

# **ECE 598 – Advanced Operating Systems Lecture 13**

Vince Weaver

`http://web.eece.maine.edu/~vweaver`

`vincent.weaver@maine.edu`

6 March 2018

# Announcements

- Homework #6  
Due after midterm. Be sure to look at memory problem.
- Trouble with vmwos. Got device-tree working.  
Uninitialized data was 0s on some pis, 5s on another.  
Accidentally working.
- Warnings on why its good to comment your code.



# HW#5 Review

- Forgot to say could work in groups.
- Shell to userspace
- Avoid using `sizeof()` where you mean `strlen()`
- Add a time system call
  - Writing to a user-supplied pointer. Dangerous?
  - Linux uses more elaborate `copy_to_user()`
- Questions
  - Nonblocking `getchar`
  - Why run in userspace?



- Changing back to kernel mode
- What is an ABI
- System call of choice
  - manpages, from section “2”
  - many operate on file descriptors
  - chmod, inotify, exit, fork, truncate, futex, stat, wait
  - surprised no one said perf\_event\_open



# Midterm Review

- Closed book/notes/computer but can bring one piece of notebook paper (front only) with notes on it
- Questions will be similar to those from homeworks
- Topics
  - Benefits of an OS / Downsides of an OS
  - Serial communication: why are we using it? What does 9600 7E1 mean? How does hardware and software flow control work?
  - Boot process



- High level, how the GPIO interface works
- Interrupts: how they switch processor mode, why FIQ is different from IRQ mode. How to switch back from userspace.
- System calls
- ABI
- Memory allocation: first vs best fit



# Idle Task Notes

- What does the system do if no jobs are ready to run?
- wfi vs msr (ARM1176 wfi is a nop)
- What happens if forget to setup a stack for the idle task? Not an issue unless you try to add a printk to track down a problem



# Advanced Memory Handling





# Security/Safety

- Want a way to mark memory regions as user only, or read-only, or no-execute
- Some processors provide “segments” for this
- Some ARM processors have a “Memory Protection Unit” (MPU)
- Most modern processors have an MMU (memory management unit) to do full virtual memory



# Using More Memory than Physically Available

- How can you have a program that accesses more RAM than available in physical memory?
- Swapping, as discussed before
- Can manually swap out small parts of a program, this technique is called overlays.
- Split program in parts. Only load the part currently



running at any given time.

- Can we have hardware do this automatically? This is part of the idea of virtual memory.



# Virtual Memory

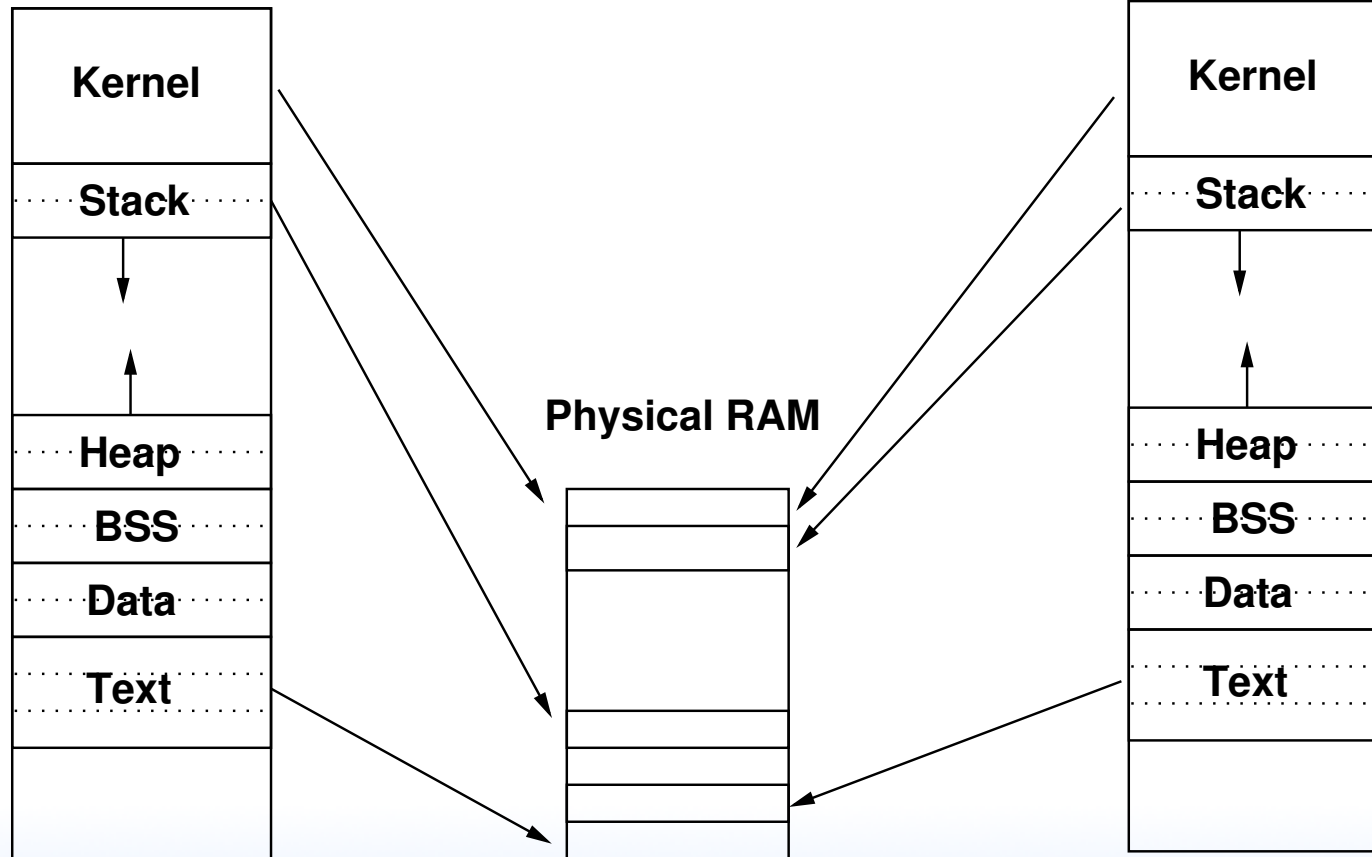
- Original purpose was to give the illusion of more main memory than available, with disk as backing store.
- Give each process own linear view of memory.
- Demand paging (no swapping out whole processes).
- Execution of processes only partly in memory, effectively a cache.
- Memory protection
- Reduces fragmentation



# Diagram

Virtual Process 1

Virtual Process 2



# Memory Management Unit

Can run without MMU. There's even MMU-less Linux.  
How do you keep processes separate? Very carefully...



# Page Table

- Collection of Page Table Entries (PTE)
- Some common components:
  - ID of owner
  - Virtual Page Number
  - valid bit, location of page (memory, disk, etc)
  - protection info (read only, etc)
  - page is dirty, age (how recent updated, for LRU)
  - Much of this info can fit in page info (lower 12 bytes free)



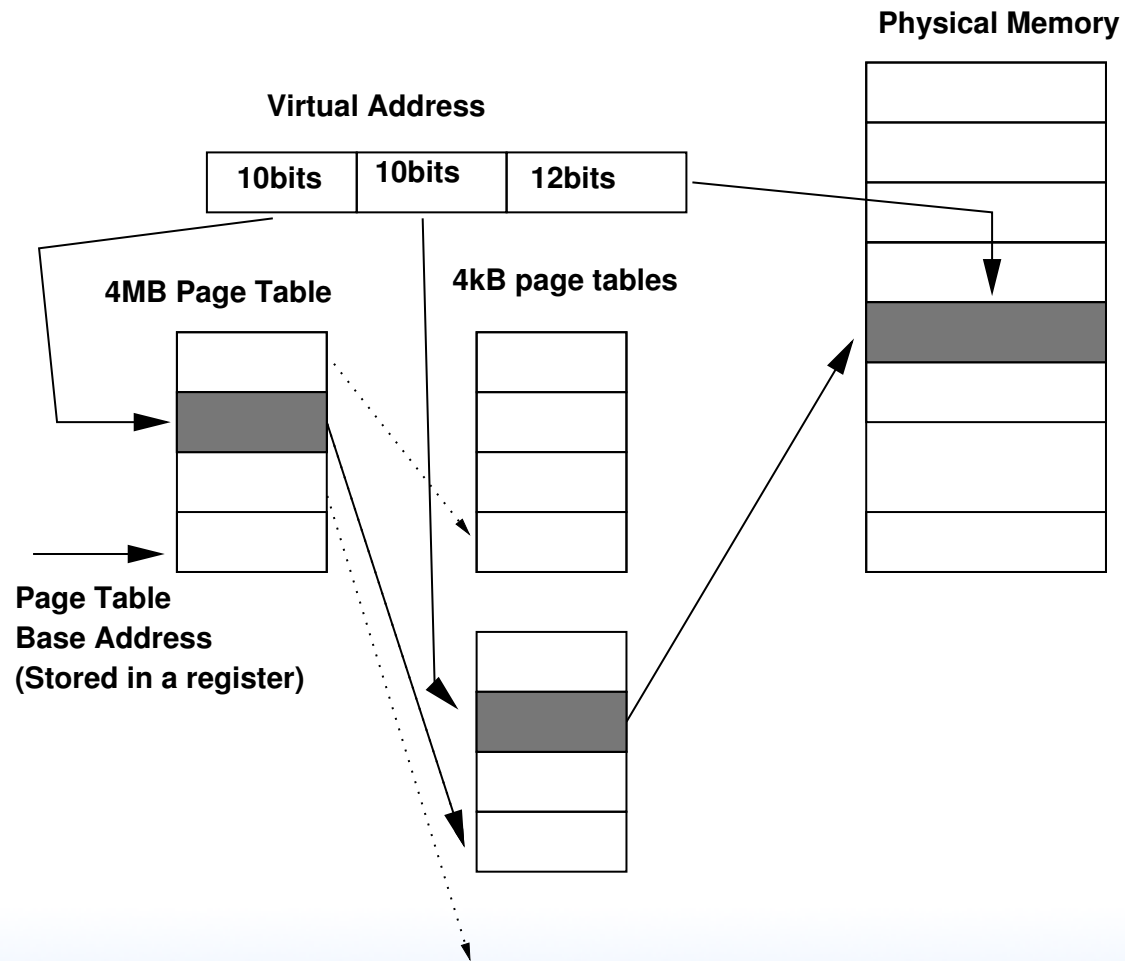
# Hierarchical Page Tables

- With 4GB memory and 4kb pages, you have 1 Million pages per process. If each has 4-byte PTE then 4MB of page tables per-process. Too big.
- It is likely each process does not use all 4GB at once. (sparse) So put page tables in swappable virtual memory themselves!  
4MB page table is 1024 pages which can be mapped in 1 4KB page.





# Hierarchical Page Table Diagram



# Hierarchical Page Table Diagram

- 32-bit x86 chips have hardware 2-level page tables
- ARM 2-level page tables

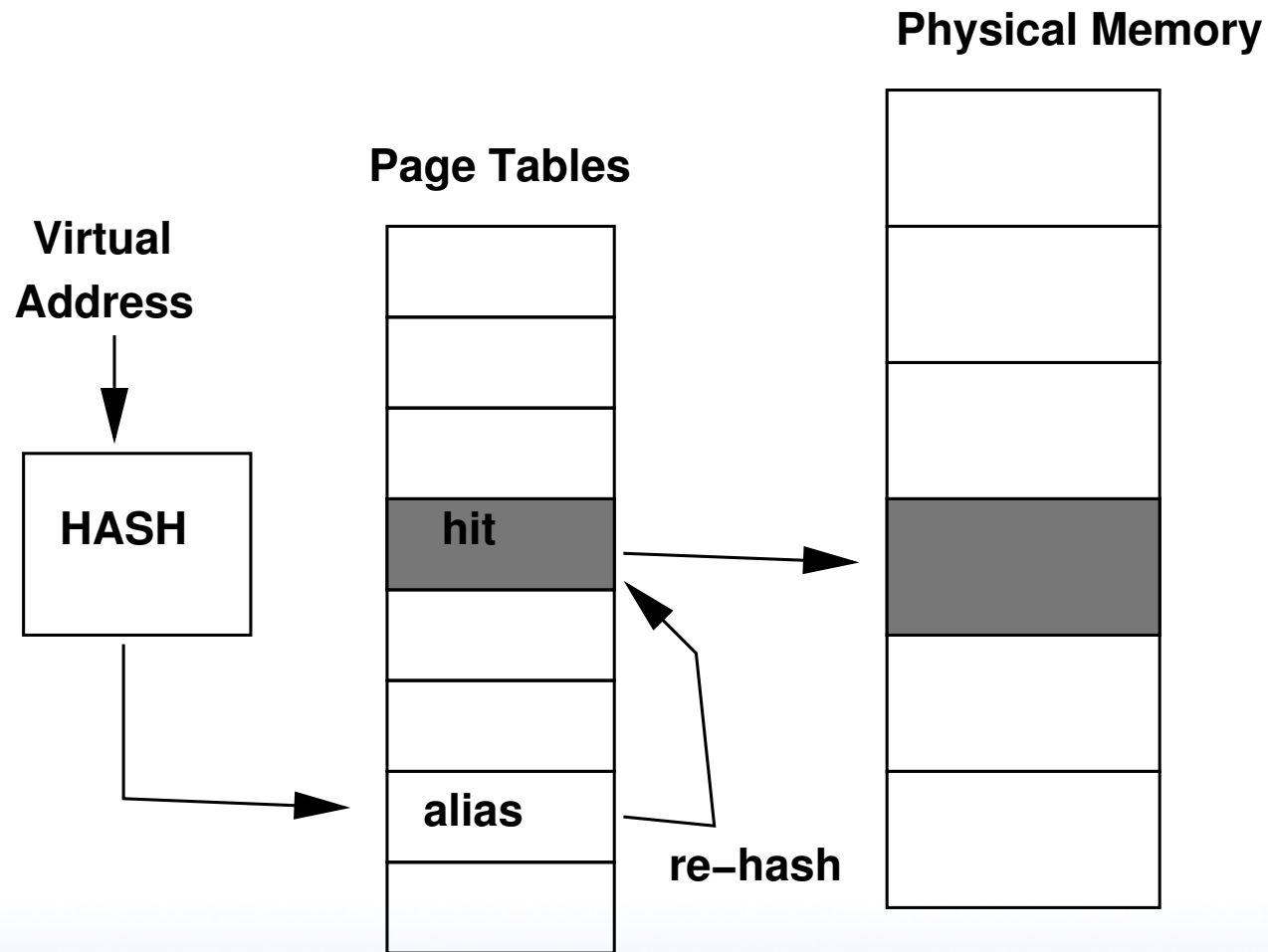


# Inverted Page Table

- How to handle larger 64-bit address spaces?
- Can add more levels of page tables (4? 5?) but that becomes very slow
- Can use hash to find page. Better best case performance, can perform poorly if hash algorithm has lots of aliasing.



# Inverted Page Table Diagram



# Walking the Page Table

- Can be walked in Hardware or Software
- Hardware is more common
- Early RISC machines would do it in Software. Can be slow. Has complications: what if the page-walking code was swapped out?



# TLB

- Translation Lookaside Buffer  
(Lookaside Buffer is an obsolete term meaning cache)
- Caches page tables
- Much faster than doing a page-table walk.
- Historically fully associative, recently multi-level multi-way
- TLB shutdown – when change a setting on a mapping and TLB invalidated on all other processors



# Flushing the TLB

- When do you need to flush?
- May need to do this on context switch if doesn't store ASID or ASIDs run out.
- Sometimes called a “TLB Shutdown”
- Hurts performance as the TLB gradually refills
- Avoiding this is why the top part is mapped to kernel under Linux



# What happens on a memory access

- If in TLB, not a problem, right page fetched from physical memory, TLB updated
- If not in TLB, then the page tables are walked
- If no physical mapping in page table, then page fault happens





# What happens on a page fault

- Walk the page table and see if the page is valid and there
- “minor” – page is already in memory, just need to point a PTE at it. For example, shared memory, shared libraries, etc.
- “major” – page needs to be created or brought in from disk. Demand paging.  
Needs to find room in physical memory. If no free space available, needs to kick something out. Disk-backed



(and not dirty) just discarded. Disk-backed and dirty, written back. Memory can be paged to disk. Eventually can OOM. Memory is then loaded, or zeroed, and PTE updated. Can it be shared? (zero page)

- “invalid” – segfault



# Uses of VM in an operating system

- Process separation, security
- Each process own view of memory
- Kernel mapped into each process address space
- Auto-growing stack
- zero page?
- Memory overcommit



- Demand paging
- Copy-on-write with fork



# What happens on a fork?

- Do you actually copy all of memory?  
Why would that be bad? (slow, also often `exec()` right away)
- Page table marked read-only, then shared
- Only if writes happen, take page fault, then copy made  
Copy-on-write (COW)

