# Virtual Memory

CS 351: Systems Programming
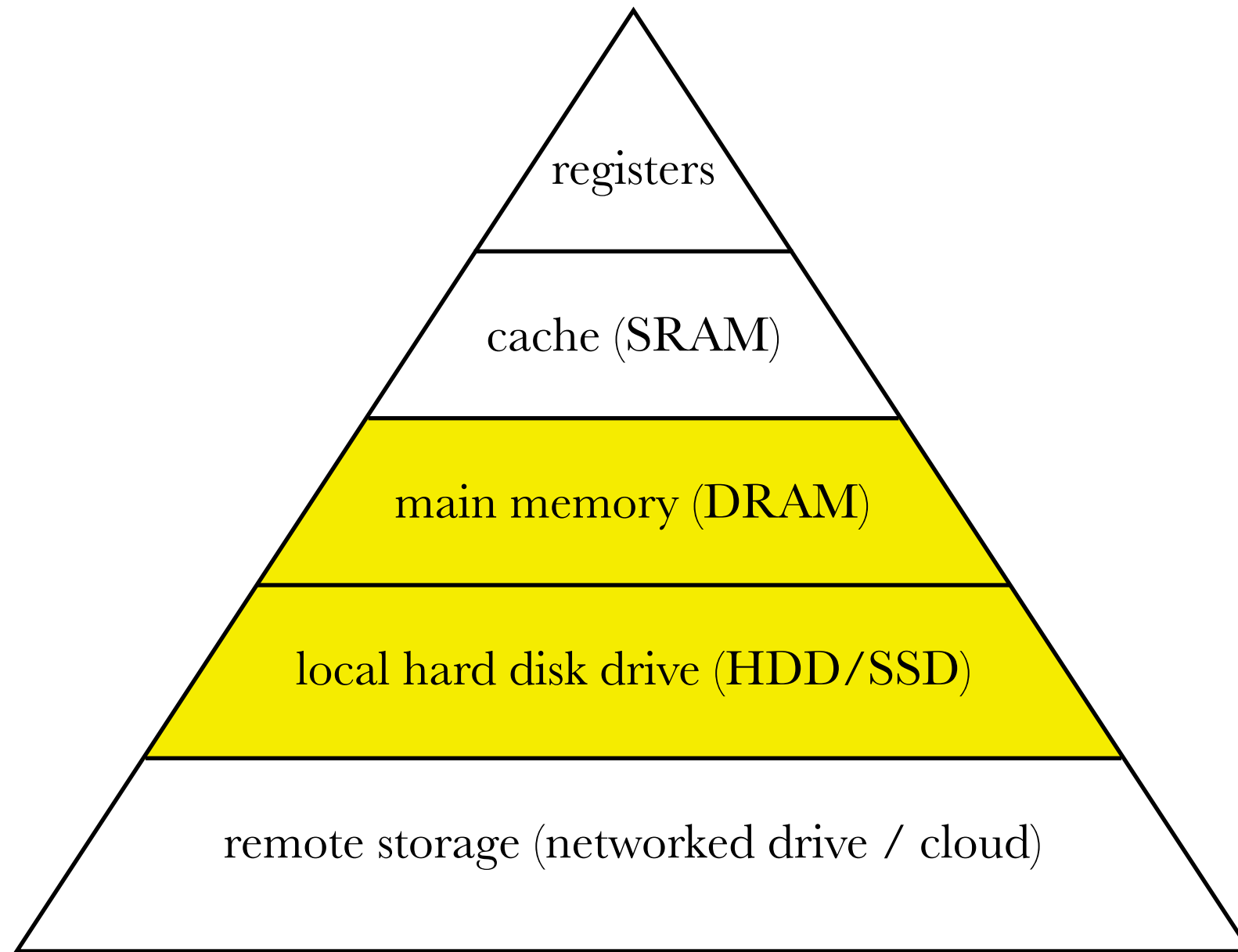Michael Lee <lee@iit.edu>

previously: SRAM ⇔ DRAM

next: DRAM ⇔ HDD, SSD, etc.

i.e., memory as a "cache" for disk

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

main goals:

1. maximize memory *throughput*

2. maximize memory *utilization*

3. provide *address space consistency* & *memory protection* to processes

*throughput* = # bytes per second

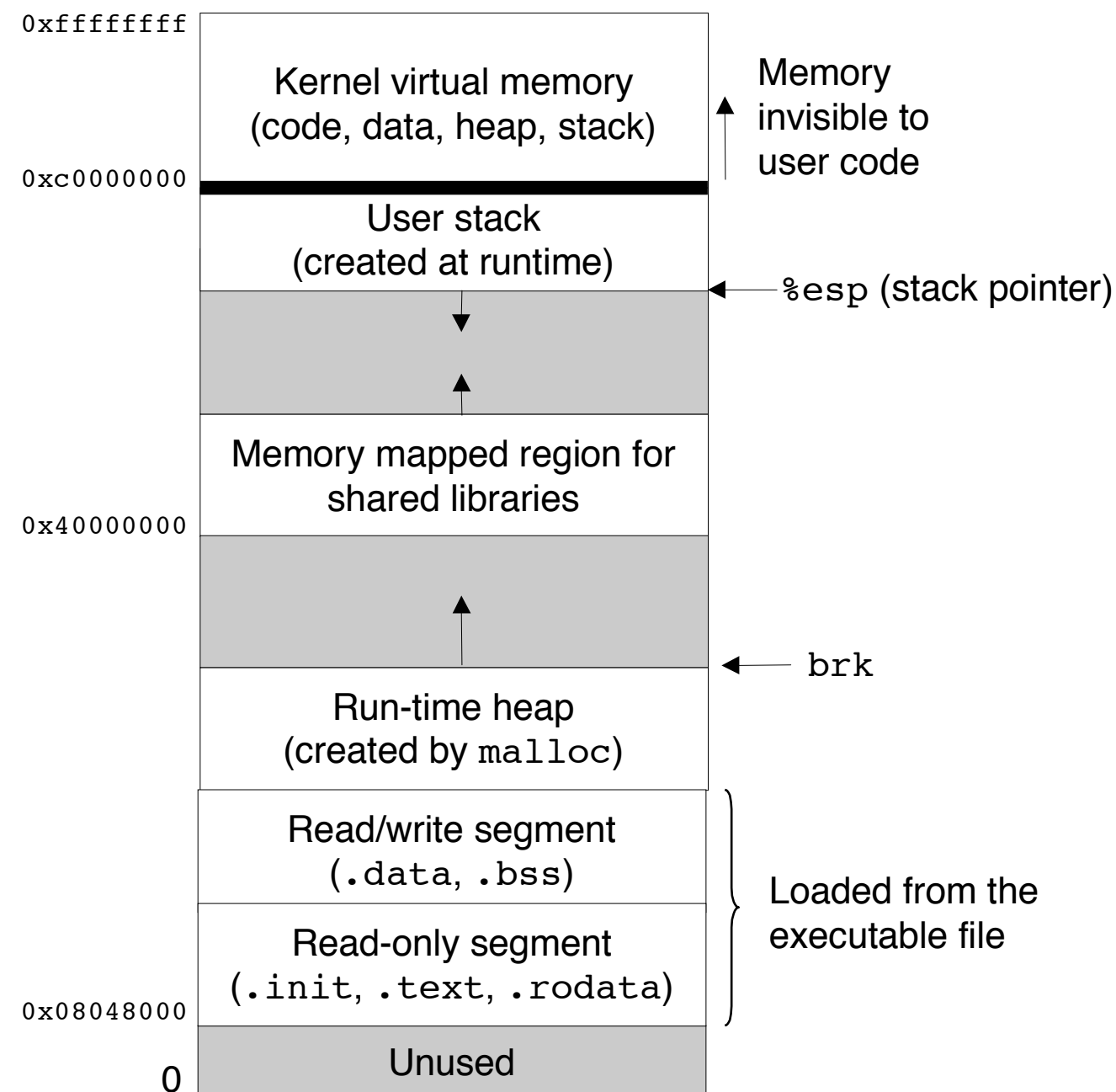- depends on access latencies (DRAM, HDD, etc.) and "hit rate"

*utilization* = fraction of allocated memory that contains "user" data (aka *payload*)

- vs. metadata and other overhead required for memory allocation

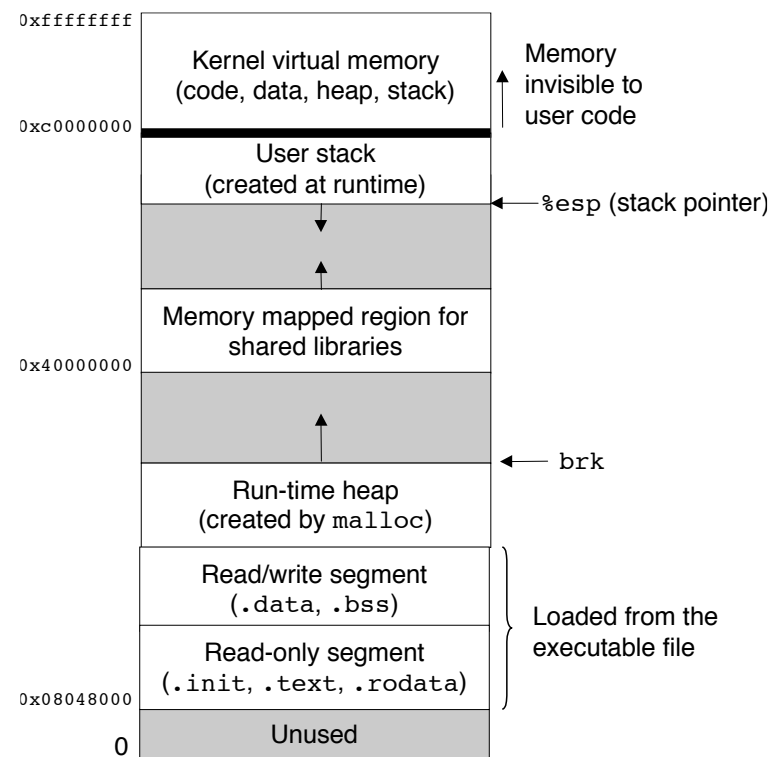*address space consistency* → provide a uniform "view" of memory to each process
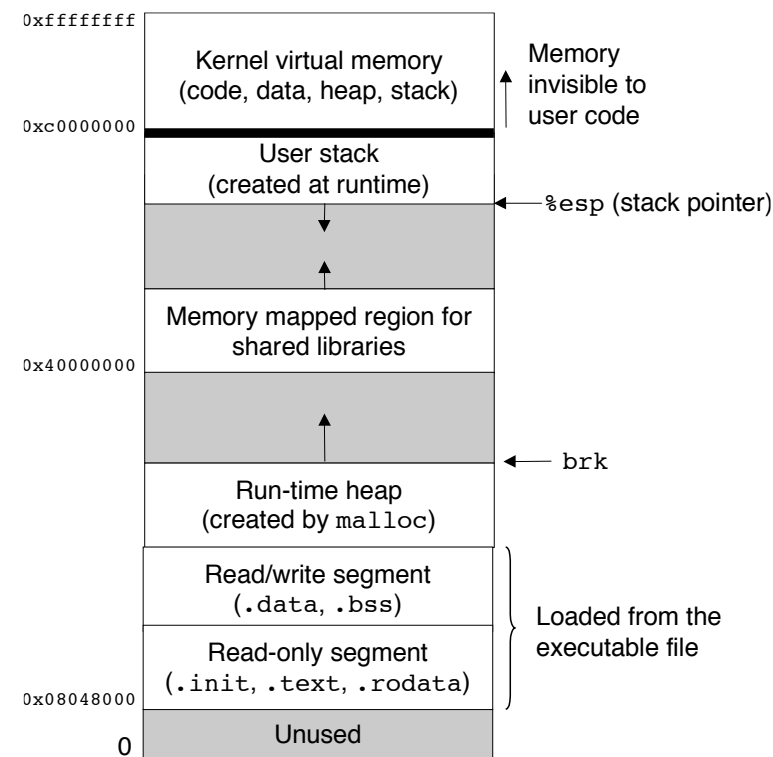
*address space consistency* → provide a uniform "view" of memory to each process

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

*memory protection* → prevent processes from directly accessing each other's address space
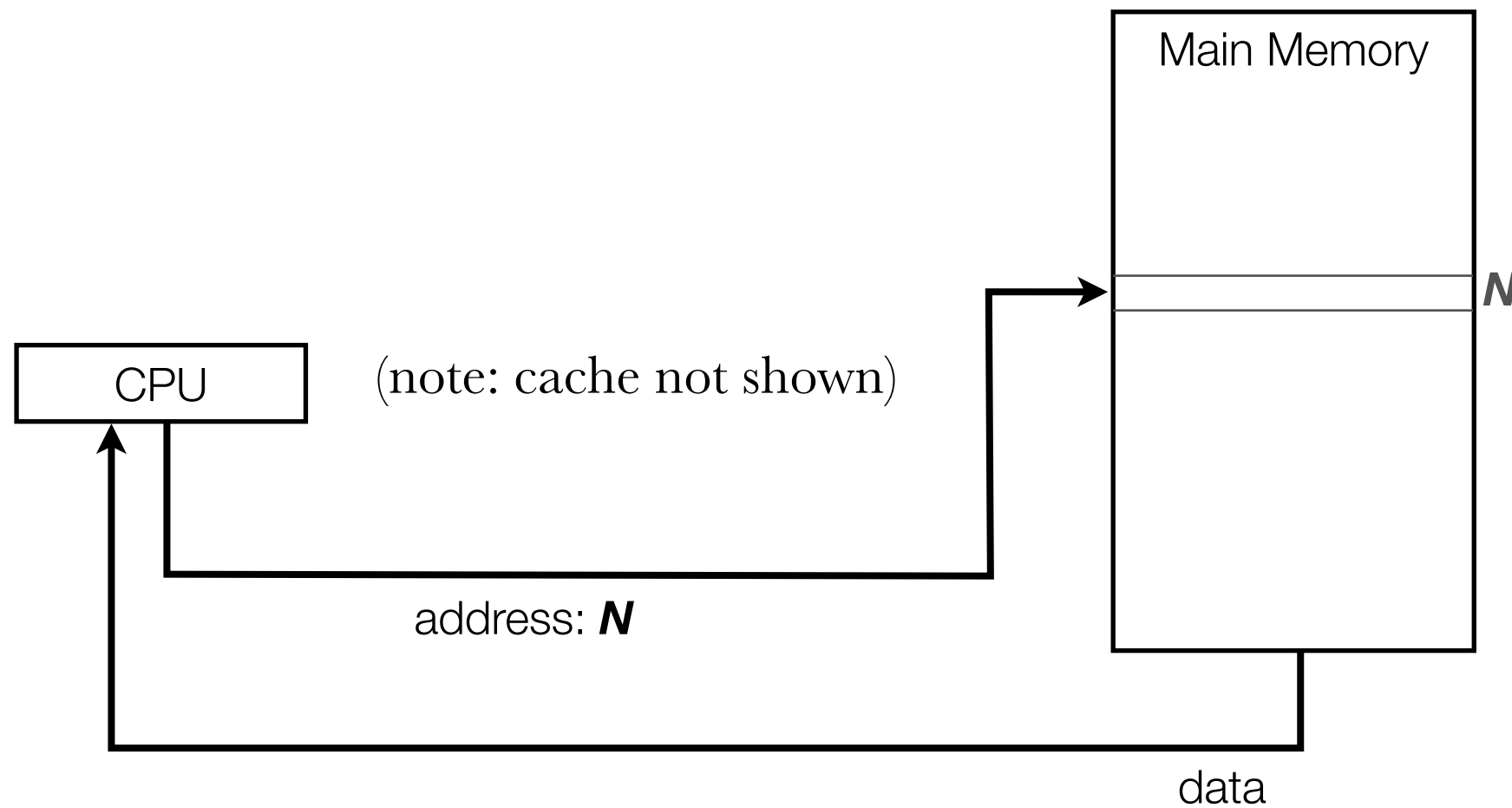
$P_0$     $P_1$     $P_2$

*memory protection* → prevent processes from
directly accessing each other's address space

i.e., every process should be provided with a managed, *virtualized* address space

"memory addresses": what are they, really?

Main Memory

*N*

CPU

(note: cache not shown)

address: ***N***

data

"physical" address: (byte) index into DRAM

```
int glob = 0xDEADBEEE;

main() {
    fork();
    glob += 1;
}
```

**(gdb) set detach-on-fork off**
**(gdb) break main**
Breakpoint 1 at 0x400508: file memtest.c, line 7.
**(gdb) run**
Breakpoint 1, main () at memtest.c:7
7        fork();
**(gdb) next**
[New process 7450]
8        glob += 1;
**(gdb) print &glob**
$1 = (int *) 0x6008d4
**(gdb) next**
9      }
**(gdb) print /x glob**
$2 = 0xdeadbeef
**(gdb) inferior 2**
[Switching to inferior 2 [process 7450]
#0  0x000000310acac49d in __libc_fork ()
131      pid = ARCH_FORK ();
**(gdb) finish**
Run till exit from #0  in __libc_fork ()
8        glob += 1;
**(gdb) print /x glob**
$4 = 0xdeadbeee
**(gdb) print &glob**
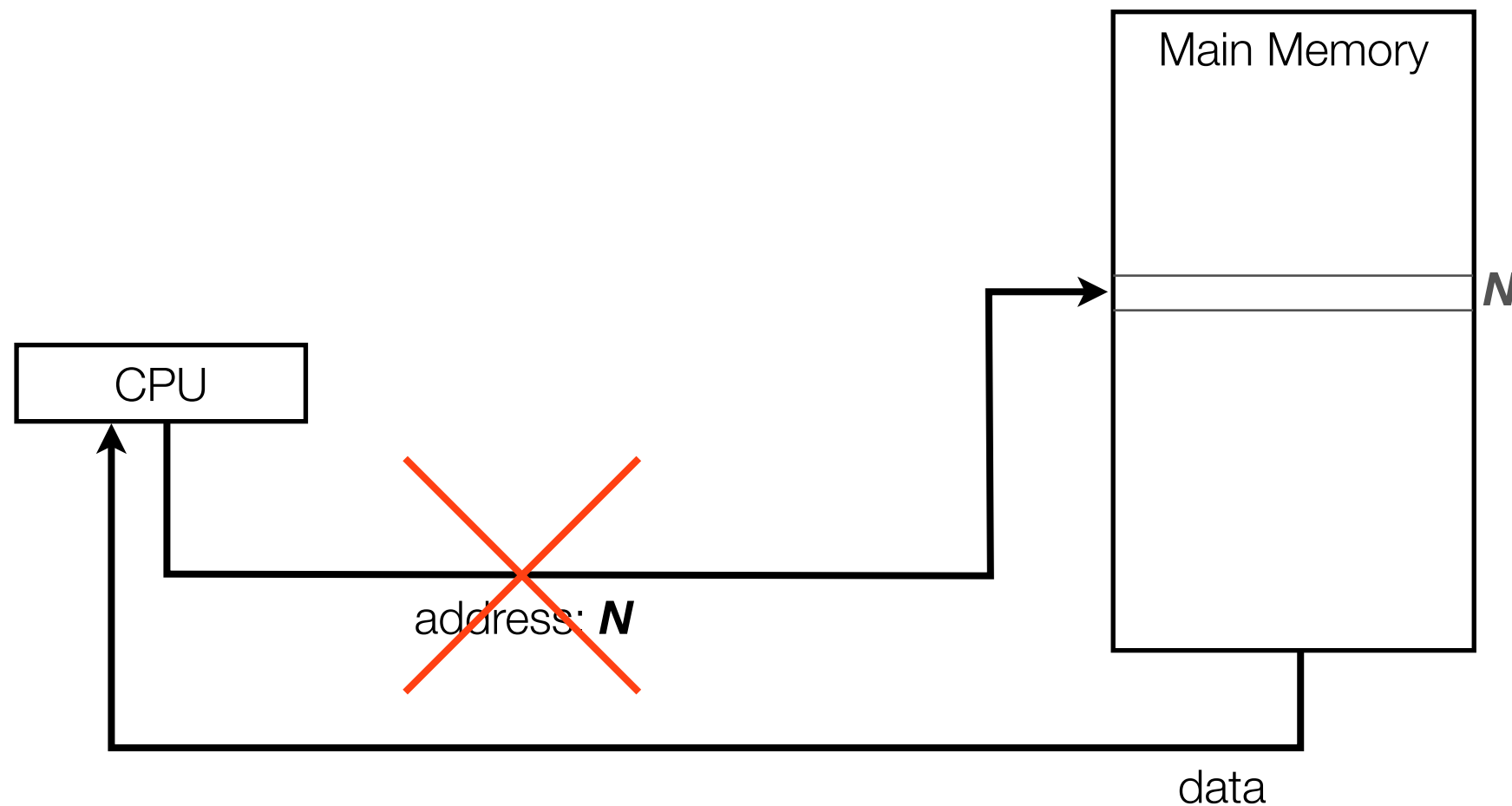$5 = (int *) 0x6008d4

*parent*

*child*

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY
```

instructions executed by the CPU *do not* refer directly to *physical* addresses!

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

processes reference *virtual* addresses,

the CPU relays virtual address requests to the *memory management unit* (MMU),

which are *translated* to physical addresses

Main Memory

**MMU**

physical address

CPU

**virtual address**

address translation unit

disk address

(note: cache not shown)

"swap" space

essential problem: translate request for a
virtual address → physical address

… this must be **FAST**, as *every* memory
access from the CPU must be translated

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

both hardware/software are involved:

- MMU (hw) handles simple and fast operations (e.g., table lookups)

- Kernel (sw) handles complex tasks (e.g., eviction policy)

# §Virtual Memory Implementations

# 1. simple relocation

# 1. simple relocation



- per-process relocation address is loaded by kernel on every context switch

# 1. simple relocation



- problem: processes may easily overextend their bounds and trample on each other

# 1. simple relocation



- incorporate a *limit* register to provide memory protection

# 1. simple relocation

Main Memory

MMU

limit reg.

*L*

CPU

relocation reg.

*B*

VA: *N*

PA: *N+B*

*B+L*

*N*

*B*

process
sandbox

data

assert (0 ≤ *N* ≤ *L*)

**-** assertion failure triggers a fault, which
summons kernel (which signals process)

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

pros:

- simple & fast!

- provides protection

but: available memory for mapping depends on value of base address

i.e., address spaces are *not consistent*!

also: all of a process *below the address limit* must be loaded in memory

i.e., memory may be *vastly under-utilized*

# 2. segmentation

- partition virtual address space into *multiple logical segments*

- individually map them onto physical memory with relocation registers

## Segmented Virtual Address Space

0  Seg #3: Stack

0  Seg #2: Heap

0  Seg #1: Data

0  Seg #0: Code

## MMU
### Segment Table

|   | Base | Limit |
|---|------|-------|
| 0 | $B_0$ | $L_0$ |
| 1 | $B_1$ | $L_1$ |
| 2 | $B_2$ | $L_2$ |
| 3 | $B_3$ | $L_3$ |

## Main Memory

$B_2+L_2$

$B_2$

$B_3+L_3$

$B_3$

$B_0+L_0$

$B_0$

$B_1+L_1$

$B_1$

# virtual address has form *seg#:offset*

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

MMU

Segment Table

| | Base | Limit |
|---|------|-------|
| 0 | $B_0$ | $L_0$ |
| 1 | $B_1$ | $L_1$ |
| 2 | $B_2$ | $L_2$ |
| 3 | $B_3$ | $L_3$ |

Main Memory

$B_2+L_2$

$B_2$

$B_3+L_3$

$B_3$

$B_0+L_0$

$B_0$

$B_1+L_1$

$B_1$

CPU

VA: *seg#:offset*

PA: *offset + $B_2$*

$\oplus$

assert (*offset* ≤ $L_2$)

data

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

Segment Table

| | Base | Limit |
|---|---|---|
| 0 | $B_0$ | $L_0$ |
| 1 | $B_1$ | $L_1$ |
| 2 | $B_2$ | $L_2$ |
| 3 | $B_3$ | $L_3$ |

\- implemented as MMU registers

\- part of kernel-maintained, per-process metadata (aka "process control block")

\- re-populated on each context switch

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

pros:

- still very *fast*

- translation = register access & addition

- memory *protection* via limits

- segmented addresses improve *consistency*

- variable segment sizes → *memory fragmentation*

- fragmentation potentially *lowers utilization*

- can fix through compaction, but expensive!

# 3. paging

- partition virtual and physical address spaces into *uniformly sized* **pages**

- virtual pages map onto physical pages

stack

heap

data

code

physical memory

- minimum mapping granularity = page

- not all of a given segment need be mapped

modified mapping problem:

- a virtual address is broken down into
  *virtual page number* & *page offset*

- determine which physical page (if any)
  a given virtual page is loaded into

- if physical page is found, use page
  offset to access data

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

# Given page size = $2^p$ bytes

$p$

**VA:**

| virtual page number | virtual page offset |
|---|---|

$p$

**PA:**

| physical page number | physical page offset |
|---|---|

VA: | *virtual page number* | *virtual page offset* |

address
translation

PA: | *physical page number* | *physical page offset* |

VA:

$n$

*virtual page number* | *virtual page offset*

*translation structure:* **page table**

*valid*  *PPN*

*index*

$2^n$ entries

*if invalid, page
is not mapped*

PA:

*physical page number* | *physical page offset*

page table entries (PTEs) typically contain additional metadata, e.g.:

- dirty (modified) bit

- access bits (shared or kernel-owned pages may be read-only or inaccessible)

e.g., 32-bit virtual address,
4KB ($2^{12}$) page size,
4-byte PTE size;

- size of page table?

e.g., 32-bit virtual address,
    4KB ($2^{12}$) pages,
    4-byte PTEs;

- # pages = $2^{32} \div 2^{12} = 2^{20} = 1M$

- page table size = 1M × 4 bytes = **4MB**

4MB is much too large to fit in the MMU — insufficient registers and SRAM!

Page table resides in **main memory**

The translation process (aka *page table walk*) is performed by hardware (MMU).

The kernel must initially populate, then continue to manage a process's page table

The kernel also populates a *page table base register* on context switches

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

# translation: *hit*



CPU

❶ VA: **N**

Address
Translator
(part of MMU)

❷ *page table
walk*

❸ PA: **N′**

Main
Memory

Page
Table

❹ *data*

translation: *miss*

kernel decides where to place page, and what to evict (if memory is full)

- e.g., using LRU replacement policy

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

this system enables **on-demand paging**

i.e., an active process need only be partly in memory (load rest from disk dynamically)

but if working set (of active processes)
exceeds available memory, we may have
**swap thrashing**

integration with caches?

Q: do caches use physical or virtual addresses for lookups?

# Virtual address based Cache

# Physical address based Cache

Q: do caches use physical or virtual addresses for lookups?

A: caches typically use *physical* addresses

**%\*@$&#!!!**

saved by hardware:

the *Translation Lookaside Buffer* (TLB) — a cache used solely for VPN→PPN lookups

# TLB + Page table

(exercise for reader: revise earlier translation diagrams!)

virtual address

n-1                                              p p-1                              0

| virtual page number (VPN) | page offset |
|---|---|

TLB

valid        tag        physical page number (PPN)

= 

**TLB Hit** ◄

physical address

byte offset

Cache

valid        tag                data

=

**Cache Hit** ◄

**Data**

TLB mappings are *process specific* — requires flush & reload on context switch

- some architectures store PID (aka "virtual space" ID) in TLB

Familiar caching problem:

- TLB caches a few thousand mappings

- vs. *millions* of virtual pages per process!

we can improve TLB hit rate by reducing the number of pages …

by increasing the size of each page

compute # pages for 32-bit memory for:

- 1KB, 512KB, 4MB pages

  - $2^{32} \div 2^{10} = 2^{22}$  = 4M pages

  - $2^{32} \div 2^{19} = 2^{13}$  = 8K pages

  - $2^{32} \div 2^{22} = 2^{10}$  = 1K pages
    (not bad!)

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

# Process **A**

Virtual Memory

# Process **B**

Virtual Memory

lots of wasted space!

Physical Memory

Process **A**

Process **B**

Virtual Memory

Virtual Memory

Physical Memory

increasing page size results in increased *internal fragmentation* and *lower utilization*

i.e., TLB effectiveness needs to be balanced against memory utilization

so what about 64-bit systems?

$2^{64}$ = 16 Exabyte address space

$\approx$ 4 billion x 4GB

most modern implementations support a
max of $2^{48}$ (256TB) addressable space

page table size (assuming 4K page size)?

- # pages $= 2^{48} \div 2^{12} = 2^{36}$

- PTE size $= 8$ bytes (64 bits)

- PT size $= 2^{36} \times 8 = 2^{39}$ bytes

$= \textbf{512GB}$

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

# 512GB

(just for the virtual memory *mapping* structure)

(and we need *one per process*)

(these things aren't going to fit in memory)

instead, use *multi-level* page tables:

- split an address translation into two (or more) separate table lookups

- unused parts of the table don't need to be in memory!

# "toy" memory system
## - 8 bit addresses
## - 32-byte pages

| 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |

├──── VPN ────┤├──── page offset ────┤

## Page Table

| | |
|---|---|
| 7 | (unmapped) |
| 6 | PPN |
| 5 | (unmapped) |
| 4 | PPN |
| 3 | (unmapped) |
| 2 | (unmapped) |
| 1 | (unmapped) |
| 0 | (unmapped) |

*all 8 PTEs must be in memory at all times*

# "toy" memory system
## - 8 bit addresses
## - 32-byte pages

| 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |

page offset

3 — (unmapped)
2 — PPN
1 — (unmapped)
0 — PPN

1
0

*page "directory"*

*all unmapped;
don't need in memory!*

3 — (unmapped)
2 — (unmapped)
1 — (unmapped)
0 — (unmapped)

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

# "toy" memory system
## - 8 bit addresses
## - 32-byte pages

Linear Address

| 31 | 22 | 21 | 12 | 11 | 0 |
|----|----|----|----|----|----|
| Directory | | Table | | Offset | |

/ 12   4-KByte Page

Physical Address

/ 10   Page Table

Page Directory

PTE

PDE with PS=0

/ 20

/ 20

/ 10

/ 32

CR3

# IA-32 paging (4KB pages)

IIT College of Science
ILLINOIS INSTITUTE OF TECHNOLOGY

Linear Address

# x86-64 paging (4KB pages)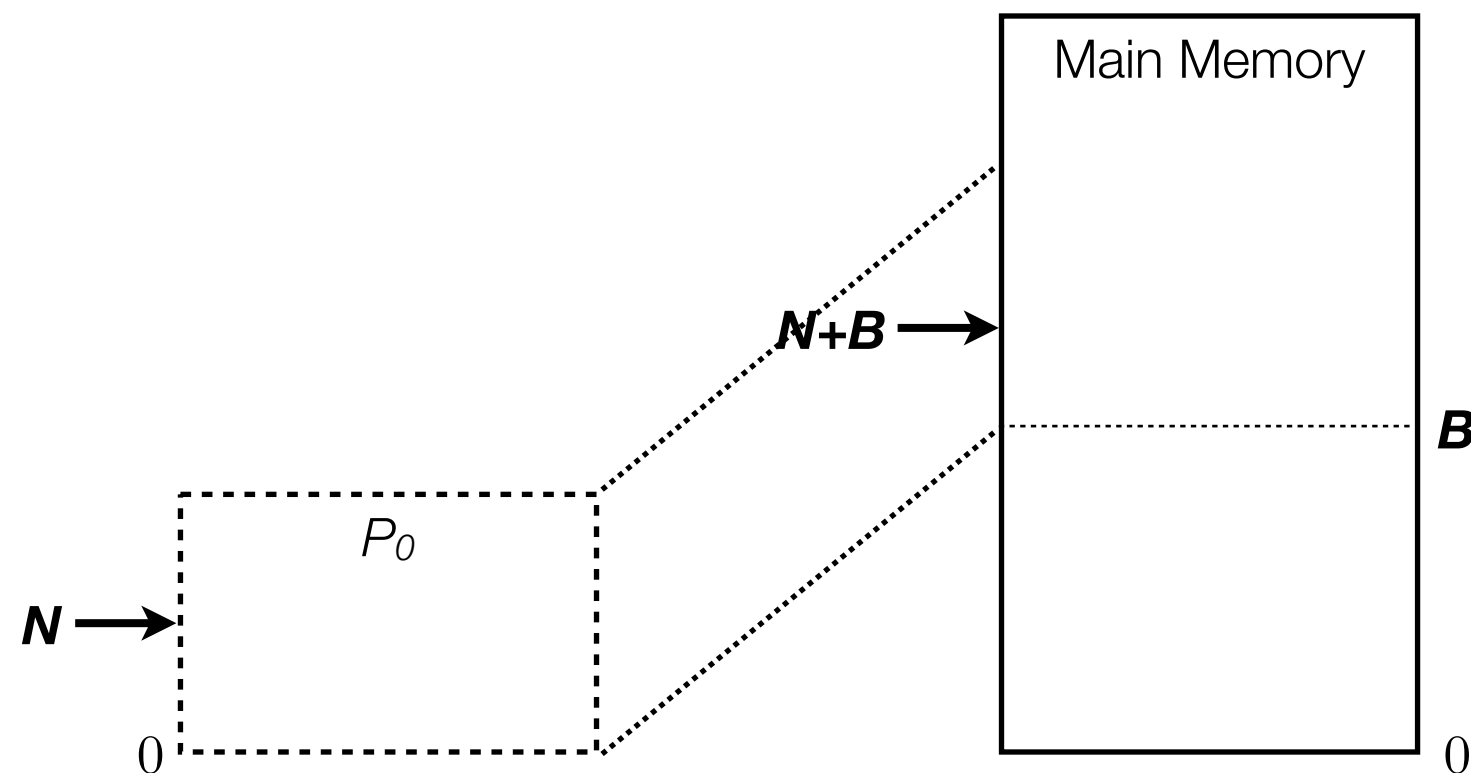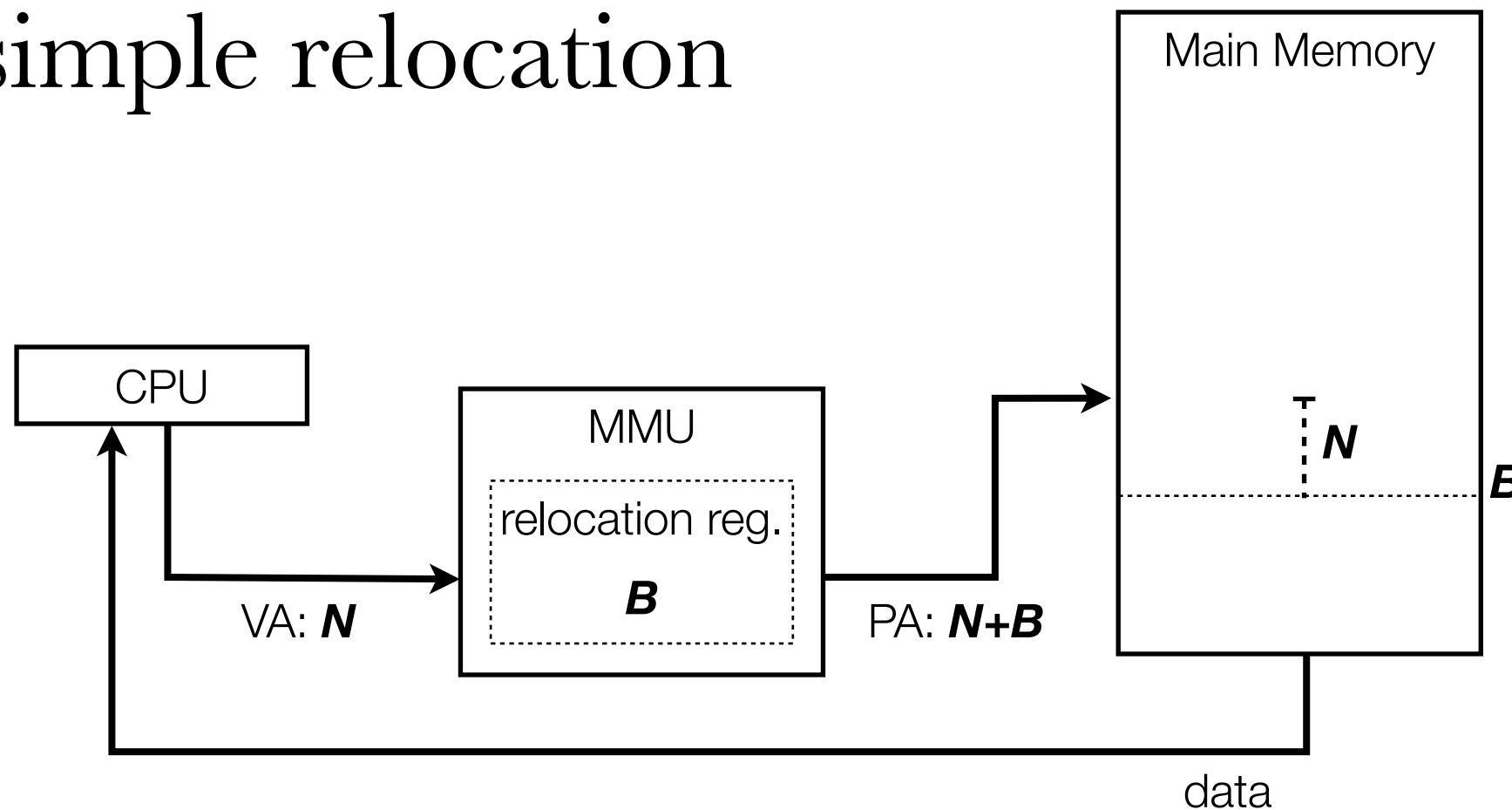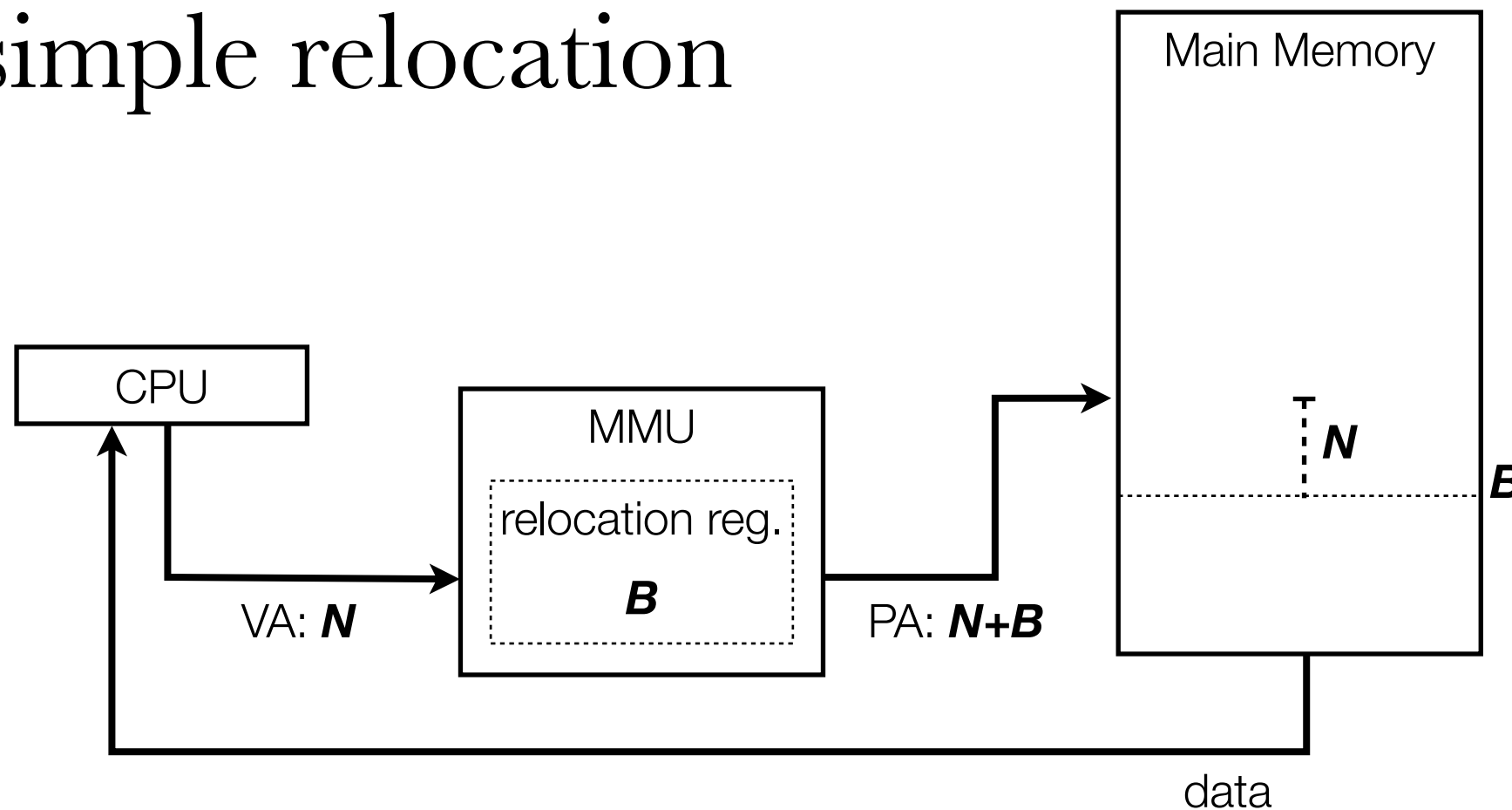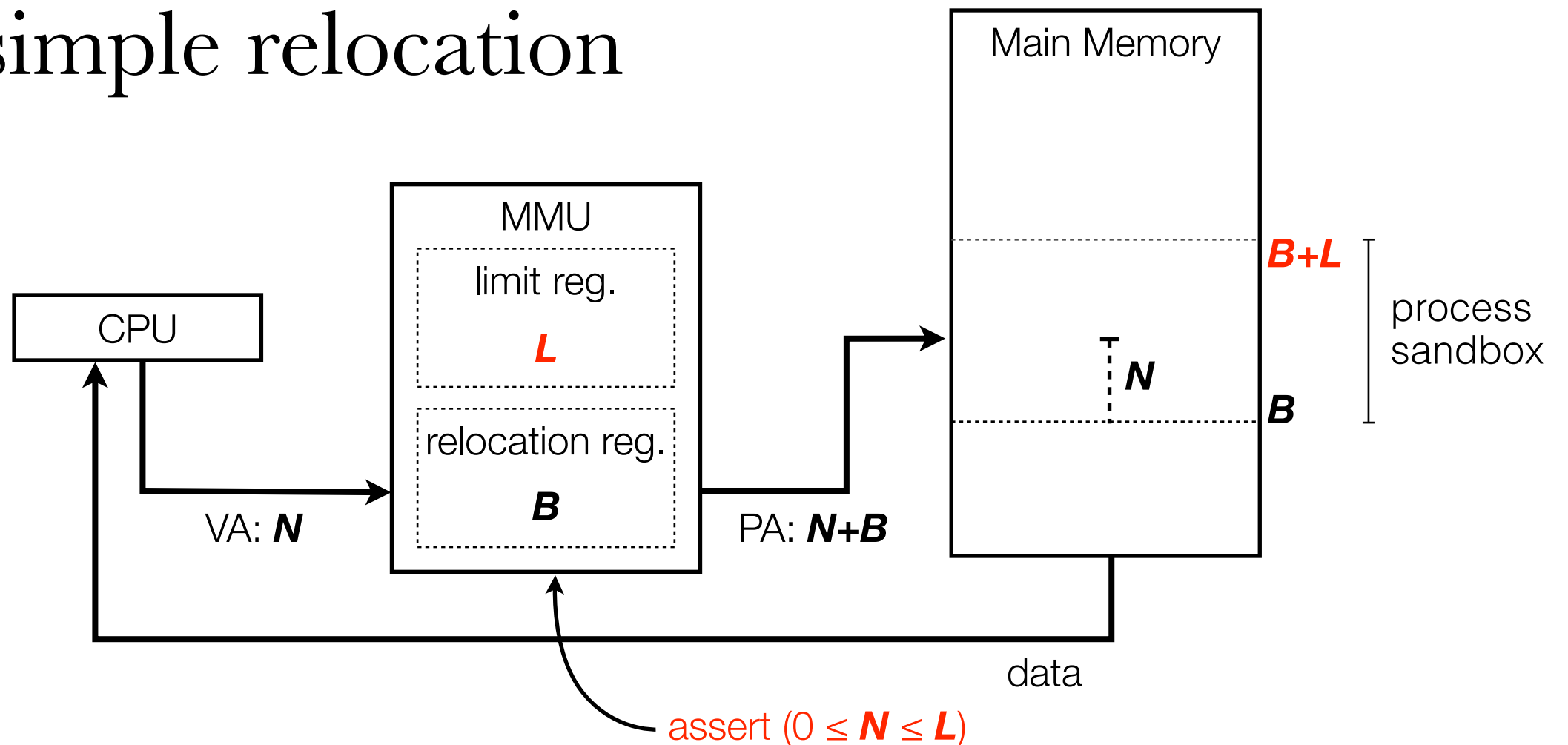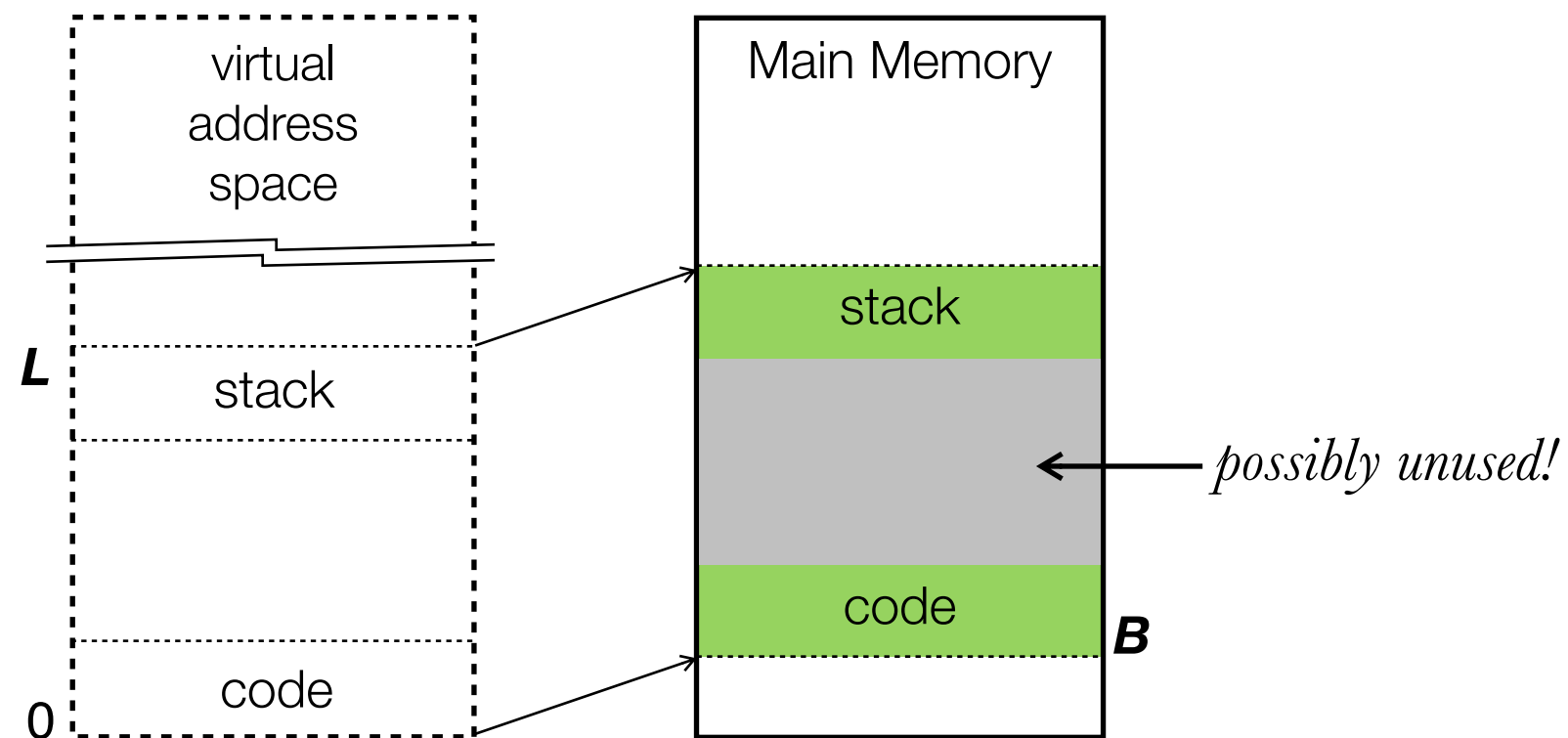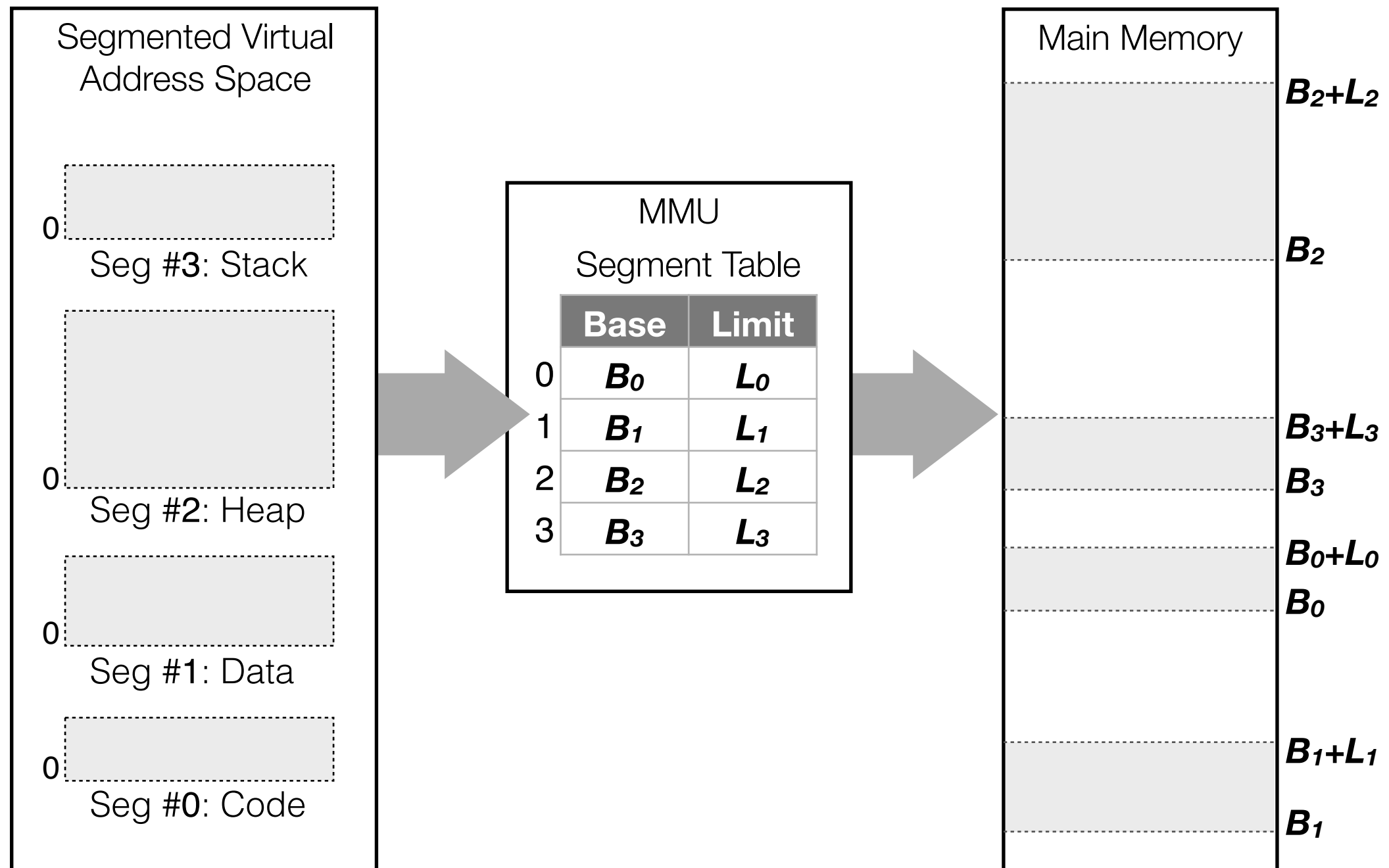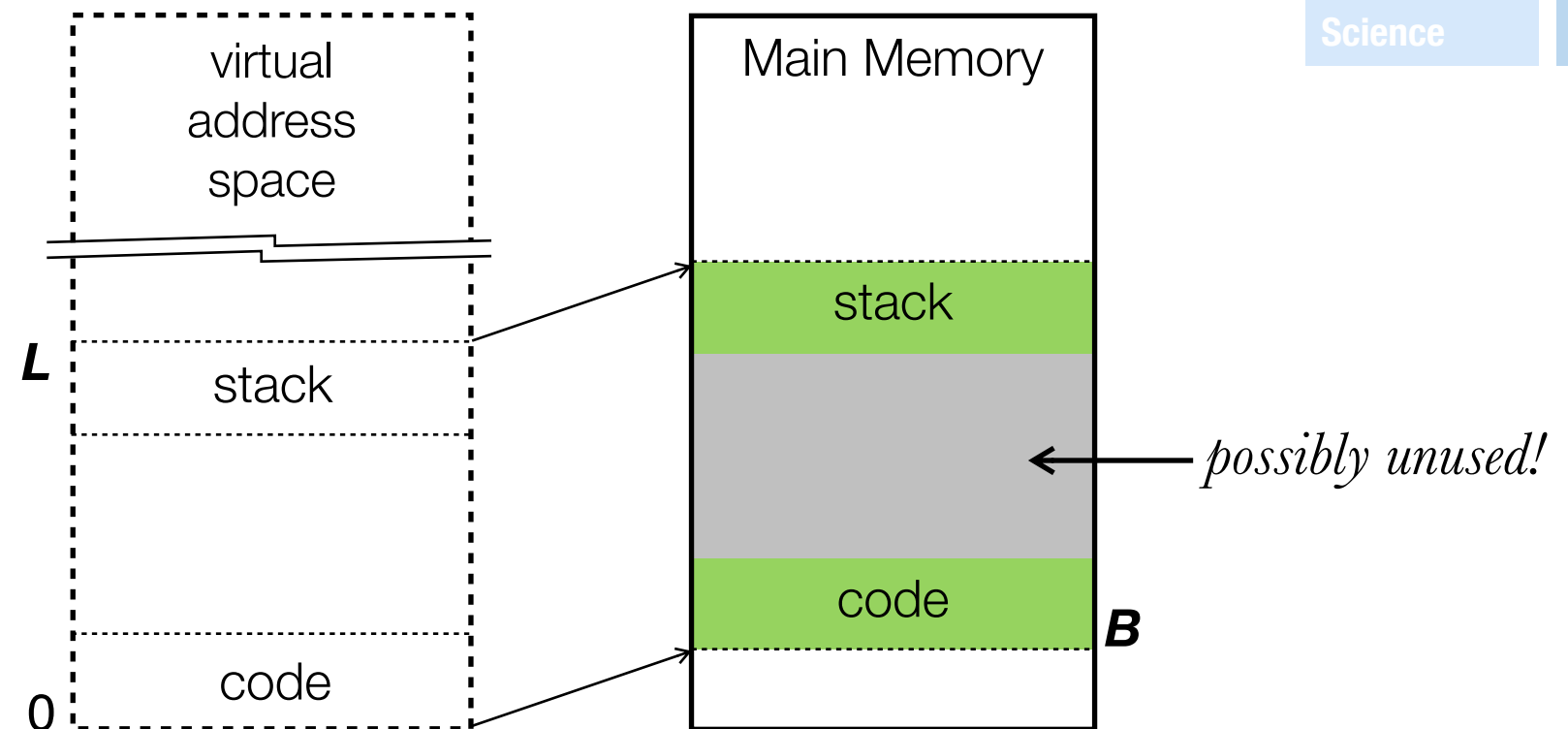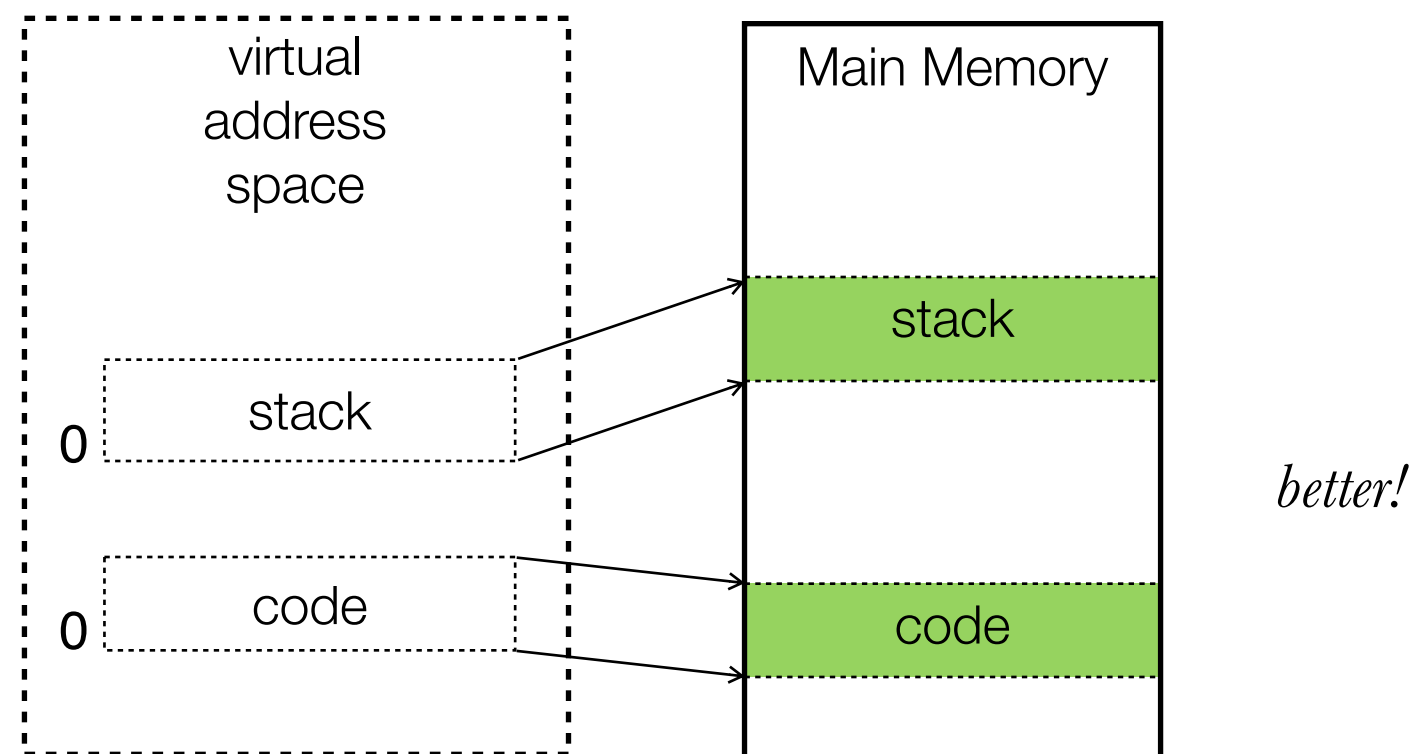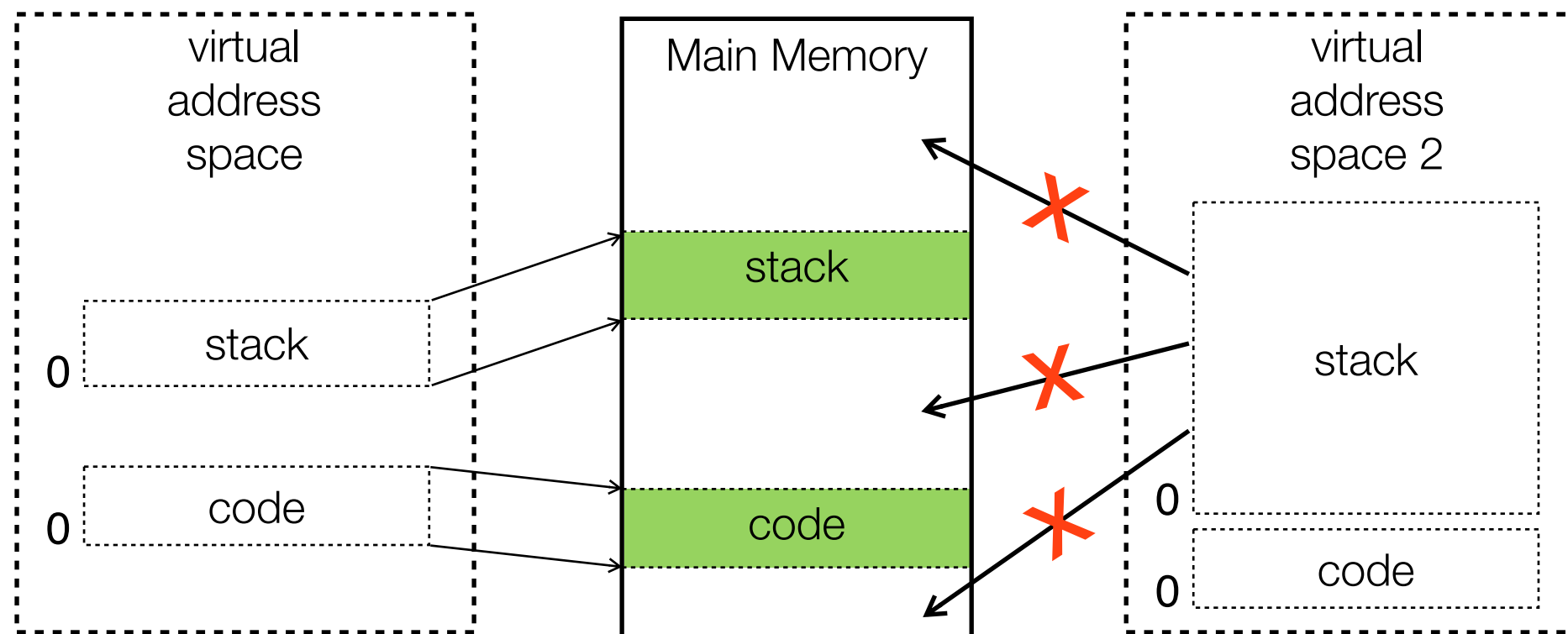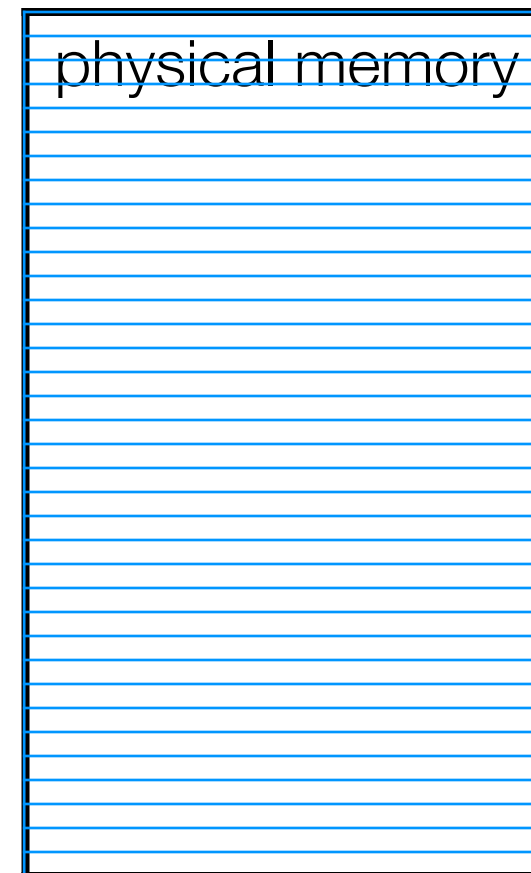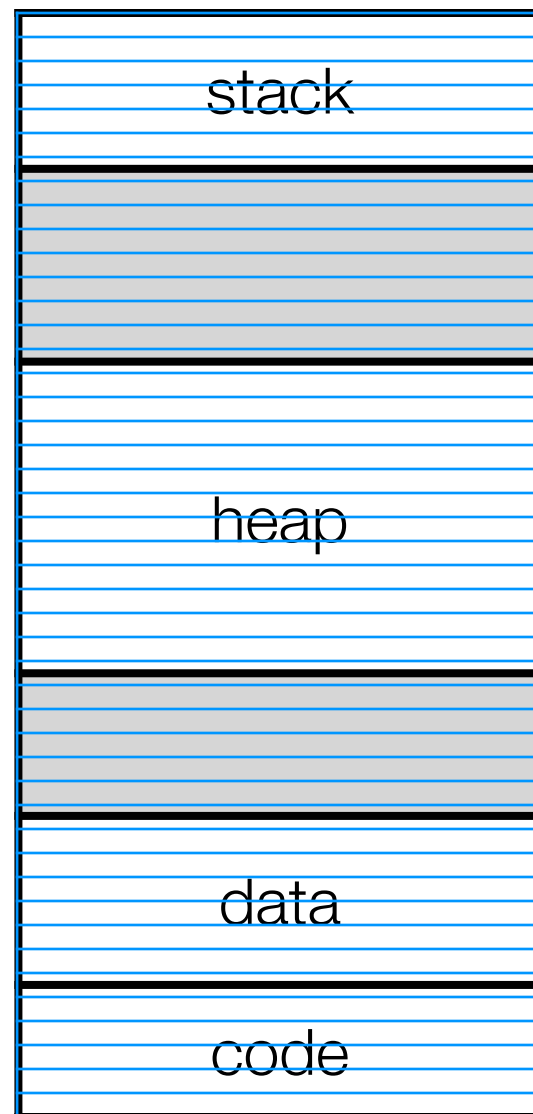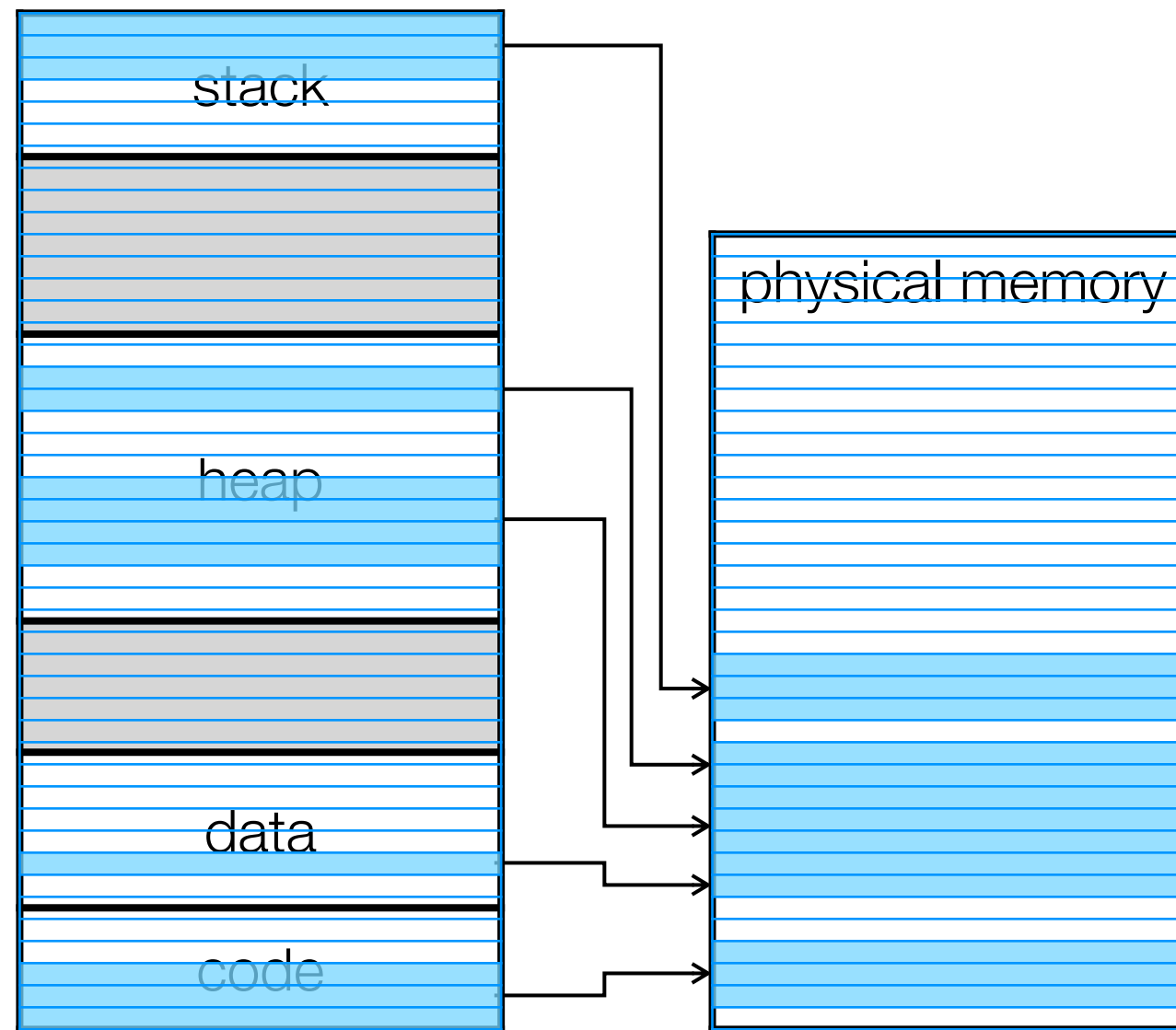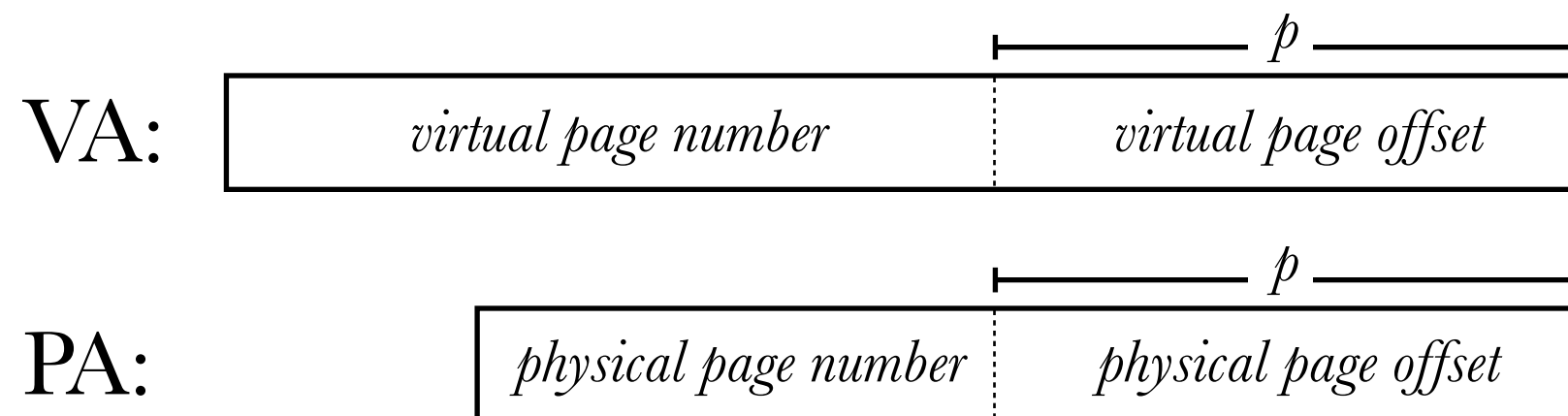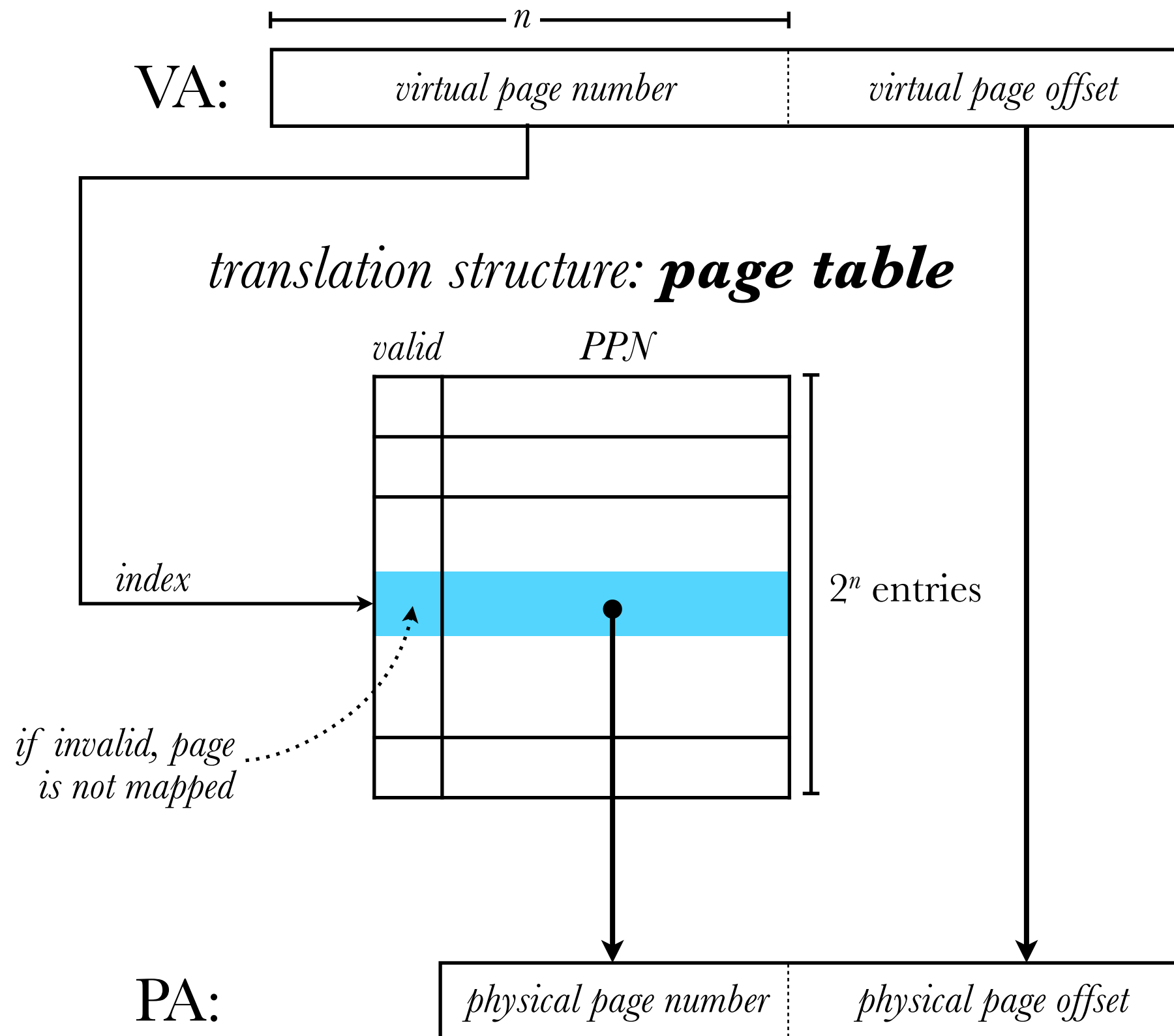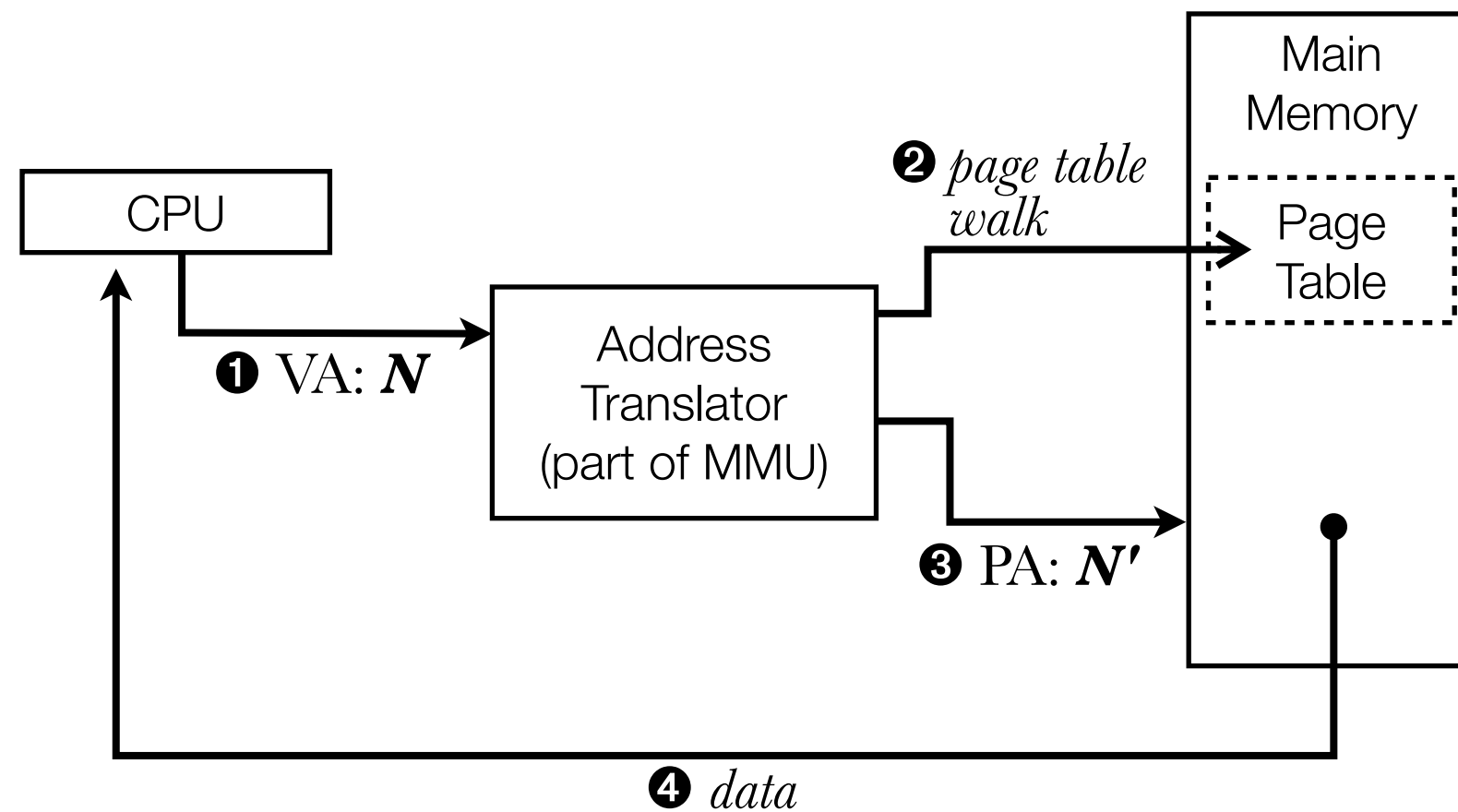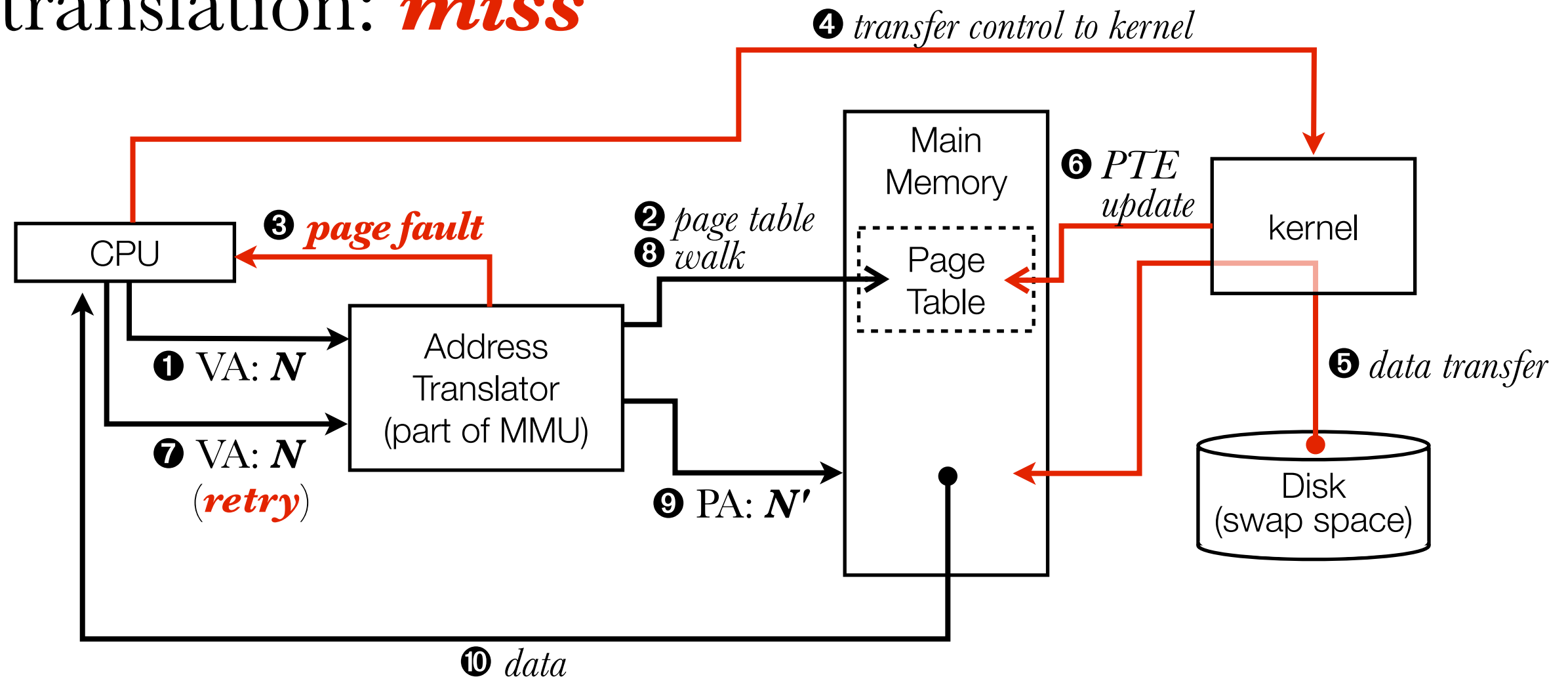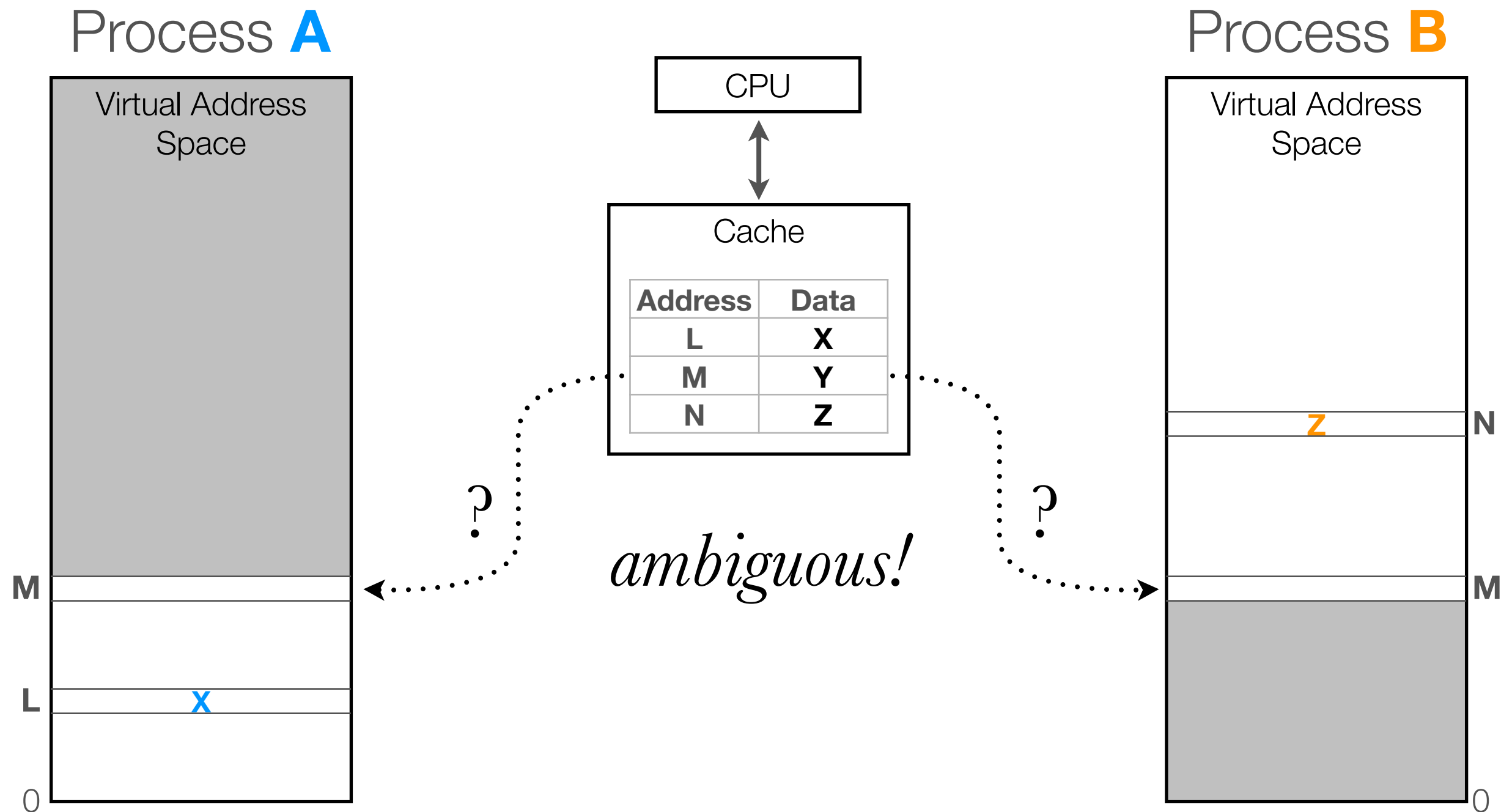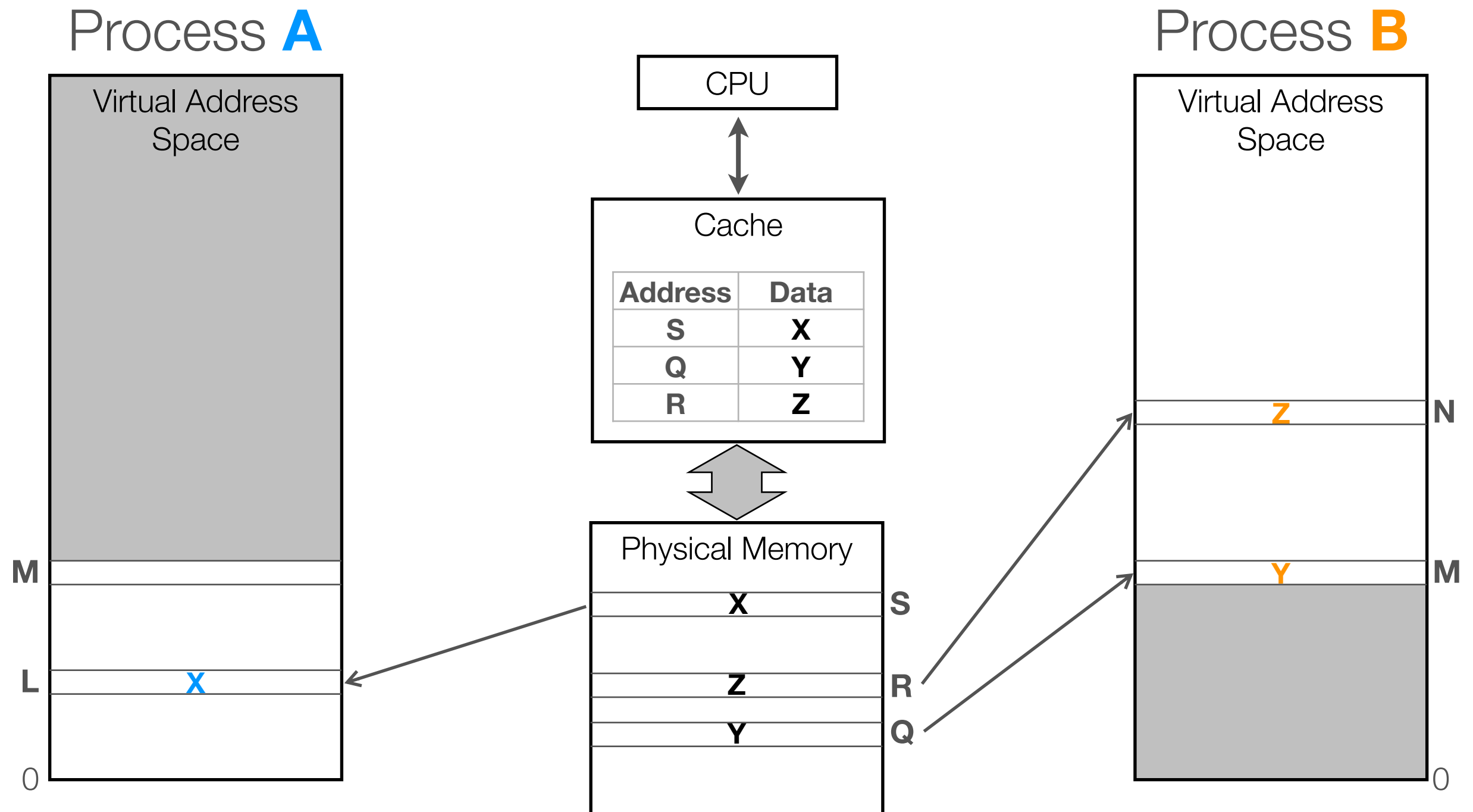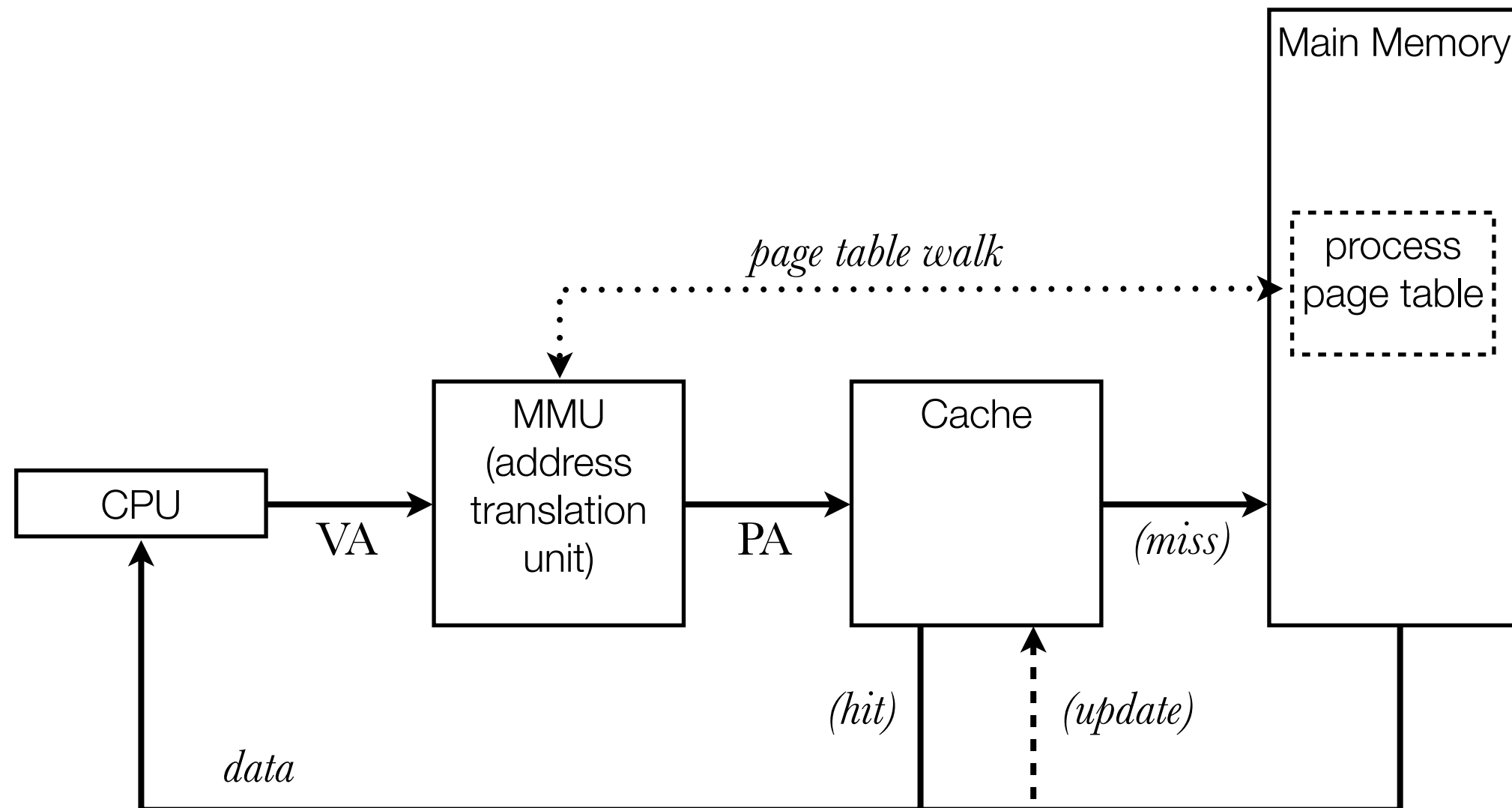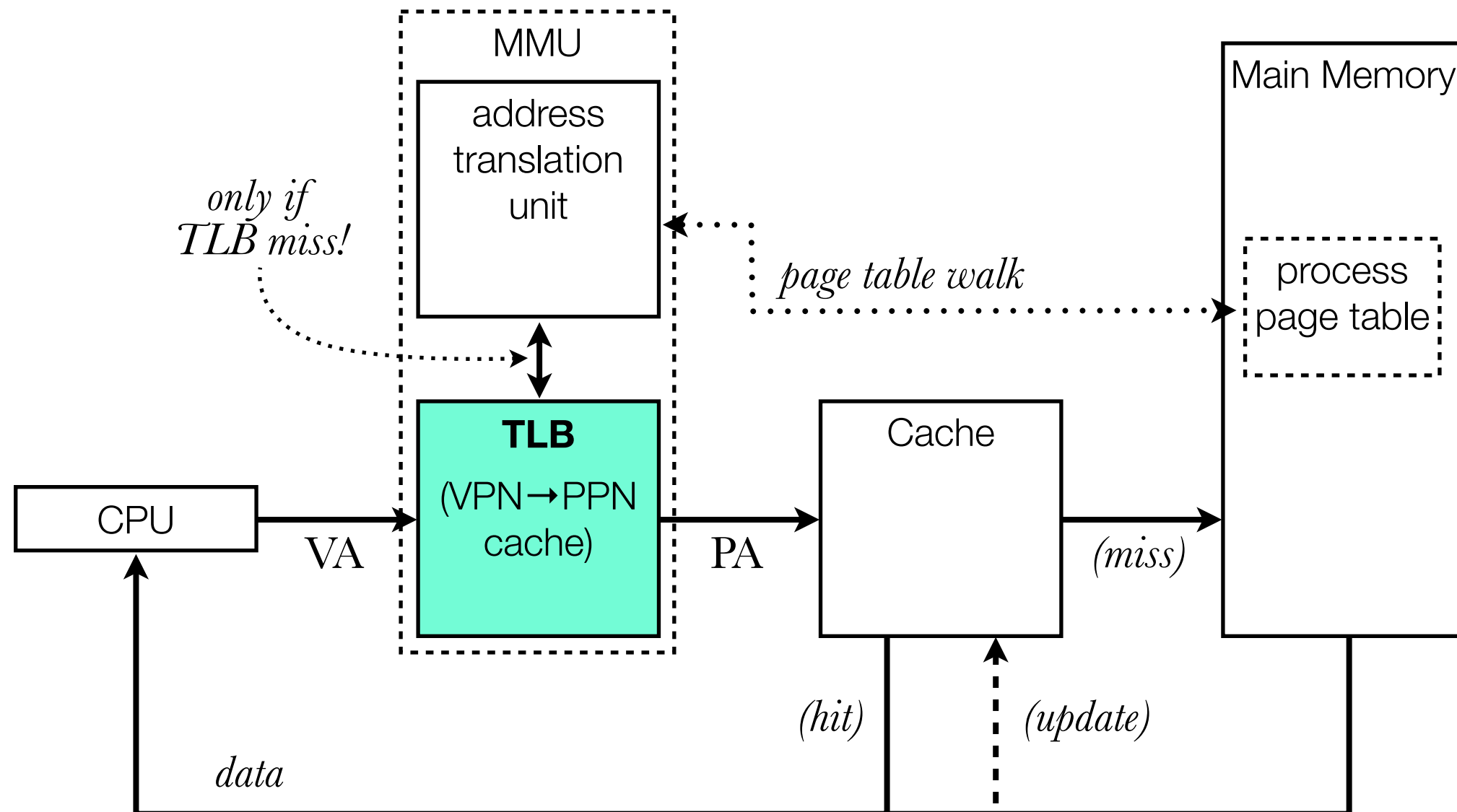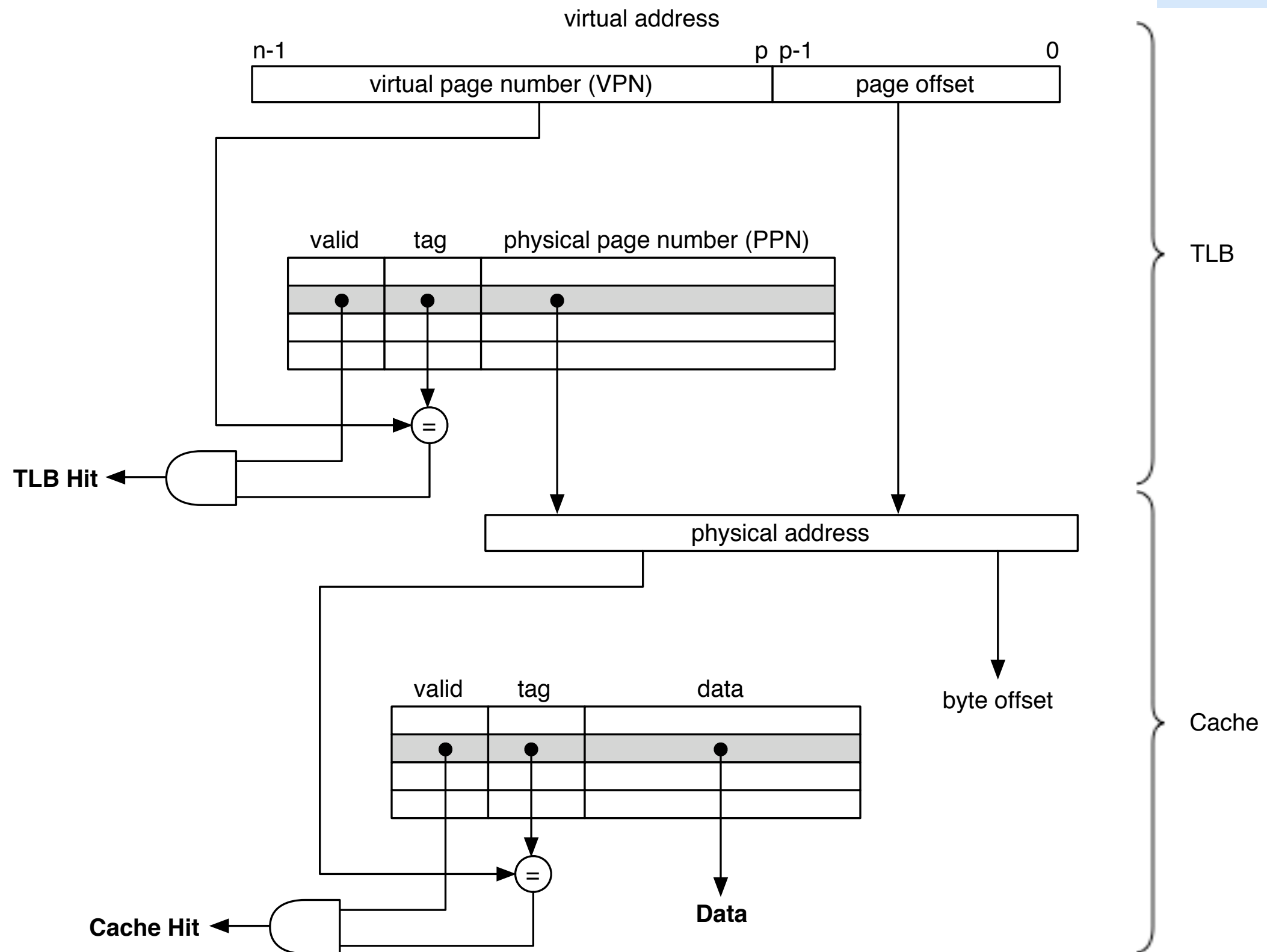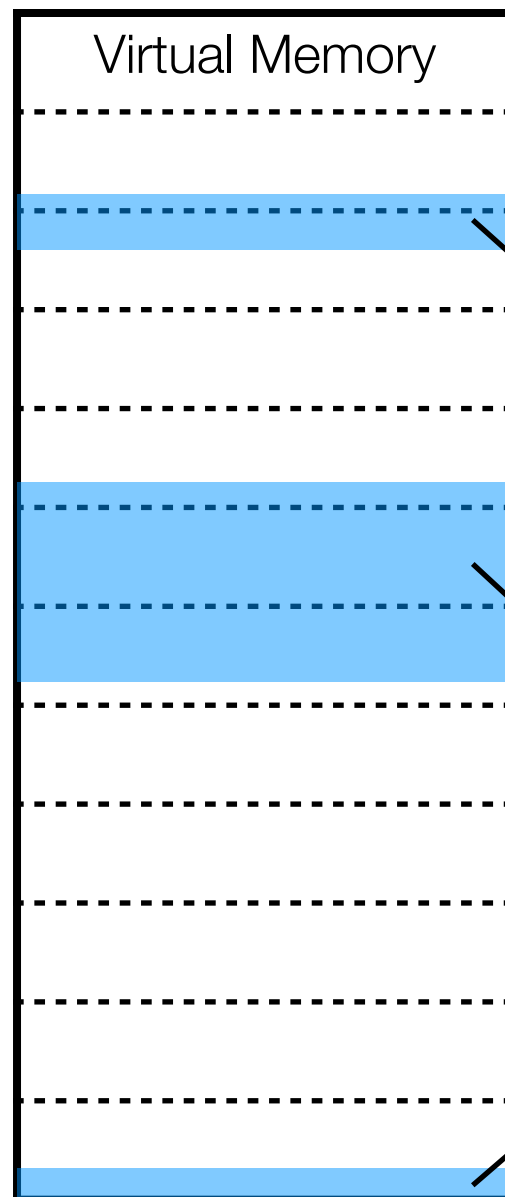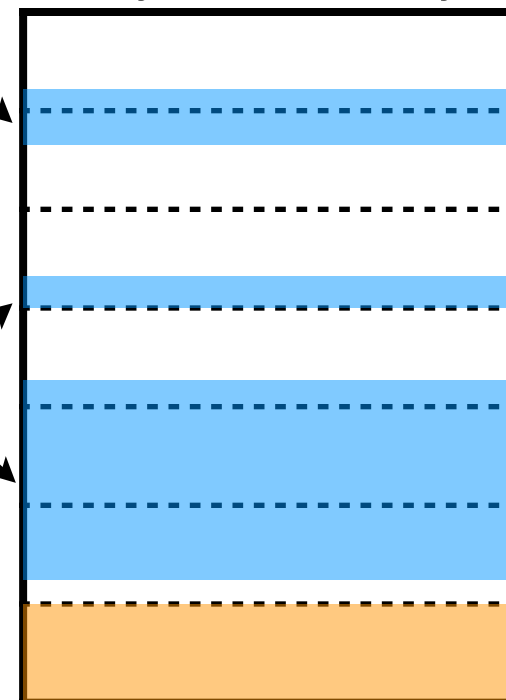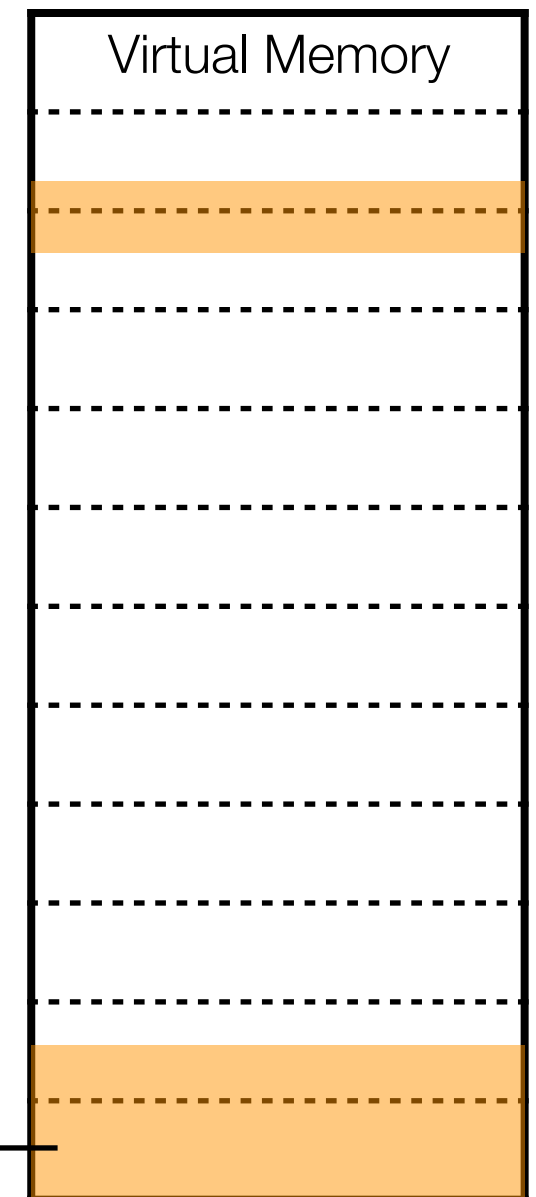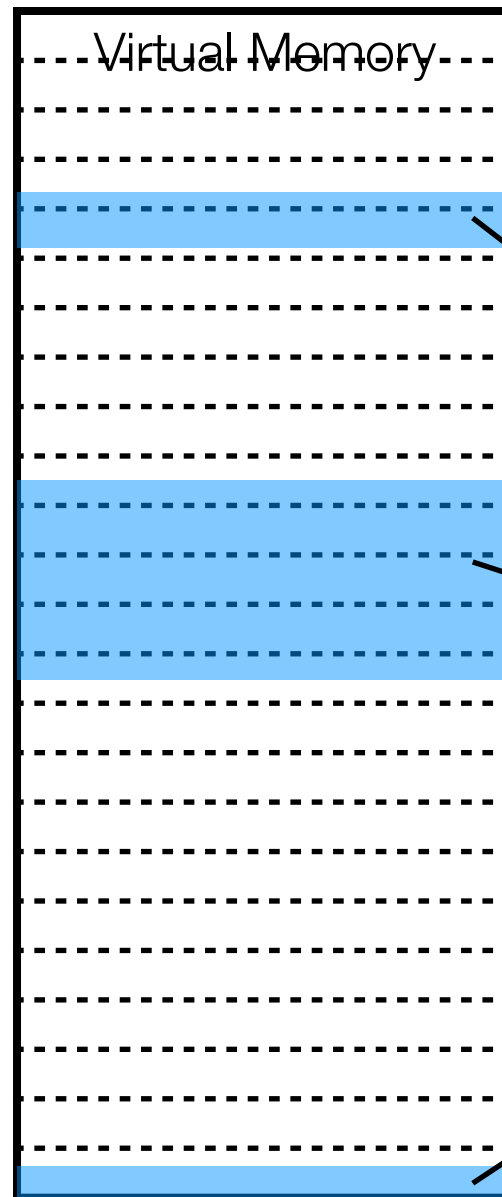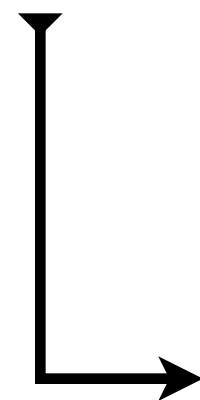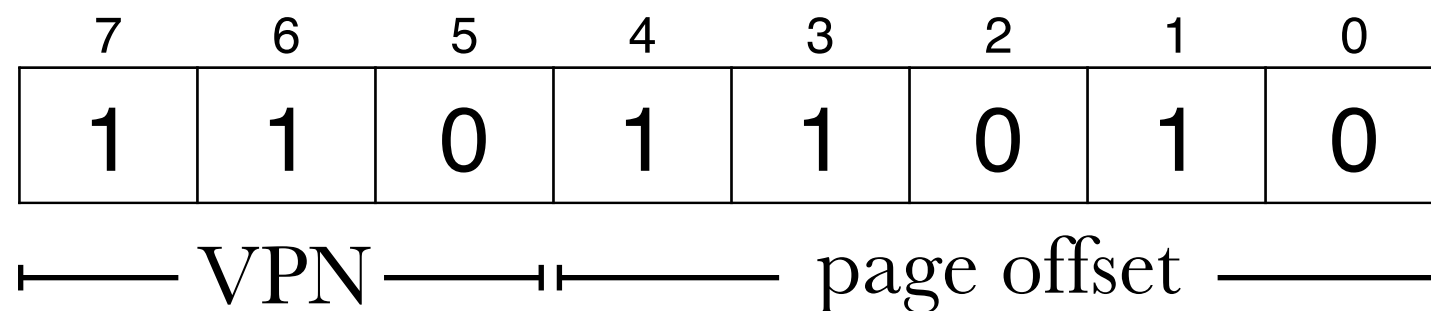