

CSCI 4061: Virtual Memory

Chris Kauffman

*Last Updated:
Thu Mar 7 09:41:41 CST 2019*

Logistics

Reading

- ▶ Stevens/Rago: Ch 14.8 on `mmap()`
- ▶ Virtual Memory Reference: Bryant/O'Hallaron, Computer Systems. Ch 9 (CSCI 2021)

Goals

- ▶ Finish up File Ops
- ▶ Virtual Memory System
- ▶ Memory Mapped Files

Lab06 on binary files

How did it go?

A2 is Coming

- ▶ Small versioning / backup system
- ▶ Use file stats and recursive directory traversal
- ▶ Use binary file format with `read()` / `write()` / `mmap()`

Exercise: The View of Memory Addresses so Far

- ▶ Every **process** (running program) has some memory, divided into roughly 4 areas (which are...?)
- ▶ Reference different data/variables through their addresses
- ▶ If only a single program could run at time, no trouble: load program into memory and go
- ▶ Running multiple programs gets interesting particularly if they both reference the *same memory location*, e.g. address 1024

PROGRAM 1

```
## load global from 1024
```

```
movq 1024, %rax
```

...

PROGRAM 2

```
# add to global at #1024
```

```
addl %esi, 1024
```

...

- ▶ Is this a **problem** (is there a conflict between these programs?)
- ▶ What are possible solutions?

Answers: The View of Addresses so Far

- ▶ 4 areas of memory are roughly: (1) Stack (2) Heap (3) Globals (4) Text/Instructions
- ▶ Both programs cannot use address #1024 so *may* conflict with each other

PROGRAM 1

```
## load global from 1024
```

```
movq 1024, %rax
```

```
...
```

PROGRAM 2

```
# add to global at #1024
```

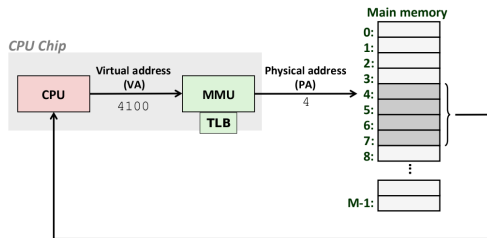
```
addl %esi, 1024
```

```
...
```

- ▶ **Solution 1:** Never let Programs 1 and 2 run together (bleck!)
- ▶ **Solution 2:** Translate every memory address in every program on loading it, run with physical addresses
 - ▶ Tough/impossible as not all addresses are known at compile/load time...
- ▶ **Solution 3:** Translate every memory address/access in every program while it runs (uh-oh)

Addresses are a Lie

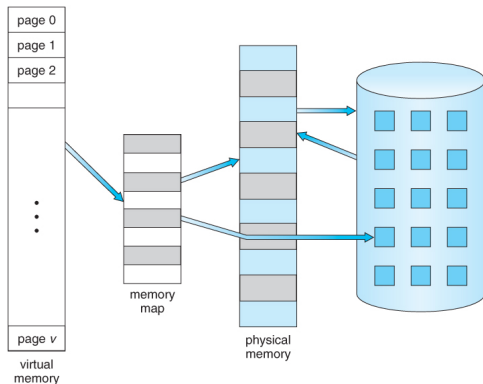
- ▶ Operating system uses tables and hardware to translate every memory address reference **on the fly**
- ▶ Processes know **virtual addresses** which are translated via the memory subsystem to physical addresses in RAM and on disk
- ▶ Translation must be **FAST** so usually involves hardware



- ▶ **MMU (Memory Manager Unit)** is a hardware element specifically designed for address translation
- ▶ Usually contains a special cache, **TLB (Translation Lookaside Buffer)**, which keeps some translated addresses

Address Translation

- ▶ OS maintains tables to translate virtual to physical addresses
- ▶ Contiguous virtual addresses may be spread all over RAM/Disk
- ▶ Address translation is NOT CONSTANT $O(1)$, has an impact on performance of real algorithms*
- ▶ Addresses usually mapped in hunks called **pages** like pages in a notebook



Source: John T. Bell Operating Systems Course Notes

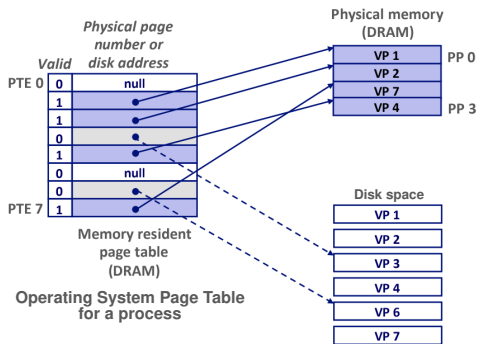
*See: [On a Model of Virtual Address Translation \(2015\)](#)

Paged Memory

- ▶ Physical memory is divided into hunks called **pages**
- ▶ Common page size supported by many OS's (Linux) and hardware MMU's is $4\text{KB} = 4096$ bytes
- ▶ Memory is usually byte addressable so need offset into page
- ▶ 12 bits for offset into page
- ▶ $A - 12$ bits for **page number** where A is the address size in bits
- ▶ Usually A is NOT 64-bits
 - > cat /proc/cpuinfo
 - vendor_id : GenuineIntel
 - cpu family : 6
 - model : 79
 - model name : Intel(R) Xeon(R) CPU E5-1620 v4 @ 3.50GHz
 - ...
 - address sizes : 46 bits physical, 48 bits virtual
- ▶ Leaves one with something like $48 - 12 = 36$ bits for page #s
- ▶ Means a **page table** may have up to 2^{36} entries (!)

Page Tables and Page Table Entries (PTE)

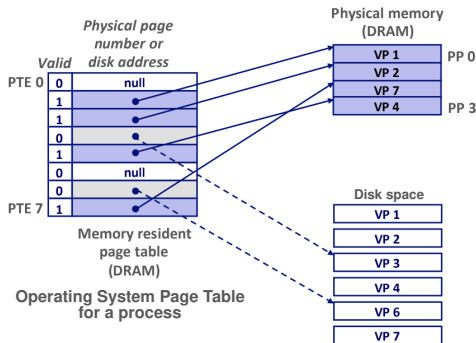
- ▶ OS **Page Table** allows translation from virtual address to physical addresses
- ▶ Each **Page Table Entry (PTE)** contains a physical page number in DRAM or Disk
- ▶ Page table contains all possible addresses for a process and where that data currently exists
- ▶ Some virtual addresses like 0x00 are **unmapped** (null page table entry)



- ▶ Accessing unmapped addresses leads to a **segmentation fault**

Translating Addresses

- ▶ On using a Virtual Memory address, hardware looks it up in the Page Table
- ▶ If valid (hit), address is already in DRAM, translates to physical DRAM address
- ▶ If not valid (miss), address is on disk, move to DRAM potentially evicting a current resident
- ▶ Miss = **Page fault**, notifies OS to move disk data to DRAM

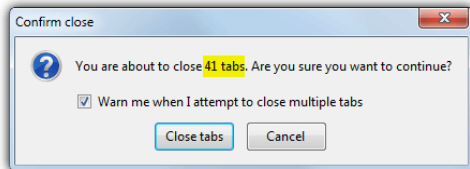


- ▶ Lookup.. Hits.. misses.. this should sound **familiar to..**

Virtual Memory Caches Physical Memory

- ▶ Virtual Memory allows illusion of 2^{48} bytes (hundreds of TBs) of memory when physical memory might only be 2^{30} to 2^{36} (few to hundreds of GBs)
- ▶ Disk space is used for space beyond main memory
- ▶ Pages that are frequently used stay in DRAM (swapped in)
- ▶ Pages that haven't been used for a while end up on disk (**swapped out**)

- ▶ DRAM (physical memory) is then thought of as a cache for Virtual Memory which can be as big as disk space allows



Like when I was writing my composition paper but then got distracted and opened 41 Youtube tabs and when I wanted to write again it took like 5 minutes for Word to load back up because it was swapped out.

The Many Other Advantages of Virtual Memory

- ▶ Caching: Seen that VirtMem can treat main memory as a cache for larger memory
- ▶ Security: Translation allows OS to check memory addresses for validity
- ▶ Debugging: Similar to above, Valgrind checks addresses for validity
- ▶ Sharing Data: Processes can share data with one another by requesting OS to map virtual addresses to same physical addresses
- ▶ Sharing Libraries: Can share same program text between programs by mapping address space to same shared library
- ▶ Convenient I/O: Map internal OS data structures for files to virtual addresses to make working with files free of `read()/write()`

But first...

Exercise: Quick Review

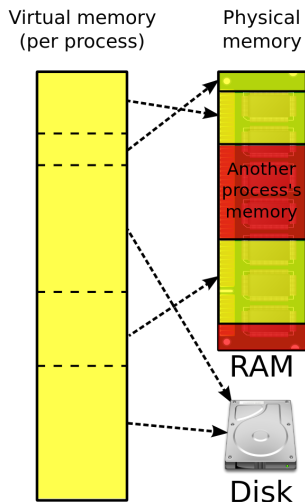
1. While running a program, memory address #1024 always refers to a physical location in DRAM (True/False: why?)
2. Two programs which both use the address #1024 cannot be simultaneously run (True/False: why?)
3. What do MMU and TLB stand for and what do they do?
4. What is a memory page? How big is it usually?
5. What is a Page Table and what is it good for?

Answers: Quick Review

1. While running a program, memory address #1024 always refers to a physical location in DRAM (True/False: why?)
 - ▶ False: #1024 is usually a **virtual address** which is translated by the OS/Hardware to a physical location which *may* be in DRAM but may instead be paged on onto disk
2. Two programs which both use the address #1024 cannot be simultaneously run (True/False: why?)
 - ▶ False: The OS/Hardware will likely translate these identical virtual addresses to **different physical locations** so that the programs do not clobber each other's data
3. What do MMU and TLB stand for and what do they do?
 - ▶ Memory Management Unit: a piece of hardware involved in translating Virtual Addresses to Physical Addresses/Locations
 - ▶ Translation Lookaside Buffer: a special cache used by the MMU to make address translation **fast**
4. What is a memory page? How big is it usually?
 - ▶ A discrete hunk of memory usually 4Kb (4096 bytes) big
5. What is a Page Table and what is it good for?
 - ▶ A table maintained by the operating system that is used to map Virtual Addresses to Physical addresses for each page

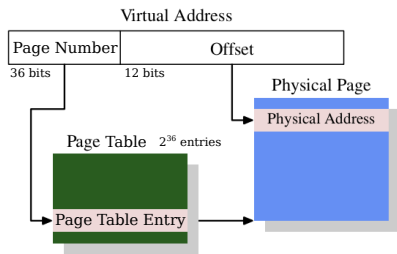
Exercise: Page Table Size

- ▶ Page tables map a virtual page to physical location
- ▶ Maintained by operating system in memory
- ▶ A **direct page** table has one entry per virtual page
- ▶ Each page is $4K = 2^{12}$ bytes, so 12 bits for offset of address into a page
- ▶ Virtual Address Space is 2^{48}
- ▶ **How many** pages of virtual memory are there?
 - ▶ How many bits specify a virtual page number?
 - ▶ How big is the page table? Is this a problem?



How big does the page table mapping virtual to physical

Answers: Page Table Size



"What Every Programmer Should Know About Memory" by Ulrich Drepper, Red Hat, Inc.

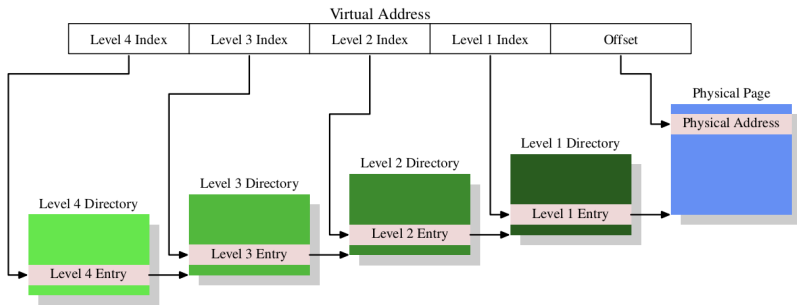
```
48 bits for virtual address
- 12 bits for offset
-----
36 bits for virtual page number
```

So, 2^{36} virtual pages...

- ▶ Every page table entry needs at least 8 bytes for a physical address
- ▶ Plus maybe 8 bytes for other stuff (on disk, permissions)
- ▶ 16 bytes per PTE = 2^4 bytes $\times 2^{36}$ PTEs =
- ▶ 2^{40} = 1 Terabyte of space for the Page Table (!!!)

You've been lying again, haven't you professor...

Page Tables Usually Have Multiple Levels

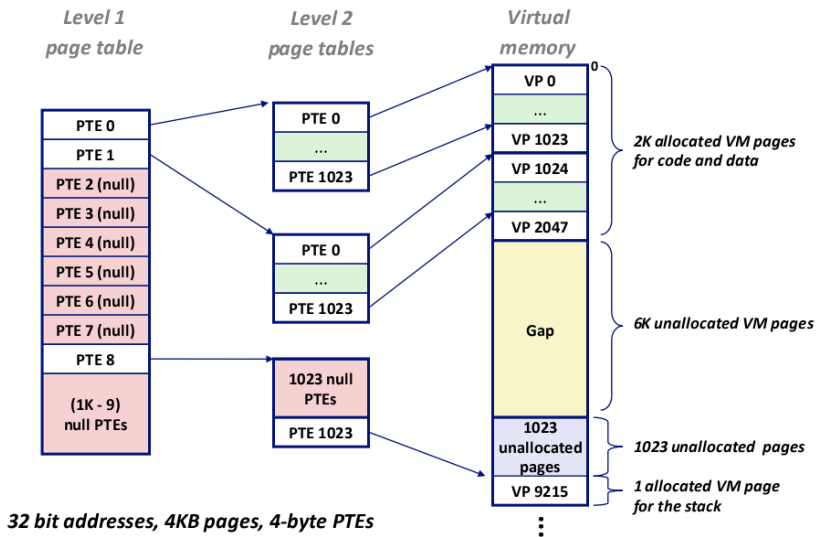


"What Every Programmer Should Know About Memory" by Ulrich Drepper, Red Hat, Inc.

- ▶ Fix this absurdity with **multi-level page tables**: a sparse tree
- ▶ Virtual address divided into sections which indicate which PTE to access at different table levels
- ▶ 3-4 level page table is common in modern architectures
- ▶ Many entries in different levels are NULL (not mapped) most of the 2^{36} virtual pages are not mapped to a physical page (see next diagram)

Textbook Example: Two-level Page Table

Space savings gained via NULL portions of the page table/tree



Source: Bryant/O'Hallaron, CSAPP 3rd Ed

Exercise: Printing Contents of file

1. Write a simple program to print all characters in a file. What are key features of this program?
2. Examine `mmap_print_file.c`: does it contain all of these key features? Which ones are missing?

Answers: Printing Contents of file

1. Write a simple program to print all characters in a file. What are key features of this program?
 - ▶ Open file
 - ▶ Read 1 or more characters into memory using `read()/fscanf()`
 - ▶ Print those characters with `write/printf()`
 - ▶ Read more characters and print
 - ▶ Stop when end of file is reached
 - ▶ Close file
2. Examine `mmap_print_file.c`: does it contain all of these key features? Which ones are missing?
 - ▶ Missing the `read()/fscanf()` portion
 - ▶ Uses `mmap()` to get **direct access** to the bytes of the file
 - ▶ Treat bytes as an array of characters and print them directly

mmap(): Mapping Addresses is Ammazing

- ▶ `ptr = mmap(NULL, size,...,fd,0)` arranges backing entity of `fd` to be mapped to be mapped to `ptr`
- ▶ `fd` often a file opened with `open()` system call

```
int fd = open("gettysburg.txt", O_RDONLY);
// open file to get file descriptor

char *file_chars = mmap(NULL, size, PROT_READ, MAP_SHARED,
                        fd, 0);
// call mmap to get a direct pointer to the bytes in file associated
// with fd; NULL indicates don't care what address is returned;
// specify file size, read only, allow sharing, offset 0

printf("%c",file_chars[0]);           // print 0th char
printf("%c",file_chars[5]);          // print 5th char
```

`mmap()` allows file reads/writes without `read()/write()`

- ▶ Memory mapped files are not just for reading
- ▶ With appropriate options, writing is also possible

```
char *file_chars =  
    mmap(NULL, size, PROT_READ | PROT_WRITE,  
        MAP_SHARED, fd, 0);
```

- ▶ Assign new value to memory, OS writes changes into the file
- ▶ Example: `mmap_tr.c` to transform one character to another

Mapping things that aren't characters

`mmap()` just gives a pointer: can assert type of what it points at

- ▶ Example `int *`: treat file as array of binary ints
- ▶ Notice changing array will write to file

```
// mmap_increment.c

int fd = open("binary_nums.dat", O_RDWR);
// open file descriptor, like a FILE *

int *file_ints = mmap(NULL, size, PROT_READ | PROT_WRITE, MAP_SHARED, fd, 0);
// get pointer to file bytes through mmap,
// treat as array of binary ints

int len = size / sizeof(int);
// how many ints in file

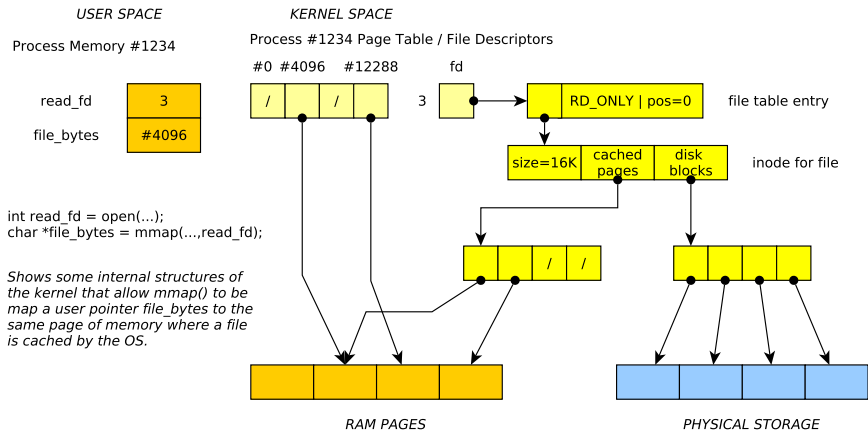
for(int i=0; i<len; i++){
    printf("%d\n",file_ints[i]); // print all ints
}

for(int i=0; i<len; i++){
    file_ints[i] += 1; // increment each file int, writes back to disk
}
```

OS usually Caches Files in RAM

- ▶ For efficiency, part of files are stored in RAM by the OS
- ▶ OS manages internal data structures to track which parts of a file are in RAM, whether they need to be written to disk
- ▶ `mmap()` alters a process Page Table to translate addresses to the cached file page
- ▶ OS tracks whether page is changed, either by file write or `mmap()` manipulation
- ▶ Automatically writes back to disk when needed
- ▶ Changes by one process to cached file page will be seen by other processes
- ▶ **See diagram on next slide**

Diagram of Kernel Structures for mmap()



Exercise: `mmap_tr.c`

Speculate on how to use `mmap()` to write the following program.

```
> gcc -o mmap_tr mmap_tr.c
> mmap_tr gettysburg.txt f p
Transforming 'f' to 'p' in gettysburg.txt
Transformation complete
```

```
> mmap_tr gettysburg.txt F P
Transforming 'F' to 'P' in gettysburg.txt
Transformation complete
```

```
> head gettysburg.txt
```

Pour score and seven years ago our pathers brought porth on this continent, a new nation, conceived in Liberty, and dedicated to the proposition that all men are created equal.

Now we are engaged in a great civil war, testing whether that nation, or any nation so conceived and so dedicated, can long endure. We are met on a great battle-pield op that war. We have come to dedicate a portion op that pield, as pinal resting place por those who here gave their lives that that nation might live. It is altogether pitting and proper that we should do this.

```
>
```

Answers in `mmap_tr.c`

`mmap()` Compared to Traditional `read()/write()` I/O

Advantages of `mmap()`

- ▶ Avoid following cycle
 - ▶ `read()/fread()/fscanf()` file contents into memory
 - ▶ Analyze/Change data
 - ▶ `write()/fwrite()/fscanf()` write memory back into file
- ▶ Saves memory and time
- ▶ Many Linux mechanisms backed by `mmap()` like processes sharing memory

Drawbacks of `mmap()`

- ▶ Always maps **pages** of memory: multiple of 4096b (4K)
- ▶ For small maps, lots of wasted space
- ▶ Cannot change size of files with `mmap()`: must used `write()` to extend or other calls to shrink
- ▶ No bounds checking, just like everything else in C

One Page Table Per Process

- ▶ OS maintains a page table for each running program (**process**)
- ▶ Each process believes its address space ranges from 0x00 to 0xBIG (0 to 2^{48}), its virtual address space
- ▶ Virtual addresses are mapped to physical locations in DRAM or on Disk via page tables

Physical Memory	
00x	H E L L
01x	R L D !
02x	0 W O
03x	H A V E
04x	F U N
05x	L O T
06x	S O F
07x	; -)

Process A	
Page Table	Virtual Memory
00x 00	00x H E L L
01x 02	01x 0 W O
02x 01	02x R L D !
03x n.a.	03x #####
04x n.a.	04x #####
05x 07	05x ; -)

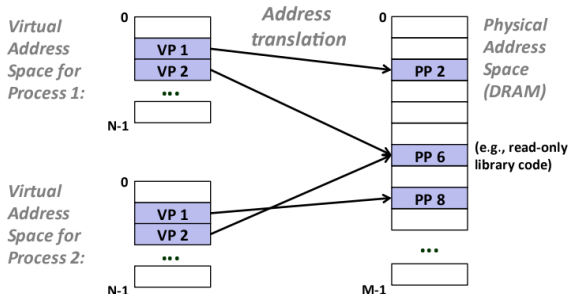
Process B	
Page Table	Virtual Memory
00x 03	00x H A V E
01x 05	01x L O T
02x 06	02x S O F
03x 04	03x F U N
04x n.a.	04x #####
05x 07	05x ; -)

Source: OSDev.org

*Two processes with their own page tables. Notice how contiguous virtual addresses are mapped to non-contiguous spots in physical memory. Notice also the **sharing** of a page.*

Pages and Mapping

- ▶ Memory is segmented into hunks called **pages**, 4Kb is common (use `page_size.c` to see your system's page size)
- ▶ OS maintains tables of which pages of memory exist in RAM, which are on disk
- ▶ OS maintains tables per process that translate process virtual addresses to physical pages
- ▶ **Shared Memory** can be arranged by mapping virtual addresses for two processes to the same memory page



Shared Memory Calls

- ▶ Using OS system calls, can usually create shared memory
- ▶ Unix System V (five) IPC includes the following

```
key_t key = ftok("crap", 'R');  
// make the SysV IPC key
```

```
int shmid = shmget(key, 1024, 0644 | IPC_CREAT);  
// connect to (and possibly create) the segment:
```

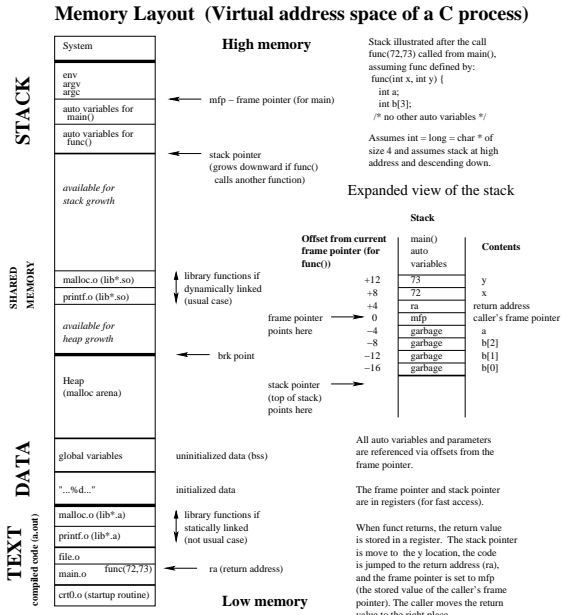
```
char *data = shmat(shmid, (void *)0, 0);  
// attach to the segment to get a pointer to it:
```

- ▶ Multiple processes can all "see" the same unit of memory
- ▶ This is an old style but still useful
- ▶ Will cover more **interprocess communication (IPC)** later
- ▶ Modern incarnations favor `mmap()` followed by `fork()`

Exercise: Process Memory Image and Libraries

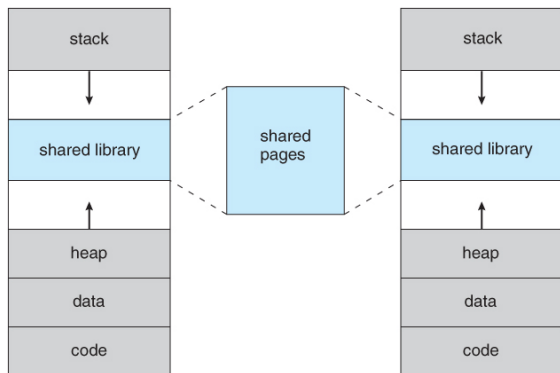
- ▶ How many programs on the system need to use malloc() and printf()?
- ▶ Where is the code for malloc() or printf() in the process memory?

Right: A detailed picture of the virtual memory image, by Wolf Holzman



Shared Libraries: *.so Files

- ▶ Code for libraries can be shared
- ▶ `libc.so`: shared library with `malloc()`, `printf()` etc in it
- ▶ OS puts into one page, maps all linked procs to it



Source: John T. Bell Operating Systems Course Notes

pmap: show virtual address space of running process

```
> ./memory_parts
0x5579a4cbe0c0 : global_arr
0x7fff96aff6f0 : local_arr
0x5579a53aa260 : malloc_arr
0x7f441f8bb000 : mmap'd file
my pid is 7986
press any key to continue
```

- ▶ While a program is running, determine its **process id**
- ▶ Call pmap to see how its virtual address space maps
- ▶ For more details of pmap output, refer to [this diagram](#) from a now defunct article by Andreas Fester

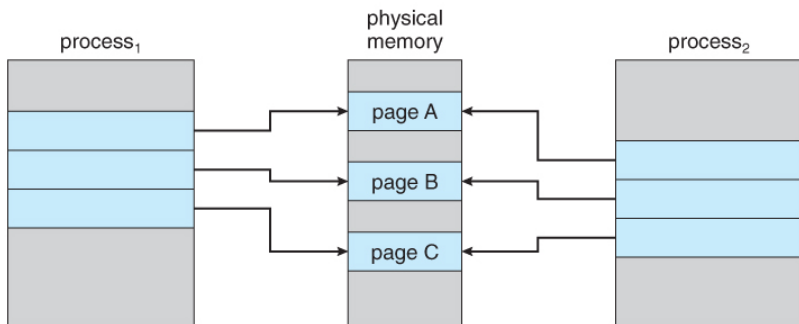
```
> pmap 7986
7986:  ./memory_parts
00005579a4abd000      4K r-x-- memory-parts
00005579a4cbd000      4K r---- memory-parts
00005579a4cbe000      4K rw--- memory-parts
00005579a4cbf000      4K rw--- [ anon ]
00005579a53aa000    132K rw--- [ heap ]
00007f441f2e1000   1720K r-x-- libc-2.26.so
00007f441f48f000   2044K ----- libc-2.26.so
00007f441f68e000     16K r---- libc-2.26.so
00007f441f692000      8K rw--- libc-2.26.so
00007f441f694000     16K rw--- [ anon ]
00007f441f698000    148K r-x-- ld-2.26.so
00007f441f88f000      8K rw--- [ anon ]
00007f441f8bb000      4K r--s- gettysburg.txt
00007f441f8bc000      4K r---- ld-2.26.so
00007f441f8bd000      4K rw--- ld-2.26.so
00007f441f8be000      4K rw--- [ anon ]
00007fff96ae1000   132K rw--- [ stack ]
00007fff96b48000     12K r---- [ anon ]
00007fff96b4b000      8K r-x-- [ anon ]
total                    4276K
```


Memory Protection

- ▶ Output of `pmap` indicates another feature of virtual memory: protection
- ▶ OS marks pages of memory with Read/Write/Execute/Share permissions
- ▶ Attempt to violate these and get segmentation violations (segfault)
- ▶ Ex: Executable page (instructions) usually marked as `r-x`: no write permission.
- ▶ Ensures program don't accidentally write over their instructions and change them
- ▶ Ex: By default, pages are not shared (no `s` permission) but can make it so with the right calls

Fork and Shared Pages

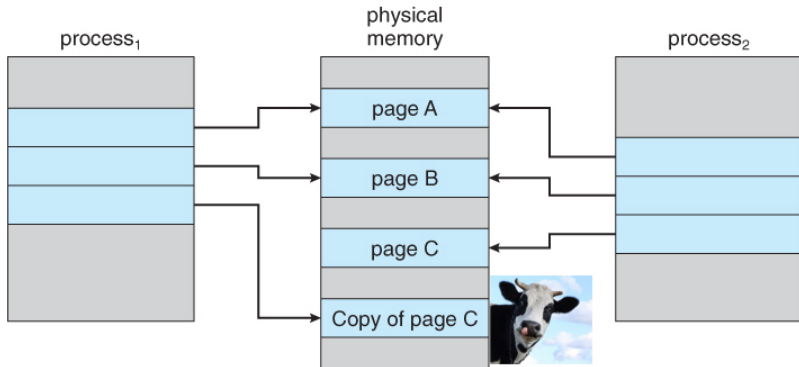
- ▶ `fork()`'ing a process creates a nearly identical copy of a process
- ▶ Might need to copy all memory from parent to child pages
- ▶ Can save a lot of time if memory pages of child process are **shared with parent** - no copying needed (initially)
- ▶ What's the major danger here?



Source: John T. Bell Operating Systems Course Notes

Fork, Shared Pages, Copy on Write (COW Pages)

- ▶ If neither process writes to the page, sharing doesn't matter
- ▶ If either process writes, OS will make a copy and remap addresses to copy so it is exclusive
- ▶ Fast if hardware Memory Management Unit and OS know what they are doing (Linux + Parallel Python/R + Big Data)



Source: John T. Bell Operating Systems Course Notes

Summary

A computer using a Virtual Memory system sees the OS and hardware cooperate to translate every program address from a virtual address to a physical RAM address.

- ▶ **Consequence 1:** All programs see a linear address space but each has its own #1024 which does not conflict
- ▶ **Consequence 2:** The OS must maintain Page Table data structures that map each process's virtual addresses to physical locations
- ▶ **Consequence 3:** Computers with small amounts of RAM can "fake" larger amounts by using disk space; RAM serves as a cache for this larger virtual memory
- ▶ **Consequence 4:** Every address deference is translated so the OS can catch out of bounds references or direct the dereference to special RAM areas like memory mapped files
- ▶ **Consequence 5:** Memory can shared and manipulated for efficiency such as `fork()` using Copy on Write