

Program Optimization

- Code motion
- Reduction in strength → Shift, add instead of multiply or divide

- Procedure calls → optimization blocker^{#1}
Compiler won't do code motion to prevent side effects.

- Memory matters

Instruction control unit
determines the destination
of jump
predicts whether will be taken

Starts fetching instruction at predicted destination

Execution unit checks whether prediction was ok.

If not → signals instruction control
instruction control invalidates any operations generated from mispredicted instructions.
begins fetching & decoding instructions at correct target.

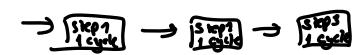
Unrolling & Accumulating :

- can unroll to any degree L
- can accumulate K results in parallel
- L must be multiple of K.

SIMD operations : (Single Inst Multiple Data)
single precision → 8 additions
double " → 4 "

Latency vs Throughput

latency cycles/issue
Integer multiply 3 1



(String example)

↑ 3 independent multiplications will be pipelined.
Sequentially dependent won't be that efficient.

$$T = CPE \times n + \text{Overhead}$$

b is a slope of line

Cycles/issue: 1
Latency: 3

$$P_1 = a \times b$$

$$P_2 = a \times c$$

$$P_3 = P_1 + P_2$$

Pipeline Functional Units							
	1	2	3	4	5	6	7
Stage 1	a+b	a+c					p1=ab
Stage 2		a+b	a+c				p2=ac
Stage 3			a+b	a+c			p3=ab+ac

Branches → CPU must work well ahead of Execution unit to generate enough operations to keep EU busy.
When encounters conditional branch, cannot reliably determine where to continue fetching.

The Memory hierarchy

SRAM : retains value indefinitely, insensitive to electrical noise, faster & expensive
RAM : DRAM : value must be refreshed, sensitive " " "(6-transistor circuit)
basic storage unit is a cell
multiple RAM chips form a memory
reading the content slows & cheaper (capacitor & transistor)
of a memory cell destroys its content.

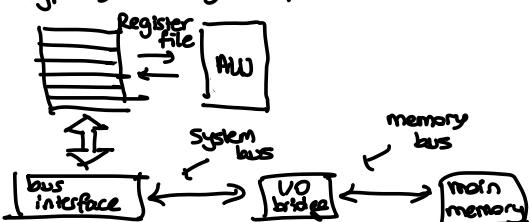
DRAM & SRAM are volatile memories

Nonvolatile memories → ROM → PROM
Firmware → Program stored in ROM → EEPROM
(boot time code)
BIOS, graphics card, disk controllers) → Flash Memory

Traditional Bus Structure

Bus is a collection of parallel wires that carry address, data & control signals.

Typically shared by multiple devices.



Memory Read Transaction :

- CPU places address A on the memory bus
- Main memory reads from the memory bus, retrieves word x and places it on the bus.
- CPU reads word x from the bus and copies it into register file.

Memory Write Transaction :

- CPU places address A on bus. Main memory needs it and waits for the corresponding data word to arrive.
- CPU places data word y on the bus
- Main memory reads data word y from the bus and stores it.

it at address A.

Locality

Principle of Locality: Programs tend to use data & instructions with addresses near or equal to those they have used recently.

Temporal Locality: Recently referenced items are likely to be referenced again in the future.



Spatial Locality: Items with nearby addresses tend to be referenced close together in time.



L0: Registers

↑ smaller, faster,
costlier per byte

L1: L1 Cache (SRAM)

↓ larger, slower,
cheaper per byte.

L2: L2 Cache (SRAM)

L3: Main Memory (DRAM)

L4: Local Secondary Storage (local disks)

L5: Remote Secondary Storage (tapes, distributed file systems, web servers)

Caches:

fundamental idea of
memory hierarchy

→ For each k , the faster, smaller device at level k
serves as a cache for the larger, slower device at level $k+1$

The memory hierarchy

- Creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.

Because of locality, programs tend to access the data at level k more often than they access the data at level $k+1$. Thus the storage at level $k+1$ can be slower, and thus larger and cheaper per bit.

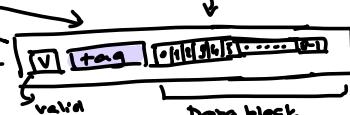
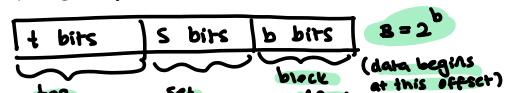
General Cache Organization (S, E, B)



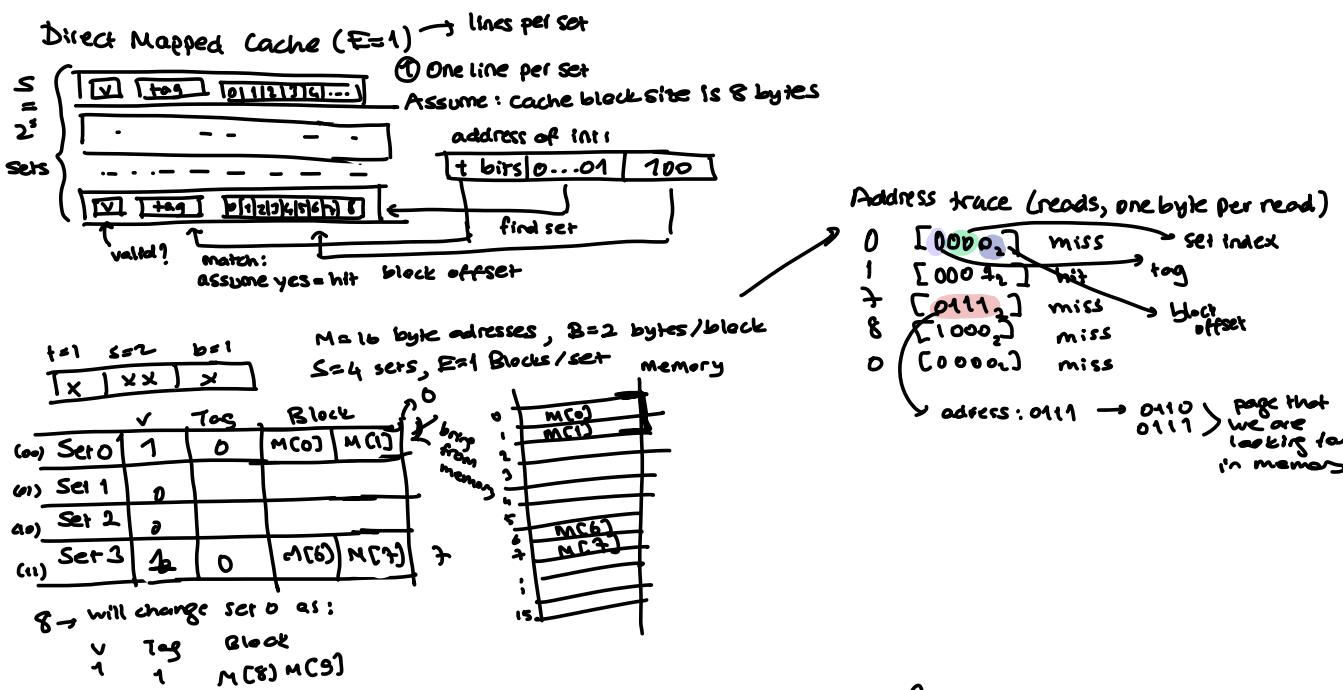
Read:

- Locate set
- Check if any line in set has matching tag
- If + line valid: hit
- Locate data starting at offset -

Address of word:



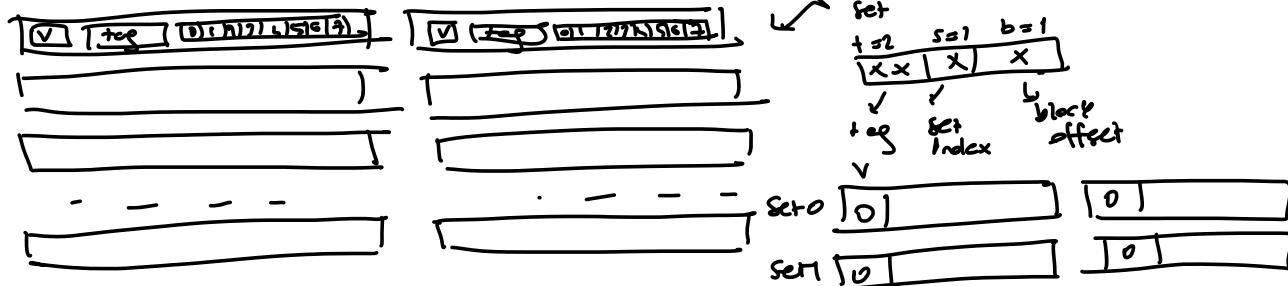
$B = 2^b$ bytes per cache block (the data)



E-way Set Associative Cache

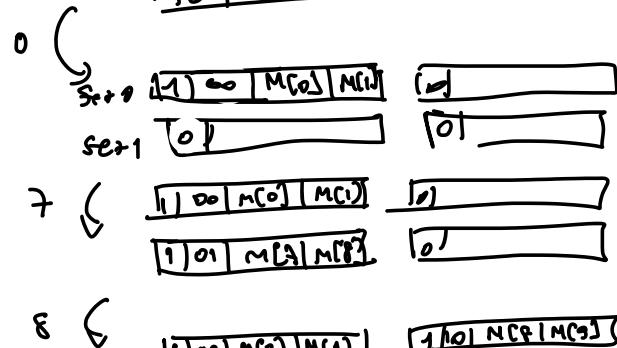
$E=2$: Two lines per set

Assume: cache block size: 8



Address Trace

- 0 [0000₂] → miss
- 1 [0001₂] → hit
- 7 [0111₂] → miss
- 8 [1000₂] → miss
- 0 [0000₂] → hit

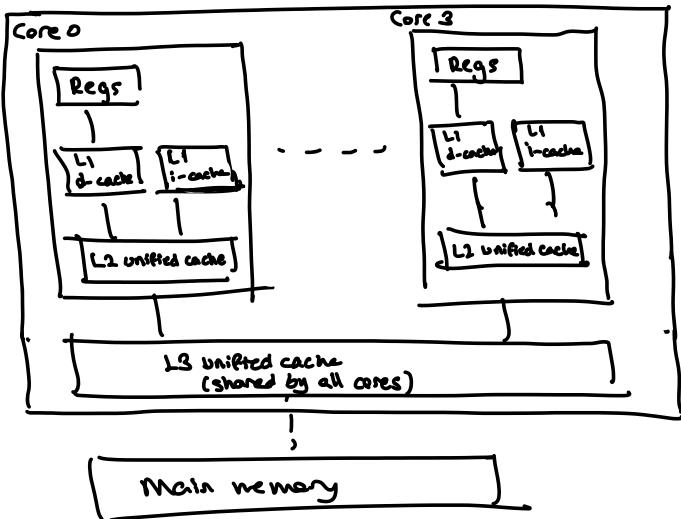


- Write Hit:
- Write through: write immediately to memory
 - Write back: defer write to memory until replacement of LRU (need a dirty bit)

- Write Miss:
- Write allocate: load into cache, update line in cache
 - No-write allocate: writes immediately to memory

Typical:
write through + no write-allocate
write back + write allocate

Intel core i7 Cache Hierarchy



Cache Summary

parameters: $C = S \times E \times B \rightarrow$ data block size
 ↴ ↴
 # of sets # of lines per set
 only the total
 ↴
 # of data bytes are stored

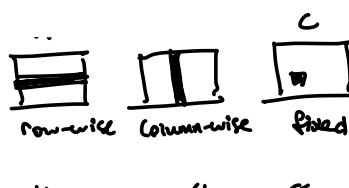
Layout of C arrays in Memory

- allocated in row-major order (each row in contiguous memory locations)
- stepping through columns in 1 row
 $\text{for } (i=0; i < N; i++)$
 $\text{sum} = a[0][i]$
- accesses successive elements
 $a_{0,0}, a_{0,1}, a_{0,2}, a_{0,3}, \dots$
 if block size > 4, exploit spatial locality
 Compulsory miss rate = 4 bytes / 1B
- stepping through rows in 1 column
 $\text{for } (i=0; i < N; i++)$
 $\text{sum} = a[i][0]$
- accesses distant elements
 no spatial locality
 Compulsory miss rate = 1 (1/100)
 $a_{0,0}, a_{1,0}, a_{2,0}, a_{3,0}, \dots$

Matrix Multiplication

		misses per inner loop
ijk loop →	$a[i][k] \rightarrow$	0.25
	$b[k][j] \rightarrow$	1.0
	$c[i][j] \rightarrow$	0.0
$jik \rightarrow$	~	0.25
		1.0
		0.0
$kij \rightarrow$	$a[i][k] \rightarrow$	0.0
	$b[k][j] \rightarrow$	0.25
	$c[i][j] \rightarrow$	0.25
$lki \rightarrow$	$c[i][j] \rightarrow$	1.0
	$a[l][k] \rightarrow$	1.0
	$b[k][j] \rightarrow$	0.0

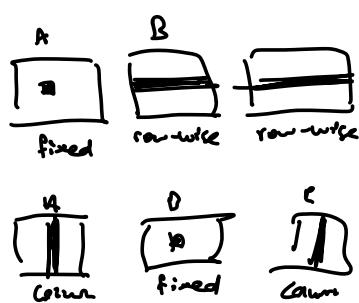
misses per inner loop



iJK ($\&$ jik): 2 loads
 0 stores
 1.25 misses

kij (ijk): 2 loads
 1 store
 0.5 misses

JKi ($\&$ jki): 2 loads
 1 store
 2.0 misses



Direct mapped \rightarrow easy comparison
 2-way \rightarrow too many top comparisons
 but less collisions

Types of cache misses

Compulsory miss:
 occur because cache is empty

Conflict miss:
 occur when the level K cache is large enough, but multiple data objects all map to the same level K block.

Capacity miss:
 occurs when the set of active cache blocks (working set) is larger than the cache.

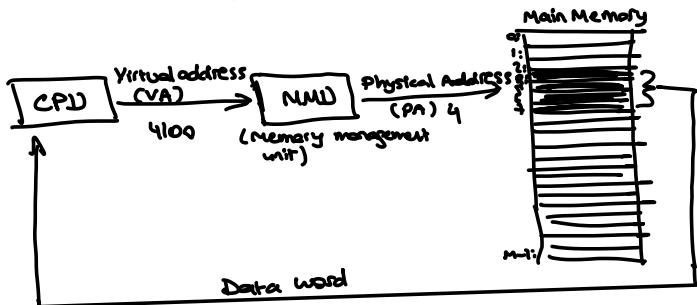


Set
 Tag | Block \rightarrow Address,
 offset
 Maximum # of tag
 comparisons needed

There are several
 cache types, but
 Aybuke didn't
 let me write
 them.

Summary

Virtual Memory



Address Spaces:

Linear address space: ordered set of contiguous non-negative integer addresses.

Virtual address space: Set of $N = 2^n$ virtual addresses $\{0, 1, 2, 3, \dots, N-1\}$

Physical address space: Set of $M = 2^m$ physical addresses $\{0, 1, 2, 3, \dots, M-1\}$

Every byte in main memory has one physical address and 0 or more virtual addresses.

Why Virtual Memory?

① Uses main memory efficiently

- use DRAM as a cache for parts of a virtual address space



② Simplifies memory management

- Each process gets the same uniform linear address space



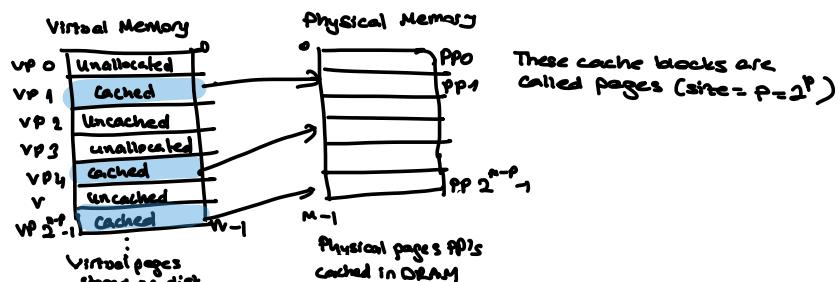
③ Isolates address spaces

- One process can't interfere with another's memory.
- User program cannot access privileged kernel information and code.

VM as a Tool for Caching

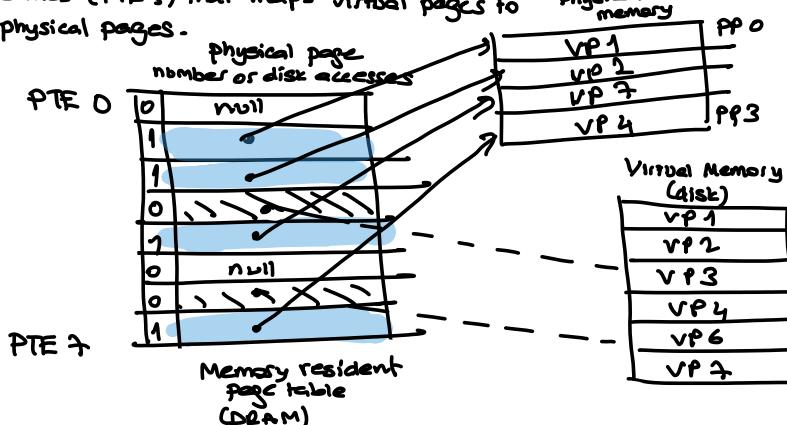
Virtual memory is an array of N contiguous bytes stored on disk.

The contents of the array on disk are cached in physical memory (DRAM cache).



Enabling Data Structure: Page Table

A page table is an array of page table entries (PTE's) that maps virtual pages to physical pages.



PageHit: reference to VM word that is in physical memory (DRAM cache hit)

PageFault: " " " " " that is not in physical memory.

(when it is mapped to a valid bit 0,

page miss causes page fault.

page fault handler selects a victim to be evicted (here VP4) and updates it with the new one.

Offending instruction is restarted → page hit

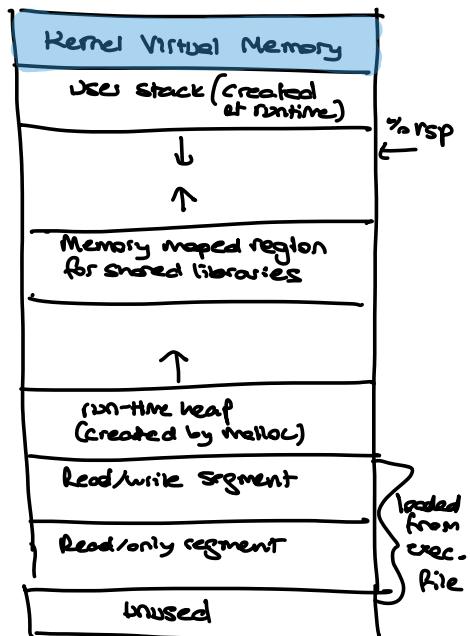
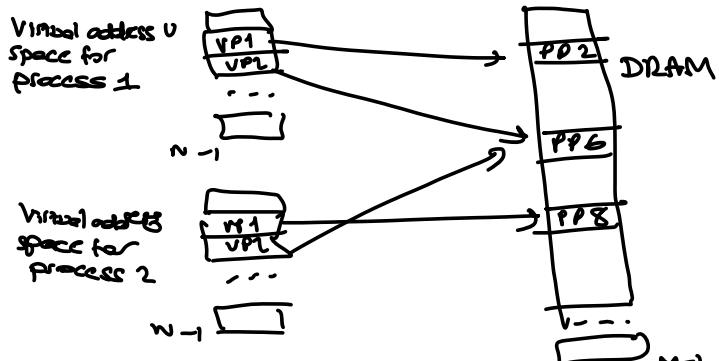
Set of active virtual pages → working set

If working set size < main memory size: good performance after compulsory misses

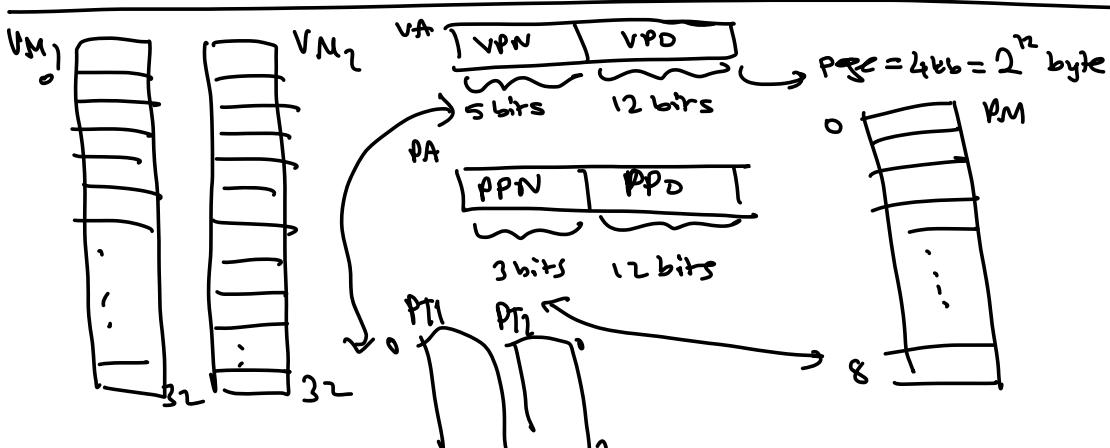
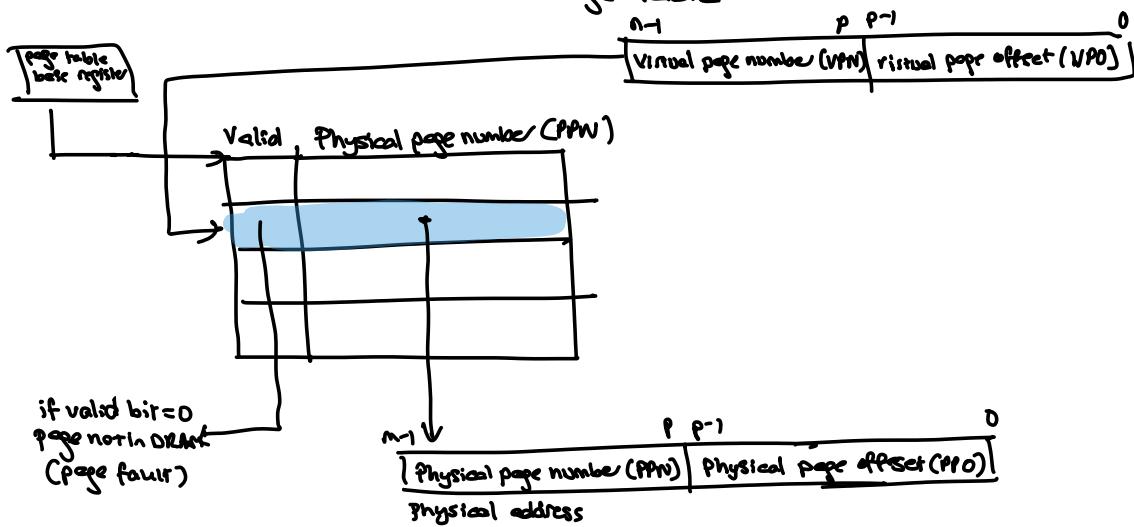
If sum(working set size) > main memory size

thrashing: performance meltdown where pages are swapped in & out continuously

VM as a Tool for Memory Management:
Key Idea: each process has its own virtual address space.



Address Translation with a Page Table



VM as a Tool for Memory Protection

- Extend PTE's with permission bits
- Page Fault Handler checks these before remapping.

	READ	WRITE	EXEC	Address
VP0	No	Yes	No	PP6
VP1	No	Yes	Yes	PP4
VP2	Yes	Yes	Yes	PP2

Speeding up Translation with a TLB \rightarrow Cache for Page Table

PTE's are cached in L1 like any other memory word.

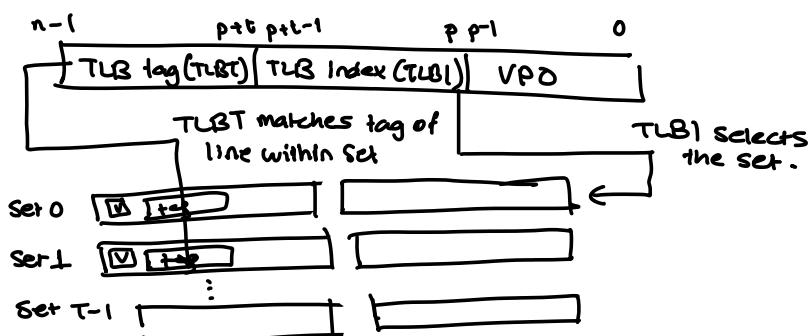
PTE's may be evicted by other data references

PTE hit still requires a small L1 delay

Solution: Translation Lookaside Buffer (TLB)

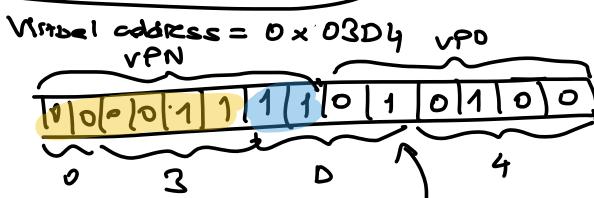
- Small set-associative hardware cache in MMU
- maps virtual page numbers to physical page numbers
- contains complete page table entries for small number of pages.

1) MMU uses the VPN portion of the virtual address to access the TLB.

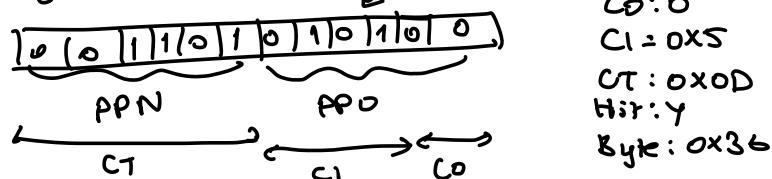


Address Translation Example:

VA \rightarrow VPN \rightarrow TLB \rightarrow PPN \rightarrow PA



Physical Address



Address Translation (Page Hit)

- Processor sends virtual address to MMU.
 - MMU fetches PTE from page table in memory
 - MMU sends physical address to cache/memory
 - Cache/memory sends data word to processor.
- 2 memory accesses

Address Translation (Page Fault)

- Processor sends virtual address to MMU.
- MMU fetches PTE from table in memory.
- Valid bit is 0, so MMU triggers page fault exception.
- Handler identifies victim.
- Handler pages in new page & updates PTE in memory.
- Handler returns to original process restarting faulting instruction.

Size of the PT

32 bit \rightarrow 1A32

1Page = 4 KB (2^{12})

(PM) = 1 GB (2^{30})

(VM) = 4 GB (2^{32})

\rightarrow 20 bits \rightarrow 2 bits

VPN | VPO

32 bits

1PTE = 2^{20} PTE $\Rightarrow 2^{20} \times 2^2 = 2^{22} = 4MB$

AT

1 | PPO

— 18 — 12 —

PTE \Rightarrow 1 | PPN

— 18 bits —

\rightarrow PTE = 32 bits

$\approx 4 bytes$

VPN = 0x0F TLB Hit: Y
TLBI = 0x3 Page Fault: N

TLBT = 0x03

PPN = 0x0D

C0: 0

C1 = 0x5

C2: 0x0

Hit: Y

Byte: 0x36

