# CS202 (003): Operating Systems Virtual Memory III, Weensy OS

Instructor: Jocelyn Chen

Most of the materials covered in this slide come from the lecture notes of Mike Walfish's CS202



### Last Time

### Page faults

A reference is illegal, either because it's not mapped in the page tables or because there is a protection violation.

This is a quite powerful mechanism! (It turns out you can build interesting functionalities by triggering page faults)

## How does OS get involved in page fault (in x86)?

Process constructs a trap frame and transfer execution to an interrupt/trap handler



%rip now points to the code handle the trap (using Interrupt Descriptor Table)

%cr2 holds the faulting virtual address

When page fault happens, the **kernel** sets up the process's page entries properly, or terminates the process

U/S: user mode fault / supervisor mode fault R/W: access was read / access was write P: not-present page / protection violation



31	15       5       4       3       2       1       0         Reserved       SQ       Reserved       P       SQ       P
Ρ	<ol> <li>The fault was caused by a non-present page.</li> <li>The fault was caused by a page-level protection violation.</li> </ol>
W/R	<ol> <li>The access causing the fault was a read.</li> <li>The access causing the fault was a write.</li> </ol>
U/S	<ol> <li>A supervisor-mode access caused the fault.</li> <li>A user-mode access caused the fault.</li> </ol>
RSVD	<ol> <li>The fault was not caused by reserved bit violation.</li> <li>The fault was caused by a reserved bit set to 1 in some paging-structure entry.</li> </ol>
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.
SGX	<ol> <li>The fault is not related to SGX.</li> <li>The fault resulted from violation of SGX-specific access-control requirements.</li> </ol>

.

#### Figure 4-12. Page-Fault Error Code

. .

.

### When does page fault occur?

Overcommitting physical memory

"Your program thinks it has 64GB of memory, but your hardware has 16 GB of physical memory"

### How does this work? Disk was (is) used to store memory pages

Advantages: address space looks huge Disadvantages: access to "paged" memory (as disk pages that live on the disk are known) are **slow** 

### **Rough Implementation**

On a page fault, the kernel reads in the faulting page. It may need to send a page to disk (when satisfy the following TWO):

1. kernel is out of memory

2. the page that it selects to write out is dirty)



### What are some other use cases of page fault?

Store memory pages across the network (Distributed Shared Memory)

Copy-on-write (fork, mmap, ...) When creating a copy of another process, don't copy its memory. Just copy its page tables, mark the pages as read-only

When a write happens, a page fault results. at that point, the kernel allocates a new page, copies the memory over, and restarts the user program to do a write Then, only do copies of memory when there is a fault as a result of a write

#### Accounting

Good way to sample what percentage of the memory pages are written to in any time slice: mark a fraction of them not present, see how often you get faults

On a page fault, the page fault handler went and retrieved the needed page from some other machine









## Paging in day-to-day use

### Demand paging

Program code is loaded into memory only when it's needed, not all at once

Growing the stack

The seemingly contiguous virtual memory can scatter across different locations in physical memory

BSS page allocation (Block Started by Symbol)

The OS can save memory by not allocating physical pages for the BSS until the program actually tries to use variables in this segment.

Shared text

Sharing the read-only parts of a program between multiple processes running the same program

Shared libraries

Multiple programs can use the same library code in memory, saving space

Shared libraries

Allowing multiple processes to access the same memory region

What does paging from the disk cost?



 $1.1 * t_M > (1-p) * t_M + p * t_D$ What does *p* need to be to ensure that  $p = 0.1 * \frac{t_M}{t_D - t_M} \approx \frac{10^1 ns}{10^7 ns} = 10^{-6}$ paging hurts performance by less than 10%?

> Page faults are **super-expensive**! "need to pay attention to the slow case if it's really slow and common enough to matter."

### Costs of page faults

(1-p) \* memory\_access\_time + p \* page\_fault\_time

#### where p is the prob of a page fault

disk\_access\_time(t<sub>D</sub>)  $\approx 10$ ms = 10<sup>7</sup>ms memeory\_access\_time( $t_M$ )  $\approx 100$ ns





## A Cache System

### A Cache System

any system that temporarily stores frequency used data the cache itself is smaller than the the storage it is cached on

- 2. If the cache is full, decide which existing entry to evict to make room

#### Cache Miss

the requested data isn't in the cache 1. fetch the missing data from the slower main storage

How to decide which entry to throw away if we get a cache miss?

## VM as a Cache System

### A Cache System

any system that temporarily stores frequency used data the cache itself is smaller than the the storage it is cached on

• Virtual memory is an **abstraction** that provides programs with the illusion of a large, contiguous memory space • Physical RAM is typically much smaller than the virtual address space • The operating system keeps only a subset of all pages (fixed-size blocks of memory) in physical RAM at any given time • The rest of the pages are **stored on disk** (in the swap space or paging file)

How to decide which page to throw away if we get a 'page-not-present in memory' fault?

#### FIFO

MIN (optimal)

throw out the oldest

throw away the entry that won't be used for the longest time

### LRU

throw out the least recently used



#### FIFO

throw out the oldest

throw away the entry that won't be used for the longest time

How do we evaluate these algorithms?

Reference string (sequence of page accesses) Input: Cache size (i.e. physical memory)

# of cache evictions (i.e. number of swaps) **Output:** 

MIN (optimal)

LRU

throw out the least recently used



#### FIFO

#### throw out the oldest

Α	В	С	Α	В	D	Α	D	В	С	В
А	А	А	А	А	D	D	D	D	С	С
-	В	В	В	В	В	А	А	А	А	А
_	-	С	С	С	С	С	С	В	В	В

Number of Hits: 4

Page Faults: 7

Hit Rate: 36.36%

MIN (optimal)

throw away the entry that won't be used for the longest time



Page Faults: 5

В	D	Α	D	В	С	В
А	А	А	А	А	А	А
В	В	В	В	В	В	В
С	D	D	D	D	С	С

Number of Hits: 6

Hit Rate: 54.55%

### LRU

#### throw out the least recently used

Α	В	С	A	В	D	А	D	В	С	В
A	А	А	В	С	А	В	В	А	D	D
-	В	В	С	А	В	D	А	D	В	С
-	_	С	А	В	D	А	D	В	С	В

- Number of Hits: 6
- Page Faults: 5
- Hit Rate: 54.55%



#### FIFO

#### throw out the oldest



Number of Hits: 0

Page Faults: 12

Hit Rate: 0.0%

MIN (optimal)

throw away the entry that won't be used for the longest time



Page Faults: 6

4	В	С	D	Α	В	С	D
4	А	А	А	А	В	В	В
3	В	С	С	С	С	С	С
)	D	D	D	D	D	D	D

Number of Hits: 6

Hit Rate: 50.0%

### LRU

#### throw out the least recently used



Number of Hits: 0

Page Faults: 12

Hit Rate: 0.0%



## Replacement policy (adding new memory)

FIFO

throw out the oldest



Number of Hits: 3

Page Faults: 9

Hit Rate: 25.0%

MIN (optimal)

throw away the entry that won't be used for the longest time



Page Faults: 7

Ą	В	E	A	В	С	D	Е
Ą	А	А	А	А	А	А	А
В	В	В	В	В	С	D	D
D	D	Е	Е	Е	Е	Е	Е

Number of Hits: 5

Hit Rate: 41.67%

### LRU

throw out the least recently used



Number of Hits: 2

Page Faults: 10

Hit Rate: 16.67%



#### FIFO

MIN (optimal)

throw out the oldest

throw away the entry that won't be used for the longest time

### LRU

throw out the least recently used

#### Pretty decent!

It approximates OPT when: principle of temporal locality holds strongly





### Implementing LRU

In hardware, it's a lot of work to timestamp each reference and keep the list ordered

Implementing LRU in OS/hardware is a lot of pain!

In OS, it doubles the memory traffic (since after every reference, have to move some structure to the head of some list)

**Approximating LRU\*** 

Periodically, sweep through all pages

- Used? Clear use bit
- Unused? reclaim
  - update core map
  - invalidate page table
  - write back if dirty
  - TLB shootdown
  - add to free list

(\*yes, LRU was already an approximation...)







## Generalizing CLOCK: Nth Chance

- On page fault, OS checks accessed bit: • If 1, then clear it, and also clear the counter.  $\circ$  If 0, then increment the counter; if count == N, replace page.

Large N implies better approximation to LRU: e.g., N = 1000 is a very good LRU approximation. However, a large N implies more work by the OS before a page can be replaced.

Decent approximations to LRU, assuming that past is a good predictor of the future

N = 1 implies the default clock algorithm.

• With each page, OS maintains a counter to indicate the number of sweeps that page has gone through.

https://www.cse.iitd.ernet.in/~sbansal/os/lec/l30.html#:~:text=Nth%20chance%3A%20The%20clock%20algorithm,N%20chances%20before%20evicting%20it.



Process requires more memory than the system has

Each time a page is brought in, another page, whose contents will soon be referenced, is thrown out

A program touches 50 pages (each equally likely) but only have 40 physical page frames

If we have enough physical pages, 100ns/ref If we have too few physical pages, assuming every 5th reference leads to a page fault, then: 4 ref \* 100 ns + 1 page fault \* 10ms for disk I/O

- - This lead to 5 refs per (10ms + 400ns) ~ 2ms/ref = **20,000x slowdown!**

### Thrashing

Process requires more memory than the system has

Each time a page is brought in, another page, whose contents will soon be referenced, is thrown out

What we want: virtual memory the size of disk with access time the speed of physical memory

What we have: memory with access time roughly at the same magnitude as disk access

**Note:** this issue is not limited to page access, but we are discussing this issue in the context of page access

### Thrashing

## Thrashing - what are the causes?

What we have: memory with access time roughly at the same magnitude as disk access

- What we want: virtual memory the size of disk with access time the speed of physical memory

process don't reuse memory (no temporal locality) OR process reuses memory but the memory that is absorbing most of the accesses doesn't fit

Each processes fit the memory individually, but too much to fit for all processes in the system!

## Thrashing - What do we do?

Each processes fit the memory individually, but too much to fit for all processes in the system!

Working Set

### The pages a process has touched over some trailing window of time

Only run a set of processes s.t. the union of their working sets fit in memory

Page fault frequency

### Track the metric (# page faults/instructions executed)

If that thing rises above a threshold, and there is not enough memory on the system, swap out the process

# Lab 4: Weensy OS (Yes, it is already released)

In Lab 4, you will write a mini OS, WeensyOS, that implements the virtual memory architecture and a few important system calls.

## Weensy OS structure

sys\_page\_alloc() for process allocating memory

(sys\_page\_alloc is analogous to brk() or mmap() in POSIX)

exception\_return() for when returning back into user space

%rax is what the application return value is

\* Process registers, process state \* Process page table - a pointer (kernel virtual address, which is the identical physical address) \* to an L1 page table L1 page table's first entry points to a page table, and so on...

virtual\_memory\_lookup(): lookup a physical page using pagetable and virtual memory.

Kernel Code

Processes

Files with p-\*

Files with k-\*

virtual\_memory\_map(): map virtual address -> physical address

#### look at process.h for

#### look at kernel.h for process control block (PCB): struct proc

typedef struct physical\_pageinfo { int8 t owner; //kernel, reserved, free, pid int8 t refcount; physical\_pageinfo;

static physical pageinfo pageinfo [PAGENUMBER(MEMSIZE PHYSICAL)];

// one physical\_pageinfo struct per \_physical\_ page

pageinfo array

## Weensy OS Memory Related

WeensyOS begins with the kernel and all processes sharing a single address space. This is defined by the kernel\_pagetable.

Kernel's pagetable is identity-mapped: Virtual address X maps to physical address X. As you work through the project, you will shift processes to use independent address space where each process can access only a subset of physical memory.

The OS supports 3MB of virtual memory on top of 2MB of physical memory. (Recall the point of virtualization, from the perspective of the process, it thinks it has 3MB of memory. But in reality, it doesn't.)

Assume page size to be 4KB, each entry in the page table is 64 bit. How to we support 3MB of virtual memory? How many L4 pagetable do we need? (2 L4 page tables)

## Weensy OS Macros and Constants

Macro	Meaning	Constant	Meaning	
PAGESIZE	Size of a memory page. Equals 4096 (or, equivalently, 1 << 12).	KERNEL_START_AD DR	Start of kernel code.	
PAGENUMBER(addr)	Page number for the page containing addr. Expands to an expression analogous to addr / PAGESIZE.	KERNEL_STACK_TO P	Top of kernel stack. The kernel stack is one page long.	
PAGEADDRESS(pn)	The initial address (zeroth byte) in page number pp. Expands to an	console	Address of CGA console memory.	
	expression analogous to pn * PAGESIZE.	PROC_START_ADD	Start of application code. Applications should not be a	
PAGEINDEX(addr, level)	The index in the levelth page table for addr. level must be between 0 and 3; 0 returns the level-1 page table index (address bits 39–47),		at console.	
	1 returns the level-2 index (bits 30–38), 2 returns the level-3 index (bits 21–29), and 3 returns the level-4 index (bits 12–20).	MEMSIZE_PHYSICA L	Size of physical memory in bytes. WeensyOS does not a physical addresses ≥ this value. Defined as 0x200000 (2)	
PTE_ADDR(pe)	The physical address contained in page table entry pe. Obtained by masking off the flag bits (setting the low-order 12 bits to zero).	MEMSIZE_VIRTUAL	Size of virtual memory. WeensyOS does not support virt addresses $\geq$ this value. Defined as 0x300000 (3MB).	





# Bring your questions next Tuesday!