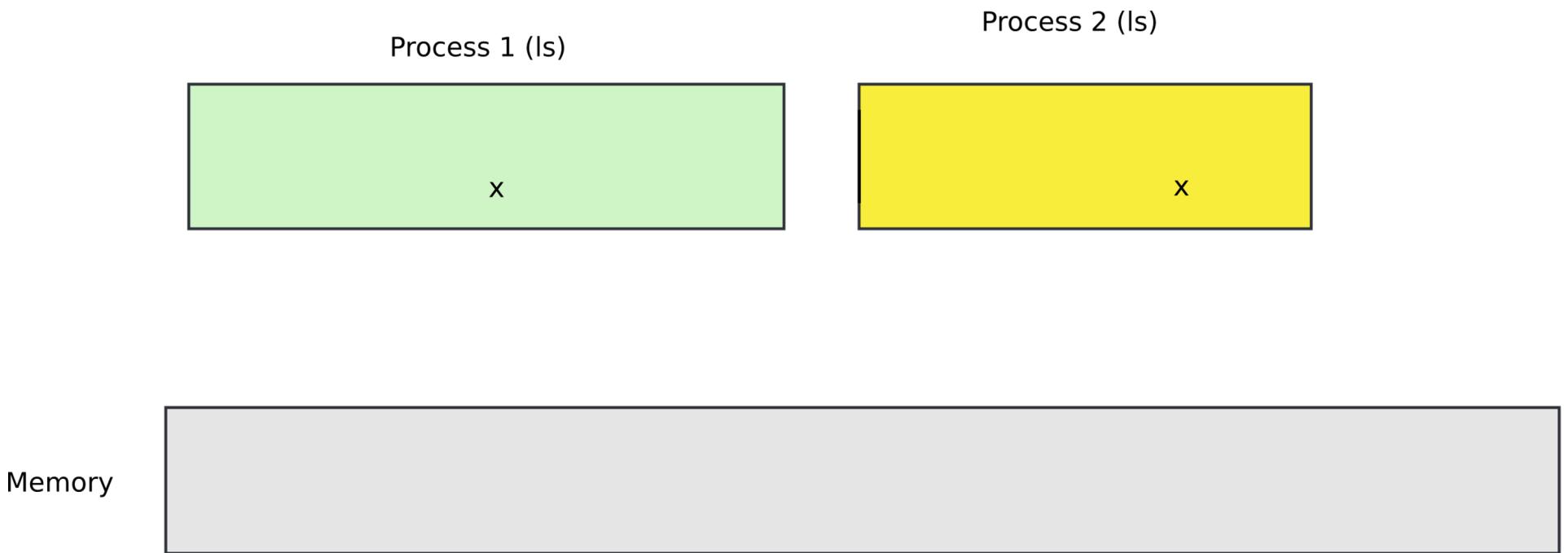
- Another way is to ask for hardware support

# Segmentation

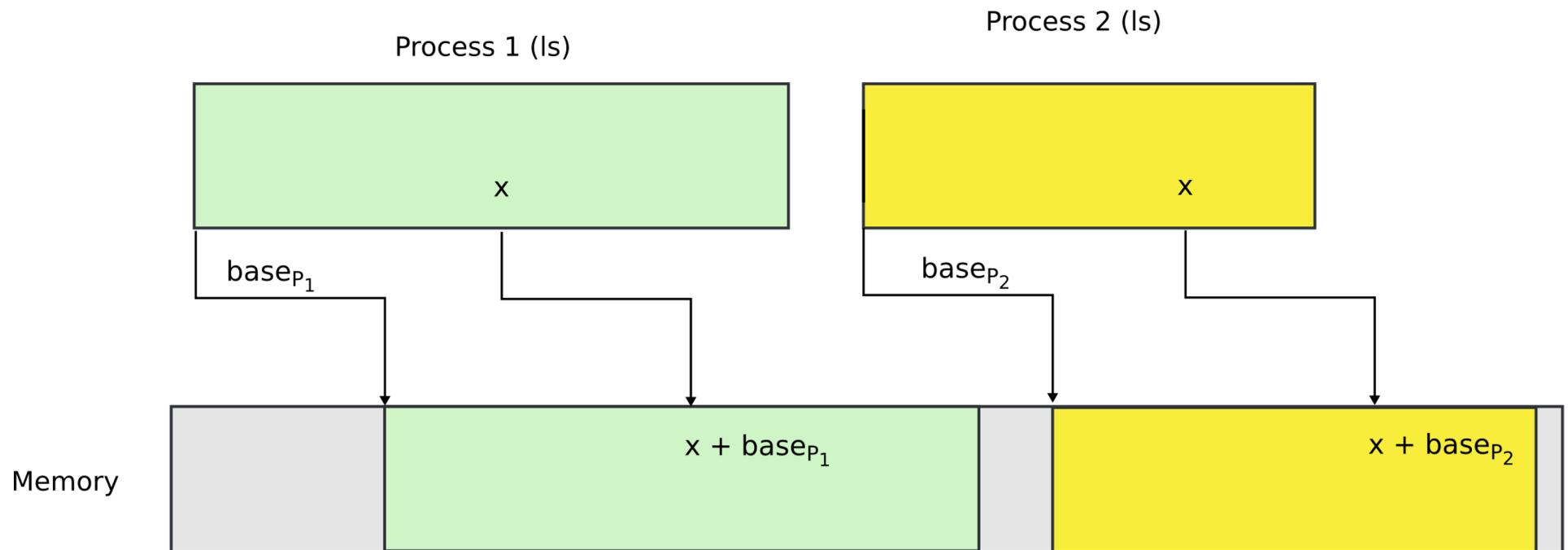
# What are we aiming for?

- Illusion of a private address space
- Identical copy of an address space in multiple programs
- 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

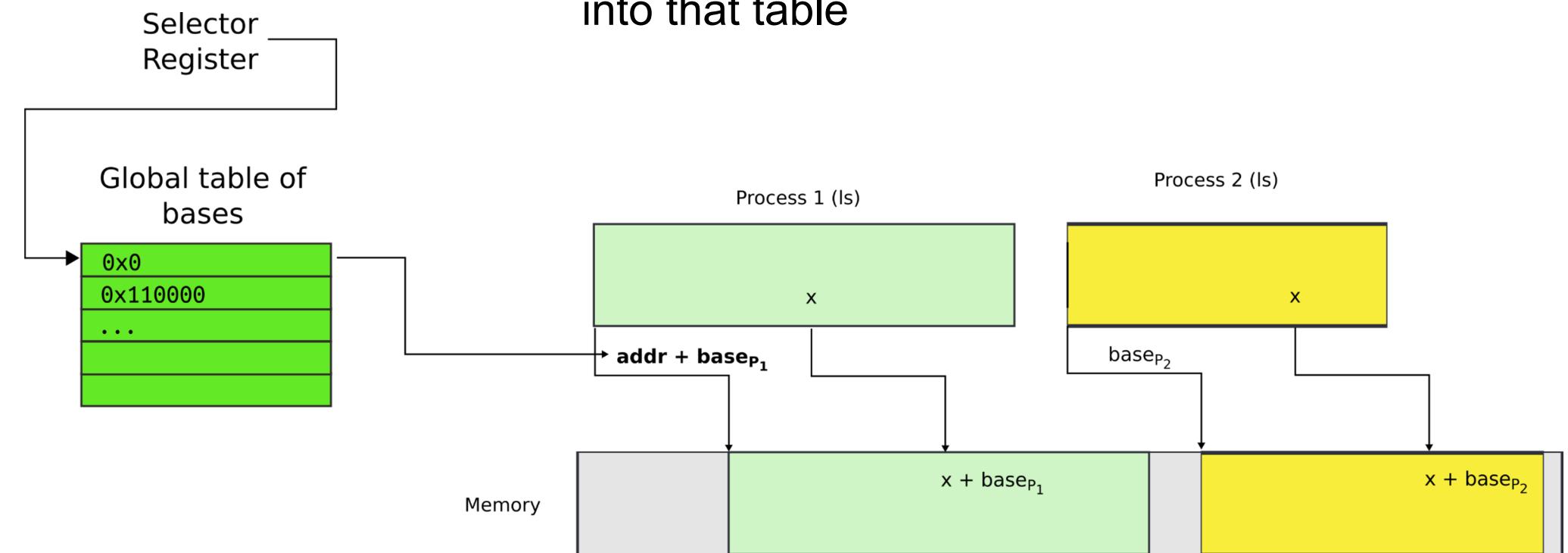
# Seems easy

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

# Seems easy

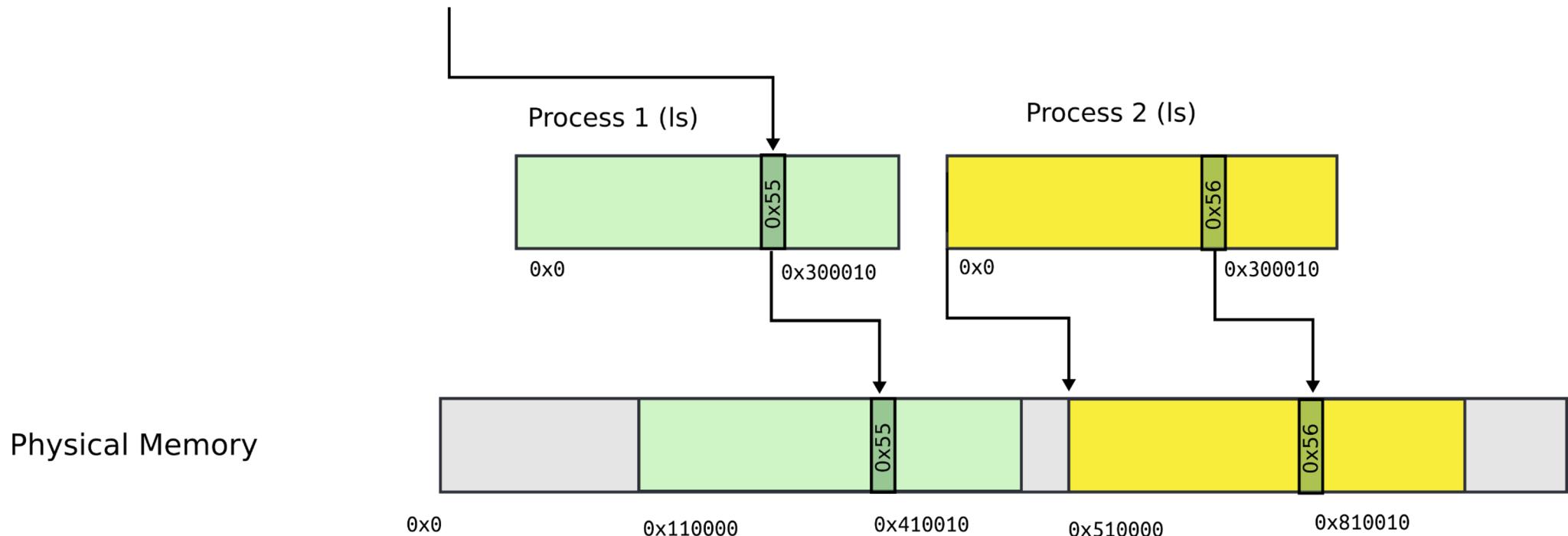
- One problem
- Where does this base address come from?
- Hardware can maintain a table of base addresses
  - One base for each process
  - Dedicate a special register to keep an index into that table

- 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



# Segmentation: example

```
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0x0  
EBX = 0x300010
```

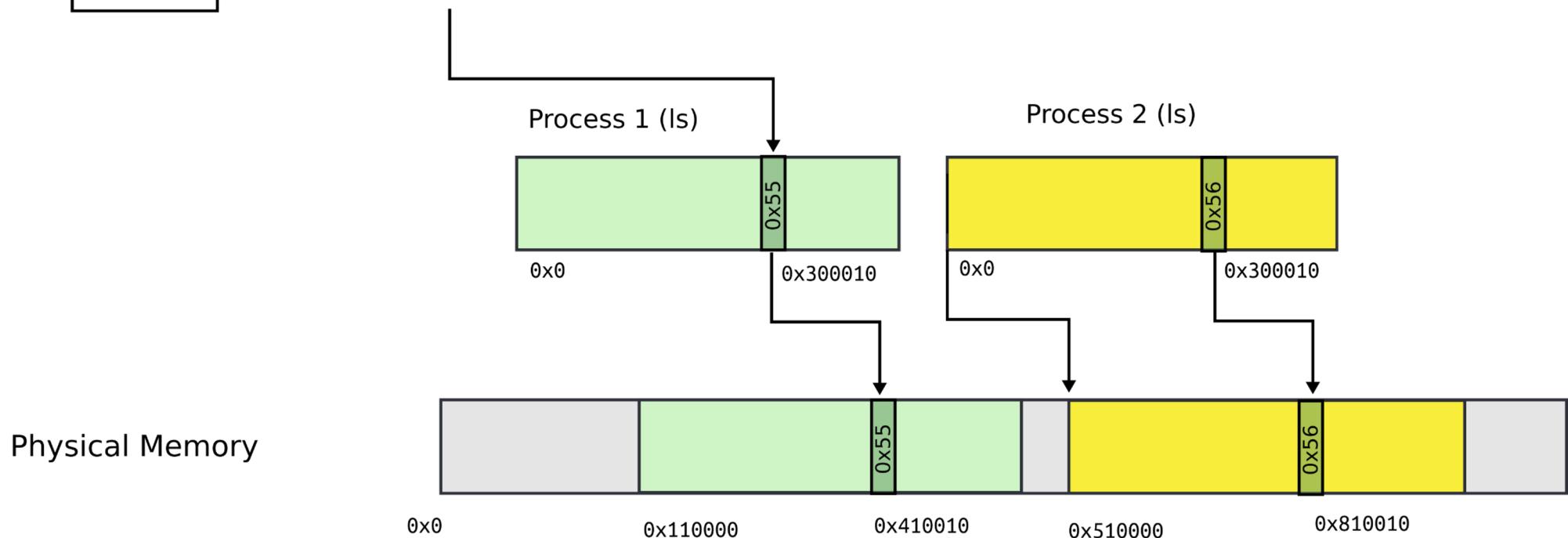


# Segmentation: address consists of two parts

Segment register  
(CS, SS, DS, ES, FS, GS)

0x1

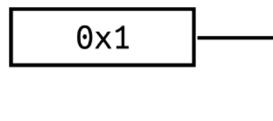
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0x0  
**EBX = 0x300010, DS = 0x1**



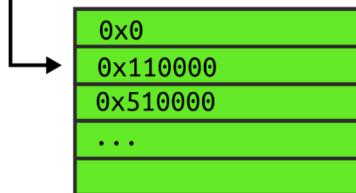
- Segment register contains segment selector
- General registers contain offsets
- Intel calls this address: “**logical address**”

# Segmentation: Global Descriptor Table

Segment register  
(CS, SS, DS, ES, FS, GS)

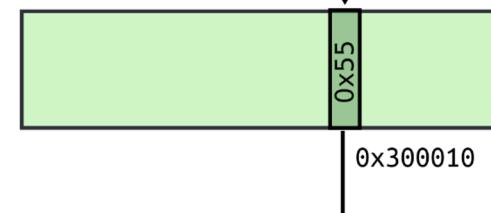


Global Descriptor Table  
(table of segment sizes and bases)

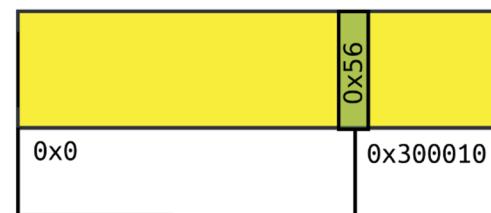


mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0x0  
EBX = 0x300010, DS = 0x1

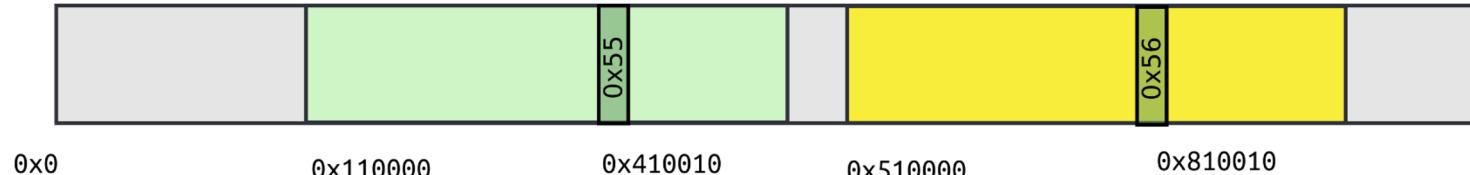
Process 1 (ls)



Process 2 (ls)

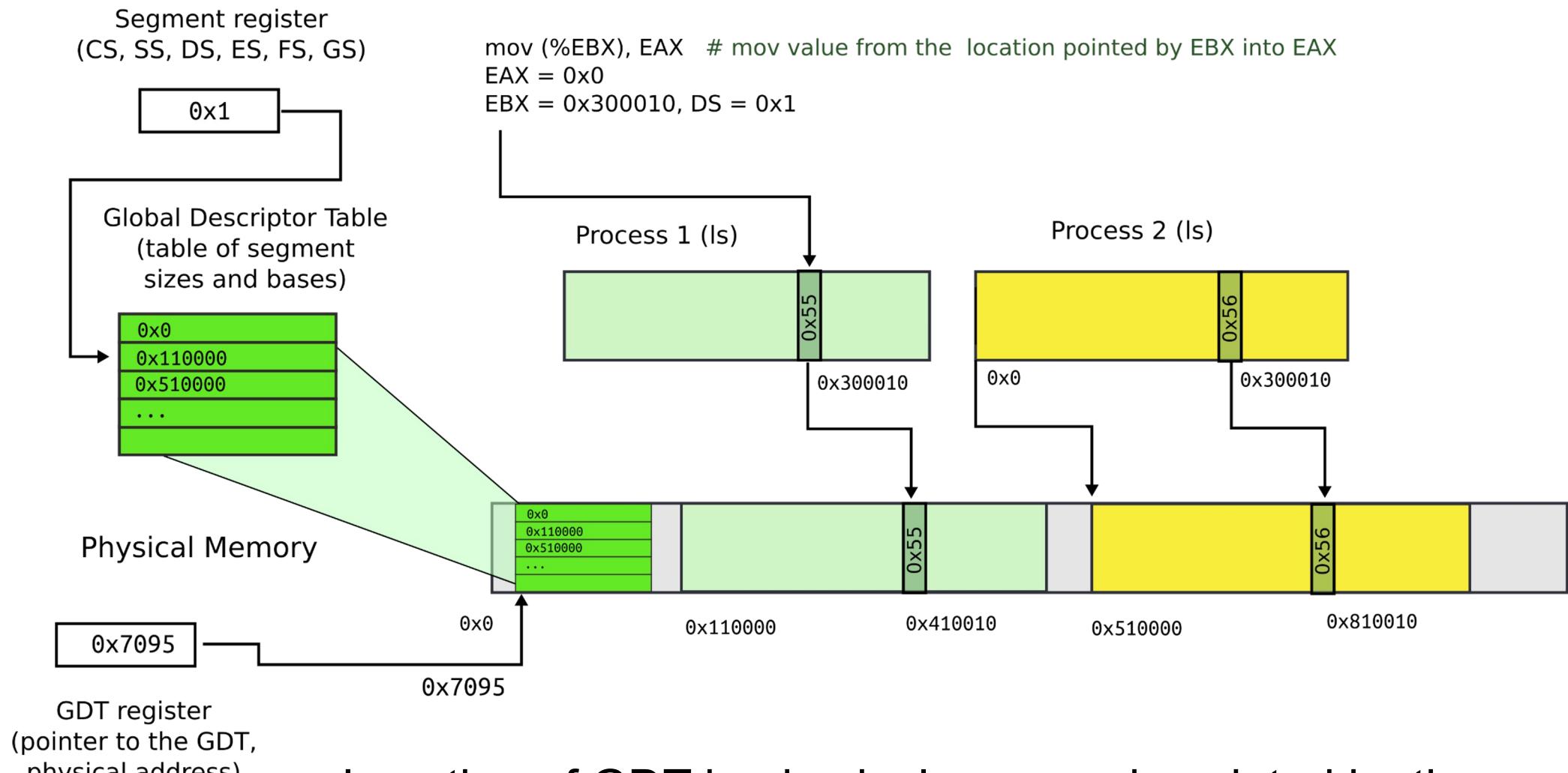


Physical Memory



- GDT is an array of segment descriptors
- Each descriptor contains base and limit for the segment
- Plus access control flags

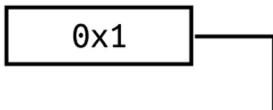
# Segmentation: Global Descriptor Table



- Location of GDT in physical memory is pointed by the GDT register

# Segmentation: base + offset

Segment register  
(CS, SS, DS, ES, FS, GS)



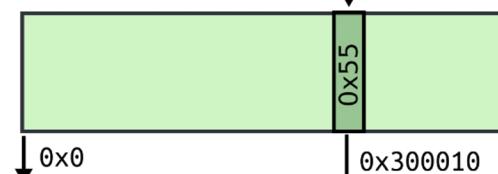
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0x0  
EBX = 0x300010, DS = 0x1

Global Descriptor Table  
(table of segment sizes and bases)

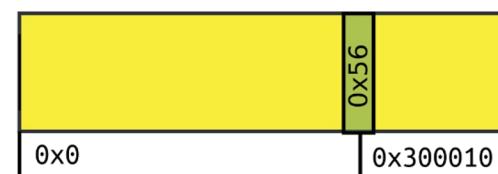
0x0
0x110000
0x510000
...
0x0

base<sub>P1</sub>

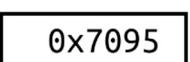
Process 1 (ls)



Process 2 (ls)



Physical Memory



0x0  
0x7095

0x110000

0x410010

0x510000

0x810010

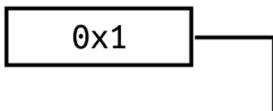
GDT register  
(pointer to the GDT,  
physical address)

0x0  
0x7095

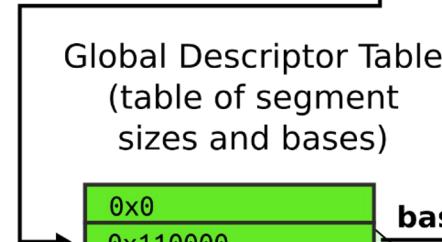
- Segment register (0x1) chooses an entry in GDT
- This entry contains base of the segment (0x110000) and limit (size) of the segment (not shown)

# Segmentation: base + offset

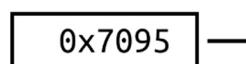
Segment register  
(CS, SS, DS, ES, FS, GS)



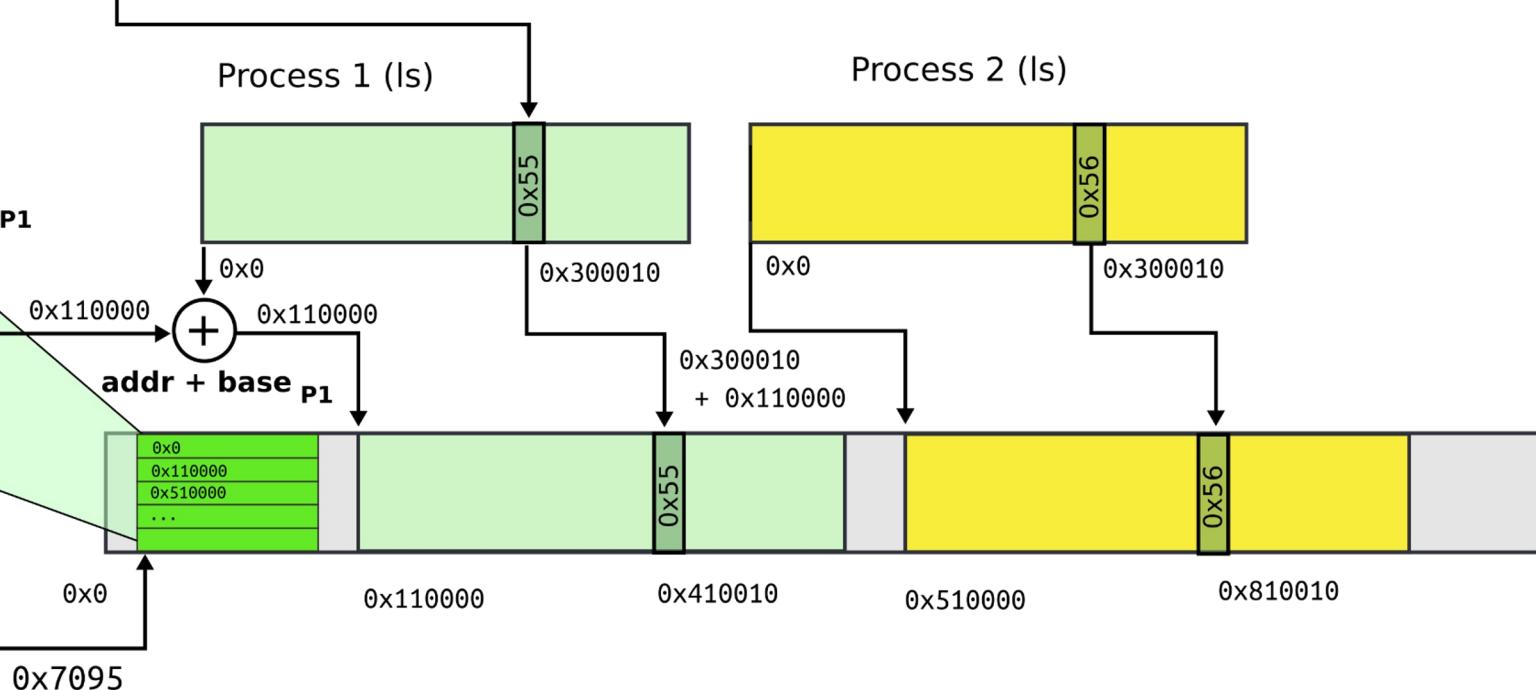
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0x0  
EBX = 0x300010, DS = 0x1



Physical Memory



GDT register  
(pointer to the GDT,  
physical address)



- Physical address:
- $0x410010 = 0x300010 \text{ (offset)} + 0x110000 \text{ (base)}$
- Intel calls this address “linear”

# Segmentation: process 2

Segment register  
(CS, SS, DS, ES, FS, GS)

0x1

mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0x0  
EBX = 0x300010, DS = 0x1

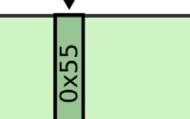
Global Descriptor Table  
(table of segment  
sizes and bases)

0x0	base <sub>P1</sub>
0x510000	
...	
...	

Process 1 (ls)

0x0

0x300010



Process 2 (ls)

0x0

0x300010



Physical Memory

0x890000

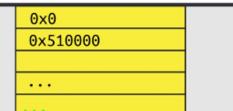
0x0  
0x890000

0x110000

0x410010

0x510000

0x810010



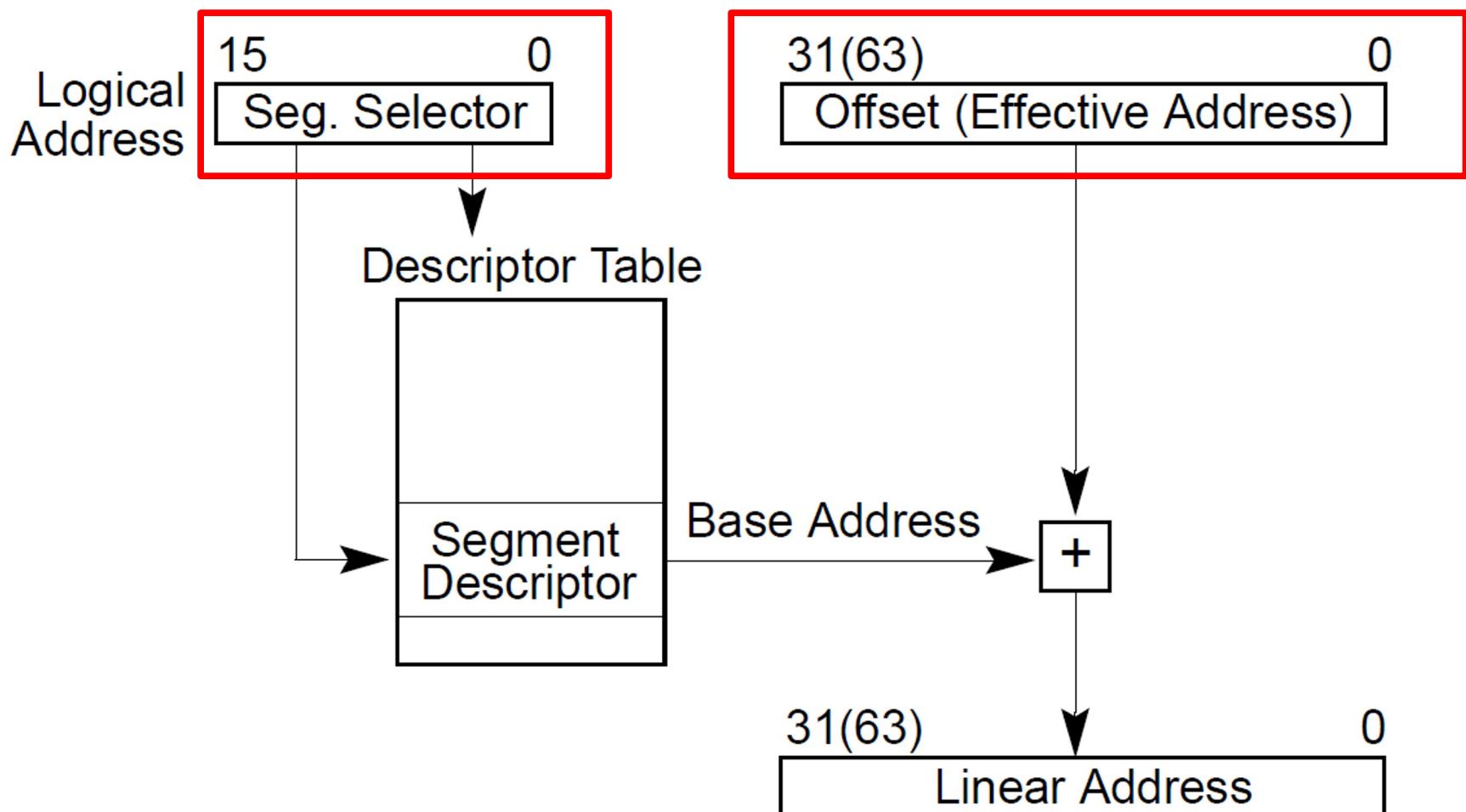
GDT register  
(pointer to the GDT,  
physical address)

- Each process has a **private GDT**
- OS switches between GDTs

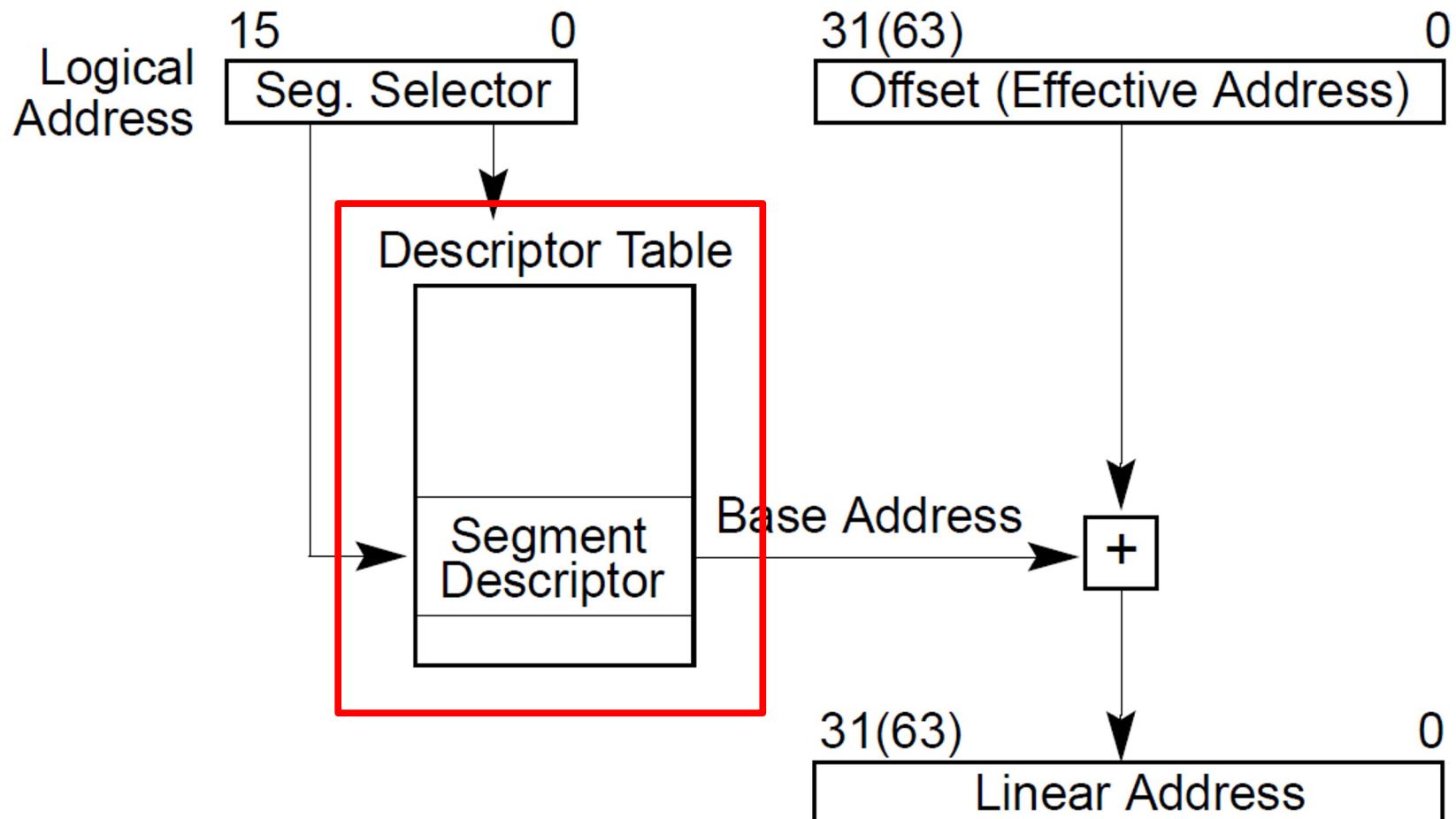
New addressing mode:  
“logical addresses”

# All addresses are logical address

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

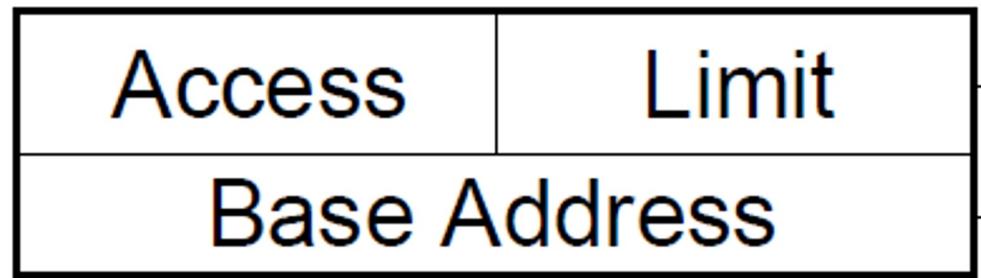


- Segment selector (16 bit)
- Is simply an index into an array (Descriptor Table)
- That holds segment descriptors
  - Base and limit (size) for each segment

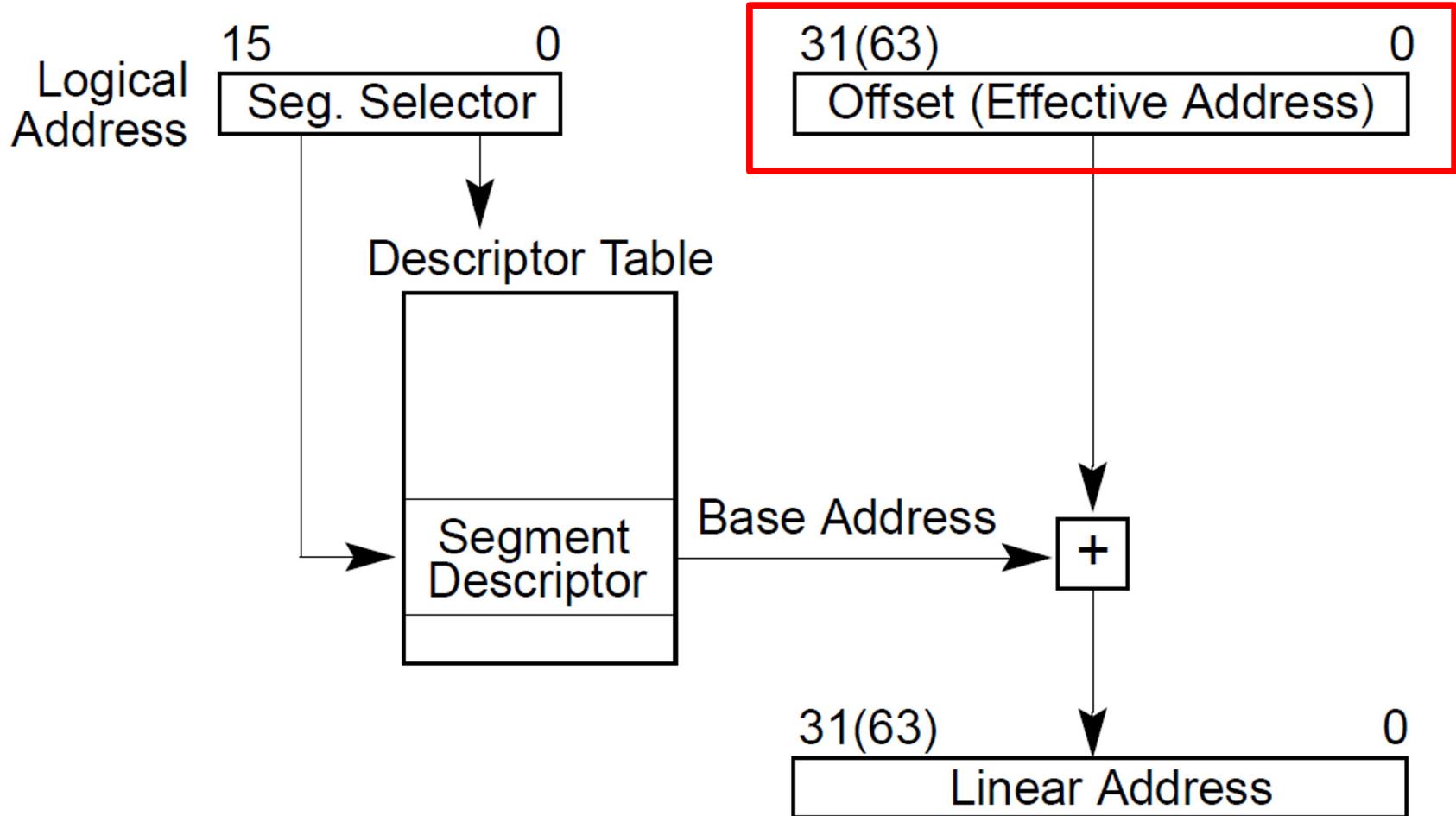


# Elements of the descriptor table are segment descriptors

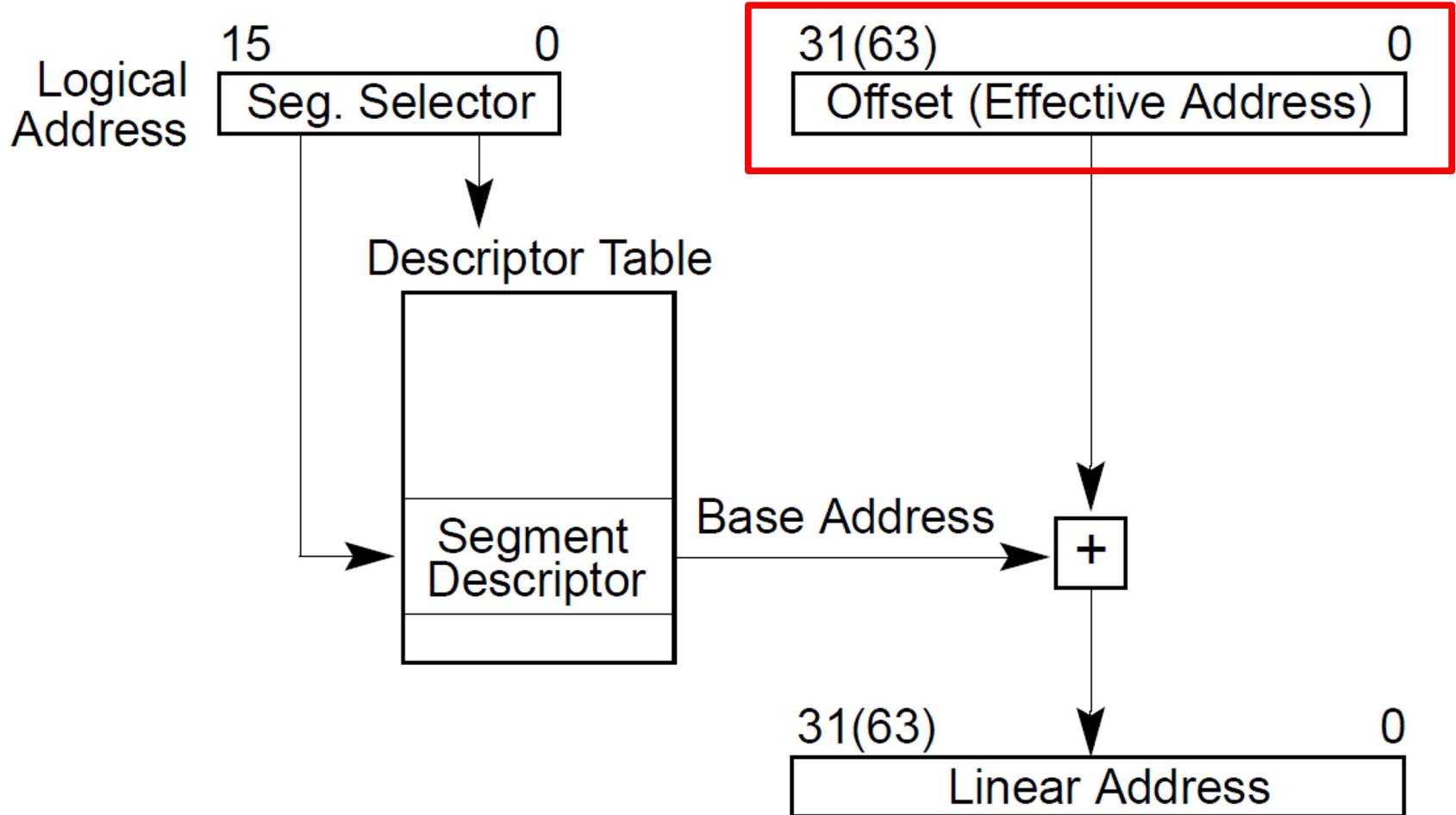
- Base address
  - 0 – 4 GB
- Limit (size)
  - 0 – 4 GB
- Access rights
  - Executable, readable, writable
- Privilege level (0 - 3)



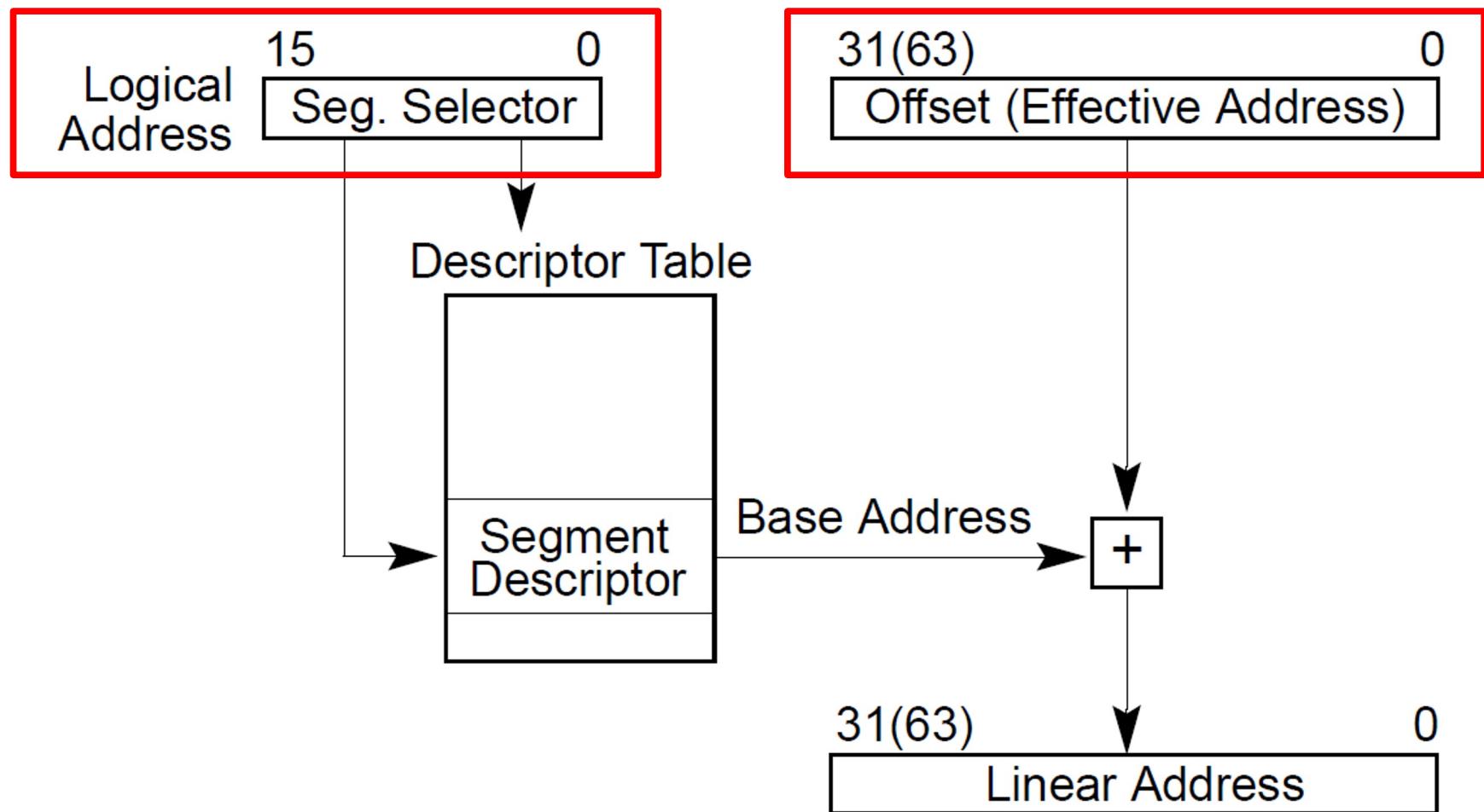
- Offsets into segments (x in our example) or “**effective addresses**” are in registers



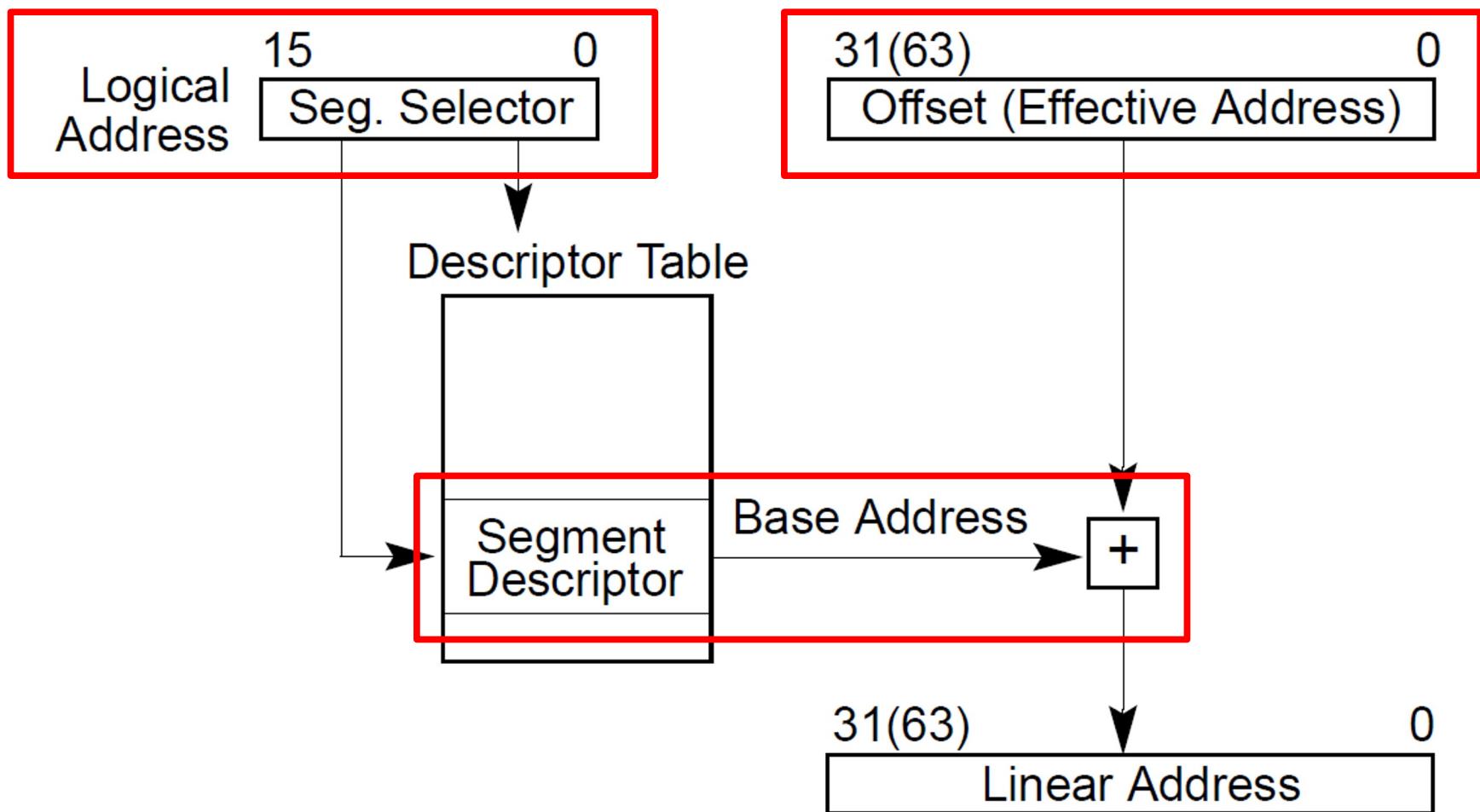
- Logical addresses are translated into physical
  - $\text{Effective address} + \text{DescriptorTable}[selector].Base$



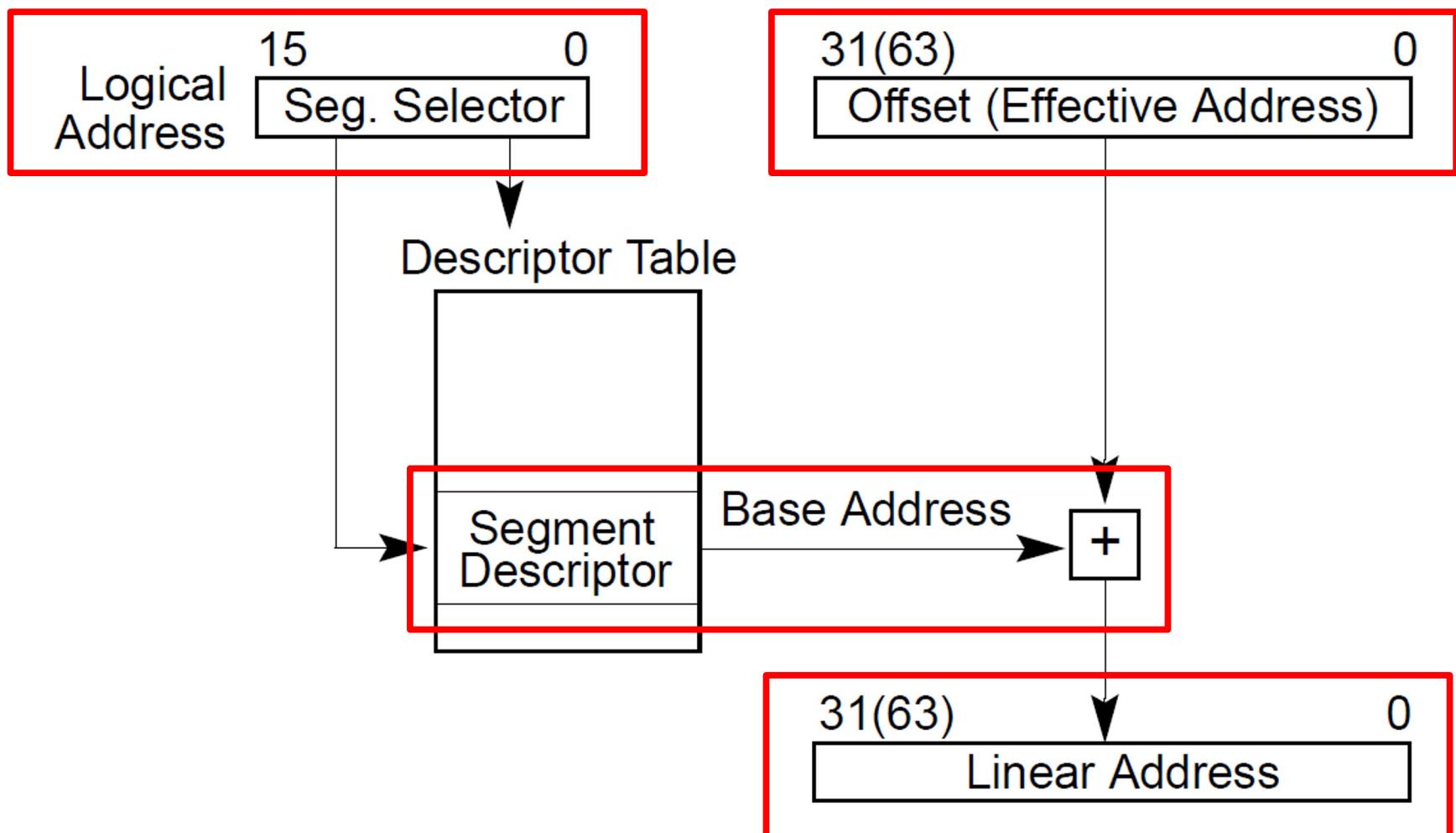
- Logical addresses are translated into physical
  - $\text{Effective address} + \text{DescriptorTable}[selector].Base$



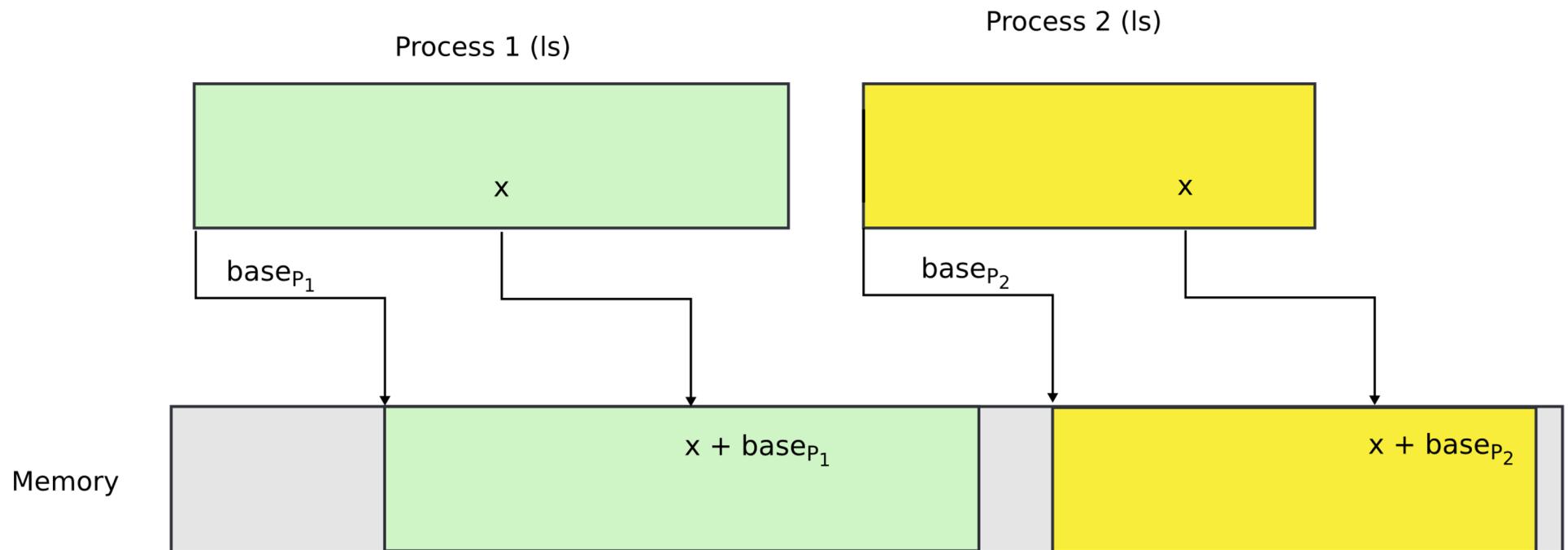
- Logical addresses are translated into physical
  - $\text{Effective address} + \text{DescriptorTable}[selector].Base$



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



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



# Segment registers

- Hold 16 bit segment selectors
  - Indexes into GDT
- Segments are associated with one of three types of storage
  - Code
  - Data
  - Stack

# Programing with segements (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 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.

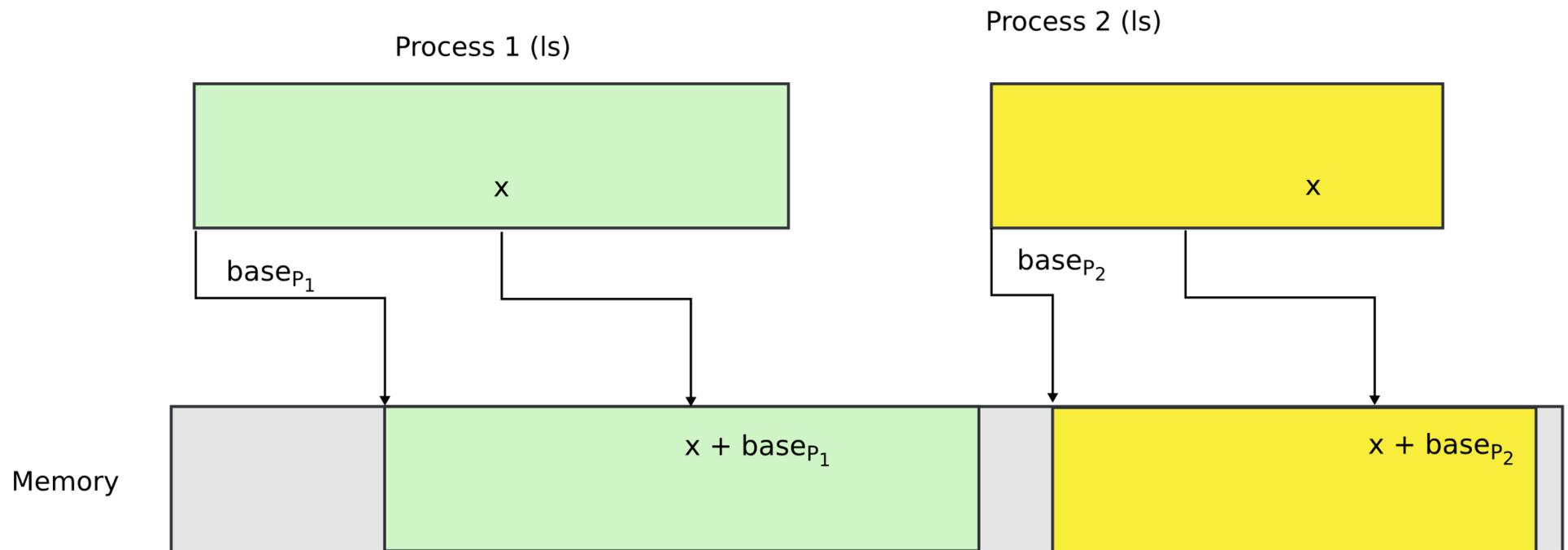
- This is cumbersome,
- Instead the idea is: infer code, data and stack segments from the instruction type
  - Control-flow instructions use code segment (jump, call)
  - Stack management (push/pop) uses stack
  - Most loads/stores use data segment
- Extra segments (es, fs, gs) must be used explicitly

# Segmentation: what did we achieve

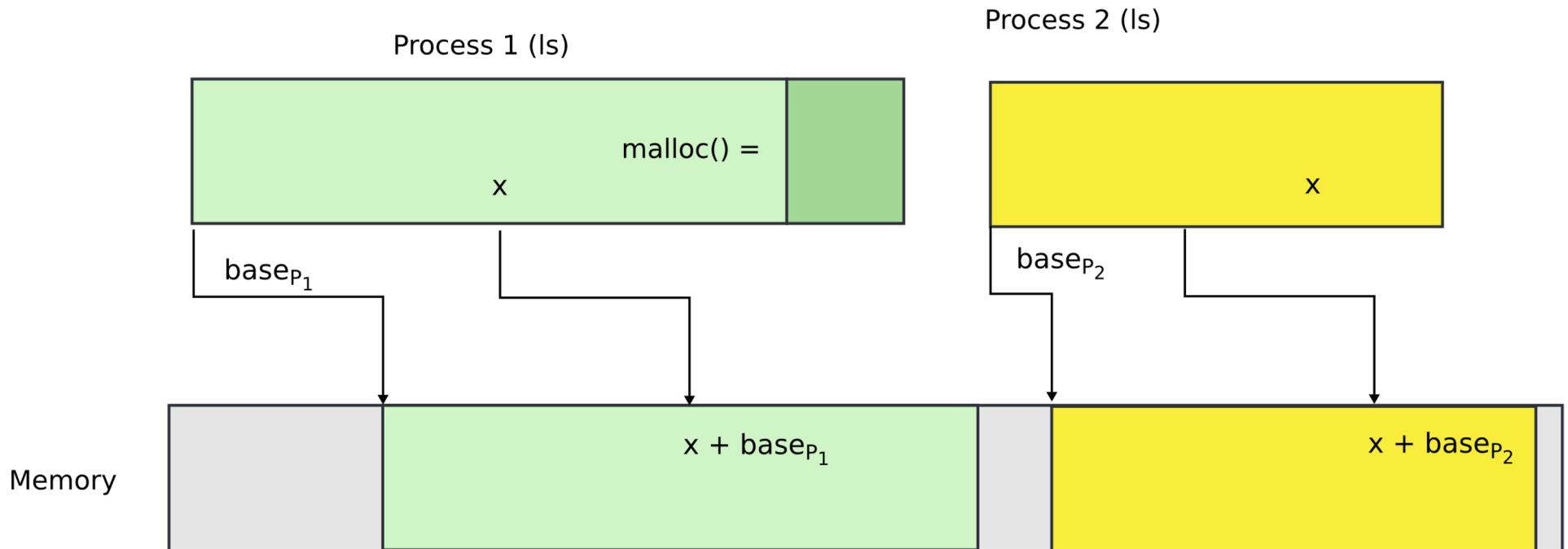
- Illusion of a private address space
- Identical copy of an address space in multiple programs
  - We can implement fork()
- Isolation
  - Processes cannot access memory outside of their segments

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

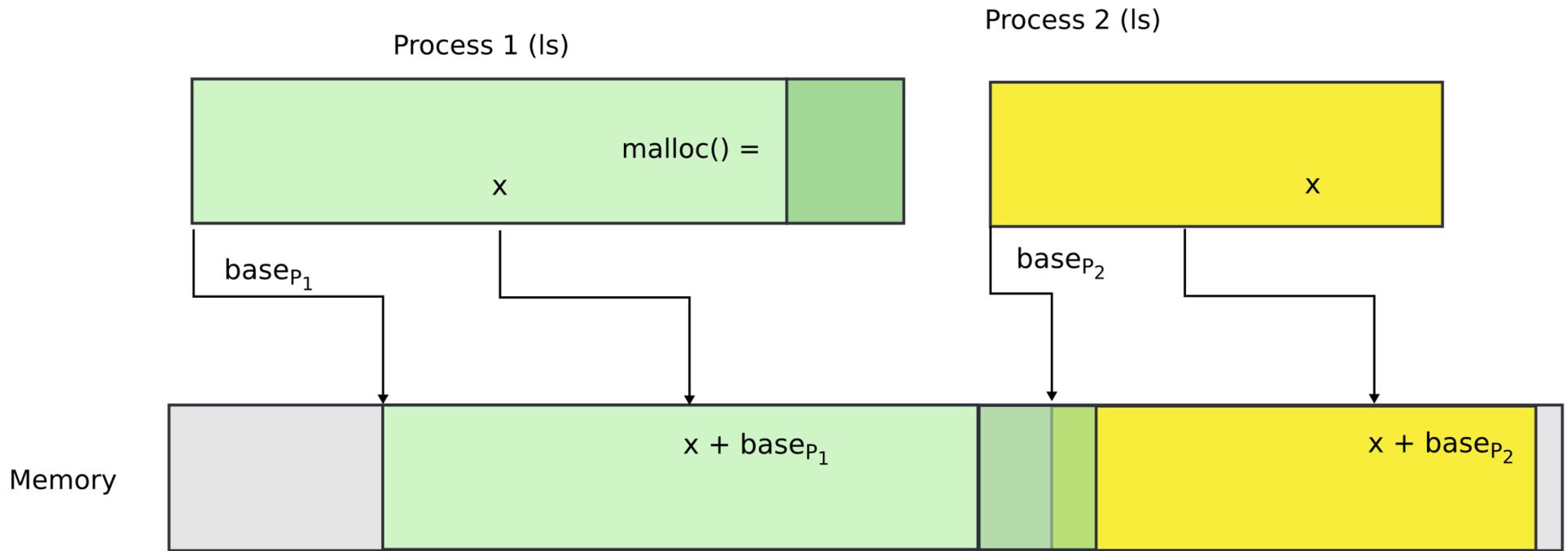
# Segmentation is ok... but



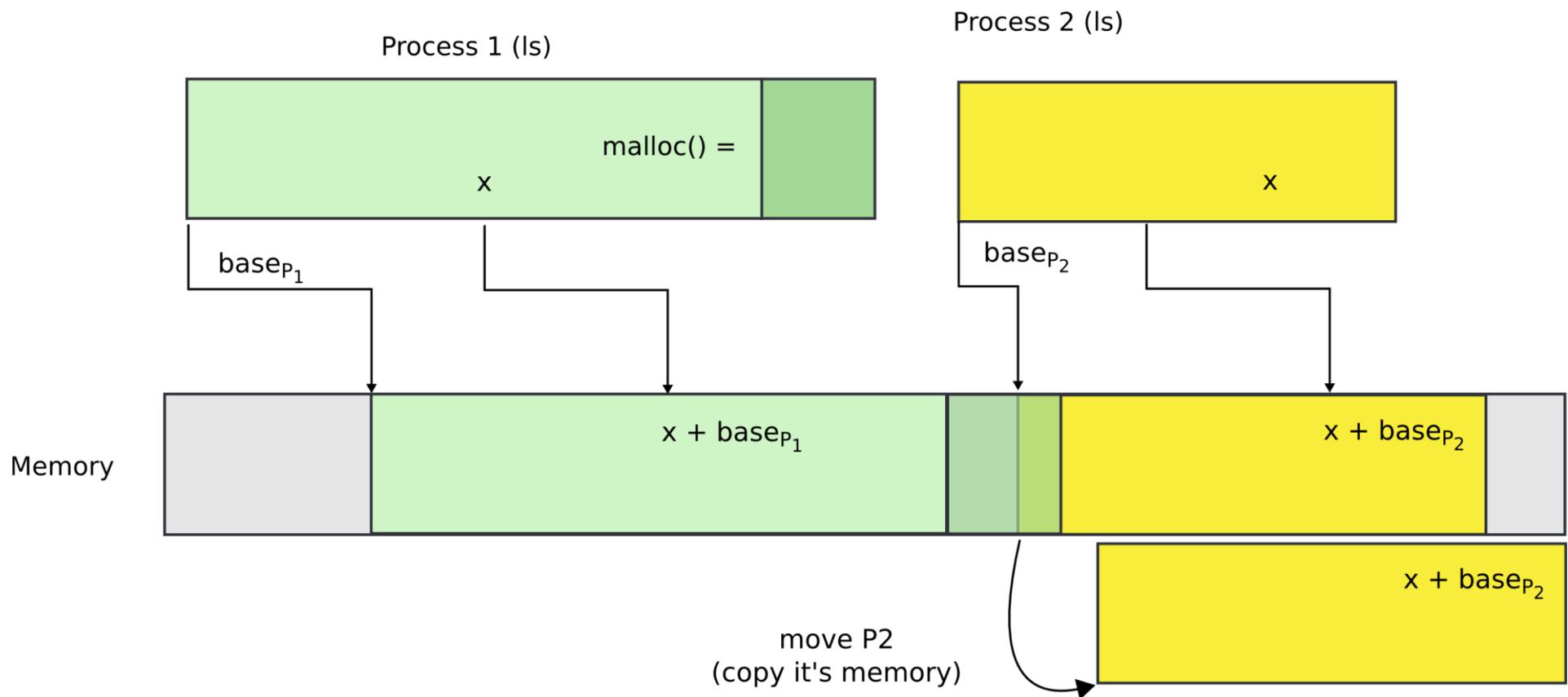
# What if process needs more memory?



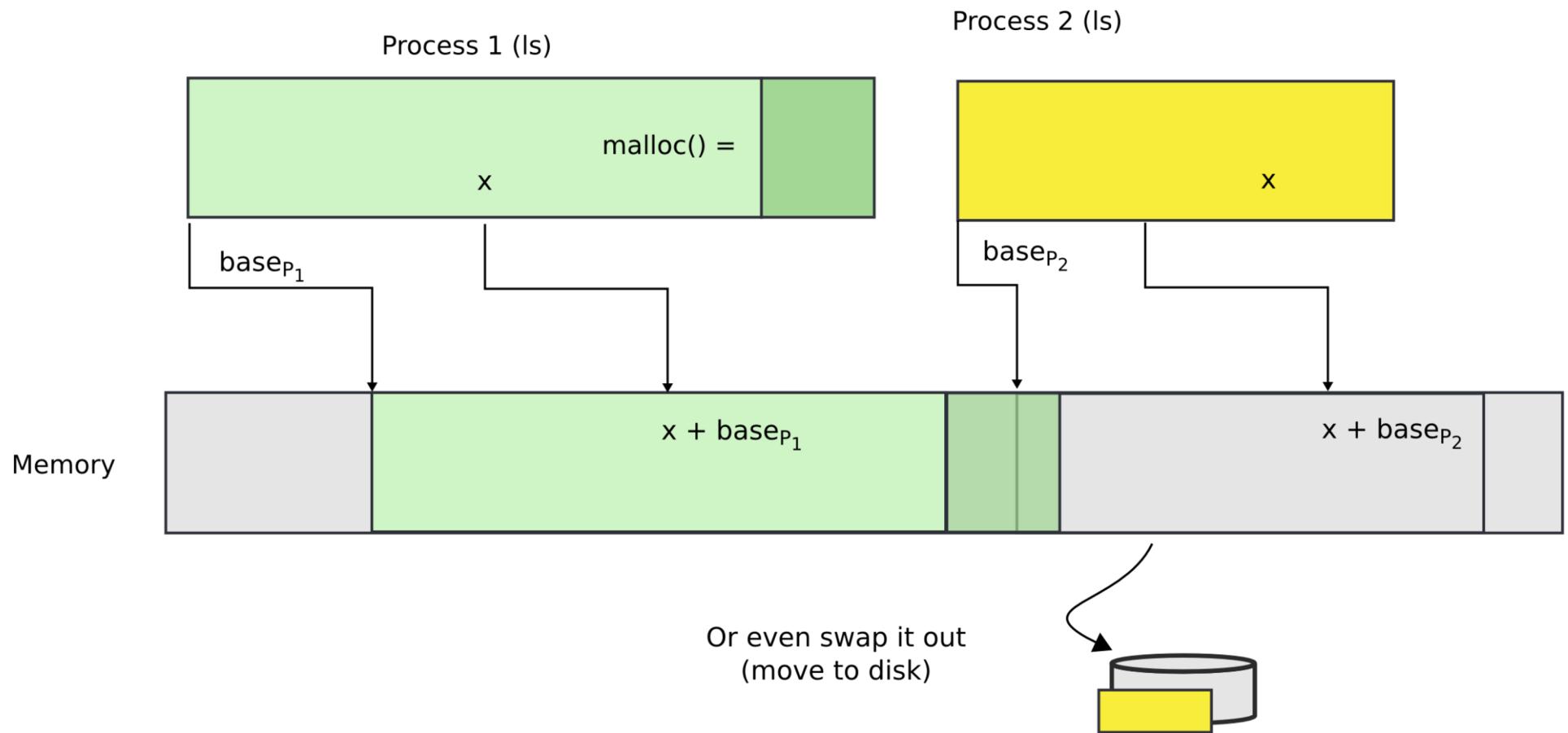
# What if process needs more memory?



# You can move P2 in memory



# Or even swap it out to disk

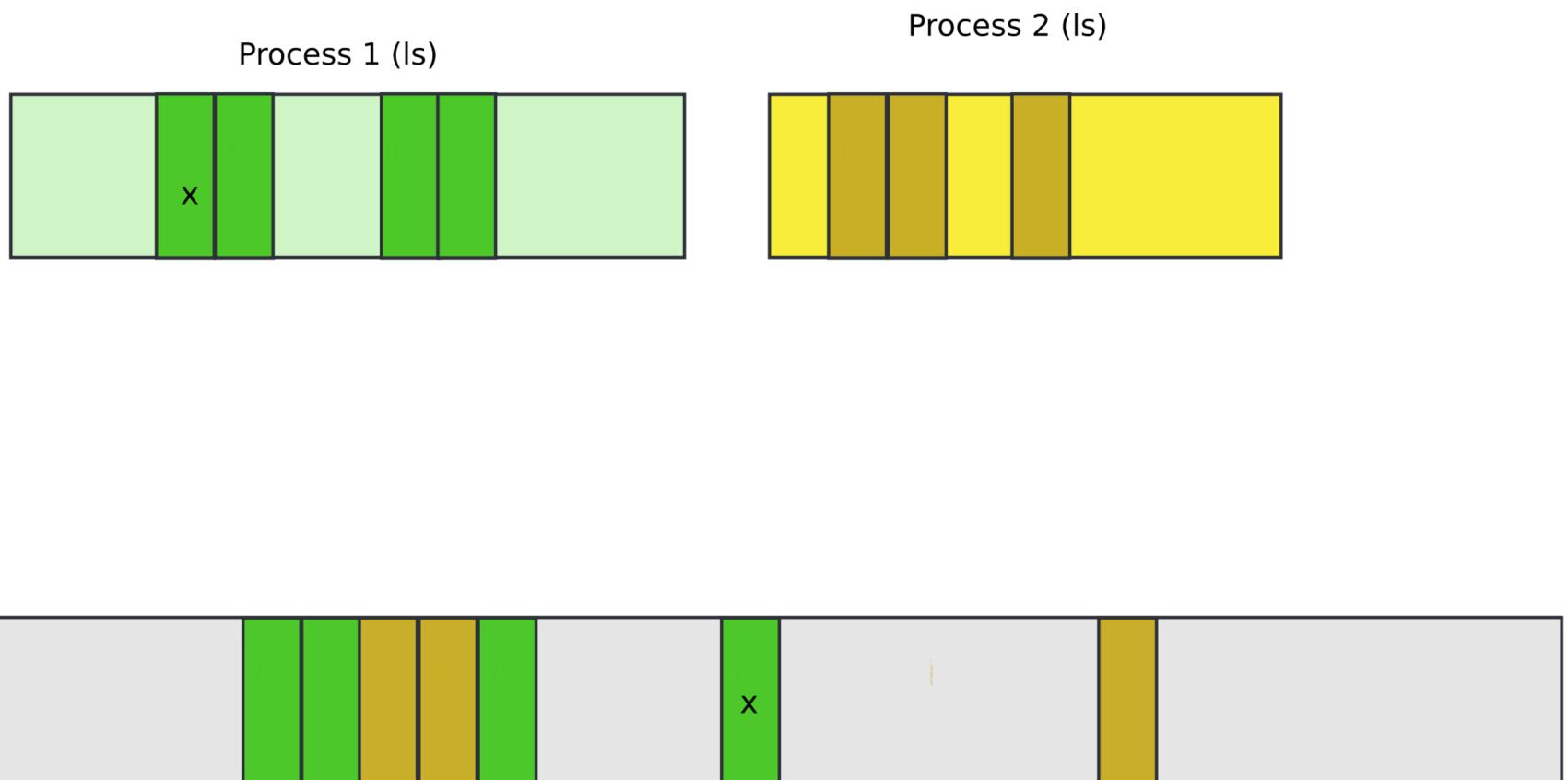


# Problems with segments

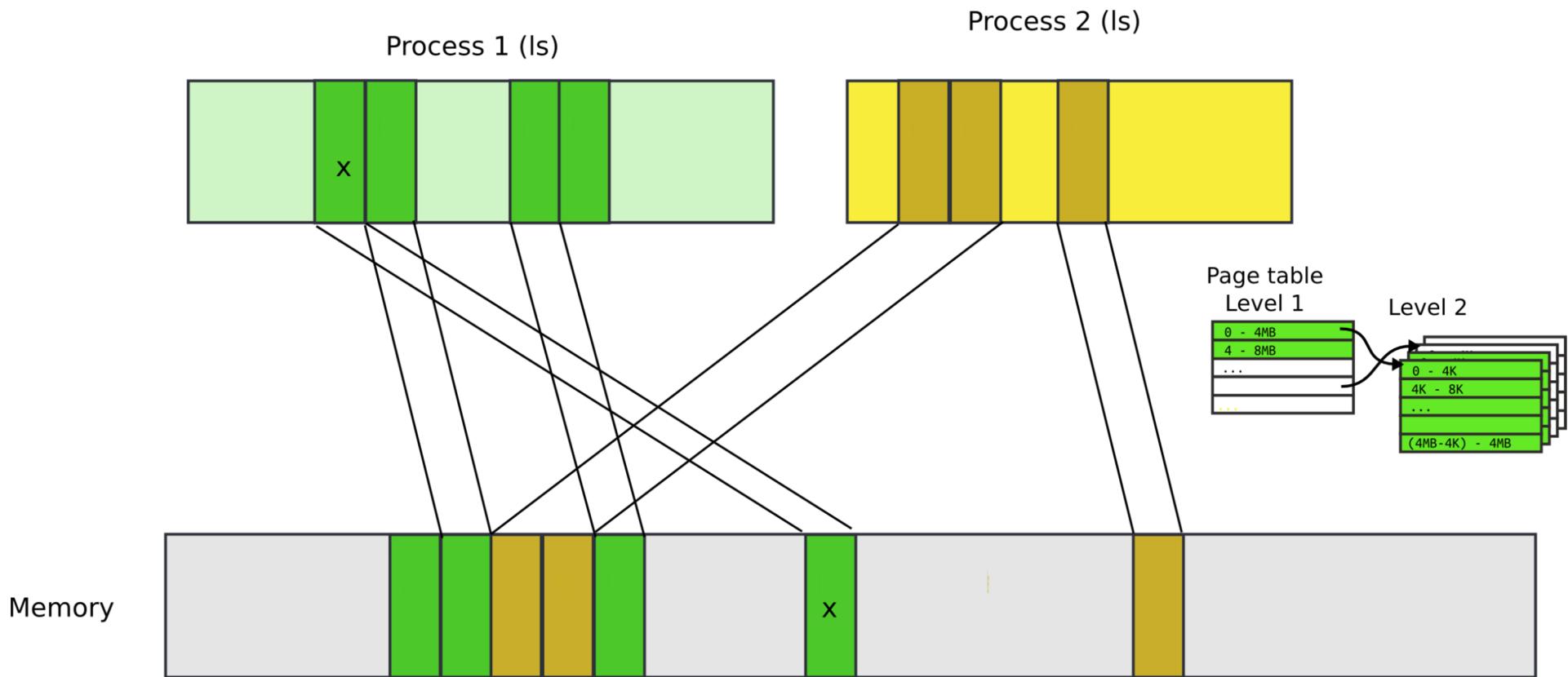
- Segments are somewhat inconvenient
  - Relocating or swapping the entire process takes time
- Memory gets fragmented
  - There might be no space (gap) for the swapped out process to come in
  - Will have to swap out other processes

# Paging

# Pages



# Pages

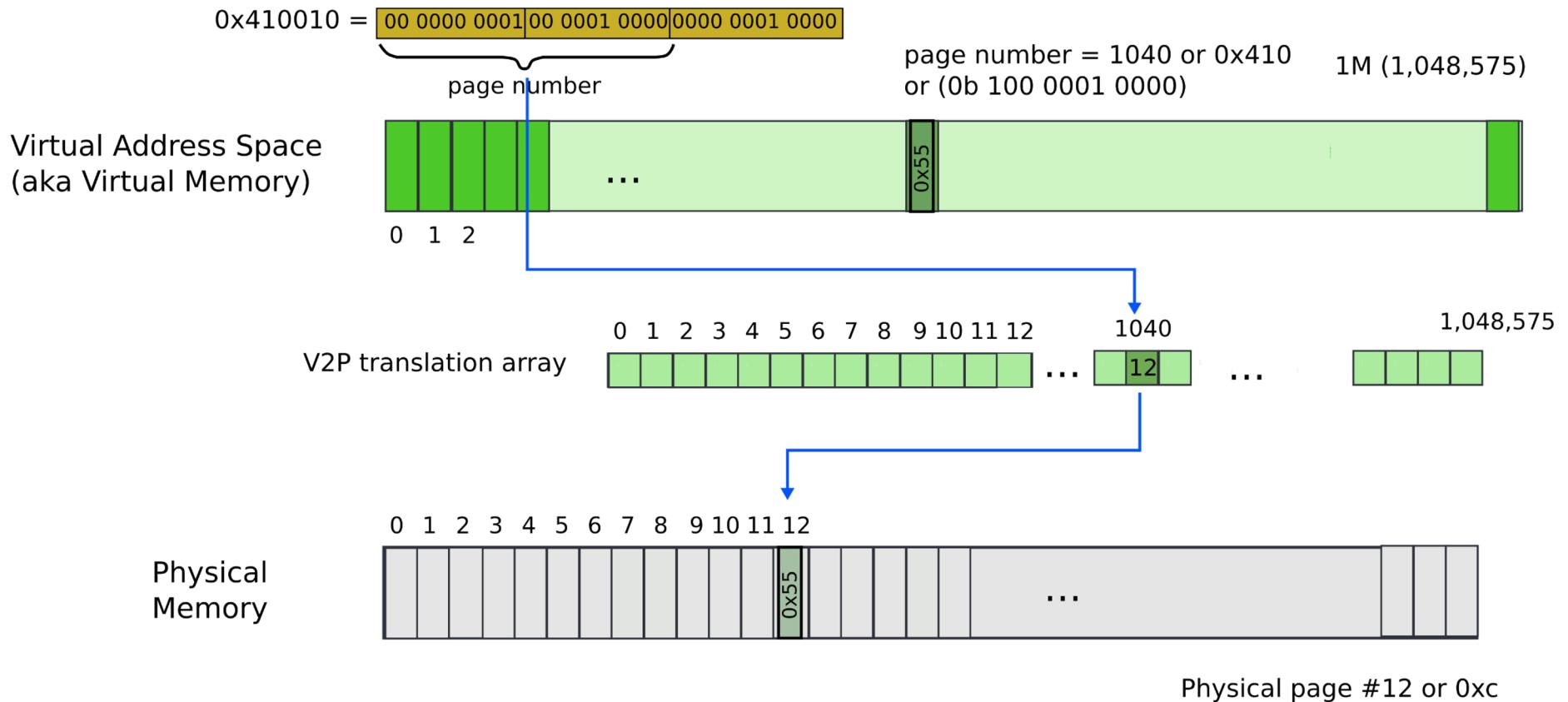


# Paging idea

- 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
- Compared with segmentation (single base + limit)
  - Much more flexibility

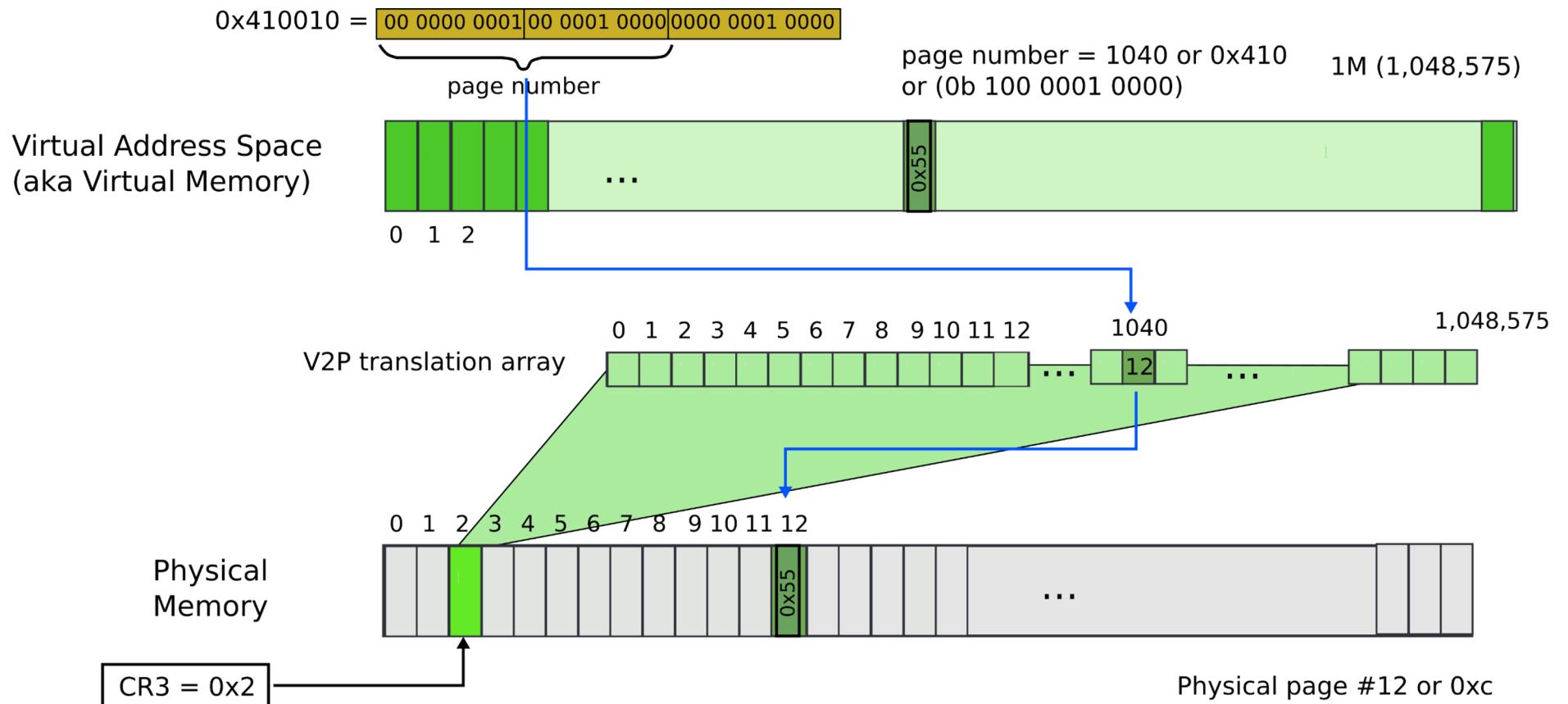
# How can we build this translation mechanism?

# Paging: naive approach: translation array



- Linear address 0x410010
- Remember it's result of logical to linear translation (aka segmentation)
  - $0x410010 = 0x300010$  (offset) +  $0x110000$  (base)

# Paging: naive approach: translation array



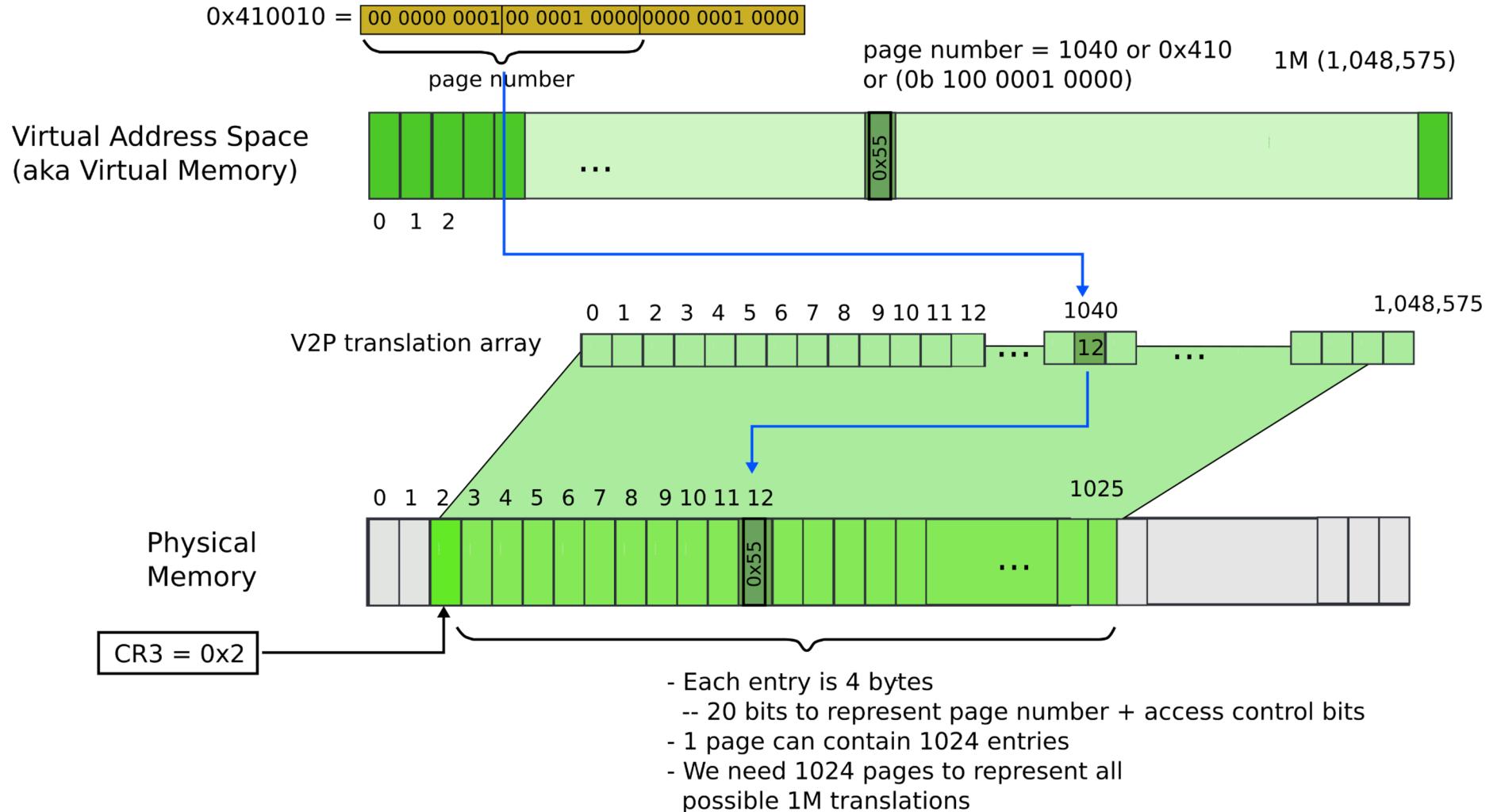
- Linear address 0x410010
- Remember it's result of logical to linear translation (aka segmentation)
  - $0x410010 = 0x300010 \text{ (offset)} + 0x110000 \text{ (base)}$

# What is wrong?

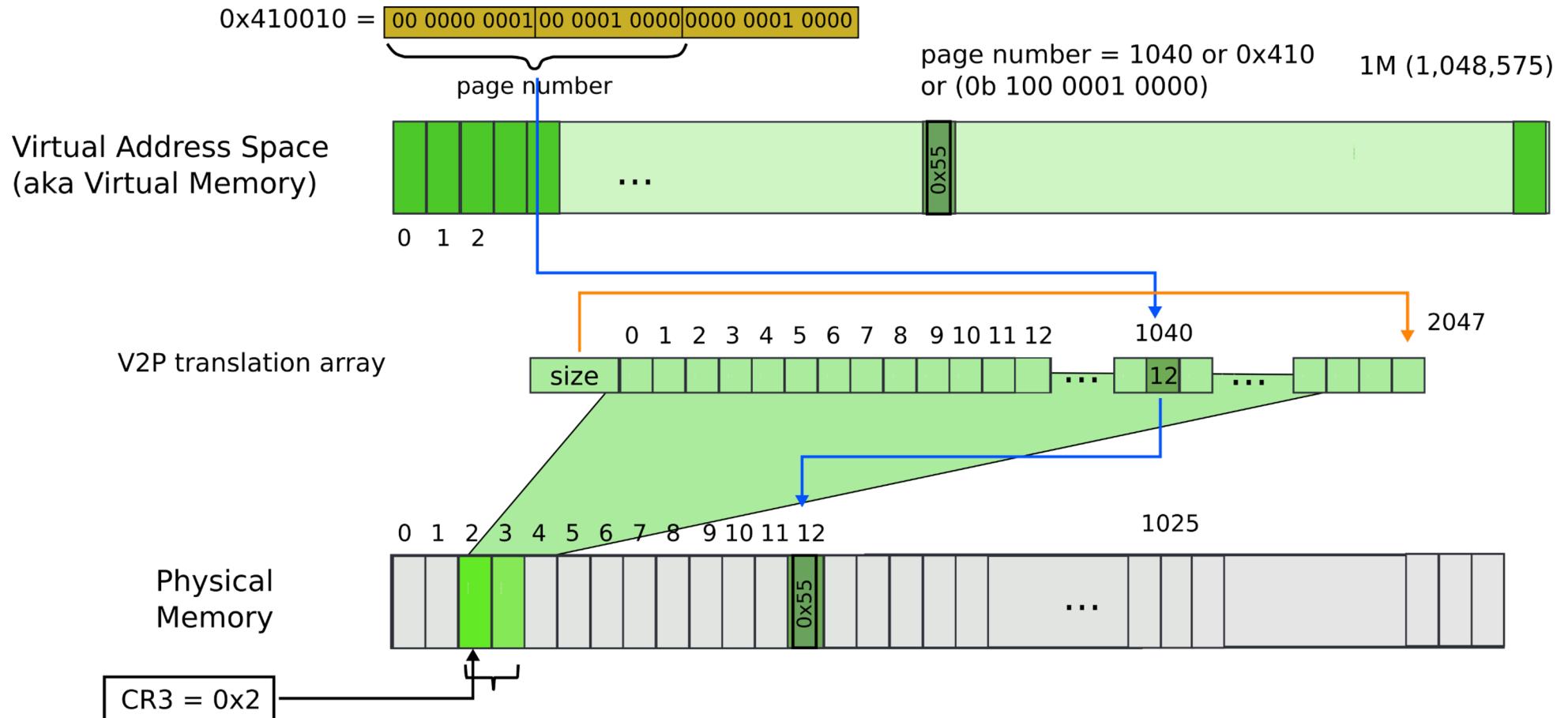
# What is wrong?

- We need 4 bytes to relocate each page
- 20 bits for physical page number
- 12 bits of access flags
- Therefore, we need array of 4 bytes x 1M entries
  - 4MBs

# Paging: naive approach: translation array



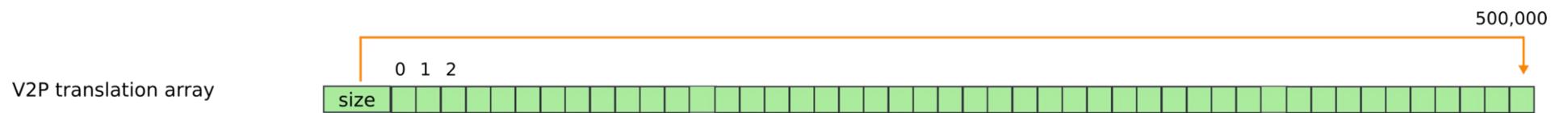
# Paging: array with size



- The size controls how many entries are required

# But still what may go wrong?

# Paging: array with size

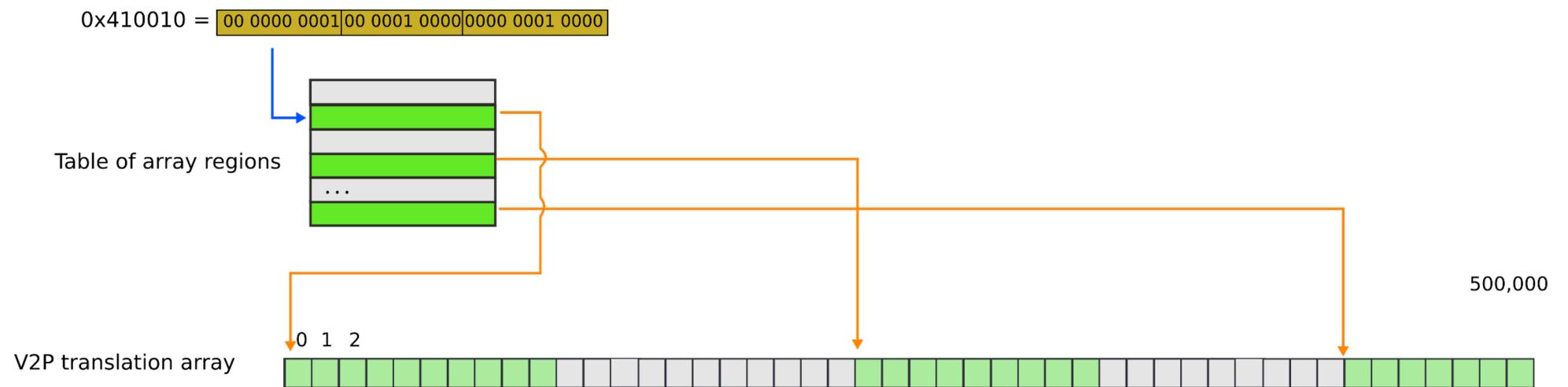


# Paging: array with size

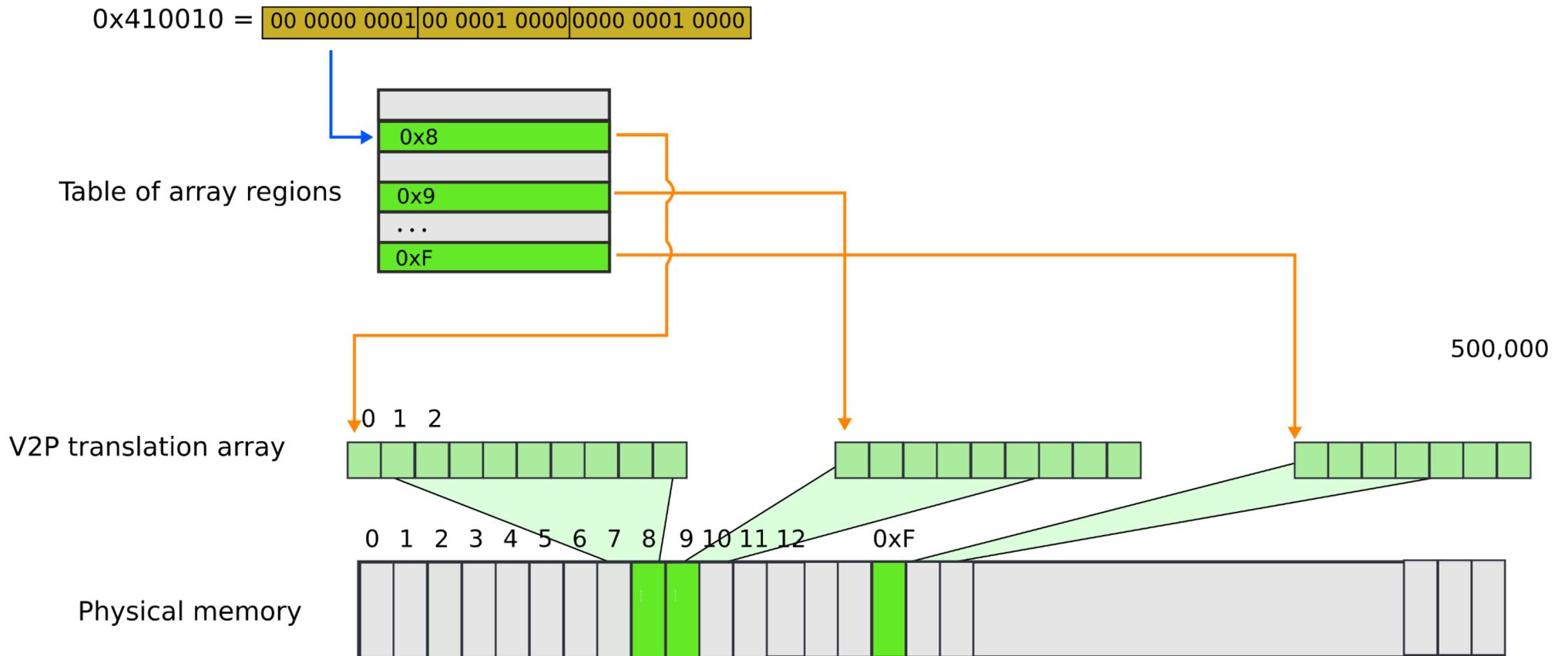


# Can we improve?

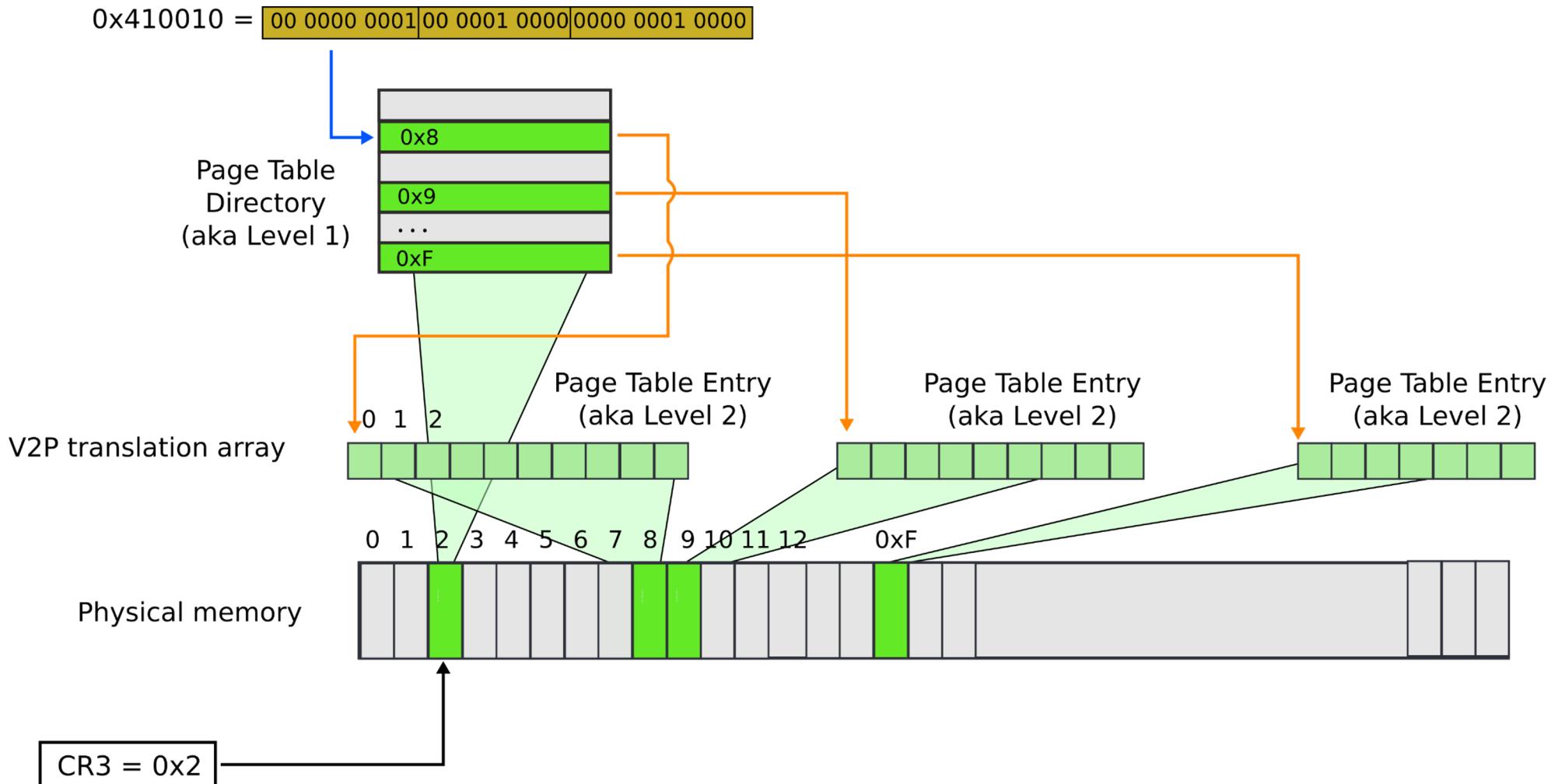
# Paging: array with chunks



# Paging: array with chunks



# Paging: page table



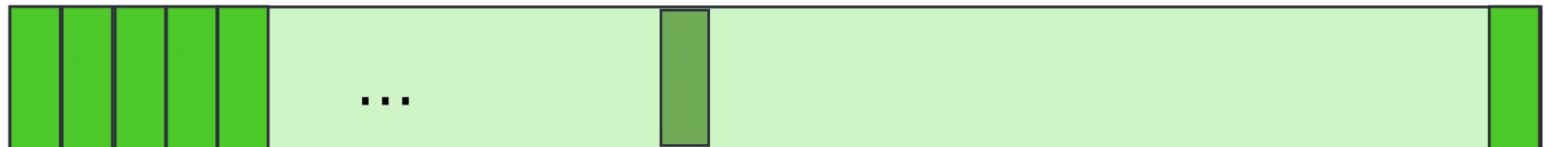
```
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX  
EAX = 0  
EBX = 20 983 809
```

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

CR3 = 0

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

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

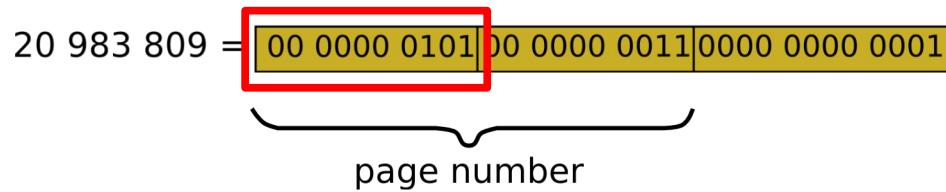
Physical  
Memory



```

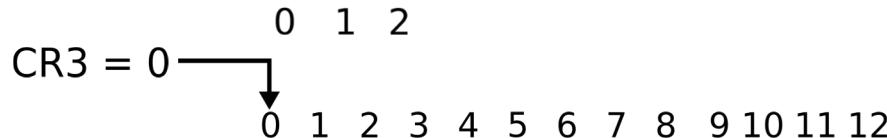
mov (%EBX), EAX # mov value from the location pointed by EBX into EAX
EAX = 0
EBX = 20 983 809

```



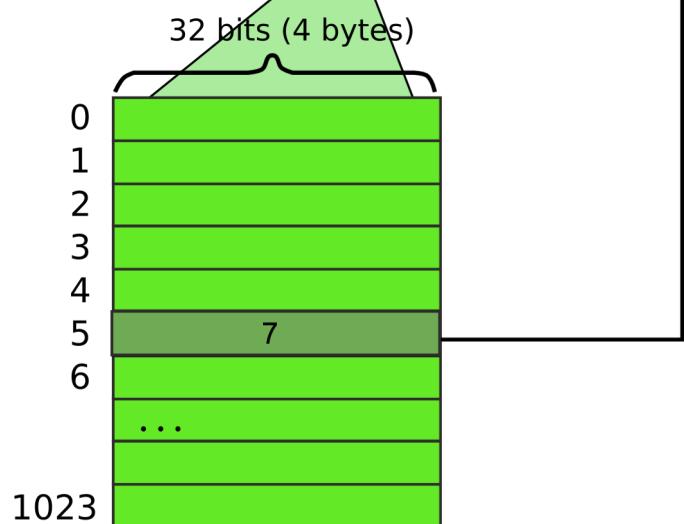
1M (1,048,575)

Virtual Address  
Space (or Memory)  
of the Process



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

Physical  
Memory

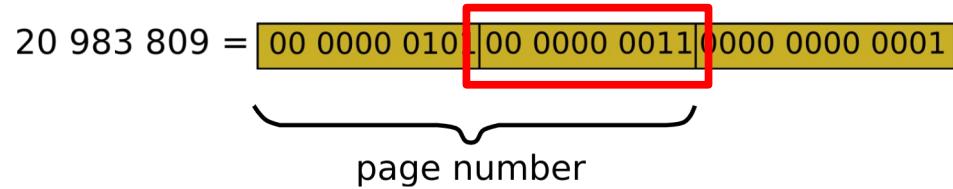


Level 1  
(Page Table  
Directory)

```

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

```



1M (1,048,575)

Virtual Address  
Space (or Memory)  
of the Process

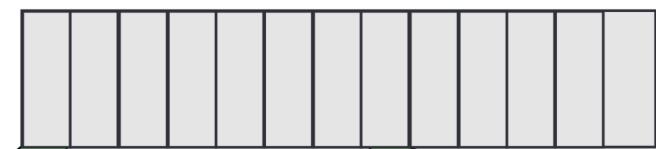


CR3 = 0

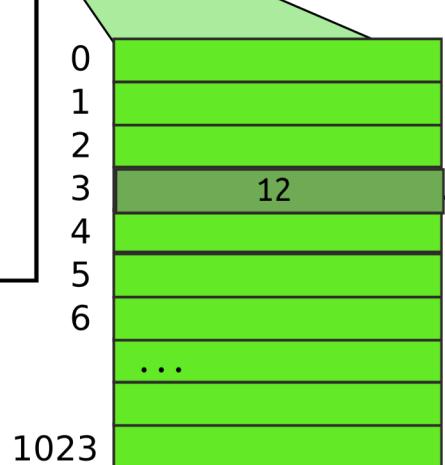
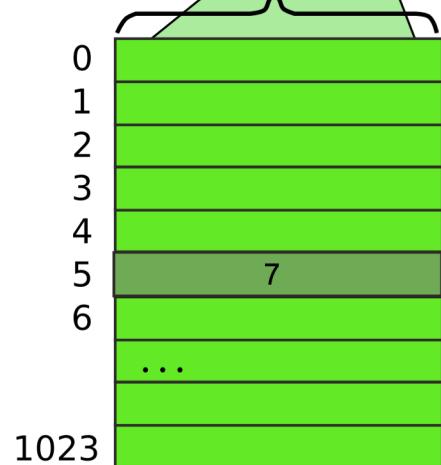
0 1 2 3 4 5 6 7 8 9 10 11 12

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

Physical  
Memory



32 bits (4 bytes)



Level 1  
(Page Table  
Directory)

Level 2  
(Page Table)

```

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)

Virtual Address  
Space (or Memory)  
of the Process



CR3 = 0 →  
0 1 2 3 4 5 6 7 8 9 10 11 12

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

Physical  
Memory



Level 1  
(Page Table  
Directory)

Level 2  
(Page Table)

Page

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

• **Result:**

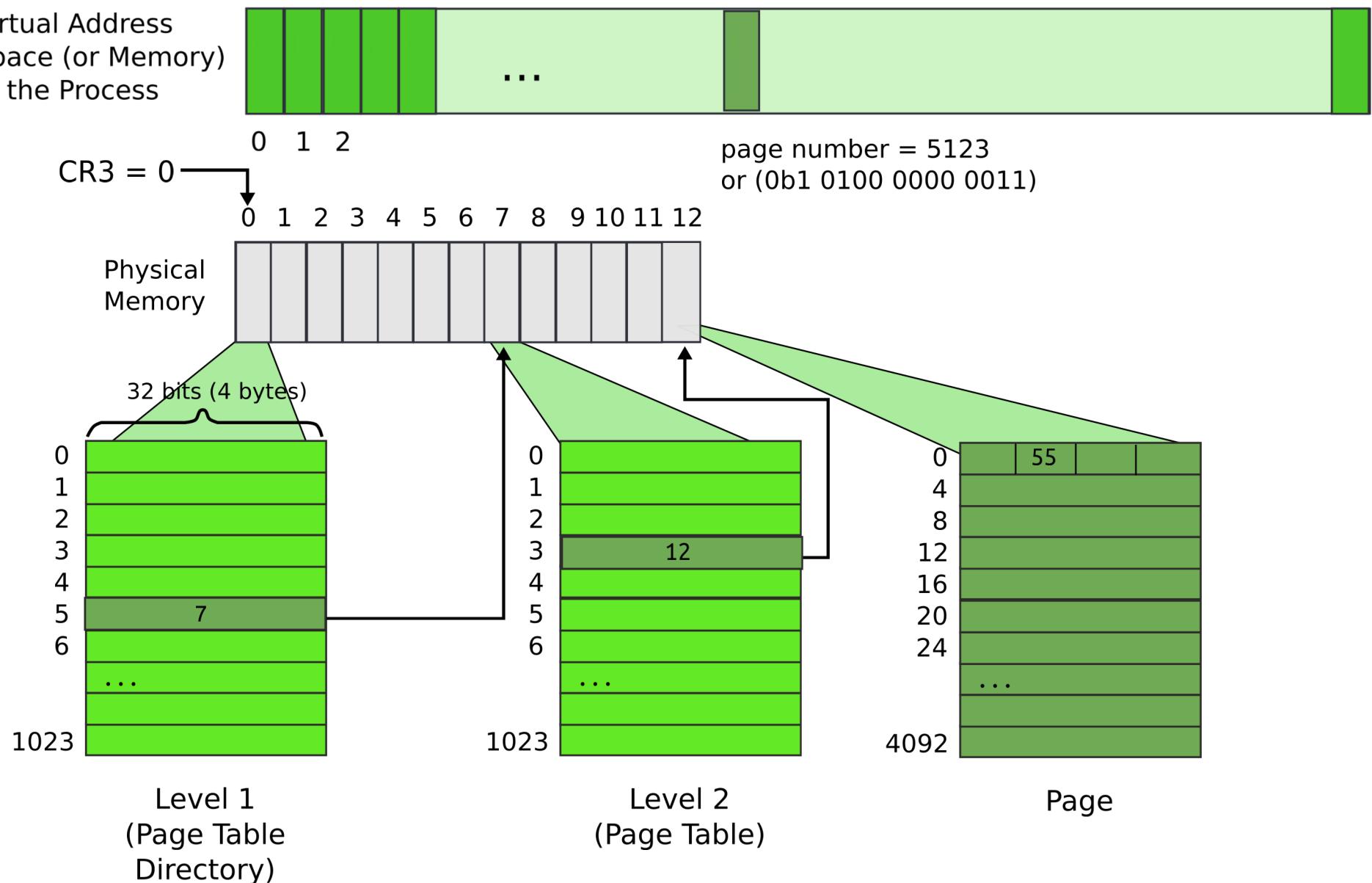
20 983 809 = **00 0000 0101|00 0000 0011|0000 0000 0001**

page number

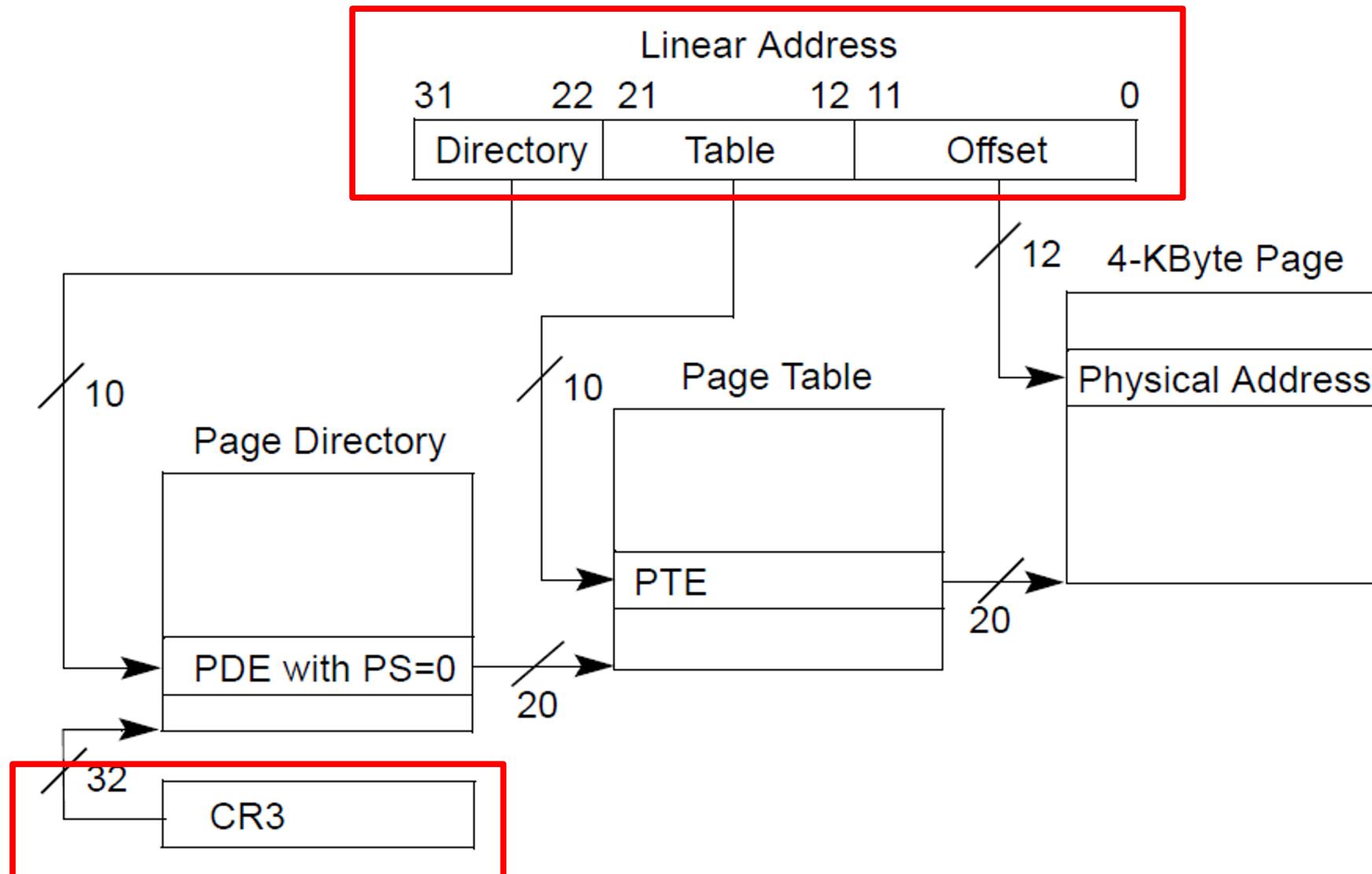
- Result:
  - EAX = 55

1M (1.048.575)

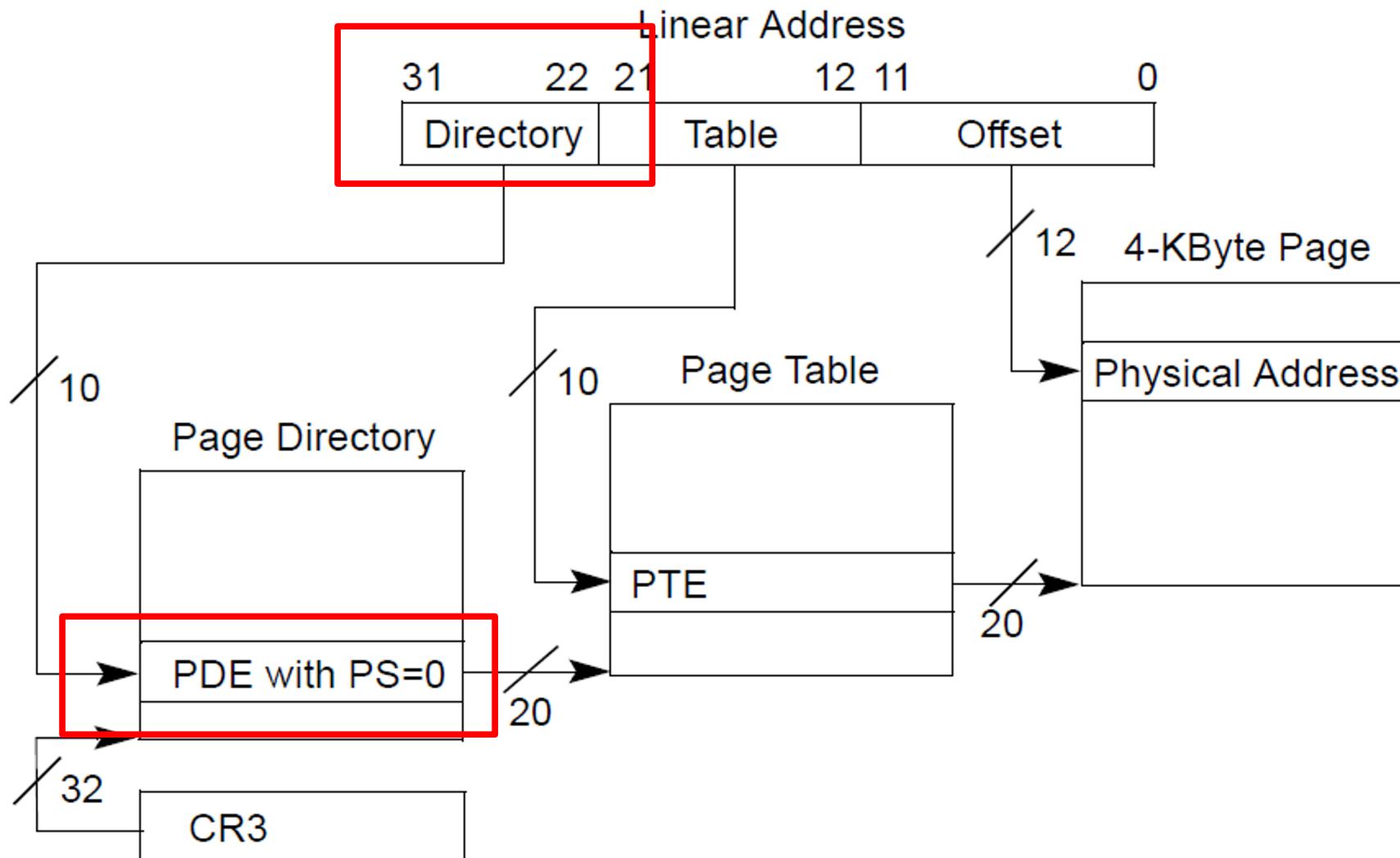
## Virtual Address Space (or Memory) of the Process



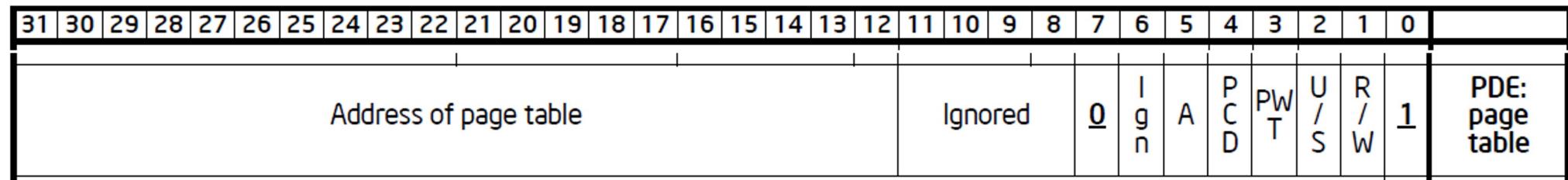
# Page translation



# Page translation

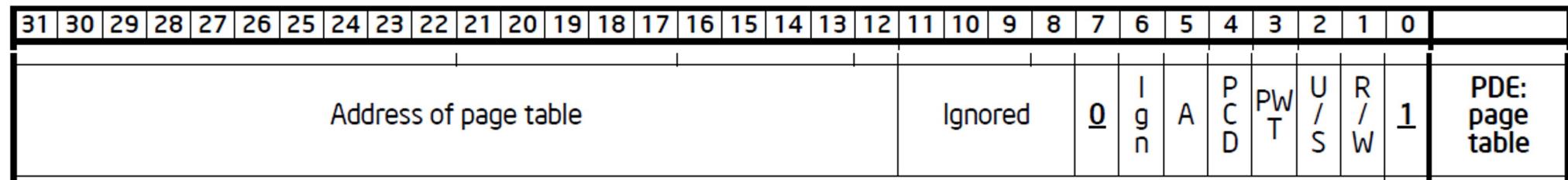


# Page directory entry (PDE)



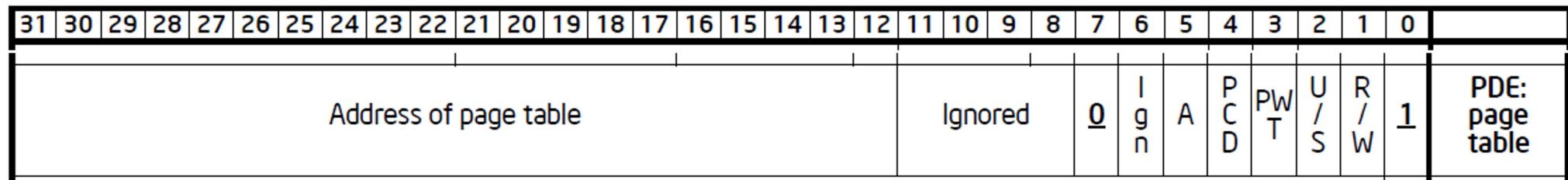
- 20 bit address of the page table

# Page directory entry (PDE)



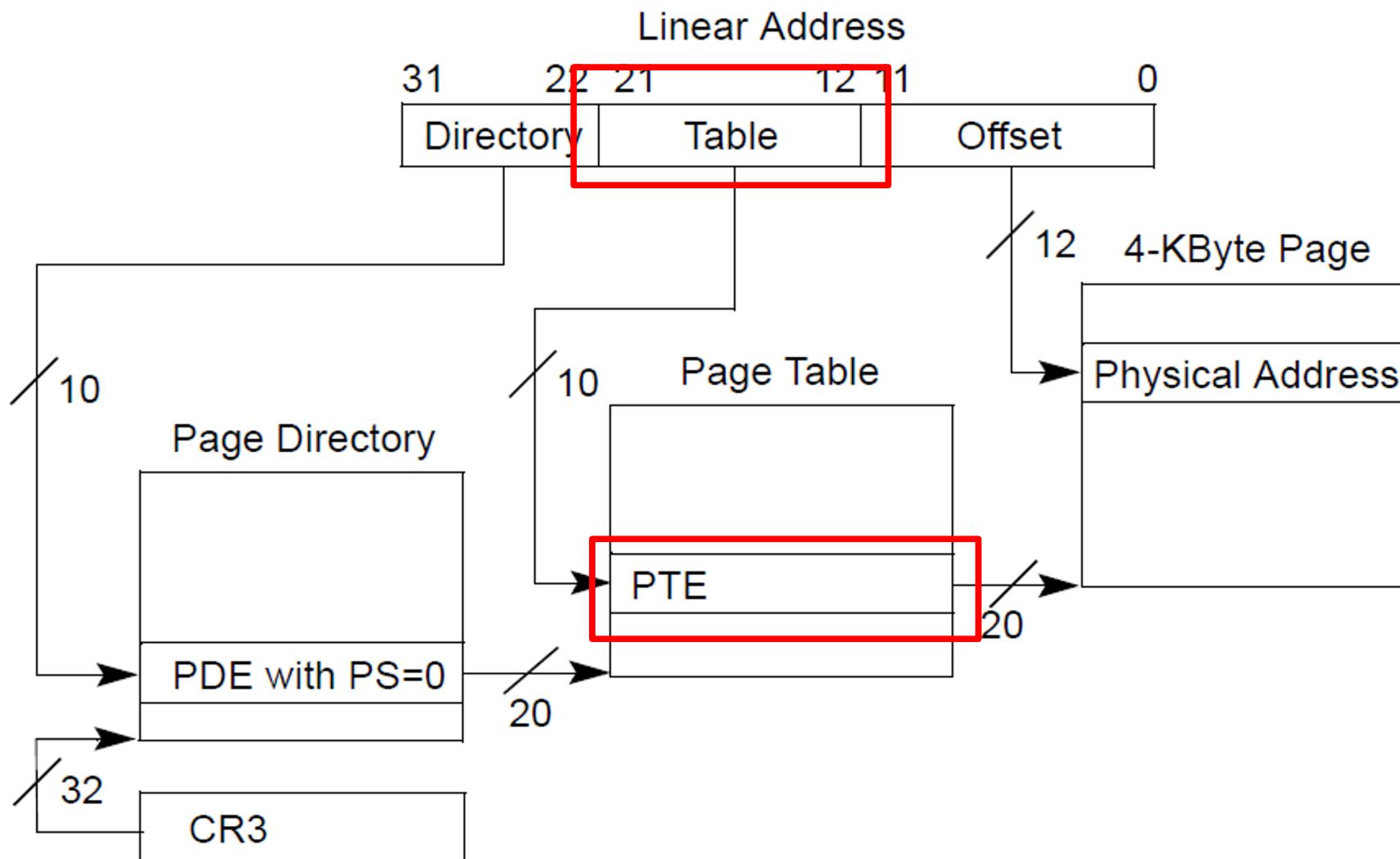
- 20 bit address of the page table
- Wait... 20 bit address, but we need 32 bits

# Page directory entry (PDE)

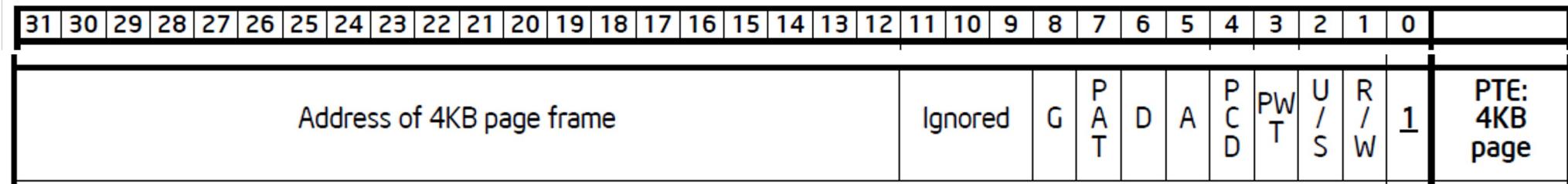


- 20 bit address of the page table
- Wait... 20 bit address, but we need 32 bits
- Pages 4KB each, we need 1M to cover 4GB
- Pages start at 4KB (page aligned boundary)

# Page translation

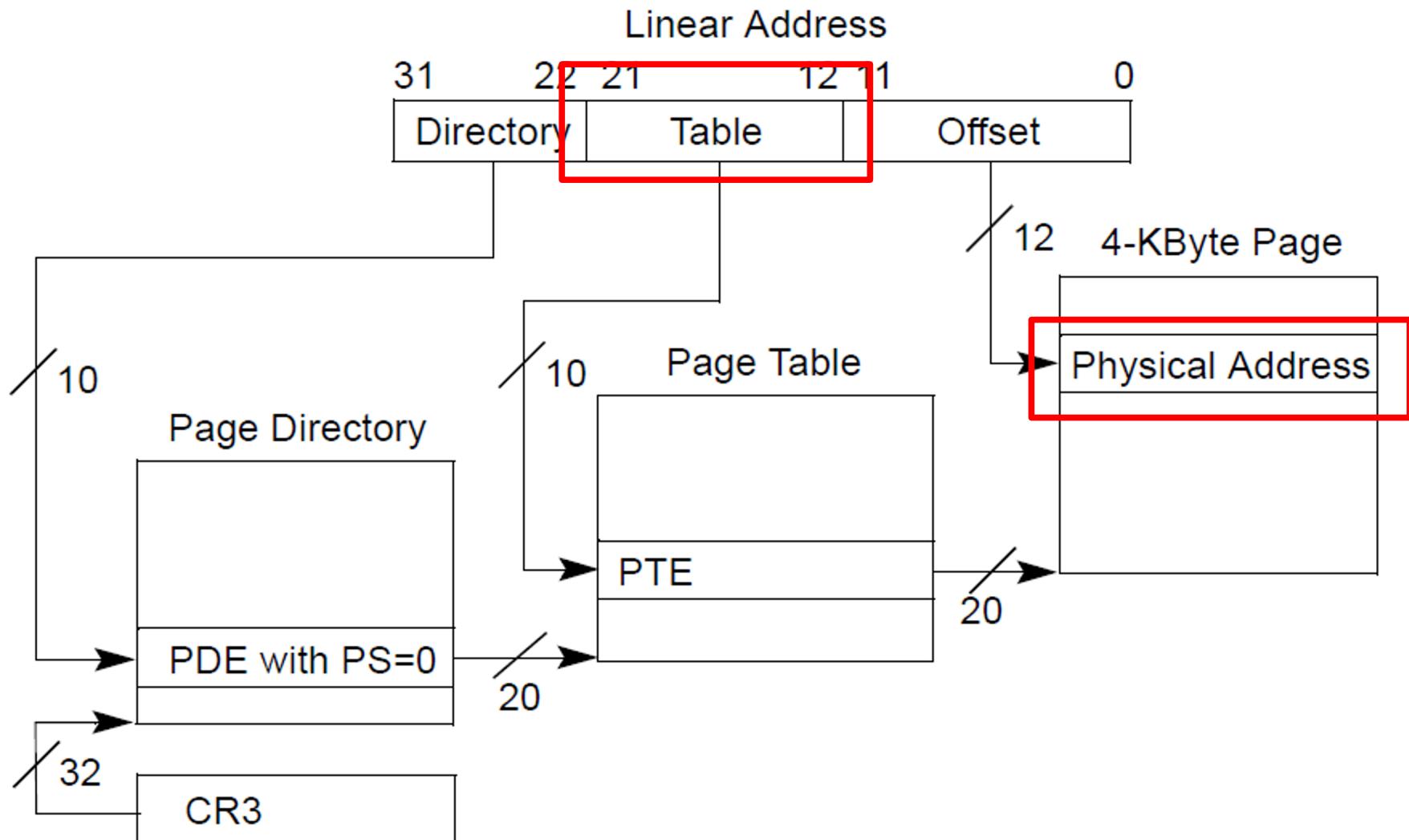


# Page table entry (PTE)



- 20 bit address of the 4KB page
- Pages 4KB each, we need 1M to cover 4GB

# Page translation

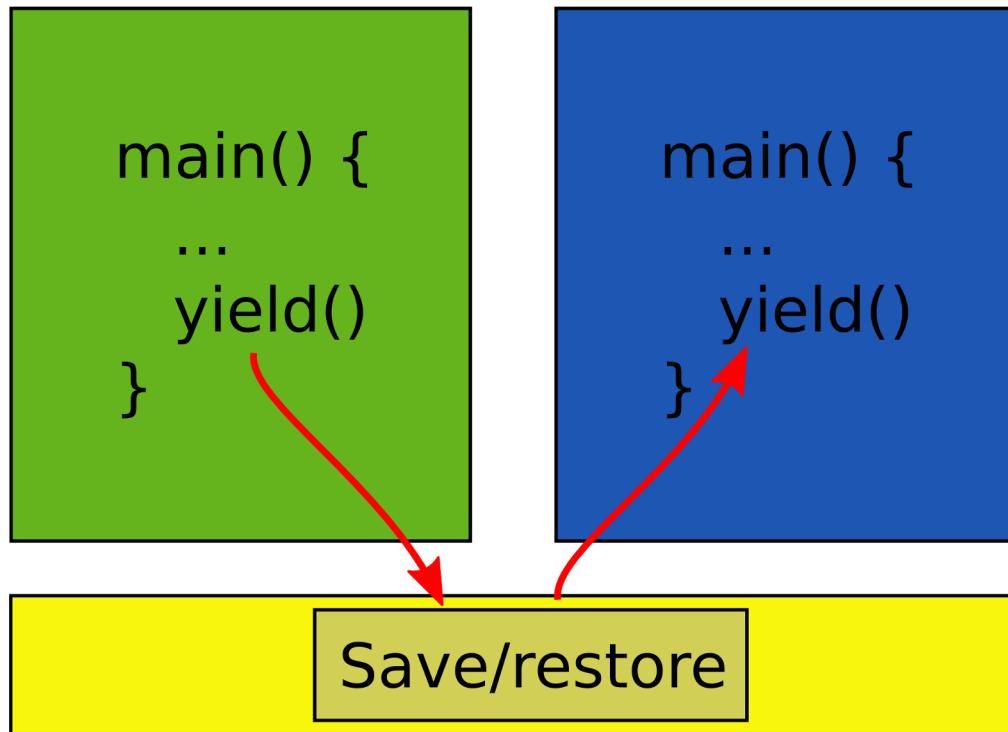


# Benefit of page tables

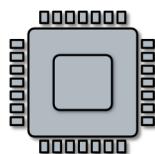
... Compared to arrays?

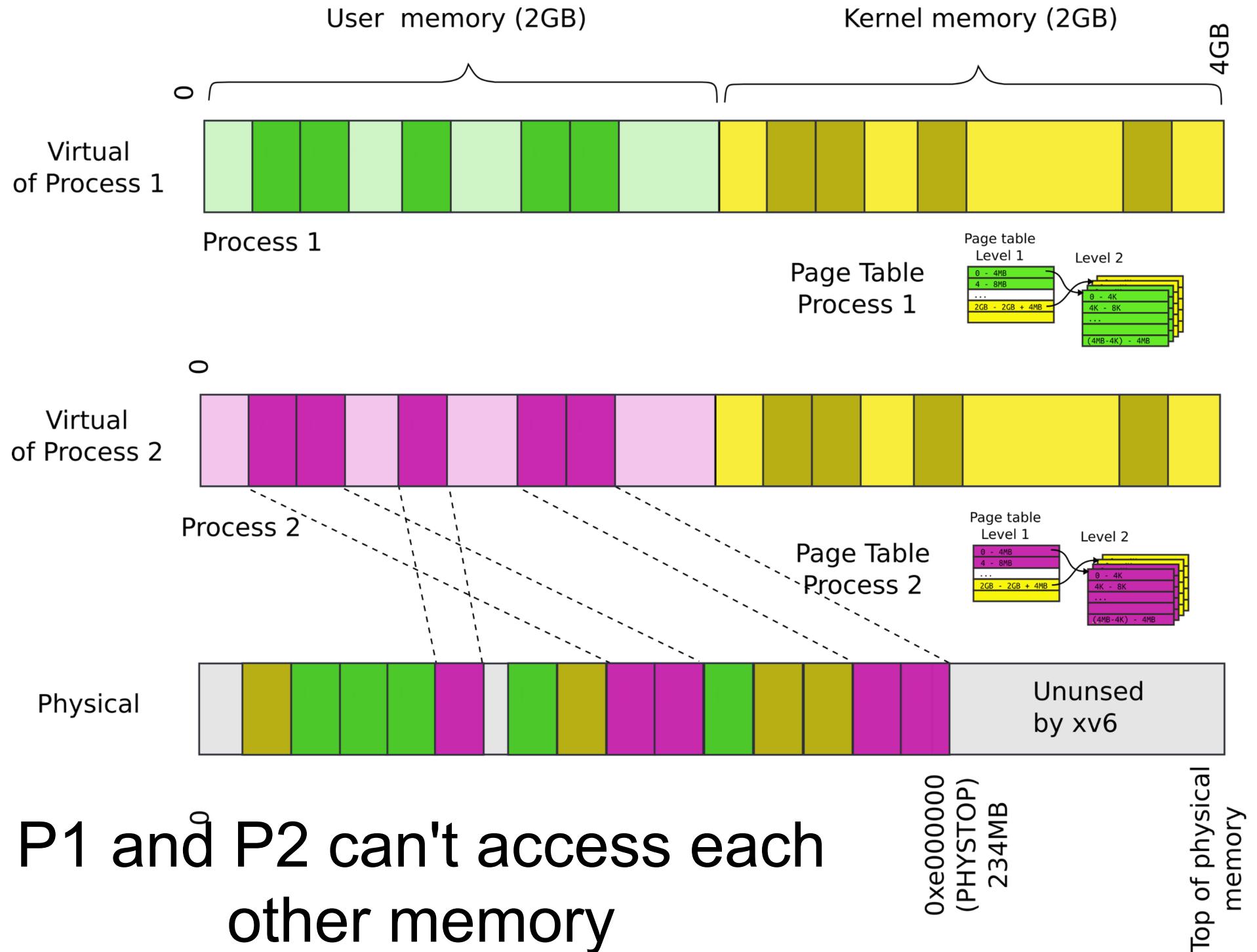
- Page tables represent sparse address space more efficiently
- An entire array has to be allocated upfront
- But if the address space uses a handful of pages
- Only page tables (Level 1 and 2 need to be allocated to describe translation)
- On a dense address space this benefit goes away
- I'll assign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

# What about isolation?



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





# Compared to segments pages allow

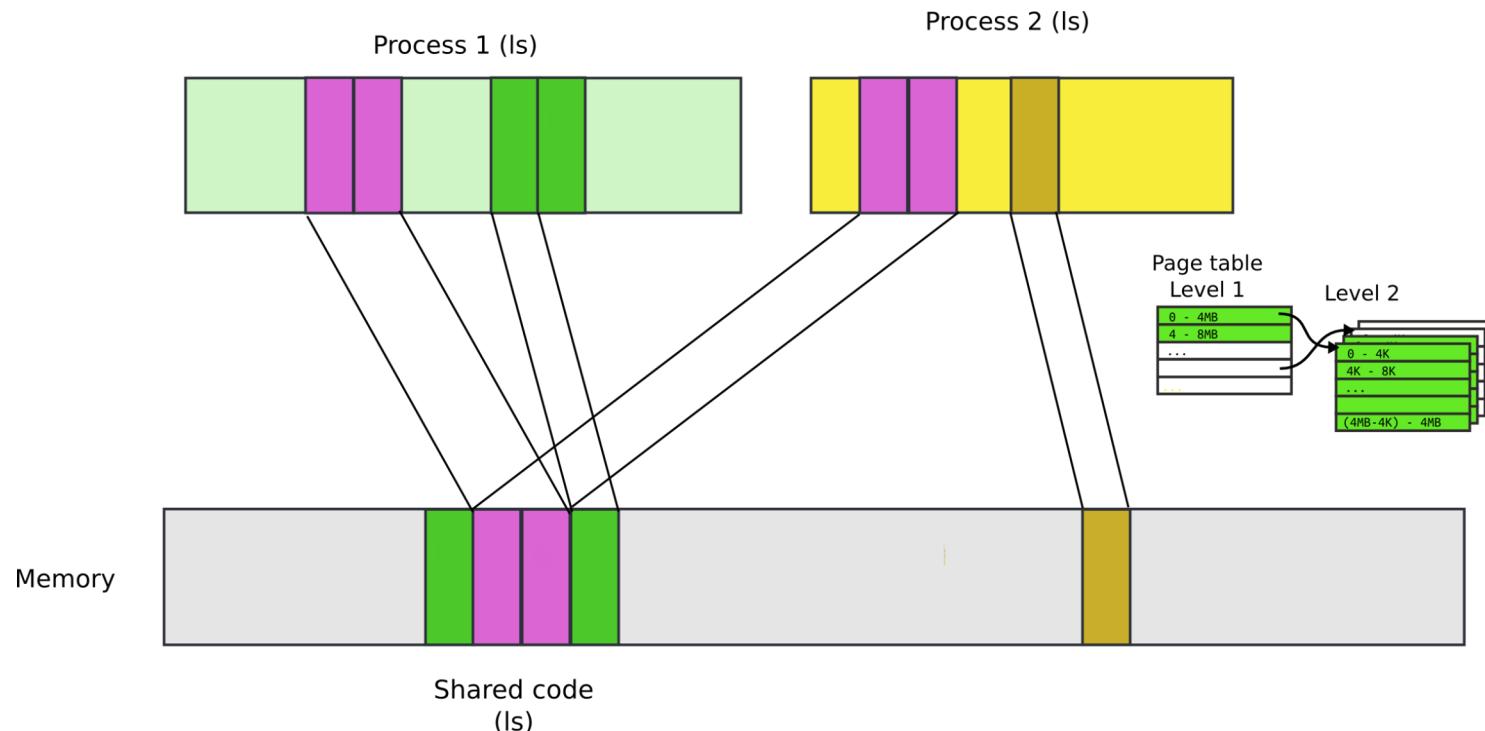
...

- 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

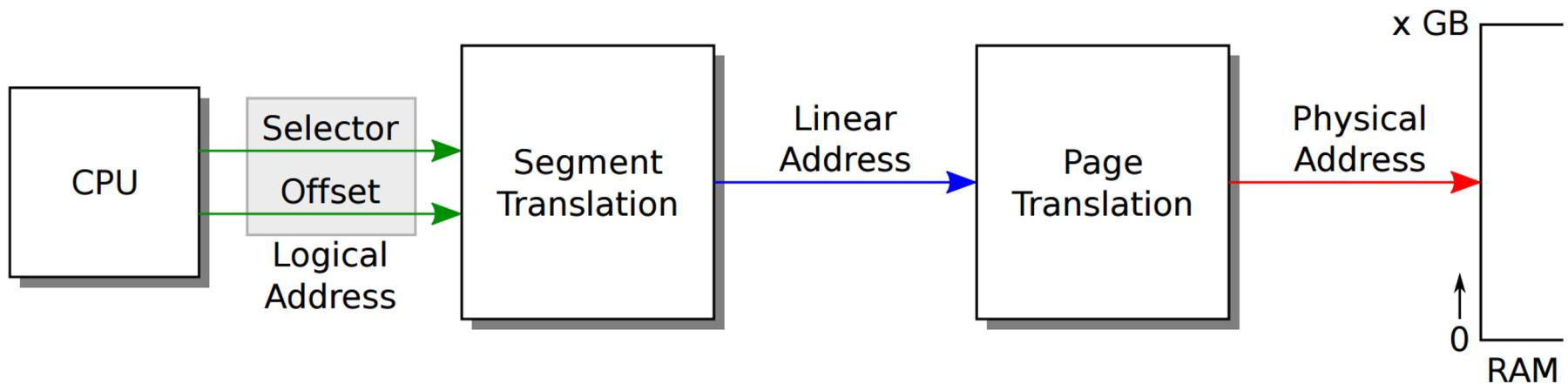
# Compared to segments pages allow

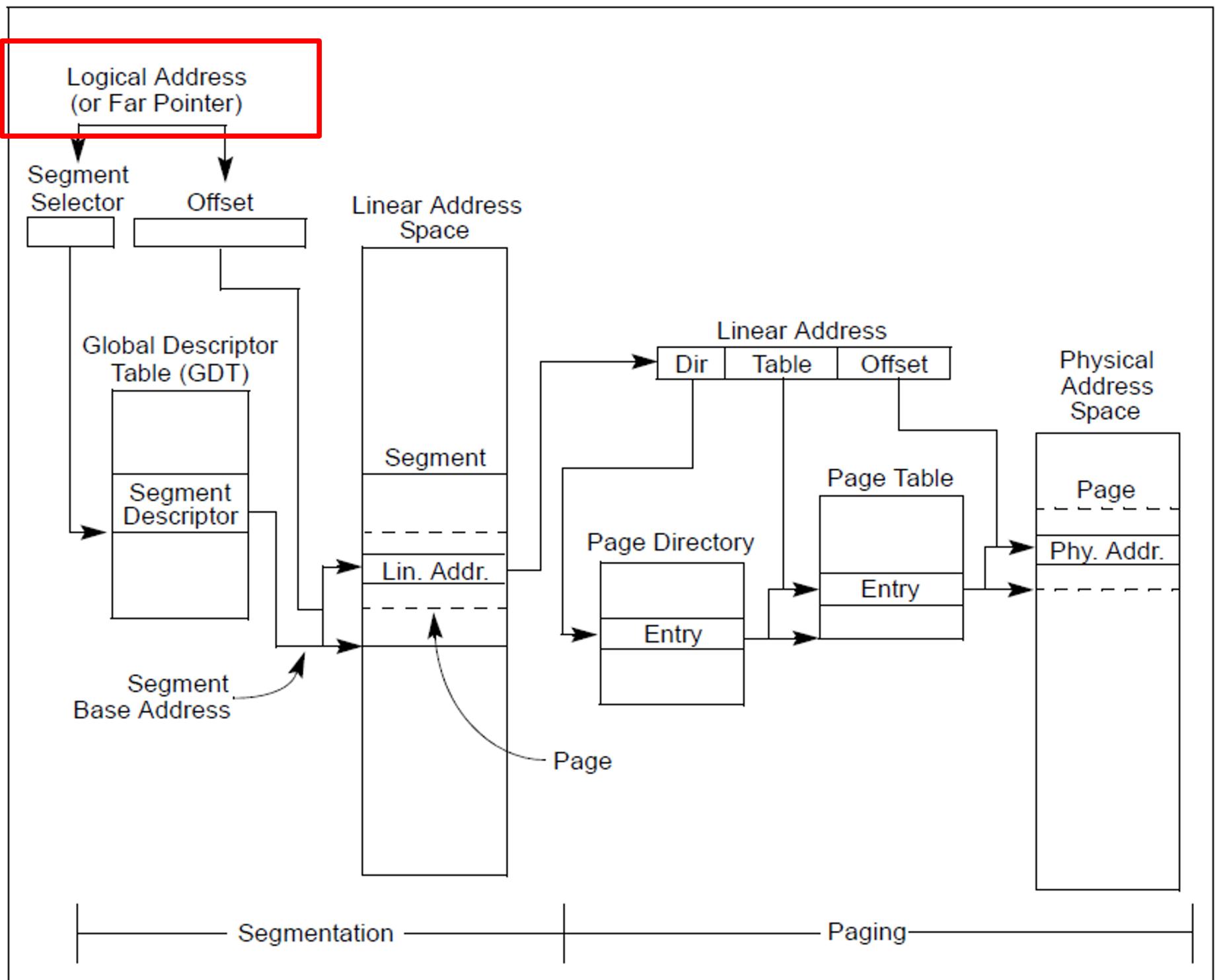
...

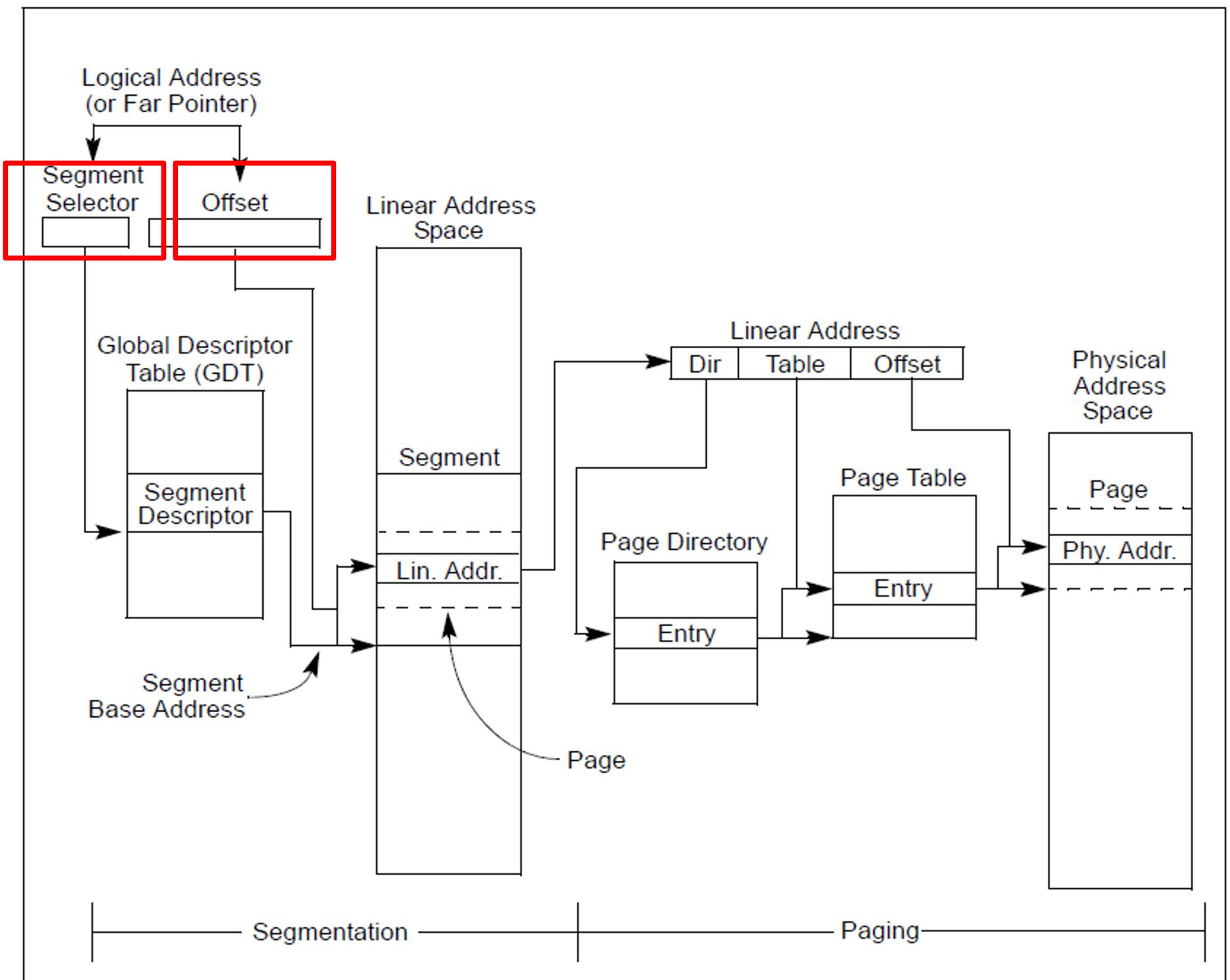
- Share a region of memory across multiple programs
- Well... segments allow this too
- Communication (shared buffer of messages)
- Shared libraries

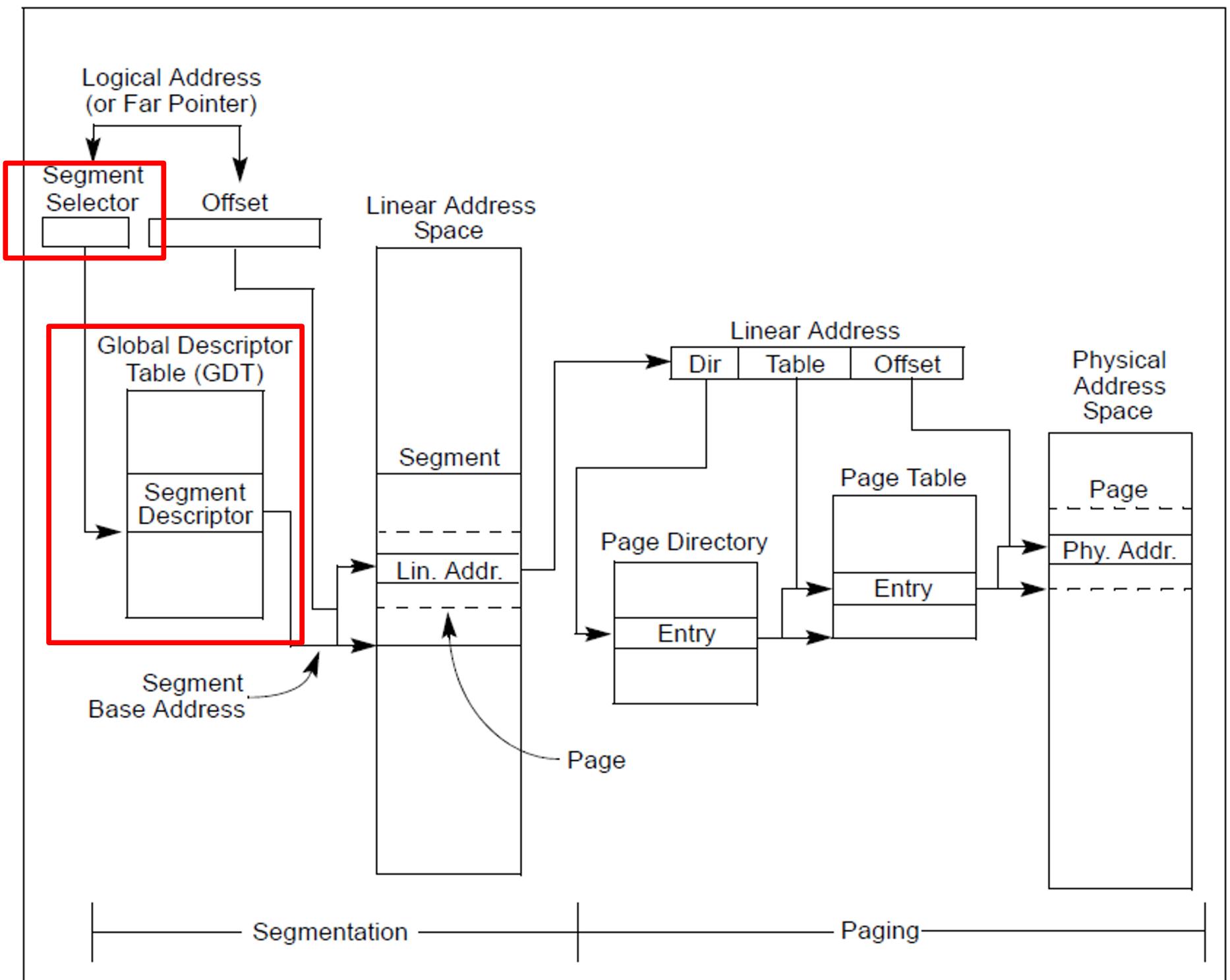


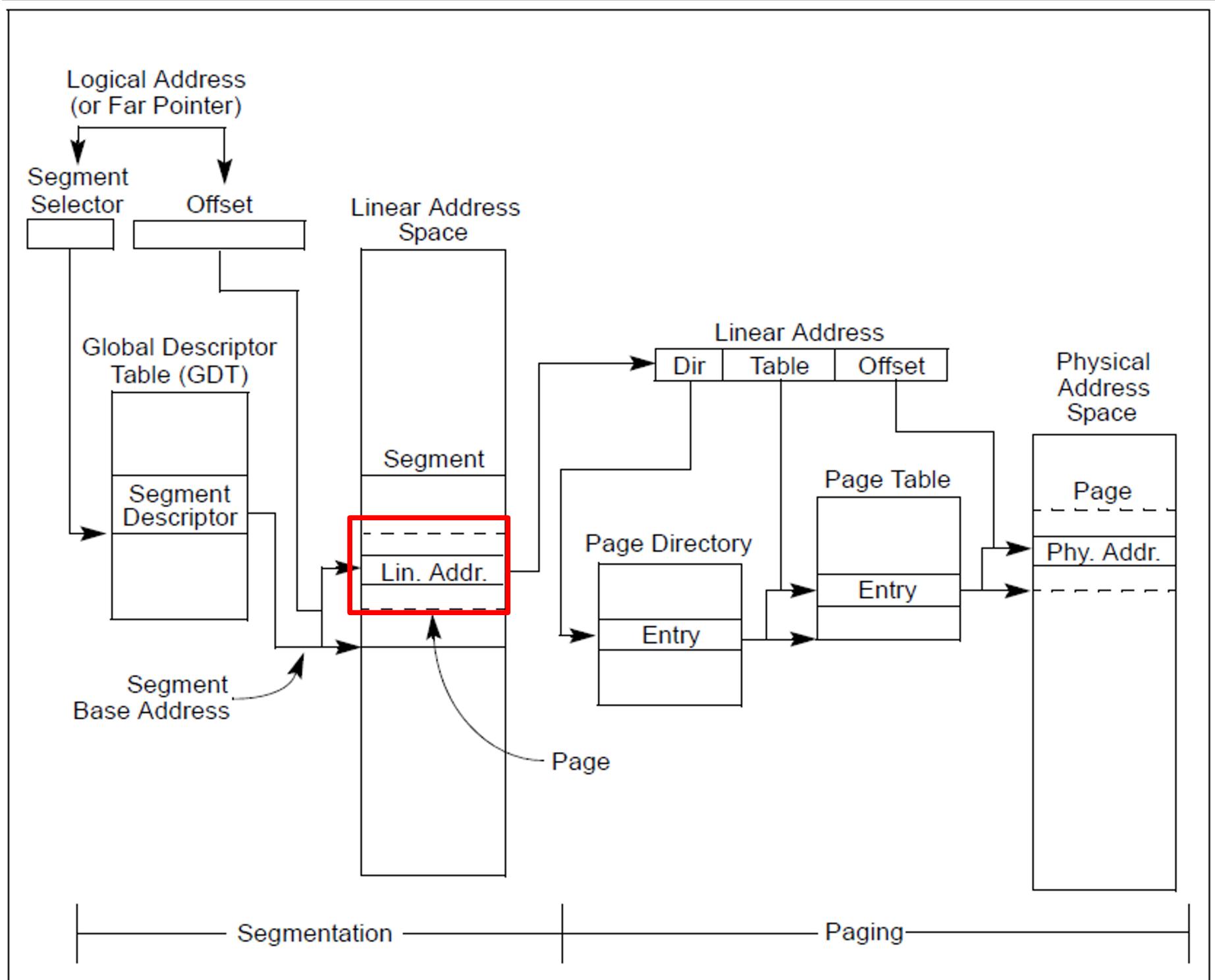
# Recap: complete address translation

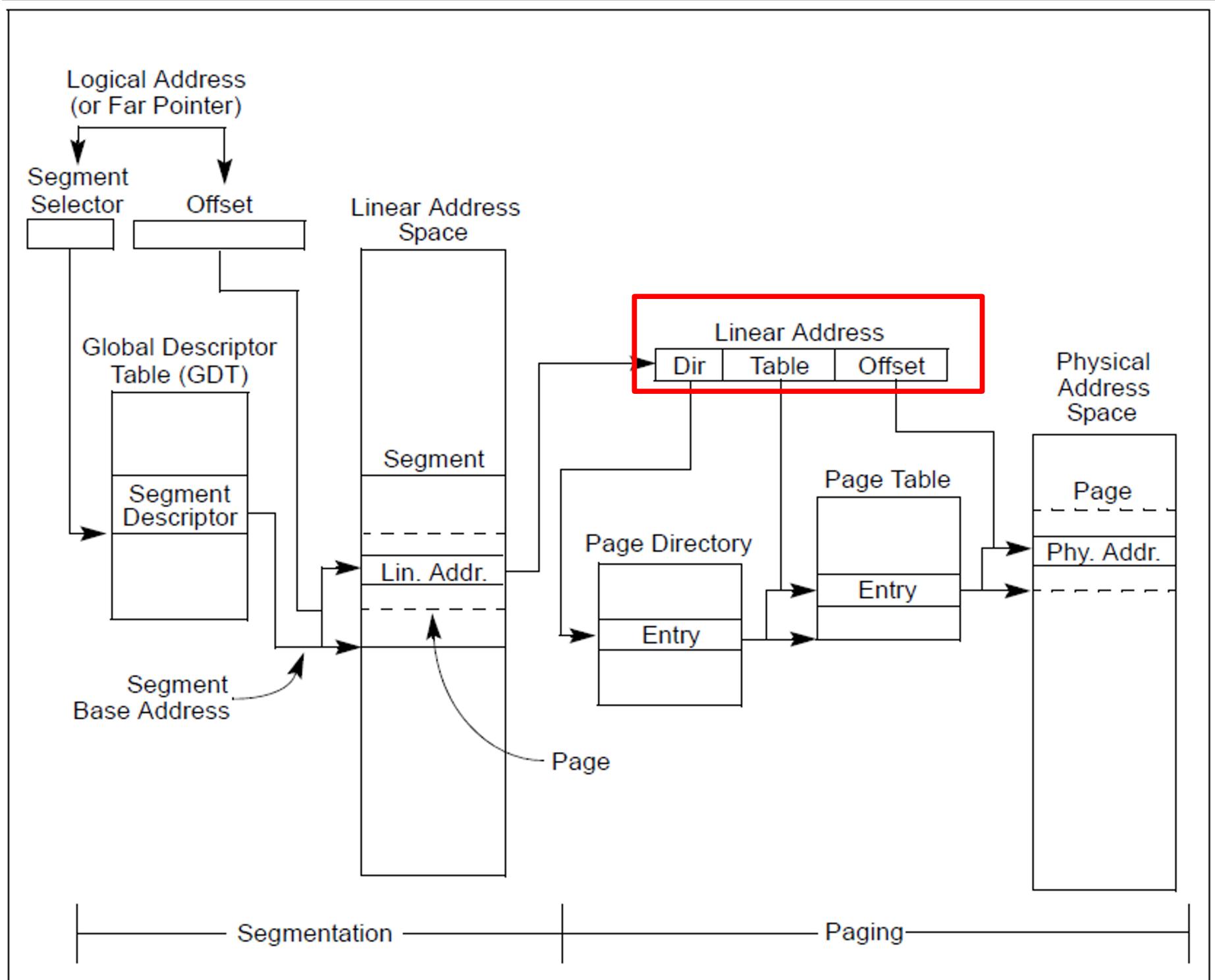


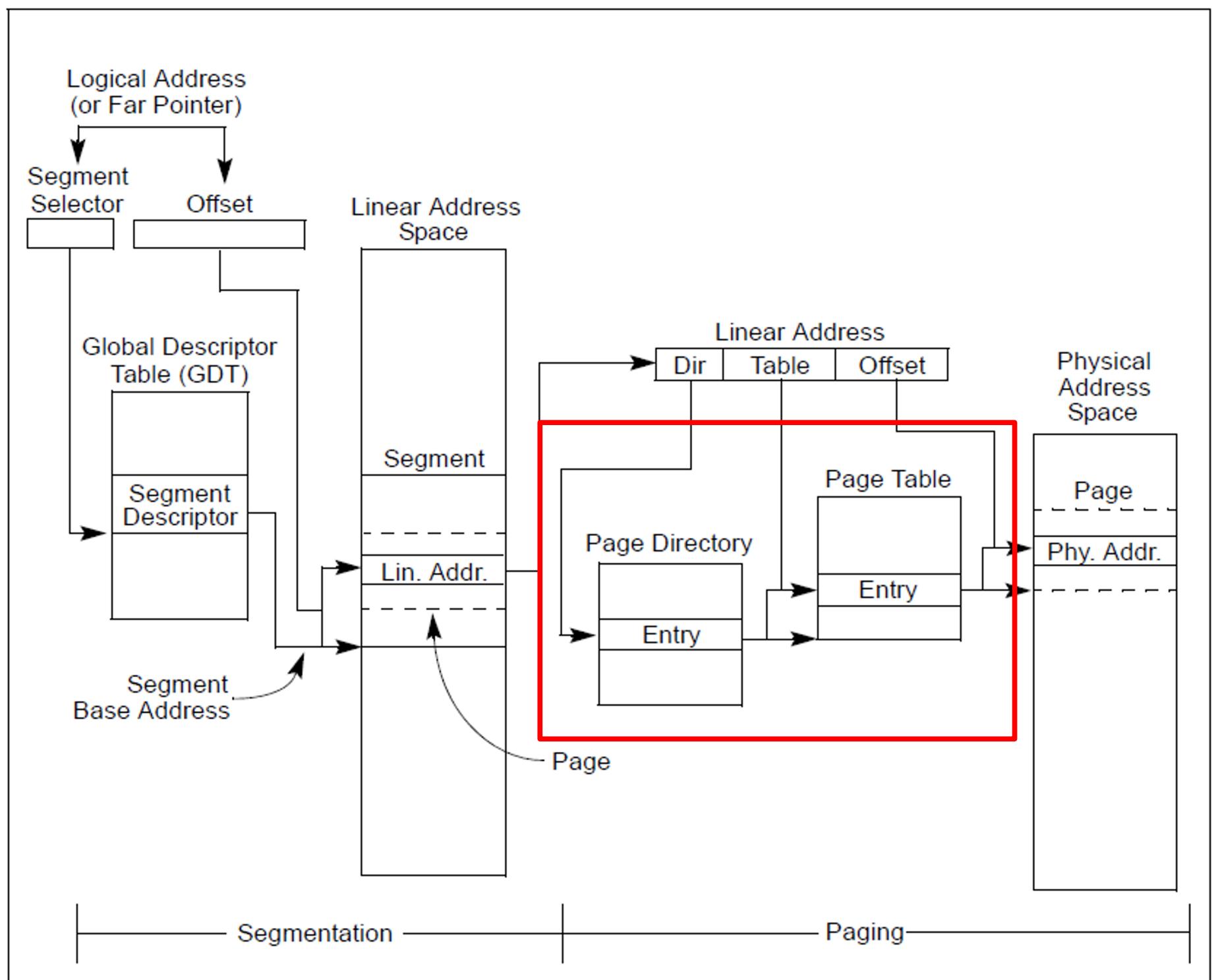


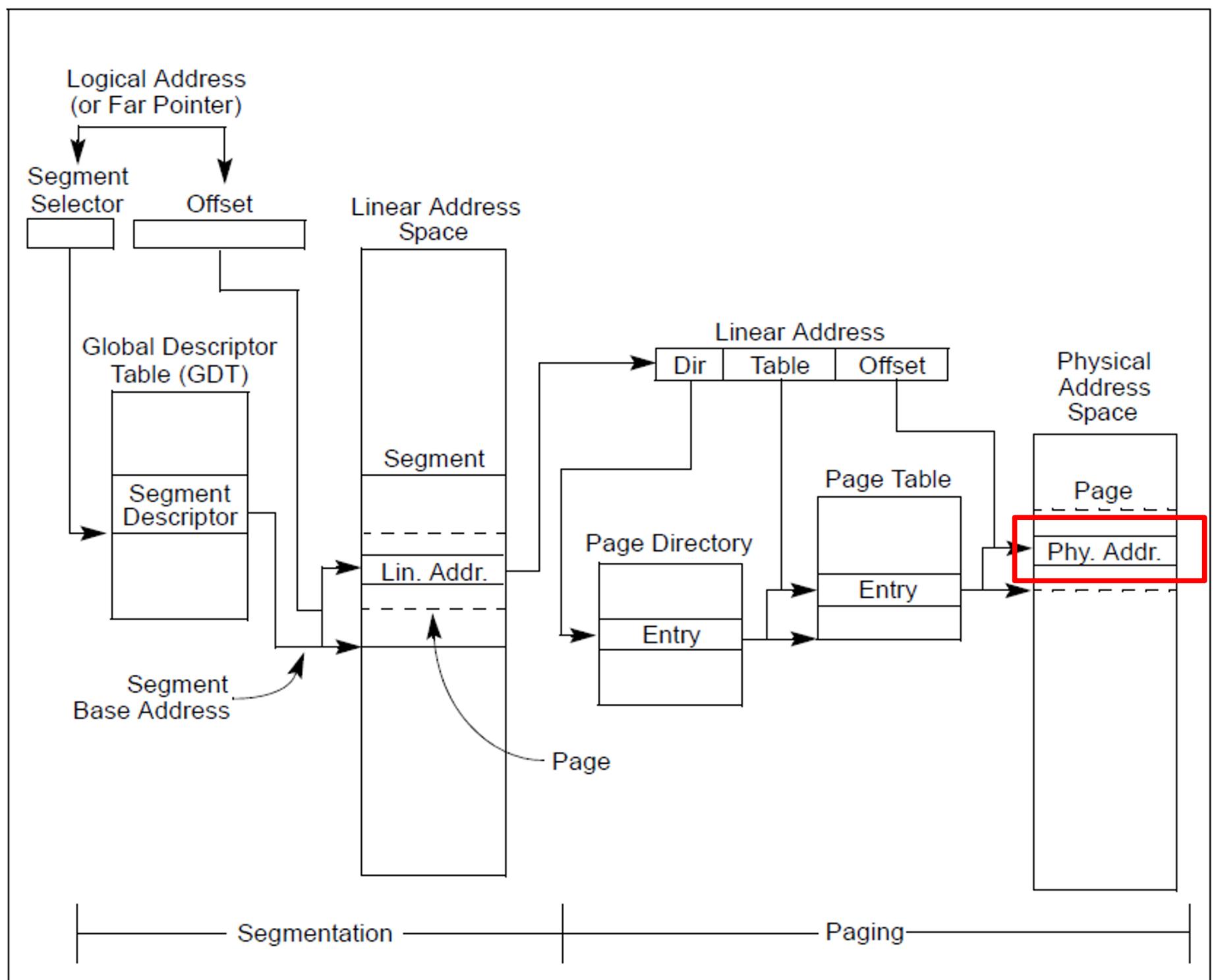


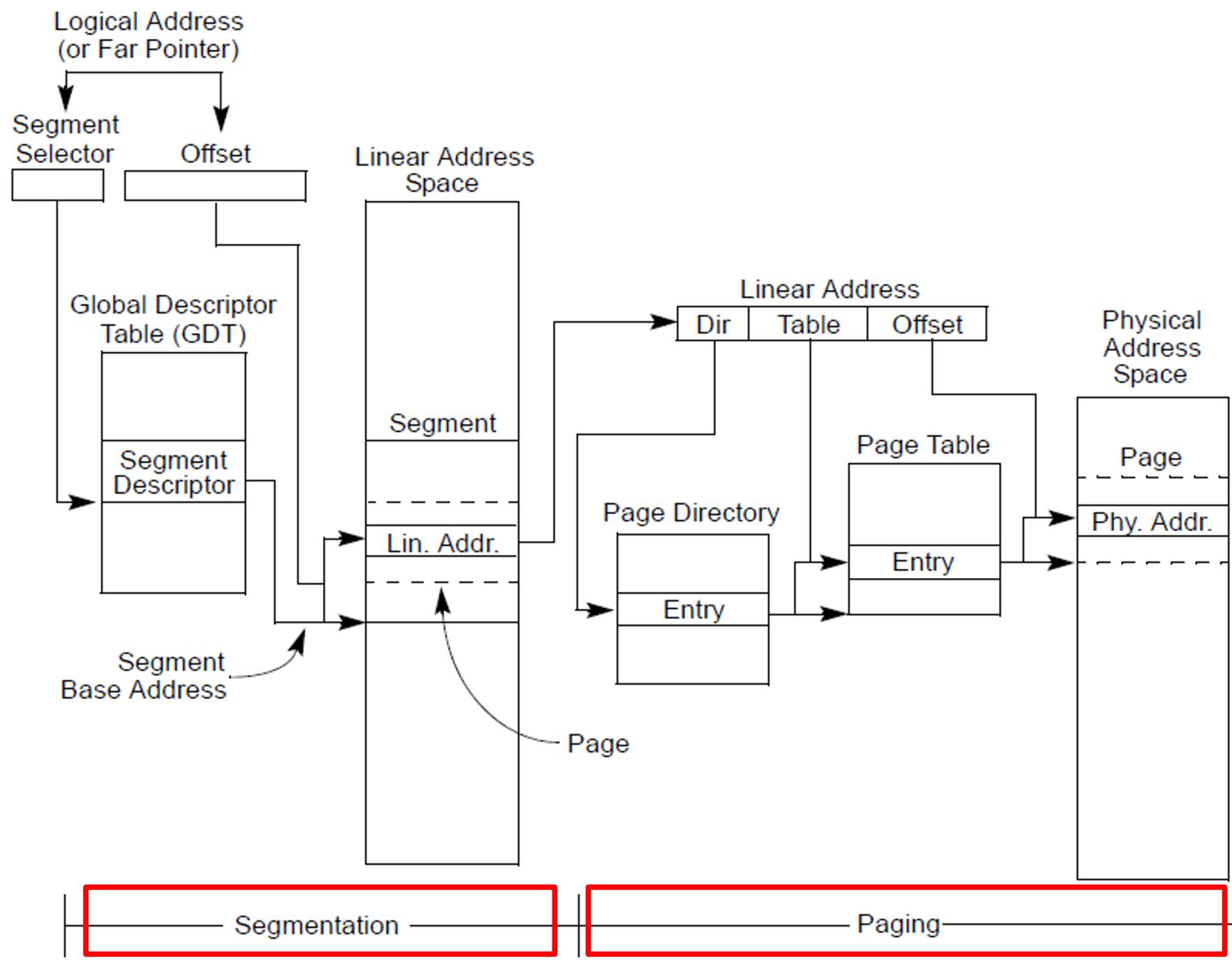








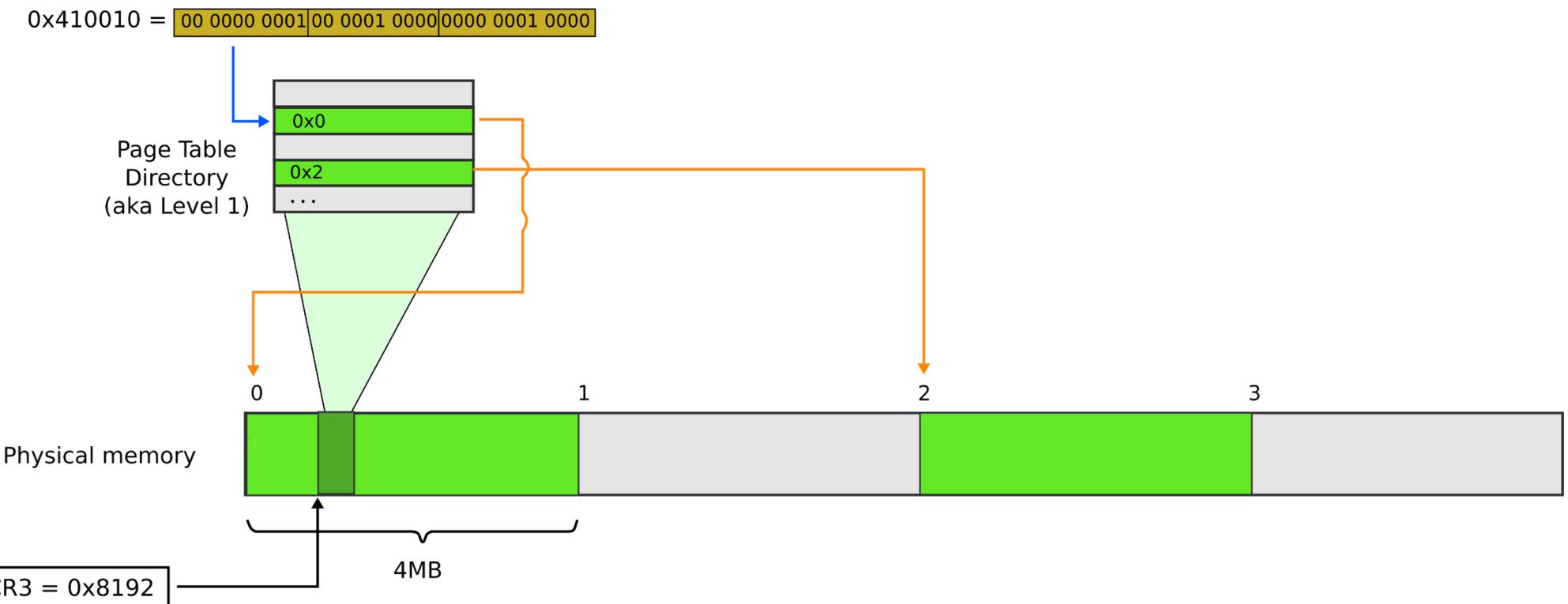




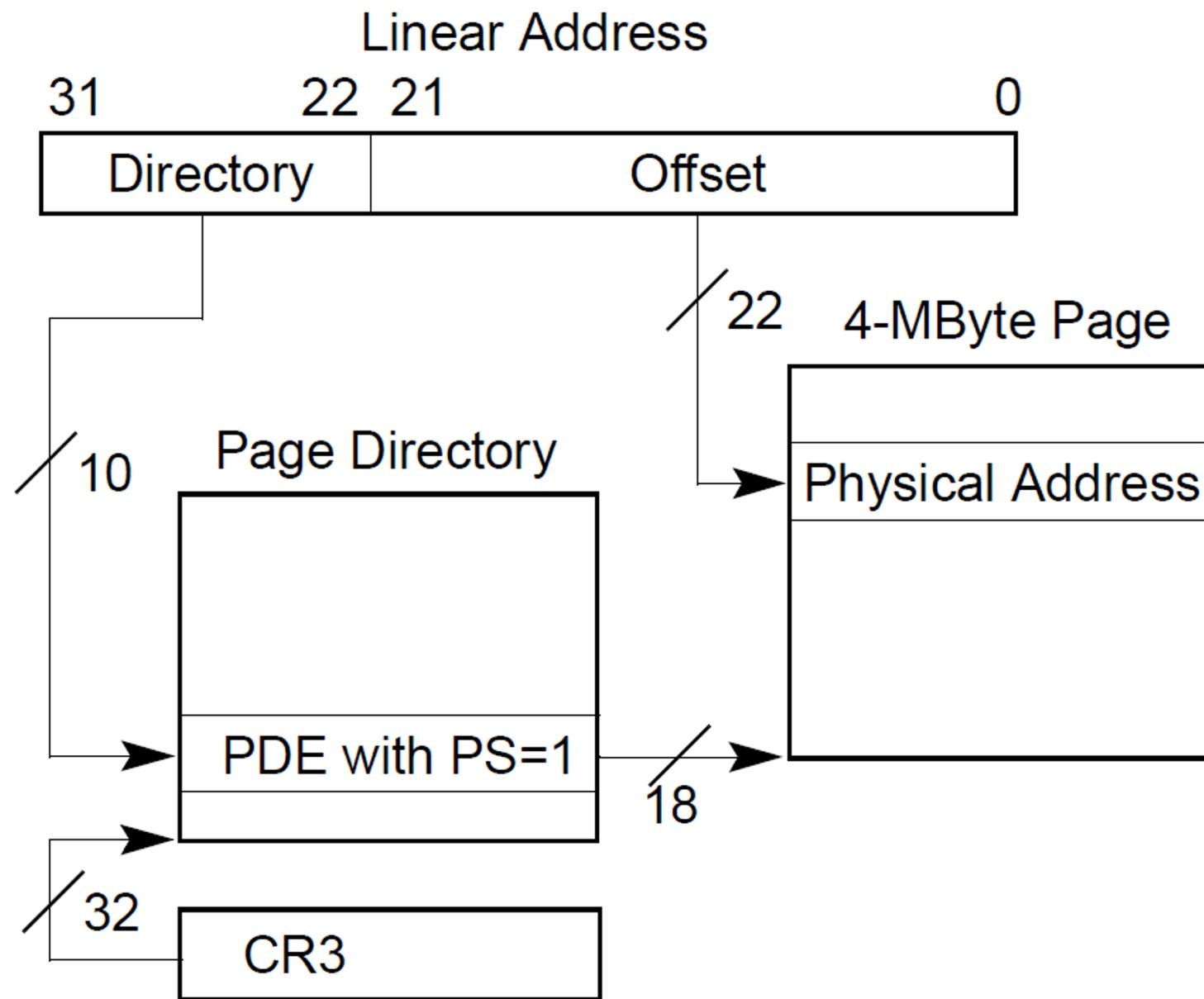
# 32bit x86 supports two page sizes

- 4KB pages
- 4MB pages

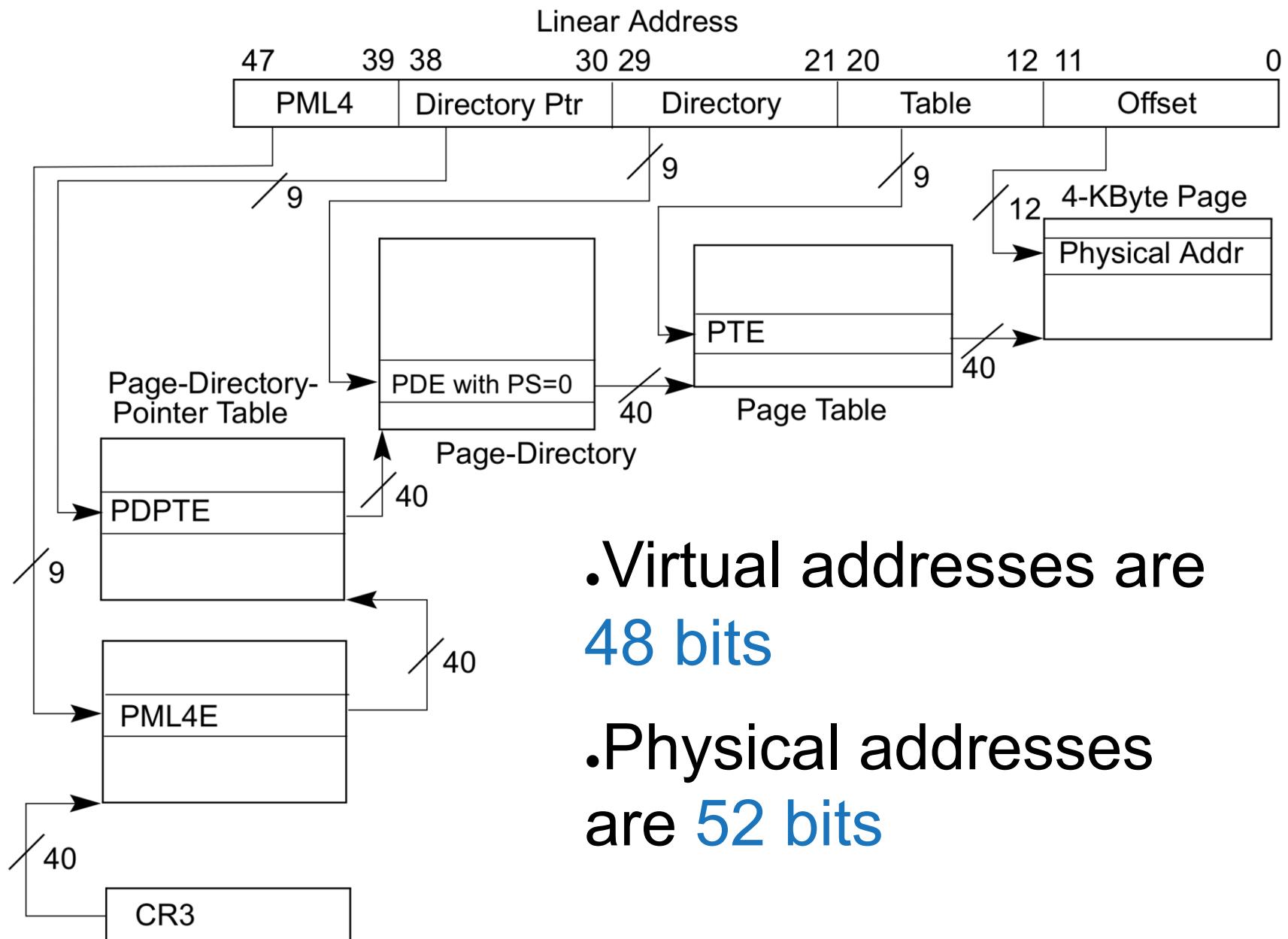
# Page translation for 4MB pages



# Page translation for 4MB pages



# Page translation in 64bit mode



# Questions?

# What pages are used for

- 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

# More paging tricks

- Determine a working set of a program?

# More paging tricks

- Determine a working set of a program?
- Use “accessed” bit

# More paging tricks

- Determine a working set of a program?
- Use “accessed” bit
- Iterative copy of a working set?
- Used for virtual machine migration

# More paging tricks

- Determine a working set of a program?
- Use “accessed” bit
- Iterative copy of a working set?
- Used for virtual machine migration
- Use “dirty” bit

# More paging tricks

- Determine a working set of a program?
- Use “accessed” bit
- Iterative copy of a working set?
- Used for virtual machine migration
- Use “dirty” bit
- Copy-on-write memory, e.g. lightweight fork()?

# More paging tricks

- Determine a working set of a program?
- Use “accessed” bit
- Iterative copy of a working set?
- Used for virtual machine migration
- Use “dirty” bit



# TLB

- .CPU caches results of page table walks
- .In translation lookaside buffer (TLB)

Virt	Phys
0xf0231000	0x1000
0x00b31000	0x1f000
0xb0002000	0xc1000
-	-



# TLB invalidation

- After every page table update, OS needs to manually invalidate cached values
- Flush TLB
  - Either one specific entry
  - Or entire TLB, e.g., when CR3 register is loaded
    - This happens when OS switches from one process to another
- This is expensive
  - Refilling the TLB with new values takes time

# Tagged TLBs

- Modern CPUs have “tagged TLBs”,
- Each TLB entry has a “tag” – identifier of a process
- No need to flush TLBs on context switch
- On Intel this mechanism is called
- Process-Context Identifiers (PCIDs)

Virt	Phys	Tag
0xf0231000	0x1000	P1
0x00b31000	0x1f000	P2
0xb0002000	0xc1000	P1

# When would you disable paging?

# When would you disable paging?

- Imagine you're running a memcached
- Key/value cache
- You serve 1024 byte values (typical) on 10Gbps connection
- 1024 byte packets can leave every 835ns, or 1670 cycles (2GHz machine)
- This is your target budget per packet

# When would you disable paging?

- Now, to cover 32GB RAM with 4K pages
- You need 64MB space
- 64bit architecture, 4-level page tables (or 5-levels now)
- Page tables do not fit in L3 cache
- Modern servers come with 32MB cache
- Every cache miss results in up to 4 cache misses due to page walk (remember 4-level page tables)
- Each cache miss is 250 cycles
- Solution: 1GB pages

# Back of the envelope

- If a page is 4K and an entry is 4 bytes, how many entries per page?

# Back of the envelope

- If a page is 4K and an entry is 4 bytes, how many entries per page?
- 1k

# Back of the envelope

- If a page is 4K and an entry is 4 bytes, how many entries per page?
- 1k
- How large of an address space can 1 page represent?

# Back of the envelope

- If a page is 4K and an entry is 4 bytes, how many entries per page?  
• 1k
- How large of an address space can 1 page represent?  
•  $1\text{k entries} * 1\text{page/entry} * 4\text{K/page} = 4\text{MB}$

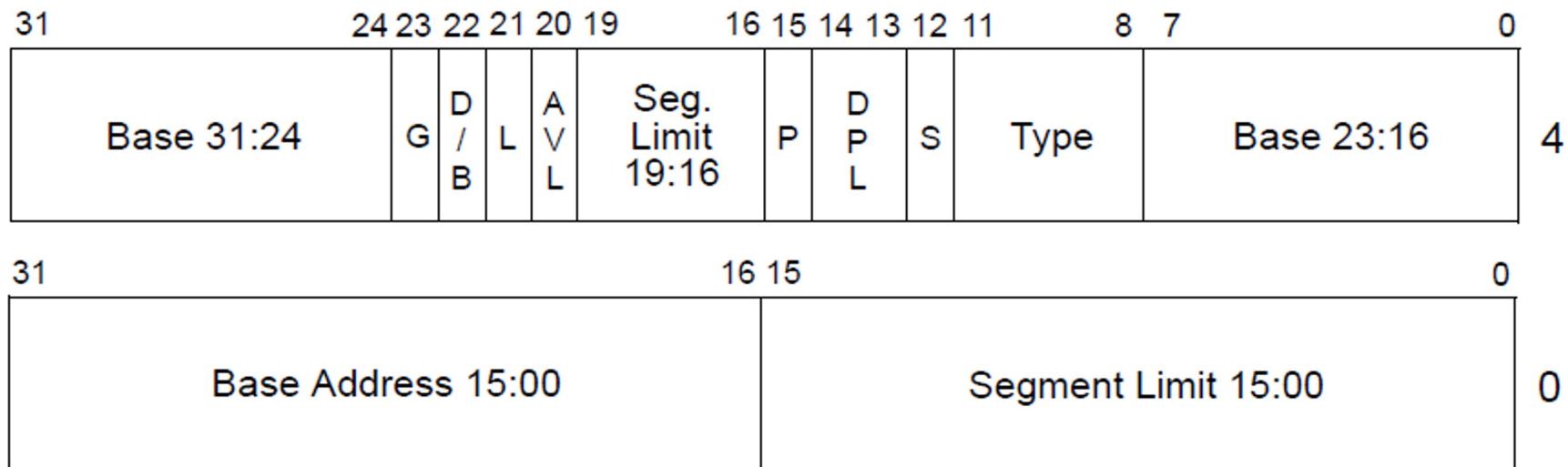
# Back of the envelope

- If a page is 4K and an entry is 4 bytes, how many entries per page?  
• 1k
- How large of an address space can 1 page represent?  
•  $1\text{k entries} * 1\text{page/entry} * 4\text{K/page} = 4\text{MB}$
- How large can we get with a second level of translation?

# Back of the envelope

- If a page is 4K and an entry is 4 bytes, how many entries per page?
  - 1k
- How large of an address space can 1 page represent?
  - $1\text{k entries} * 1\text{page/entry} * 4\text{K/page} = 4\text{MB}$
- How large can we get with a second level of translation?
  - $1\text{k tables/dir} * 1\text{k entries/table} * 4\text{k/page} = 4 \text{ GB}$
- Nice that it works out that way!

# Segment descriptors



L — 64-bit code segment (IA-32e mode only)

AVL — Available for use by system software

BASE — Segment base address

D/B — Default operation size (0 = 16-bit segment; 1 = 32-bit segment)

DPL — Descriptor privilege level

G — Granularity

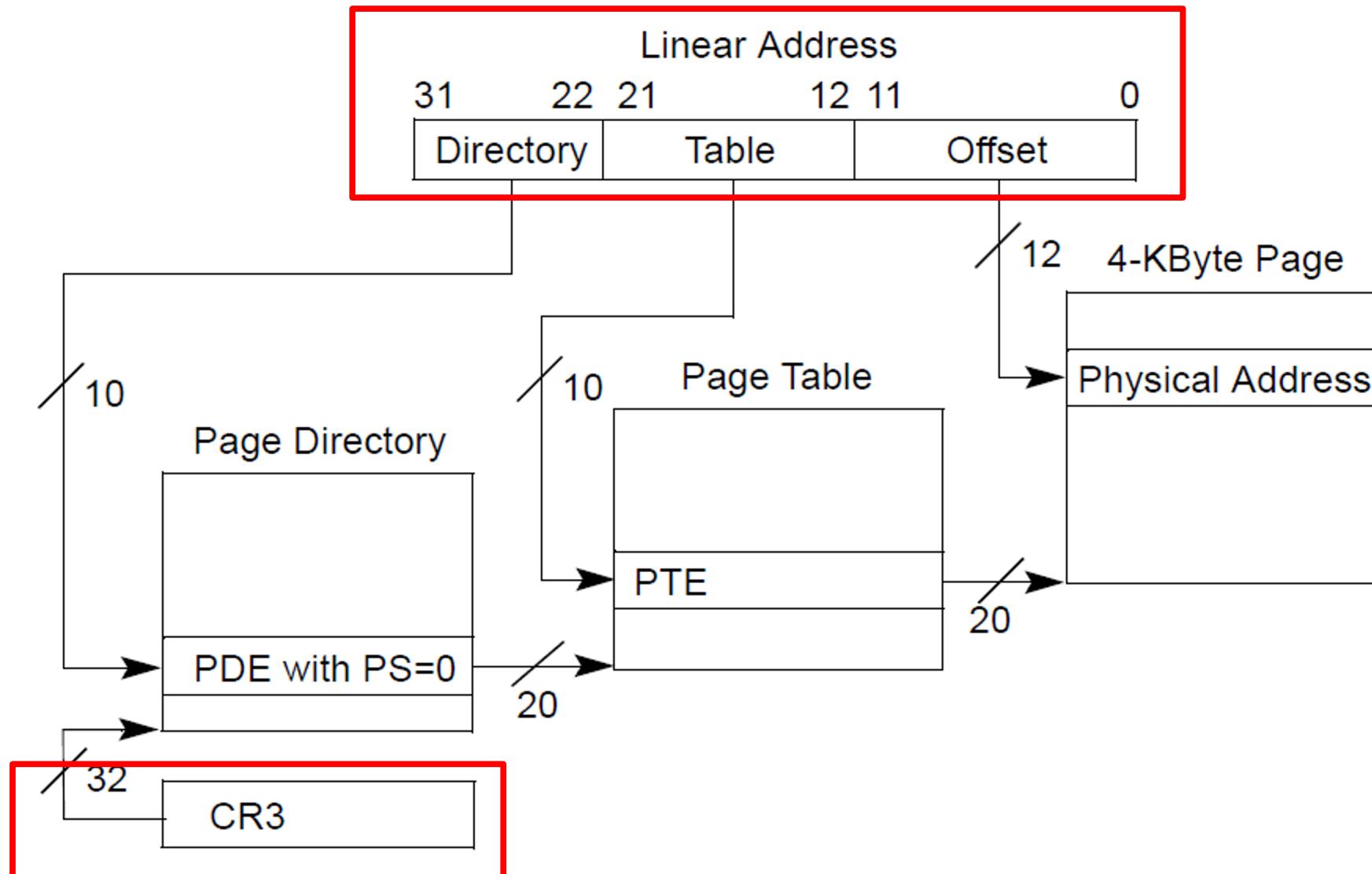
LIMIT — Segment Limit

P — Segment present

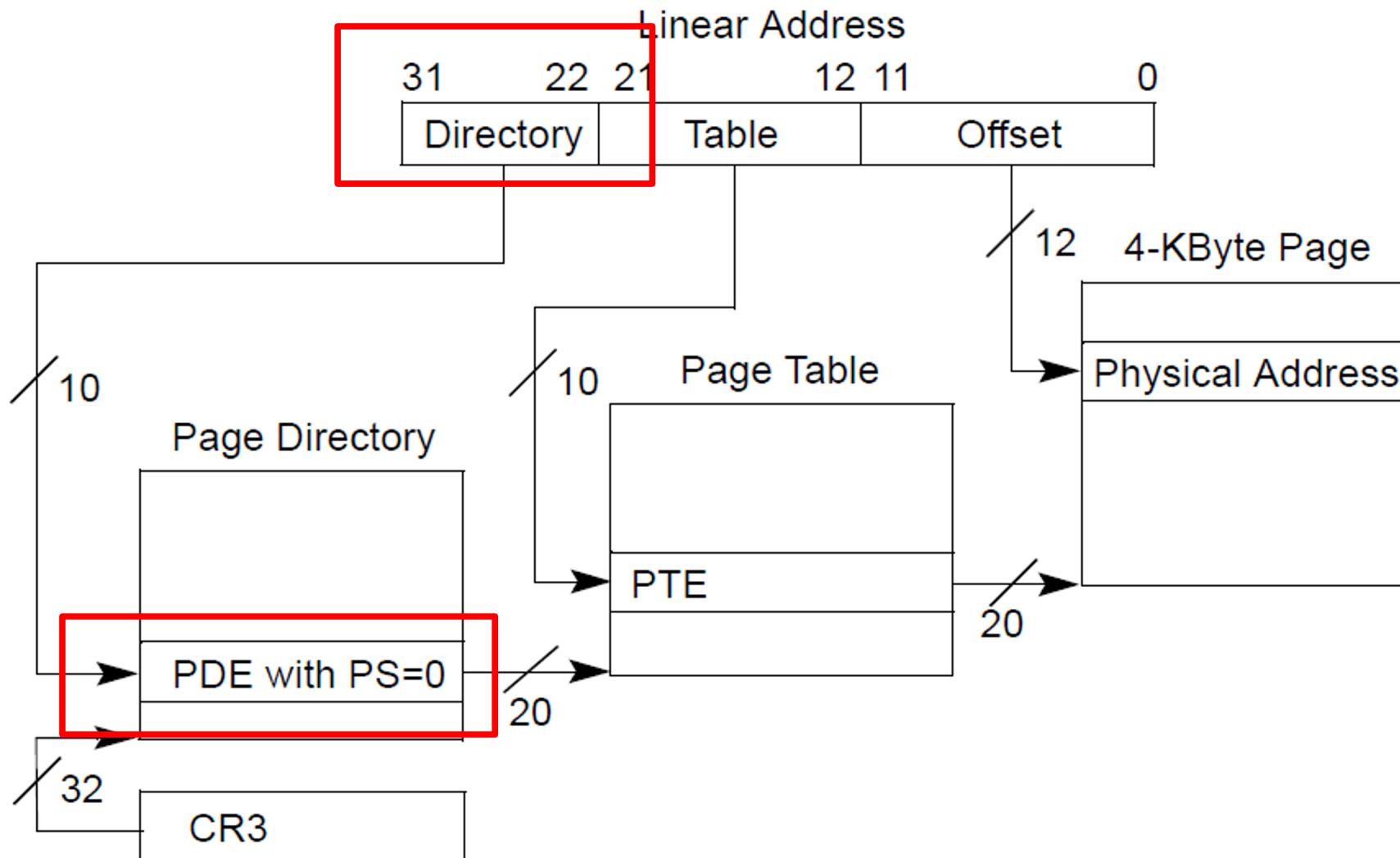
S — Descriptor type (0 = system; 1 = code or data)

TYPE — Segment type

# Page translation



# Page translation

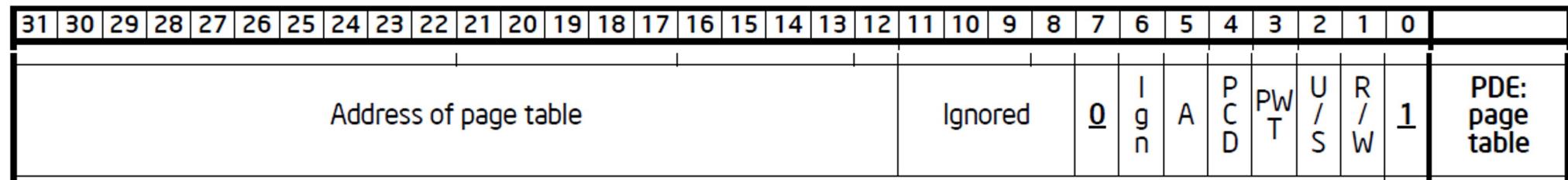


# Page directory entry (PDE)

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
Address of page table												Ignored		0	I	g	n	A	P	C	P	W	U	/	R	/	1	PDE: page table				

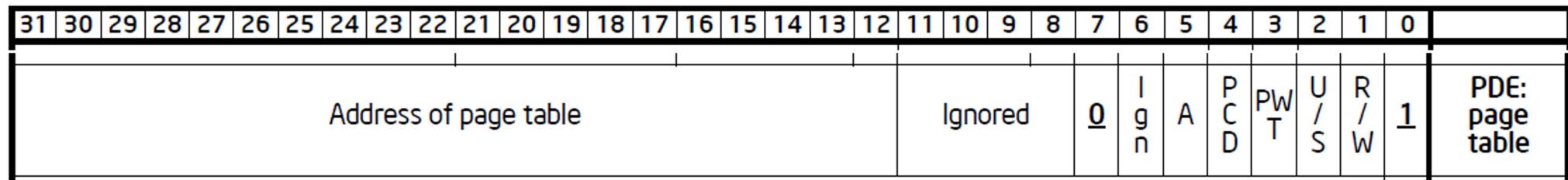
.20 bit address of the page table

# Page directory entry (PDE)



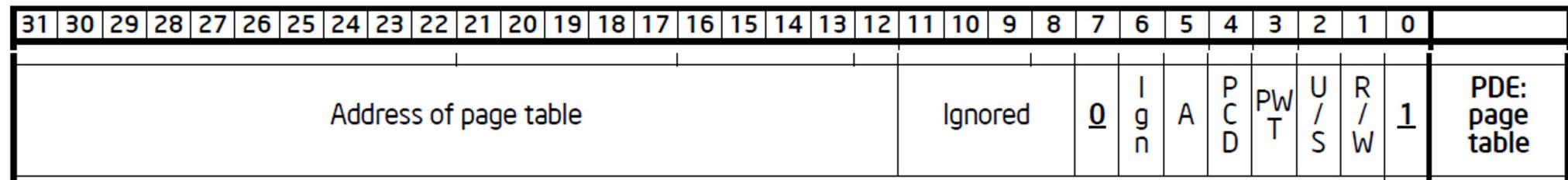
- .20 bit address of the page table
- .Wait... 20 bit address, but we need 32 bits

# 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?

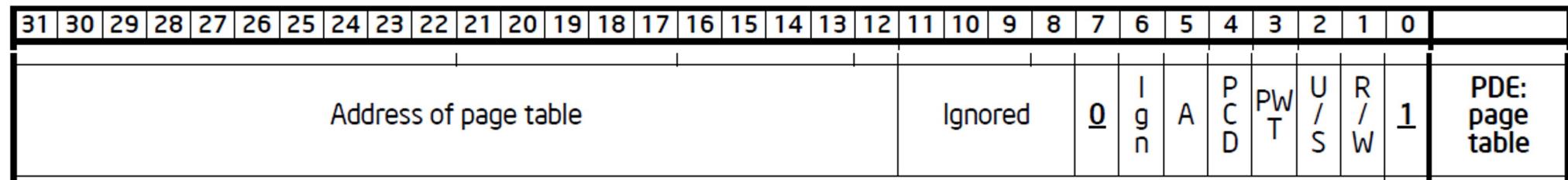
# Page directory entry (PDE)

The picture can't be displayed.

Address of page table	Ignored	0	I	g	n	A	P	C	PW	U	S	R	/	W	1	PDE: page table
-----------------------	---------	---	---	---	---	---	---	---	----	---	---	---	---	---	---	-----------------------

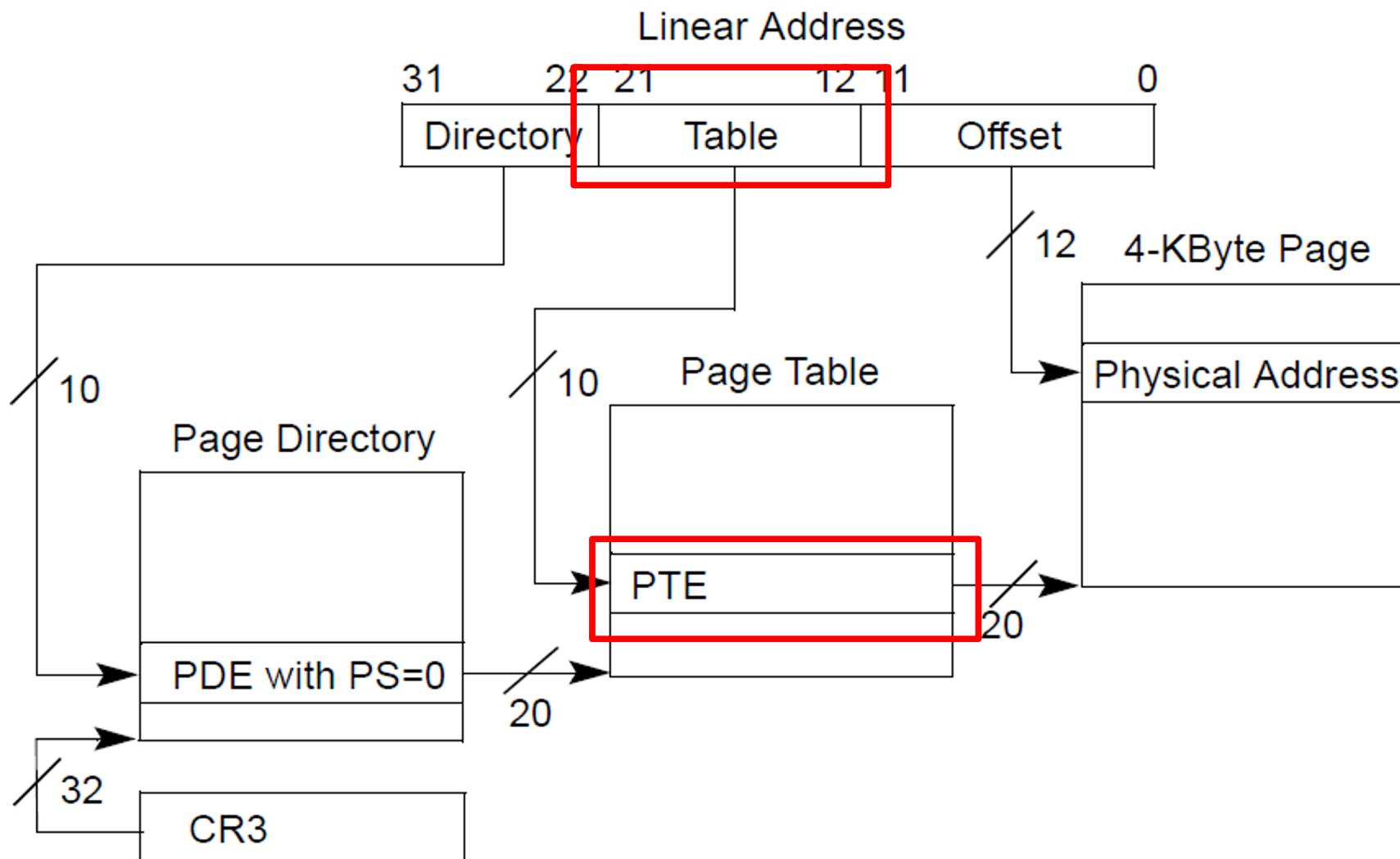
- .Bit #1: R/W – writes allowed?
- .But allowed where?
- .One page directory entry controls 1024 Level 2 page tables
  - Each Level 2 maps 4KB page
- .So it's a region of  $4\text{KB} \times 1024 = 4\text{MB}$

# Page directory entry (PDE)

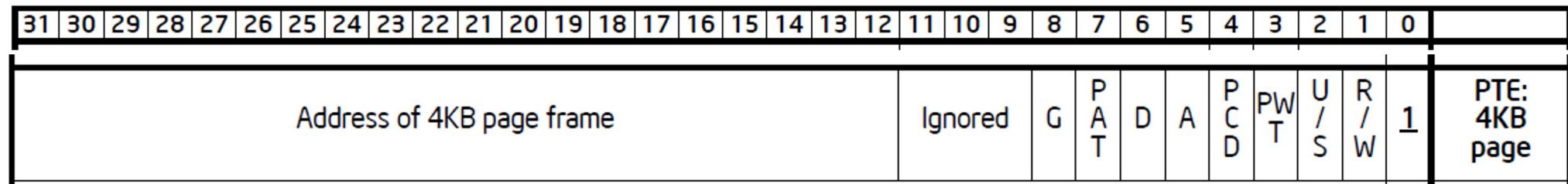


- .Bit #2: U/S – user/supervisor
- .If 0 – user-mode access is not allowed
- .Allows protecting kernel memory from user-level applications

# 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

 The picture can't be displayed.



Questions?