Computer Architecture

Lecture 07 – Address Translation & Virtual Mem

Tian Xia

Institute of Artificial Intelligence and Robotics Xi'an Jiaotong University

http://gr.xjtu.edu.cn/web/pengjuren

Any problem in computer science can be solved with another layer of indirection.

--David Wheeler British Computer Scientist

(1927--2004)



- Fellow of the Royal Society (1981)
- Computer Pioneer Award (1985)
- Fellow, Computer History Museum (2003)

Bare Machine



In a bare machine, the only kind of address is a physical address, corresponding to address lines of actual hardware memory.

Managing Memory in Bare Machines

- Early machines only ran one program at a time, with this program having unrestricted access to all memory and all I/O devices
 - This simple memory management model was also used in turn by the first minicomputer and first microcomputer systems
- Subroutine libraries became popular, were written in location-independent form
 - Different programs use different combination of routines
- To run program on bare machines, use linker or loader program to relocate library modules to actual locations in physical memory

Dynamic Address Translation

- Motivation
 - Each process limited to a non-overlapping contiguous physical memory region (space doesn't start from addr 0...) 1.ITU 2023
 - Everything must fit in the region
- Location-independent programs
 - Programming and storage management ease
 - → need for a *base* register
- Protection
 - Independent programs should not affect each other inadvertently
 - → need for a *bound* register
- Multiprogramming drives requirement for resident supervisor software (e.g. OS) to manage context switches between multiple programs

Program 1 Physical Memory Program 2 OS

Simple Base and Bound Translation

(Scheme used on all Cray vector supercomputers prior to X1, 2002)



Base(Staring address) and bounds(size of region) registers are visible/accessible only when processor is running in the supervisor mode (privileged control registers)

Separate Areas for Program and Data

(Scheme used on all Cray vector supercomputers prior to X1, 2002)



- Shared program segment
- Very fast

What about more base/bound pairs?

Base and Bound Machine



Can fold addition of base register into (register+immediate) address calculation using a carry-save adder (sums three numbers with only a few gate delays more than adding two numbers)

External Fragmentation with Segments



As users come and go, the storage is "fragmented" (external)

Plan ahead to avoid bubbles

Programs are moved around to compact the storage

Reason 1: Adding Disks to Hierarchy



 Simplified Computer Memory Hierarchy Illustration: Ryan J. Leng

Need to devise a mechanism to "Connect" memory and disk in the memory hierarchy.

Reason 2: Simplifying Memory for Apps



Reason 3: Protection Between Processes

- With a bare system, addresses issued with loads/stores are real physical addresses
- This means any program can issue any address, therefore can access any part of memory, even areas which it doesn't own
 - Ex: The OS data structures
- We should send all addresses through a mechanism that the OS controls, before they make it out to DRAM
 - a translation and protection mechanism

2 Parts to Modern Virtual Memory

In a multi-tasking system, virtual memory supports the **illusion** of a **large**, **private**, and **uniform** memory space to each process

- Ingredient A: naming and protection

 each process sees a large, contiguous address space without holes (*for convenience*)
 each process's memory is private, i.e., protected from access by other processes (*for sharing and protection, Readable? Writeable? Executable?*)
 - Ingredient B: demand paging (for hierarchy and efficiency)
 - location-independent programs
 - capacity of secondary storage (swap space on disk)
 - speed of primary storage (DRAM)



Names for Memory Locations



- Machine Language address
 - As specified in machine code
- Virtual Address (VA)
 - ISA specifies translation of machine code address into virtual address of program variable (sometime called effective address)
- Physical Address (PA)
 - Operating System specifies mapping of virtual address into name for a physical memory location

The Common Denominator : Address Translation

- Large, private, and uniform abstraction achieved through address translation
 - user process operates on virtual address (VA)
 - HW translates VA to physical address (PA) on every memory cia@XJ1 reference
- Through address translation
 - control which physical locations (DRAM and/or swap disk) can be referred to by a process
 - allow dynamic allocation and relocation of physical backing store (where in DRAM and/or swap disk)
- Address translation HW and policies controlled by the OS and protected from user

Paged Memory Systems (How)

 Fixed-sized pages (mostly 4KB) in virtual address space are mapped to physical address using Page Table (PT)



- For each program, an independent Page Table is maintained
- Paging makes it possible to store a large contiguous virtual memory space using non-contiguous physical memory pages 17

Private Address Space per User



Virtual Address Space Pages for Job 1

| 0 | |
|---|--|
| 1 | |
| 2 | |
| 3 | |

Virtual Address Space Pages for Job 2

> 0 1 2 3

Virtual Address Space Pages for Job 3





- Each user has a Page Table contains an entry for each user page
- **Persistent OS** residing in memory to control all page tables

Paging Simplifies Memory Allocation

- Fixed-size pages can be kept on OS free list and allocated as needed to any process
- Process memory usage can easily grow and shrink dynamically
- Paging suffers from internal fragmentation (inside Page) where not all bytes on a page are used
 - Much less of an issue than external fragmentation or compaction for common page sizes (4-8KB)
 - But one reason that many oppose move to larger page sizes

Coping with Limited Primary Storage

- Paging reduces fragmentation, but still many problems would not fit into primary memory, have to move data to/from secondary storage (drum, disk) 12023
- Early approach:
 - Manual overlays, programmer explicitly copies code and data in and out of primary memory
 - Tedious coding, error-prone (jumping to non-resident code?)
 - IBM Cell microprocessor using in Playstation-3 had explicitly managed local store!
 - Many new "deep learning" accelerators have similar arch.
- Using virtual memory pages:
 - Put a physical page in primary or secondary storage wherever suitable
 - Maintain its position in a virtual memory page table entry.

Where should Page Tables Reside ?

- Space required by the page table (PT) is proportional to the address space, number of users, ...
 - Space requirement is large (4GB space = 1M PT Entries of 4KB pages for each user = 4MB for 32-bit entry)
- Bad Idea: Keep PT of current user in special registers
 - May not be feasible for large page tables
 - Too expensive to keep in registers
 - Increases the cost of context swap
- Good Idea: Keep PTs in the main memory
 - Use one Page Table Base Register (PTBR) to hold PT's location in the main memory
 - Needs one additional reference to retrieve the page base address and another to access the data word
 - > Double the number of memory references!

Page Tables Live in Memory



Simple linear page tables are too large, so hierarchical page tables are commonly used (see later)

Common for modern OS to place page tables in kernel's **virtual memory** (page tables can be swapped to secondary storage)

Linear Page Table

- Page Table Entry (PTE) contains:
 - A bit to indicate if a page exists
 - PPN (physical page number) for a memory-resident page
 - DPN (disk page number) for a page on the disk
 - Status bits for protection and in usage Pengiu Rena The second second



Hierarchical Page Table Walk: SPARC v8



Hierarchical Page Table is a tree structure, PTBR is the root.

Only create page table when **necessary**, reduces memory footprint
 Termed as **page table walk**, usually performed in **hardware unit**.

Two-Level Page Tables in Physical Memory



Address Translation & Protection Check



Address translation is very expensive! In a two-level page table, each reference requires several memory accesses

Translation-Lookaside Buffers (TLB)

A good VM design needs to be **fast** (~ one cycle) and space **efficient**

Idea: Cache the address translation of frequently used pages



TLB Designs

- Typically 32-128 entries, usually fully associative
 - Each entry maps a large page, hence less spatial locality across pages → more likely that two entries conflict 23
 - Sometimes larger TLBs (256-512 entries) are 4-8 way setassociative
 - Larger systems sometimes have multi-level (L1 and L2) TLBs
- Replacement policy: Random or FIFO
- TLB Reach: Size of largest virtual address space that can be simultaneously mapped by TLB
 - Example: 64 TLB entries, 4KB pages, one page per entry

64 entries * 4 KB = 256 KB (if contiguous)

– TLB Reach =

28

Variable-Size Page TLB

Some Systems support multiple page sizes



MIPS using Pagemask mark different Page size (OS manage)

Handling a TLB Miss

- TLB misses when:
 - The page exists in memory \rightarrow just add the missing entry in TLB.
 - The page doesn't exist in memory \rightarrow transfer control to the OS
- Software (MIPS, Alpha)
 - TLB miss causes an TLB miss exception and the operating system walks the page tables and reloads TLB.
 - Very expensive on out-of-order superscalar processor as requires a flush of pipeline to jump to trap handler.
- Hardware (SPARC v8, x86, PowerPC, RISC-V)
 - A memory management unit (MMU) walks the page tables and reloads the TLB.
 - If a missing (data or PT) page is encountered during the TLB reloading, MMU gives up and signals a Page Fault exception for the original instruction.

Handling a TLB Miss



Handling a TLB Miss



Page Fault Exception

Occurs when an instruction references a memory page that is **not in main memory**.



Page Fault Exception

- Since it takes a long time to transfer a page (msecs), page faults are handled completely in software by OS
- Page Fault Handler (of OS) does the following:
 - Assign an unused page in DRAM
 - If no unused page is left, a page currently in DRAM is swapped out
 - □ Replacement policy: Pseudo-LRU, implemented in software
 - □ If the replaced page is 'dirty', write it back to
 - page table entry that maps that VPN->PPN is marked as invalid/DPN
 - If virtual page exist in disk, Initiate transfer of the requested page from disk to DRAM, assigning to an unused page
 - Another job may be run on the CPU while the first job waits for the requested page to be read from disk
 - Page table entry of the requested page is updated with a (now) valid PPN
 - Return and re-execute the exception-causing instruction
 - Need for precise exceptions

Handling VM-related exceptions



- Handling TLB miss needs a hardware or software mechanism to refill TLB (usually hardware as MMU)
- Handling Page Fault (e.g., page is on disk) needs restartable exception so software handler can resume after retrieving page
 - Precise exceptions are easy to restart
 - Can be imprecise but restartable, but this complicates OS software
- A protection violation may abort process
 - But often handled the same as a page fault


Address Translation – Putting it all together



Page-Based Virtual-Memory Machine (Hardware Page-Table Walk)



MMU uses untranslated physical memory to access page tables

Sequential Access to TLB & Cache (Physical Index/Physical Tag, PIPT)



Adding one more stage for TLB access will increase:

- For I-Cache the miss prediction penalty
- D-Cache load latency (critical path!)

Address Translation in CPU Pipeline



TLB miss? Page Fault?TLB miss? Page Fault?Protection violation?Protection violation?

- Need to cope with additional latency of TLB:
 - slow down the clock?
 - Unacceptable for modern CPUs
 - pipeline the TLB and cache access?
 - Sub-optimal solution (still long latency)
 - virtual address caches
 - parallel TLB/cache access

Virtual-Address Caches



- one-step process in case of a hit (+)
- Homonym and aliasing problems (-)

Homonym in Virtual Address



- Conflicting virtually-tagged entry (both TLB and VIVT cache)
- Software (OS): Cache and TLB needs to be flushed on a context switch
- Hardware: add Address Space Identifier (ASID) into Tags

Aliasing in Virtual-Address Caches





Virtual cache can have two copies of same physical data. Writes to one copy not visible to reads of other!

Two virtual pages share one physical page

General Solution: Prevent aliases coexisting in cache

Software (i.e., OS) solution for direct-mapped cache: VAs of shared pages must **agree in cache index bits**; this ensures all VAs accessing same PA will conflict in direct-mapped cache (early SPARCs)

Concurrent Access to TLB & Cache (Virtual Index/Physical Tag, VIPT)

L=Cache Index, b=Cache Block, k=Page Size,



Tag comparison is made after both accesses are completed.

if Cases: <u>*L* + *b* ≤ *k*</u>

Index L is available without consulting the TLB.

- → Cache and TLB accesses can begin simultaneously!
- → Actually, these are **physical indexed Cache**

Concurrent Access to TLB & Large L1 The problem with L1 Cache > Page size

L=Cache Index, b=Cache Block, k=Page Size,

If Cases: L + b > k



VA₁ and VA₂ both map to same PA

Virtual-Index Physical-Tag Caches: Associative Organization



A solution via Second-Level Cache



Physical-index physical-tag (PIPT) , Inclusive L2 Cache:

- PIPT: Applied for most L2/L3/LLC cache design
- Unified: L2 cache backs up both Instruction and Data L1 caches
- Inclusive: L2 has copies of all cache lines in both L1 D-Cache and I-Cache

Anti-Aliasing (VIPT) Using L2 [MIPS R10000,1996]



Virtually Addressed Cache (Virtual Index/Virtual Tag, VIVT)



- Directly use full virtual address for L1 cache access
- Only check TLB on L1 cache miss
- Use **physical address** for L2 cache access

Anti-Aliasing (VIVT) using L2



Summary: Caching v.s Demand Paging



Summary: TLB, Cache and Page

TLB、 Page Table、 D-Cache and Physical Page (Total 2^4=16 cases)
■ IF PT hit → Physical Page exists in DDR

■ IF TLB hit or D\$ hit → PT hit and physical page exists in DDR

| Situation | TLB | Page Table | D-Cache | Physical Page | Operations |
|----------------------|------|-----------------|---------|-----------------|--|
| D\$ hit | Hit | TLB hit, it hit | Hitn X | D\$ hit, it hit | No need to check Phys. Mem |
| D\$ miss | Hit | TLB hit, it bit | Miss | TLB hit, it hit | Update D\$ |
| TLB miss | Miss | D\$ hit, it Hit | Hit | D\$ hit, it hit | Update TLB, then access TLB again |
| TLB+D\$ both miss | Miss | Hit | Miss | PT hit, it hit | Update TLB and D\$, then access TLB again |
| Page Fault | Miss | Miss | Miss | Miss | Page Fault process |

Cost

Summary : TLB, Cache and Page

| TLB | PAGE TABLE | CACHE | POSSIBLE? | NOTES |
|------|---------------|-------|-----------|---|
| Hit | Hit | Miss | Yes | Page Table is never checked if TLB hits |
| Miss | Hit | Hit | Yes | TLB misses, entry found in Page Table. After retry data is found in Cache |
| Miss | Hit | Miss | Yes | TLB misses ontry found in Page Table. After retry data misses in Cache |
| Miss | Miss | Miss | Yesen | TLB misses, followed by a Page Fault. After retry data misses in Cache |
| Hit | Miss | Miss | GIV No | Cannot have translation in TLB if Page Table not in memory |
| Hit | Miss | Hit | No | Cannot have translation in TLB if Page Table not in memory |
| Miss | Miss | Hit | No | Data cannot be allowed in Cache if Page Table is not memory |

Summary : Modern Virtual Memory Systems

Illusion of a large, private, uniform store

OS

user,

mapping

ΓΙΒ

PA

54

VA

Secondary

Storage

Protection & Privacy Several users, each with their private address space and one or more shared address spaces Kia@XJTL **Demand Paging** Provides the ability to run programs larger Primary than the primary memory Memory Hides differences in physical memory layouts Cost The price is address translation on

each memory reference

Use TLB and Virtual Indexed cache

Why a Privileged Architecture?

- Profiles (Simple Embedded w/wo Protection, Unix-like OS, Cloud OS) Privileges and Modes Ren&Tian Xia@XJTU 2023 Privileged Features CSRs Instructions • Memory Addressing • Translation Protection Trap Handling. Exceptions • Interrupts • Counters Time ullet
 - Performance

RISC-V Privilege Modes

- Machine mode (M-mode, highest privileges)
 - -A.K.A monitor mode, microcode mode, ...
- Hypervisor-Extended Supervisor Mode (HS-Mode)
- Supervisor Mode (S-mode)
- User Mode (U-mode, lowest privileges)
- Supported combinations of modes:
 - M (simple embedded systems)
 - M, U (embedded systems with security)
 - M, S, U (systems running Unix-like OS)
 - M, S, HS, U (systems running hypervisors, Cloud OS Capable)

RISC-V System State (Privileged)

- Processor registers
 - Compute registers
 - General-purpose (x0-x31)
 - (ia@XJTU 2023 Optional floating-point (f0-f31)
 - Optional vector (v0-v31)
 - Optional custom
 - Control and status registers (CSRs)
 - Accessibility controlled by privilege mode or higher
- System main memory
- System I/O devices
- All system memory and device control registers mapped into flat machine physical address space

Extended Page Tables (EPT) in HS Mode



Extended Page Tables (EPT) in HS Mode



Physical Memory Protection (PMP)



Machine Physical Address Space

An optional **physical memory protection (PMP)** unit provides per-hart machinemode control registers to allow physical memory access privileges (read, write, execute) to be specified for each physical memory region

M-Mode controls PMPs

- M-mode has access to entire machine after reset
- Configures PMPs and ioPMPs to contain each active context inside a physical partition
- Can even restrict M-mode access to regions until next reset
- M-mode can dynamically swap PMP settings to run different security contexts on a hart
- RISC-V hardware thread (hart)

Multiple Concurrent Security Contexts



Machine Physical Address Space

PMP checks are applied to all accesses when the hart is running in S or U modes

Memory Protection for RISC-V Modes



Privilege Level Change

Combine privilege level change with interrupt/exception transfer

- switch to next higher privilege level on interrupt/exception
- privilege level restored on return from interrupt/exception
- Interrupt/exception control transfer is only gateway to privileged mode

 lower-level code can never escape into privileged mode

 lower-level code don't even need to know there is a privileged mode



Exceptions

- An exception is a transfer of control to the OS kernel in response to some event (i.e., change in processor state)
 - Kernel is the memory-resident part of the OS
 - Examples of events: Divide by 0, arithmetic overflow, page fault, I/O request completes, typing Ctrl-C



Exception Tables

Exception numbers



Asynchronous Exceptions (Interrupts)

- Caused by events external to the processor
 - Indicated by setting the processor's *interrupt pip*
 - Handler returns to "next" instruction Xia[®]
- Examples:
 - Timer interrupt
 - Every few ms, an external timer chip triggers an interrupt
 - Used by the kernel to take back control from user programs
 - I/O interrupt from external device
 - Hitting Ctrl-C at the keyboard
 - Arrival of a packet from a network
 - Arrival of data from a disk

Synchronous Exceptions

Caused by events that occur as a result of executing an instruction:

- Traps

- Intentional, set program up to "trip the trap" and do something
- Examples: software interrupts, system calls, gdb breakpoints
- Returns control to "next" instruction
- Faults
 - Unintentional but possibly recoverable
 - Examples: page faults (recoverable), protection faults (unrecoverable), floating point exceptions
 - Either re-executes faulting ("current") instruction or aborts
- Aborts
 - Unintentional and unrecoverable
 - Examples: illegal instruction, parity error, machine check
 - Aborts current program

System Calls

- Each x86-64 system call has a unique ID number
- **Examples:**

| impics. | | ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ |
|---------|--------------|---|
| Number | Name | Description |
| 0 | read | Read file |
| 1 | write | Write file |
| 2 | open rian r. | Open file |
| 3 | close 8 | Close file |
| 4 | stat | Get info about file |
| 57 00 | fork | Create process |
| 59 | execve | Execute a program |
| 60 | _exit | Terminate process |
| 62 | kill | Send signal to process |

RISC-V Secure Embedded Systems (M, U modes)

- M-mode runs secure boot and runtime monitor
- Embedded code runs in U-mode
- Physical memory protection (PMP) on U-mode accesses
- Interrupt handling can be delegated to U-mode code
 - User-level interrupt support (N-extension)
- Provides arbitrary number of isolated security contexts



RISC-V Virtual Memory Architectures (M, S, U modes)

- Designed to support current Unix-style operating systems
- Sv32 (RV32)
- ian Xia@XJTU 2023 Demand-paged 32-bit virtual-address spaces
 - 2-level page table
 - 4 KiB pages, 4 MiB megapages
- Sv39 (RV64)
 - Demand-paged 39-bit virtual-address spaces
 - 3-level page table
 - 4 KiB pages, 2 MiB megapages, 1 GiB gigapages
- Sv48, Sv57, Sv64 (RV64)
 - Sv39 + 1/2/3 more page-table levels

S-Mode runs on top of M-mode

- M-mode runs secure boot and monitor
- S-mode runs OS
- U-mode runs application on top of OS or M-mode



- PMP checks are also applied to page-table accesses for virtual-address translation, for which the effective privilege mode is S. Optionally, PMP checks may additionally apply to M-mode accesses.
- PMP can grant permissions to S and U modes, which by default have none, and can revoke permissions from M-mode, which by default has full permissions.
Virtual Memory Use Today - 1

- Servers/desktops/laptops/smartphones have full demand-paged virtual memory
 - Portability between machines with different memory sizes
 - Protection between multiple users or multiple tasks
 - Share small physical memory among active tasks
 - Simplifies implementation of some OS features
- Vector supercomputers have translation and protection but rarely complete demand-paging
- (Older Crays: base&bound, Japanese & Cray X1/X2: pages)
 - Don't waste expensive CPU time thrashing to disk (make jobs fit in memory)
 - Mostly run in batch mode (run set of jobs that fits in memory)
 - Difficult to implement restartable vector instructions

Virtual Memory Use Today - 2

- Most embedded processors and DSPs provide physical addressing only
 - Can't afford area/speed/power budget for virtual memory support
 - Often there is no secondary storage to swap to!
 - Programs custom written for particular memory configuration in product
 - Difficult to implement restartable instructions for exposed architectures

Next Lecture: Branch Prediction

Acknowledgements

Some slides contain material developed and copyright by:

- Arvind (MIT)
- Krste Asanovic (MIT/UCB)
- Joel Emer (Intel/MIT)
- James Hoe (CMU)
- David Patterson (UCB)
- 12@XJTU 2023 David Wentzlaff (Princeton University)
- MIT material derived from course 6.823
- UCB material derived from course CS252 and CS 61C