

Virtual Memory

- Use main memory as a “cache” for secondary (disk) storage
 - Managed jointly by CPU hardware and the operating system (OS)

Virtualization & Protection

- Programs share main memory
 - Each gets a private virtual address space holding its frequently used code and data
 - Protected from other programs

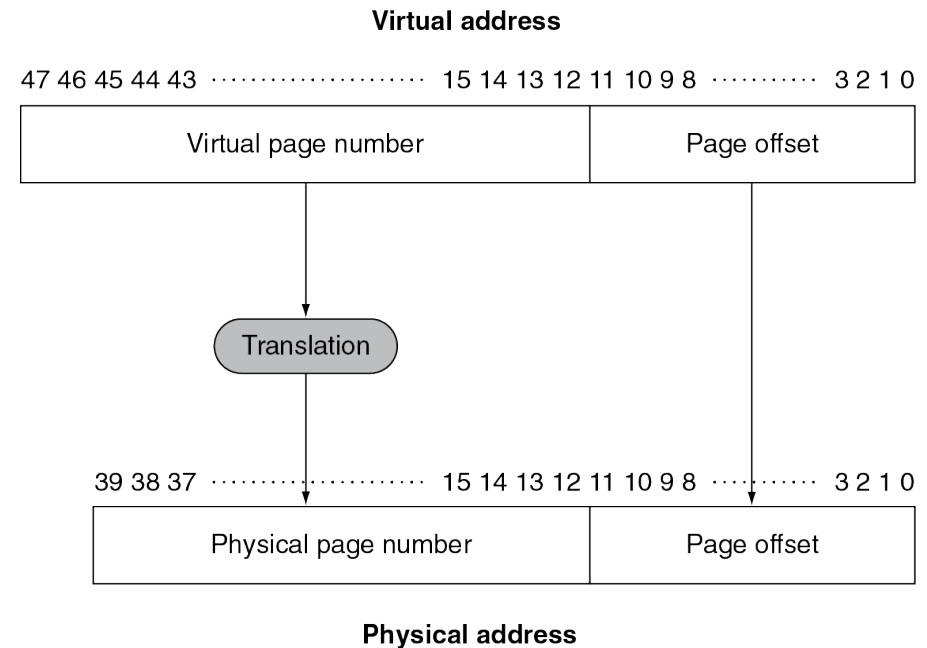
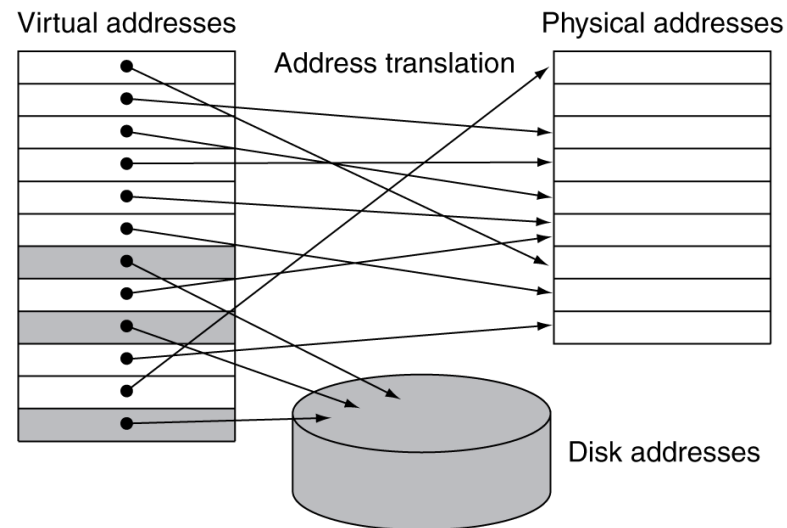
- CPU and OS translate virtual addresses to physical addresses

- VM “block” is called a page
- VM translation “miss” is called a page fault

Also solve the
Fragmentation
problem:
available mem.
for new
process is
fragmented

Address Translation

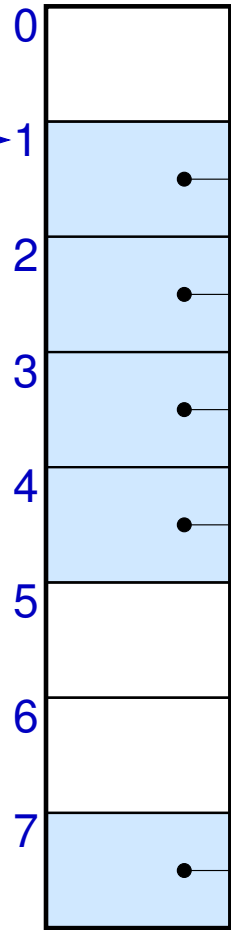
- Fixed-size pages (e.g., 4K)



Page Fault Penalty

- On page fault, the page must be fetched from disk
 - Takes millions of clock cycles
 - Handled by OS code
- Try to minimize page fault rate
 - Fully associative placement
 - Smart replacement algorithms

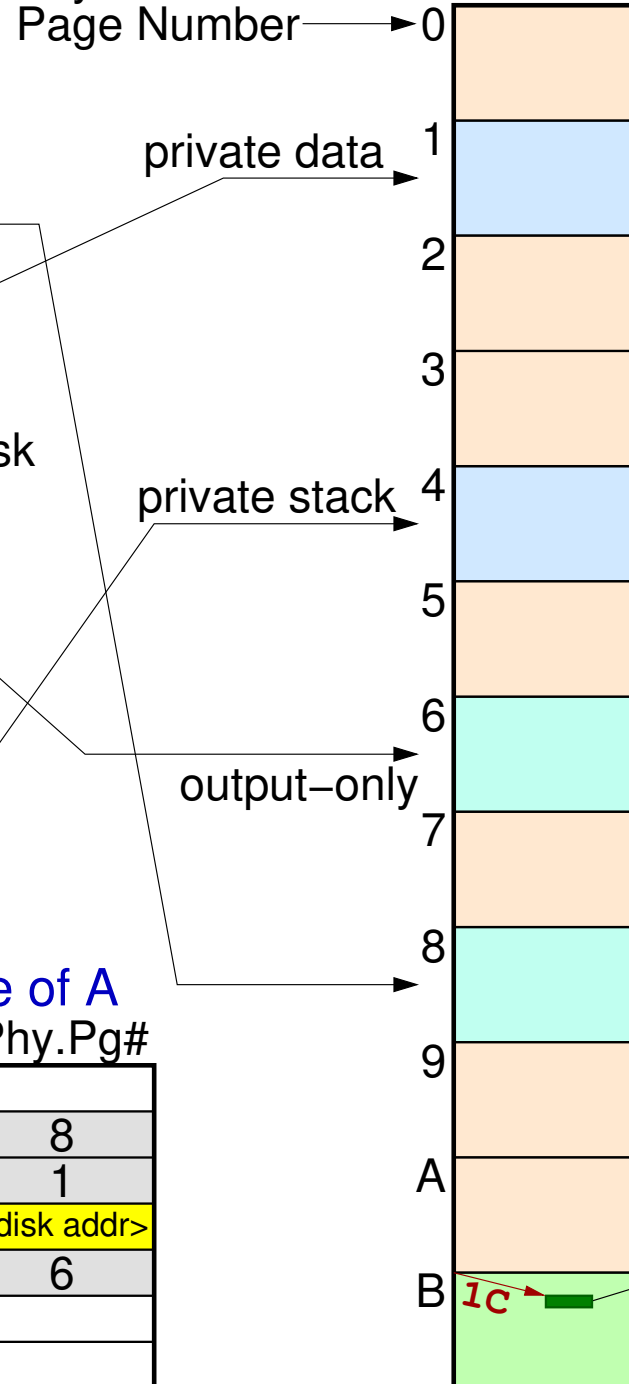
Process A Virtual Memory



Page Table of A

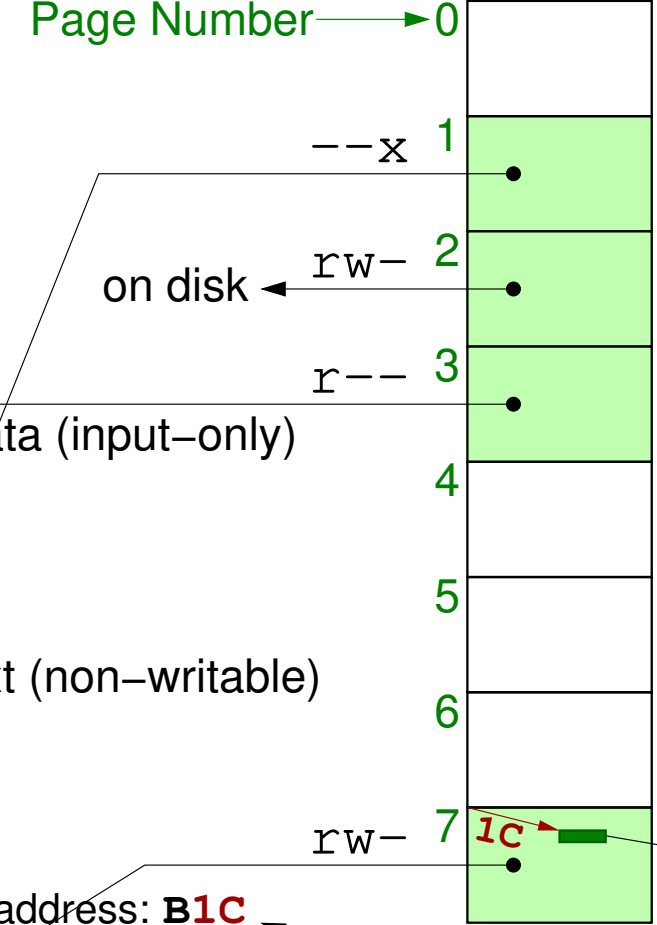
	V	Prot	D	R	Phy.Pg#
0	0	---			
1	1	--x			8
2	1	rw-			1
3	0	r--			<disk addr>
4	1	-w-			6
5	0	---			
6	0	---			
7	1	rw-			4

Physical Memory



Two Processes, A and B, instances of a same program, isolated from each other, except for a shared data page

Process B Virtual Memory



shared data (input-only)

shared text (non-writable)

physical address: **B1C**

private stack

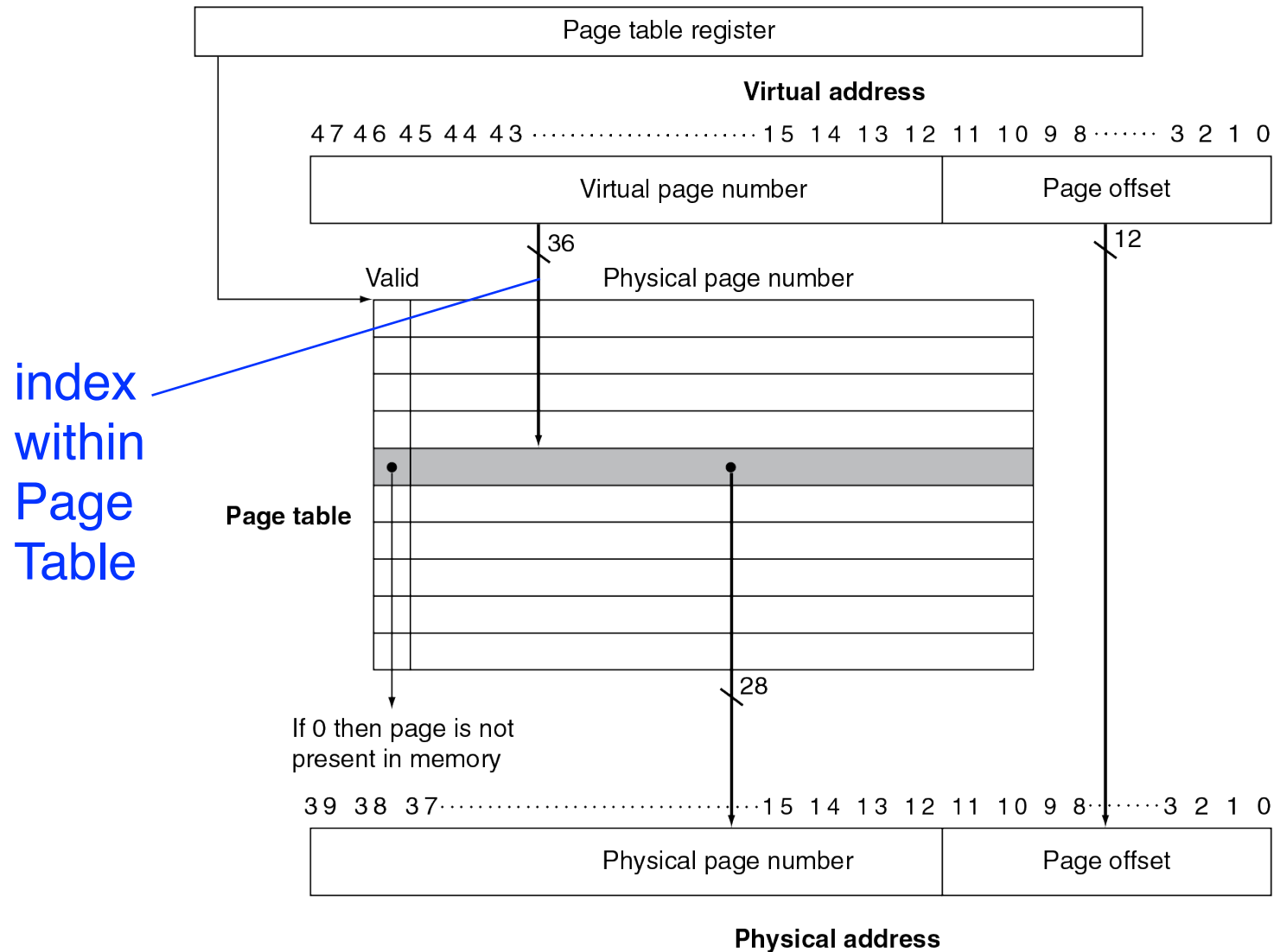
virtual address: **71C**

1C is the offset-within-page

Page Tables (per Process!)

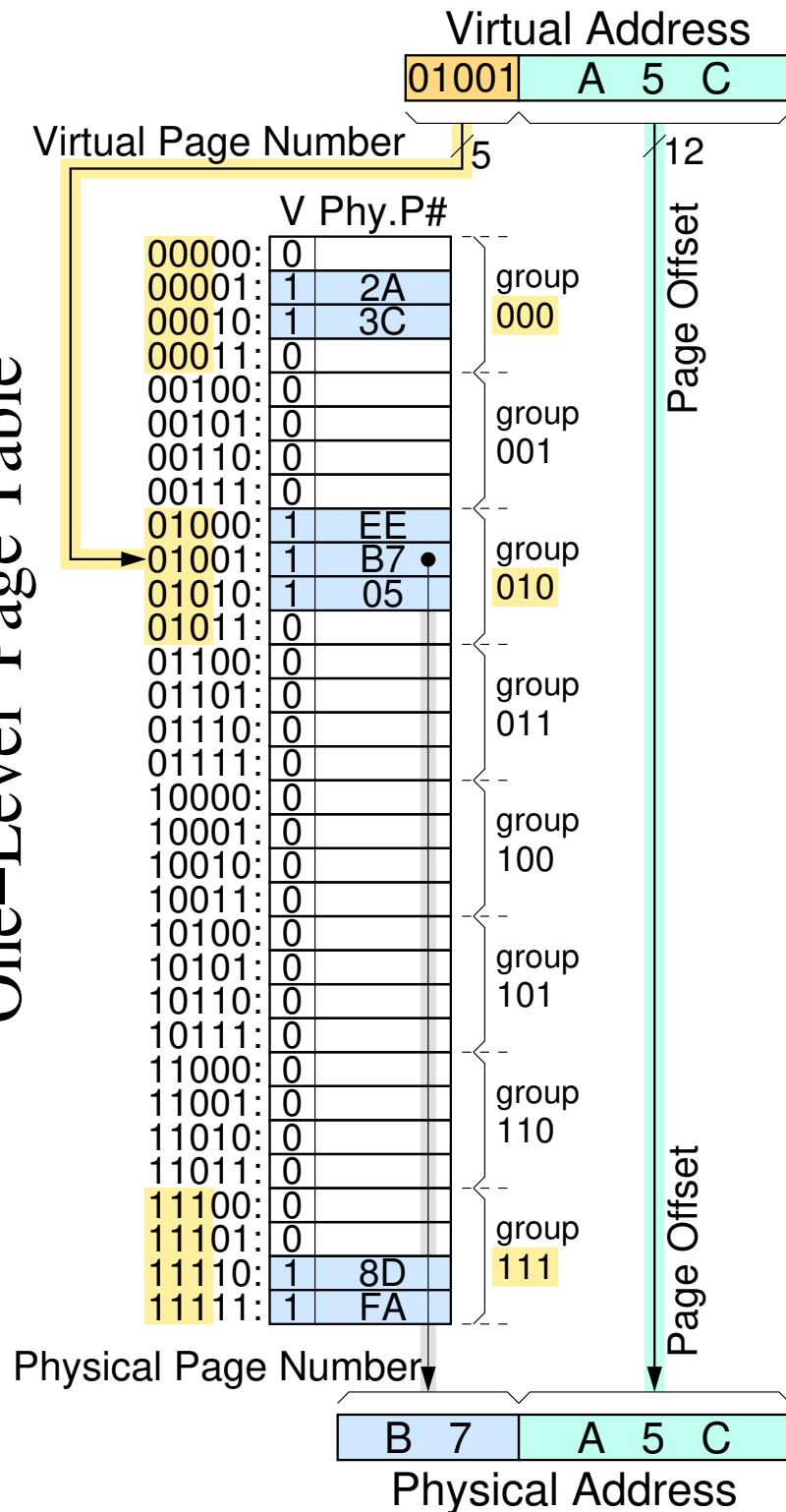
- Stores placement information
 - Array of page table entries, indexed by virtual page number
 - Page table register in CPU points to page table in physical memory (to the page table of the currently running process!)
- If page is present in memory
 - PTE stores the physical page number
 - Plus other status bits (referenced, dirty, ...)
- If page is not present
 - PTE can refer to location in swap space on disk

Translation Using a Page Table

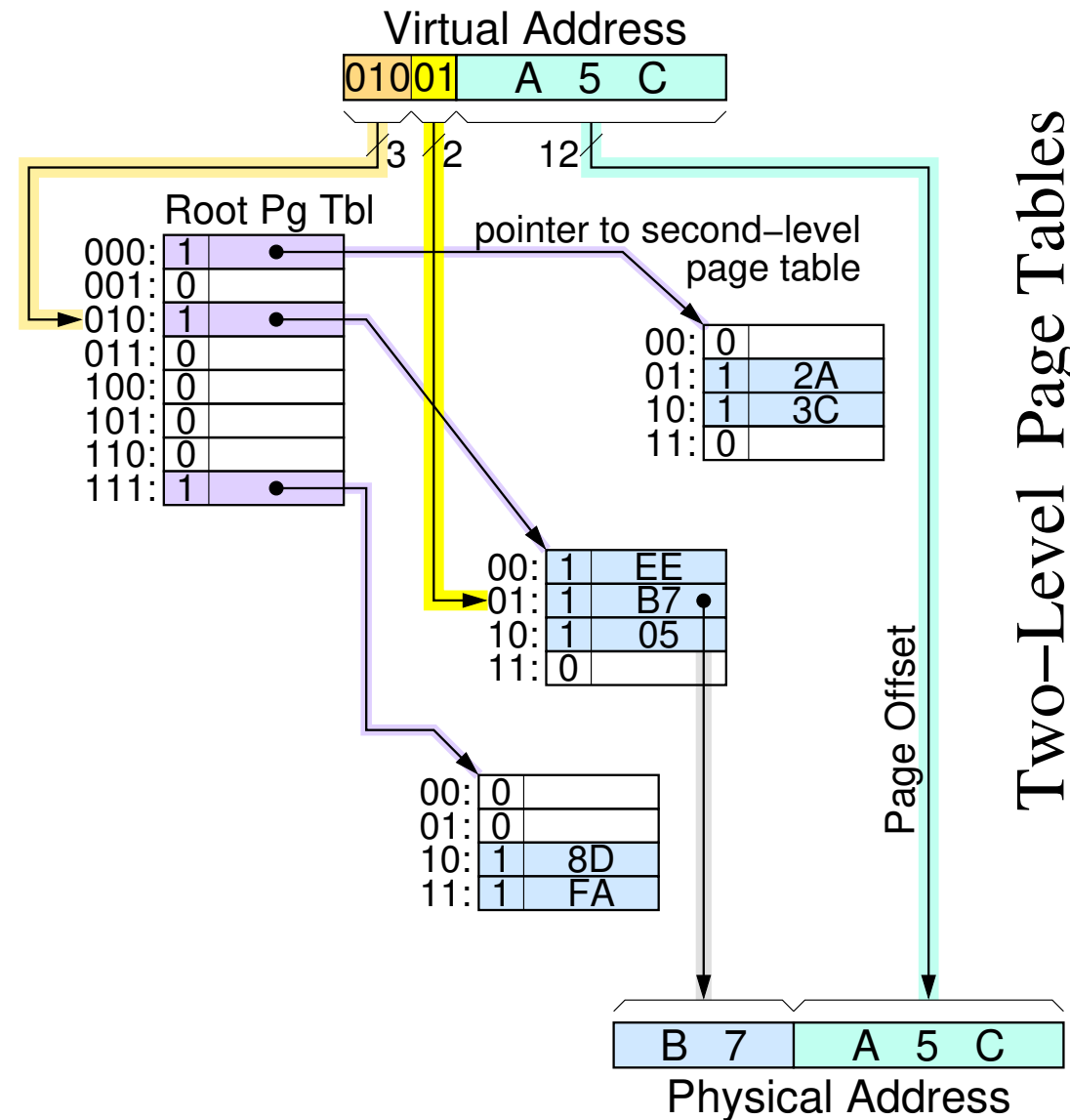


Problem:
Very Large Size
of single-level
Page Table;
Solution:
Multi-Level
Page Tables.

One-Level Page Table



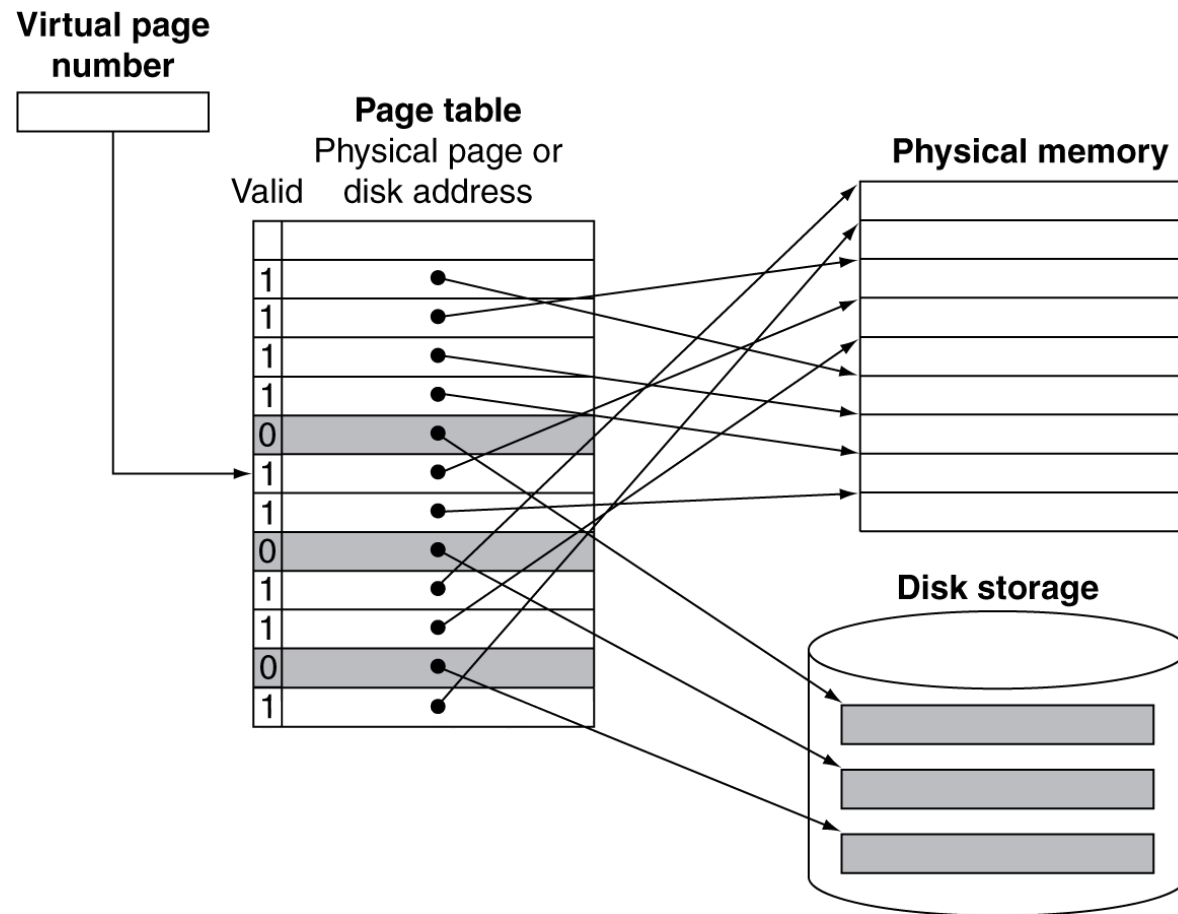
Two-Level Page Tables



In this example:

- Page size: 4 KBytes
- Virtual Address Space: 128 KBytes
=> 32 virtual pages per process
- Physical Address Space: 1 MByte
=> 256 physical pages

Mapping Pages to Storage



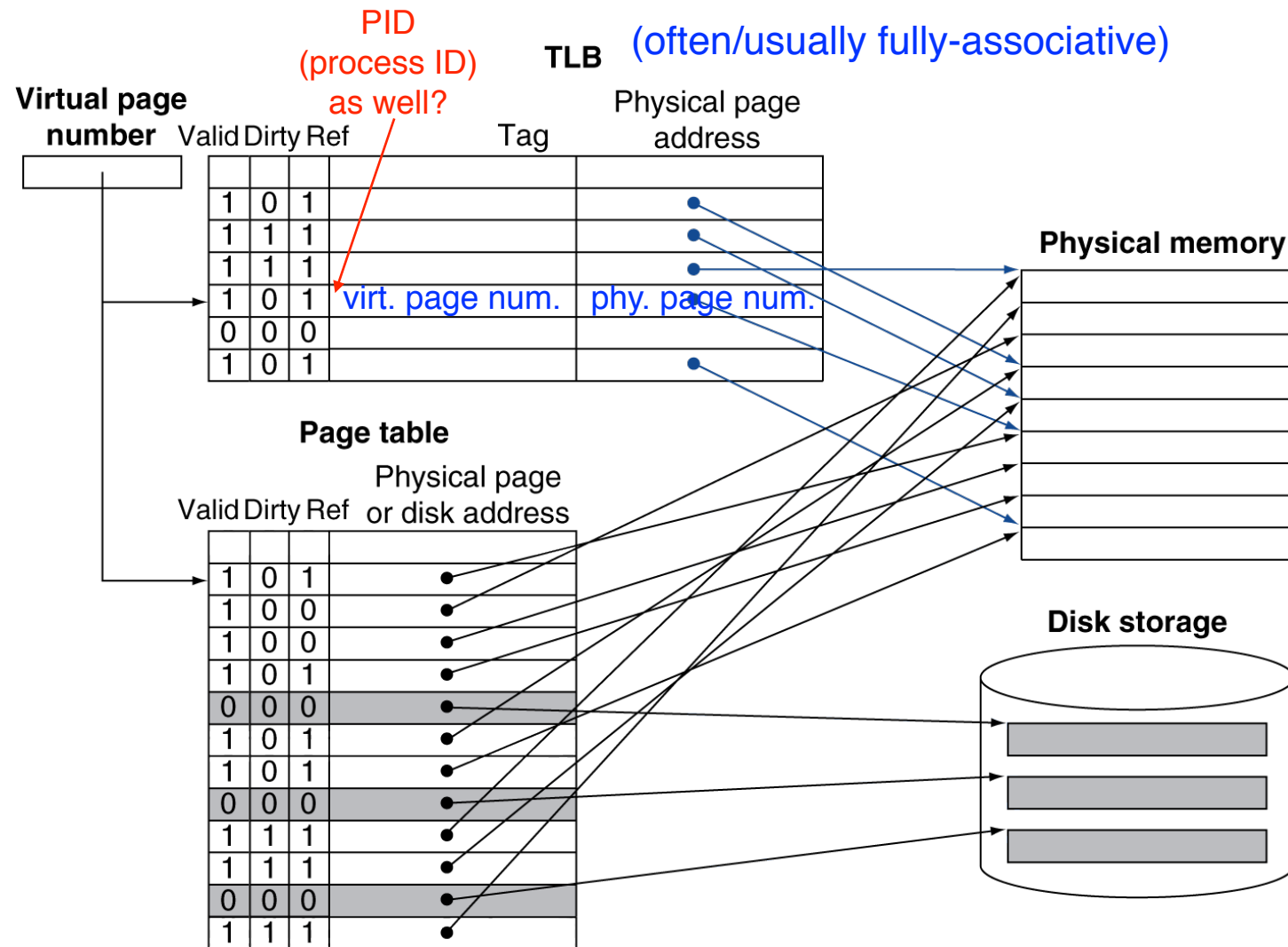
Replacement and Writes

- To reduce page fault rate, prefer least-recently used (LRU) replacement
 - Reference bit (aka use bit) in PTE set to 1 on access to page
 - Periodically cleared to 0 by OS
 - A page with reference bit = 0 has not been used recently
- Disk writes take millions of cycles
 - Block at once, not individual locations
 - Write through is impractical
 - Use write-back
 - Dirty bit in PTE set when page is written

Fast Translation Using a TLB

- Address translation would appear to require extra memory references
 - One to access the PTE (and more for multi-level tables)
 - Then the actual memory access
- But access to page tables has good locality
 - So use a fast cache of PTEs within the CPU
 - Called a Translation Look-aside Buffer (TLB)
 - Typical: 16–512 PTEs, 0.5–1 cycle for hit, 10–100 cycles for miss, 0.01%–1% miss rate
 - Misses could be handled by hardware or software

Fast Translation Using a TLB



TLB Misses

- If page is in memory
 - Load the PTE from memory and retry
 - Could be handled in hardware
 - Can get complex for more complicated page table structures
 - Or in software
 - Raise a special exception, with optimized handler
- If page is not in memory (page fault)
 - OS handles fetching the page and updating the page table
 - Then restart the faulting instruction

Page Table structure
fixed in hardware

TLB Miss Handler

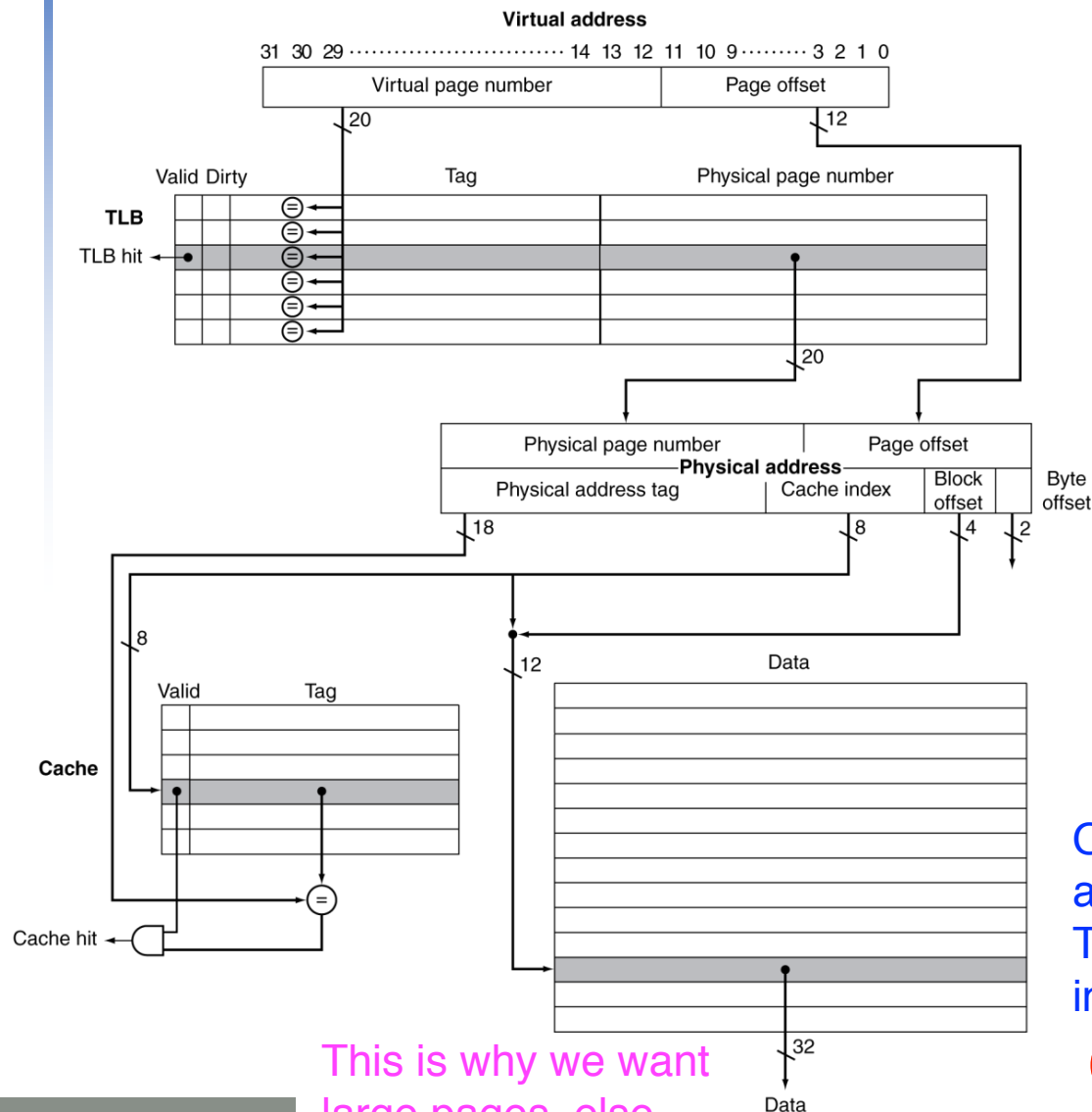
(when TLB misses
are handled in software)

- TLB miss indicates
 - ^{either} Page present, but PTE not in TLB
 - ^{or} Page not present
- Must recognize TLB miss before destination register overwritten ^(in stage 4 of our pipeline, before stage 5)
 - Raise exception
- Handler copies PTE from memory to TLB
 - Then restarts instruction
 - If page not present, page fault will occur

Page Fault Handler

- Use faulting virtual address to find PTE
- Locate page on disk
- Choose page to replace
 - If dirty, write to disk first
- Read page into memory and update page table
- Make process runnable again
 - Restart from faulting instruction

TLB and Cache Interaction



This is why we want large pages, else forced to increase associativity of L1

- If cache tag uses physical address
 - Need to translate before cache lookup
- Alternative: use virtual address tag
 - Complications due to aliasing
 - Different virtual addresses for shared physical address

Often we want: physical addr. cache, and TLB access in parallel with tag read from cache. This requires cache index to be fully contained in page offset bits, which means:

Cache Way Size \leq Page Size

Memory Protection

- Different tasks can share parts of their virtual address spaces

- But need to protect against errant access
- Requires OS assistance

■ Hardware support for OS protection

- Privileged supervisor mode (aka kernel mode)
- Privileged instructions *only executable in supervisor mode*
- Page tables and other state information only accessible in supervisor mode
- System call exception (e.g., ecall in RISC-V)

Supervisor mode can only be entered at (hardwired) exception handler address, only through ecall or exception

Exceptions and Interrupts

- “Unexpected” events requiring change in flow of control
 - Different ISAs use the terms differently
- Exception
 - Arises within the CPU
 - e.g., undefined opcode, syscall, ...
- Interrupt
 - From an external I/O controller
- Dealing with them without sacrificing performance is hard

Handling Exceptions

- Save PC of offending (or interrupted) instruction
 - In RISC-V: Supervisor Exception Program Counter (SEPC)
- Save indication of the problem
 - In RISC-V: Supervisor Exception Cause Register (SCAUSE)
 - 64 bits, but most bits unused
 - Exception code field: 2 for undefined opcode, 12 for hardware malfunction, ...
- Jump to handler
 - Assume at 0000 0000 1C09 0000_{hex}

Many RISC-V computers store the exception entry address in a special register named Supervisor Trap Vector (STVEC), which the OS can load with a value of its choosing.

An Alternate Mechanism

- Vectored Interrupts
 - Handler address determined by the cause
- Exception vector address to be added to a vector table base register:
 - Undefined opcode 00 0100 0000_{two}
 - Hardware malfunction: 01 1000 0000_{two}
 -: ...
- Instructions either
 - Deal with the interrupt, or
 - Jump to real handler

Handler Actions

- Read cause, and transfer to relevant handler
- Determine action required
- If restartable
 - Take corrective action
 - use SEPC to return to program
- Otherwise
 - Terminate program
 - Report error using SEPC, SCAUSE, ...

Exceptions in a Pipeline

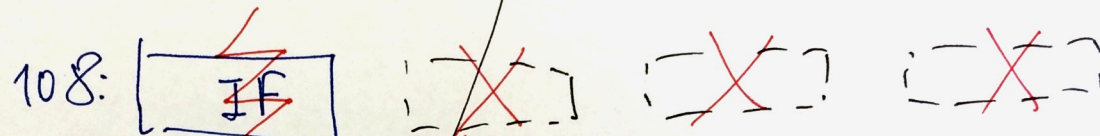
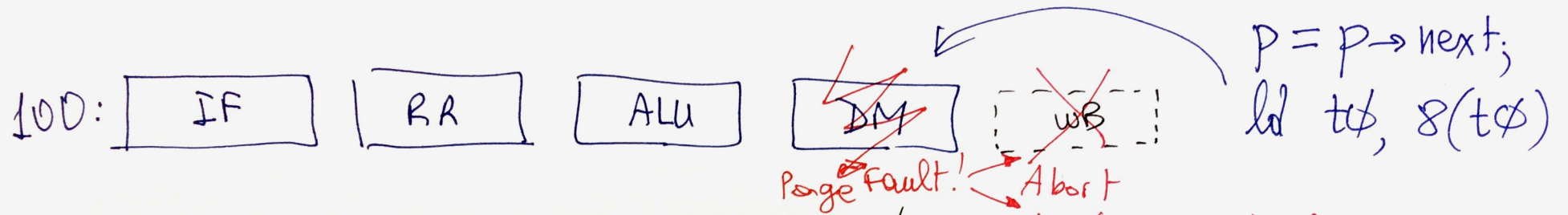
- Another form of control hazard
- Consider malfunction on add in EX stage
add x1, x2, x1
 - Prevent x1 from being clobbered
 - Complete previous instructions
 - Flush add and subsequent instructions
 - Set SEPC and SCAUSE register values
 - Transfer control to handler
- Similar to mispredicted branch
 - Use much of the same hardware

Exception Properties

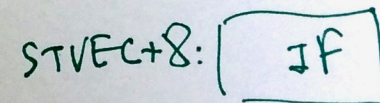
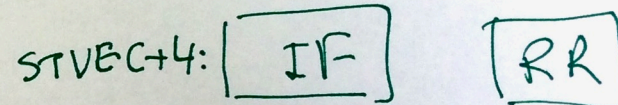
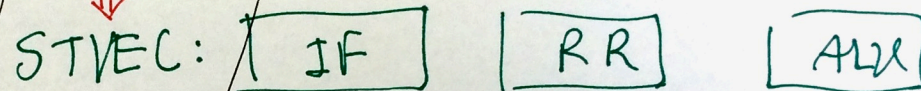
- Restartable exceptions
 - Pipeline can flush the instruction
 - Handler executes, then returns to the instruction
 - Refetched and executed from scratch
- PC saved in SEPC register
 - Identifies causing instruction

Multiple Exceptions

- Pipelining overlaps multiple instructions
 - Could have multiple exceptions at once
- Simple approach: deal with exception from earliest instruction
 - Flush subsequent instructions
 - “Precise” exceptions
- In complex pipelines
 - Multiple instructions issued per cycle
 - Out-of-order completion
 - Maintaining precise exceptions is difficult!



Page Fault! Abort



~~$SEPC \leftarrow 108$~~
 $SEPC \leftarrow 100$

Multiple
Exceptions
(exception from
earlier instruction
occurs after
exception from
later instruction)
in this example