

# Virtual Memory: Systems

CSci 2021: Machine Architecture and Organization  
 April 20th-22nd, 2020

Your instructor: Stephen McCamant

Based on slides originally by:  
 Randy Bryant, Dave O'Hallaron

# Today

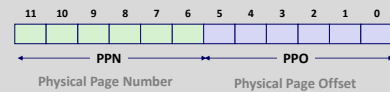
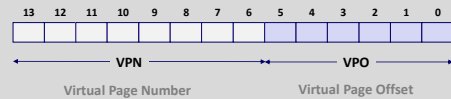
- Simple memory system example
- Case study: Core i7/Linux memory system
- Memory mapping

# Review of Symbols

- Basic Parameters
  - $N = 2^n$ : Number of addresses in virtual address space
  - $M = 2^m$ : Number of addresses in physical address space
  - $P = 2^p$ : Page size (bytes)
- Components of the virtual address (VA)
  - TLBI: TLB index
  - TLBT: TLB tag
  - VPO: Virtual page offset
  - VPN: Virtual page number
- Components of the physical address (PA)
  - PPO: Physical page offset (same as VPO)
  - PPN: Physical page number
  - CO: Byte offset within cache line
  - CI: Cache index
  - CT: Cache tag

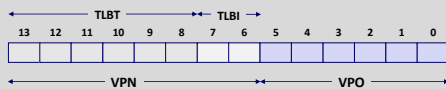
# Simple Memory System Example

- Addressing
  - 14-bit virtual addresses
  - 12-bit physical address
  - Page size = 64 bytes



# 1. Simple Memory System TLB

- 16 entries
- 4-way associative



Set	Tag	PPN	Valid	Tag	PPN	Valid	Tag	PPN	Valid	Tag	PPN	Valid
0	03	-	0	09	0D	1	00	-	0	07	02	1
1	03	2D	1	02	-	0	04	-	0	0A	-	0
2	02	-	0	08	-	0	06	-	0	03	-	0
3	07	-	0	03	0D	1	0A	34	1	02	-	0

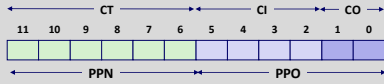
# 2. Simple Memory System Page Table

Only show first 16 entries (out of 256)

VPN	PPN	Valid	VPN	PPN	Valid
00	28	1	08	13	1
01	-	0	09	17	1
02	33	1	0A	09	1
03	02	1	0B	-	0
04	-	0	0C	-	0
05	16	1	0D	2D	1
06	-	0	0E	11	1
07	-	0	0F	0D	1

### 3. Simple Memory System Cache

- 16 lines, 4-byte block size
- Physically addressed
- Direct mapped



Idx	Tag	Valid	B0	B1	B2	B3
0	19	1	99	11	23	11
1	15	0	-	-	-	-
2	18	1	00	02	04	08
3	36	0	-	-	-	-
4	32	1	43	6D	8F	09
5	0D	1	36	72	F0	1D
6	31	0	-	-	-	-
7	16	1	11	C2	DF	03

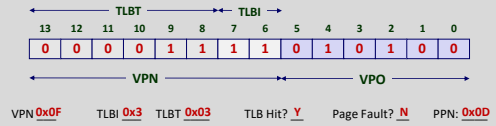
Idx	Tag	Valid	B0	B1	B2	B3
8	24	1	3A	00	51	89
9	2D	0	-	-	-	-
A	2D	1	93	15	DA	38
B	0B	0	-	-	-	-
C	12	0	-	-	-	-
D	16	1	04	96	34	15
E	13	1	83	77	1B	D3
F	14	0	-	-	-	-

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

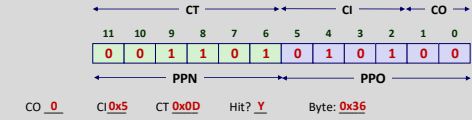
7

### Address Translation Example #1

Virtual Address: 0x03D4



Physical Address



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

8

### Address Translation Example #2

Virtual Address: 0x0020



VPN 0x00 TLBI 0 TLBT 0x00 TLB Hit? N Page Fault? N PPN: 0x28

Physical Address



CO 0 CI 0x8 CT 0x28 Hit? N Byte: Mem

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

9

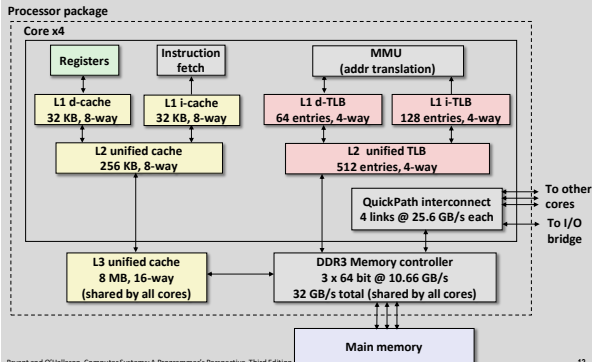
### Today

- Simple memory system example
- Case study: Core i7/Linux memory system
- Memory mapping

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

11

### Intel Core i7 Memory System



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

12

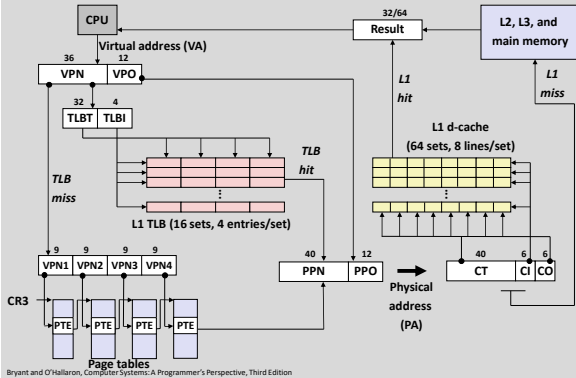
### Review of Symbols

- Basic Parameters**
  - $N = 2^n$ : Number of addresses in virtual address space
  - $M = 2^m$ : Number of addresses in physical address space
  - $P = 2^p$ : Page size (bytes)
- Components of the virtual address (VA)**
  - TLBI: TLB index
  - TLBT: TLB tag
  - VPO: Virtual page offset
  - VPN: Virtual page number
- Components of the physical address (PA)**
  - PPO: Physical page offset (same as VPO)
  - PPN: Physical page number
  - CO: Byte offset within cache line
  - CI: Cache index
  - CT: Cache tag

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

13

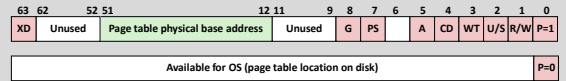
## End-to-end Core i7 Address Translation



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

14

## Core i7 Level 1-3 Page Table Entries



Each entry references a 4K child page table. Significant fields:

**P:** Child page table present in physical memory (1) or not (0).

**R/W:** Read-only or read-write access permission for all reachable pages.

**U/S:** user or supervisor (kernel) mode access permission for all reachable pages.

**WT:** Write-through or write-back cache policy for the child page table.

**A:** Reference bit (set by MMU on reads and writes, cleared by software).

**PS:** Page size either 4 KB or 4 MB (defined for Level 1 PTEs only).

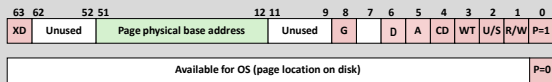
**Page table physical base address:** 40 most significant bits of physical page table address (forces page tables to be 4KB aligned)

**XD:** Disable or enable instruction fetches from all pages reachable from this PTE.

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

15

## Core i7 Level 4 Page Table Entries



Each entry references a 4K child page. Significant fields:

**P:** Child page is present in memory (1) or not (0)

**R/W:** Read-only or read-write access permission for child page

**U/S:** User or supervisor mode access

**WT:** Write-through or write-back cache policy for this page

**A:** Reference bit (set by MMU on reads and writes, cleared by software)

**D:** Dirty bit (set by MMU on writes, cleared by software)

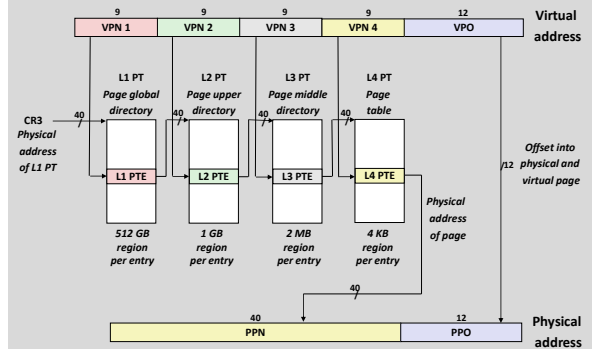
**Page physical base address:** 40 most significant bits of physical page address (forces pages to be 4KB aligned)

**XD:** Disable or enable instruction fetches from this page.

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

16

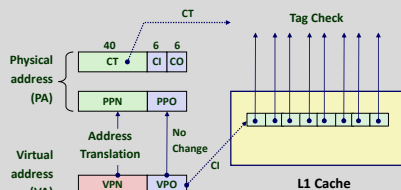
## Core i7 Page Table Translation



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

17

## Cute Trick for Speeding Up L1 Access



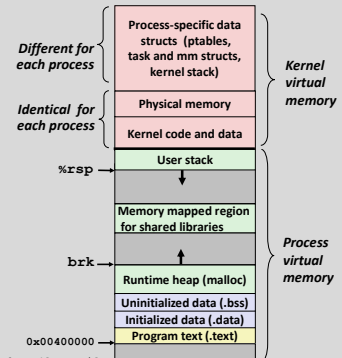
### Observation

- Bits that determine CI identical in virtual and physical address
- Can index into cache while address translation taking place
- Generally we hit in TLB, so PPN bits (CT bits) available next
- "Virtually indexed, physically tagged"
- Cache carefully sized to make this possible

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

18

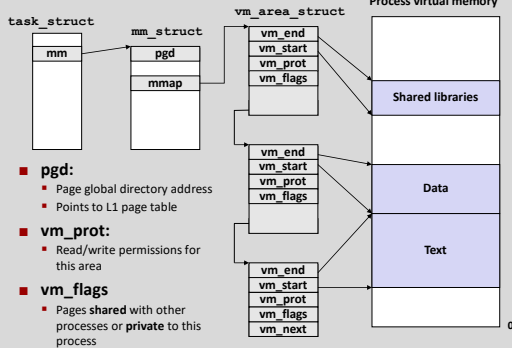
## Virtual Address Space of a Linux Process



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

19

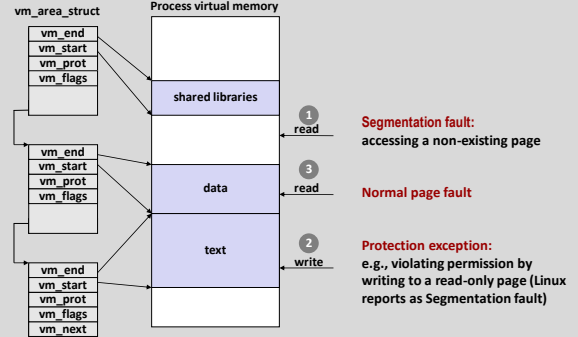
## Linux Organizes VM as Collection of "Areas"



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

20

## Linux Page Fault Handling



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

21

## Today

- Simple memory system example
- Case study: Core i7/Linux memory system
- Memory mapping

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

22

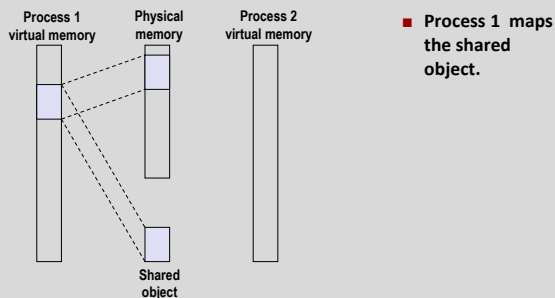
## Memory Mapping

- VM areas initialized by associating them with disk objects.
  - Process is known as **memory mapping**.
- Area can be **backed by** (i.e., get its initial values from) :
  - **Regular file** on disk (e.g., an executable object file)
    - Initial page bytes come from a section of a file
  - **Anonymous file** (e.g., nothing)
    - First fault will allocate a physical page full of 0's (**demand-zero page**)
    - Once the page is written to (**dirty**), it is like any other page
- Dirty pages are copied back and forth between memory and a special **swap file** (or partition).

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

23

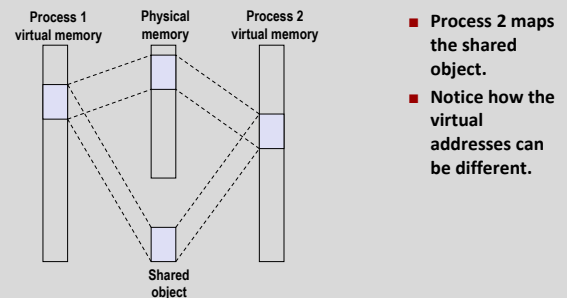
## Sharing Revisited: Shared Objects



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

24

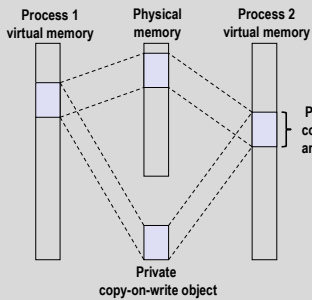
## Sharing Revisited: Shared Objects



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

25

## Sharing Revisited: Private Copy-on-write (COW) Objects

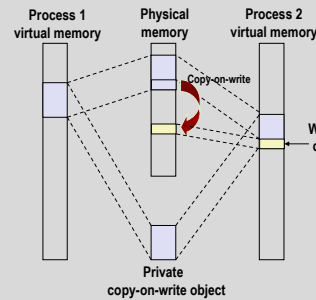


- Two processes mapping a **private copy-on-write (COW)** object.
- Area flagged as private copy-on-write
- PTEs in private areas are flagged as read-only

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

26

## Sharing Revisited: Private Copy-on-write (COW) Objects



- Instruction writing to private page triggers protection fault.
- Handler creates new R/W page.
- Instruction restarts upon handler return.
- Copying deferred as long as possible!

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

27

## User-Level Memory Mapping

```
void *mmap(void *start, int len,
           int prot, int flags, int fd, int offset)
```

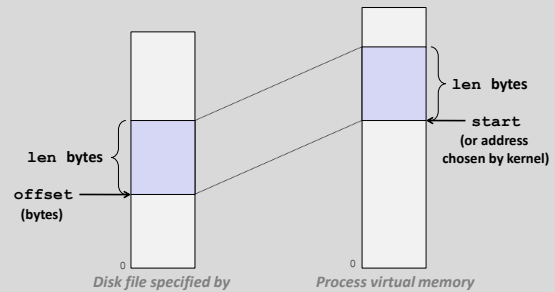
- Map `len` bytes starting at `offset` of the file specified by file descriptor `fd`, preferably at address `start`
  - `start`: may be 0 for "pick an address"
  - `prot`: `PROT_READ`, `PROT_WRITE`, ...
  - `flags`: `MAP_ANON`, `MAP_PRIVATE`, `MAP_SHARED`, ...
- Return a pointer to start of mapped area (may not be `start`)

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

30

## User-Level Memory Mapping

```
void *mmap(void *start, int len,
           int prot, int flags, int fd, int offset)
```



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

31

## Example: Using mmap to Copy Files

- Copying a file to `stdout` without transferring data through other program memory.

```
#include "csapp.h"

void mmapcopy(int fd, int size)
{
    /* Ptr to memory mapped area */
    char *bufp;

    bufp = mmap(NULL, size,
               PROT_READ,
               MAP_PRIVATE,
               fd, 0);
    write(1, bufp, size);
    return;
}

/* mmapcopy driver */
int main(int argc, char **argv)
{
    struct stat stat;
    int fd;

    /* Check for required cmd line arg */
    if (argc != 2) {
        printf("usage: %s <filename>\n",
              argv[0]);
        exit(0);
    }

    /* Copy input file to stdout */
    fd = open(argv[1], O_RDONLY, 0);
    fstat(fd, &stat);
    mmapcopy(fd, stat.st_size);
    exit(0);
}

mmapcopy.c
```

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

32