

K22 - Operating Systems: Design Principles and Internals

Fall 2025 @dit

Vaggelis Atlidakis

Lecture 17

References: Similar OS courses @Columbia, @Stanford, @UC San Diego, @Brown, @di (previous years);
and textbooks: Operating Systems: Three Easy Pieces, Operating Systems: Principles and Practice, Operating
System Concepts, Linux Kernel Development, Understanding the Linux Kernel

TOP 20 RESEARCH UNIVERSITIES IN EASTERN EUROPE BY CITATION RANK (OUT OF 796)¹

DATA SOURCES: **GOOGLE SCHOLAR** AND **WEBOMETRICS RANKING**,² JAN. 2025

CITATION RANK		HIGHER EDUCATION INSTITUTION (HEI)	COUNTRY	CITATIONS
EASTERN EUROPE (796 HEIs)	WORLD (6,191 HEIs)			
1	85	National and Kapodistrian University of Athens / Εθνικό και Καποδιστριακό Πανεπιστήμιο Αθηνών	Greece	4,864,139
2	264	Aristotle University of Thessaloniki / Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης	Greece	2,182,896
3	388	University of Tartu / Tartu Ülikool	Estonia	1,498,293
4	395	University of Patras / Πανεπιστήμιο Πατρών	Greece	1,476,112
5	425	National Technical University of Athens / Εθνικό Μετσόβιο Πολυτεχνείο	Greece	1,346,248
6	436	University of Crete / Πανεπιστήμιο Κρήτης	Greece	1,308,661
7	451	Hacettepe University / Hacettepe Üniversitesi	Turkey	1,248,796
8	456	National Research University Higher School of Economics / Национальный исследовательский университет Высшая школа экономики	Russia	< 1,240,618
9	483	Eötvös Loránd University Budapest / Eötvös Loránd Tudományegyetem ELTE	Hungary	1,148,651
10	484	University of Ljubljana / Univerza v Ljubljani	Slovenia	1,147,135
11	487	University of Thessaly / Πανεπιστήμιο Θεσσαλίας	Greece	1,134,957
12	501	Charles University in Prague / Univerzita Karlova v Praze UK	Czechia	1,089,590
13	513	University of Ioannina / Πανεπιστήμιο Ιωαννίνων	Greece	1,054,433

Overview

- We'll start from hardware and follow a question-oriented approach

~~— Intro [Q: What is an OS?]~~

~~— Events [Q: When does the OS run?]~~

~~— Runtime [Q: How does a program look like in memory?]~~

~~— Processes [Q: What is a process?]~~

~~— IPC [Q: How do processes communicate?]~~

~~— Threads [Q: What is a thread?]~~

~~— Synchronization [Q: What goes wrong w/o synchronization?]~~

~~— Time Management [Q: What is scheduling?]~~

- **Memory Management [Q: What is virtual memory?]**

- **Files [Q: What is a file descriptor?]**

- **Storage Management [Q: How do we allocate disk space to files?]**

- * Basic (H/W & S/W)
- * **Abstractions**
- * **Primitives**
- * **Mechanisms**

Overview

- **Memory management**

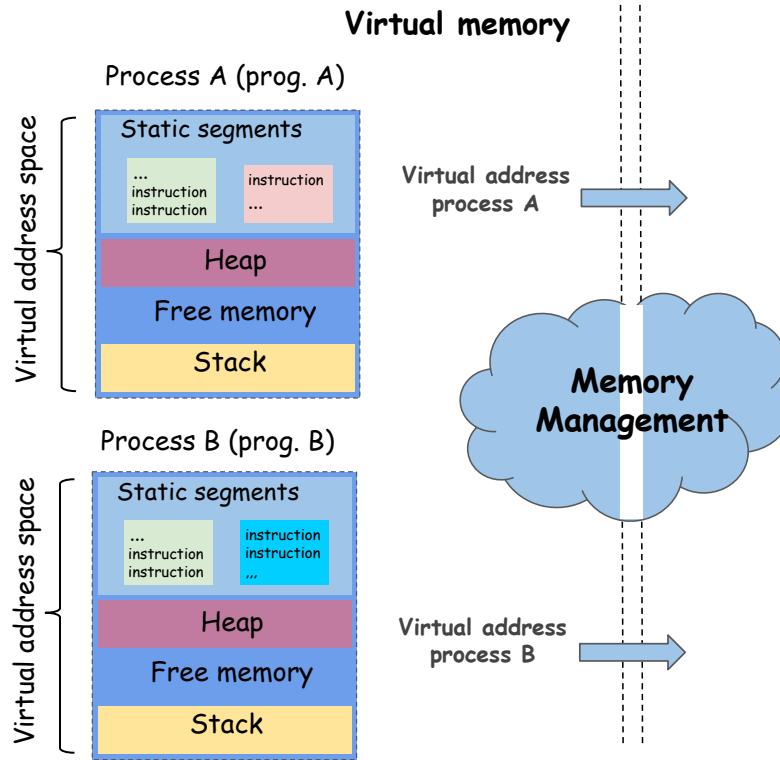
- Q1: What is virtual memory?
- Q2: Why is virtual memory necessary?
- Q3: How is virtual memory implemented?
- ...

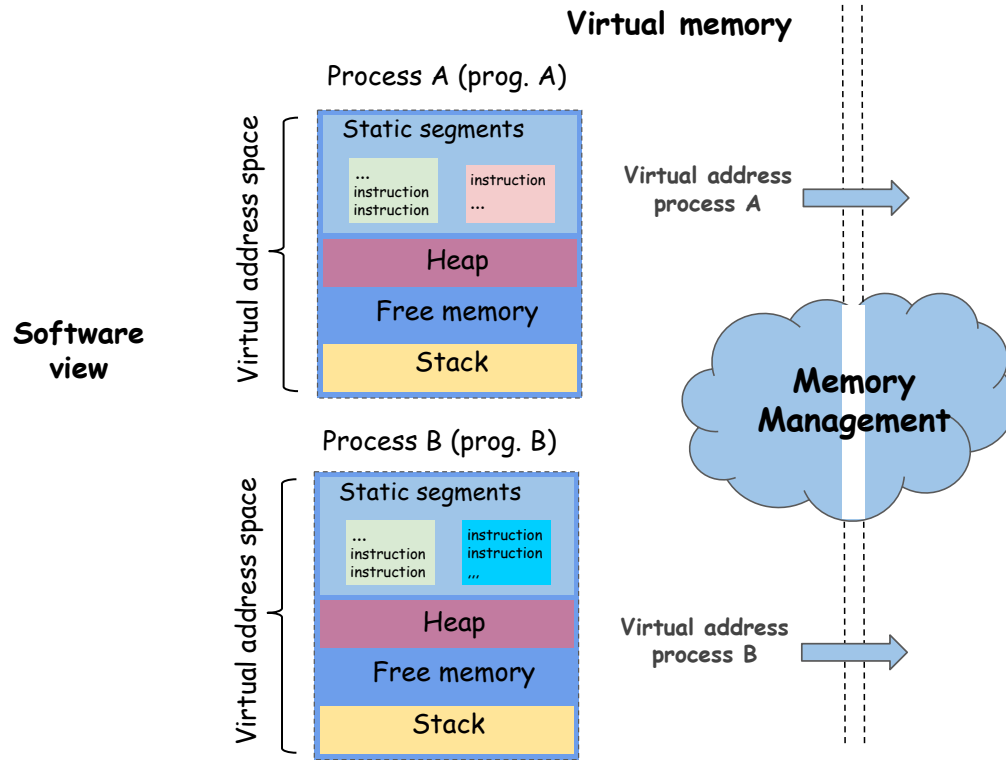
What is virtual memory?

>> A layer of abstraction

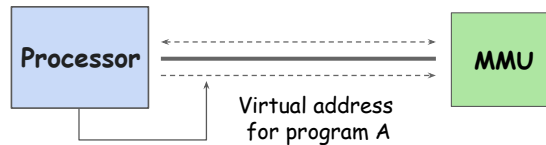
- Translates virtual addresses into physical addresses
- **Virtual addr.:** the language processes talk to the processor about memory
- **Physical addr.:** the language actual contents of memory are accessed
- The way programmers imagine programs interact with memory is an illusion on top of the virtual memory abstraction layer

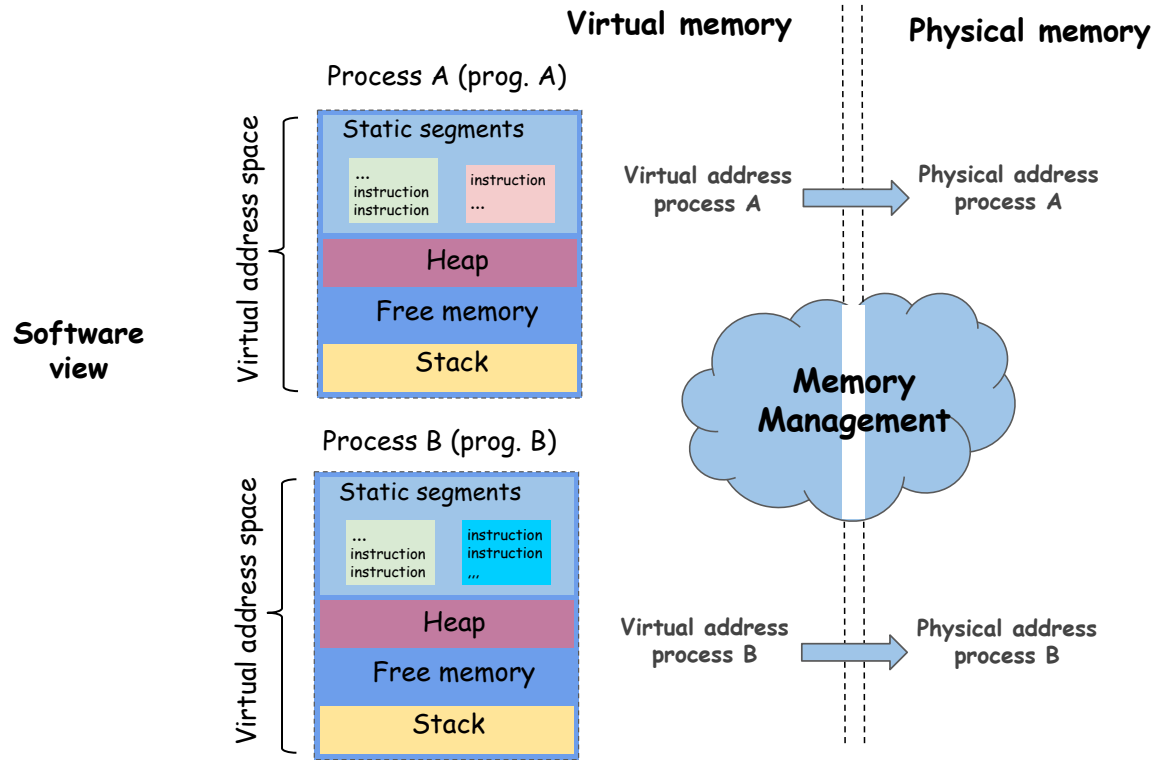
**Software
view**



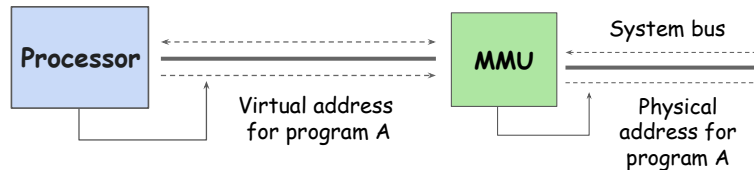


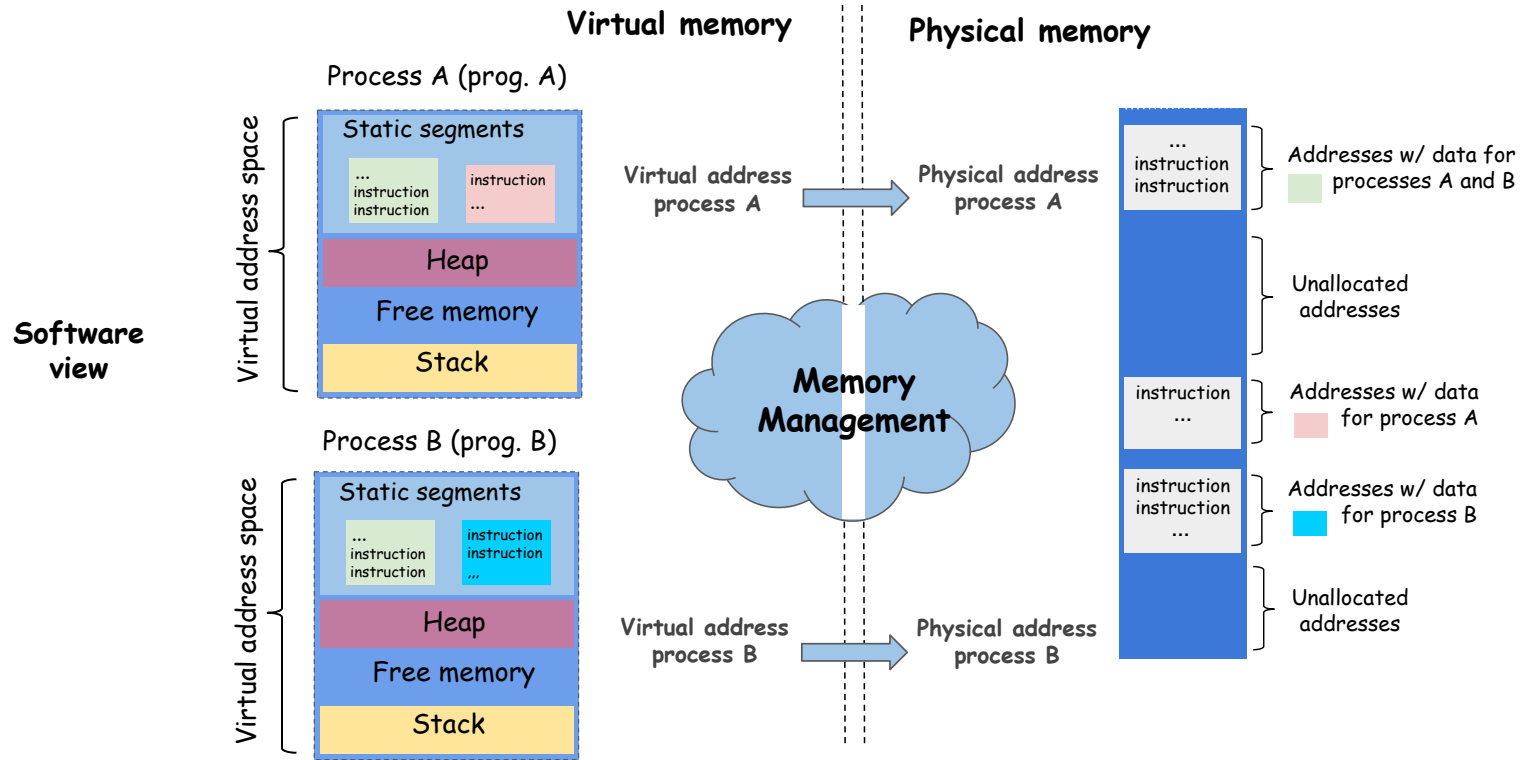
Hardware view



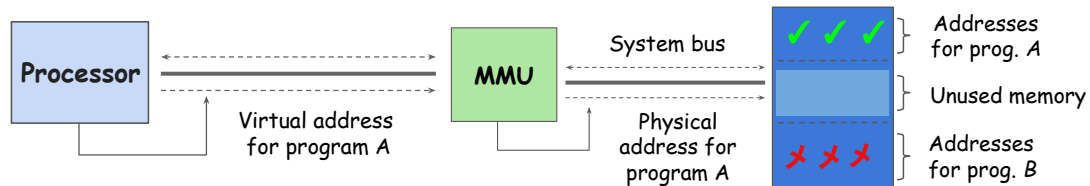


Hardware view





Hardware view



What is virtual memory? (demo)

What is virtual memory? (demo)

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

int main(int argc, char **argv) {
    int i;
    pid_t pid = getpid();

    i = atoi(argv[1]);

    printf("[PID: %d]; &i: %p; i: %d\n", pid, &i, i);
    printf("[PID: %d]; Sleep for %d sec\n", pid, i);
    sleep(i);
    printf("[PID: %d]; I am done now\n", pid, i);
    return 0;
}
```

What is virtual memory? (demo)

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

int main(int argc, char **argv) {
    int i;
    pid_t pid = getpid();

    i = atoi(argv[1]);

    printf("[PID: %d]; &i: %p; i: %d\n", pid, &i, i)
    printf("[PID: %d]; Sleep for %d sec\n", pid, i);
    sleep(i);
    printf("[PID: %d]; I am done now\n", pid, i);
    return 0;
}
```

→ git:(master) X echo 0 | sudo tee
/proc/sys/kernel/randomize_va_space
0

→ vm git:(master) X ./hello_vm 10 &
[1] 186394
[PID: 186394]; &i: 0xffffffff210; i: 10
[PID: 186394]: Seep for 10 sec

→ vm git:(master) X ./hello_vm 3 &
[2] 186432
[PID: 186432]; &i: 0xffffffff210; i: 3
[PID: 186432]: Sleep for 3 sec

→ vm git:(master) X [PID: 186432]; I am done now
[2] + 186432 done ./hello_vm 3

→ vm git:(master) X [PID: 186394]; I am done now
[1] + 186394 done ./hello_vm 10

What is virtual memory? (demo)

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

int main(int argc, char **argv) {
    int i;
    pid_t pid = getpid();

    i = atoi(argv[1]);

    printf("[PID: %d]; &i: %p; i: %d\n", pid, &i, i)
    printf("[PID: %d]; Sleep for %d sec\n", pid, i);
    sleep(i);
    printf("[PID: %d]; I am done now\n", pid, i);
    return 0;
}
```

→ git:(master) X echo 0 | sudo tee
/proc/sys/kernel/randomize_va_space
0

→ vm git:(master) X ./hello_vm 10 &
[1] 186394

[PID: 186394]; &i: 0xffffffff210; i: 10
[PID: 186394]: Seep for 10 sec

→ vm git:(master) X ./hello_vm 3 &
[2] 186432

[PID: 186432]; &i: 0xffffffff210; i: 3
[PID: 186432]: Sleep for 3 sec

} Same address
different value,
both programs
running...

→ vm git:(master) X [PID: 186432]; I am done now
[2] + 186432 done ./hello_vm 3

→ vm git:(master) X [PID: 186394]; I am done now
[1] + 186394 done ./hello_vm 10

Why need virtual memory?

Why need virtual memory?

- **Fault isolation:** Anything that goes wrong in the address space of one process does not affect any other process

Why need virtual memory?

- **Fault isolation:** Anything that goes wrong in the address space of one process does not affect any other process
- **Illusion of continuous memory:** Programs are written, compiled, assembled, linked, and loaded assuming access to **continuous** memory

Why need virtual memory?

- **Fault isolation:** Anything that goes wrong in the address space of one process does not affect any other process
- **Illusion of continuous memory:** Programs are written, compiled, assembled, linked, and loaded assuming access to **continuous** memory
- **Frugal use of resources**
 - >> Demand paging: Postpone allocating physical memory until necessary

Why need virtual memory?

- **Fault isolation:** Anything that goes wrong in the address space of one process does not affect any other process
- **Illusion of continuous memory:** Programs are written, compiled, assembled, linked, and loaded assuming access to **continuous** memory
- **Frugal use of resources**
 - >> Demand paging: Postpone allocating physical memory until necessary
 - >> Copy-On-Write: One physical replica of common readed-only code

Why need virtual memory?

- **Fault isolation:** Anything that goes wrong in the address space of one process does not affect any other process
- **Illusion of continuous memory:** Programs are written, compiled, assembled, linked, and loaded assuming access to **continuous** memory
- **Frugal use of resources**
 - >> Demand paging: Postpone allocating physical memory until necessary
 - >> Copy-On-Write: One physical replica of common readed-only code
- **Performant use of resources**
 - >> Programs have a memory footprint as much as the size of storage and run at speed close to the speed of CPU caches (how?...)

Why need virtual memory?

- **Fault isolation:** Anything that goes wrong in the address space of one process does not affect any other process
- **Illusion of continuous memory:** Programs are written, compiled, assembled, linked, and loaded assuming access to **continuous** memory
- **Frugal use of resources**
 - >> Demand paging: Postpone allocating physical memory until necessary
 - >> Copy-On-Write: One physical replica of common readed-only code
- **Performant use of resources**
 - >> Programs have a memory footprint as much as the size of storage and run at speed close to the speed of CPU caches (how?...)
 - >> **CoW:** Laziness -> defering Vs. **Locality:** Proactivity -> prefetching

How is virtual memory implemented?

"We've rewritten the VM several times in the last ten years, and I expect it will be changed several more times in the next few years. Within five years, we'll almost certainly have to make the current three-level page tables be four levels, etc."

—Linus Torvald, 2001.

How is virtual memory implemented?

- > Usually considered the most complex kernel subsystem
 - It's simple if you view it as an index

Ask the right questions

How is virtual memory implemented?

- > Usually considered the most complex kernel subsystem
- It's simple if you view it as an index

Ask the right questions

Q1: What are the entries of the index?

Q2: How are the entries of the index used?

Q3: How are the entries of the index allocated?

Q4: How are the entries of the index replaced?

How is virtual memory implemented?

- > Usually considered the most complex kernel subsystem
- It's simple if you view it as an index

Ask the right questions

Q1: What are the entries of the index?

Q2: How are the entries of the index used?

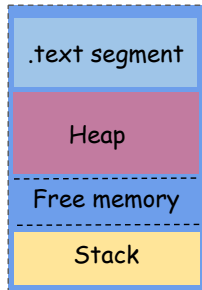
Q3: How are the entries of the index allocated?

Q4: How are the entries of the index replaced?

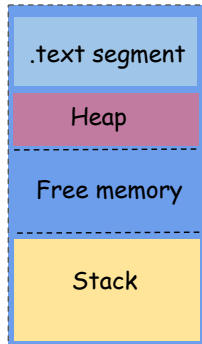
Mechanism: Segmentation

Virtual memory

Process A



Process B

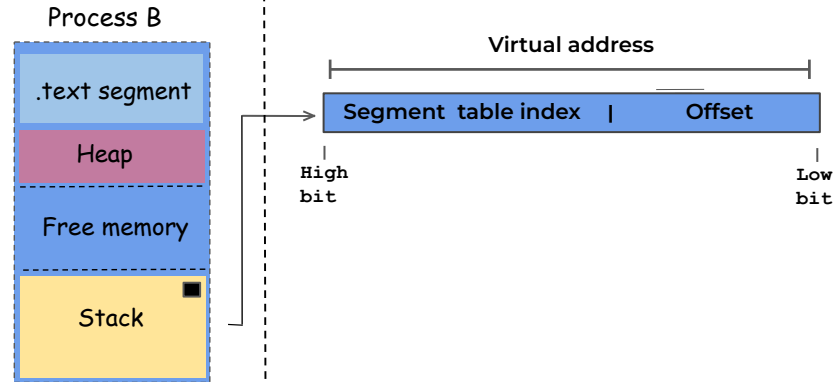
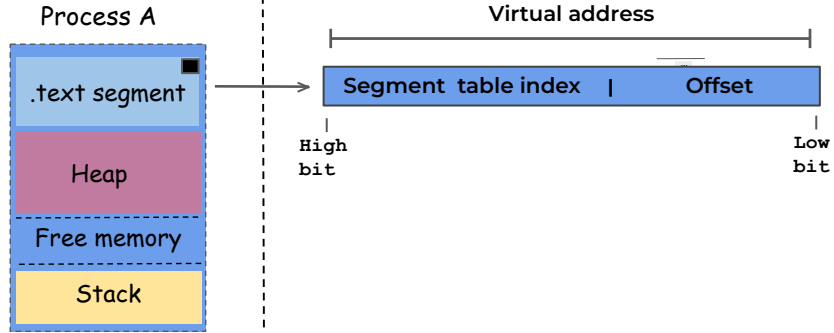


Physical memory



Mechanism: Segmentation

Virtual memory

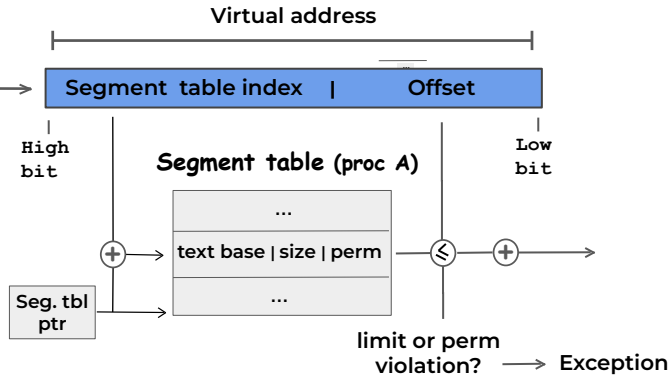
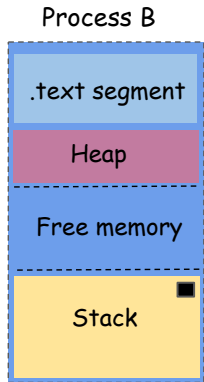
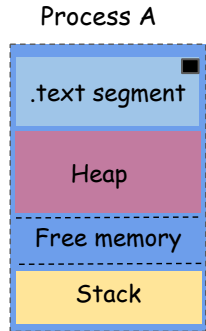


Physical memory



Mechanism: Segmentation

Virtual memory

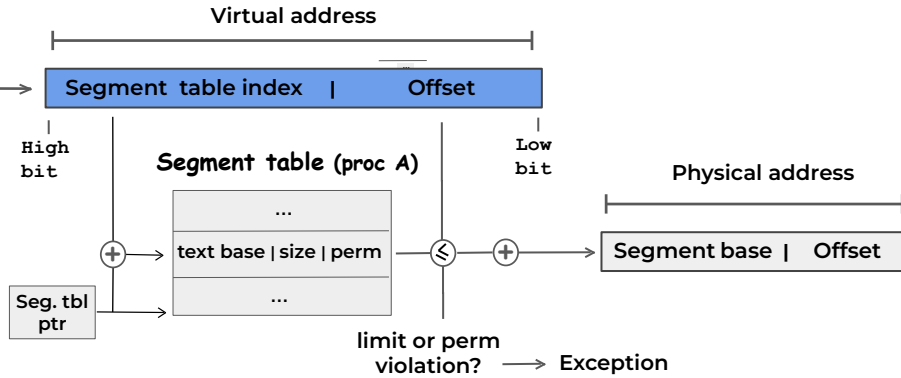
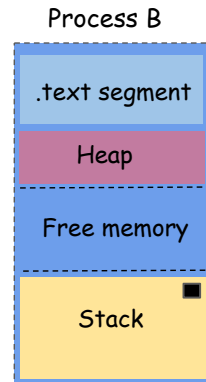
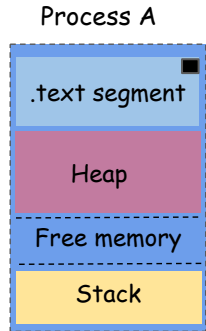


Physical memory

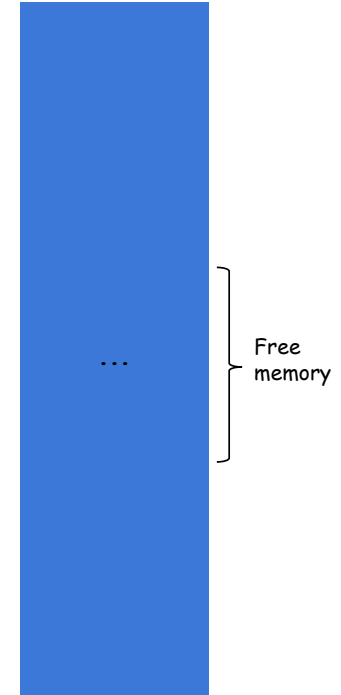


Mechanism: Segmentation

Virtual memory

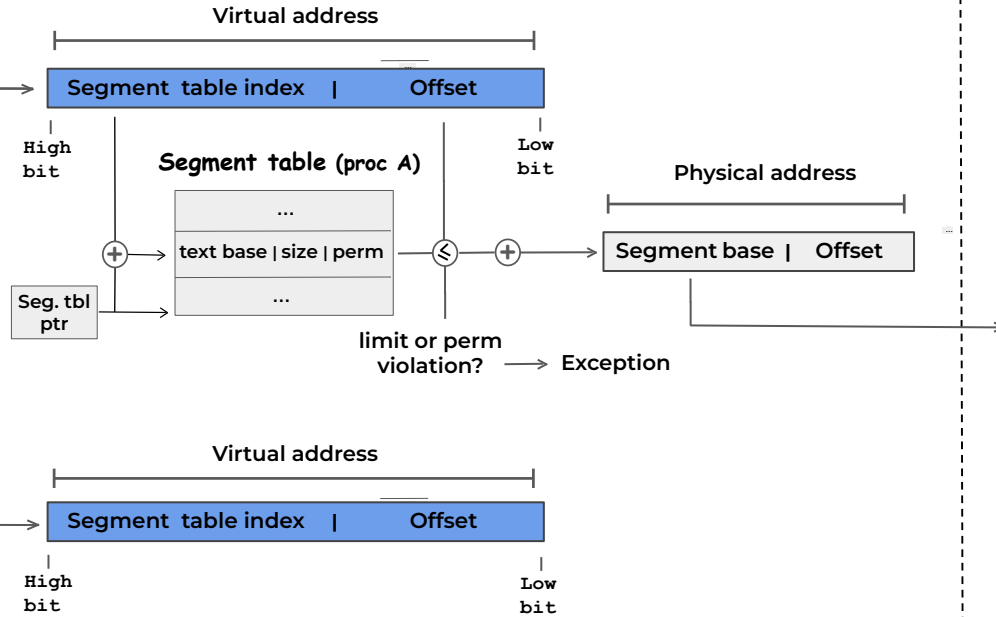
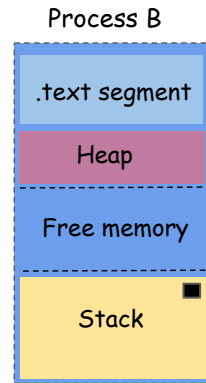
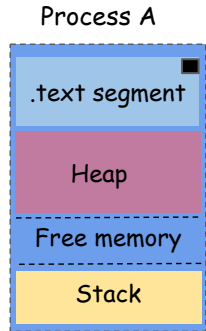


Physical memory

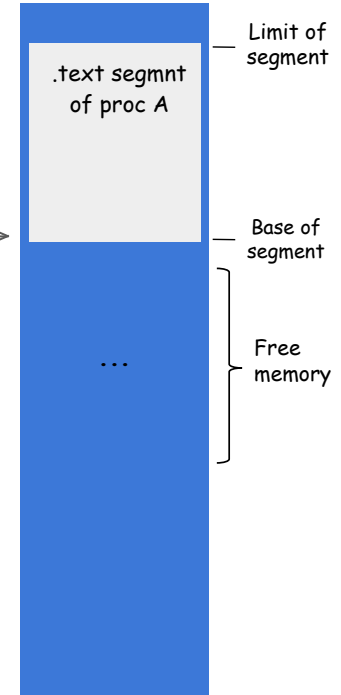


Mechanism: Segmentation

Virtual memory

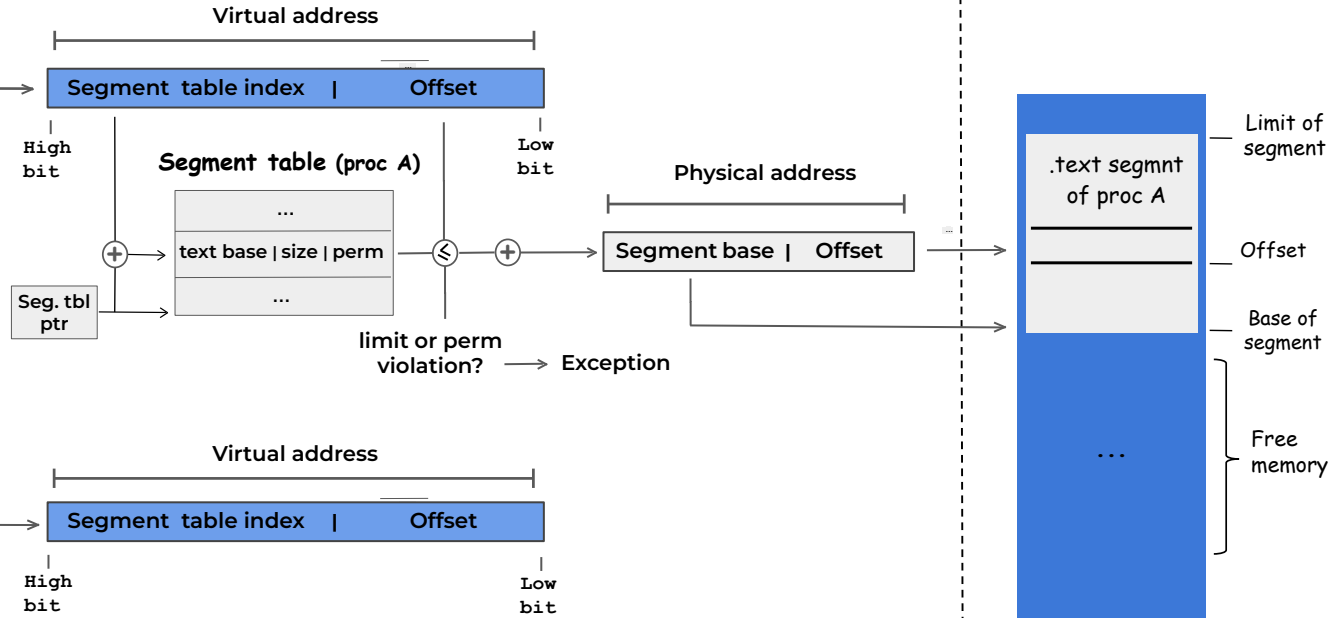
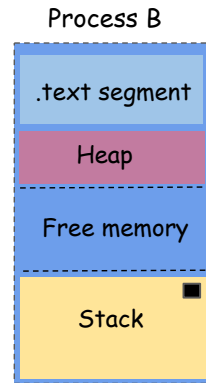
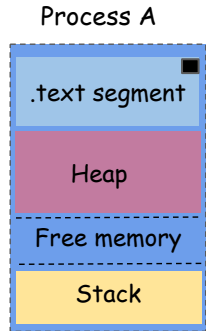


Physical memory



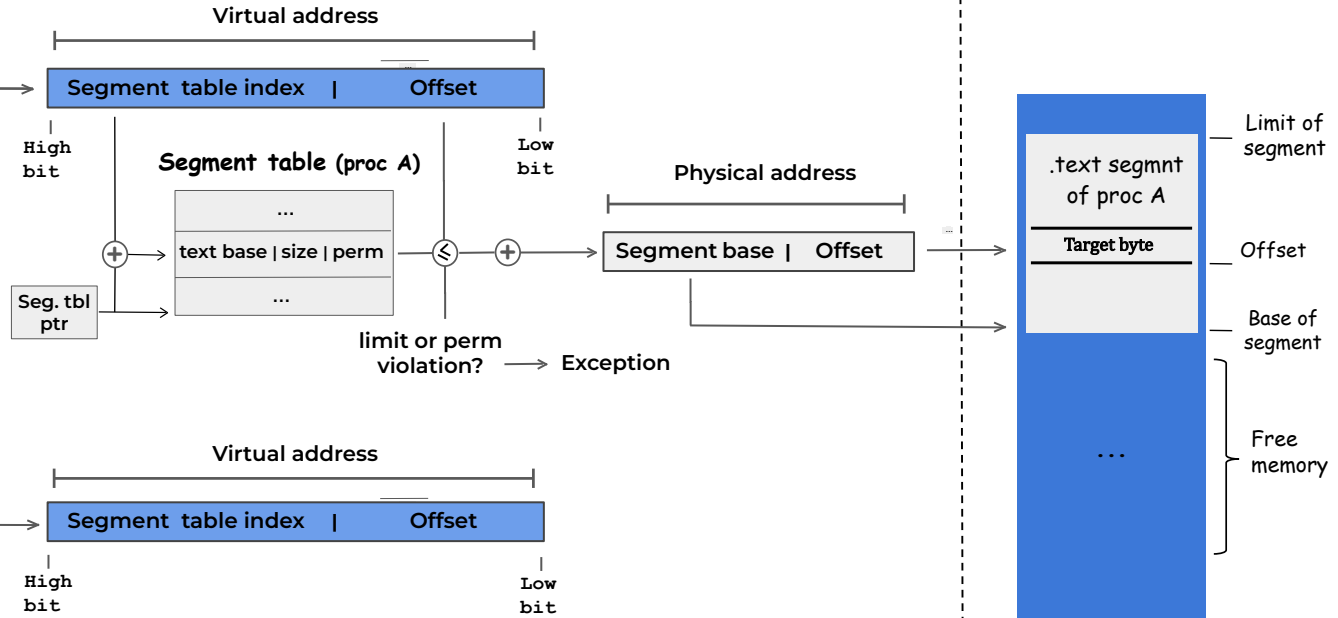
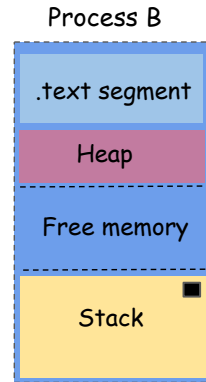
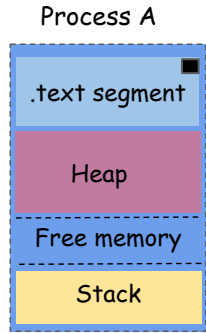
Mechanism: Segmentation

Virtual memory



Mechanism: Segmentation

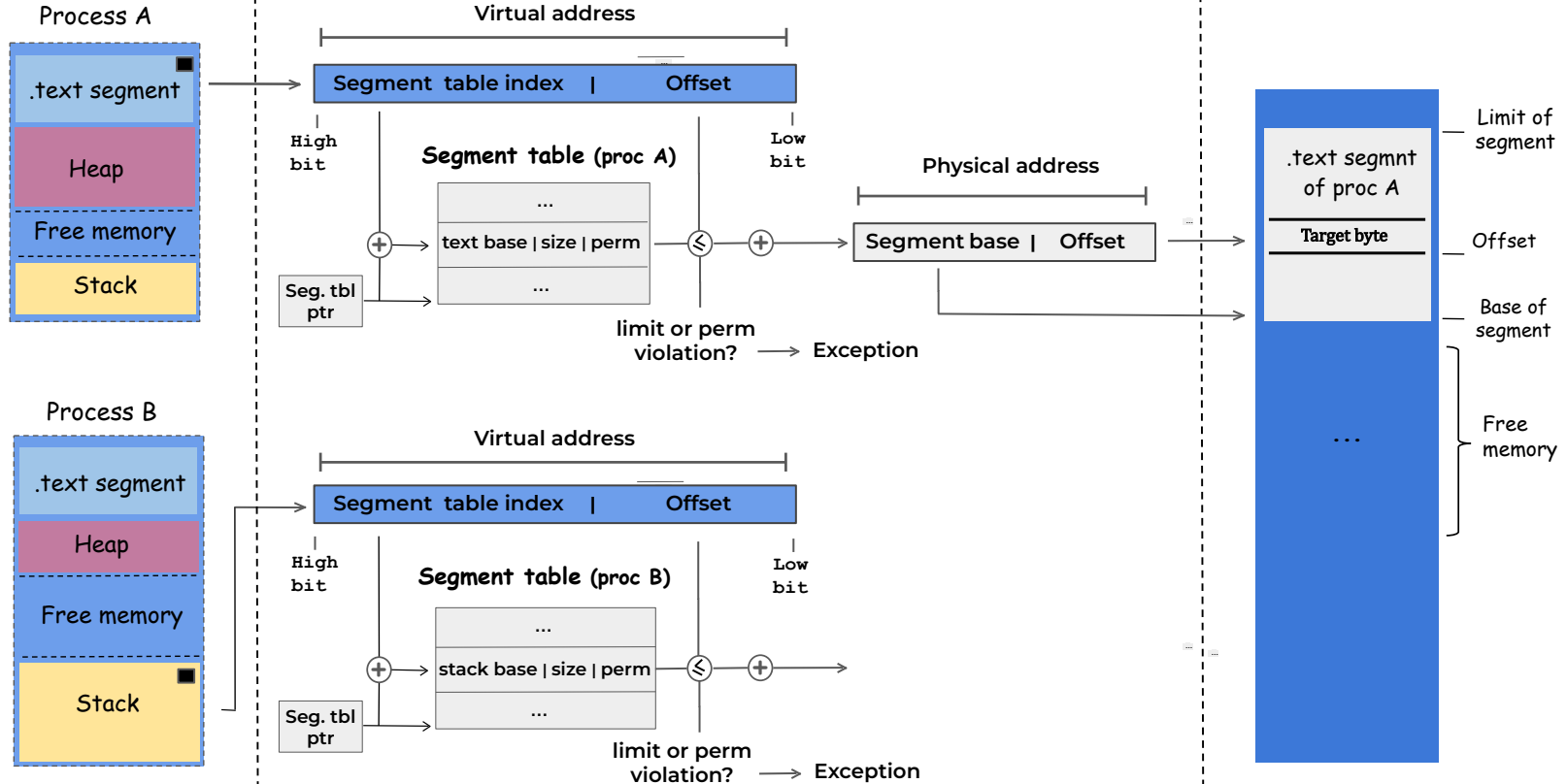
Virtual memory



Mechanism: Segmentation

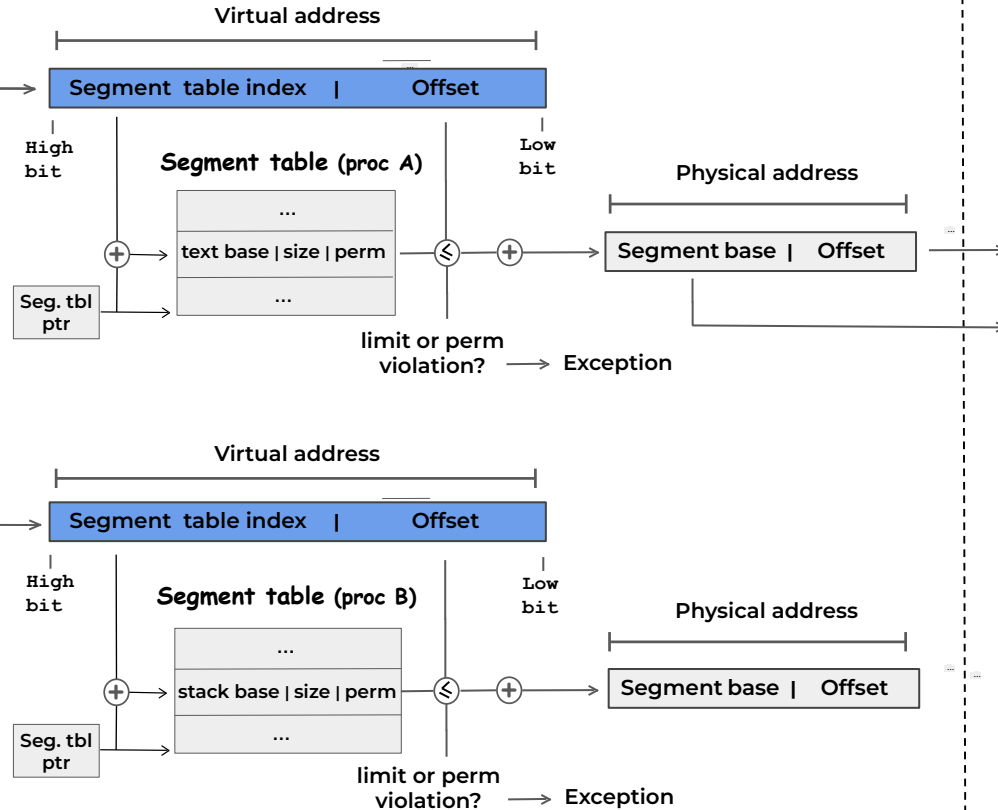
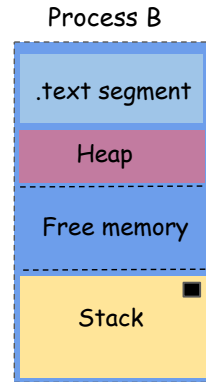
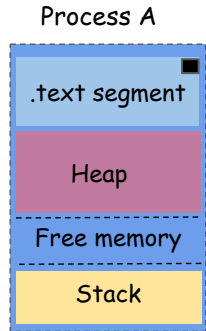
Virtual memory

Physical memory

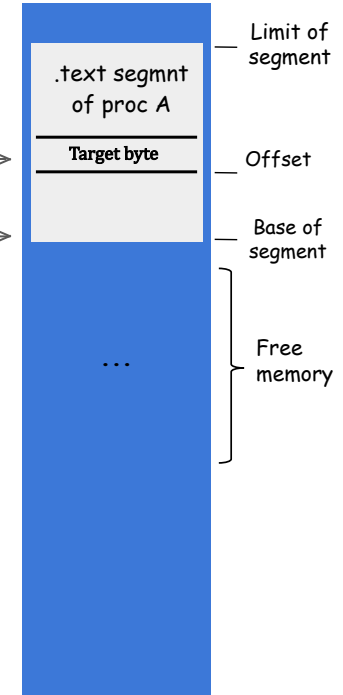


Mechanism: Segmentation

Virtual memory

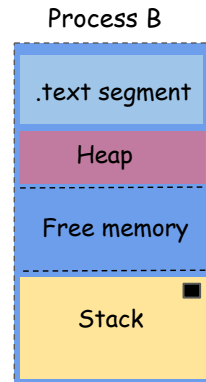
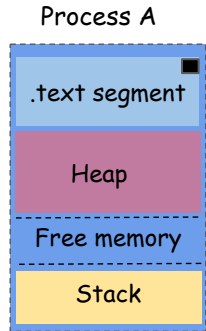


Physical memory

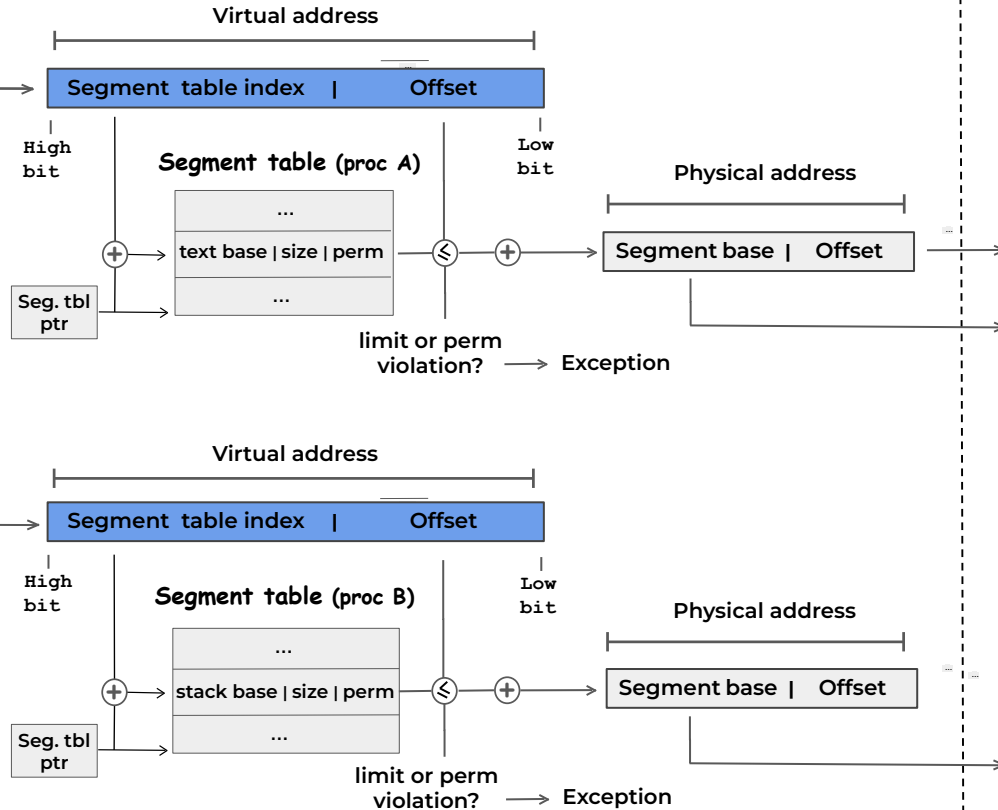
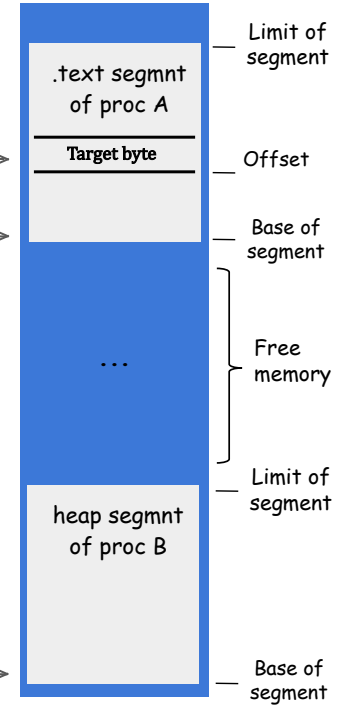


Mechanism: Segmentation

Virtual memory

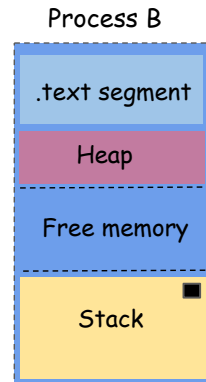
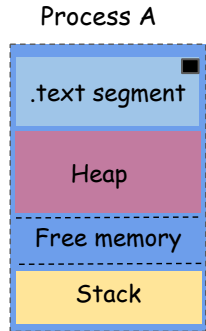


Physical memory

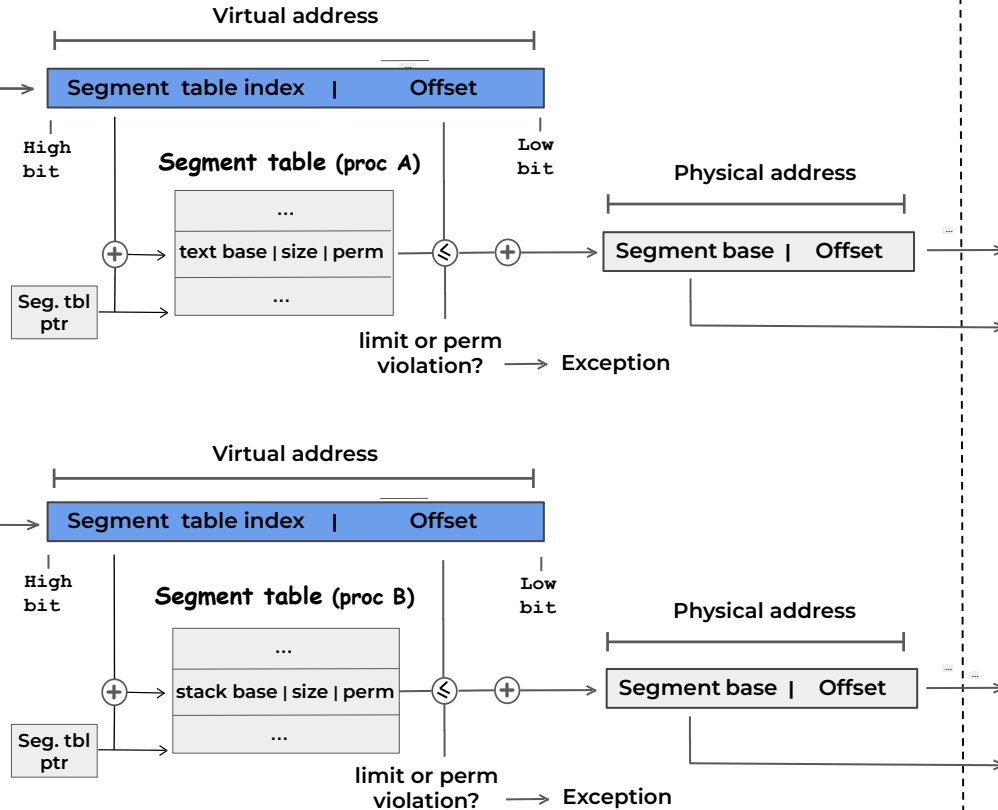
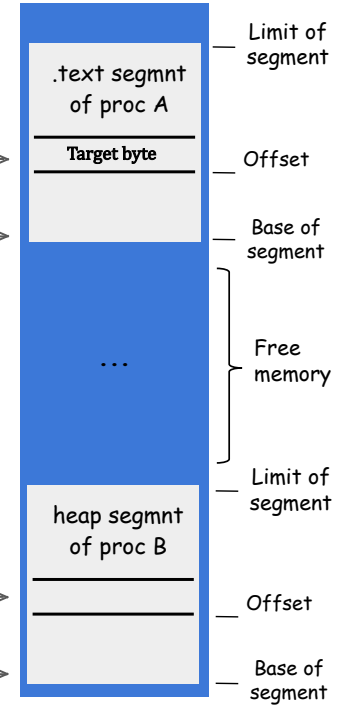


Mechanism: Segmentation

Virtual memory

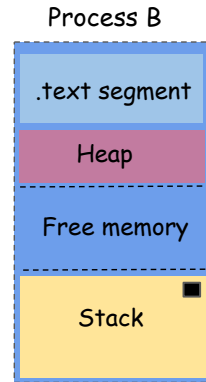
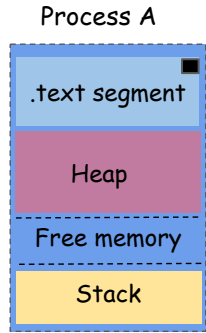


Physical memory

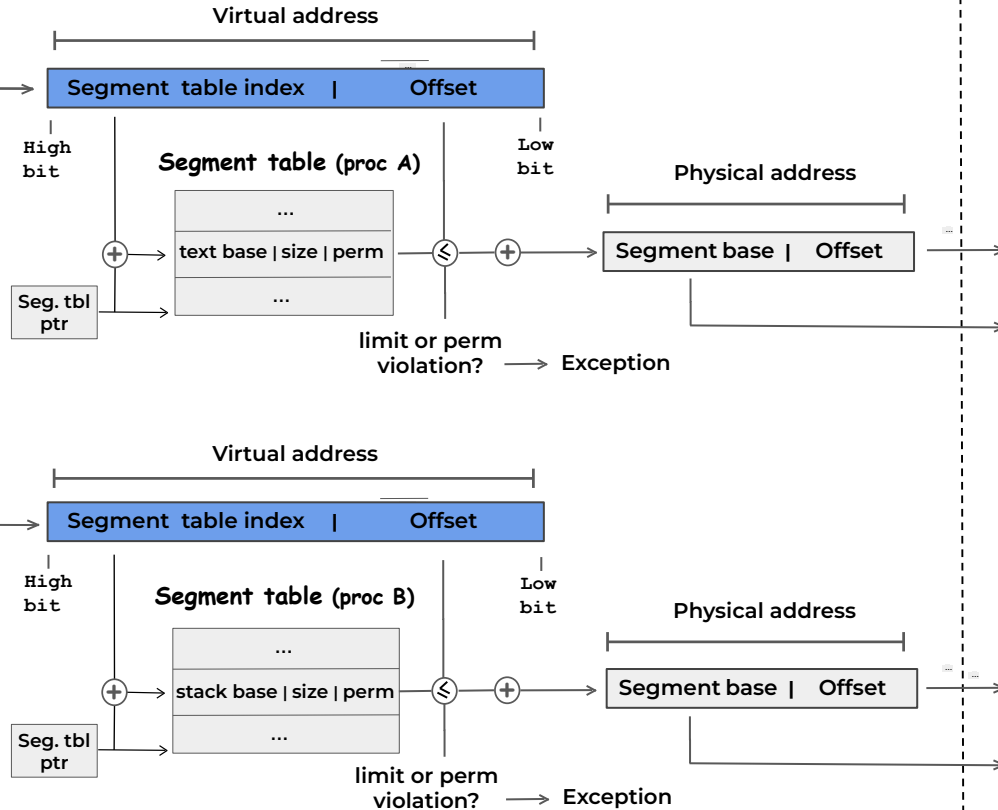
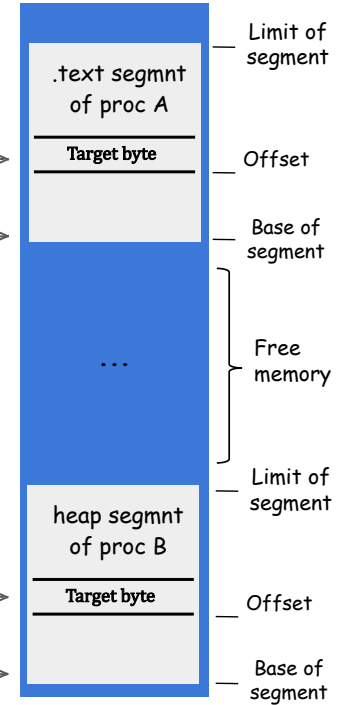


Mechanism: Segmentation

Virtual memory



Physical memory



Mechanism: Segmentation

- › Assume 32-bit virtual and physical address space
w/ 8-bits for the segment table index (segment selector) and 24-bits for offset

Q1: How many segments can a process have?

Mechanism: Segmentation

- › Assume 32-bit virtual and physical address space
w/ 8-bits for the segment table index (segment selector) and 24-bits for offset

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Mechanism: Segmentation

- › Assume 32-bit virtual and physical address space
w/ 8-bits for the segment table index (segment selector) and 24-bits for offset

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

Mechanism: Segmentation

- › Assume 32-bit virtual and physical address space
w/ 8-bits for the segment table index (segment selector) and 24-bits for offset

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 \times (2^{20}) = 16$ MiB

Mechanism: Segmentation

- Assume **32-bit** virtual and physical address space
w/ **8-bits** for the segment table index (segment selector) and **24-bits** for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- $0x00001030 =$

Mechanism: Segmentation

- Assume **32-bit** virtual and physical address space
w/ **8-bits** for the segment table index (segment selector) and **24-bits** for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- $0x00001030 = [0000\ 0000] [0000\ 0000] [0001\ 0000] [0011\ 0000]$

Mechanism: Segmentation

- Assume **32-bit** virtual and physical address space
w/ **8-bits** for the segment table index (segment selector) and **24-bits** for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- $0x00001030 = [0000\ 0000] [0000\ 0000] [0001\ 0000] [0011\ 0000]$
= $[0100\ 1111] [0000\ 0000] [0001\ 0000] [0011\ 0000]$

Mechanism: Segmentation

- Assume 32-bit virtual and physical address space
w/ 8-bits for the segment table index (segment selector) and 24-bits for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- $0x00001030 = [0000\ 0000] [0000\ 0000] [0001\ 0000] [0011\ 0000]$
= $[0100\ 1111] [0000\ 0000] [0001\ 0000] [0011\ 0000]$
= $0x4f\ 00\ 10\ 30 \Rightarrow 00\ 10\ 30 < 9f\ ff\ ff ?$ OK
= $0x4f001030$

Mechanism: Segmentation

- Assume **32-bit** virtual and physical address space
w/ **8-bits** for the segment table index (segment selector) and **24-bits** for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- 0x02a20357 = [0000 0010] [1010 0010] [0000 0011] [0101 1111]

Mechanism: Segmentation

- Assume **32-bit** virtual and physical address space
w/ **8-bits** for the segment table index (segment selector) and **24-bits** for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- 0x02a20357 = [0000 0010] [1010 0010] [0000 0011] [0101 1111]
= [0011 0000] [1010 0010] [0000 0011] [0101 1111]

Mechanism: Segmentation

- Assume 32-bit virtual and physical address space
w/ 8-bits for the segment table index (segment selector) and 24-bits for offset

Segment table (process A)

Seg No	Base	Size	Perm
0	0x4f00000	0xffffffff	rw-
1	0x0100000	0x1ffffff	r-x
2	0x3000000	0x9ffffff	r--

Q1: How many segments can a process have?

- $2^8 = 64$ segments

Q2: How large can a segment be in bytes?

- 2^{24} bytes = $16 * (2^{20}) = 16$ MiB

Q3: What's the physical address for virtual address?

- $0x02a20357 = [0000\ 0010] [1010\ 0010] [0000\ 0011] [0101\ 1111]$
= $[0011\ 0000] [1010\ 0010] [0000\ 0011] [0101\ 1111]$
= $0x30\ a2\ 03\ 5f \Rightarrow a2\ 03\ 5f < 9f\ ff\ ff ?$ NOT_OK

Mechanism: Segmentation

Limitations of Segmentation

Mechanism: Segmentation

Limitations of Segmentation

> **Fragmentation:** Inability to use available memory

Mechanism: Segmentation

Limitations of Segmentation

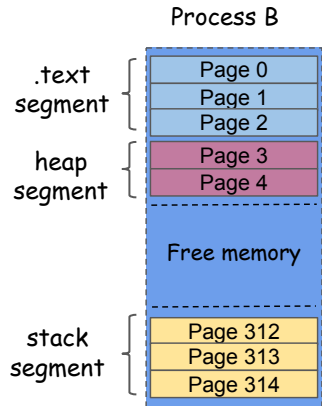
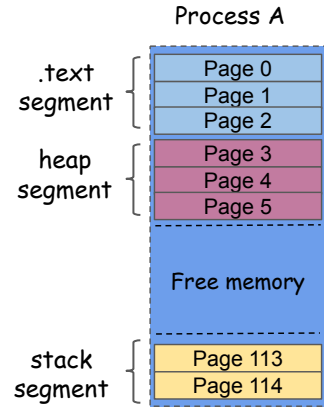
- > **Fragmentation:** Inability to use available memory
 - >> **Internal fragmentation:** Unused portions internally on each segment

Mechanism: Segmentation

Limitations of Segmentation

- > **Fragmentation:** Inability to use available memory
 - >> **Internal fragmentation:** Unused portions internally on each segment
 - >> **External fragmentation:** Free segments not usable due to unfit sizes

Virtual memory

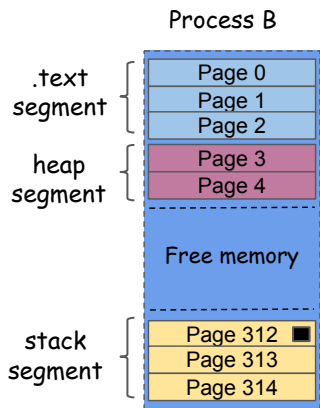
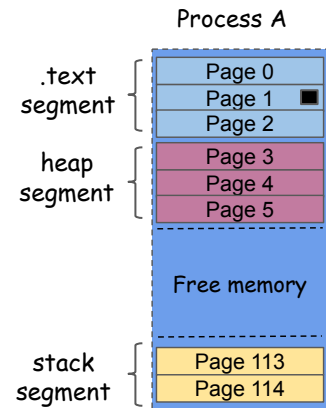


Mechanism: Paging

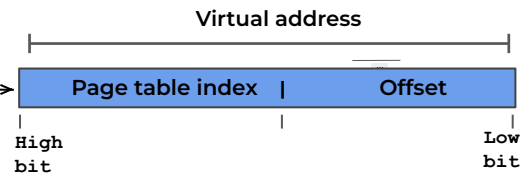
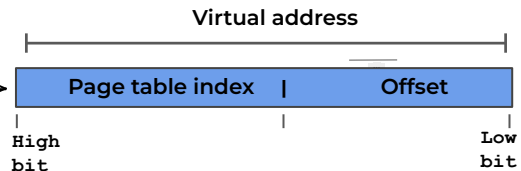
Physical memory



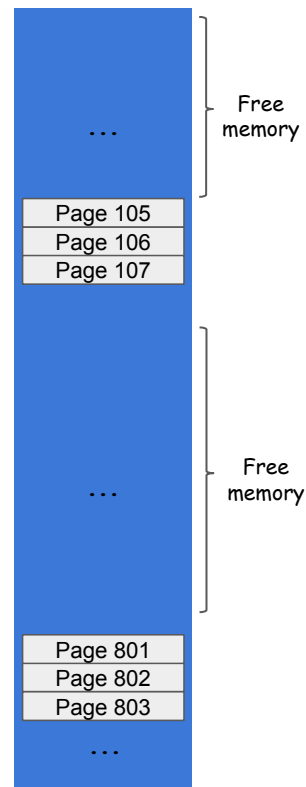
Virtual memory



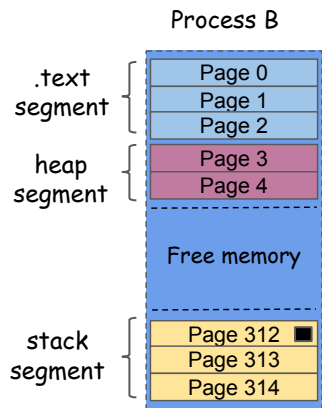
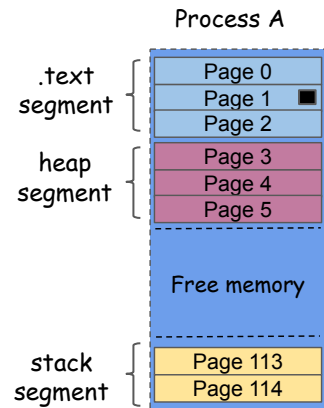
Mechanism: Paging



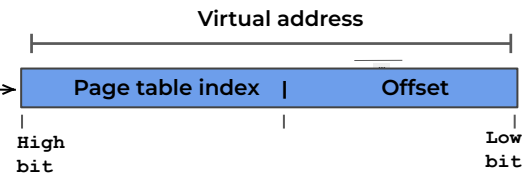
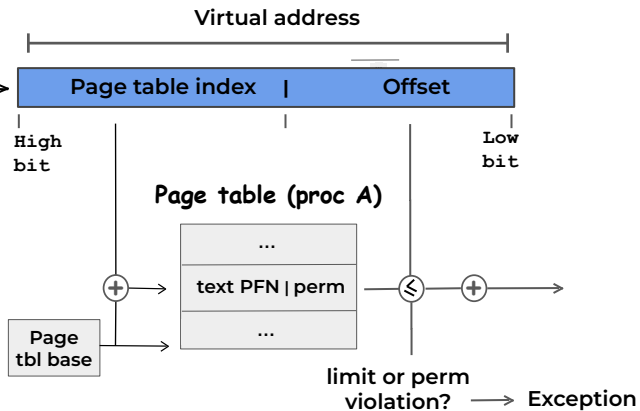
Physical memory



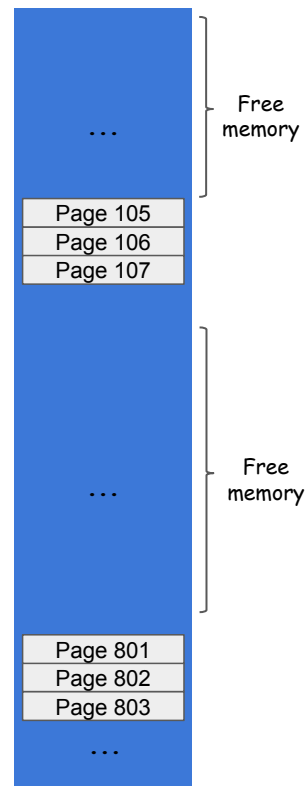
Virtual memory



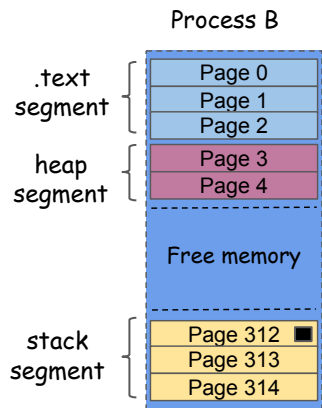
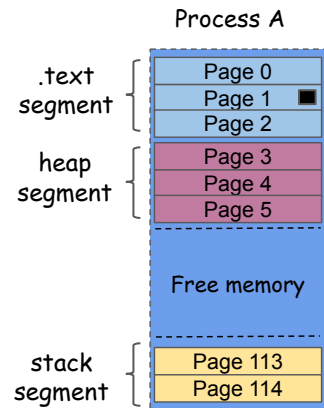
Mechanism: Paging



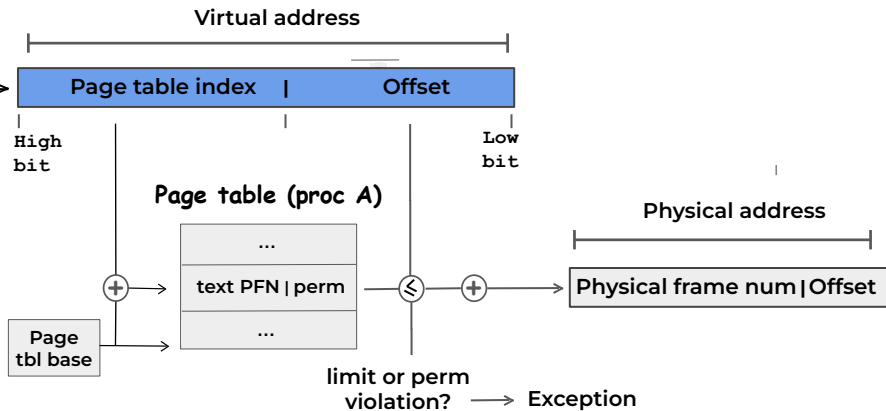
Physical memory



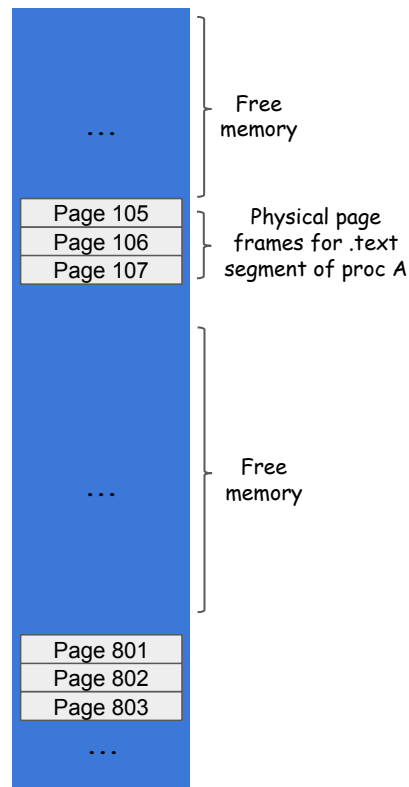
Virtual memory



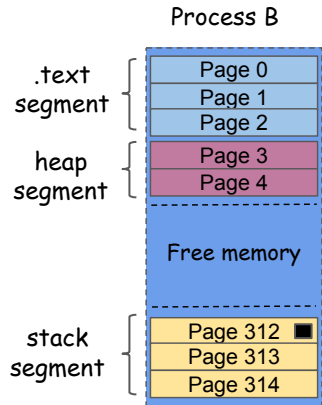
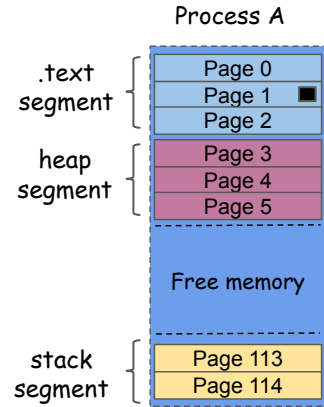
Mechanism: Paging



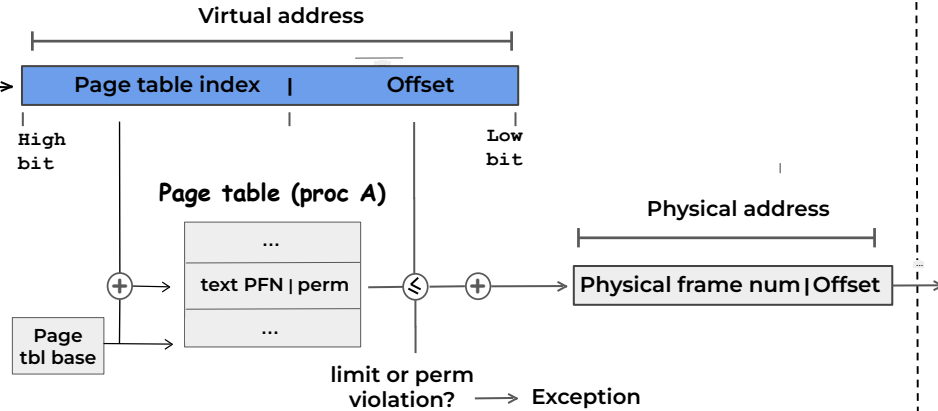
Physical memory



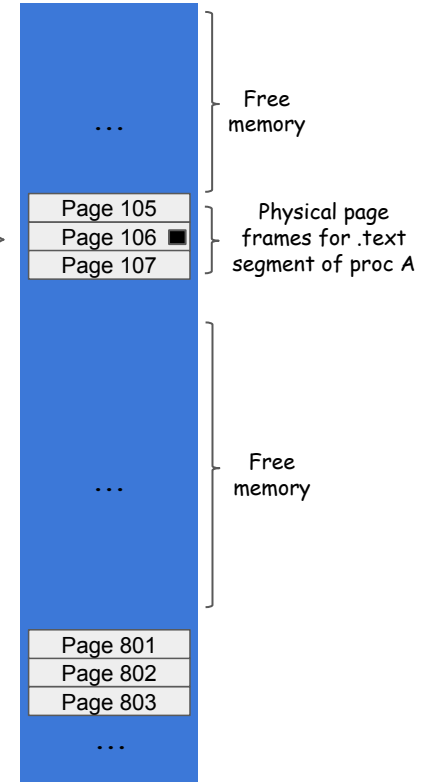
Virtual memory



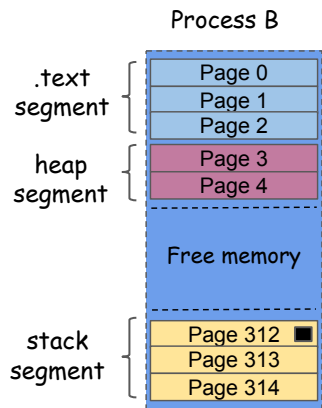
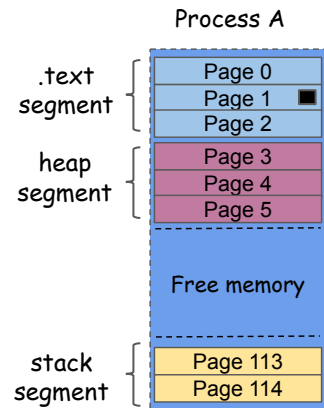
Mechanism: Paging



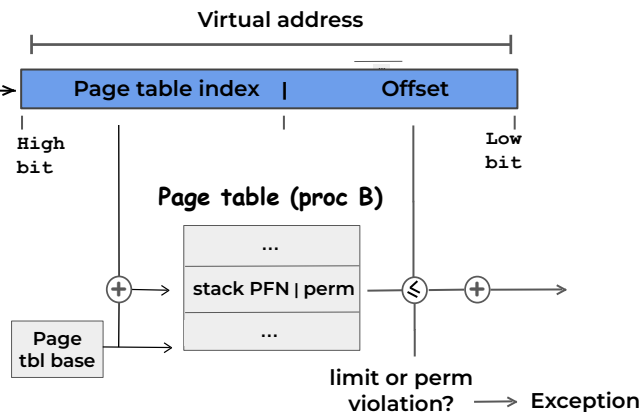
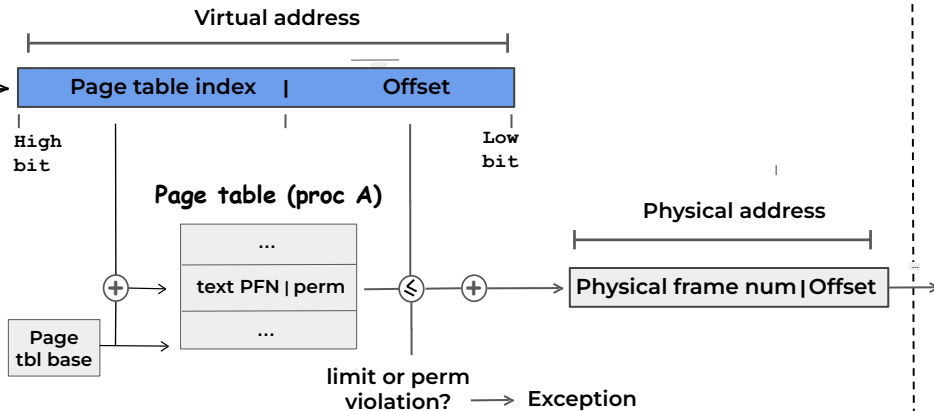
Physical memory



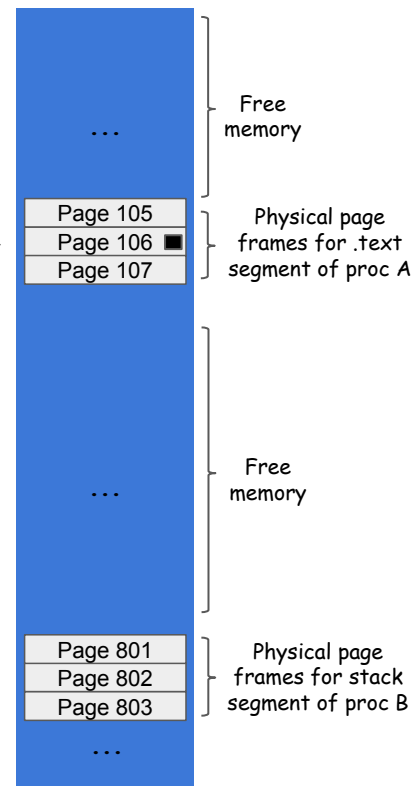
Virtual memory



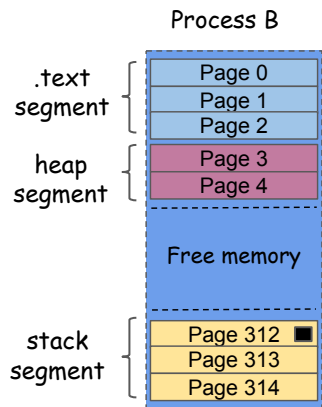
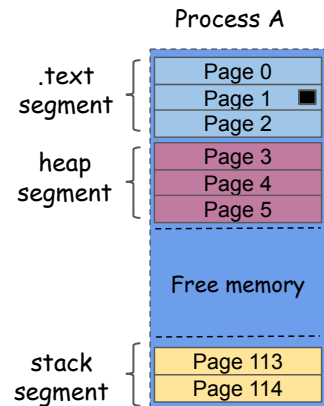
Mechanism: Paging



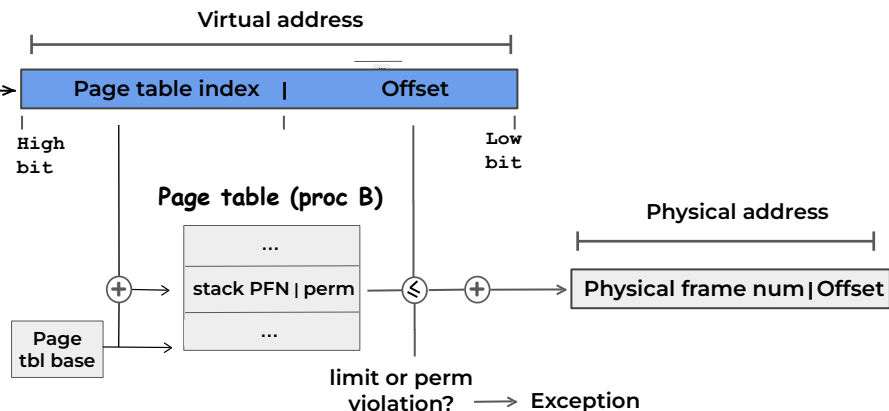
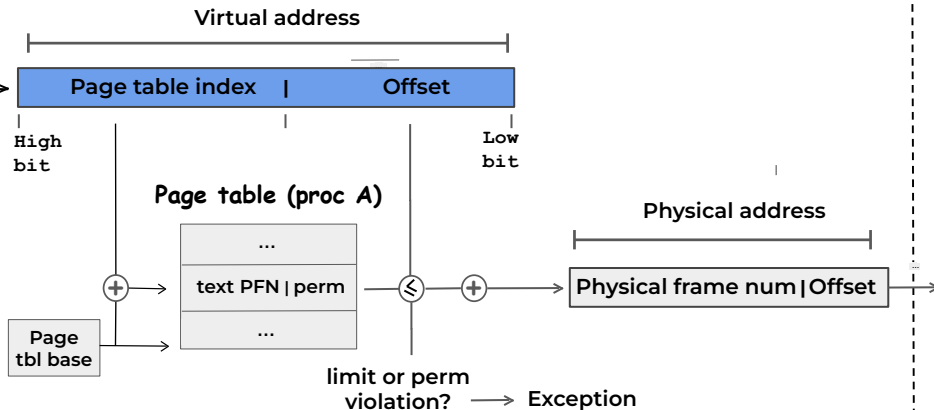
Physical memory



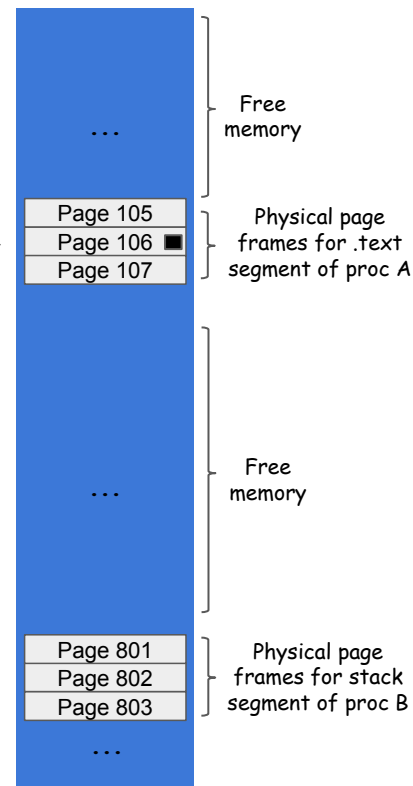
Virtual memory



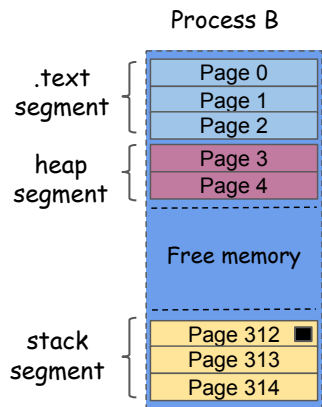
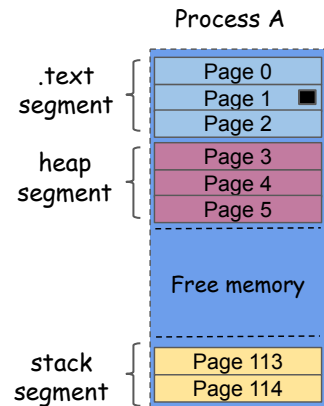
Mechanism: Paging



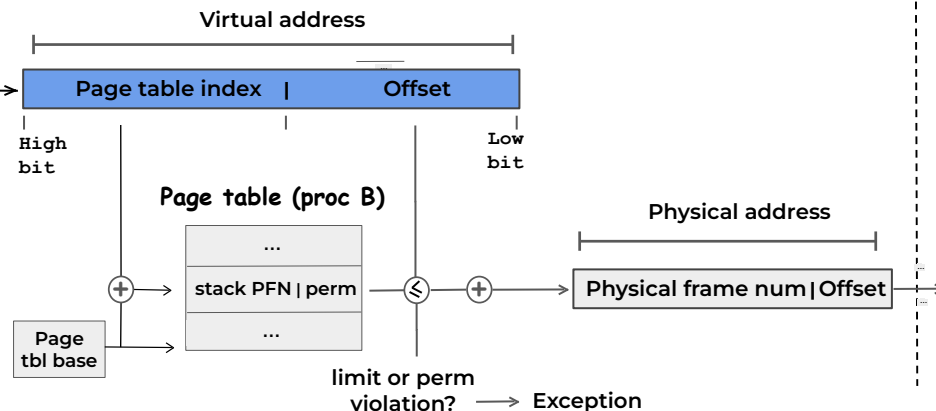
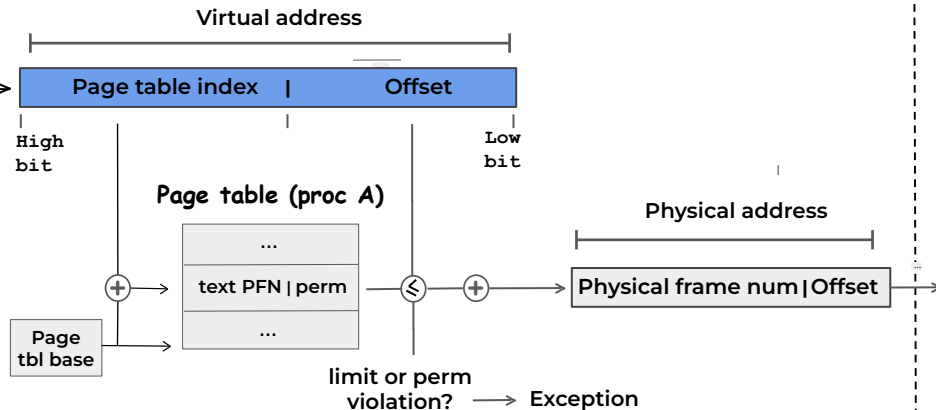
Physical memory



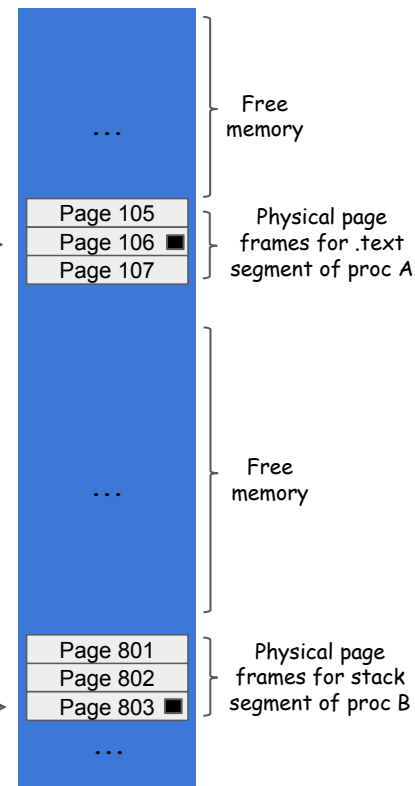
Virtual memory



Mechanism: Paging



Physical memory



Mechanism: Paging

➤ Assume **32-bit** virtual and physical address space, **20-bits** for page table index and **4KB pages** (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Mechanism: Paging

➤ Assume **32-bit** virtual and physical address space, **20-bits** for page table index and **4KB pages** (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

- 0x00000400 =

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Mechanism: Paging

➤ Assume **32-bit** virtual and physical address space, **20-bits** for page table index and **4KB pages** (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

- 0x00000400 = [0000 0000] [0000 0000] [0000 0100] [0000 0000]

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Mechanism: Paging

- Assume 32-bit virtual and physical address space, 20-bits for page table index and 4KB pages (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

$$\begin{aligned} - 0x00000400 &= [0000\ 0000] [0000\ 0000] [0000\ 0100] [0000\ 0000] \\ &= [0000\ 0000] [0001\ 0010] [1110\ 0100] [0000\ 0000] \end{aligned}$$

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Mechanism: Paging

➤ Assume **32-bit** virtual and physical address space, **20-bits** for page table index and **4KB pages** (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

$$\begin{aligned} - 0x00000400 &= [0000\ 0000] [0000\ 0000] [0000\ 0100] [0000\ 0000] \\ &= [0000\ 0000] [0001\ 0010] [1110\ 0100] [0000\ 0000] \\ &= 0x00\ 12\ e4\ 00 = 0x0012e400 \end{aligned}$$

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Mechanism: Paging

➤ Assume **32-bit** virtual and physical address space, **20-bits** for page table index and **4KB pages** (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

- 0x00000400 = [0000 0000] [0000 0000] [0000 0100] [0000 0000]
= [0000 0000] [0001 0010] [1110 0100] [0000 0000]
= 0x00 12 e4 00 = 0x0012e400

- 0x00001402 = [0000 0000] [0000 0000] [0001 0100] [0000 0010]

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Mechanism: Paging

- Assume **32-bit** virtual and physical address space, **20-bits** for page table index and **4KB pages** (i.e., 12 bits for offset)

Q: What's the physical address for virtual address?

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

$$\begin{aligned} - 0x00000400 &= [0000\ 0000] [0000\ 0000] [0000\ 0100] [0000\ 0000] \\ &= [0000\ 0000] [0001\ 0010] [1110\ 0100] [0000\ 0000] \\ &= 0x00\ 12\ e4\ 00 = 0x0012e400 \end{aligned}$$

$$\begin{aligned} - 0x00001402 &= [0000\ 0000] [0000\ 0000] [0001\ 0100] [0000\ 0010] \\ &= [0000\ 0000] [0000\ 0110] [1010\ 0100] [0000\ 0010] \end{aligned}$$

Mechanism: Paging

- Assume 32-bit virtual and physical address space, 20-bits for page table index and 4KB pages (i.e., 12 bits for offset)

Page table (process A)

Index	PFN	Perm
0	302	rw-
1	106	r-x
2	19	r--

Q: What's the physical address for virtual address?

$$\begin{aligned} - 0x00000400 &= [0000\ 0000] [0000\ 0000] [0000\ 0100] [0000\ 0000] \\ &= [0000\ 0000] [0001\ 0010] [1110\ 0100] [0000\ 0000] \\ &= 0x00\ 12\ e4\ 00 = 0x0012e400 \end{aligned}$$

$$\begin{aligned} - 0x00001402 &= [0000\ 0000] [0000\ 0000] [0001\ 0100] [0000\ 0010] \\ &= [0000\ 0000] [0000\ 0110] [1010\ 0100] [0000\ 0010] \\ &= 0x00\ 06\ a4\ 02 = 0x0006a402 \end{aligned}$$

Mechanism: Paging

Limitations of segmentation

- › Fragmentation: Inability to use available memory
 - ›› Internal fragmentation: Unused portions internally on each segment
 - ›› External fragmentation: Free segments not usable due to unfit sizes

Limitations of single-level paging

Mechanism: Paging

Limitations of segmentation

- › Fragmentation: Inability to use available memory
 - ›› Internal fragmentation: Unused portions internally on each segment
 - ›› External fragmentation: Free segments not usable due to unfit sizes

Limitations of single-level paging

- › Size of page table: E.g., with 32-bit virtual addresses, need a 4MB flat array per process to map a sparse 4GB address space

Mechanism: Paging

Limitations of segmentation

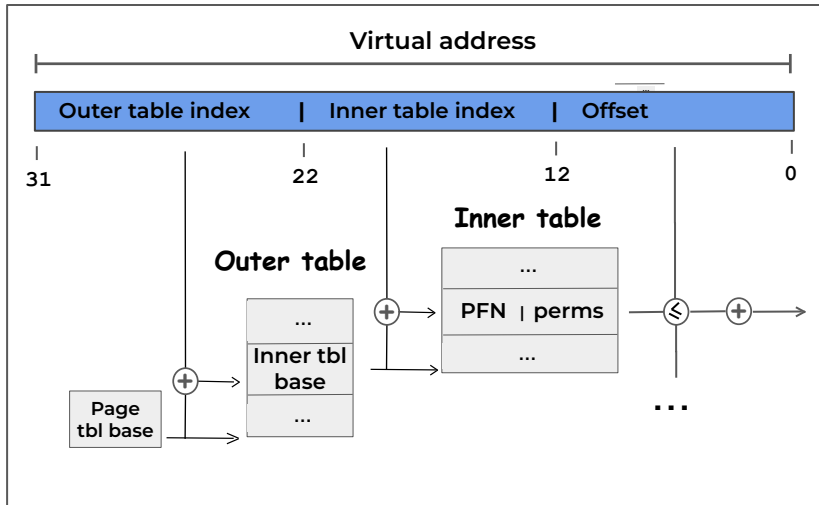
- › Fragmentation: Inability to use available memory
 - ›› Internal fragmentation: Unused portions internally on each segment
 - ›› External fragmentation: Free segments not usable due to unfit sizes

Limitations of single-level paging

- › Size of page table: E.g., with 32-bit virtual addresses, need a 4MB flat array per process to map a sparse 4GB address space
- › Internal fragmentation: Still an issue, but not as prominent

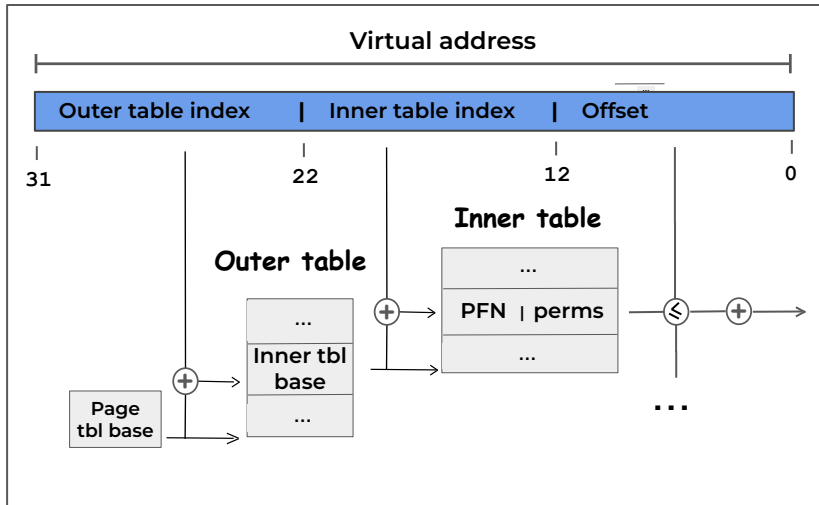
Mechanism: 2-level hierarchical paging

- > Assume **32-bit** virtual and physical address space, where the **12 lower bits** are for the offset, the **10 middle bits** are the index for the inner table, and the **10 higher bits** are the index for the outer table



Mechanism: 2-level hierarchical paging

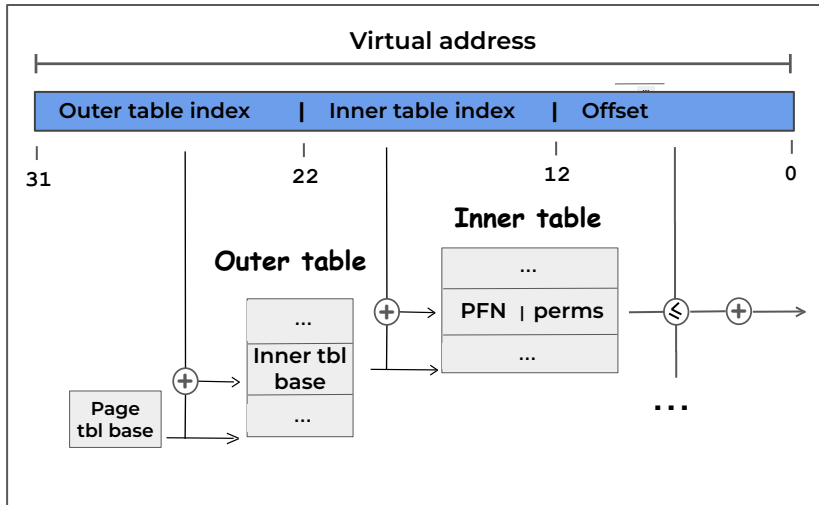
- > Assume **32-bit** virtual and physical address space, where the **12 lower bits** are for the offset, the **10 middle bits** are the index for the inner table, and the **10 higher bits** are the index for the outer table



Q: Size of "map-able" address space?

Mechanism: 2-level hierarchical paging

- > Assume **32-bit** virtual and physical address space, where the **12 lower bits** are for the offset, the **10 middle bits** are the index for the inner table, and the **10 higher bits** are the index for the outer table

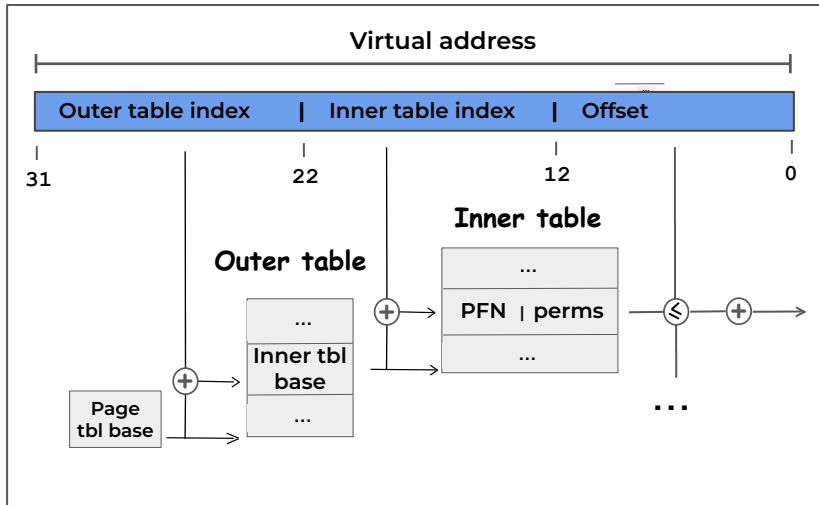


Q: Size of "map-able" address space?

- > **With two pages** (one for each of the outer and inner table): $1 * 1024 \text{ PTEs} * 4\text{KB} = 4\text{MiB}$

Mechanism: 2-level hierarchical paging

- > Assume **32-bit** virtual and physical address space, where the **12 lower bits** are for the offset, the **10 middle bits** are the index for the inner table, and the **10 higher bits** are the index for the outer table

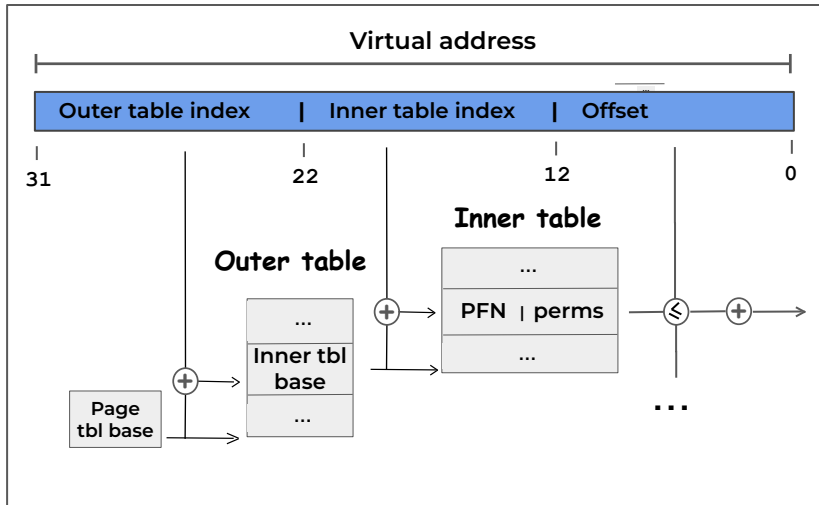


Q: Size of "map-able" address space?

- > **With two pages** (one for each of the outer and inner table): $1 * 1024 \text{ PTEs} * 4\text{KB} = 4\text{MiB}$
- > **With each additional page in the inner table** (i.e., 1024 PTEs), 4 additional MiB can be mapped

Mechanism: 2-level hierarchical paging

- > Assume **32-bit** virtual and physical address space, where the **12 lower bits** are for the offset, the **10 middle bits** are the index for the inner table, and the **10 higher bits** are the index for the outer table

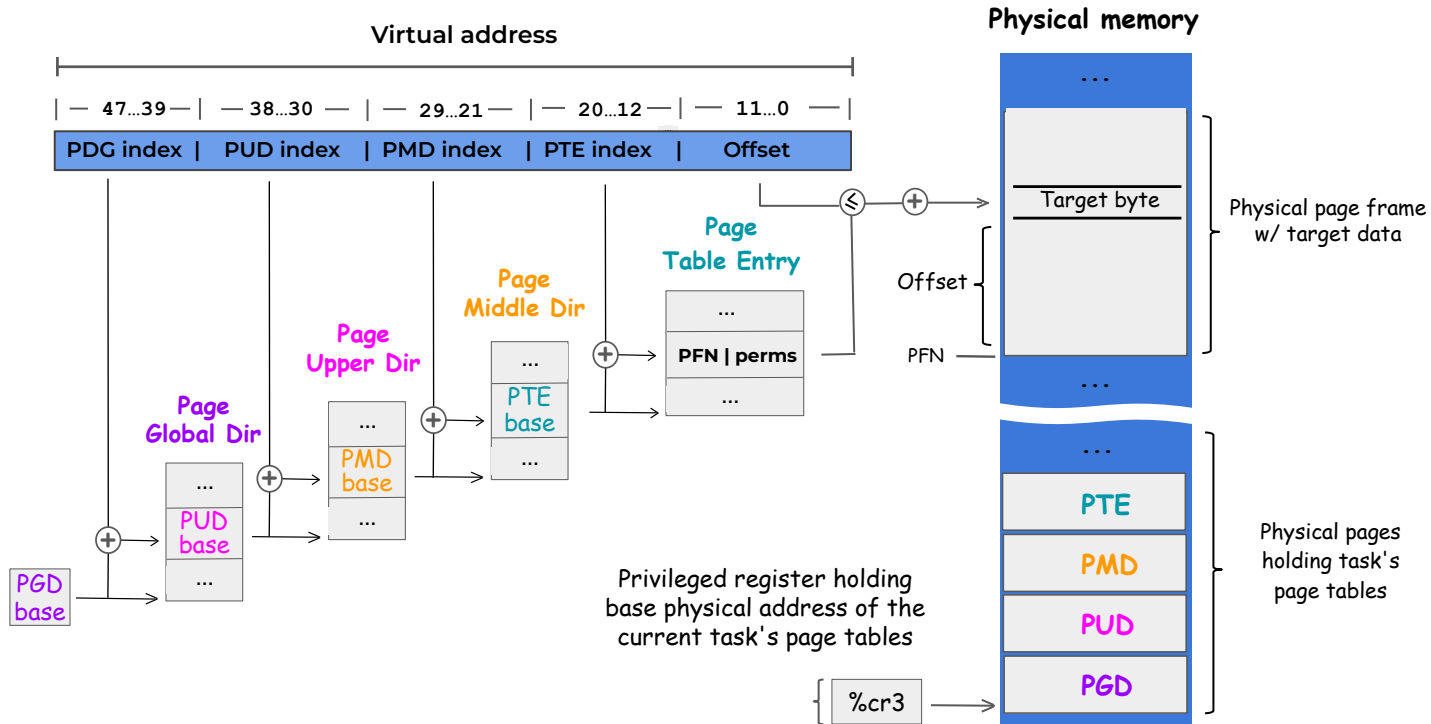


Q: Size of "map-able" address space?

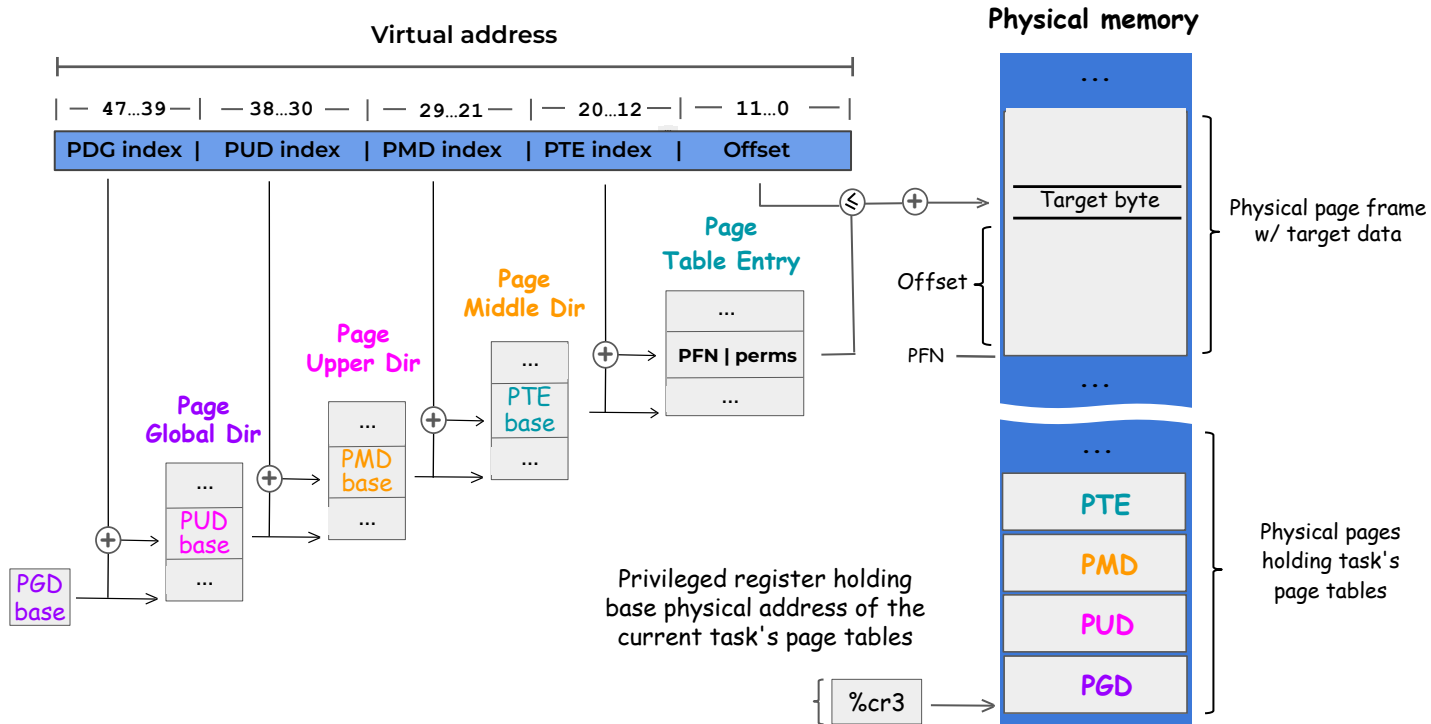
- > **With two pages** (one for each of the outer and inner table): $1 * 1024 \text{ PTEs} * 4\text{KB} = 4\text{MiB}$
- > **With each additional page in the inner table** (i.e., 1024 PTEs), 4 additional MiB can be mapped

There is another positive side-effect...hang on...

Mechanism: Multi-level paging

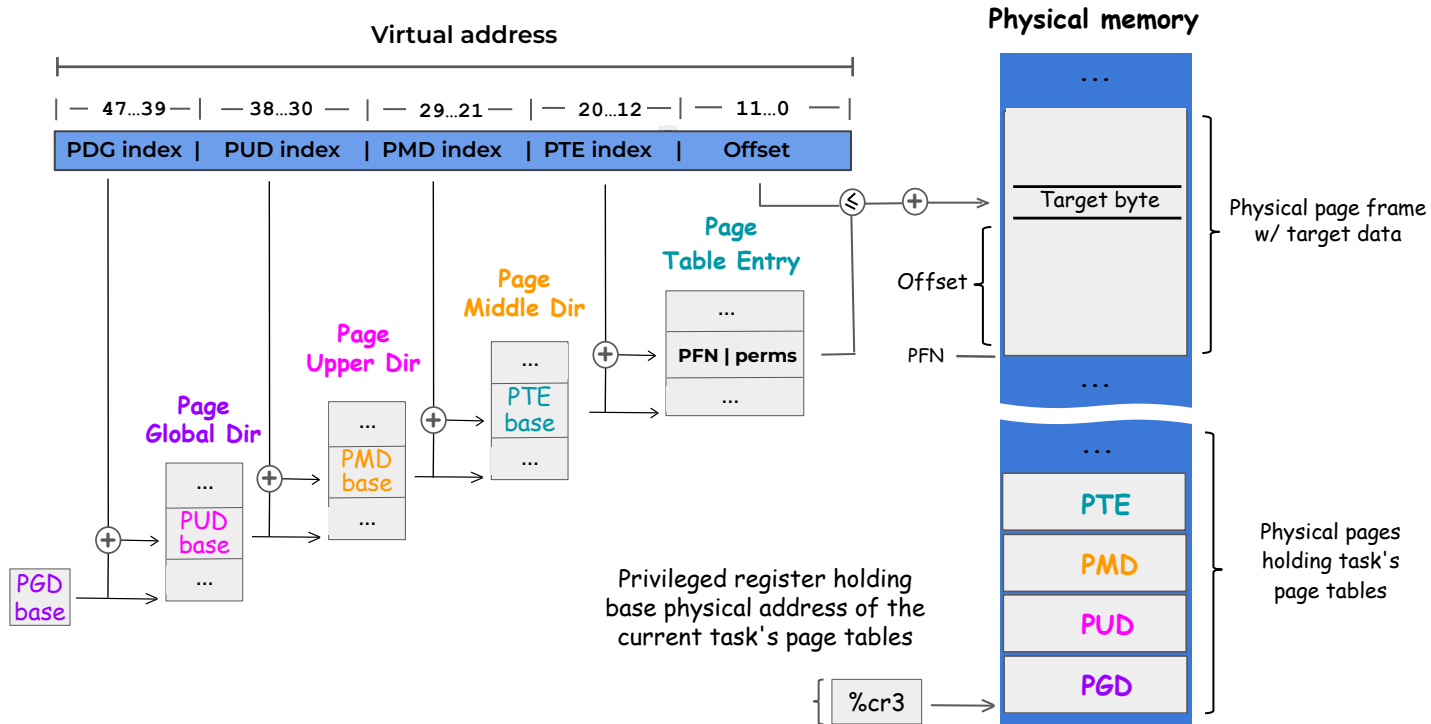


Mechanism: Multi-level paging



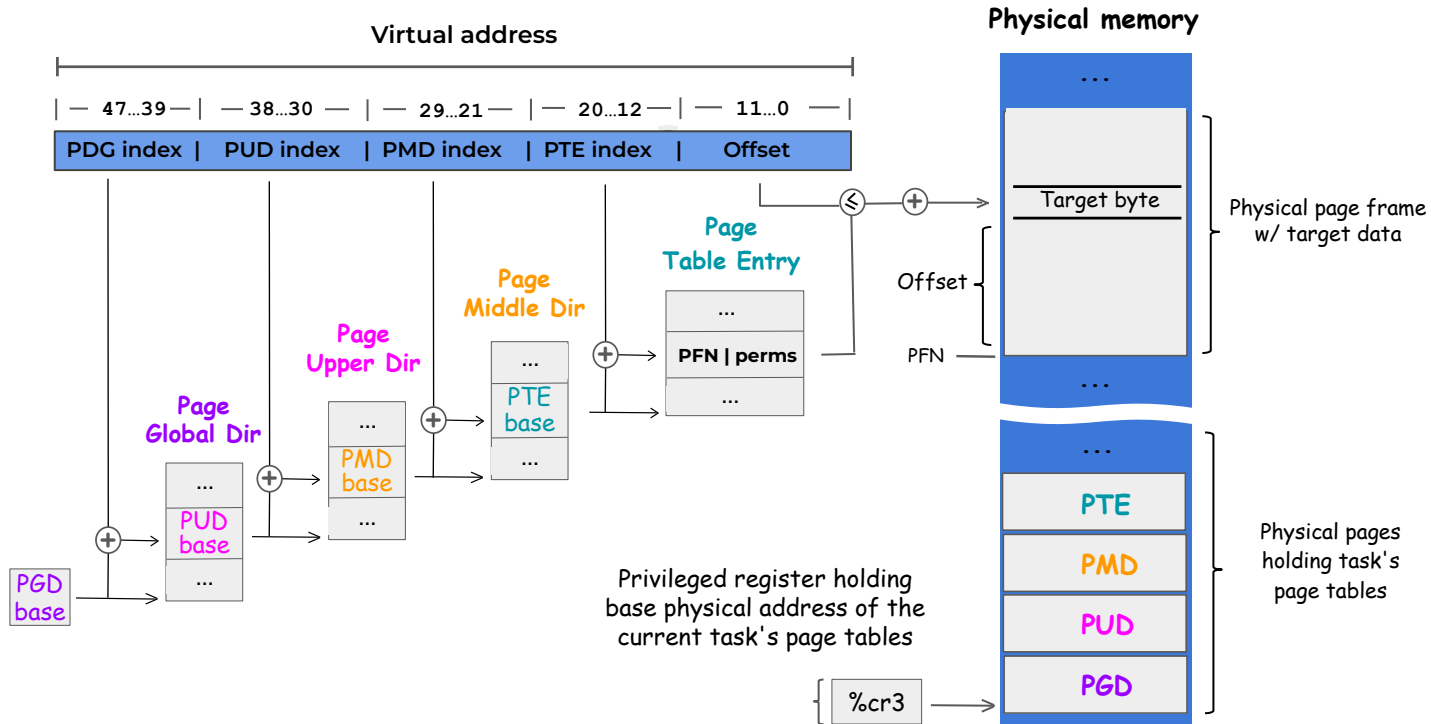
➤ See relevant linux definitions: [here](#)

Mechanism: Multi-level paging



➤ What is the reach of one page of PTE/PMD/PUD/PGD entries?

Mechanism: Multi-level paging



➤ What is the reach of one page of PTE/PMD/PUD/PGD entries? Why care?

x86 PTE format that maps a 4KB page

66665555555555		3210987654321		M ¹ M-1		33332222222222		2109876543210987654321098765432109876543210																				
X D	Prot. Key		Ignored		Rsvd.		Address of 4KB page frame										R	Ign.	G	P A T	D	A	P C D	P W T	U / S	R / W	1	PTE: 4KB page
Ignored																								Q	PTE: not present			
0 (P)		Present; must be 1 to map a 4-KByte page																										
1 (R/W)		Read/write; if 0, writes may not be allowed to the 4-KByte page referenced by this entry (see Section 5.6)																										
2 (U/S)		User/supervisor; if 0, user-mode accesses are not allowed to the 4-KByte page referenced by this entry (see Section 5.6)																										
3 (PWT)		Page-level write-through; indirectly determines the memory type used to access the 4-KByte page referenced by this entry (see Section 5.9.2)																										
4 (PCD)		Page-level cache disable; indirectly determines the memory type used to access the 4-KByte page referenced by this entry (see Section 5.9.2)																										
5 (A)		Accessed; indicates whether software has accessed the 4-KByte page referenced by this entry (see Section 5.8)																										
6 (D)		Dirty; indicates whether software has written to the 4-KByte page referenced by this entry (see Section 5.8)																										

- See [Intel Software Developer's Manual](#)

x86 Page fault error codes @%cr2 (faulting addr. @%cr3)

31		15		7	6	5	4	3	2	1	0
Reserved		SGX	Reserved	HLAT	SS	PK	I/D	RSVD	U/S	W/R	P

- P
- 0 The fault was caused by a non-present page.
 - 1 The fault was caused by a page-level protection violation.
- W/R
- 0 The access causing the fault was a read.
 - 1 The access causing the fault was a write.
- U/S
- 0 A supervisor-mode access caused the fault.
 - 1 A user-mode access caused the fault.
- RSVD
- 0 The fault was not caused by reserved bit violation.
 - 1 The fault was caused by a reserved bit set to 1 in some paging-structure entry.
- I/D
- 0 The fault was not caused by an instruction fetch.
 - 1 The fault was caused by an instruction fetch.
- PK
- 0 The fault was not caused by protection keys.
 - 1 There was a protection-key violation.
- SS
- 0 The fault was not caused by a shadow-stack access.
 - 1 The fault was caused by a shadow-stack access.
- HLAT
- 0 The fault occurred during ordinary paging or due to access rights.
 - 1 The fault occurred during HLAT paging.
- SGX
- 0 The fault is not related to SGX.
 - 1 The fault resulted from violation of SGX-specific access-control requirements.

➤ See [Intel Software Developer's Manual](#)

Checking the address space size of your machine

Checking the address space size of your machine

```
→ ~ cat /proc/$$/maps | head
aaaaab8d0000-aaaaab9a3000 r-xp 00000000 08:02 2375001
aaaaab9b3000-aaaaab9b5000 r--p 000d3000 08:02 2375001
aaaaab9b5000-aaaaab9bb000 rw-p 000d5000 08:02 2375001
aaaaab9bb000-aaaaab9cf000 rw-p 00000000 00:00 0
aaaaac024000-aaaaac40b000 rw-p 00000000 00:00 0
ffffb82b0000-ffffb8530000 r--s 00000000 08:02 2890570
ffffb8530000-ffffb853f000 r-xp 00000000 08:02 2495184
ffffb853f000-ffffb854e000 ---p 0000f000 08:02 2495184
ffffb854e000-ffffb854f000 r--p 0000e000 08:02 2495184
ffffb854f000-ffffb8550000 rw-p 0000f000 08:02 2495184
```

Checking the address space size of your machine

```
→ ~ cat /proc/$$/maps | head
aaaaab8d0000-aaaaab9a3000 r-xp 00000000 08:02 2375001
aaaaab9b3000-aaaaab9b5000 r--p 000d3000 08:02 2375001
aaaaab9b5000-aaaaab9bb000 rw-p 000d5000 08:02 2375001
aaaaab9bb000-aaaaab9cf000 rw-p 00000000 00:00 0
aaaaac024000-aaaaac40b000 rw-p 00000000 00:00 0
ffffb82b0000-ffffb8530000 r--s 00000000 08:02 2890570
ffffb8530000-ffffb853f000 r-xp 00000000 08:02 2495184
ffffb853f000-ffffb854e000 ---p 0000f000 08:02 2495184
ffffb854e000-ffffb854f000 r--p 0000e000 08:02 2495184
ffffb854f000-ffffb8550000 rw-p 0000f000 08:02 2495184
```

> Why am I seeing 12 hexadecimal digits?

Checking the address space size of your machine

```
→ ~ cat /proc/$$/maps | head
aaaaab8d0000-aaaaab9a3000 r-xp 00000000 08:02 2375001
aaaaab9b3000-aaaaab9b5000 r--p 000d3000 08:02 2375001
aaaaab9b5000-aaaaab9bb000 rw-p 000d5000 08:02 2375001
aaaaab9bb000-aaaaab9cf000 rw-p 00000000 00:00 0
aaaaac024000-aaaaac40b000 rw-p 00000000 00:00 0
fffffb82b0000-fffffb8530000 r--s 00000000 08:02 2890570
fffffb8530000-fffffb853f000 r-xp 00000000 08:02 2495184
fffffb853f000-fffffb854e000 ---p 0000f000 08:02 2495184
fffffb854e000-fffffb854f000 r--p 0000e000 08:02 2495184
fffffb854f000-fffffb8550000 rw-p 0000f000 08:02 2495184
→ ~ cat /boot/config-`uname -r` | grep "VA_BITS\|PA_BITS"
# CONFIG_ARM64_VA_BITS_39 is not set
CONFIG_ARM64_VA_BITS_48=y
CONFIG_ARM64_VA_BITS=48
CONFIG_ARM64_PA_BITS_48=y
CONFIG_ARM64_PA_BITS=48
```

How is virtual memory implemented?

- > Usually considered the most complex kernel subsystem
- It's simple if you view it as an index

Ask the right questions

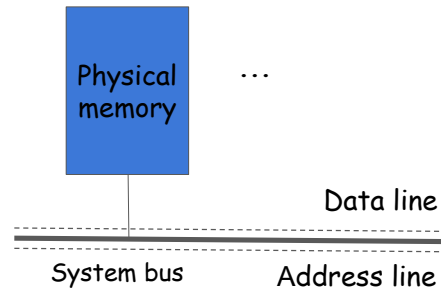
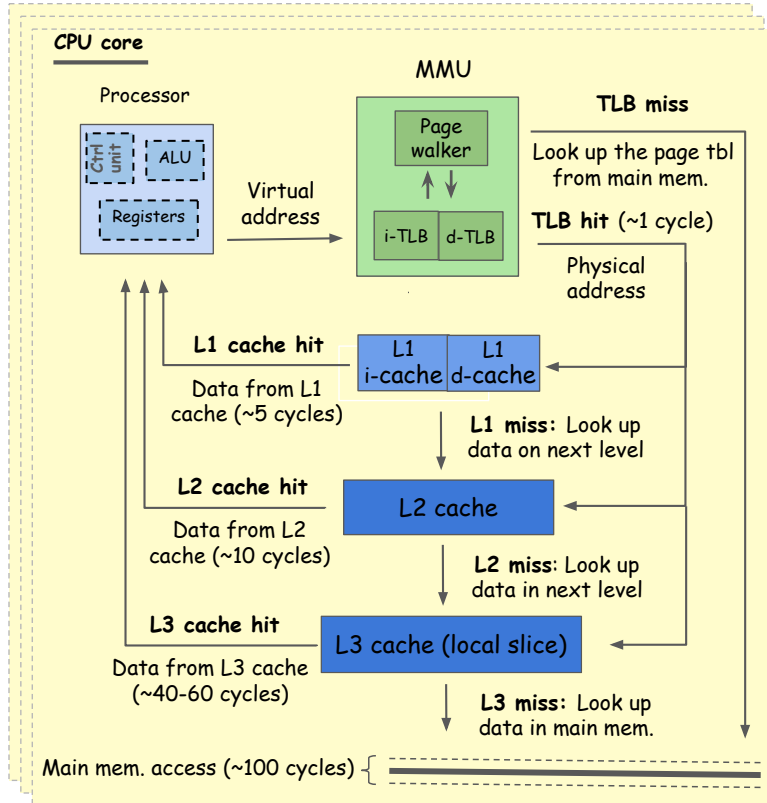
~~Q1: What are the entries of the index?~~

~~Q2: How are the entries of the index used?~~

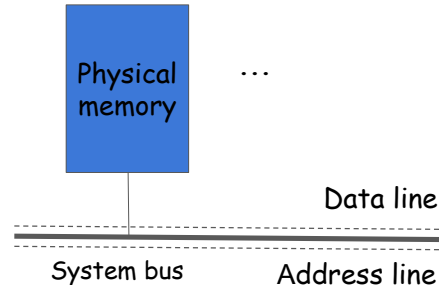
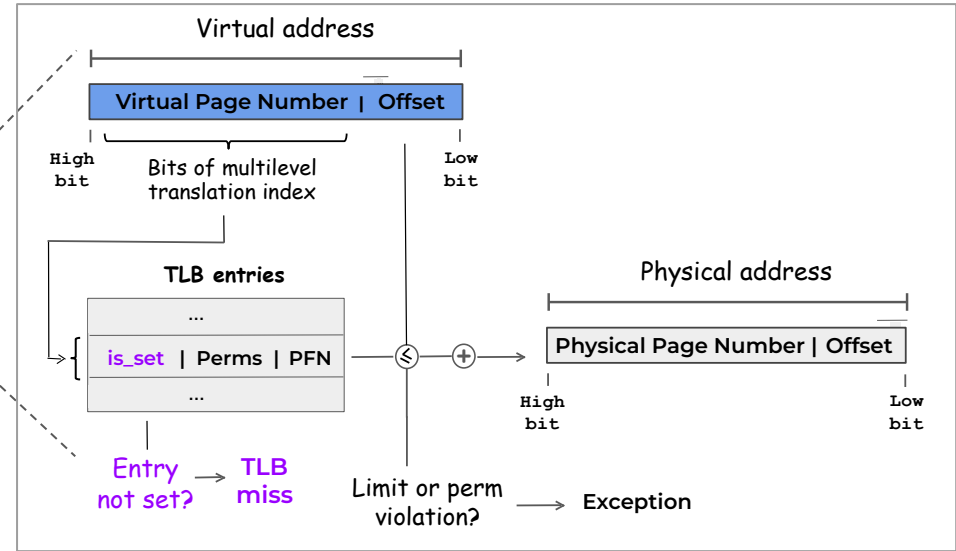
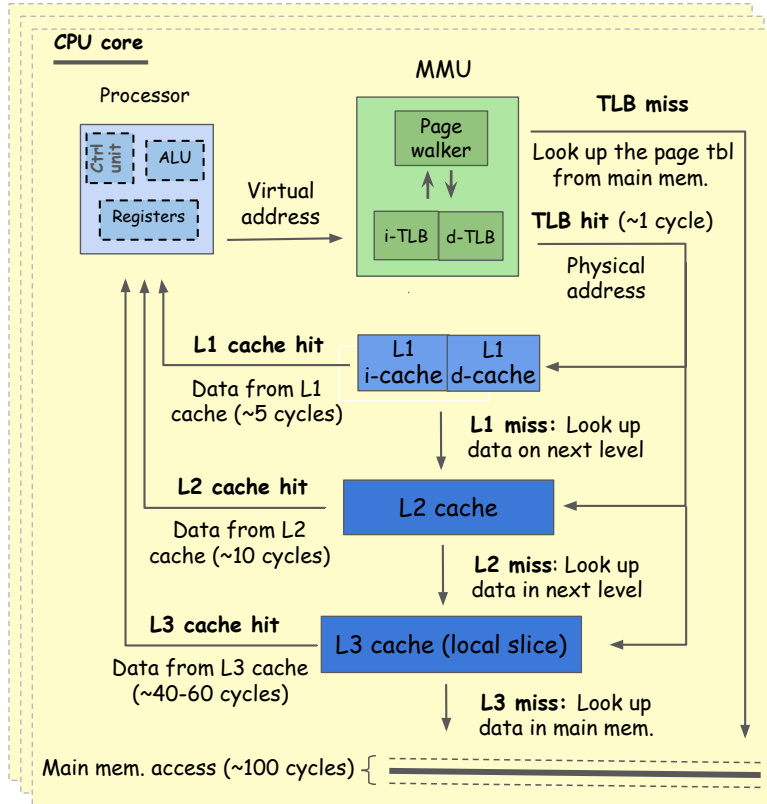
~~Q3: How are the entries of the index allocated?~~

~~Q4: How are the entries of the index replaced?~~

Translation Lookaside Buffer (TLB)

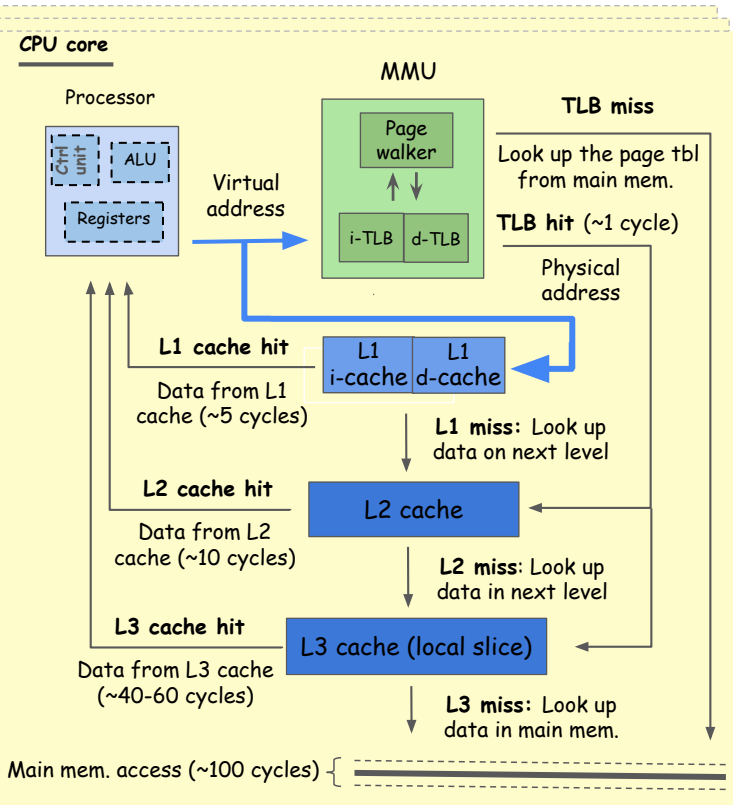


Translation Lookaside Buffer (TLB)

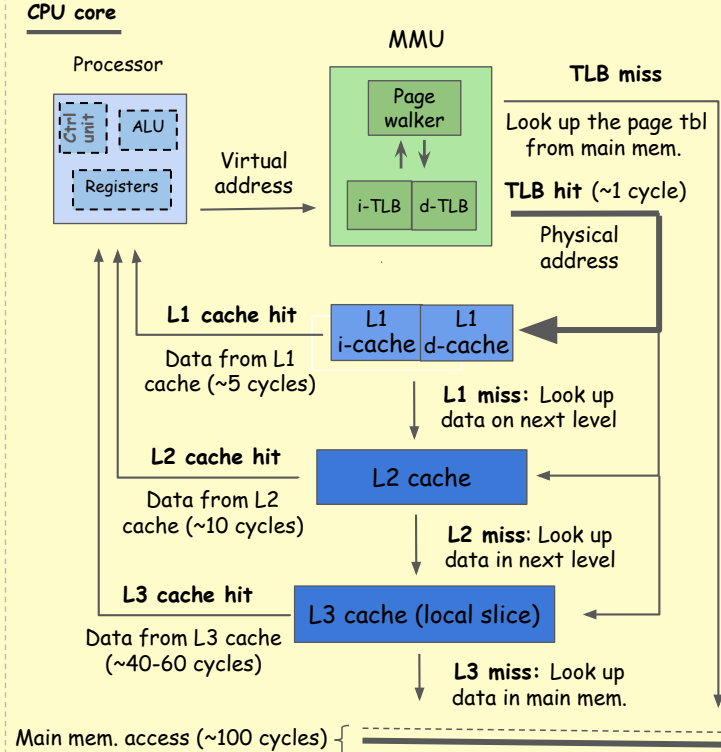
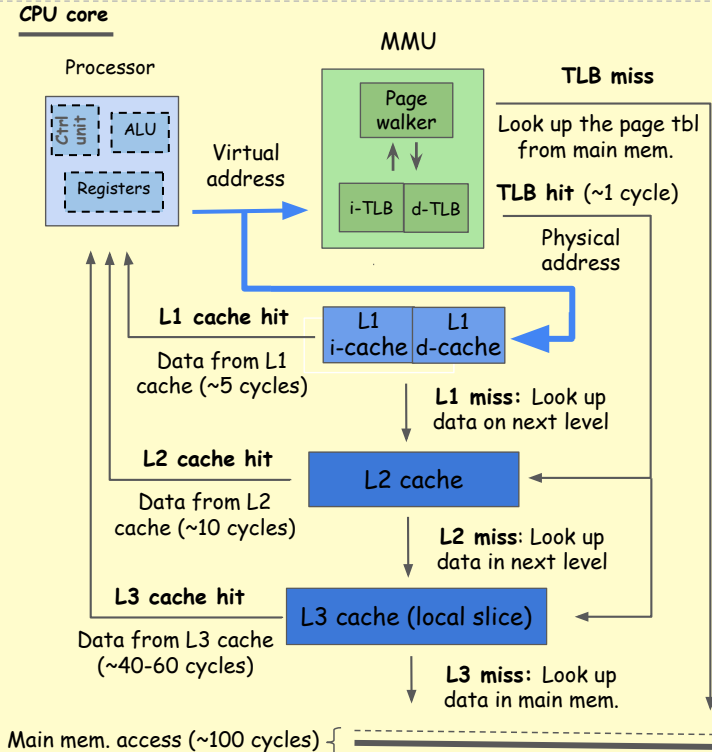


L1 cache: Why not virtually indexed and tagged?

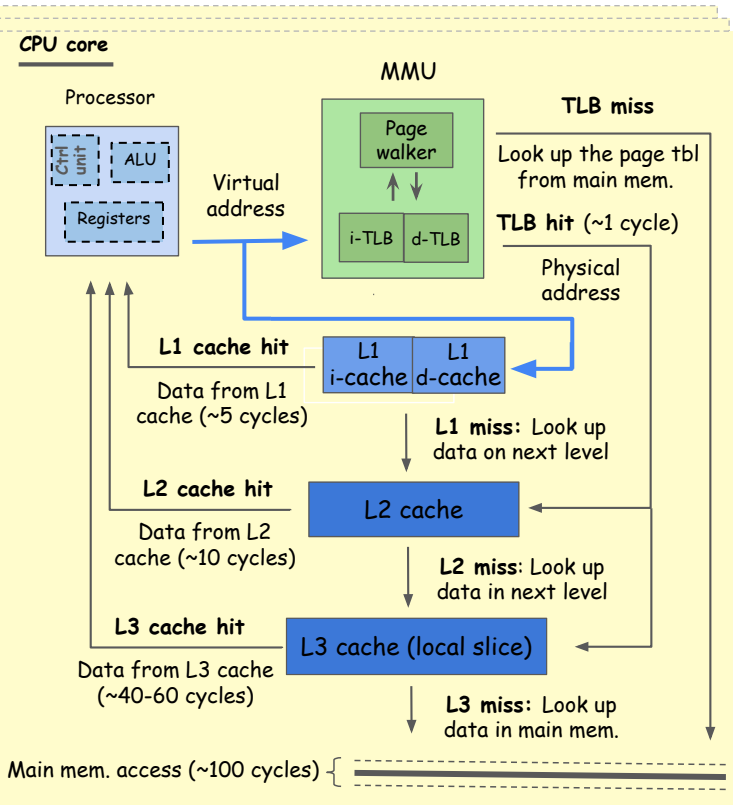
L1 cache: Why not virtually indexed and tagged?



L1 cache: Why not virtually indexed and tagged?

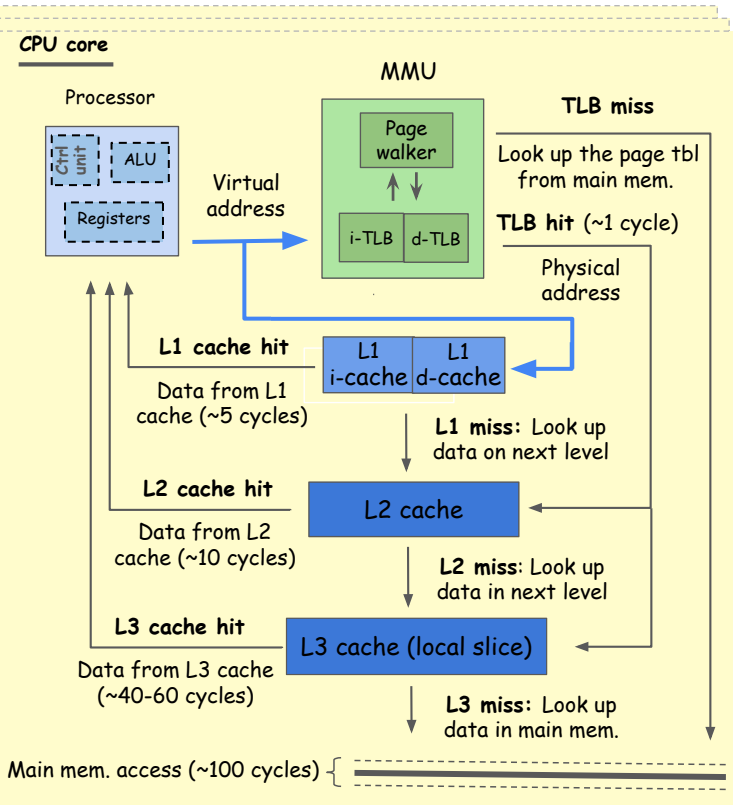


L1 cache: Why not virtually indexed and tagged?



Homonym problem: After a context switch a proc can access leftover cache contents not mapped to its virtual address space (**data leakage**)

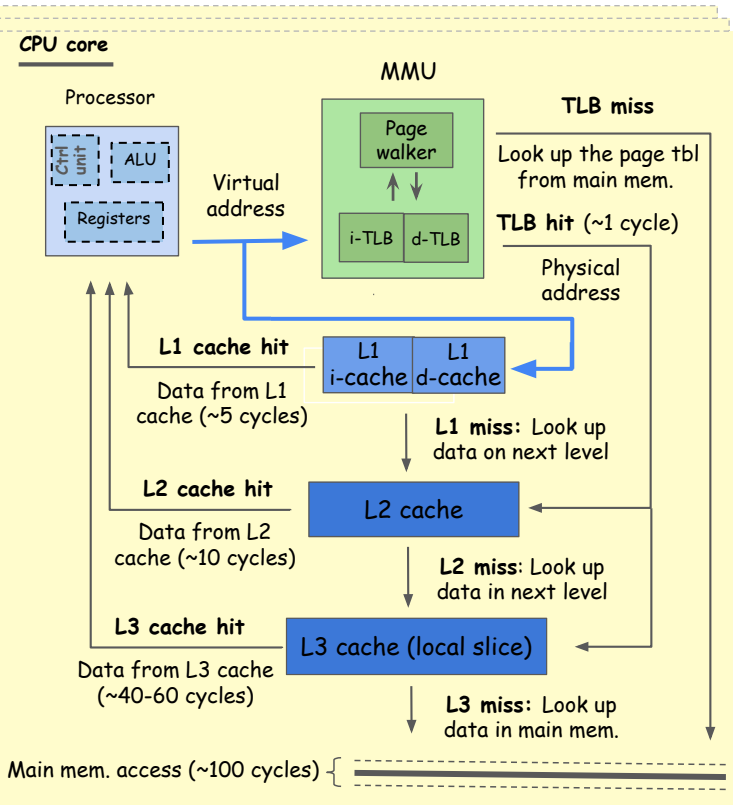
L1 cache: Why not virtually indexed and tagged?



Homonym problem: After a context switch a proc can access leftover cache contents not mapped to its virtual address space (**data leakage**)

» Flush caches on ctx switch? Slow

L1 cache: Why not virtually indexed and tagged?

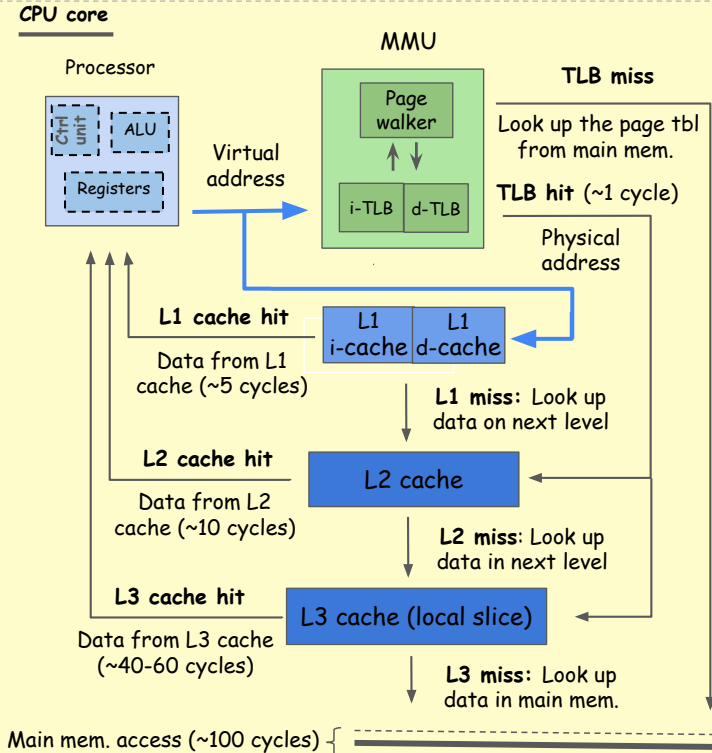


Homonym problem: After a context switch a proc can access leftover cache contents not mapped to its virtual address space (**data leakage**)

» Flush caches on ctx switch? Slow

» Use Addr. Space Ids in cache tags? Wasteful

L1 cache: Why not virtually indexed and tagged?



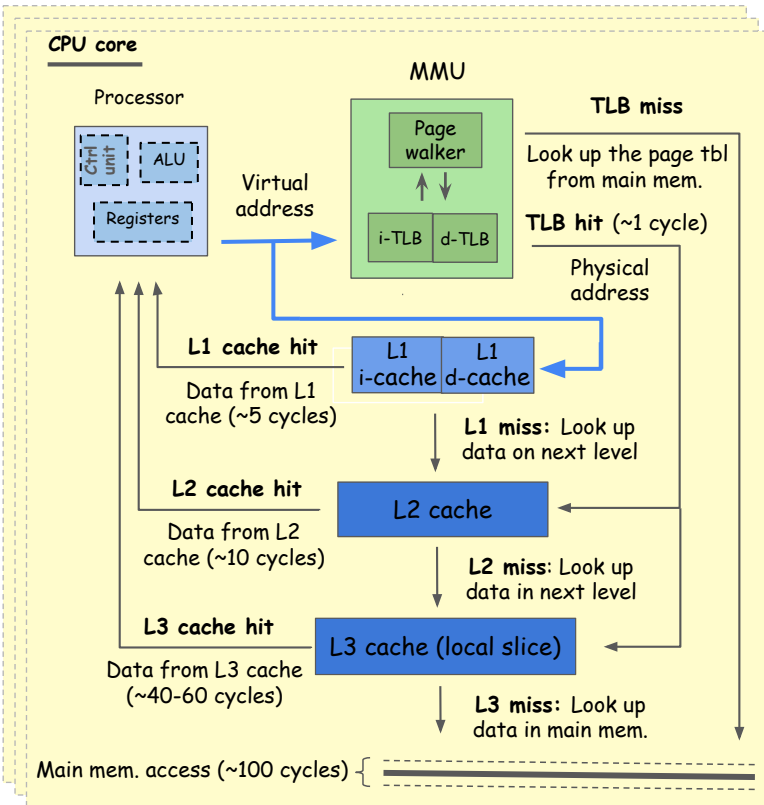
Homonym problem: After a context switch a proc can access leftover cache contents not mapped to its virtual address space (**data leakage**)

» Flush caches on ctx switch? Slow

» Use Addr. Space Ids in cache tags? Wasteful

Synonym problem (a.k.a aliasing): Different virtual addresses reference the same data => OS decision => Invisible to hardware's cache coherency protocols...

L1 cache: Why not virtually indexed and tagged?



Homonym problem: After a context switch a proc can access leftover cache contents not mapped to its virtual address space (**data leakage**)

» Flush caches on ctx switch? Slow

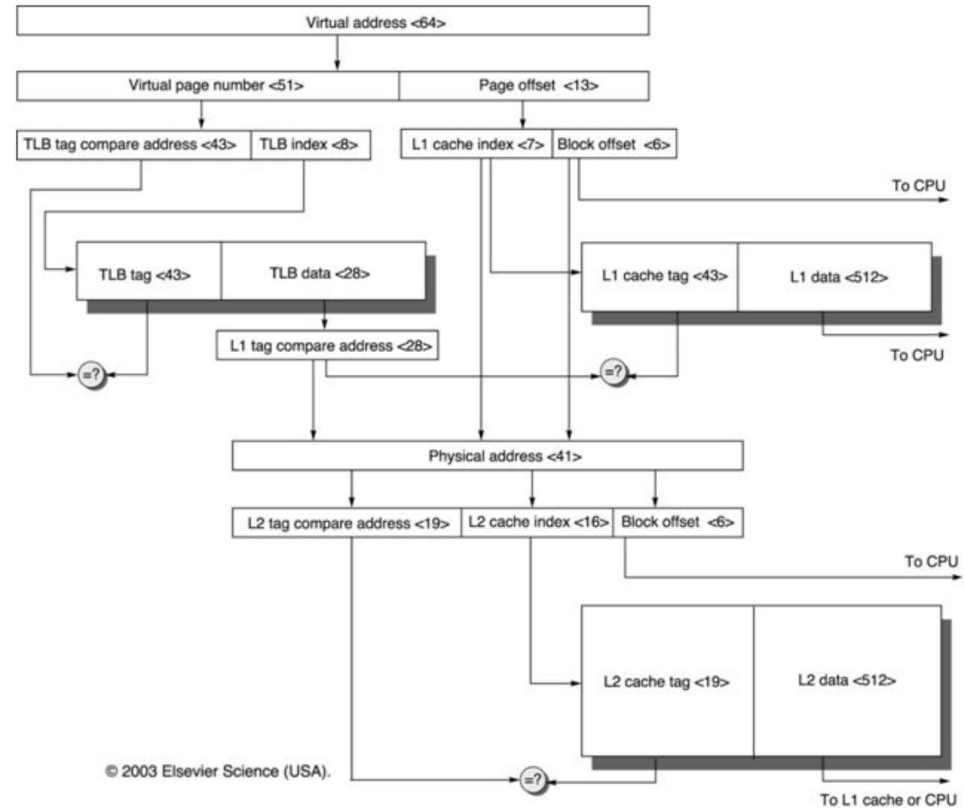
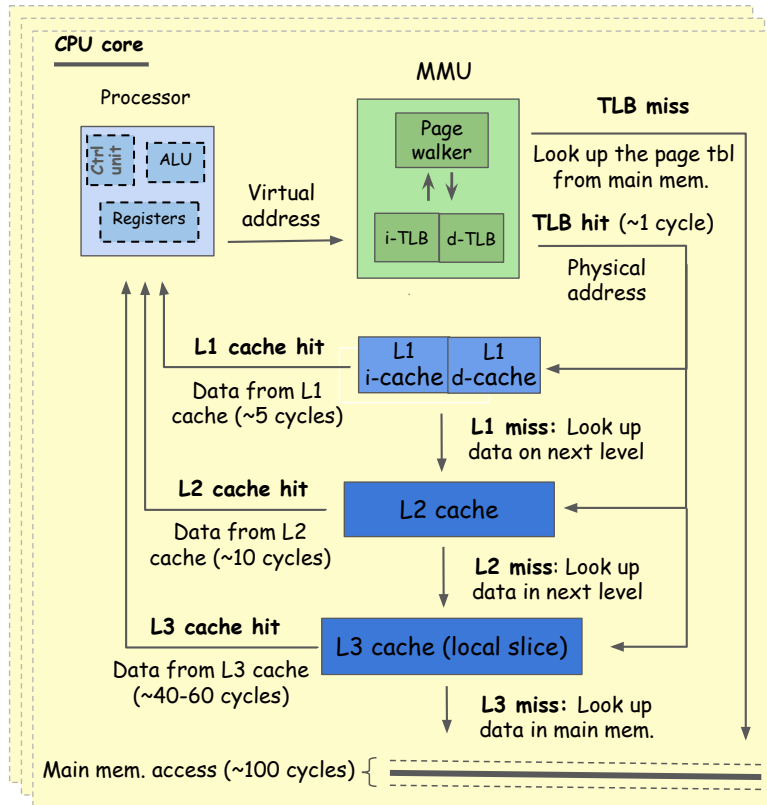
» Use Addr. Space Ids in cache tags? Wasteful

Synonym problem (a.k.a aliasing): Different virtual addresses reference the same data => OS decision => Invisible to hardware's cache coherency protocols...

» Page coloring (beyond scope)

» Use VIPT or PIPT caches

L1 cache: Virtually-indexed / Physically-tagged



How is virtual memory implemented?

- > Usually considered the most complex kernel subsystem
- It's simple if you view it as an index

Ask the right questions

~~Q1: What are the entries of the index?~~

~~Q2: How are the entries of the index used?~~

~~Q3: How are the entries of the index allocated?~~

~~Q4: How are the entries of the index replaced?~~

The working set model for program behavior

- > A process can be in main memory iff all the pages it is currently using can be in main memory, by P. Denning (1968)

The working set model for program behavior

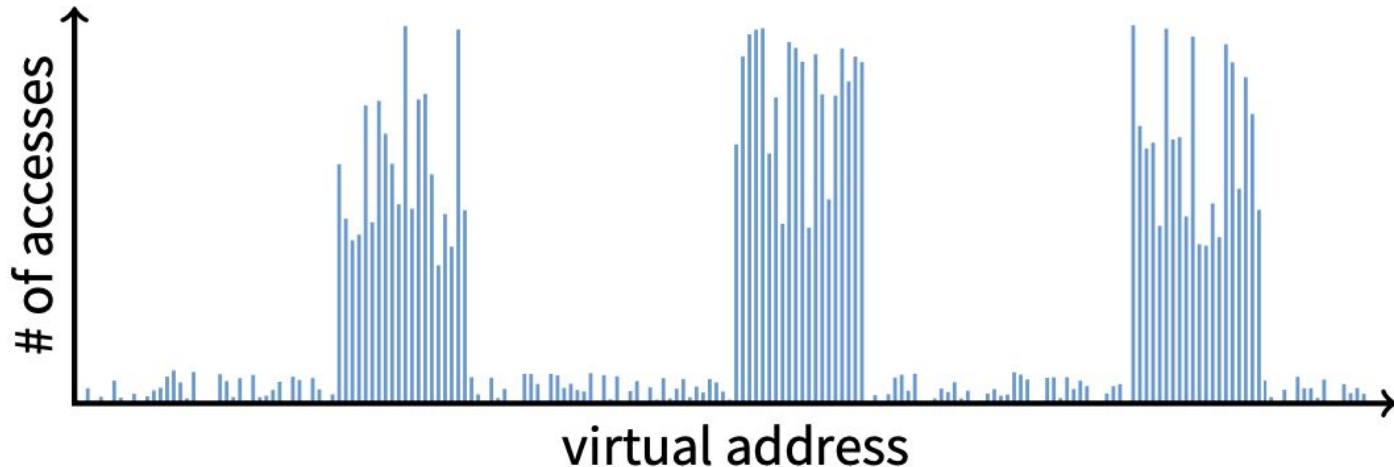
> A process can be in main memory iff all the pages it is currently using can be in main memory, by P. Denning (1968)

Rule of thumb: 20% of memory gets 80% of accesses

The working set model for program behavior

> A process can be in main memory iff all the pages it is currently using can be in main memory, by P. Denning (1968)

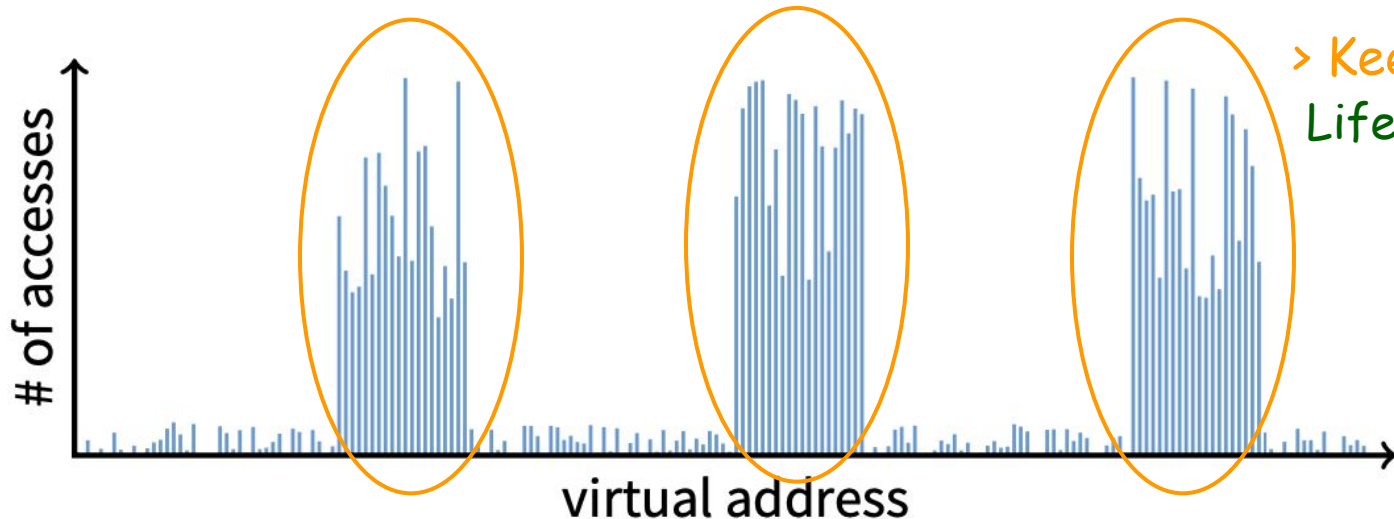
Rule of thumb: 20% of memory gets 80% of accesses



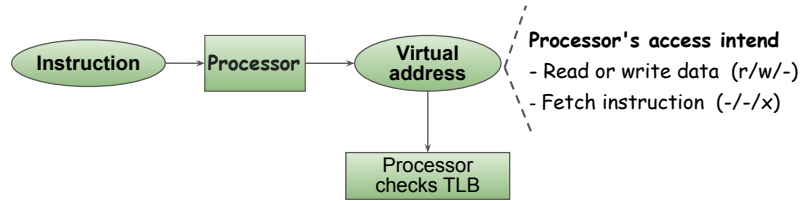
The working set model for program behavior

> A process can be in main memory iff all the pages it is currently using can be in main memory, by P. Denning (1968)

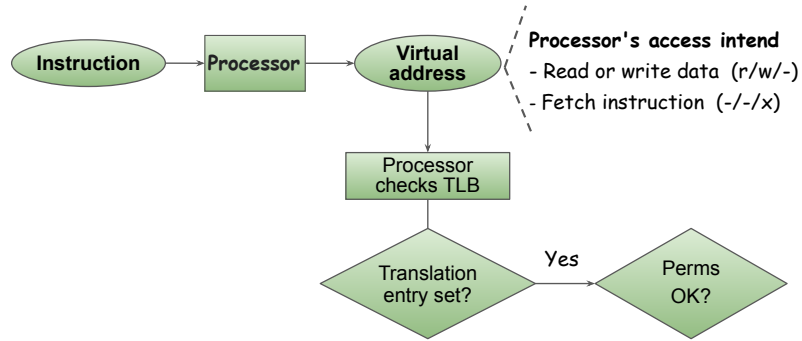
Rule of thumb: 20% of memory gets 80% of accesses



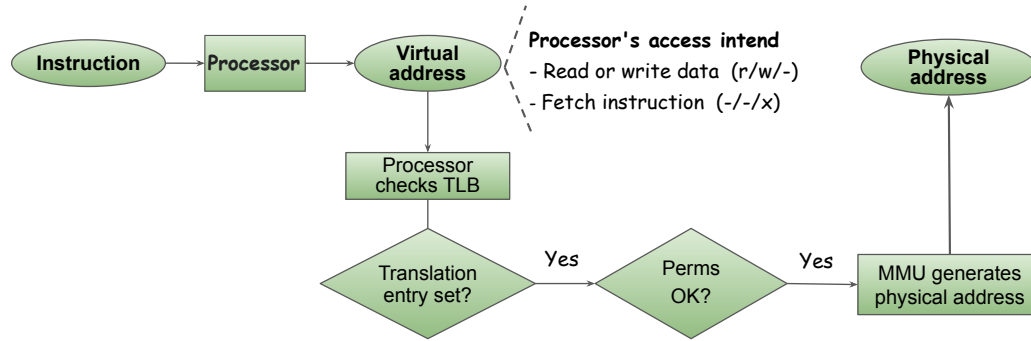
Demand paging via page fault handling



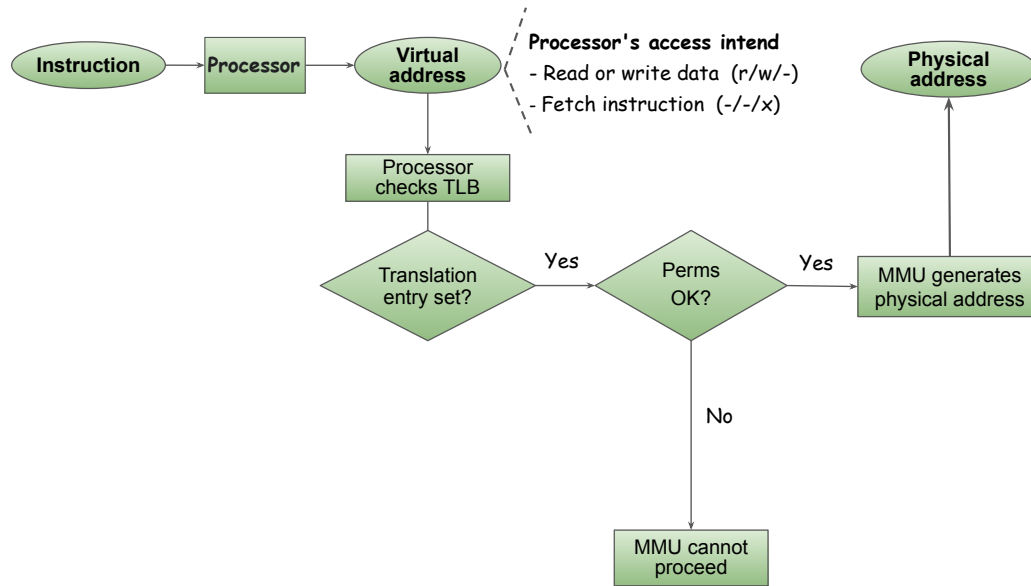
Demand paging via page fault handling



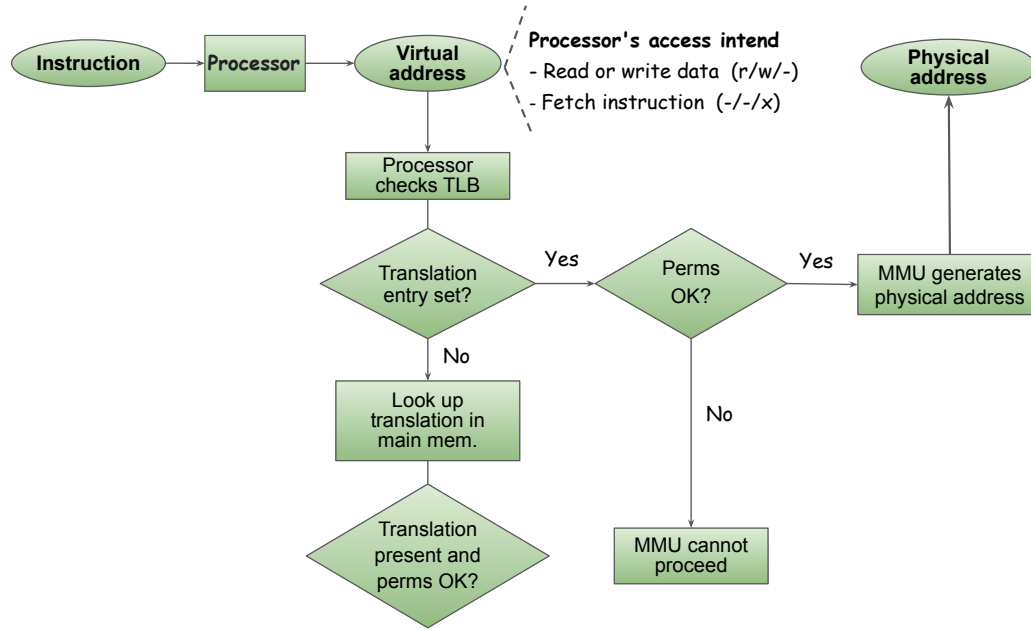
Demand paging via page fault handling



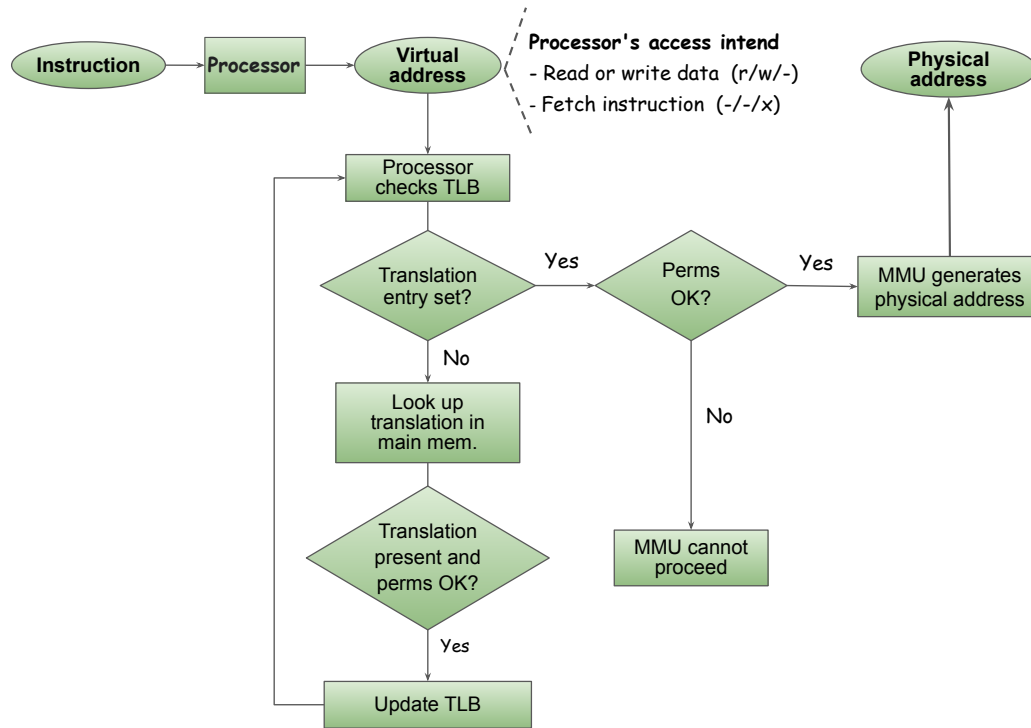
Demand paging via page fault handling



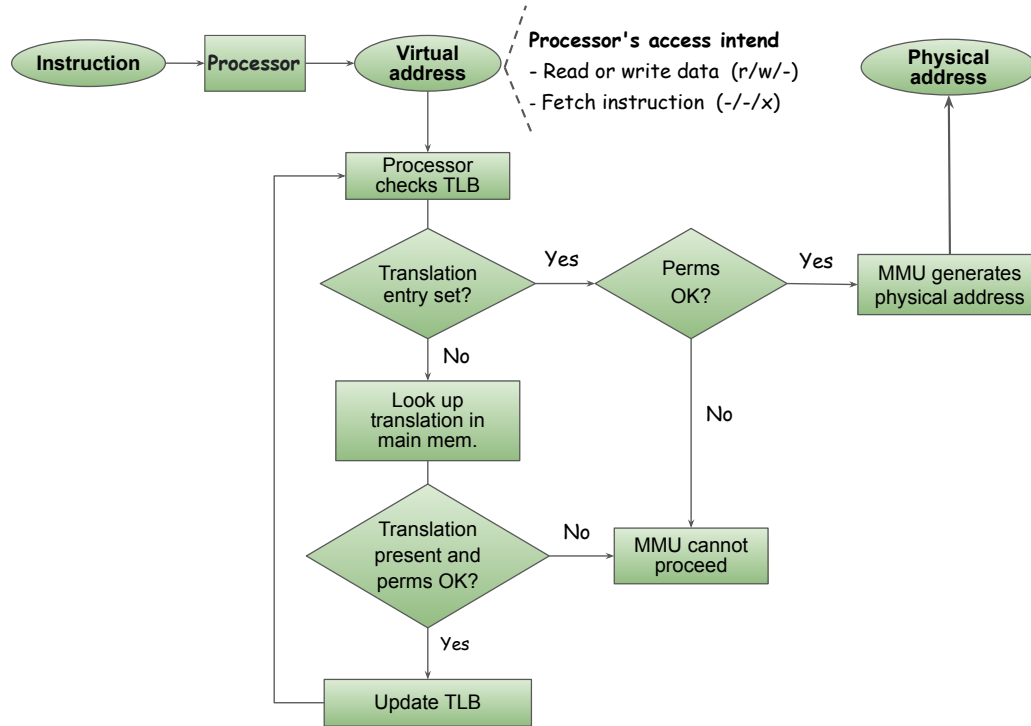
Demand paging via page fault handling



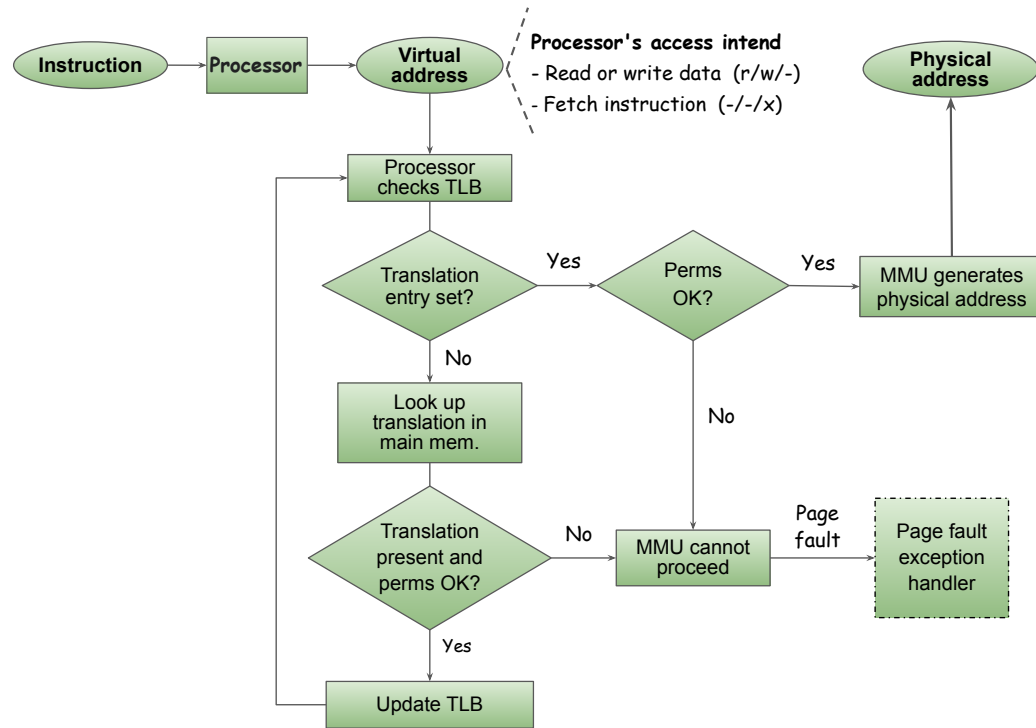
Demand paging via page fault handling



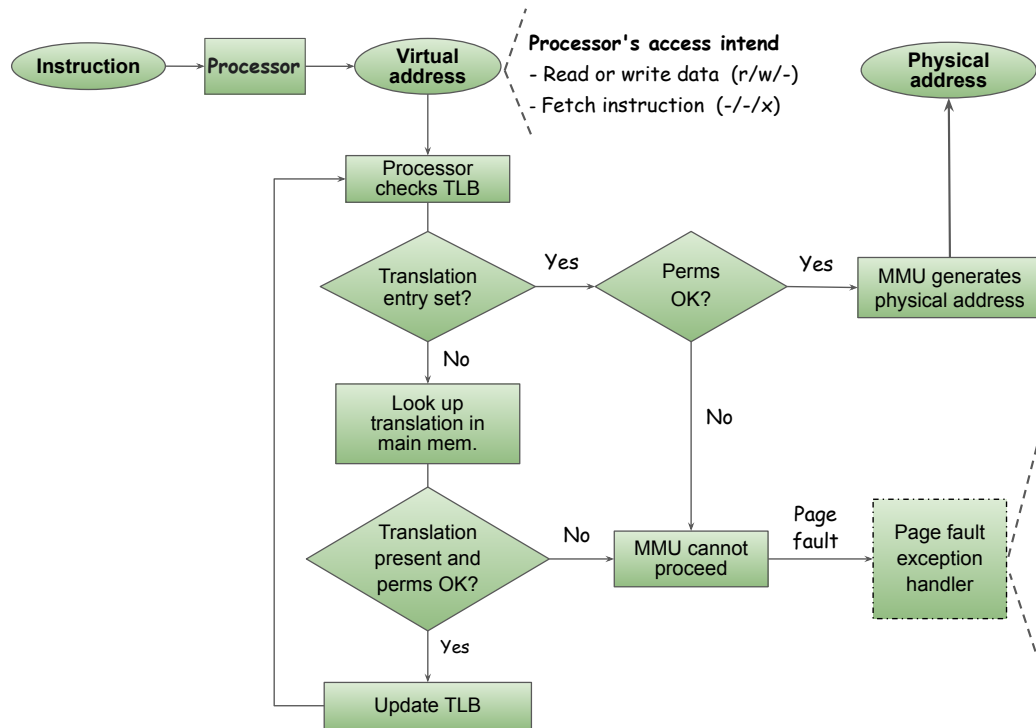
Demand paging via page fault handling



Demand paging via page fault handling



Demand paging via page fault handling



Page fault handling (see [here](#) and [here](#))

> Access not permitted (see [here](#)) => SIGSEGV

>> Translation entry not set (see [here](#))

>>> File-backed VMA (see [here](#))

Minor p.f. => [do read fault](#)

=> [do cow fault](#)

=> [do shared fault](#)

>>> Not file-backed VMA (see [here](#))

Minor p.f => [do anonymous page](#)

>> Translation entry set (see [here](#))

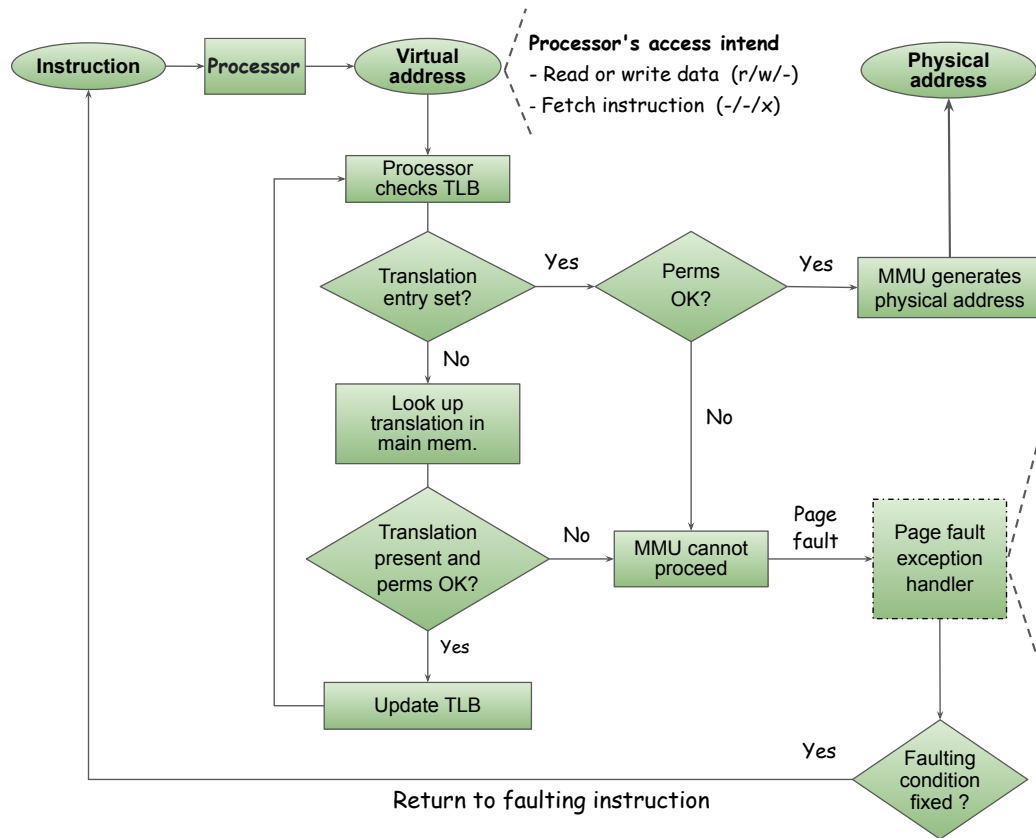
>>> "present" bit off (see [here](#))

Major p.f. => [do swap page](#)

>>> "write" bit off (see [here](#))

Minor p.f. => [do wp page](#)

Demand paging via page fault handling



Page fault handling (see [here](#) and [here](#))

> Access **not** permitted (see [here](#)) => SIGSEGV

>> Translation entry **not** set (see [here](#))

>>> File-backed VMA (see [here](#))

Minor p.f. => [do read fault](#)

=> [do cow fault](#)

=> [do shared fault](#)

>>> Not file-backed VMA (see [here](#))

Minor p.f => [do anonymous page](#)

>> Translation entry **set** (see [here](#))

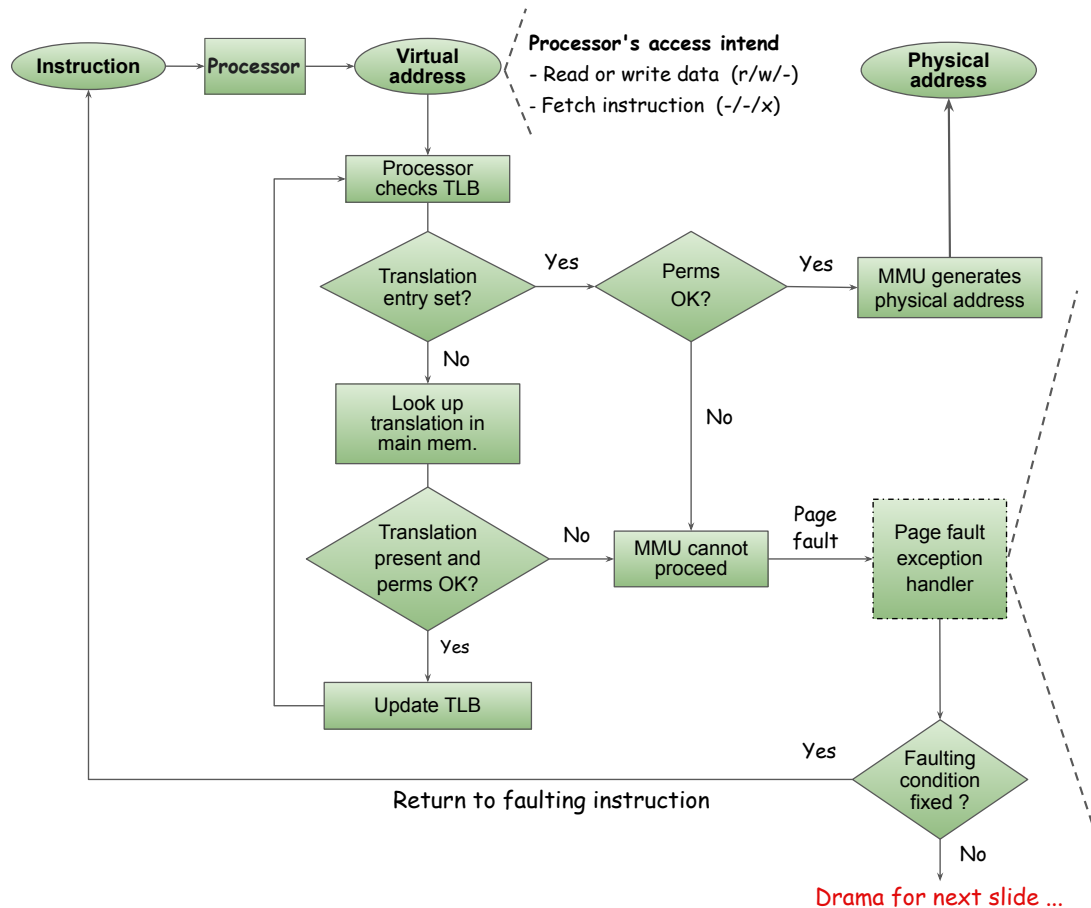
>>> "present" bit off (see [here](#))

Major p.f. => [do swap page](#)

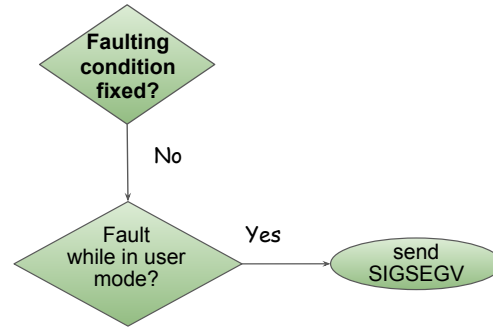
>>> "write" bit off (see [here](#))

Minor p.f. => [do wp page](#)

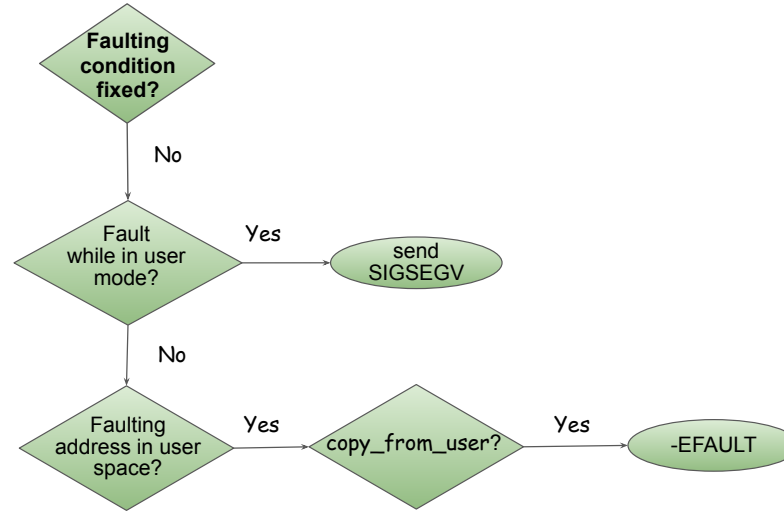
Demand paging via page fault handling



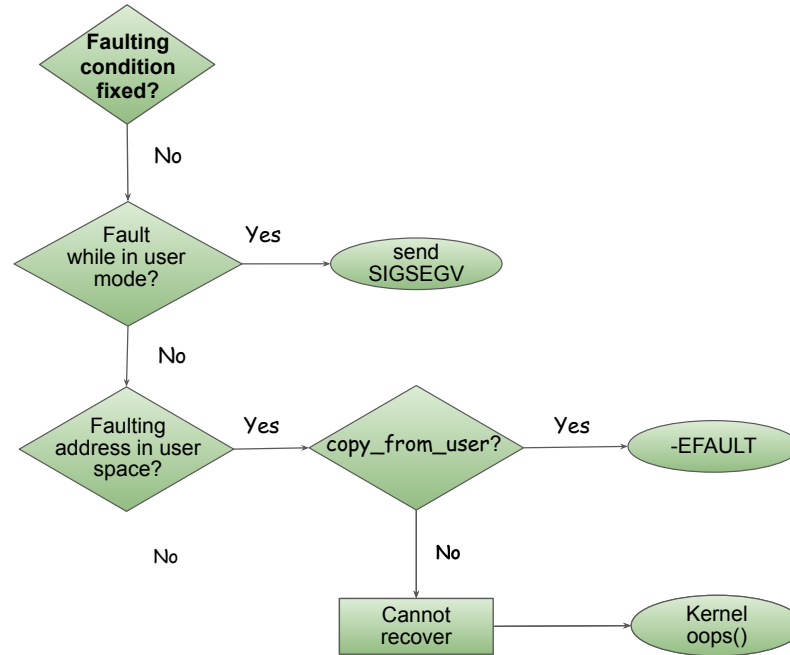
Page fault handling (cont'ed)



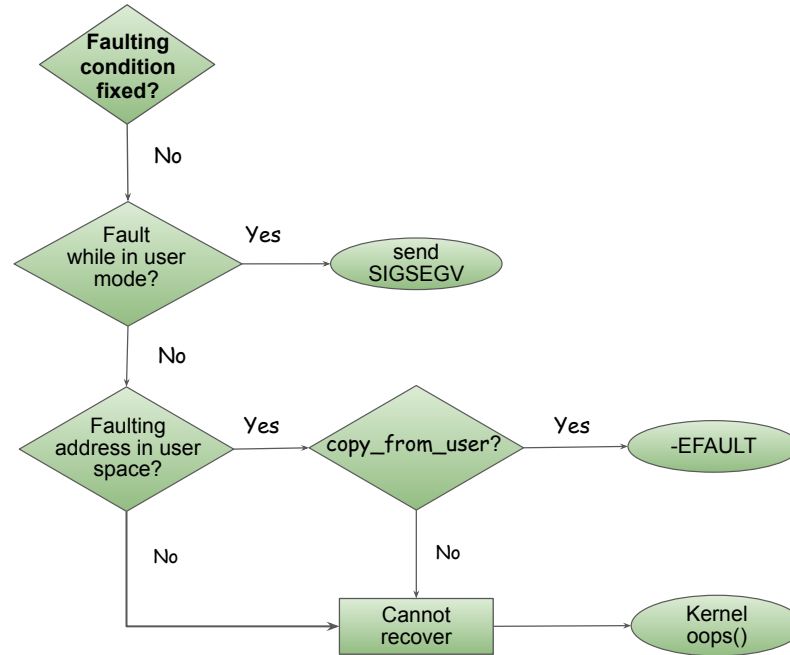
Page fault handling (cont'ed)



Page fault handling (cont'ed)



Page fault handling (cont'ed)



Page fault demo: read->write vs. write->read

```
int main(int argc, char **argv) {
    char a;
    int vma_size = 2 * 4096;
    char *buffer = mmap(NULL, vma_size,
        PROT_READ | PROT_WRITE,
        MAP_PRIVATE | MAP_ANONYMOUS, -1, 0);

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st read)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd read)\n",
            i / 4096, end_time - start_time);
    }

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st write)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd write)\n",
            i / 4096, end_time - start_time);
    }
}
```

→ [git:\(master\)](#) x ./read_write_page_faults

(1) minor page faults: 87, major page faults: 0

page-0: Time elapsed: 2633 nanoseconds (1st read)
page-0: Time elapsed: 78 nanoseconds (2nd read)
page-1: Time elapsed: 1956 nanoseconds (1st read)
page-1: Time elapsed: 78 nanoseconds (2nd read)

} Read pg faults

(2) minor page faults: 89, major page faults: 0

Page fault demo: read->write vs. write->read

```
int main(int argc, char **argv) {
    char a;
    int vma_size = 2 * 4096;
    char *buffer = mmap(NULL, vma_size,
        PROT_READ | PROT_WRITE,
        MAP_PRIVATE | MAP_ANONYMOUS, -1, 0);

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st read)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd read)\n",
            i / 4096, end_time - start_time);
    }

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st write)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd write)\n",
            i / 4096, end_time - start_time);
    }
}
```

→ [git:\(master\)](#) X ./read_write_page_faults

(1) minor page faults: 87, major page faults: 0
page-0: Time elapsed: 2633 nanoseconds (1st read)
page-0: Time elapsed: 78 nanoseconds (2nd read)
page-1: Time elapsed: 1956 nanoseconds (1st read)
page-1: Time elapsed: 78 nanoseconds (2nd read) } Read pg faults

(2) minor page faults: 89, major page faults: 0
page-0: Time elapsed: 4131 nanoseconds (1st write)
page-0: Time elapsed: 113 nanoseconds (2nd write)
page-1: Time elapsed: 3694 nanoseconds (1st write)
page-1: Time elapsed: 58 nanoseconds (2nd write) } Write pg faults

(3) minor page faults: 91, major page faults: 0

Page fault demo: read->write vs. write->read

```
int main(int argc, char **argv) {
    char a;
    int vma_size = 2 * 4096;
    char *buffer = mmap(NULL, vma_size,
        PROT_READ | PROT_WRITE,
        MAP_PRIVATE | MAP_ANONYMOUS, -1, 0);

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st read)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd read)\n",
            i / 4096, end_time - start_time);
    }

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st write)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd write)\n",
            i / 4096, end_time - start_time);
    }
}
```

→ [git:\(master\)](#) X ./read_write_page_faults

(1) minor page faults: 87, major page faults: 0
page-0: Time elapsed: 2633 nanoseconds (1st read)
page-0: Time elapsed: 78 nanoseconds (2nd read)
page-1: Time elapsed: 1956 nanoseconds (1st read)
page-1: Time elapsed: 78 nanoseconds (2nd read) } Read pg faults

(2) minor page faults: 89, major page faults: 0
page-0: Time elapsed: 4131 nanoseconds (1st write)
page-0: Time elapsed: 113 nanoseconds (2nd write)
page-1: Time elapsed: 3694 nanoseconds (1st write)
page-1: Time elapsed: 58 nanoseconds (2nd write) } Write pg faults

→ [git:\(master\)](#) X ./write_read_page_faults

(1) minor page faults: 88, major page faults: 0
page-0: Time elapsed: 5868 nanoseconds (1st write)
page-0: Time elapsed: 115 nanoseconds (2nd write)
page-1: Time elapsed: 5487 nanoseconds (1st write)
page-1: Time elapsed: 48 nanoseconds (2nd write) } Write pg faults

(2) minor page faults: 90, major page faults: 0

Page fault demo: read->write vs. write->read

```
int main(int argc, char **argv) {
    char a;
    int vma_size = 2 * 4096;
    char *buffer = mmap(NULL, vma_size,
        PROT_READ | PROT_WRITE,
        MAP_PRIVATE | MAP_ANONYMOUS, -1, 0);

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st read)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        a = buffer[i];
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd read)\n",
            i / 4096, end_time - start_time);
    }

    for (int i = 0; i < vma_size; i += 4096) {
        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (1st write)\n",
            i / 4096, end_time - start_time);

        start_time = clock_gettime_ns();
        buffer[i] = 'A';
        end_time = clock_gettime_ns();
        printf("page-%d: Time elapsed: %lu ns (2nd write)\n",
            i / 4096, end_time - start_time);
    }
}
```

→ [git:\(master\)](#) X ./read_write_page_faults

(1) minor page faults: 87, major page faults: 0
page-0: Time elapsed: 2633 nanoseconds (1st read)
page-0: Time elapsed: 78 nanoseconds (2nd read)
page-1: Time elapsed: 1956 nanoseconds (1st read)
page-1: Time elapsed: 78 nanoseconds (2nd read) } Read pg faults

(2) minor page faults: 89, major page faults: 0
page-0: Time elapsed: 4131 nanoseconds (1st write)
page-0: Time elapsed: 113 nanoseconds (2nd write)
page-1: Time elapsed: 3694 nanoseconds (1st write)
page-1: Time elapsed: 58 nanoseconds (2nd write) } Write pg faults

(3) minor page faults: 91, major page faults: 0

→ [git:\(master\)](#) X ./write_read_page_faults

(1) minor page faults: 88, major page faults: 0
page-0: Time elapsed: 5868 nanoseconds (1st write)
page-0: Time elapsed: 115 nanoseconds (2nd write)
page-1: Time elapsed: 5487 nanoseconds (1st write)
page-1: Time elapsed: 48 nanoseconds (2nd write) } Write pg faults

(2) minor page faults: 90, major page faults: 0
page-0: Time elapsed: 90 nanoseconds (1st read)
page-0: Time elapsed: 92 nanoseconds (2nd read)
page-1: Time elapsed: 59 nanoseconds (1st read)
page-1: Time elapsed: 46 nanoseconds (2nd read) } No read pg faults

(3) minor page faults: 90, major page faults: 0

How is virtual memory implemented?

- > Usually considered the most complex kernel subsystem
- It's simple if you view it as an index

Ask the right questions

~~Q1: What are the entries of the index?~~

~~Q2: How are the entries of the index used?~~

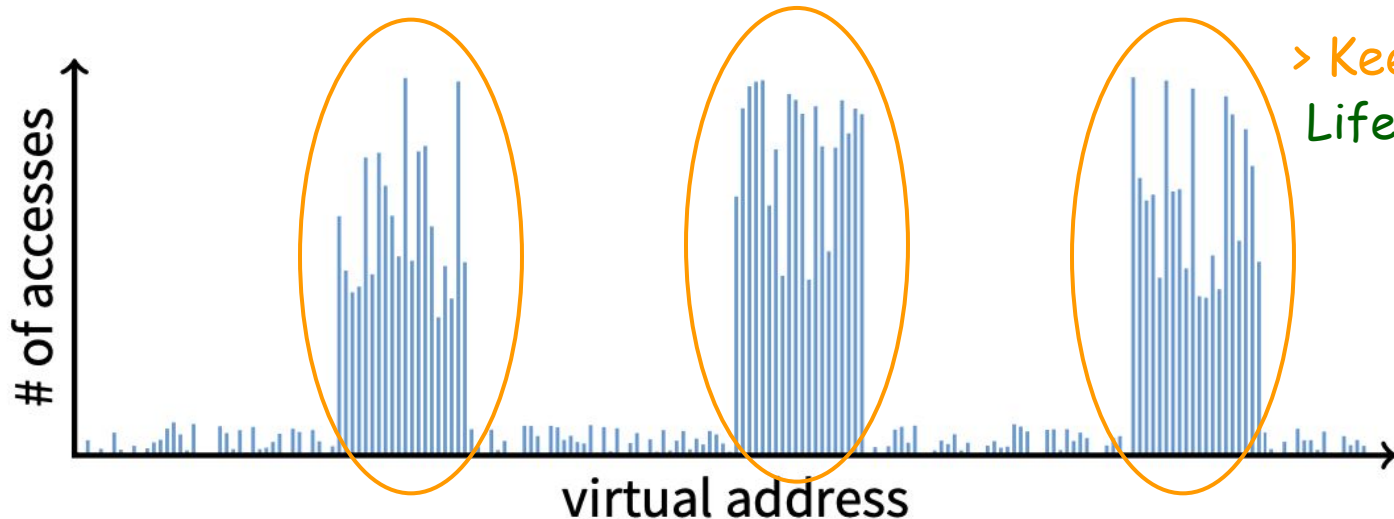
~~Q3: How are the entries of the index allocated?~~

~~Q4: How are the entries of the index replaced?~~

The working set model for program behavior

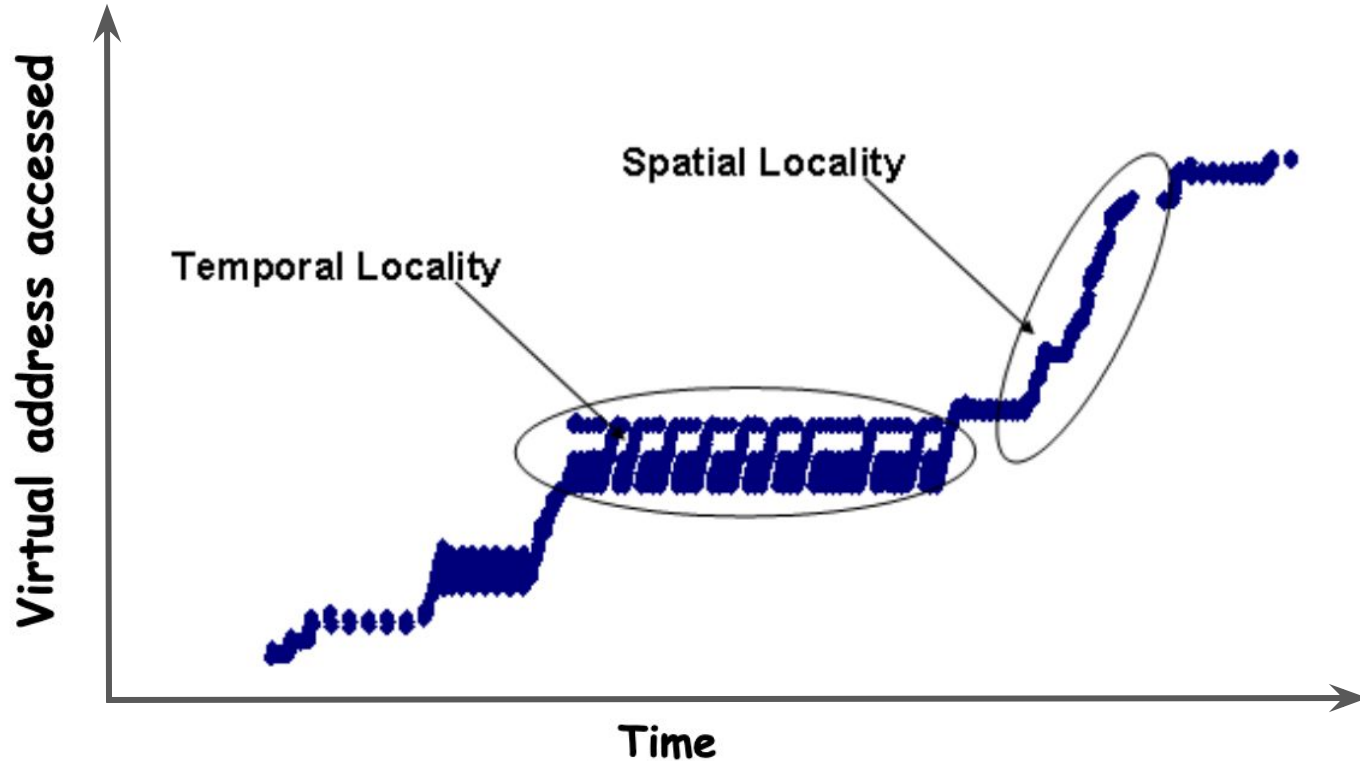
> A process can be in main memory iff all the pages that it is currently using can be in main memory, by P. Denning (1968)

Rule of thumb: 20% of memory gets 80% of accesses



> Keep these "hot"?
Life is good...

Temporal and spatial locality



Replacement policies

- > When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

> Approximate LRU (a.k.a. CLOCK algorithm)

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

> **Approximate LRU (a.k.a. CLOCK algorithm)**

- LRU is expensive to implement in software => **Hardware support?**

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

> **Approximate LRU (a.k.a. CLOCK algorithm)**

- LRU is expensive to implement in software => **Hardware support?**
- **Remember PTE "referenced" bit?**

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

> **Approximate LRU (a.k.a. CLOCK algorithm)**

- LRU is expensive to implement in software => **Hardware support?**
- **Remember PTE "referenced" bit?** The OS periodically swipes through PTEs and resets "accessed" bits

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

> **Approximate LRU (a.k.a. CLOCK algorithm)**

- LRU is expensive to implement in software => **Hardware support?**
- **Remember PTE "referenced" bit?** The OS periodically swipes through PTEs and resets "accessed" bits
- Translation entries w/ "accessed" bits off

Replacement policies

> When all physical memory is in use, the OS must reclaim some physical page frames in order to serve new translations

Which translations are "good candidates" to be evicted?

> **Converse of locality property:** The Last Recently Used (LRU) page is the least to be used soon ("best eviction candidate")

> **Approximate LRU (a.k.a. CLOCK algorithm)**

- LRU is expensive to implement in software => **Hardware support?**
- **Remember PTE "referenced" bit?** The OS periodically swipes through PTEs and resets "accessed" bits
- **Translation entries w/ "accessed" bits off => Not recently accessed**

Replacement policies

> Not all pages are equal: Reclaiming decision is made based on the the type of the occupied pages

Replacement policies

- > **Not all pages are equal:** Reclaiming decision is made based on the type of the occupied pages
 - **File-backed pages:** Dropped from the kernel's page cache **first** since they can be reloaded later from the corresponding file

Replacement policies

> **Not all pages are equal:** Reclaiming decision is made based on the type of the occupied pages

- **File-backed pages:** Dropped from the kernel's page cache **first** since they can be reloaded later from the corresponding file
- **Anonymous pages w/o file backing:** Must be written to the swap space in order to be available for reloading; **dropped next**

Replacement policies

- > **Not all pages are equal:** Reclaiming decision is made based on the type of the occupied pages
 - **File-backed pages:** Dropped from the kernel's page cache **first** since they can be reloaded later from the corresponding file
 - **Anonymous pages w/o file backing:** Must be written to the swap space in order to be available for reloading; **dropped next**
 - **File-system-related kernel caches:** **Dropped last** (next chapter)

Replacement policies

- > **Not all pages are equal:** Reclaiming decision is made based on the the type of the occupied pages
 - **File-backed pages:** Dropped from the kernel's page cache **first** since they can be reloaded later from the corresponding file
 - **Anonymous pages w/o file backing:** Must be written to the swap space in order to be available for reloading; **dropped next**
 - **File-system-related kernel caches:** **Dropped last** (next chapter)
 - **Last resort:** Out-Of-Memory (OOM) killer terminates the process

Thrashing (or, "when nothing works")

Thrashing (or, "when nothing works")

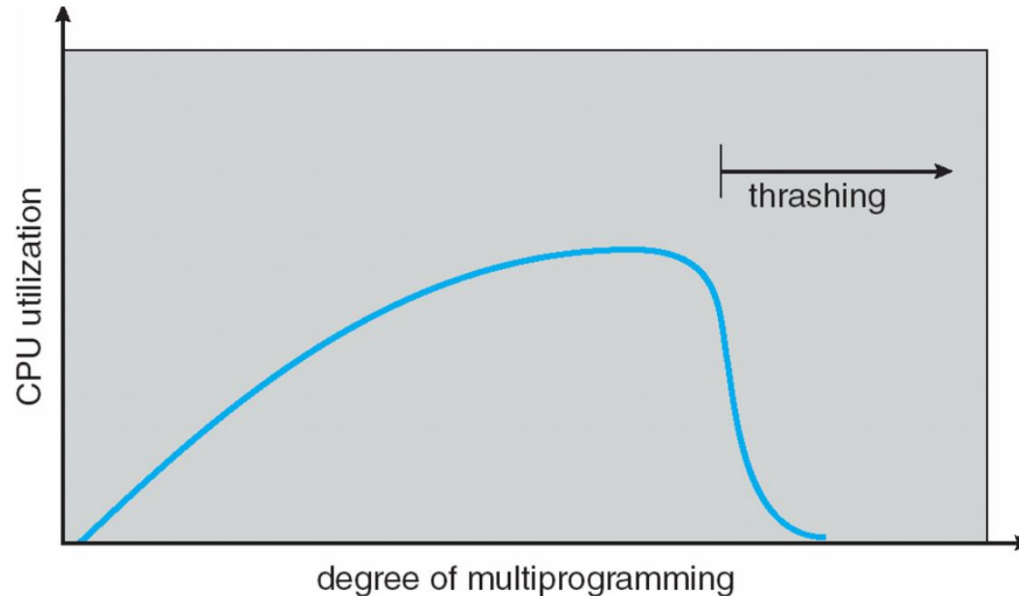
- I/O devices at 100% utilization but CPU utilization drops to 0%
- Processes remain blocked waiting for pages to be fetched in

Thrashing (or, "when nothing works")

- I/O devices at 100% utilization but CPU utilization drops to 0%
- Processes remain blocked waiting for pages to be fetched in
- The system does no useful work

Thrashing (or, "when nothing works")

- > I/O devices at 100% utilization but CPU utilization drops to 0%
- Processes remain blocked waiting for pages to be fetched in
- The system does no useful work



Thrashing (or, "when nothing works")

> Reasons for thrashing

- Memory accesses w/o temporal locality \Rightarrow Past \neq Future

Thrashing (or, "when nothing works")

> Reasons for thrashing

- Memory accesses w/o temporal locality \Rightarrow Past \neq Future
- Processes' hot memory does not fit in main memory

Thrashing (or, "when nothing works")

> Reasons for thrashing

- Memory accesses w/o temporal locality => Past \neq Future
- Processes' hot memory does not fit in main memory

> What we ordered? Memory with the size of disk and the speed of CPU caches

> What we got?

Thrashing (or, "when nothing works")

> Reasons for thrashing

- Memory accesses w/o temporal locality \Rightarrow Past \neq Future
- Processes' hot memory does not fit in main memory

> What we ordered? Memory with the size of disk and the speed of CPU caches

> What we got? Memory with the access time of the disk :-(

POSIX Memory management syscalls (cont'ed)

`void *mmap (void *addr, ...)`

> Creates a new mapping in the virtual address space of the calling process

- `addr`: If NULL, the kernel chooses the address at which to create the mapping
- On success, `mmap(...)` returns a pointer to the mapped area

POSIX Memory management syscalls (cont'ed)

`void *mmap (void *addr, ...)`

➤ Creates a new mapping in the virtual address space of the calling process

- `addr`: If NULL, the kernel chooses the address at which to create the mapping
- On success, `mmap(...)` returns a pointer to the mapped area

`int munmap (void *addr, size_t length)`

➤ Deletes the mappings for the specified address range

- `addr`: Start of the address range
- `length`: Size of the address range
- On success, `munmap(...)` returns 0

POSIX Memory management syscalls (cont'ed)

`void *mmap (void *addr, ...)`

➤ Creates a new mapping in the virtual address space of the calling process

- `addr`: If NULL, the kernel chooses the address at which to create the mapping
- On success, `mmap(...)` returns a pointer to the mapped area

`int munmap (void *addr, size_t length)`

➤ Deletes the mappings for the specified address range

- `addr`: Start of the address range
- `length`: Size of the address range
- On success, `munmap(...)` returns 0

`int mprotect (void *addr, size_t len, int prot)`

➤ Updates the protections for page(s) in range `[addr, addr+len)` to "prot"

- `addr`: Must be aligned to a page boundary
- On success, `mprotect(...)` returns 0

POSIX Memory management syscalls (cont'ed)

`int msync (void *addr, ...):`

> Flushes to disk changes made to the in-kernel copy of a file that was mapped into memory using mmap

- **Note:** Without use of msync(...) there is no guarantee that changes made to the in-kernel copy of a file will be written back to disk before munmap(...) is called
- On success, msync(...) returns 0

POSIX Memory management syscalls (cont'ed)

`int msync (void *addr, ...):`

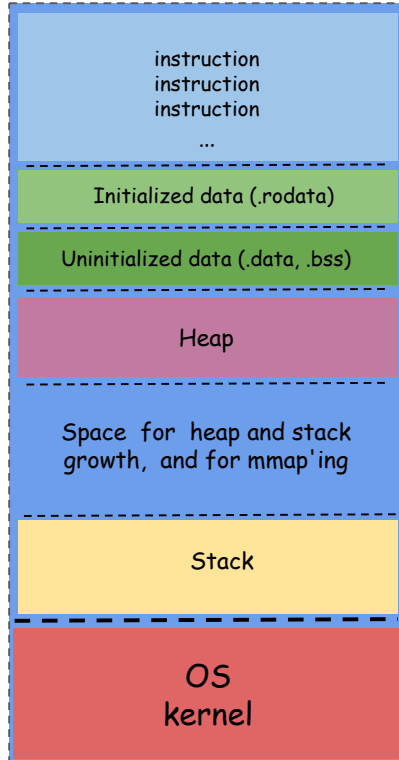
> Flushes to disk changes made to the in-kernel copy of a file that was mapped into memory using mmap

- **Note:** Without use of msync(...) there is no guarantee that changes made to the in-kernel copy of a file will be written back to disk before munmap(...) is called
- On success, msync(...) returns 0

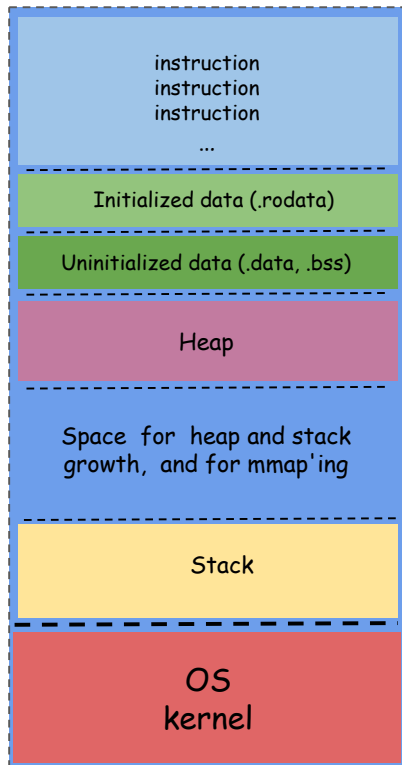
`void *mlock/unlock (void *addr, ...)`

- > Locks/unlocks part or all of the calling process's virtual address space into main mem. preventing that memory from being paged to the swap area
- On success, mlock(...)/unlock(...) returns 0

Putting it all together: 1) Per-process index of VAS segments



Putting it all together: 1) Per-process index of VAS segments



```
struct task_struct {  
    ...  
    char comm[TASK_COMM_LEN];  
    struct mm_struct *active_mm;  
    struct mm_struct *mm;  
    struct mm_struct {  
        ...  
        pgd_t *pgd;  
        unsigned long start_code, end_code;  
        unsigned long start_data, end_data;  
        unsigned long start_brk, brk;  
        unsigned long start_stack;  
        struct maple_tree mm_mt;  
        struct vm_area_struct {  
            unsigned long vm_start;  
            unsigned long vm_end;  
            pgprot_t vm_page_prot;  
            vm_flags_t vm_flags;  
            ...  
        }  
    }  
}
```

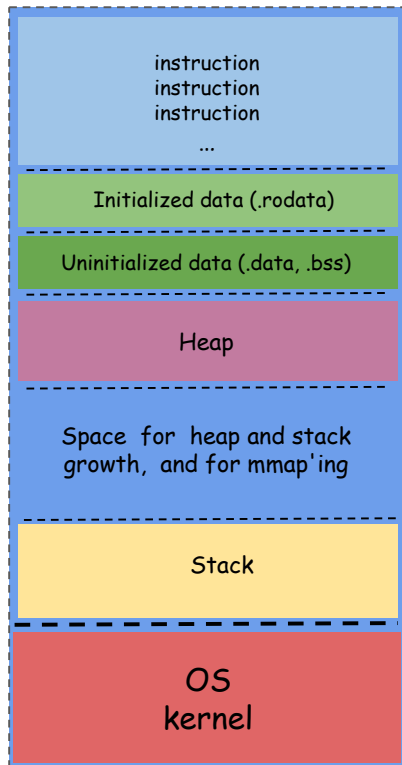
> Mapple tree: Read [here](#)

→ git:(master) X cat /proc/12453/maps

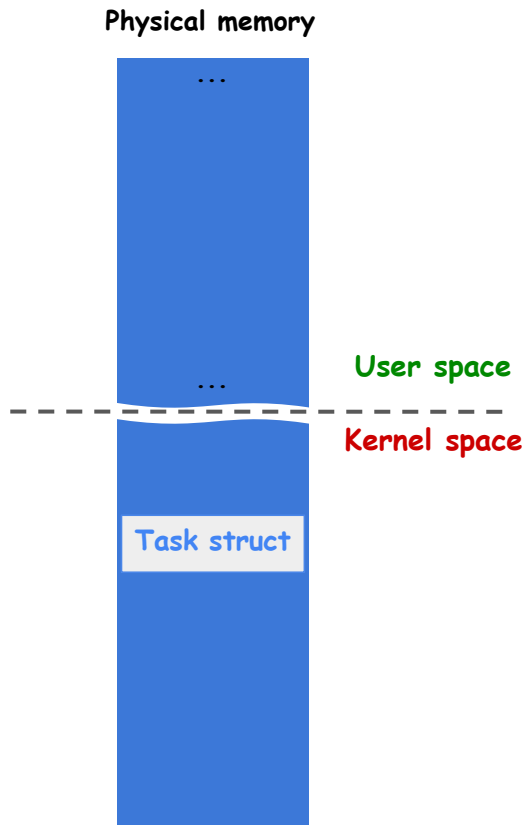
```
aaaadcda0000-aaaadcda1000 r-xp ... /foo  
aaaadcdb0000-aaaadcdb1000 r--p ... /foo  
aaaadcdb1000-aaaadcdb2000 rw-p ... /foo  
...  
aaaaed20d000-aaaaed22e000 rw-p ... [heap]  
ffffbe610000-ffffbe798000 r-xp ... libc.so  
ffffbe798000-ffffbe7a7000 ---p ... libc.so  
...  
fffff6461000-fffff6482000 rw-p ... [stack]
```

Three arrows point from the `mm_mt` field in the struct to the first three lines of the `/proc/12453/maps` output.

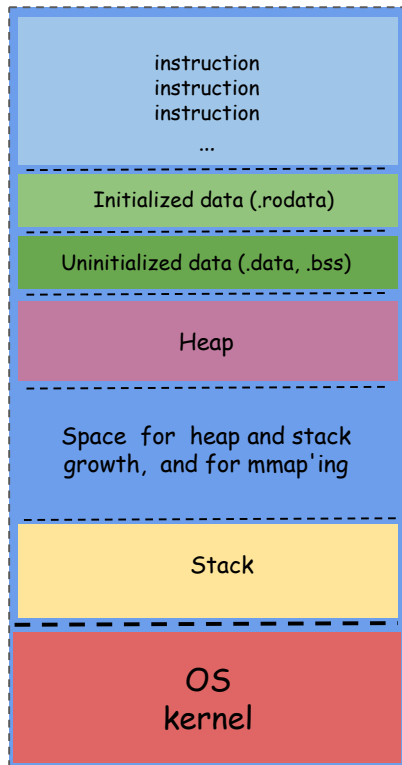
Putting it all together: (1/3) Per-process index of VAS segments



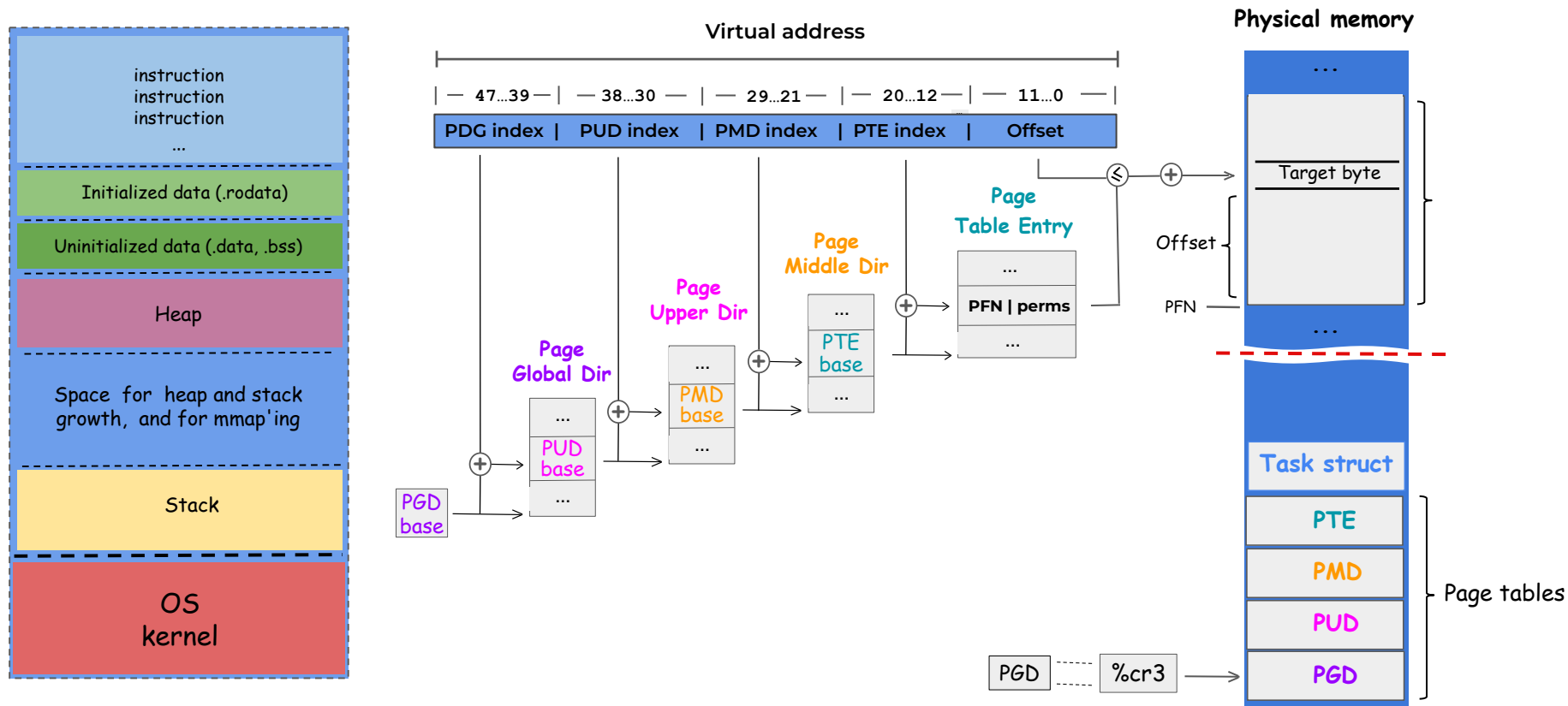
```
struct task_struct {  
    ...  
    char comm[TASK_COMM_LEN];  
    struct mm_struct *active_mm;  
    struct mm_struct *mm;  
    struct mm_struct {  
        ...  
        pgd_t *pgd;  
        unsigned long start_code, end_code;  
        unsigned long start_data, end_data;  
        unsigned long start_brk, brk;  
        unsigned long start_stack;  
        struct maple_tree mm_mt;  
        struct vm_area_struct {  
            unsigned long vm_start;  
            unsigned long vm_end;  
            pgprot_t vm_page_prot;  
            vm_flags_t vm_flags;  
            ...  
        }  
    }  
}
```



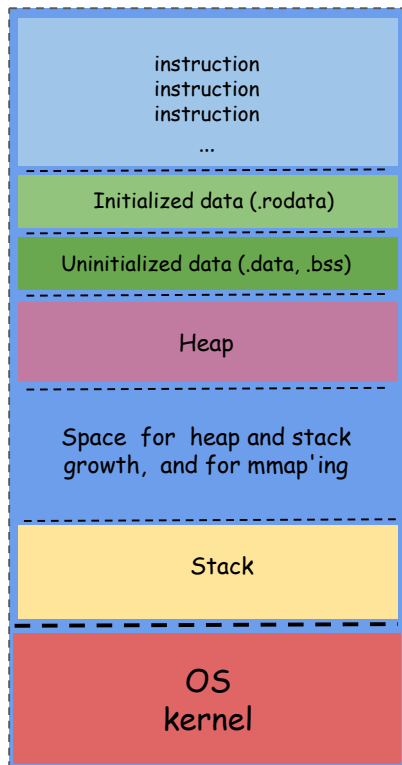
Putting it all together: (2/3) Per-process page table



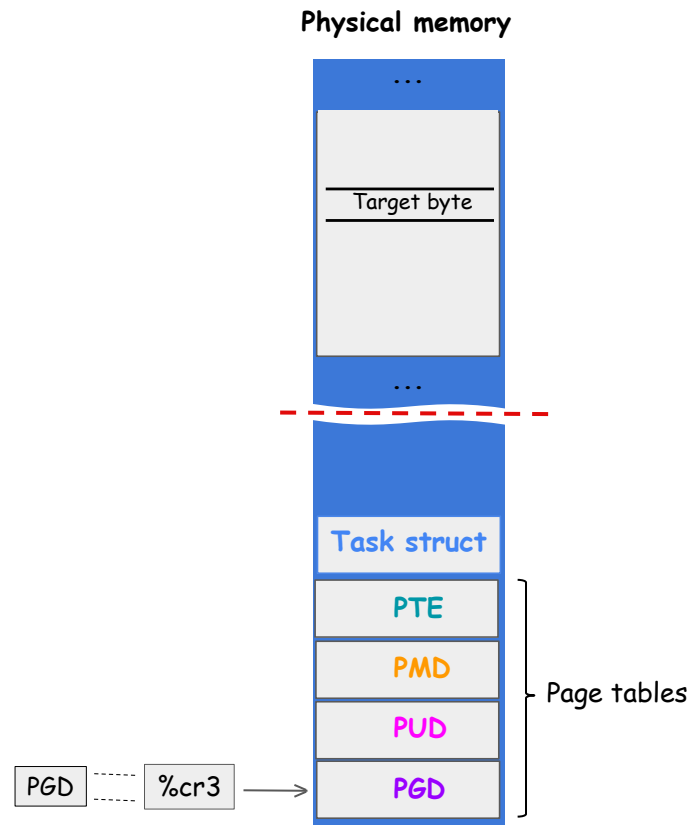
Putting it all together: (2/3) Per-process page table



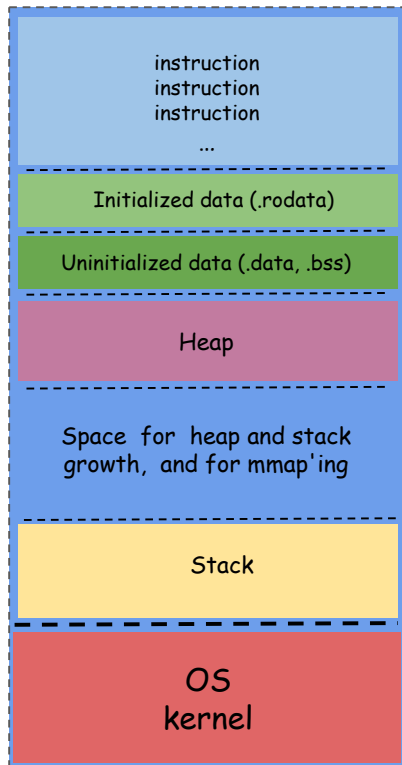
Putting it all together: (2/3) Per-process page table



```
struct task_struct {  
    ...  
    char comm[TASK_COMM_LEN];  
    struct mm_struct *active_mm;  
    struct mm_struct *mm;  
    struct mm_struct {  
        ...  
        pgd_t *pgd;  
        unsigned long start_code, end_code;  
        unsigned long start_data, end_data;  
        unsigned long start_brk, brk;  
        unsigned long start_stack;  
        struct maple_tree mm_mt;  
        struct vm_area_struct {  
            unsigned long vm_start;  
            unsigned long vm_end;  
            pgprot_t vm_page_prot;  
            vm_flags_t vm_flags;  
            ...  
        }  
    }  
}
```



Putting it all together: (3/3) System-wide page frame allocators



> Buddy memory allocator

- Physical memory is divided into "zones"
 - Buddy allocator: Range-based, power-of-2 allocator
 - Each zone has its own buddy allocator instance to manage free physical pages
- >> Request comes in
- Allocator finds the smallest block that fits
 - Only larger blocks available? Splits them into "buddies" until the right size is reached

Physical memory

