

Paging

Abhilash Jindal

Overview

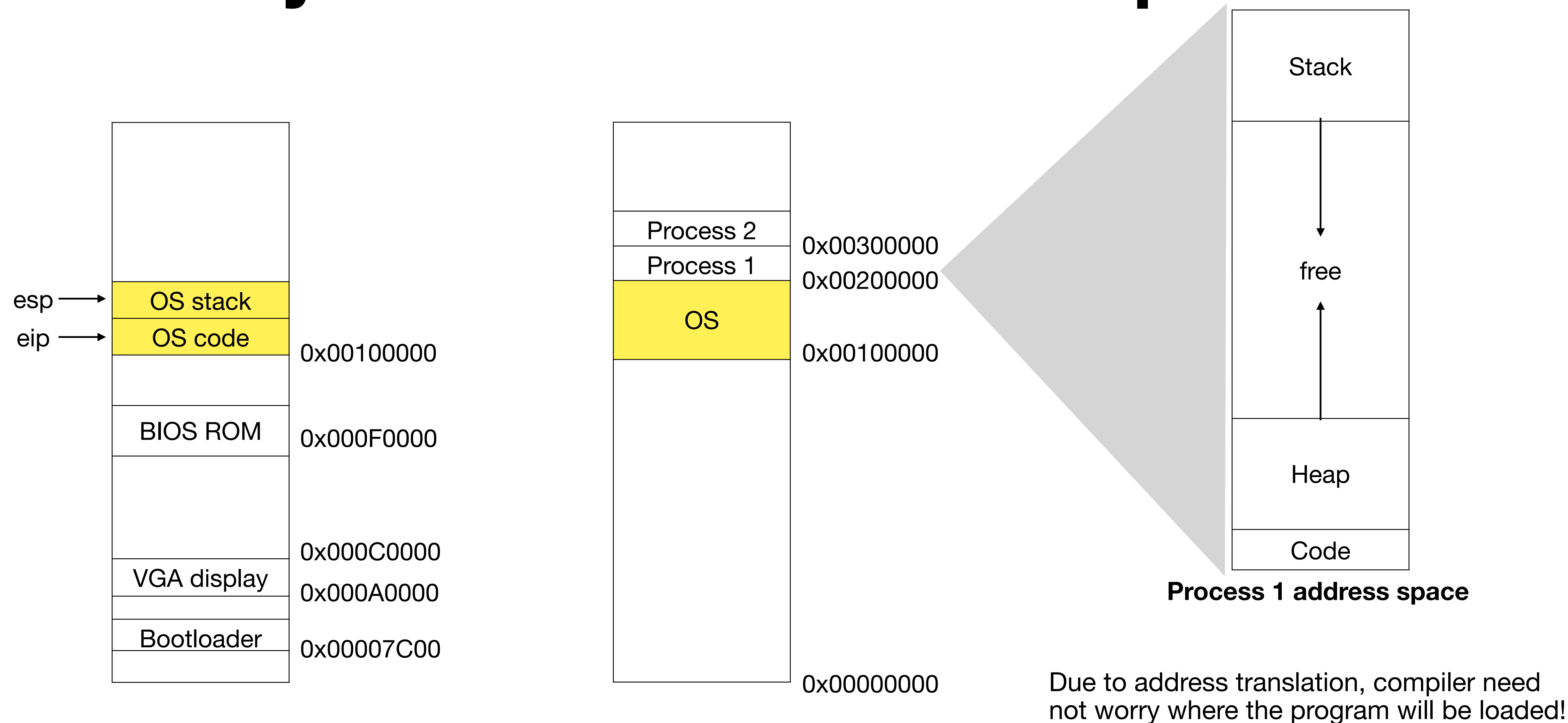
- More flexible address translation with paging (OSTEP Ch 18-20)
 - Paging hardware
- Demand paging: swapping pages to disk when memory becomes full (OSTEP Ch 21-22)
 - Swapping mechanisms
 - Page replacement algorithms
- Paging in action (xv6 book Ch 2, OSTEP Ch 23)
 - Paging on xv6
 - Fork with Copy-on-write, Guard pages

Paging Hardware

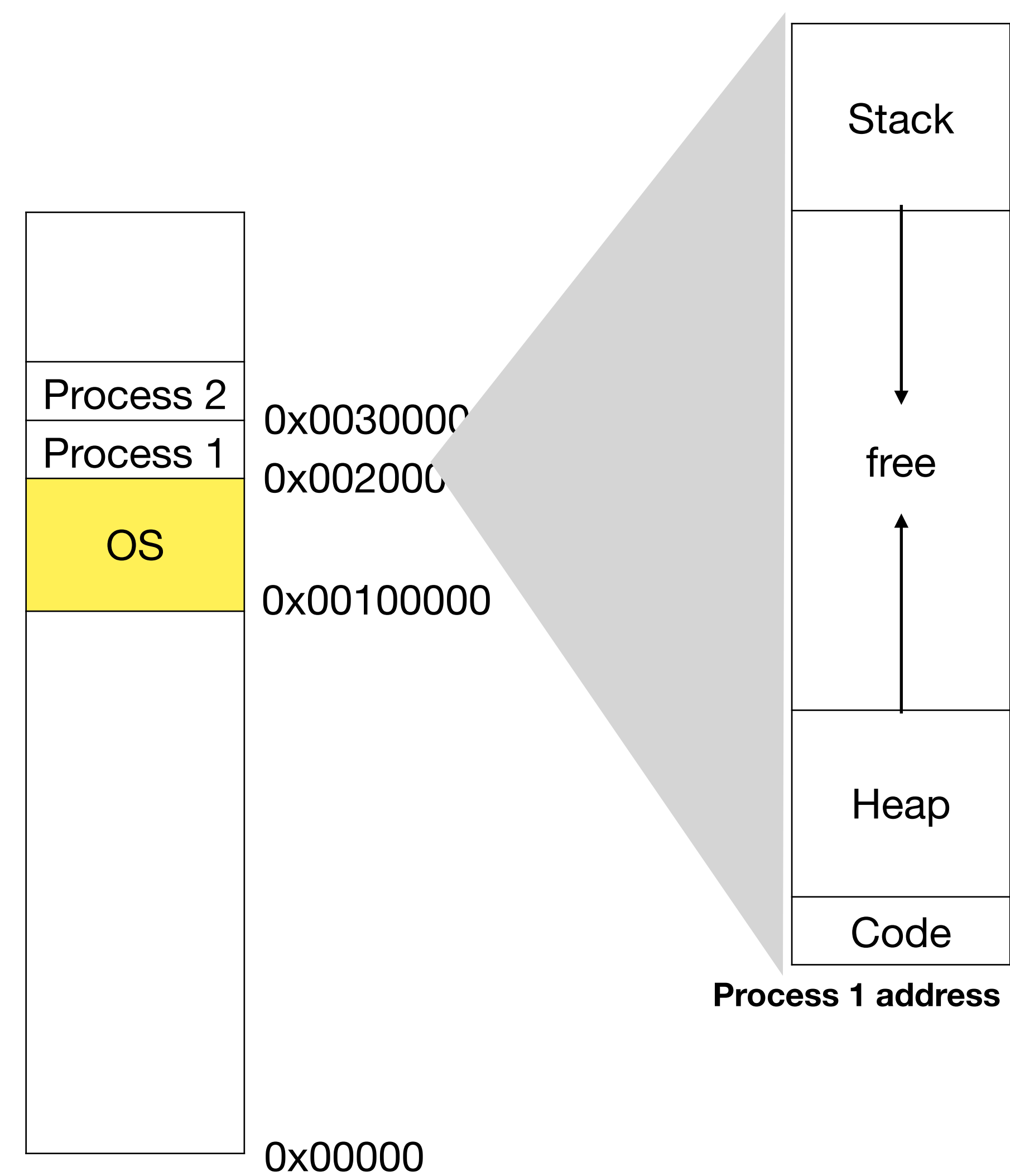
OSTEP Ch 18-20

Intel SDM Volume 3A Ch 4

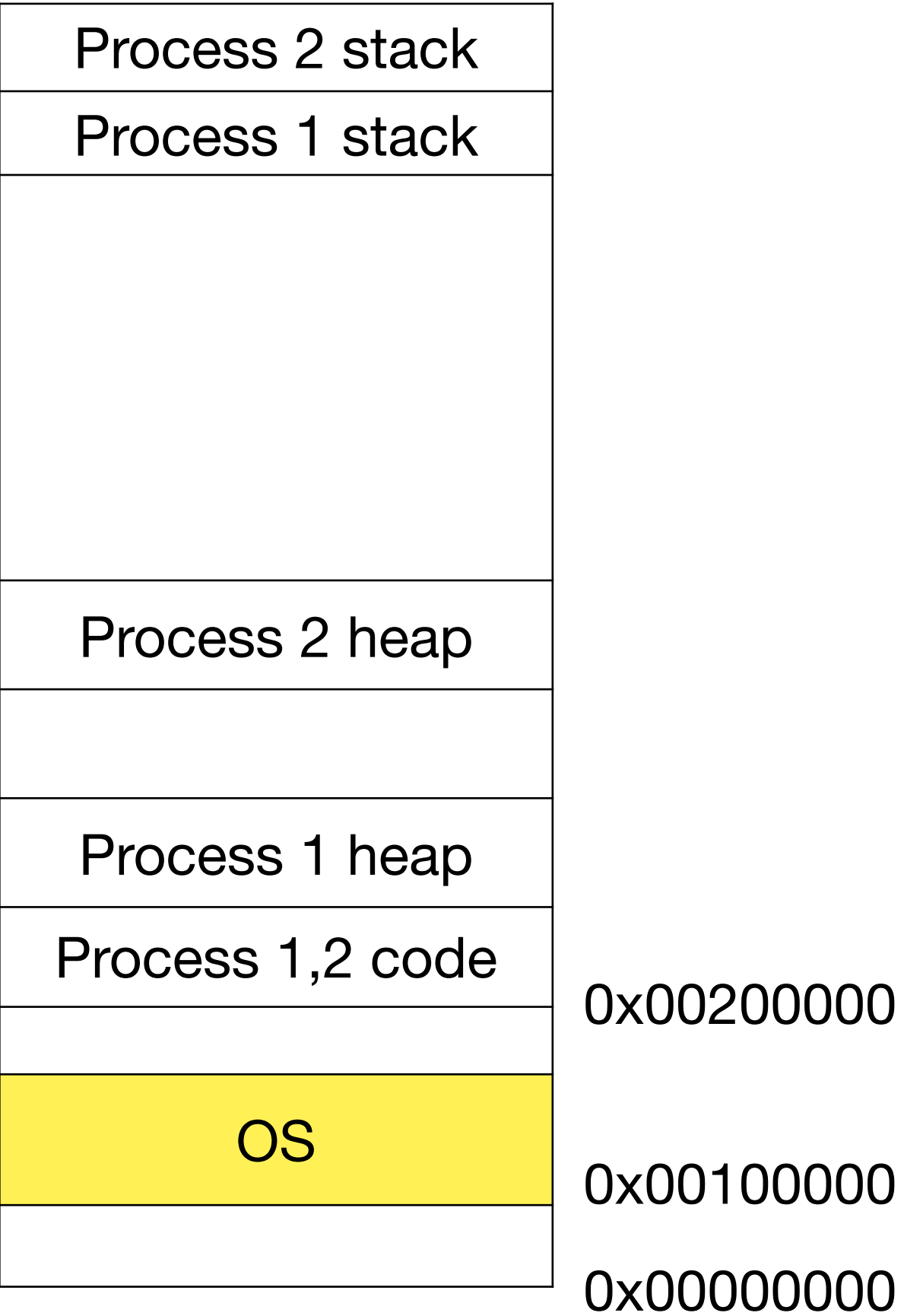
Memory isolation and address space



Segmentation

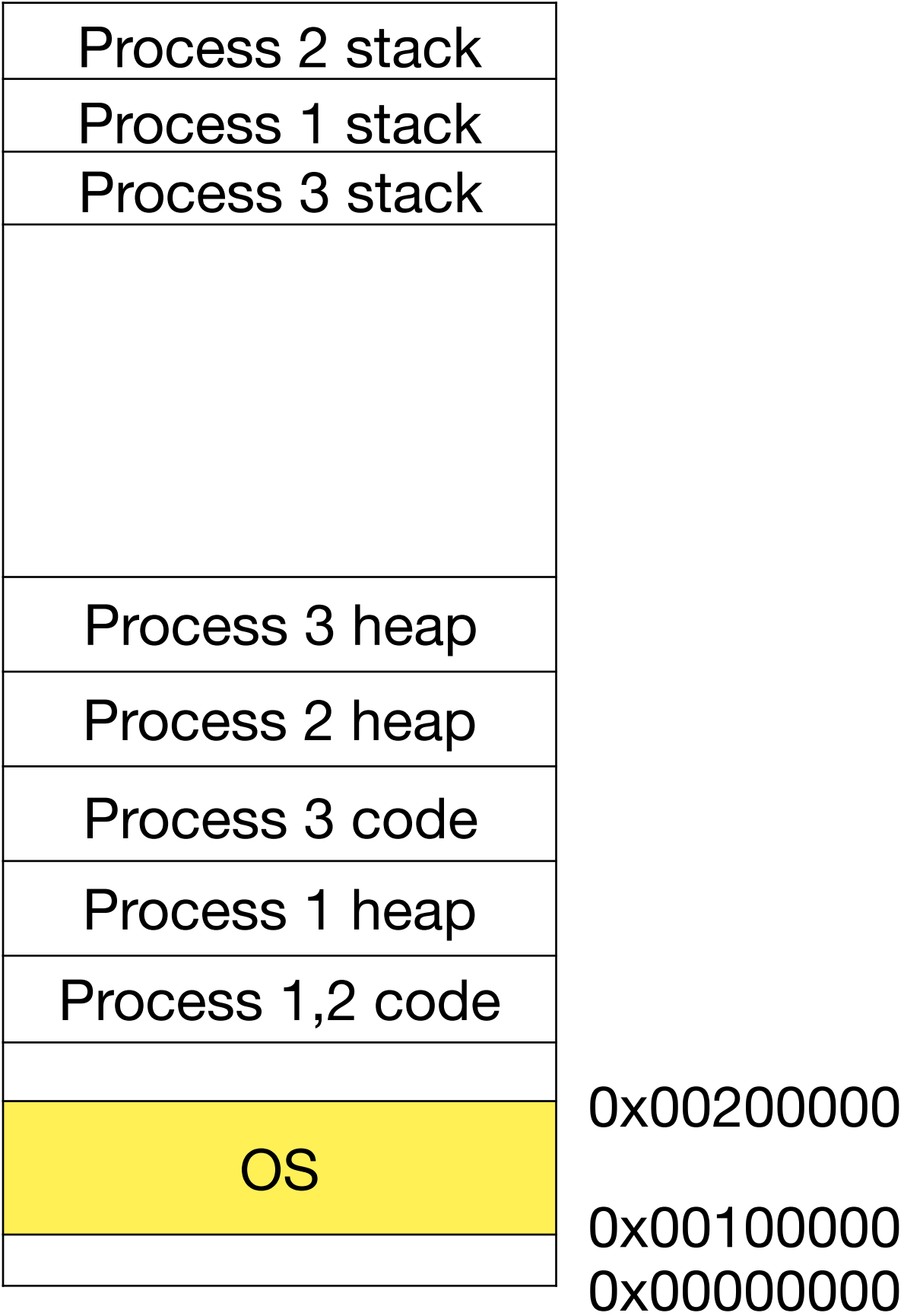


- Mapping large address spaces
- Place each segment independently to not map free space



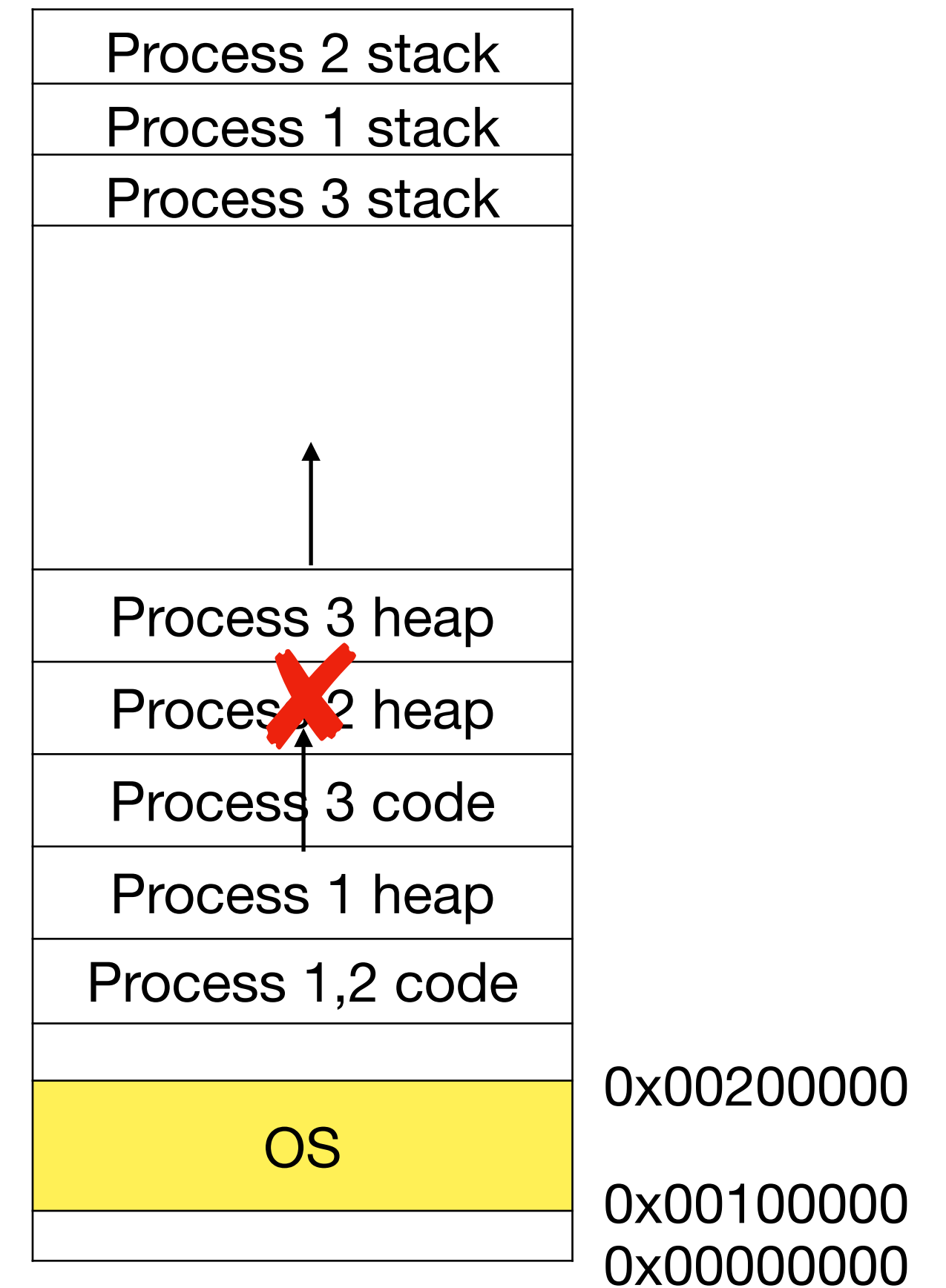
Allocating memory to a new process

- Find free spaces in physical memory
 - Difficult because segments can be of arbitrary sizes



Growing and shrinking address space

- Segments need to be contiguous in memory
- Growing might not succeed if there are other segments next to heap segment



External fragmentation

- After many processes start and exit, memory might become “fragmented” (similar to disk)
 - Example: cannot allocate 20 KB segment
- Compaction: copy all allocated regions contiguously, update segment base and bound registers
 - Copying is expensive
 - Growing heap becomes not possible

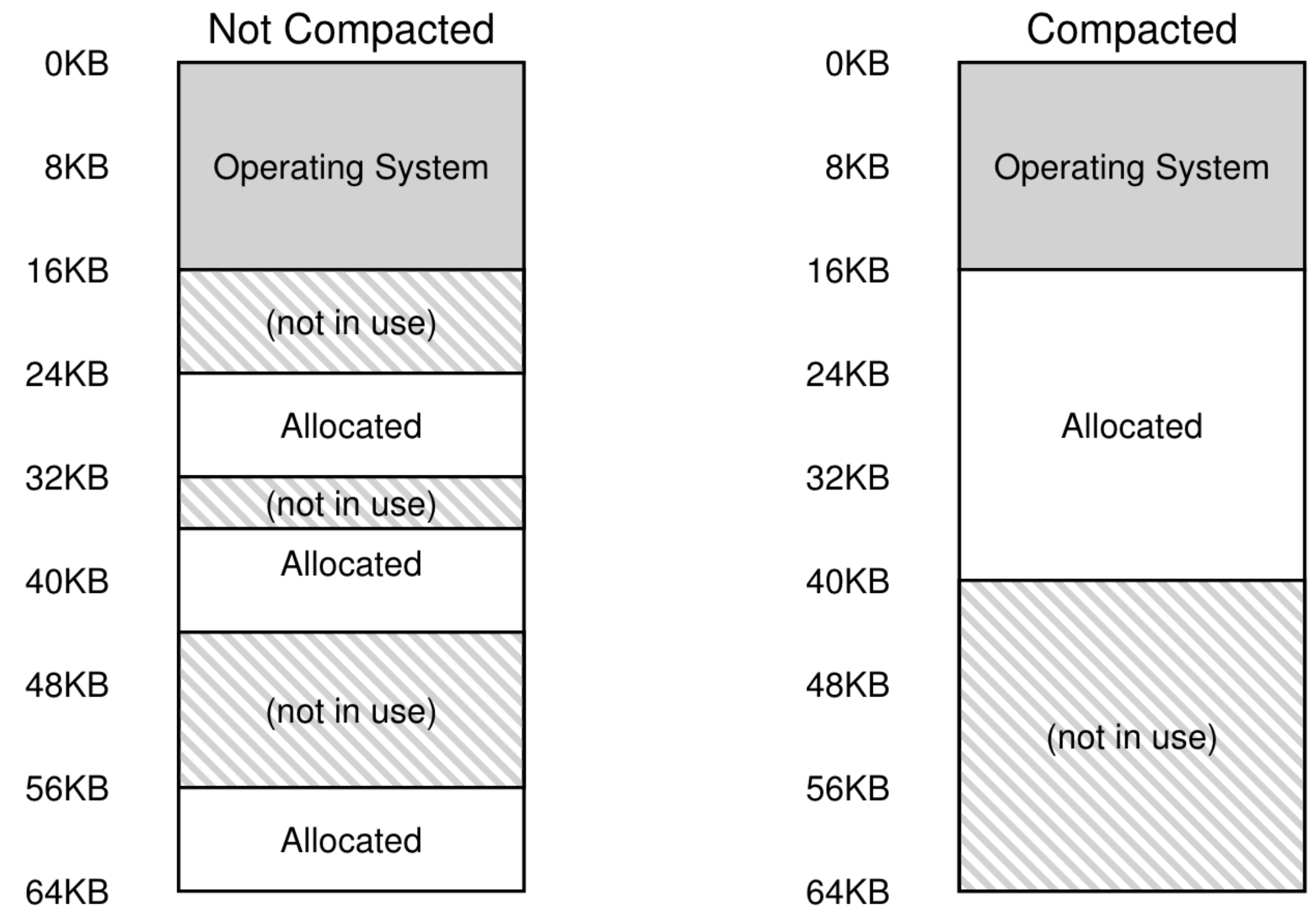


Figure 16.6: Non-compacted and Compacted Memory

Limitations of segmentation

Limited flexibility:

- To support large address spaces, burden falls on programmer/compiler to manage multiple segments
- Only an entire segment can be shared. Example: cannot share some part of CS (both processes use the same library)

Different sized segments, segments need to be contiguous in physical memory

- complicates physical memory allocator
- lead to external fragmentation
- growing/shrinking segments is awkward

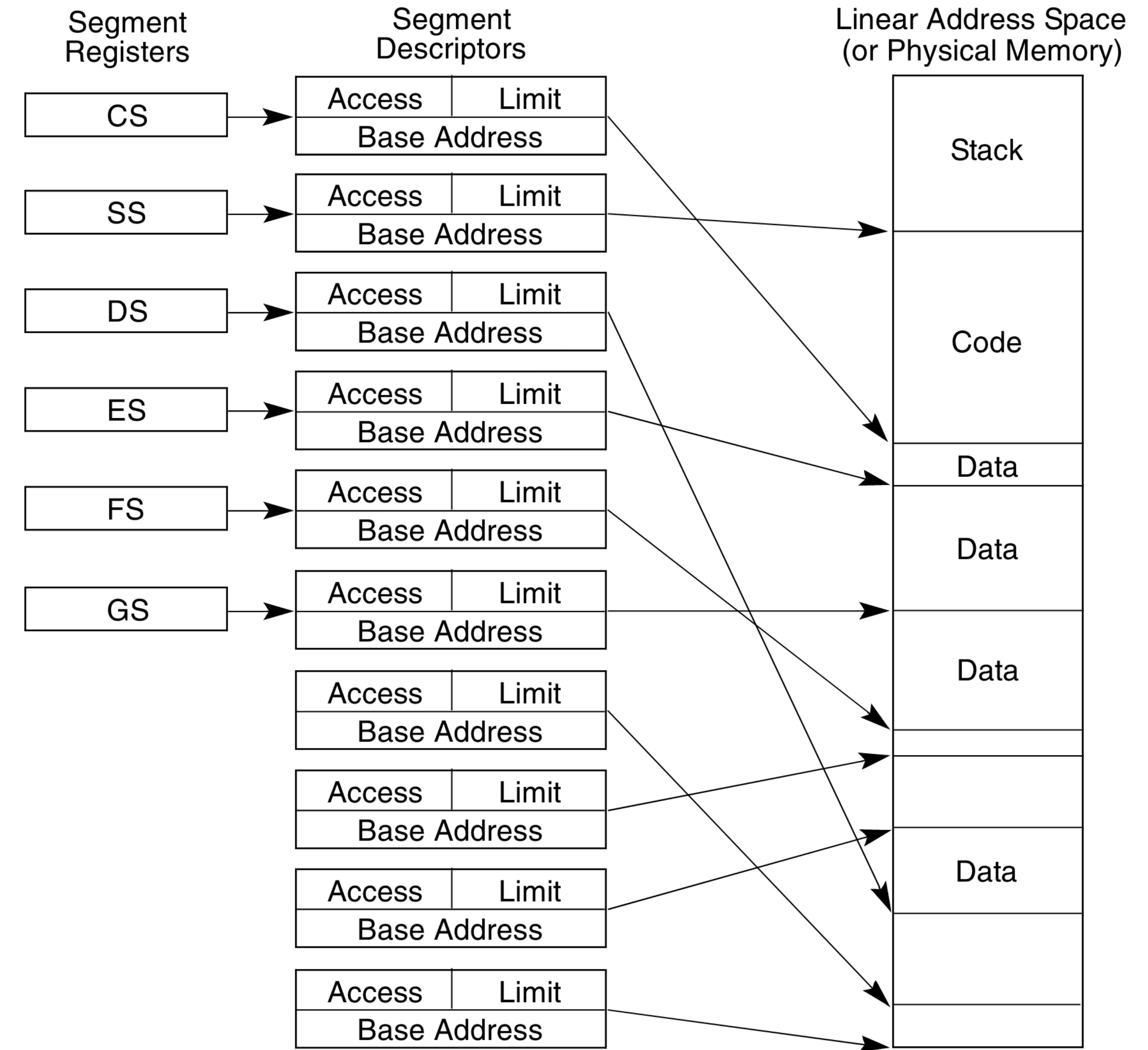


Figure 3-4. Multi-Segment Model

Paging

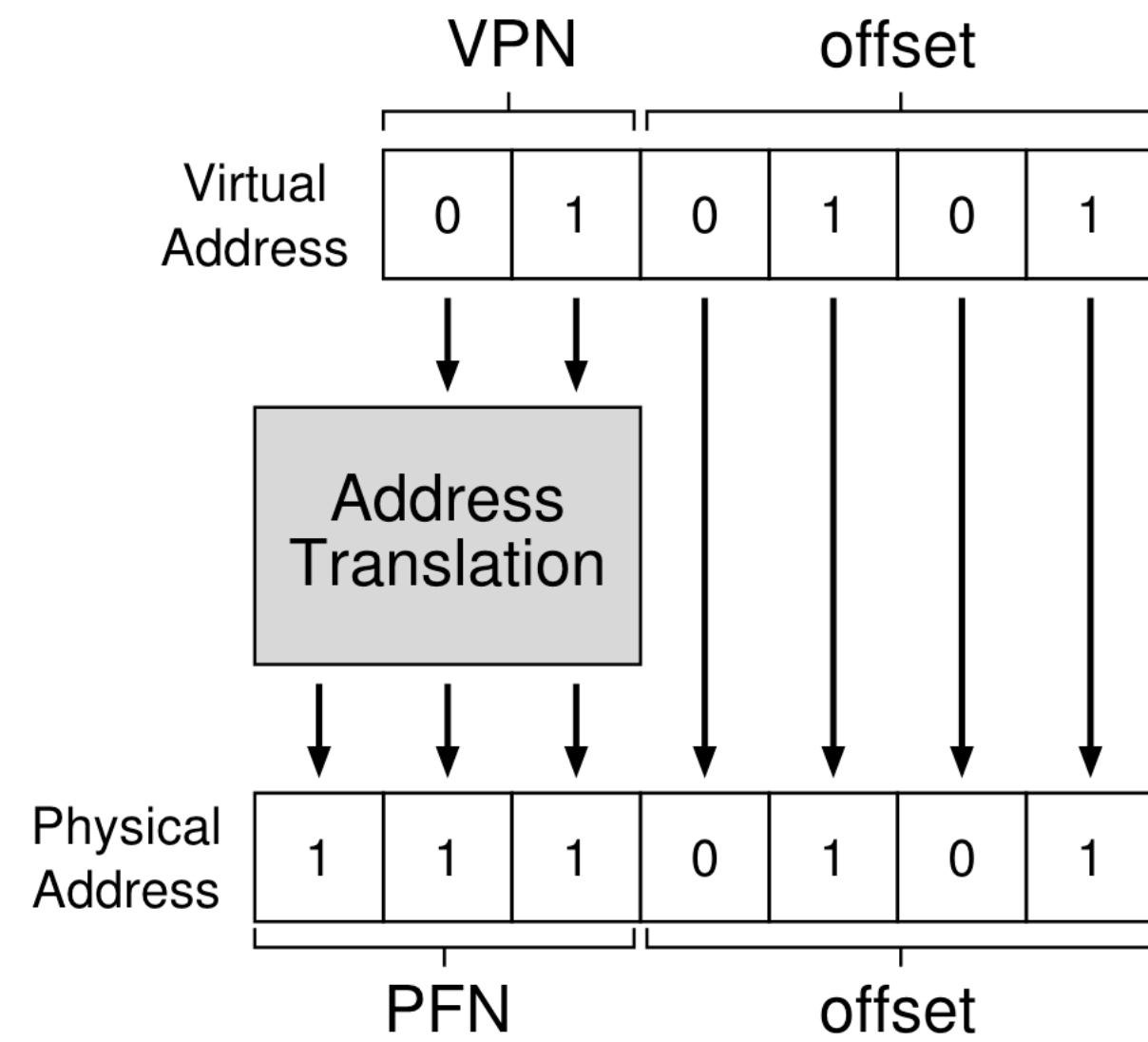
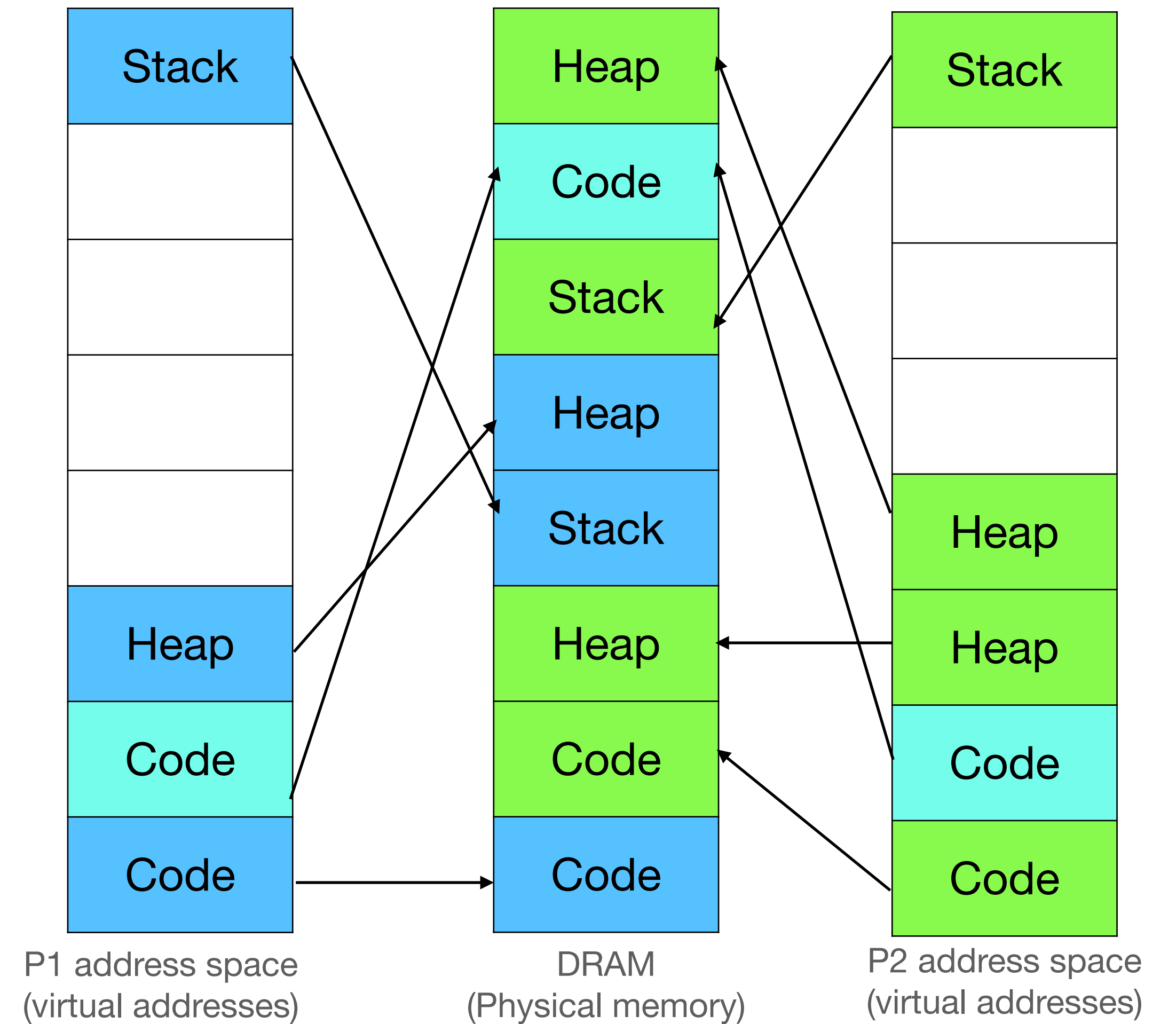


Figure 18.3: The Address Translation Process

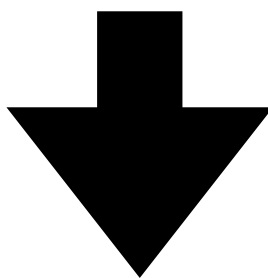


Notebook analogy

Segmentation

Preparing for OS exam:

- Read second letter from 3rd page



- Read second letter from 4th page



0	1	2	3	4	5
xv6	is	an	OS	for	x86



0	1	2	3
Write	an	SQL	query



Notebook

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
OS:1 DB:7	xv6	is	an	OS	for	x86	Write	an	SQL	query				

Notebook analogy

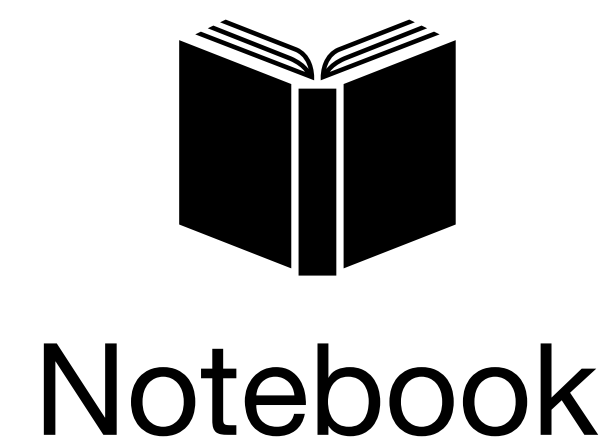
Paging



0	1	2	3	4	5
xv6	is	an	OS	for	x86



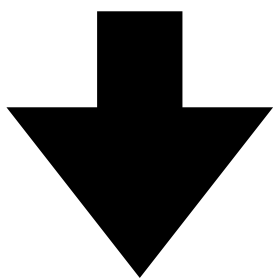
0	1	2	3
Write	an	SQL	query



0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
OS:1 DB:5	0:3,1:4, 2:2,3:8, 4:6,5:10	an	xv6	is	0:9,1:2, 2:7,3:12	for	SQL	OS	Write	x86		query		

Preparing for OS exam:

- Read second letter from 3rd page



- Read second letter from 8th page

Paging

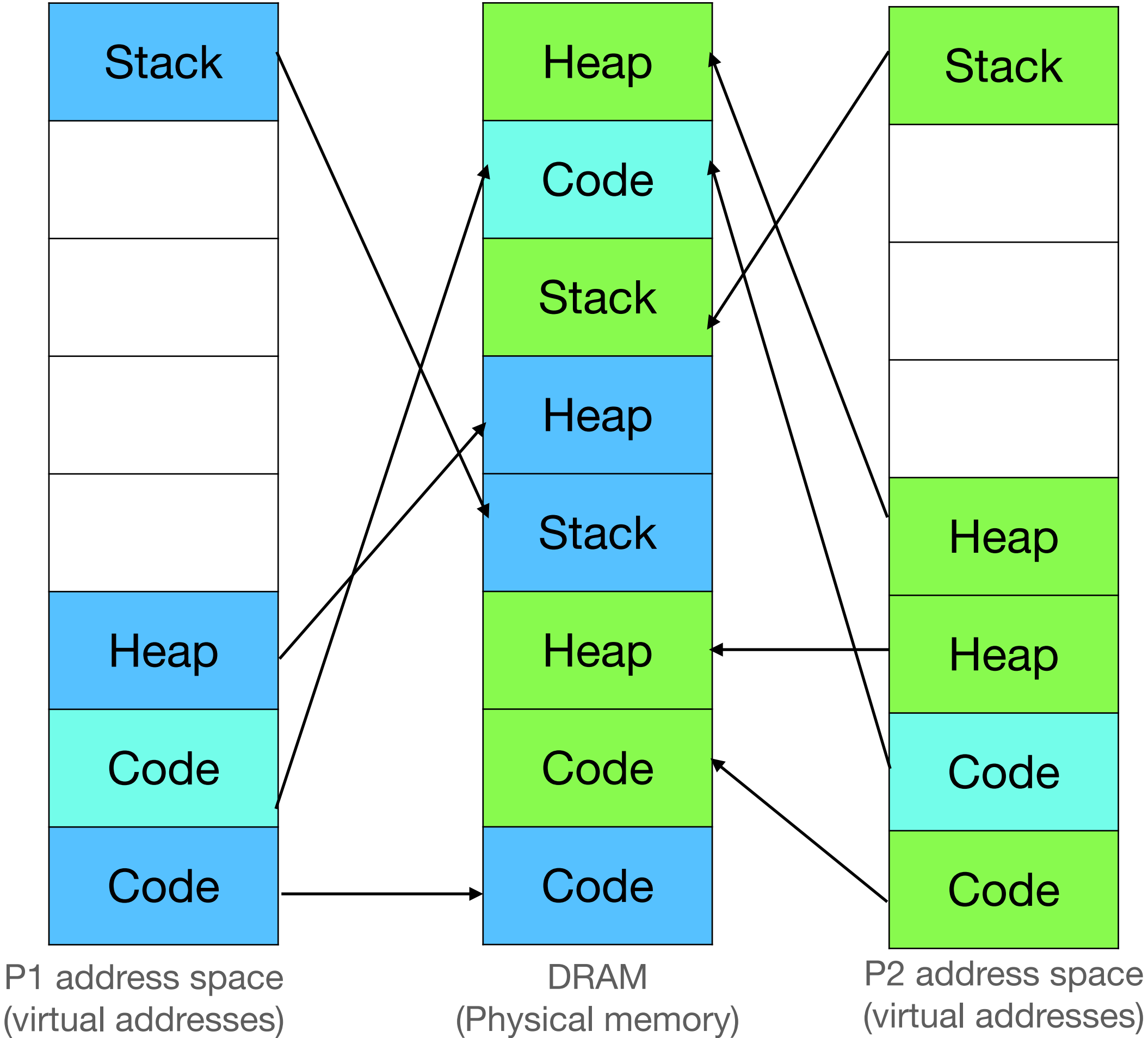
Segmentation	Paging
Large address spaces need multi segment model. Burden falls on programmer/compiler to manage multiple segments	Transparently support large address spaces. Programmer/compiler work with a flat virtual address space
Different sized segments lead to external fragmentation	Fixed sized pages (4KB, 4MB)
Different sized segments complicate physical memory allocator	
Full segment needs to be contiguous in physical memory	Page is contiguous. Neighbouring addresses (in different pages) may not be contiguous
Growing/shrinking segments is awkward	Simple: allocate another page, free a page
Address translation hardware has an adder ($va + base$) and a comparator ($va < limit$)	Address translation hardware is much more complicated

How to do page-based address translation?

Page table: Maintain a lookup table for each process

Virtual page number	Physical page number
7	3
6	x
5	x
4	x
3	x
2	4
1	6
0	0

Virtual page number	Physical page number
7	5
6	x
5	x
4	x
3	7
2	3
1	6
0	1

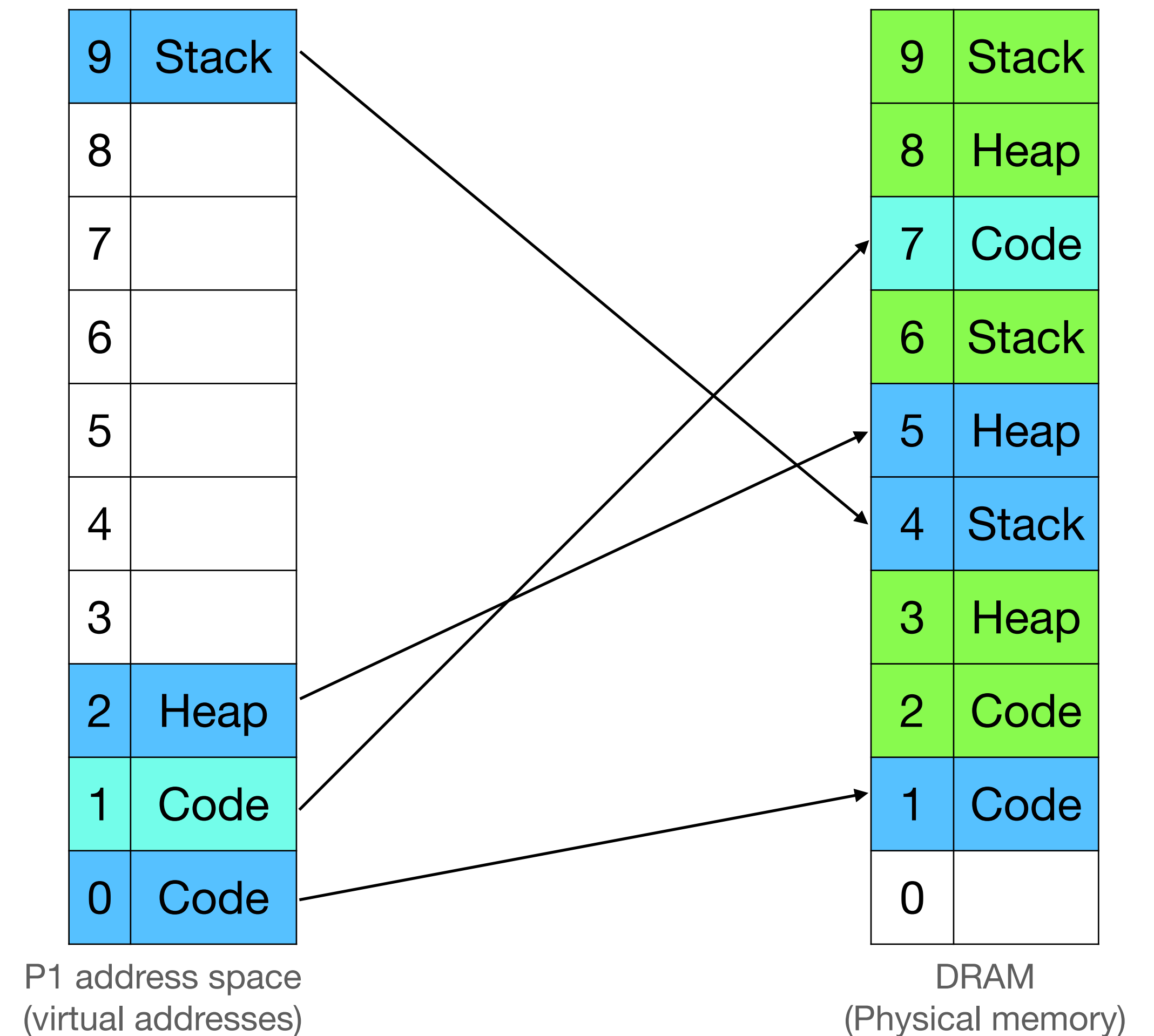


How to do page-based address translation? (2)

Page table bits

- Present bit: valid mapping
- Permission bits: read only, writeable, executable

Virtual page number	Physical page number	Present	Permission
9	4	Y	rw
8	x	N	
7	x	N	
6	x	N	
5	x	N	
4	x	N	
3	x	N	
2	5	Y	rw
1	7	Y	rx
0	1	Y	rx



Page table size

- Virtual addresses: 2^{32}
- Size of page = (4KB) = 2^{12}
- Number of page table entries = 2^{20}
- Number of pages in a 4GB DRAM = $2^{32}/2^{12} = 2^{20}$
- Size of each page table entry = 20 bits ~ 3 bytes
- Size of page table = $3 \cdot 2^{20} = 3\text{MB}$!
- 1000 processes => ~3GB memory!

Virtual page number	Physical page number
7	3
6	x
5	x
4	x
3	x
2	4
1	6
0	0

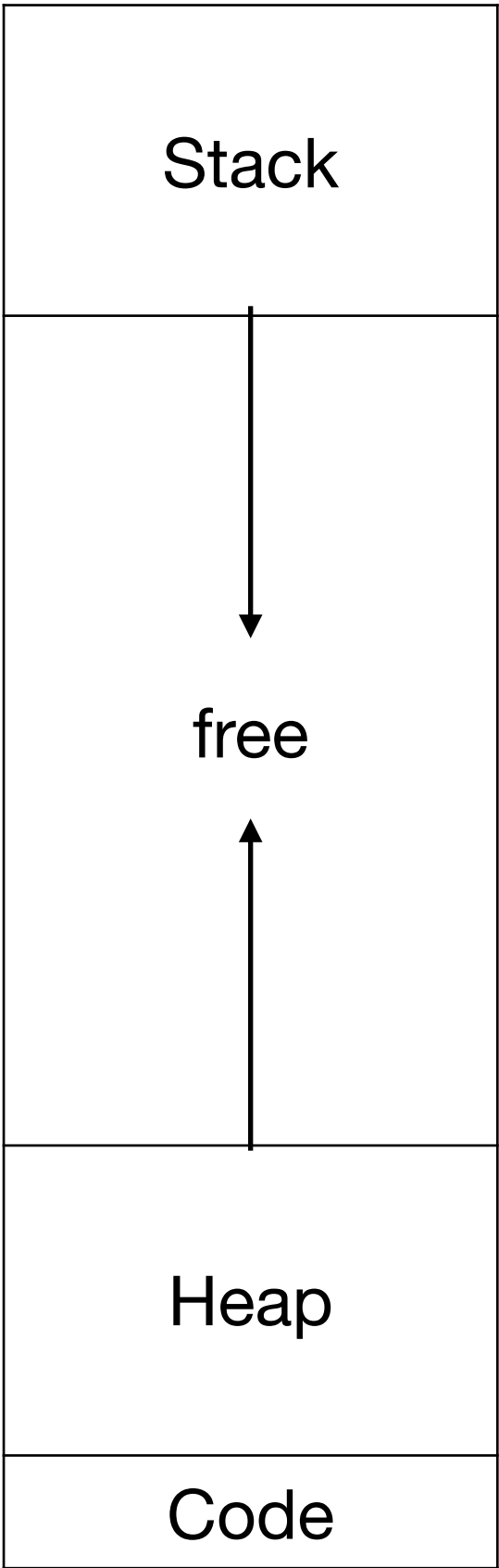
Reducing page table size

- Bigger pages
 - Virtual addresses: 2^{32}
 - Size of page = (**4MB**) = 2^{22}
 - Number of page table entries = 2^{10}
 - Number of pages in a 4GB DRAM = $2^{32}/2^{22} = 2^{10}$
 - Size of each page table entry = 10 bits ~ 2 bytes
 - Size of page table = $2 \cdot 2^{10} = 2\text{KB}$!
- Bigger pages increase internal fragmentation



Observation

- Lot of page table entries are invalid



Process 1 address space

Virtual page number	Physical page number
7	3
6	x
5	x
4	x
3	x
2	4
1	6
0	0

Multi-level page table

Virtual page number	Physical page number
7	3
6	x
5	x
4	x
3	x
2	4
1	6
0	0

PPN: 8

Virtual page number	Physical page number
6,7	9
4,5	x
2,3	11
0,1	10

PPN: 9

Virtual page number	Physical page number
7	3
6	x

PPN: 11

Virtual page number	Physical page number
3	x
2	4

PPN: 10

Virtual page number	Physical page number
1	6
0	0

- Page directory entries point to page table pages
- Unused portions of virtual address space is skipped!

Notebook analogy

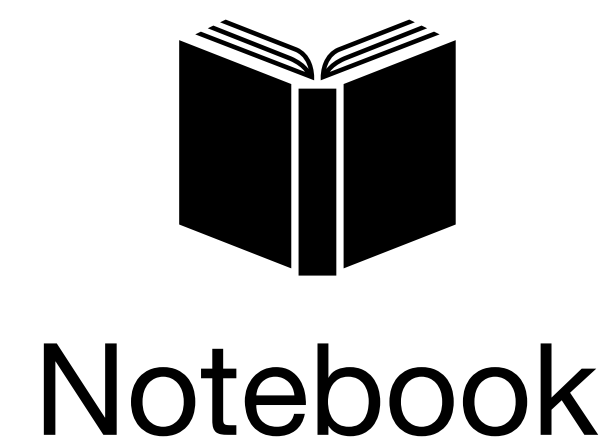
Paging



0	1	2	3	4	5
xv6	is	an	OS	for	x86



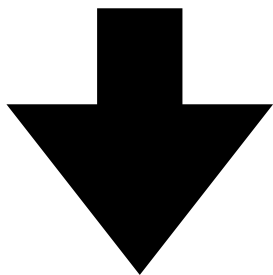
0	1	2	3
Write	an	SQL	query



0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
OS:1 DB:5	0:3,1:4, 2:2, 3:8 , 4:6,5:10	an	xv6	is	0:9,1:2, 2:7,3:12	for	SQL	OS	Write	x86		query		

Preparing for OS exam:

- Read second letter from 3rd page



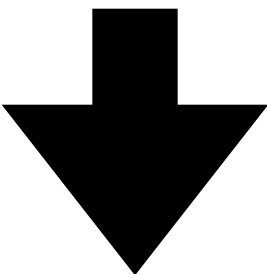
- Read second letter from 8th page

Notebook analogy

Page directories: call 4 pages a “section”

Preparing for OS exam:

- Read second letter from 3rd page in Section 0



- Read second letter from 8th page



OS

0	1	2	3	0	1	2	3
xv6	is	an	OS	for	x86		
Section 0				Section 1			



DB

0	1	2	3
Write	an	SQL	query
Section 0			



Notebook

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
OS:1 DB:5	0:11, 1:14	an	xv6	is	0:13	for	SQL	OS	Write	x86	0:3,1:4, 2:2, 3:8	query	0:9,1:2, 2:7,3:12	0:6, 1:10

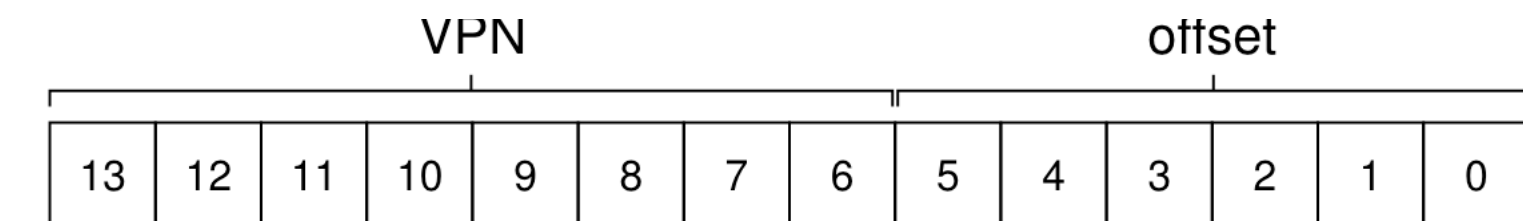
Address translation

Simple address space

- 16KB address space has 2^{14} addresses
- Each page has 64 ($= 2^6$) bytes
- Number of pages = 2^8
- First 8 bits are page number, last 6 bits are offset within the page

0000 0000	code
0000 0001	code
0000 0010	(free)
0000 0011	(free)
0000 0100	heap
0000 0101	heap
0000 0110	(free)
0000 0111	(free)
.....	... all free ...
1111 1100	(free)
1111 1101	(free)
1111 1110	stack
1111 1111	stack

Figure 20.4: A 16KB Address Space With 64-byte Pages

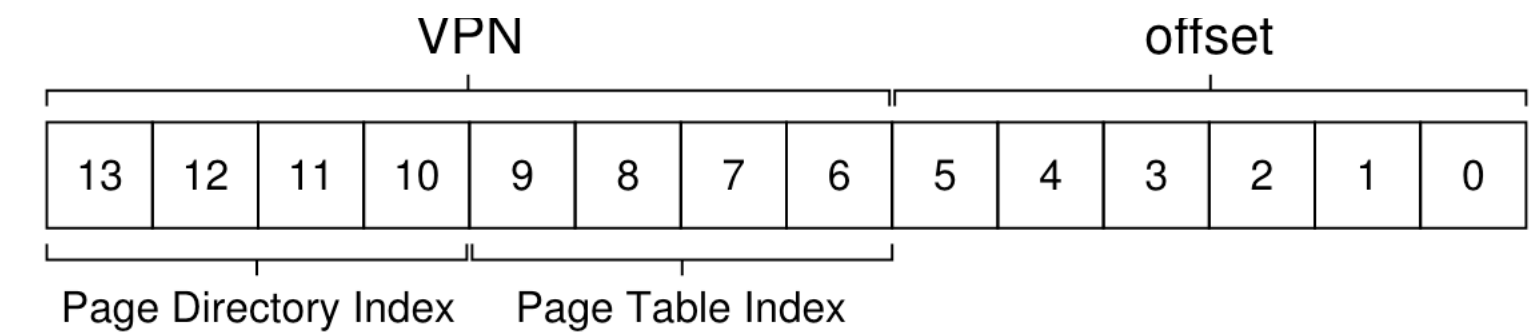


Simple address space

- Example: 0x3F81 (VA) => 0x0DC1 (PA)

Page dir idx 15				Page table idx 14				Offset					
1	1	1	1	1	1	1	0	0	0	0	0	0	1
3		F				8			1				

Physical page number 55								Offset					
0	0	1	1	0	1	1	1	0	0	0	0	0	1
0		D				C			1				



Page Directory		Page of PT (@PFN:100)			Page of PT (@PFN:101)		
PFN	valid?	PFN	valid	prot	PFN	valid	prot
100	1	10	1	r-x	—	0	—
—	0	23	1	r-x	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	80	1	rw-	—	0	—
—	0	59	1	rw-	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	55	1	rw-
101	1	—	0	—	45	1	rw-

Figure 20.5: A Page Directory, And Pieces Of Page Table

x86 segmentation and paging

- Segmentation:
 - Virtual address (logical address) => “linear address”
- Paging:
 - Linear address => Physical address

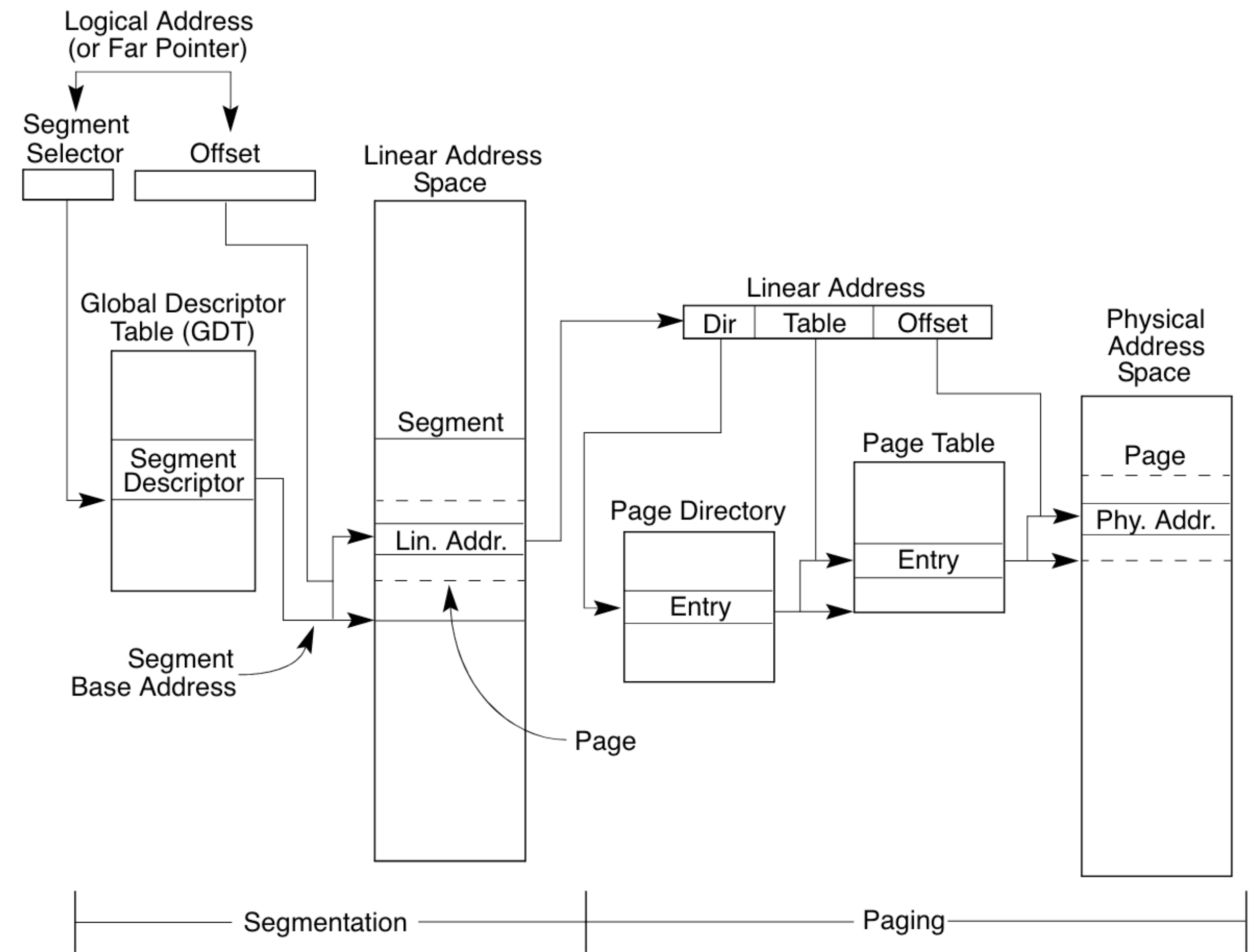
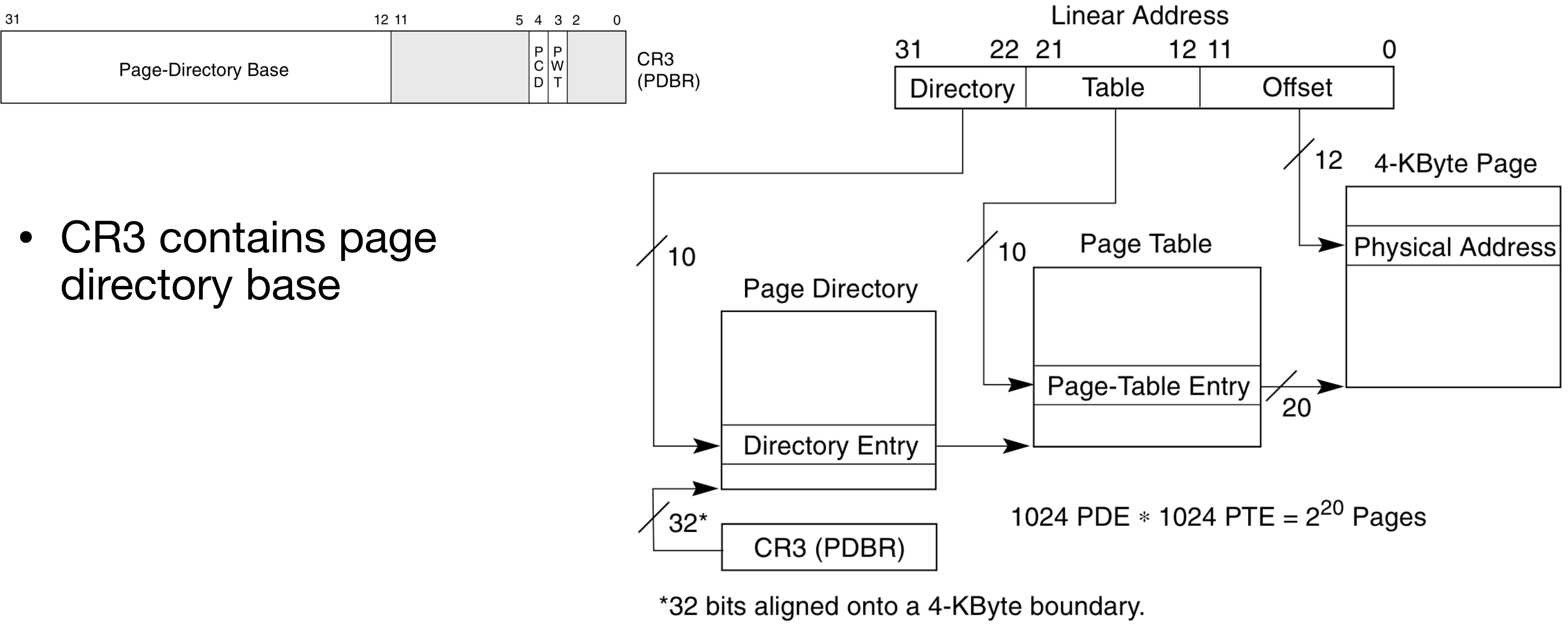


Figure 3-1. Segmentation and Paging

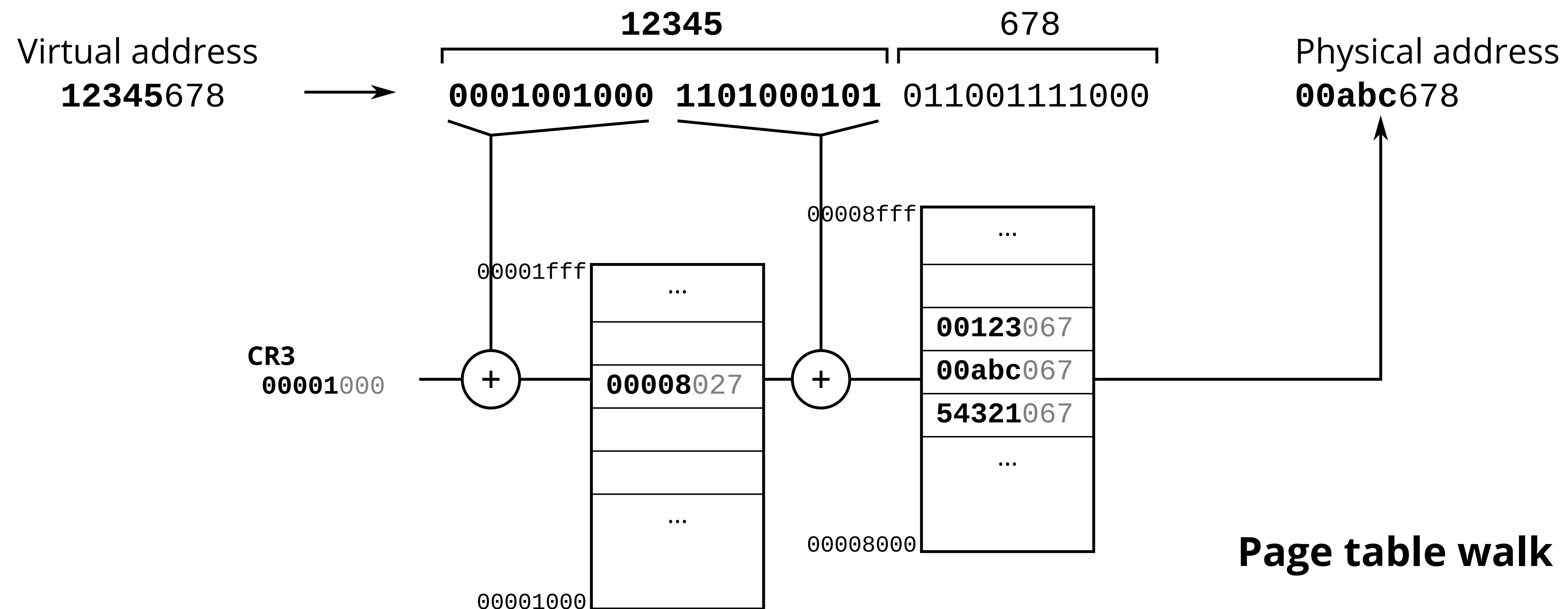
Address translation with paging on x86



- CR3 contains page directory base

Figure 3-12. Linear Address Translation (4-KByte Pages)

Address translation with paging on x86



x86 PTEs, PDEs

- 2^{20} 4KB pages in a 4GB DRAM
 - Page base address, page table base address are 20 bits
- Bits set by OS, used by hardware:
 - **Present**: It is a valid entry
 - **Read/write**: Can write if 1
 - **User/supervisor**: CPL=3 can access if 1
- Bits set by hardware, used/cleared by OS:
 - **Accessed**: Hardware accessed this page
 - **Dirty**: Hardware wrote to this page

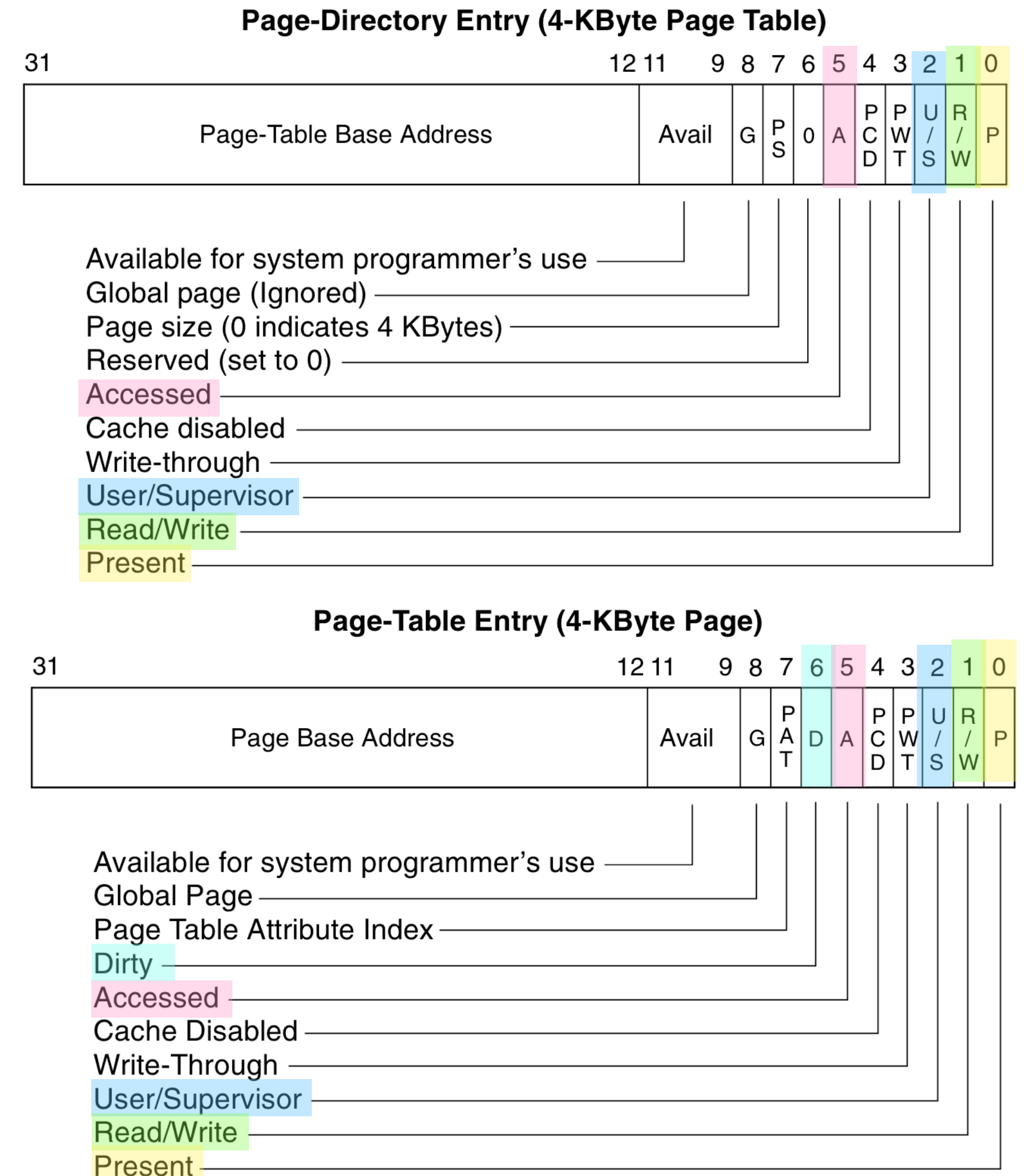
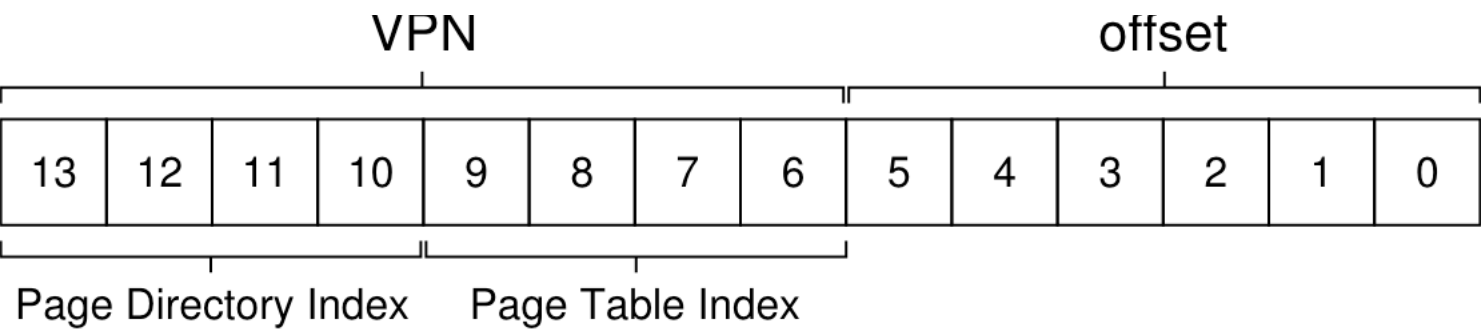


Figure 3-14. Format of Page-Directory and Page-Table Entries for 4-KByte Pages and 32-Bit Physical Addresses

Performance degradation!



- Accessing 1 memory location requires accessing 3 memory locations!

Page dir idx 15				Page table idx 14				Offset					
1	1	1	1	1	1	1	0	0	0	0	0	0	1
3		F				8			1				

Physical page number 55								Offset					
0	0	1	1	0	1	1	1	0	0	0	0	0	1
0		D				C			1				

Page Directory		Page of PT (@PFN:100)			Page of PT (@PFN:101)		
PFN	valid?	PFN	valid	prot	PFN	valid	prot
100	1	10	1	r-x	—	0	—
—	0	23	1	r-x	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	80	1	rw-	—	0	—
—	0	59	1	rw-	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	—	0	—
—	0	—	0	—	55	1	rw-
101	1	—	0	—	45	1	rw-

Figure 20.5: A Page Directory, And Pieces Of Page Table

Translation-lookaside buffer (TLB)

- First check page translation in TLB before walking the page table
- TLB hit: ~0.5-1 cycle
- TLB miss: ~10-100 cycles

Page dir idx 15				Page table idx 14				Offset					
1	1	1	1	1	1	1	0	0	0	0	0	0	1
3		F				8			1				

Physical page number 55								Offset					
0	0	1	1	0	1	1	1	0	0	0	0	0	1
0		D				C			1				

TLB

VPN								PPN						
1	1	1	1	1	1	1	0	0	0	1	1	0	1	1

Which programs will run faster?

Which programs will have lesser TLB misses?

```
int sum = 0;
for (i = 0; i < 10; i++) {
    sum += a[i];
}
```

	Offset				
	00	04	08	12	16
VPN = 00					
VPN = 01					
VPN = 02					
VPN = 03					
VPN = 04					
VPN = 05					
VPN = 06		a[0]	a[1]	a[2]	
VPN = 07	a[3]	a[4]	a[5]	a[6]	
VPN = 08	a[7]	a[8]	a[9]		
VPN = 09					
VPN = 10					
VPN = 11					
VPN = 12					
VPN = 13					
VPN = 14					
VPN = 15					

- High **spatial locality**: after the program accessed a memory location, it will access a nearby memory location
- High **temporal locality**: after the program accessed a memory location, it will soon access it again

Example bad programs

- Low spatial and temporal locality: most accesses lead to TLB miss

- Large hash table with random access

```
int get(int I)
```

```
    return a[I];
```

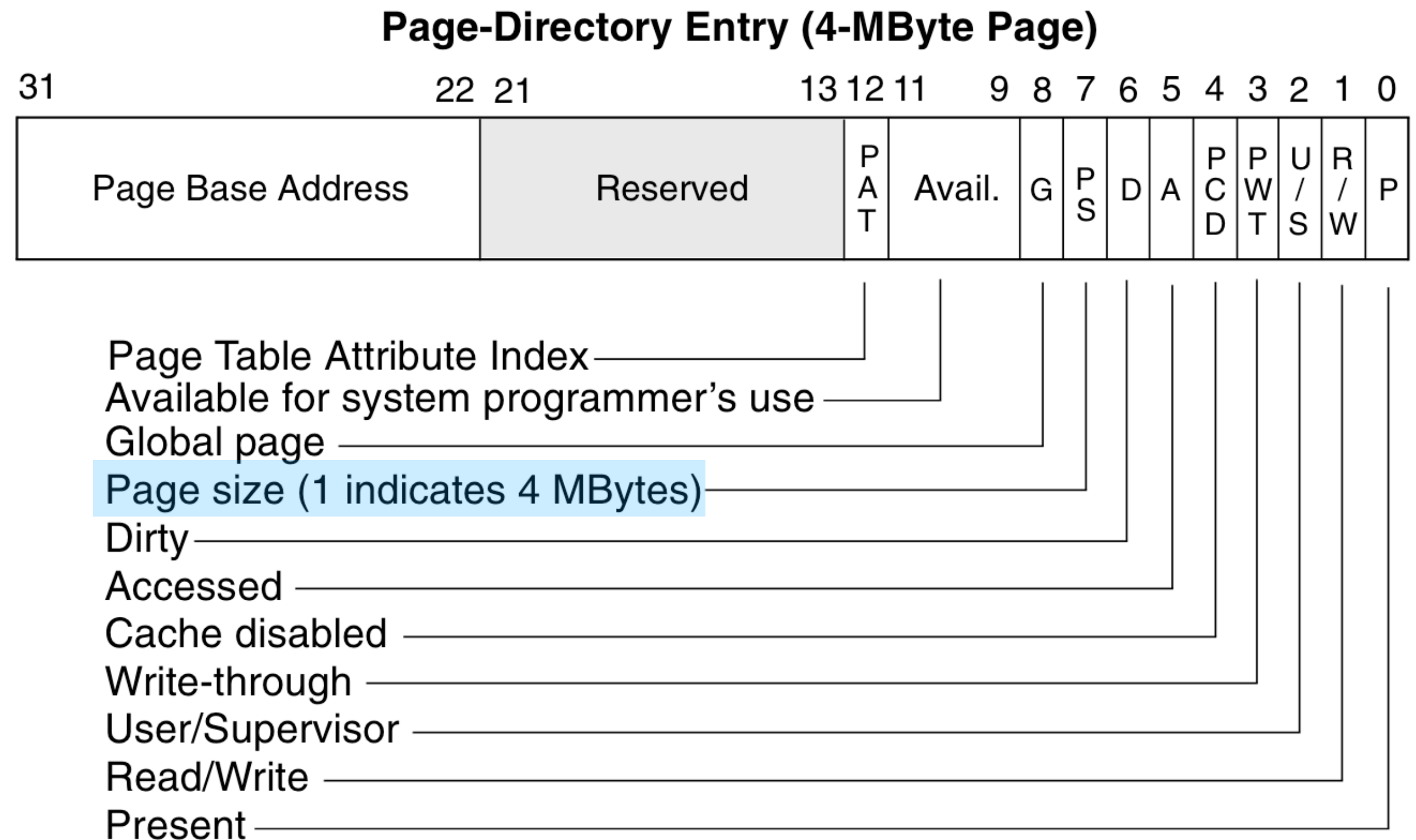
VPN=1	a[0]	a[1]	a[2]	a[3]
VPN=2	a[4]	a[5]	a[6]	a[7]
VPN=3	a[8]	a[9]	a[10]	a[11]
VPN=4	a[12]	a[13]	a[14]	a[15]
VPN=5	a[16]	a[17]	a[18]	a[19]

- Traversing large graphs: two neighbours can be on different pages
- Ok for small hash tables, small graphs.
 - Working set of the program: Amount of memory that it touches
- $\text{TLB reach} = (\text{number of TLB entries}) * (\text{page size})$
- $\text{Working set} > \text{TLB reach}$

Increasing TLB reach

- Larger TLBs
 - 32 -> 64 entries.
 - Larger caches => slower hits
- Larger pages
 - 1 TLB entry for 4MB of addresses (with a 4MB page)
 - 1024 TLB entries for 4MB of addresses (with 4KB pages)
- Allocating large pages on Linux:

```
posix_memalign(void **memptr, size_t alignment, size_t size);  
int madvise(void addr, size_t length, int advice= MADV_HUGEPAGE);
```



TLB on my x86-64 machine

```
$ cpuid
```

```
..  
(simple synth) = Intel Core (unknown type) (Kaby Lake / Coffee Lake) {Skylake}, 14nm
```

```
..  
cache and TLB information (2):
```

```
0x63: data TLB: 2M/4M pages, 4-way, 32 entries
```

```
data TLB: 1G pages, 4-way, 4 entries
```

```
0x03: data TLB: 4K pages, 4-way, 64 entries
```

```
0x76: instruction TLB: 2M/4M pages, fully, 8 entries
```

```
0xb5: instruction TLB: 4K, 8-way, 64 entries
```

```
0xc3: L2 TLB: 4K/2M pages, 6-way, 1536 entries
```

```
..
```

Separate TLBs for different page sizes

Separate TLBs for instruction and data

Much larger (slower) L2 TLB

Data TLB reach:

- $4\text{KB} * 64 = 256\text{KB}$
- $4\text{MB} * 32 = 128\text{MB}$
- $1\text{GB} * 4 = 4\text{GB}$

Instruction TLB reach:

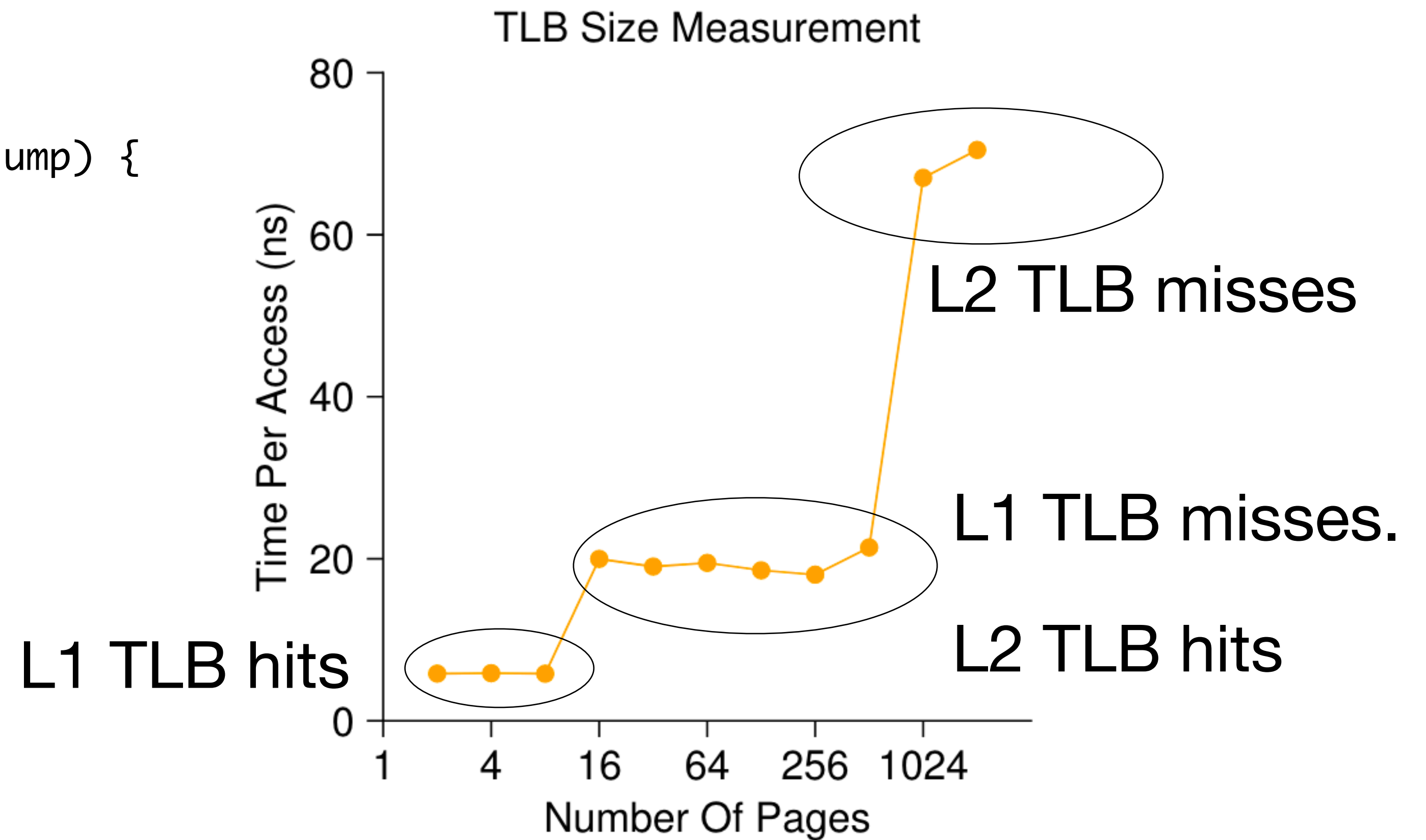
- $4\text{KB} * 64 = 256\text{KB}$
- $4\text{MB} * 8 = 32\text{MB}$

L2 TLB reach:

- $4\text{KB} * 1536 = 6\text{MB}$
- $2\text{MB} * 1536 = 3\text{GB}$

TLB size, multiple TLBs example

```
int jump = PAGE_SIZE / sizeof(int);
for (i = 0; i < NUMPAGES * jump; i += jump) {
    a[i] += 1;
}
```



TLB replacement policies

- Need to replace an entry after TLB is full. Which entry to replace to minimise TLB miss rate?
- Least recently used (LRU):
 - If an entry hasn't been used recently, it is unlikely to be used soon => assumes spatial and temporal locality of access.
 - Corner case behaviours: Program cycling over $N+1$ pages.
- Random replacement
 - Just pick an entry randomly

Context switching

- During context switch, OS changes CR3 register to change page table

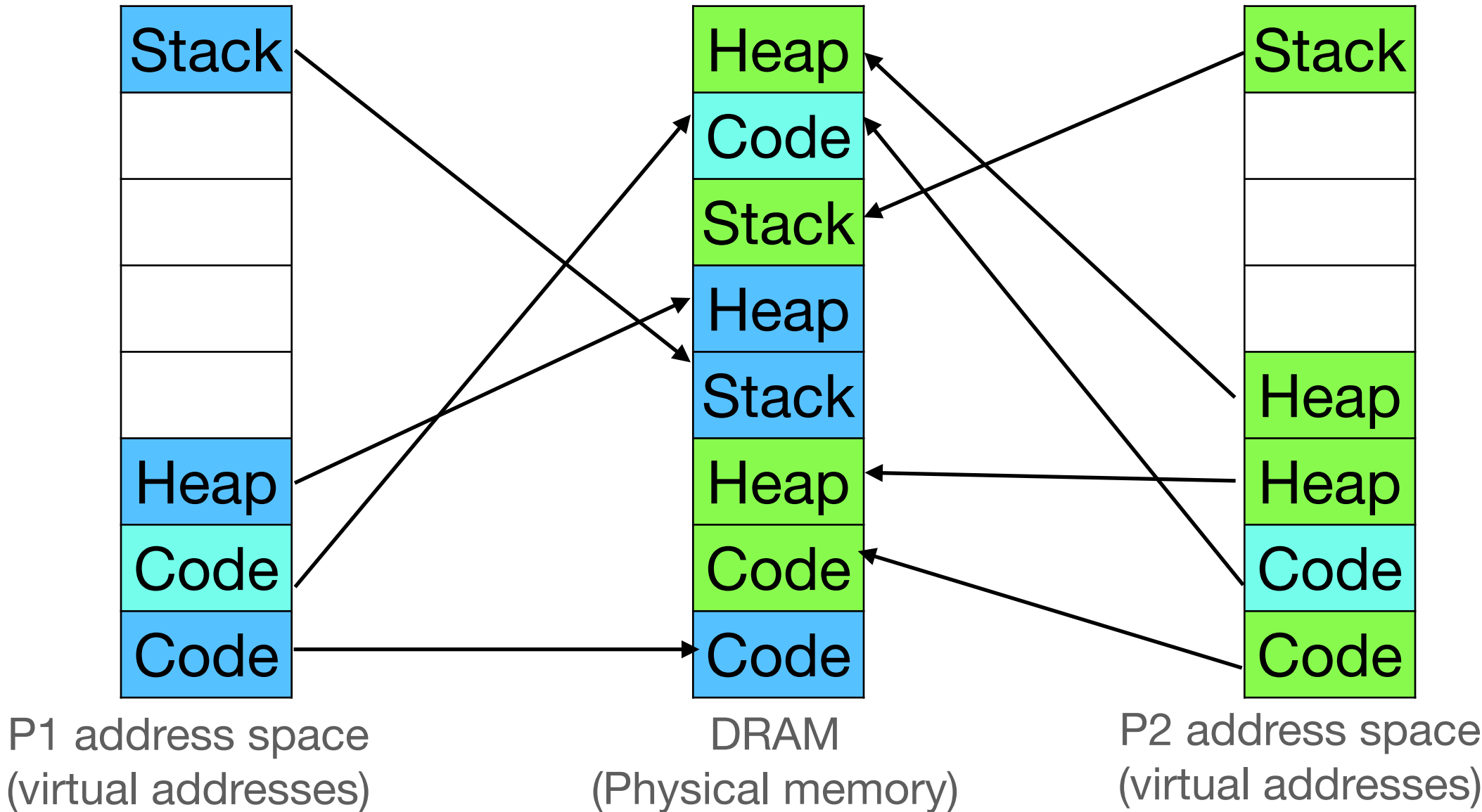
```
movl %eax, %cr3
```

- Privileged operation
- Marks each TLB entry invalid. Every memory access after context switch causes TLB miss!

Virtual page number	Physical page number
0	0
2	4
0	1

Virtual page number	Physical page number
7	3
6	x
5	x
4	x
3	x
2	4
1	6
0	0

Virtual page number	Physical page number
7	5
6	x
5	x
4	x
3	7
2	3
1	6
0	1



Tagged TLBs

- TLB entries are “tagged” with a *process context identifier*
- Additional bits in CR3 register tells hardware about the current PCID
- Upon context switch:
 - OS changes CR3: PCID, page directory base
 - Hardware *need not* invalidate TLB entries.
~5us
 - When process gets back control, some of its TLB entries might still be present!

TLB

Virtual page number	Physical page number	Process context identifier
0	0	P1
2	4	P1
0	1	P2

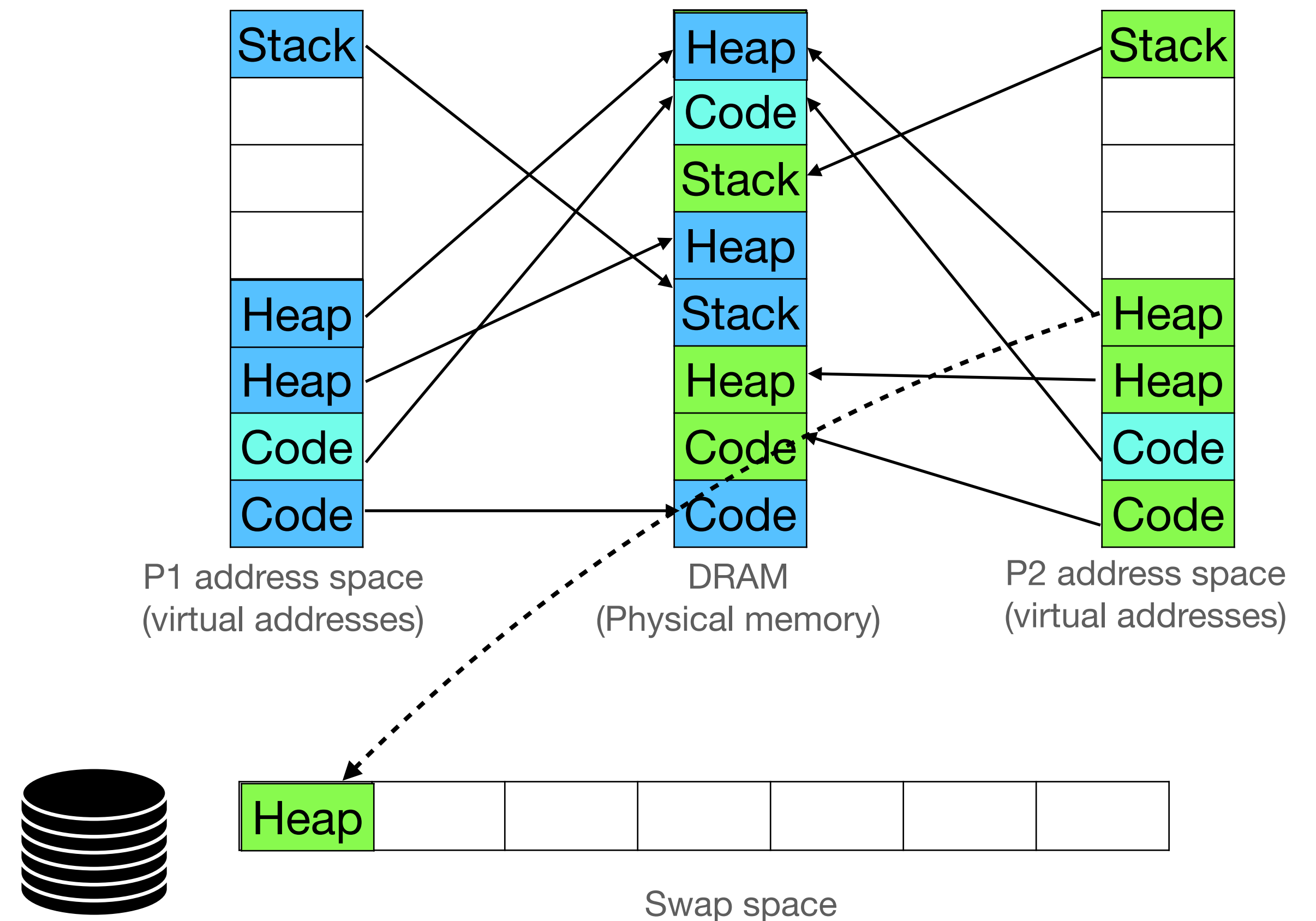
Demand Paging

OSTEP Ch 21-22

Demand paging

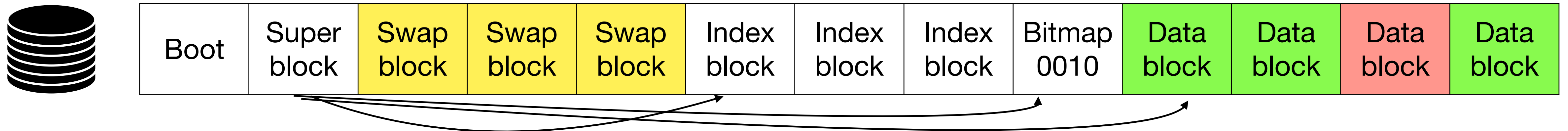
Providing illusion of large virtual address spaces

- Swap out a page to disk when physical memory is about to run out of space:
 - Mark page as not present in page table
 - Remember where page is swapped out on disk
- Add swapped out physical page to free list



Disk layout with swap space

- Reserved swap blocks, not touched by file system
- Swap space does not require crash consistency: anyways garbage after restart (all processes are dead)



Swapping out a page

- Find a free swap block on disk
- Copy page to the free block
- Run INVLPG instruction to remove page from TLB
- Mark not present, remember swap block number in PTE
- Add page to free list

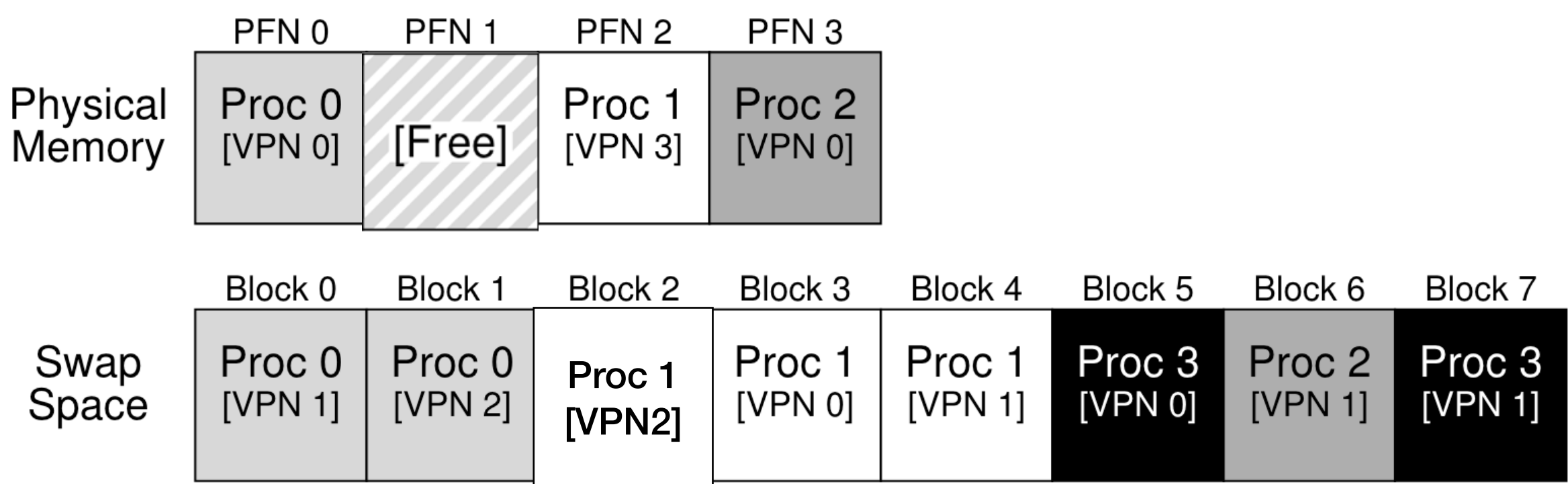


Figure 21.1: Physical Memory and Swap Space

Page table entry for Proc 1 [VPN2]

31												12 11				9	8	7	6	5	4	3	2	1	0
PFN=1, block # 2												Avail				G		P A T	D	A	P C D	P W T	U / S	R / W	1

Swapping in a page

- Hardware does page table walk to find that the page is not present. It raise page fault.
- OS handles page fault:
 - Copies page to a free physical page
 - Updates page table entry
- Hardware retries instruction. This time finds the page, adds entry to TLB, continues as normal

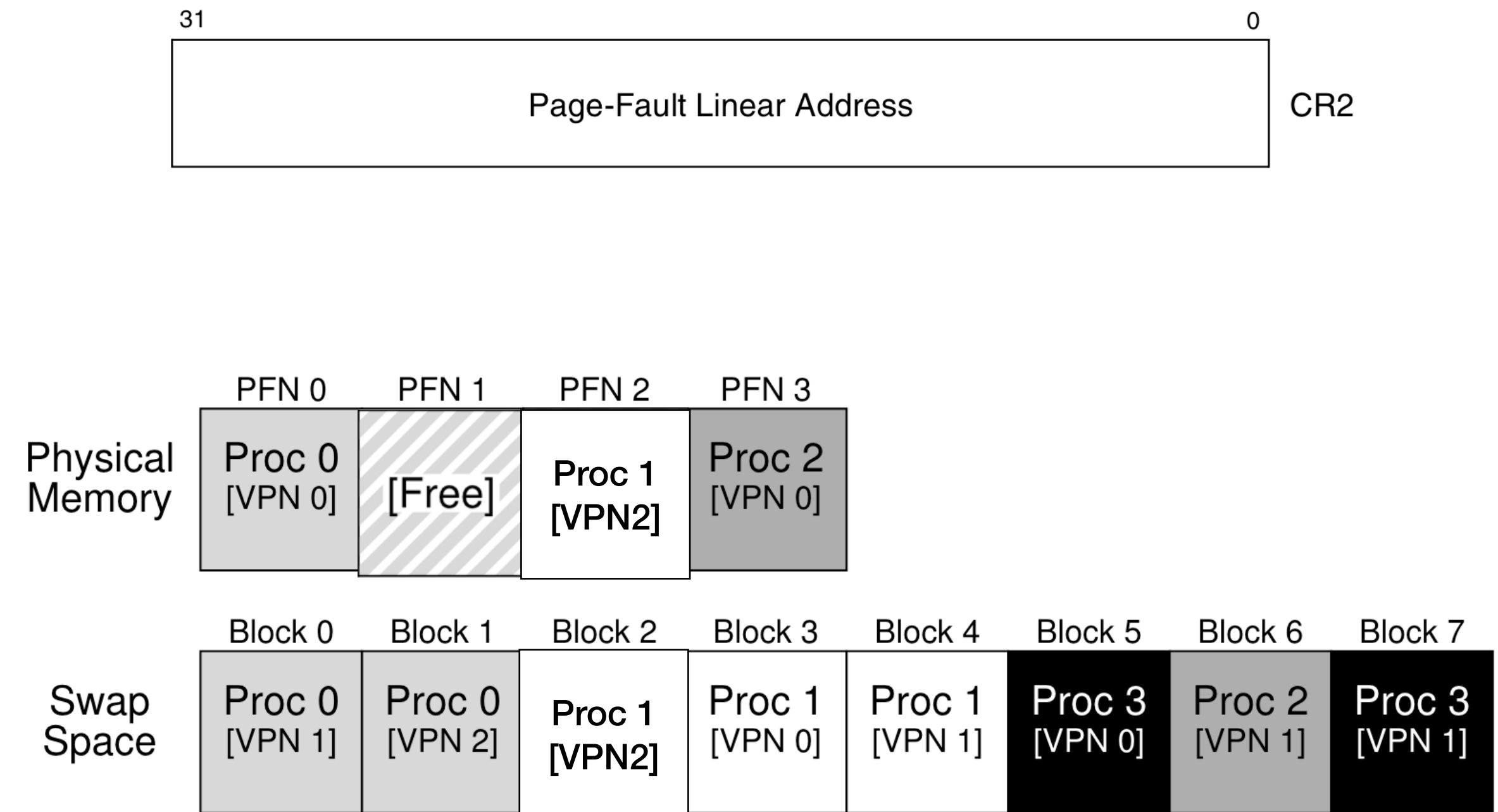


Figure 21.1: Physical Memory and Swap Space

Page table entry for Proc 1 [VPN2]



Page replacement policies

- When to evict pages?
- How many pages to evict?
- Which page to evict?

When to evict pages? How many pages to evict?

- Swap out one page when we completely run out of physical memory
 - What if OS itself needed a new page?
- Start swapping out before we completely run out
 - When there are less than N free pages left
- Swap out multiple pages in one shot until we have $M (> N)$ free pages left
 - Sends multiple disk writes in one shot reduces seek delay

Which page to evict?

- Goal: minimize number of swap ins/outs
- Belady's algorithm for optimal page replacement
 - Evict the page required furthest in the future
 - Optimal because all other pages will be required sooner
- Future is unknown!

Memory access sequence:

0, 1, 2, 0, 1, 3, 0, 3, 1, 2, 1

Access	Hit/Miss?	Evict	Resulting Cache State
0	Miss		0
1	Miss		0, 1
2	Miss		0, 1, 2
0	Hit		0, 1, 2
1	Hit		0, 1, 2
3	Miss	2	0, 1, 3
0	Hit		0, 1, 3
3	Hit		0, 1, 3
1	Hit		0, 1, 3
2	Miss	3	0, 1, 2
1	Hit		0, 1, 2

Figure 22.1: Tracing The Optimal Policy

FIFO

- Evict the page that came first to the cache
- OS appends the page to a queue when it swaps in a page (or when it allocates a new page)

Memory access sequence:

0, 1, 2, 0, 1, 3, 0, 3, 1, 2, 1

Access	Hit/Miss?	Evict	Resulting Cache State	
0	Miss		First-in→	0
1	Miss		First-in→	0, 1
2	Miss		First-in→	0, 1, 2
0	Hit		First-in→	0, 1, 2
1	Hit		First-in→	0, 1, 2
3	Miss	0	First-in→	1, 2, 3
0	Miss	1	First-in→	2, 3, 0
3	Hit		First-in→	2, 3, 0
1	Miss	2	First-in→	3, 0, 1
2	Miss	3	First-in→	0, 1, 2
1	Hit		First-in→	0, 1, 2

Figure 22.2: Tracing The FIFO Policy

Belady's anomaly

Bigger caches can have lower hit rates!

Access	Hit/miss	Resulting cache	Access	Hit/miss	Resulting cache
1	Miss	1	1	Miss	1
2	Miss	1, 2	2	Miss	1, 2
3	Miss	1, 2, 3	3	Miss	1, 2, 3
4	Miss	2, 3, 4	4	Miss	1, 2, 3, 4
1	Miss	3, 4, 1	1	Hit	1, 2, 3, 4
2	Miss	4, 1, 2	2	Hit	1, 2, 3, 4
5	Miss	1, 2, 5	5	Miss	2, 3, 4, 5
1	Hit	1, 2, 5	1	Miss	3, 4, 5, 1
2	Hit	1, 2, 5	2	Miss	4, 5, 1, 2
3	Miss	2, 5, 3	3	Miss	5, 1, 2, 3
4	Miss	5, 3, 4	4	Miss	1, 2, 3, 4
5	Hit	5, 3, 4	5	Miss	2, 3, 4, 5

FIFO does not follow “stack property”. Cache of size 4 may not contain elements in cache of size 3.

Fairness in page replacement

- Someone had **lots of pages**. I had **very little**. My page was evicted
- OS maintains “resident size” per process: 1, 7, 9
- First select a **victim process** with highest resident size, remove its pages



Least Recently Used (LRU)

- Most programs exhibit temporal locality:
 - If a page was accessed recently, it shall be accessed soon
- Keep list according to recency

Access	Hit/Miss?	Evict	Resulting Cache State	
0	Miss		LRU→	0
1	Miss		LRU→	0, 1
2	Miss		LRU→	0, 1, 2
0	Hit		LRU→	1, 2, 0
1	Hit		LRU→	2, 0, 1
3	Miss	2	LRU→	0, 1, 3
0	Hit		LRU→	1, 3, 0
3	Hit		LRU→	1, 0, 3
1	Hit		LRU→	0, 3, 1
2	Miss	0	LRU→	3, 1, 2
1	Hit		LRU→	3, 2, 1

Figure 22.5: Tracing The LRU Policy

Difficulty in implementing LRU

- In FIFO, list is updated by the OS when a new page is allocated, or when a page is swapped out
- In LRU, list needs to be updated at every access
 - OS is not running during page accesses :-/

Access	Hit/Miss?	Evict	Resulting Cache State	
0	Miss		LRU→	0
1	Miss		LRU→	0, 1
2	Miss		LRU→	0, 1, 2
0	Hit		LRU→	1, 2, 0
1	Hit		LRU→	2, 0, 1
3	Miss	2	LRU→	0, 1, 3
0	Hit		LRU→	1, 3, 0
3	Hit		LRU→	1, 0, 3
1	Hit		LRU→	0, 3, 1
2	Miss	0	LRU→	3, 1, 2
1	Hit		LRU→	3, 2, 1

Figure 22.5: Tracing The LRU Policy

Implementation options

- Option 1: Hardware maintains the LRU list
 - 4GB / 4KB $\sim 2^{20}$ pages
 - List size: 20 bits * 2^{20} pages $\sim 3\text{MB}$
 - List cannot be in CPU \Rightarrow must be in memory
 - Each memory access causes another set of memory accesses to update list

Implementation options

- Option 2: Hardware updates timestamp in page table entries. OS scans PTEs to find page with oldest timestamp
- PTEs live in memory. Updating timestamp again touches memory at every access. Defeats TLB.
- Lazily update timestamp:
 - when mapping is brought to TLB
 - when mapping is evicted from TLB
 - Once in a while
- Need many more bits in PTE to store timestamp
- Victim process has highest resident size. Scanning (worst case 2^{20}) PTEs will be slow

Page table

Physical page number	Access Timestamp
10	8:11am
11	7:05am
15	7:00am
12	8:21am

TLB

Virtual page number	Physical page number
0	10
1	11
2	15
3	12

Approximating LRU

- Give up on finding *least* recently used. OK to evict a less recently used
- Hardware just lazily sets 1 access bit

Page table

Physical page number	Access Timestamp
10	8:11am
11	7:05am
15	7:00am
12	8:21am

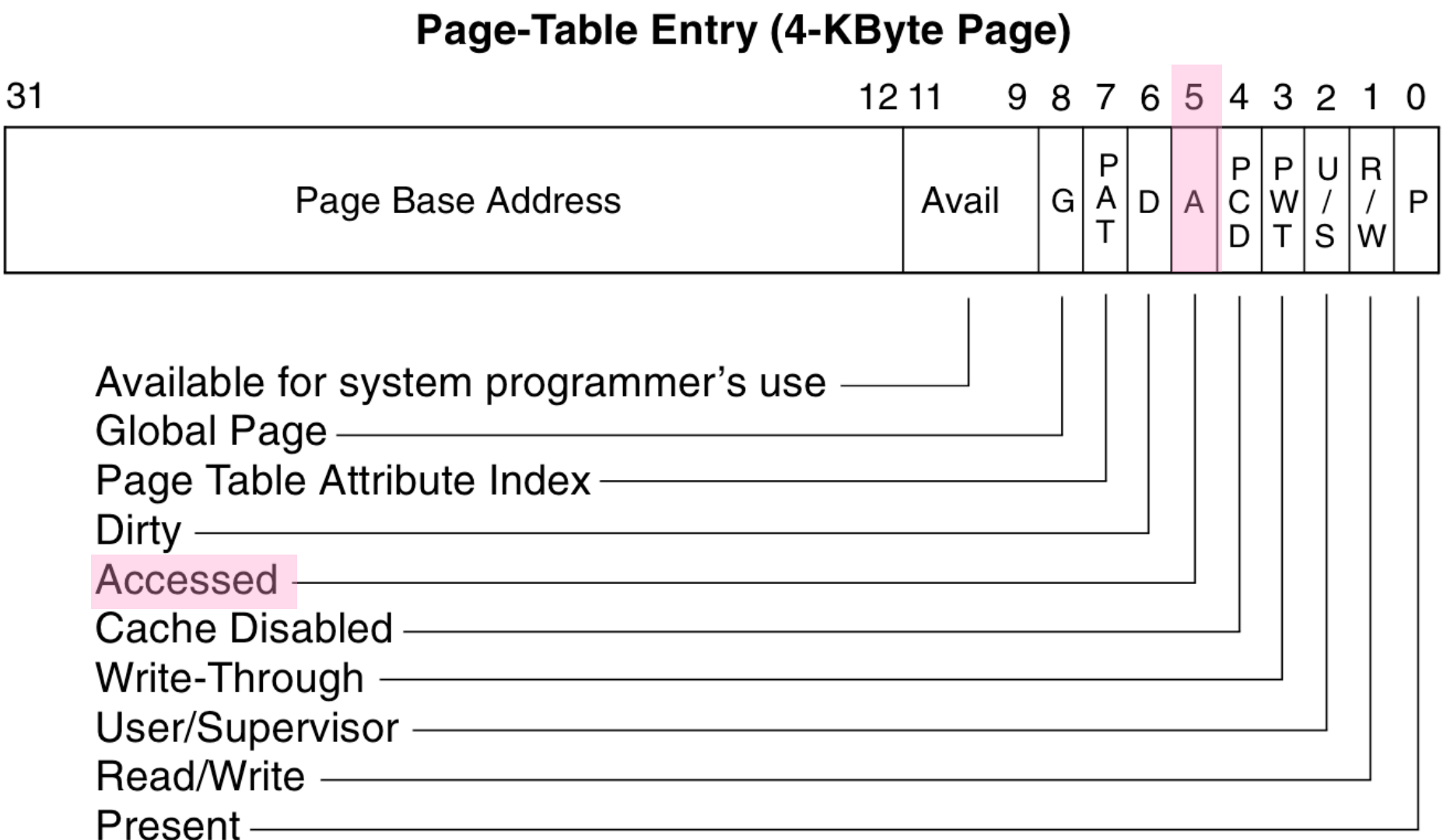


Figure 3-14. Format of Page-Directory and Page-Table Entries for 4-KByte Pages and 32-Bit Physical Addresses

Clock algorithm

- OS clears access bit. Evicts page with access bit = 0
- Hardware sets access bit to 1
- Evicted page was “not recently used”

Victim proc's pages

	AB
	0
	0
	0
	0
OS →	1
	0
	1
	0

Clock algorithm (2)

- Optimisation:
 - We should prefer evicting page that has not changed since it was brought back from disk
 - Such pages can just be deleted without doing a copy
- OS clears dirty bit when it brings a page to memory from disk
- Evict dirty pages iff not able to find a clean page

OS →

	DB	AB
	0	0
	0	0
	0	0
	1	0
	0	1
	1	0
	1	1
	0	0