CS162 Operating Systems and Systems Programming Lecture 15

Virtual Memory (2)

Professor Natacha Crooks & Matei Zaharia https://cs162.org/

Slides based on prior slide decks from David Culler, Ion Stoica, John Kubiatowicz, Alison Norman and Lorenzo Alvisi

Recall: Memory Management Wishlist

Memory Protection

Memory Sharing

Flexible Memory Placement

Support for Sparse Addresses

Runtime Lookup Efficiency

Compact Translation Table

Crooks & Zaharia CS162 © UCB Fall 2025

Recall: Increasingly powerful mechanisms



Divide logical address space of process into fixed sized chunks called pages

View physical memory as an array of fixed-sized slots called page frames

Each page frame can contain a single virtual-memory page

Pages should be small to minimise internal fragmentation (1K-16k)

How to Implement Simple Paging?

Interpret virtual address as two components



How to Implement Simple Paging?

Interpret virtual address as two components



A (Simplified) Page Table

A page table stores

virtual-to-physical address translations

One page table per process. Lives in memory.

Address stored in the in the Page Table Base Register

PTBR value saved/restored in PCB on context switch

How to access a byte?

Extract page number (first p bits)

Map virtual page number into a frame number (also called physical page number) using a page table

Extract offset (last o bits)

Convert to physical memory location: access byte at offset in frame

A (Simplified) Page Table



Assume we have a 64 bytes (2^6) of physical memory

Assume we want pages of 4 bytes (2²)

How long should our addresses be?

6 bits

How many offset bits should we assign? 2 bits

How many virtual pages can we have? 6 bit addresses: 2 bit for offsets, 4 bits for VPN. 2^4 = 16 pages

Crooks & Zaharia CS162 © UCB Fall 2025







Step 1: Extract Virtual Page Number



Step 2: Identify Physical Page Number



Step 3: Extract Frame Offset



Step 3: Extract Frame Offset



Step 3: Extract Frame Offset



Step 4: Convert to Physical Address



What is a page table entry? (32 bits)

Page Frame Number (Physical Page Number)	Free (OS)	0	PS	D	Α	PCD	PWT	U	w	Ρ
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: page has been modified recently
 - PS: Page Size

Size of page table entry: PFN (20 bits) + 12 bits for access control/caching

4 bytes

The Great Power of the PTE

Demand Paging

Keep only active pages in memory Place others on disk and mark their PTEs invalid

Copy-on-Write

UNIX fork gives *copy* of parent address space to child. Use combination of page sharing + marking pages as read-only

Zero Fill On Demand

New data pages must carry no information

Mark PTEs as invalid; page fault on use gets zeroed page

Data Breakpoints

For debugger, mark instruction page as read-only. Will trigger page-fault when try to execute Processes share a page by each mapping a page of their own virtual address space to the same frame

Use protection bits for fine-sharing



Kernel region of every process has the same page table entries

Different processes running same binary! Do not need to duplicate code segments

Shared-memory segments between different processes

Memory Layout for Linux 32-bit



From the paper:

Meltdown is a novel attack that allows overcoming memory isolation completely by providing a simple way for any user process to read the entire kernel memory of the machine it executes on, including all physical memory mapped in the kernel region. Meltdown does not exploit any software vulnerability, i.e., it works on all major operating systems.



- 1. raise_exception();
- 2. // the line below is never reached
- 3. access(probe_array[data * 4096]);]

Are we done?

How big can a page table get on x86 (32 bits)?

4KB page => 2^12 2^32/2^12 => 2^20 pages 2^20 * 4 bytes = 4 MB (approx.) That's (not) a lot per process!!

How big can a page table get on x86 (64 bits)?

4KB page => 2^12 2^64/2^12 => 2^52 pages 2^20 * 8 bytes = 36 petabytes (approx.) That's a lot per process!!

Crooks & Zaharia CS162 © UCB Fall 2025

Space overhead

With a 64-bit address space, size of page table can be huge

Time overhead

Accessing data now requires two memory accesses must also access page table, to find mapped frame

Internal Fragmentation

4KB pages

The Secret to the Whole of CS

Batching

Caching

Indirection

Specialised Hardware



Address space is sparse, i.e. has holes that are not mapped to physical memory

Most this space is taken up by page tables mapped to nothing

Process has access to full 2^64 bytes (virtually)

Physically, that would be 17,179,869,184 gigabytes

Paging the page table: 2-level paging

Tree of Page Tables



V2: What is a page table entry? (32 bits)

Inner Page Table VAddress or PFN	Free (OS)	0	PS	D	Α	PCD	PWT	U	w	Ρ
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: page has been modified recently
 - PS: Page Size

Paging the page table: 2-level paging

Tree of Page Tables

Outer Page # Inner Page #	Offset
---------------------------	--------

Number of top-Ensure that fitsDefines size of a pagelevel pageson a single page

Paging the page table: 2-level paging

Tree of Page Tables

Outer Page #	Inner Page #	Offset
10 bits	10 bits	4 KB 12 bits
	that s in a 0	

Example: x86 classic 32-bit address translation



Figure 4-2. Linear-Address Translation to a 4-KByte Page using 32-Bit Paging

Top-level page-table: Page Directory

Inner page-table: Page Directory Entries

Crooks & Zaharia CS162 © UCB Fall 2025

Example Address Space View



Sharing with multilevel page tables

Entire regions of the address space can be efficiently shared



Marking entire regions as invalid!

If region of address space unused, can mark entire inner region as invalid



Marking entire regions as invalid!

If region of address space unused, can mark entire inner region as invalid



Has this helped?

```
Assuming 10/10/12 split:
Size of Page Table
Outer: (2^10 * 4 bytes) +
Inner: 2^10 * (2^10 * 4 bytes)
```

Overhead of indirection! BUT Marking inner pages as invalid helps when address spaces are sparse

Downside: now have to do two memory accesses for translation

Paged Segmentation

Use segments for top level. Paging within each segment.

Used in x86 (32 bit). Code Segment, Data Segment, etc.



X86 64 bits has a four-level page table!



Inverted Page Table

A single page table that has an entry for each physical page of the system

Each entry contains process ID + which virtual page maps to physical page

Physical memory much smaller than virtual memory

Size proportional to size of physical memory



Inverted Page Table

Don't we have it backwards?

Add a hash table. Virtual memory can only map to specific physical frames



Address Translation Comparison

	Advantages	Disadvantages
Simple Segmentation	Fast context switching (segment map maintained by CPU)	External fragmentation
Paging (Single-Level)	No external fragmentation Fast and easy allocation	Large table size (~ virtual memory) Internal fragmentation
Paged Segmentation	Table size ~ # of pages in virtual	Multiple memory references
Multi-Level Paging	memory Fast and easy allocation	per page access
Inverted Page Table	Table size ~ # of pages in physical memory	Hash function more complex No cache locality of page table

How is the Translation Accomplished?



MMU must translate virtual address to physical address on every instruction fetch, load or store

What does the MMU need to do to translate an address? Read, check, and update PTE (set accessed bit/dirty bit on write)

How can we speedup translation?

MMU must make at least 2 memory reads to walk page table. Slow!

Use specialized hardware to cache virtual-physical memory translations!

Introducing the Translation Lookaside Buffer (TLB)

Cache: a repository for copies that can be accessed more quickly than the original

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)

Recall: In Machine Structures (eg. 61C) ...

Caching is the key to memory system performance



Recall: In Machine Structures (eg. 61C) ...

Average Memory Access Time (AMAT) = (Hit Rate x HitTime) + (Miss Rate x MissTime) Where HitRate + MissRate = 1

HitRate = $90\% => AMAT = (0.9 \times 1) + (0.1 \times 101) = 11 \text{ ns}$ HitRate = $99\% => AMAT = (0.99 \times 1) + (0.01 \times 101) = 2.01 \text{ ns}$

 $MissTime_{L1} includes$ $HitTime_{L1} + MissPenalty_{L1} \equiv HitTime_{L1} + AMAT_{L2}$

Temporal Locality (Locality in Time): Keep recently accessed data items closer to processor

Spatial Locality (Locality in Space): Move contiguous blocks to the upper levels Take advantage of the principle of locality to:

- 1) Present the illusion of having as much memory as in the cheapest technology
 - 2) Provide average speed similar to that offered by the fastest technology

Recall: fast but small/expensive. Slow but large!

Recall: Memory Hierarchy



How do we make Address Translation Fast?

Cache results of recent translations ! Cache Page Table Entries using Virtual Page # as the key



Crooks & Zaharia CS162 © UCB Fall 2025

Record recent Virtual Page # to Physical Frame # translation

If present, have the physical address without reading any of the page tables !!!

Caches the end-to-end result

Caching Applied to Address Translation



Does page locality exist?

Instruction accesses spend a lot of time on the same page (since accesses sequential) Stack accesses have definite locality of reference