# CS194-24
# Advanced Operating Systems Structures and Implementation
# Lecture 11

## TLBs, SLAB allocator
## File Systems

March 3st, 2014

Prof. John Kubiatowicz

http://inst.eecs.berkeley.edu/~cs194-24

---

## Goals for Today

- TLBs
- Paging
- SLAB allocator

Interactive is important!
Ask Questions!

Note: Some slides and/or pictures in the following are adapted from slides ©2013

---

## Recall: two-level page table



- Tree of Page Tables
- Tables fixed size (1024 entries)
  - On context-switch: save single PageTablePtr register
- Valid bits on Page Table Entries
  - Don't need every 2nd-level table
  - Even when exist, 2nd-level tables can reside on disk if not in use

---

## Recall: Multi-level Translation: Segments + Pages

- What about a tree of tables?
  - Lowest level page table⇒memory still allocated with bitmap
  - Higher levels often segmented
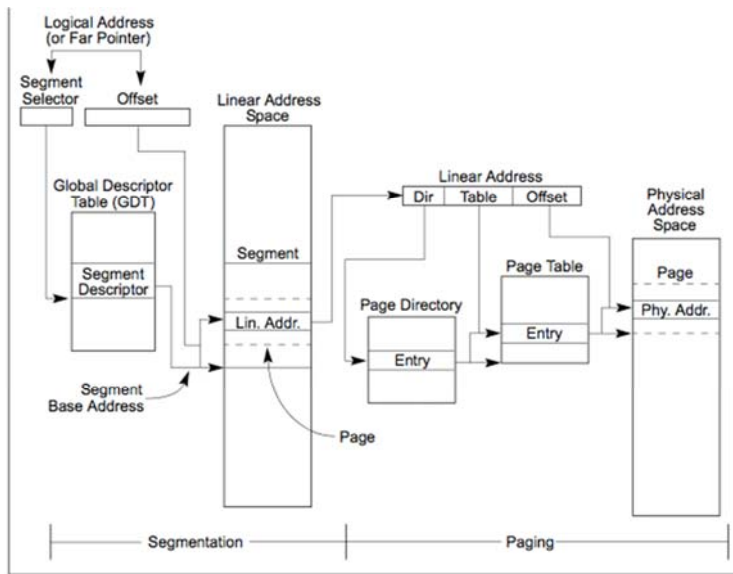- Could have any number of levels. Example (top segment):



- What must be saved/restored on context switch?
  - Contents of top-level segment registers (for this example)
  - Pointer to top-level table (page table)

## Recall: X86 Memory model with segmentation (16/32-bit)
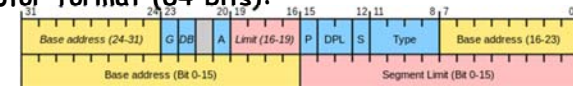
## Recall: X86 Segment Descriptors (32-bit Protected Mode)

- Segments are either implicit in the instruction (say for code segments) or actually part of the instruction
  - There are 6 registers: SS, CS, DS, ES, FS, GS
- **What is in a segment register?**
  - *A pointer to the actual segment description:*

| Segment selector [13 bits] | G/L | RPL |
|---|---|---|

  **G/L selects between GDT and LDT tables (global vs local descriptor tables)**
- Two registers: GDTR and LDTR hold pointers to the global and local descriptor tables in memory
  - Includes length of table (for < $2^{13}$) entries
- Descriptor format (64 bits):



G: Granularity of segment (0: 16bit, 1: 4KiB unit)
DB: Default operand size (0: 16bit, 1: 32bit)
A: Freely available for use by software
P: Segment present
DPL: Descriptor Privilege Level
S: System Segment (0: System, 1: code or data)
Type: Code, Data, Segment
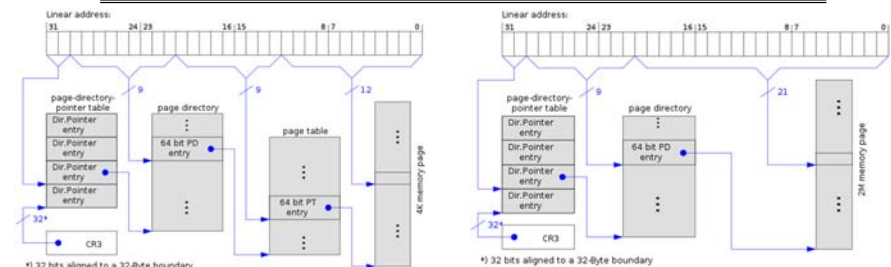
## Recall: How are segments used?

- One set of global segments (GDT) for everyone, different set of local segments (LDT) for every process
- In legacy applications (16-bit mode):
  - Segments provide protection for different components of user programs
  - Separate segments for chunks of code, data, stacks
  - Limited to 64K segments
- Modern use in 32-bit Mode:
  - Segments "flattened", i.e. every segment is 4GB in size
  - One exception: Use of GS (or FS) as a pointer to "Thread Local Storage"
    » A thread can make accesses to TLS like this:
      mov eax, gs(0x0)
- Modern use in 64-bit ("long") mode
  - Most segments (SS, CS, DS, ES) have zero base and no length limits
  - Only FS and GS retain their functionality (for use in TLS)

## Slightly More than 4GB RAM: PAE mode on x86



PAE with 4K pages          PAE with 2MB pages

- Physical Address Extension (PAE)
  - Poor-man's large memory extensions
  - More than 4GB physical memory
  - Every process still can have only 32-bit address space
- 3-Level page table
  - 64-bit PTE format
- How do processes use more than 4GB memory?
  - OS Support for mapping and unmapping physical memory into virtual address space
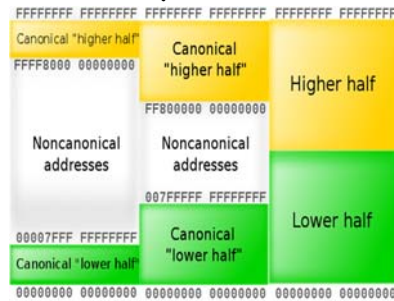  - Application Windowing Extensions (AWE)

## What about 64-bit x86-64 ("long mode")?

- **X86 long mode: 64 bit virtual addresses, 40-52 bits of physical memory**
  - Not all 64-bit virtual addresses translated
  - Virtual Addresses must be "cannonical": top n bits of must be equal
    - » n here might be 48
    - » Non-cannonical addresses will cause a protection fault

- **Using PAE scheme with 64-bit PTE can map 48-bits of virtual memory ($9 \times 4 + 12 = 48$)**
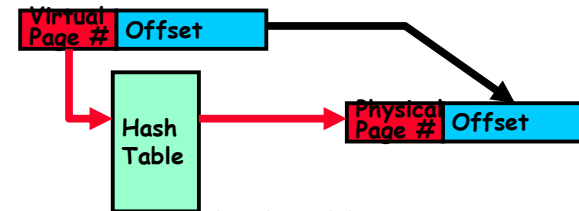- **As mentioned earlier, segments other than FS/GS disabled in long mode**

---

## Inverted Page Table

- **With all previous examples ("Forward Page Tables")**
  - Size of page table is at least as large as amount of virtual memory allocated to processes
  - Physical memory may be much less
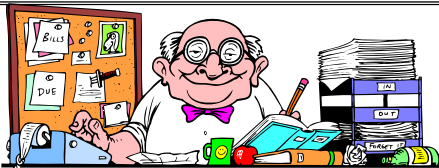    - » Much of process space may be out on disk or not in use



- **Answer: use a hash table**
  - Called an "Inverted Page Table"
  - Size is independent of virtual address space
  - Directly related to amount of physical memory
  - Very attractive option for 64-bit address spaces
- **Cons: Complexity of managing hash changes**
  - Often in hardware!
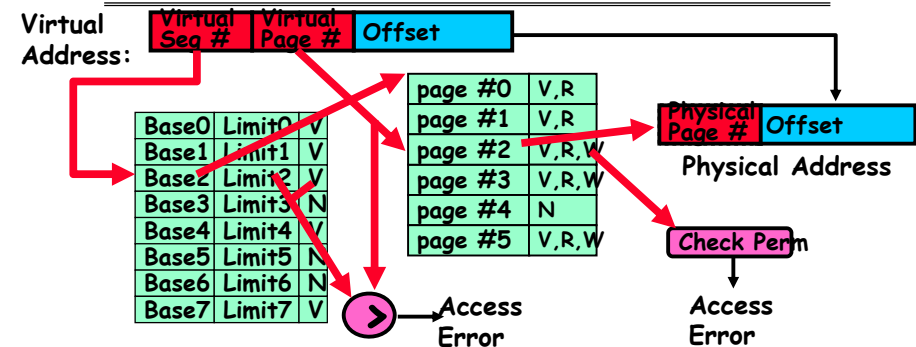
---

## Recall: Caching Concept

- **Cache: a repository for copies that can be accessed more quickly than the original**
  - Make frequent case fast and infrequent case less dominant
- **Caching underlies many of the techniques that are used today to make computers fast**
  - Can cache: memory locations, address translations, pages, file blocks, file names, network routes, etc...
- **Only good if:**
  - Frequent case frequent enough and
  - Infrequent case not too expensive
- **Important measure: Average Access time =**
  
  (Hit Rate x **Hit Time**) + (Miss Rate x **Miss Time**)

---

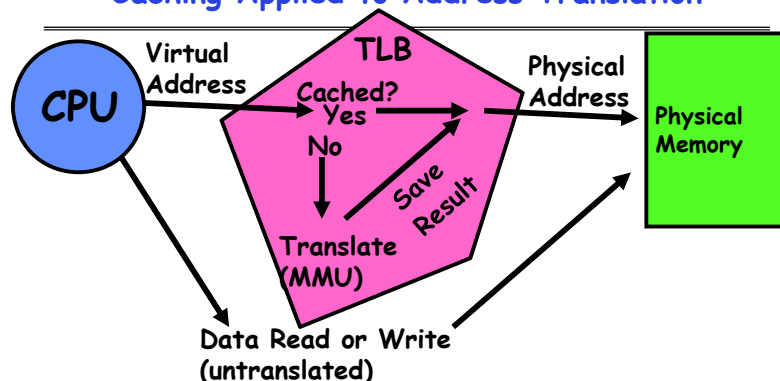## Why does caching matter for Virtual Memory?



- **Cannot afford to translate on every access**
  - At least three DRAM accesses per actual DRAM access
  - Or: perhaps I/O if page table partially on disk!
- **Even worse: What if we are using caching to make memory access faster than DRAM access???**
- **Solution? Cache translations!**
  - Translation Cache: TLB ("Translation Lookaside Buffer")

## Caching Applied to Address Translation



- **Question is one of page locality: does it exist?**
  - Instruction accesses spend a lot of time on the same page (since accesses sequential)
  - Stack accesses have definite locality of reference
  - Data accesses have less page locality, but still some…
- **Can we have a TLB hierarchy?**
  - Sure: multiple levels at different sizes/speeds

---

## What Actually Happens on a TLB Miss?

- **Hardware traversed page tables:**
  - On TLB miss, hardware in MMU looks at current page table to fill TLB (may walk multiple levels)
    » If PTE valid, hardware fills TLB and processor never knows
    » If PTE marked as invalid, causes Page Fault, after which kernel decides what to do afterwards
- **Software traversed Page tables (like MIPS)**
  - On TLB miss, processor receives TLB fault
  - Kernel traverses page table to find PTE
    » If PTE valid, fills TLB and returns from fault
    » If PTE marked as invalid, internally calls Page Fault handler
- **Most chip sets provide hardware traversal**
  - Modern operating systems tend to have more TLB faults since they use translation for many things
  - Examples:
    » shared segments
    » user-level portions of an operating system

---

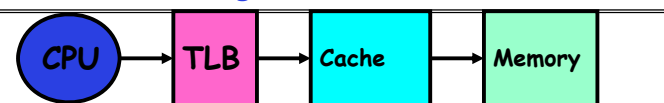## What happens on a Context Switch?

- **Need to do something, since TLBs map virtual addresses to physical addresses**
  - Address Space just changed, so TLB entries no longer valid!
- **Options?**
  - Invalidate TLB: simple but might be expensive
    » What if switching frequently between processes?
  - Include ProcessID in TLB
    » This is an architectural solution: needs hardware
- **What if translation tables change?**
  - For example, to move page from memory to disk or vice versa…
  - Must invalidate TLB entry!
    » Otherwise, might think that page is still in memory!

---

## What TLB organization makes sense?



- **Needs to be really fast**
  - Critical path of memory access
    » In simplest view: before the cache
    » Thus, this adds to access time (reducing cache speed)
  - Seems to argue for Direct Mapped or Low Associativity
- **However, needs to have very few conflicts!**
  - With TLB, the Miss Time extremely high!
  - This argues that cost of Conflict (Miss Time) is much higher than slightly increased cost of access (Hit Time)
- **Thrashing: continuous conflicts between accesses**
  - What if use low order bits of page as index into TLB?
    » First page of code, data, stack may map to same entry
    » Need 3-way associativity at least?
  - What if use high order bits as index?
    » TLB mostly unused for small programs

## TLB organization: include protection

- **How big does TLB actually have to be?**
  - Usually small: 128-512 entries
  - Not very big, can support higher associativity
- **TLB usually organized as fully-associative cache**
  - Lookup is by Virtual Address
  - Returns Physical Address + other info
- **What happens when fully-associative is too slow?**
  - Put a small (4-16 entry) direct-mapped cache in front
  - Called a "TLB Slice"
- **Example for MIPS R3000:**

| Virtual Address | Physical Address | Dirty | Ref | Valid | Access | ASID |
|---|---|---|---|---|---|---|
| 0xFA00 | 0x0003 | Y | N | Y | R/W | 34 |
| 0x0040 | 0x0010 | N | Y | Y | R | 0 |
| 0x0041 | 0x0011 | N | Y | Y | R | 0 |

---

## Example: R3000 pipeline includes TLB "stages"

**MIPS R3000 Pipeline**

| Inst Fetch | Dcd/ Reg | ALU / E.A | Memory | Write Reg |
|---|---|---|---|---|
| TLB   I-Cache | RF | Operation | | WB |
| | | E.A.  TLB | D-Cache | |

**TLB**
  64 entry, on-chip, fully associative, software TLB fault handler

**Virtual Address Space**

| ASID | | | V. Page Number | Offset |
|---|---|---|---|---|
| 6 | | | 20 | 12 |

0xx User segment (caching based on PT/TLB entry)
100 Kernel physical space, cached
101 Kernel physical space, uncached
11x Kernel virtual space

Allows context switching among
64 user processes without TLB flush

---

## Reducing translation time further

- **As described, TLB lookup is in serial with cache lookup:**



- **Machines with TLBs go one step further: they overlap TLB lookup with cache access.**
  - Works because offset available early

---

## Overlapping TLB & Cache Access

- **Here is how this might work with a 4K cache:**



- **What if cache size is increased to 8KB?**
  - Overlap not complete
  - Need to do something else.  See CS152/252
- **Another option: Virtual Caches**
  - Tags in cache are virtual addresses
  - Translation only happens on cache misses

## Modern Pipelines: SandyBridge Pipeline

## Fetch and Branch Prediction
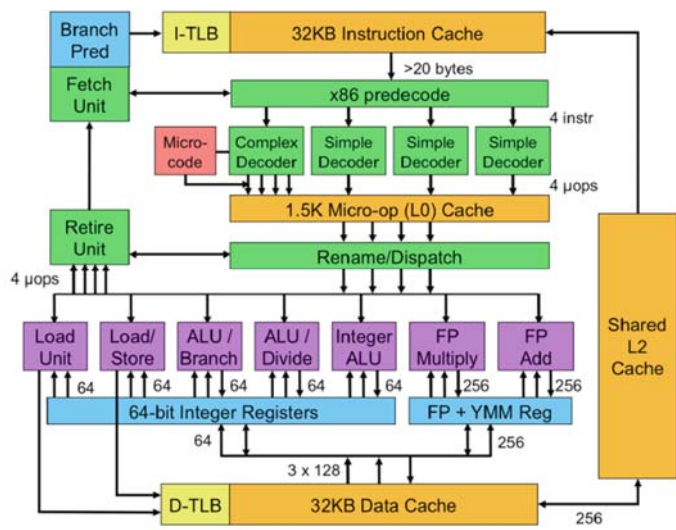


- x86 instructions turned into micro-ops
  - Cached translations are reused
- Branch prediction
  - Not entirely clear, but seems to have some combination of bi-mode, multi-level BTB, stack-based prediction for CALL/RETURN
- Predecoder
  - Finds instruction boundaries
  - Passes at least 6 instructions onto decoding infrastructure

## Out-of-Order execution: Data TLB (DTLB)



- Unified Reservation Unit
  - Full OOO execution
  - Pick 6 ready μops/cycle
  - Can have two loads or stores/cycle
    » 2 address generation units (AGUs) + store data
  - Simultaneous 256-bit Multiply and Add
  - Can have 3 regular integer ops/cycle

## Use of Mapping as a Cache: Demand Paging

- **Modern programs require a lot of physical memory**
  - **Memory per system growing faster than 25%-30%/year**
- **But they don't use all their memory all of the time**
  - **90-10 rule: programs spend 90% of their time in 10% of their code**
  - **Wasteful to require all of user's code to be in memory**
- **Solution: use main memory as cache for disk**

## Illusion of Infinite Memory



**Virtual Memory 4 GB** → TLB → **Page Table** → **Physical Memory 512 MB** → **Disk 500GB**

- **Disk is larger than physical memory ⇒**
  - In-use virtual memory can be bigger than physical memory
  - Combined memory of running processes much larger than physical memory
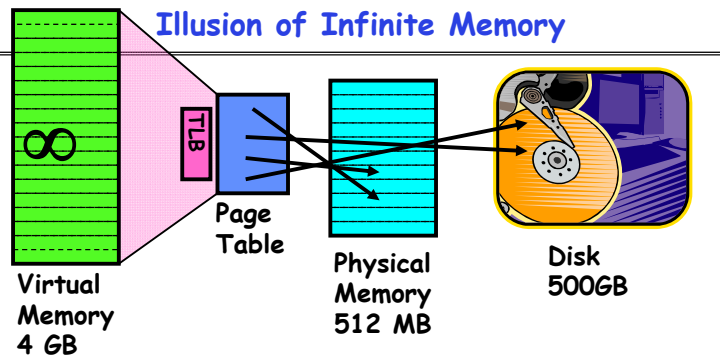    » More programs fit into memory, allowing more concurrency
- Principle: **Transparent Level of Indirection** (page table)
  - Supports flexible placement of physical data
    » Data could be on disk or somewhere across network
  - Variable location of data transparent to user program
    » Performance issue, not correctness issue

---

## Review: What is in a PTE?

- **What is in a Page Table Entry (or PTE)?**
  - Pointer to next-level page table or to actual page
  - Permission bits: valid, read-only, read-write, write-only
- **Example: Intel x86 architecture PTE:**
  - Address same format previous slide (10, 10, 12-bit offset)
  - Intermediate page tables called "Directories"

| Page Frame Number (Physical Page Number) | Free (OS) | 0 | L | D | A | PCD | PWT | U | W | P |
|---|---|---|---|---|---|---|---|---|---|---|
| 31-12 | 11-9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

  - P: Present (same as "valid" bit in other architectures)
  - W: Writeable
  - U: User accessible
  - PWT: Page write transparent: external cache write-through
  - PCD: Page cache disabled (page cannot be cached)
  - A: Accessed: page has been accessed recently
  - D: Dirty (PTE only): page has been modified recently
  - L: L=1⇒4MB page (directory only).
       Bottom 22 bits of virtual address serve as offset

---

## Demand Paging Mechanisms

- **PTE helps us implement demand paging**
  - Valid ⇒ Page in memory, PTE points at physical page
  - Not Valid ⇒ Page not in memory; use info in PTE to find it on disk when necessary
- **Suppose user references page with invalid PTE?**
  - Memory Management Unit (MMU) traps to OS
    » Resulting trap is a "Page Fault"
  - **What does OS do on a Page Fault?:**
    » Choose an old page to replace
    » If old page modified ("D=1"), write contents back to disk
    » Change its PTE and any cached TLB to be invalid
    » Load new page into memory from disk
    » Update page table entry, invalidate TLB for new entry
    » Continue thread from original faulting location
  - TLB for new page will be loaded when thread continued!
  - While pulling pages off disk for one process, OS runs another process from ready queue
    » Suspended process sits on wait queue

---

## Demand Paging Example

- **Since Demand Paging like caching, can compute average access time! ("Effective Access Time")**
  - EAT = Hit Rate x Hit Time + Miss Rate x Miss Time
  - EAT = Hit Time + Miss Rate x Miss Penalty
- **Example:**
  - Memory access time = 200 nanoseconds
  - Average page-fault service time = 8 milliseconds
  - Suppose p = Probability of miss, 1-p = Probably of hit
  - Then, we can compute EAT as follows:

    EAT = 200ns + p x 8 ms
        = 200ns + p x 8,000,000ns
- **If one access out of 1,000 causes a page fault, then EAT = 8.2 μs:**
  - This is a slowdown by a factor of 40!
- **What if want slowdown by less than 10%?**
  - 200ns x 1.1 < EAT ⇒ $p < 2.5 \times 10^{-6}$
  - This is about 1 page fault in 400000!
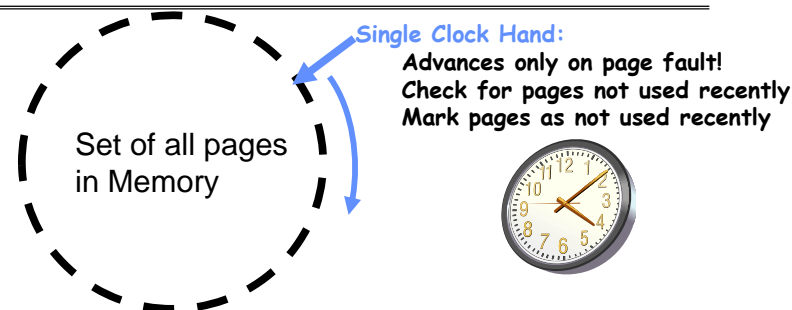
## Implementing LRU

- **Perfect:**
  - Timestamp page on each reference
  - Keep list of pages ordered by time of reference
  - Too expensive to implement in reality for many reasons
- **Clock Algorithm:** Arrange physical pages in circle with single clock hand
  - Approximate LRU (approx to approx to MIN)
  - Replace **an** old page, not **the oldest** page
- **Details:**
  - Hardware "use" bit per physical page:
    - » Hardware sets use bit on each reference
    - » If use bit isn't set, means not referenced in a long time
    - » Nachos hardware sets use bit in the TLB; you have to copy this back to page table when TLB entry gets replaced
  - On page fault:
    - » Advance clock hand (not real time)
    - » Check use bit: 1→used recently; clear and leave alone
                        0→selected candidate for replacement
  - Will always find a page or loop forever?
    - » Even if all use bits set, will eventually loop around⇒FIFO

---

## Clock Algorithm: Not Recently Used



**Single Clock Hand:**
Advances only on page fault!
Check for pages not used recently
Mark pages as not used recently

Set of all pages in Memory

- **What if hand moving slowly?**
  - Good sign or bad sign?
    - » Not many page faults and/or find page quickly
- **What if hand is moving quickly?**
  - Lots of page faults and/or lots of reference bits set
- **One way to view clock algorithm:**
  - Crude partitioning of pages into two groups: young and old
  - Why not partition into more than 2 groups?

---

## Nth Chance version of Clock Algorithm

- **Nth chance algorithm:** Give page N chances
  - OS keeps counter per page: # sweeps
  - On page fault, OS checks use bit:
    - » 1⇒clear use and also clear counter (used in last sweep)
    - » 0⇒increment counter; if count=N, replace page
  - Means that clock hand has to sweep by N times without page being used before page is replaced
- **How do we pick N?**
  - Why pick large N? Better approx to LRU
    - » If N ~ 1K, really good approximation
  - Why pick small N? More efficient
    - » Otherwise might have to look a long way to find free page
- **What about dirty pages?**
  - Takes extra overhead to replace a dirty page, so give dirty pages an extra chance before replacing?
  - Common approach:
    - » Clean pages, use N=1
    - » Dirty pages, use N=2 (and write back to disk when N=1)

---

## Clock Algorithms: Details

- **Which bits of a PTE entry are useful to us?**
  - **Use:** Set when page is referenced; cleared by clock algorithm
  - **Modified:** set when page is modified, cleared when page written to disk
  - **Valid:** ok for program to reference this page
  - **Read-only:** ok for program to read page, but not modify
    - » For example for catching modifications to code pages!
- **Do we really need hardware-supported "modified" bit?**
  - No. Can emulate it (BSD Unix) using read-only bit
    - » Initially, mark all pages as read-only, even data pages
    - » On write, trap to OS. OS sets software "modified" bit, and marks page as read-write.
    - » Whenever page comes back in from disk, mark read-only

## Clock Algorithms Details (continued)
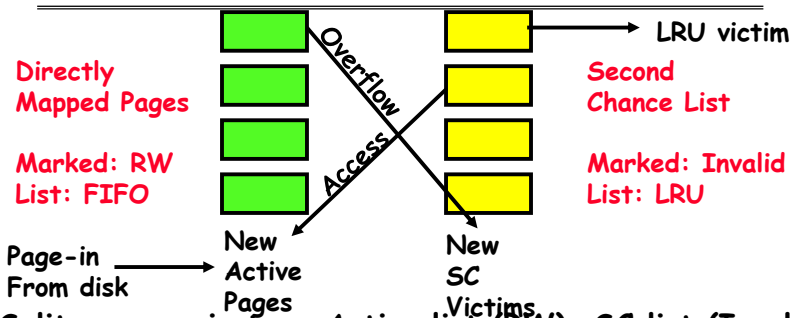
- **Do we really need a hardware-supported "use" bit?**
  - No. Can emulate it similar to above:
    - » Mark all pages as invalid, even if in memory
    - » On read to invalid page, trap to OS
    - » OS sets use bit, and marks page read-only
  - Get modified bit in same way as previous:
    - » On write, trap to OS (either invalid or read-only)
    - » Set use and modified bits, mark page read-write
  - When clock hand passes by, reset use and modified bits and mark page as invalid again
- **Remember, however, that clock is just an approximation of LRU**
  - Can we do a better approximation, given that we have to take page faults on some reads and writes to collect use information?
  - Need to identify an old page, not oldest page!
  - Answer: second chance list

## Second-Chance List Algorithm (VAX/VMS)



- **Split memory in two: Active list (RW), SC list (Invalid)**
- **Access pages in Active list at full speed**
- **Otherwise, Page Fault**
  - Always move overflow page from end of Active list to front of Second-chance list (SC) and mark invalid
  - Desired Page On SC List: move to front of Active list, mark RW
  - Not on SC list: page in to front of Active list, mark RW; page out LRU victim at end of SC list
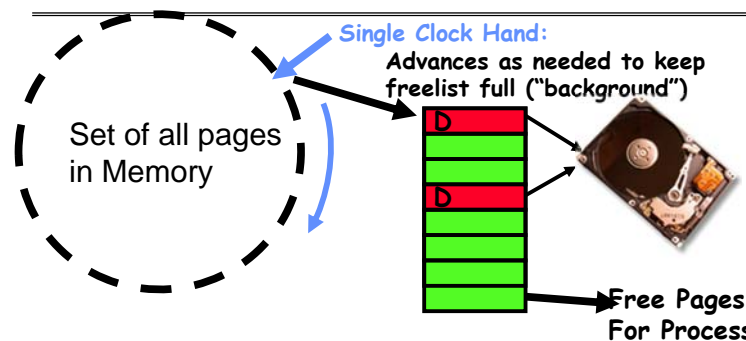
## Second-Chance List Algorithm (con't)

- **How many pages for second chance list?**
  - If 0 ⇒ FIFO
  - If all ⇒ LRU, but page fault on every page reference
- **Pick intermediate value. Result is:**
  - Pro: Few disk accesses (page only goes to disk if unused for a long time)
  - Con: Increased overhead trapping to OS (software / hardware tradeoff)
- **With page translation, we can adapt to any kind of access the program makes**
  - Later, we will show how to use page translation / protection to share memory between threads on widely separated machines
- **Question: why didn't VAX include "use" bit?**
  - Strecker (architect) asked OS people, they said they didn't need it, so didn't implement it
  - He later got blamed, but VAX did OK anyway

## Free List



- **Keep set of free pages ready for use in demand paging**
  - Freelist filled in background by Clock algorithm or other technique ("Pageout demon")
  - Dirty pages start copying back to disk when enter list
- **Like VAX second-chance list**
  - If page needed before reused, just return to active set
- **Advantage: Faster for page fault**
  - Can always use page (or pages) immediately on fault

## Summary: Examples of how to exploit a PTE

- **How do we use the PTE?**
  - Invalid PTE can imply different things:
    » Region of address space is actually invalid or
    » Page/directory is just somewhere else than memory
  - Validity checked first
    » OS can use other (say) 31 bits for location info
- **Usage Example: Demand Paging**
  - Keep only active pages in memory
  - Place others on disk and mark their PTEs invalid
- **Usage Example: Copy on Write**
  - UNIX fork gives *copy* of parent address space to child
    » Address spaces disconnected after child created
  - How to do this cheaply?
    » Make copy of parent's page tables (point at same memory)
    » Mark entries in both sets of page tables as read-only
    » Page fault on write creates two copies
- **Usage Example: Zero Fill On Demand**
  - New data pages must carry no information (say be zeroed)
  - Mark PTEs as invalid; page fault on use gets zeroed page
  - Often, OS creates zeroed pages in background

---

## Reverse Page Mapping (Sometimes called "Coremap")

- **Physical page frames often shared by many different address spaces/page tables**
  - All children forked from given process
  - Shared memory pages between processes
- **Whatever reverse mapping mechanism that is in place must be very fast**
  - Must hunt down all page tables pointing at given page frame when freeing a page
  - Must hunt down all PTEs when seeing if pages "active"
- **Implementation options:**
  - For every page descriptor, keep linked list of page table entries that point to it
    » Management nightmare – expensive
  - Linux 2.6: Object-based reverse mapping
    » Link together memory region descriptors instead (much coarser granularity)

---

## Linux Memory Details?

- **Memory management in Linux considerably more complex that the previous indications**
- **Memory Zones: physical memory categories**
  - ZONE_DMA: < 16MB memory, DMAable on ISA bus
  - ZONE_NORMAL: 16MB $\Rightarrow$ 896MB (mapped at 0xC0000000)
  - ZONE_HIGHMEM: Everything else (> 896MB)
- **Each zone has 1 freelist, 2 LRU lists (Active/Inactive)**
- **Many different types of allocation**
  - SLAB allocators, per-page allocators, mapped/unmapped
- **Many different types of allocated memory:**
  - Anonymous memory (not backed by a file, heap/stack)
  - Mapped memory (backed by a file)
- **Allocation priorities**
  - Is blocking allowed/etc

---

## Recall: Linux Virtual memory map



32-Bit Virtual Address Space          64-Bit Virtual Address Space

## Virtual Map (Details)

- **Kernel memory not generally visible to user**
  - Exception: special VDSO facility that maps kernel code into user space to aid in system calls (and to provide certain actual system calls such as gettimeofday().
- **Every physical page described by a "page" structure**
  - Collected together in lower physical memory
  - Can be accessed in kernel virtual space
  - Linked together in various "LRU" lists
- **For 32-bit virtual memory architectures:**
  - When physical memory < 896MB
    - » All physical memory mapped at 0xC0000000
  - When physical memory >= 896MB
    - » Not all physical memory mapped in kernel space all the time
    - » Can be temporarily mapped with addresses > 0xCC000000
- **For 64-bit virtual memory architectures:**
  - All physical memory mapped above 0xFFFF800000000000

## Internal Interfaces: Allocating Memory

- **One mechanism for requesting pages: everything else on top of this mechanism:**
  - Allocate contiguous group of pages of size $2^{order}$ bytes given the specified mask:
    ```
    struct page * alloc_pages(gfp_t gfp_mask,
                              unsigned int order)
    ```
  - Allocate one page:
    ```
    struct page * alloc_page(gfp_t gfp_mask)
    ```
  - Convert page to logical address (assuming mapped):
    ```
    void * page_address(struct page *page)
    ```
- **Also routines for freeing pages**
- **Zone allocator uses "buddy" allocator that tries to keep memory unfragmented**
- **Allocation routines pick from proper zone, given flags**

## Allocation flags

- **Possible allocation type flags:**
  - **GFP_ATOMIC:** Allocation high-priority and must never sleep. Use in interrupt handlers, top halves, while holding locks, or other times cannot sleep
  - **GFP_NOWAIT:** Like GFP_ATOMIC, except call will not fall back on emergency memory pools. Increases likely hood of failure
  - **GFP_NOIO:** Allocation can block but must not initiate disk I/O.
  - **GFP_NOFS:** Can block, and can initiate disk I/O, but will not initiate filesystem ops.
  - **GFP_KERNEL:** Normal allocation, might block. Use in process context when safe to sleep. This should be default choice
  - **GFP_USER:** Normal allocation for processes
  - **GFP_HIGHMEM:** Allocation from ZONE_HIGHMEM
  - **GFP_DMA** Allocation from ZONE_DMA. Use in combination with a previous flag

## Page Frame Reclaiming Algorithm (PFRA)

- **Several entrypoints:**
  - Low on Memory Reclaiming: The kernel detects a "low on memory" condition
  - Hibernation reclaiming: The kernel must free memory because it is entering in the suspend-to-disk state
  - Periodic reclaiming: A kernel thread is activated periodically to perform memory reclaiming, if necessary
- **Low on Memory reclaiming:**
  - Start flushing out dirty pages to disk
  - Start looping over all memory nodes in the system
    - » try_to_free_pages()
    - » shrink_slab()
    - » pdflush kenel thread writing out dirty pages
- **Periodic reclaiming:**
  - Kswapd kernel threads: checks if number of free page frames in some zone has fallen below pages_high watermark
  - Each zone keeps two LRU lists: Active and Inactive
    - » Each page has a last-chance algorithm with 2 count
    - » Active page lists moved to inactive list when they have been idle for two cycles through the list
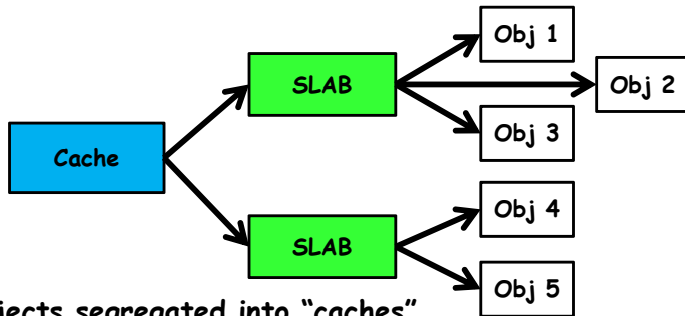    - » Pages reclaimed from Inactive list

## SLAB Allocator

- **Replacement for free-lists that are hand-coded by users**
  - Consolidation of all of this code under kernel control
  - Efficient when objects allocated and freed frequently



- **Objects segregated into "caches"**
  - Each cache stores different type of object
  - Data inside cache divided into "slabs", which are continuous groups of pages (often only 1 page)
  - Key idea: avoid memory fragmentation

---

## SLAB Allocator Details

- **Based on algorithm first introduced for SunOS**
  - Observation: amount of time required to initialize a regular object in the kernel exceeds the amount of time required to allocate and deallocate it
  - Resolves around object caching
    » Allocate once, keep reusing objects
- **Avoids memory fragmentation:**
  - Caching of similarly sized objects, avoid fragmentation
  - Similar to custom freelist per object
- **Reuse of allocation**
  - When new object first allocated, constructor runs
  - On subsequent free/reallocation, constructor does not need to be reexecuted

---

## SLAB Allocator: Cache Construction

- **Creation of new Caches:**
```
struct kmem_cache * kem_cache_create(const char *name,
                                     size_t size,
                                     size_t align,
                                     unsigned long flags,
                                     void (*ctor)(void *));
```
  - **name:** name of cache
  - **size:** size of each element in the cache
  - **align:** alignment for each object (often 0)
  - **flags:** possible flags about allocation
    » SLAB_HWCACHE_ALIGN: Align objects to cache lines
    » SLAB_POISON: Fill slabs to known value (0xa5a5a5a5) in order to catch use of uninitialized memory
    » SLAB_RED_ZONE: Insert empty zones around objects to help detect buffer overruns
    » SLAB_PANIC: Allocation layer panics if allocation fails
    » SLAB_CACHE_DMA: Allocations from DMA-able memory
    » SLAB_NOTRACK: don't track uninitialized memory
  - **ctor:** called whenever new pages are added to the cache

---

## SLAB Allocator: Cache Use

- **Example:**
```
task_struct_cachep =
    kmem_cache_create("task_struct",
                      sizeof(struct task_struct),
                      ARCH_MIN_TASKALIGN,
                      SLAB_PANIC | SLAB_NOTRACK,
                      NULL);
```

- **Use of example:**
```
struct task_struct *tsk;

tsk = kmem_cache_alloc(task_struct_cachep, GFP_KERNEL);
if (!tsk)
    return NULL;

kmem_free(task_struct_cachep,tsk);
```

## SLAB Allocator Details (Con't)

- **Caches can be later destroyed with:**
  `int kmem_cache_destroy(struct kmem_cache *cachep);`
  - Assuming that all objects freed
  - No one ever tries to use cache again
- **All caches kept in global list**
  - Including global caches set up with objects of powers of 2 from $2^5$ to $2^{17}$
  - General kernel allocation (kmalloc/kfree) uses least-fit for requested cache size
- **Reclamation of memory**
  - Caches keep sorted list of empty, partial, and full slabs
    - » Easy to manage – slab metadata contains reference count
    - » Objects within slabs linked together
  - Ask individual caches for full slabs for reclamation

## Recall: Kmalloc/Kfree: The "easy interface" to memory

- **Simplest kernel interface to manage memory: kmalloc()/kfree()**
  - **Allocate chunk of memory in kernel's address space (will be physically contiguous and virtually contiguous):**
    `void * kmalloc(size_t size, gfp_t flags);`
  - **Example usage:**
    ```
    struct dog *p;
    p = kmalloc(sizeof(struct dog), GFP_KERNEL);
    if (!p)
            /* Handle error! */
    ```
  - **Free memory:** `void kfree(const void *ptr);`
  - **Important restrictions!**
    - » Must call with memory previously allocated through `kmalloc()` interface!!!
    - » Must not free memory twice!

## Alternatives for allocation

- **According to Robert Love, "SLAB" has become a name for any allocator with a similar API**
  - Kinda like "Kleenex" has become a generic noun
- **A number of options in the kernel for object allocation:**
  - **SLAB: original allocator based on Bonwick's paper from SunOS**
  - **SLUB: Newer allocator with same interface but better use of metadata (Default since Linux 2.6.23)**
    - » Keeps SLAB metadata in the page data structure (for pages that happen to be in kernel caches)
    - » Debugging options compiled in by default, just need to be enabled
  - **SLOB: low-memory footprint allocator for embedded systems**

## Summary (1/2)

- **Memory is a resource that must be shared**
  - Controlled Overlap: only shared when appropriate
  - Translation: Change Virtual Addresses into Physical Addresses
  - Protection: Prevent unauthorized Sharing of resources
- **Segment Mapping**
  - Segment registers within processor
  - Segment ID associated with each access
    - » Often comes from portion of virtual address
    - » Can come from bits in instruction instead (x86)
  - Each segment contains base and limit information
    - » Offset (rest of address) adjusted by adding base
- **Page Tables**
  - Memory divided into fixed-sized chunks of memory
  - Virtual page number from virtual address mapped through page table to physical page number
  - Offset of virtual address same as physical address
  - Large page tables can be placed into virtual memory
- **Multi-Level Tables**
  - Virtual address mapped to series of tables
  - Permit sparse population of address space
- **Inverted page table**
  - Size of page table related to physical memory size

# Summary (2/2)

- **PTE: Page Table Entries**
  - **Includes physical page number**
  - **Control info (valid bit, writeable, dirty, user, etc)**
- **A cache of translations called a "Translation Lookaside Buffer" (TLB)**
  - **Relatively small number of entries (< 512)**
  - **Fully Associative (Since conflict misses expensive)**
  - **TLB entries contain PTE and optional process ID**
- **On TLB miss, page table must be traversed**
  - **If located PTE is invalid, cause Page Fault**
- **On context switch/change in page table**
  - **TLB entries must be invalidated somehow**
- **TLB is logically in front of cache**
  - **Thus, needs to be overlapped with cache access to be really fast**