# CS 3204
# Operating Systems

Lecture 19

Godmar Back

Virginia Tech

---

## Announcements

---

# Virtual Memory

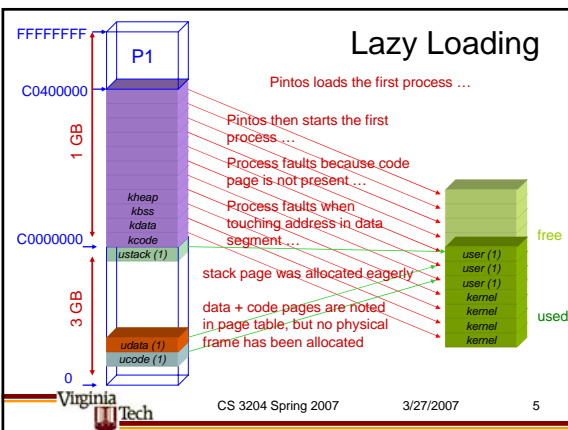Paging Techniques (cont'd)

Virginia Tech

---

## Demand paging

- Idea: only keep data in memory that's being used
  - Needed for virtualization – don't use up physical memory for data processes don't access
- Requires that actual allocation of physical page frames be delayed until first access
- Many variations
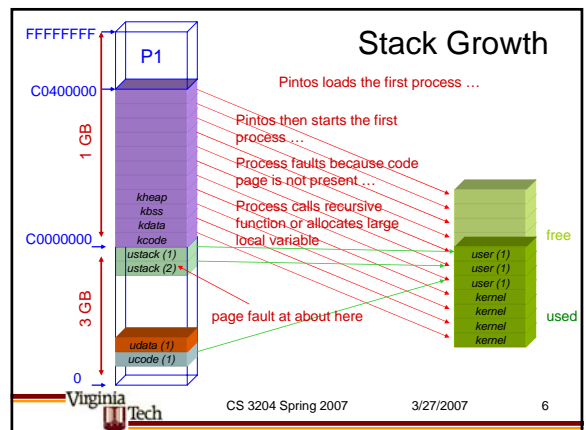  - Lazy loading of text & data, mmapped pages & newly allocated heap pages
  - Copy-on-write

---

## Lazy Loading



FFFFFFFF
C0400000
1 GB
C0000000
3 GB
0

P1

kheap
kbss
kdata
kcode
ustack (1)

udata (1)
ucode (1)

Pintos loads the first process …

Pintos then starts the first process …

Process faults because code page is not present …

Process faults when touching address in data segment …

stack page was allocated eagerly

data + code pages are noted in page table, but no physical frame has been allocated

free

user (1)
user (1)
kernel
kernel
kernel
kernel

used

---

## Stack Growth



FFFFFFFF
C0400000
1 GB
C0000000
3 GB
0

P1

kheap
kbss
kdata
kcode
ustack (1)
ustack (2)

udata (1)
ucode (1)

Pintos loads the first process …

Pintos then starts the first process …

Process faults because code page is not present …

Process calls recursive function or allocates large local variable

page fault at about here

free

user (1)
user (1)
user (1)
kernel
kernel
kernel
kernel

used

1

## mmap()

FFFFFFFF
P1
C0400000

1 GB

Pintos loads the first process …

Pintos then starts the first process …

Process faults because code page is not present …

kheap
kbss
kdata
kcode

Process opens file, calls mmap(fd, addr)

C0000000

ustack (1)

3 GB

ummap (1)

Process faults when touching mapped file

Page fault handler allocs page, maps it, reads data from disk.

udata (1)
ucode (1)

0

free

user (1)
user (1)
user (1)
kernel
kernel
kernel
kernel
kernel

used

---
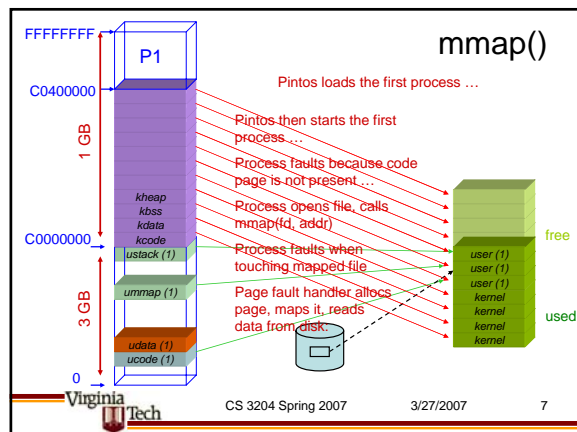
## Lazy Loading & Prefetching

- Typically want to do some prefetching when faulting in page
  - Reduces latency on subsequent faults
- Q.: how many pages? which pages?
  - Too much: waste time & space fetching unused pages
  - Too little: pay (relatively large) page fault latency too often
- Predict which pages the program will access next (how?)
- Let applications give hints to OS
  - If applications knows
  - Example: madvise(2)
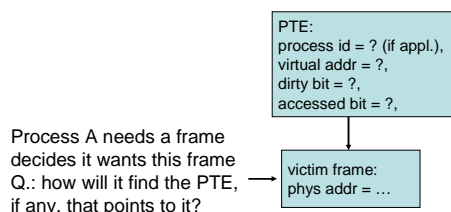  - Usual conflict: what's best for application vs what's best for system as a whole

---

## Copy-On-Write

- Sometimes, want to create a copy of a page:
  - Example: Unix fork() creates copies of all parent's pages in the child
- Optimization:
  - Don't copy pages, copy PTEs – now have 2 PTEs pointing to frame
  - Set all PTEs read-only
  - Read accesses succeed
  - On Write access, copy the page into new frame, update PTEs to point to new & old frame
- Looks like each have their own copy, but postpone actual copying until one is writing the data
  - Hope is at most one will ever touch the data – never have to make actual copy

---

## Page Eviction

- Suppose page fault occurs, but no free physical frame is there to allocate
- Must evict frame
  - Find victim frame (how – later)
  - Find & change old page table entry pointing to the victim frame
  - If data in it isn't already somewhere on disk, write to special area on disk ("swap space")
  - Install in new page table entry
  - Resume
- Requires check on page fault if page has been swapped out – fault in if so
- Some subtleties with locking:
  - How do you prevent a process from writing to a page some other process has chosen to evict from its frame?
  - What do you do if a process faults on a page that another process is in the middle of paging out?

---

## Page Eviction Example

PTE:
process id = ? (if appl.),
virtual addr = ?,
dirty bit = ?,
accessed bit = ?,

Process A needs a frame decides it wants this frame
Q.: how will it find the PTE, if any, that points to it?

victim frame:
phys addr = …

Linux uses a so-called "rmap" for that that links frames to PTE

---

## Managing Swap Space

- Continuous region on disk
  - Preferably on separate disk, but typically a partition on same disk
- Different allocation strategies are possible
  - Simplest: when page must be evicted, allocate swap space for page; deallocate when page is paged back in
  - Or: allocate swap space upfront
  - Should page's position in swap space change? What if same page is paged out multiple times?
- Can be managed via bitmap 0100100000001
  - Free/used bits for each page that can be stored
  - Pintos: note 1 page == 8 sectors

## Locking Frames

- Aka "pinned" or "wired" pages or frames
- If another device outside the CPU (e.g., DMA by network controller) accesses a frame, it cannot be paged out
  – Device driver must tell VM subsystem about this
- Also useful if you want to avoid a page fault while kernel code is accessing a user address, such as during a system call.

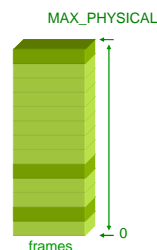## Accessing User Pointers & Paging

- Kernel must check that user pointers are valid
  – P2: easy, just check range & page table
- Harder when swapping:
  – validity of a pointer may change between check & access (if another process sneaks in and selects frame mapped to an already checked page for eviction)
- Possible solution:
  – verify & lock, then access, then unlock

```
if (verify_user(addr))
    process_terminate();
// what if addr's frame is just now
// swapped out by another process?
*addr = value;
```
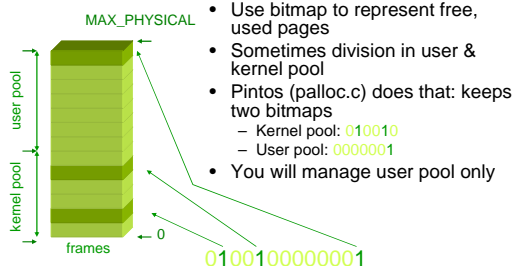
## Physical Memory Management

## Physical Memory Management

- Aka frame table management
- Task: keep efficiently track of which physical frames are used
- Allocate a frame when paging in, or eager loading
- Deallocate a frame when process exits or when page is evicted (later)

MAX_PHYSICAL

frames

## Approach 1: Bitmaps

MAX_PHYSICAL

user pool

kernel pool

frames

- Use bitmap to represent free, used pages
- Sometimes division in user & kernel pool
- Pintos (palloc.c) does that: keeps two bitmaps
  – Kernel pool: 010010
  – User pool: 0000001
- You will manage user pool only

010010000001

## Approach 2: Buddy Allocator

- Logically subdivide memory in power-of-two blocks
- Round up on allocation to next power of 2
- Split block on allocation (if necessary)

1M
512K
256K
128K
64K

A = 128 K  c = sha  64 K      256 K      D = 256 K      256 K

- Coalesce on deallocation (if possible)
  – Coalescing can be delayed
- Used in Linux: allocation requests are always multiple of pages, max blocksize is 4MB

3

## Buddy Example - Allocation

| 64 KB |
|---|

Alloc (16KB)

| 16KB | 16KB | 32KB |
|---|---|---|

Alloc (32KB)

| 16KB | 16KB | 32KB |
|---|---|---|

Alloc (4KB)

| 16KB | 4KB | 4KB | 8KB | 32KB |
|---|---|---|---|---|

Alloc (4KB)

| 16KB | 4KB | 4KB | 8KB | 32KB |
|---|---|---|---|---|

Alloc (4KB)

| 16KB | 4KB | 4KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|---|

---

## Buddy Example - Deallocation

| 16KB | 4KB | 4KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|---|

Free()

| 16KB | 4KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|

Free()

| 16KB | 8KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|

Free()

| 16KB | 16KB | 32KB |
|---|---|---|

Free()

| 32KB | 32KB |
|---|---|

Free()

| 64 KB |
|---|

---

## Fragmentation

- Def: *The inability to use memory that is unused.*
- Internal fragmentation:
  - Not all memory inside an allocated unit is used; rest can't be allocated to other users
- External fragmentation:
  - Impossible to satisfy allocation request even though total amount of memory > size requested
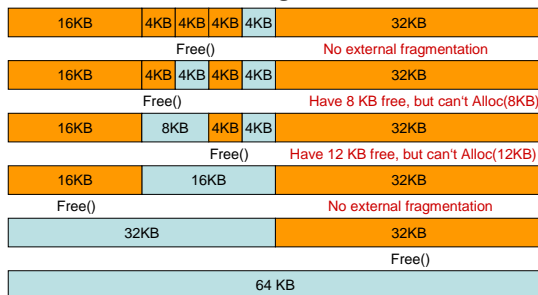
---

## Internal Fragmentation

| 64 KB |
|---|

Alloc (12KB)

| 12KB | | 16KB | 32KB |
|---|---|---|---|

Alloc (24KB)

| 12KB | | 16KB | 24KB | |
|---|---|---|---|---|

Alloc (4KB)

| 12KB | | 4KB | 4KB | 8KB | 24KB | |
|---|---|---|---|---|---|---|

Alloc (4KB)

| 12KB | | 4KB | 4KB | 8KB | 24KB | |
|---|---|---|---|---|---|---|

Alloc (3KB)

| 12KB | | 4KB | 4KB | 3KB | 4KB | 24KB | |
|---|---|---|---|---|---|---|---|

---

## External Fragmentation

| 16KB | 4KB | 4KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|---|

Free() — No external fragmentation

| 16KB | 4KB | 4KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|---|

Free() — Have 8 KB free, but can't Alloc(8KB)

| 16KB | 8KB | 4KB | 4KB | 32KB |
|---|---|---|---|---|

Free() — Have 12 KB free, but can't Alloc(12KB)

| 16KB | 16KB | 32KB |
|---|---|---|

Free() — No external fragmentation

| 32KB | 32KB |
|---|---|

Free()

| 64 KB |
|---|

---

## Buddy Allocator & Fragmentation

- Q.: what is the average internal fragmentation (per allocated object) for
  - buddy allocator with size $2^n$?
  - in bitmap allocator for objects of size n*s, where each bit represents a unit of size s?
  - in first-fit allocator from project 0?
- Q.: what external fragmentation can you expect from buddy allocator scheme?
- Q.: what's a good way to measure fragmentation in general?

# Page Size & Fragmentation

- How should a system's architect choose the page size? – Trade-Off
- Large pages:
  - Larger internal fragmentation
  - (not an issue if most pages are full…)
  - Page-in & write-back cost larger
- Small pages:
  - Higher overhead to store page table (more entries to maintain)
- Modern architectures provide support for "super pages" – 2MB or 4MB