

Page faults.

- $\text{Execution time} = \text{Instruction} \cdot \text{cycles} \cdot \text{seconds}$

Pipeline hazards:

- Structural hazards; attempt to use the same resource by two different instructions at the same time

- Data hazards; attempt to use data before its ready

- An instruction's source operands are produced by prior instructions still in the pipeline

- Control hazards; attempt to make a decision about a program's control flow before the condition has been evaluated as the new PC changes, address calculated

- Branch as jump instructions, exceptions

- Can always resolve hazards by stalling

- Pipeline control must detect the hazard

- And take action to resolve hazards

- Reducing control hazards: move branch point

- Forwarding: delay branching

- or stalling the pipeline

- or forwarding with a buffer

BFB = Branch target buffer

- Cache of target addresses

- Taken by PC when instruction is fetched

- If not a branch, branch predictor

- If not a branch, branch predictor

- If not a branch, branch predictor

Control Hazards:

- You can find a pipeline by introducing

- Flow when jumping, but could impact

- performance. Early in the pipeline

- as possible

- Branch delay slot;

- Branch does not take effect until

- the 8 cycle after its execution

- Compiler can fill the slot with

- a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

- or a NOP

Replacements and writing;

- To reduce page fault rate, prefer least recently

- used (LRU) replacement

- Reference bit (also used) in PTE set as an access

- To page

- primarily chosen to be by OS

- A page with reference bit set has not been used recently

- Disk makes sense in terms of cycles

- Block as one, not individual locations

- write through is important

- use write-back

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

- Dirty bit in PTE set page in memory

memory stall cycles;

- (memory access) (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

- (miss rate) (miss penalty)

Page table size

- Collectible in one of a page table

- 32 bits virtual addresses

- 4 KB pages/page offset = 12 bits

- 4 byte page table entries (PTEs)

- 4 bits for virtual address

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

- 2 bits for page offset

Cache

- 1 cycle hits

- 3 cycle misses

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

- 14 cycle miss penalty

-

Definers:

- PRAMs: Program Access Memory
 - fast, one, low power, cheap
 - but slow
 - connects to the bus when no power
- SRAMs: Static Random Access Memory
 - lower power
 - fast
 - going to be used for cache
 - only when power is on

Flash types:

- NOR: Read is bit cell like a NOR Gate
 - Random read/write access
 - used for instruction memory in embedded systems
- NAND: Read is bit cell like a NAND Gate
 - faster (bit/cell), but block-at-a-time access
 - cheaper per GB
- Flash: Bits wear out after 100k of access
- Not suitable for direct RAM or disk replacement

Caches:

- cache: a small amount of fast, expensive memory
 - gets between processor and dynamic main memory
 - keeps a copy of the most recently used data
 - minimizes speed in access overall, because more than the common case further
 - Reads and writes to the most frequently used addresses will be served by the cache

Locality:

- Temporal locality:
 - if program accesses are many addresses, there is a good chance they will access the same address again
- Spatial locality:
 - if a program accesses are many addresses are many others, there is a good chance it will access other nearby addresses
- Cache hit: occurs if the cache contains the data for the memory location. Hits are good, because the cache can return the data much faster than main memory.
- Cache miss: occurs if the cache does not contain the requested data. This is bad, since the CPU must then wait for the slower memory.
- blocks: caches are divided

- Direct mapped: cache is the simplest approach. each main memory address maps exactly to one cache block.
- tags: supply the rest of the address bits to let us distinguish between different memory locations that map to the same cache block.

- Valid bit: adding for each cache block
- write-back: to help put old data back into the cache and verify its validity

- Cache coherency: when the CPU tries to read from memory, the address will be sent there.

- Least recently used: if the cache fills up, older data is less likely to be requested than newer data.

- block address: if the cache block size is 2 bytes, we can conceptually split the main memory into 2-byte chunks.

- block offset: the lowest 2 bits that indexes into the 2 bytes in the cache block will store the data.

Cache Design:

- fully-associative cache: permits data to be stored in any cache block, instead of having each memory address in one particular block.

- can be placed in any unused block of the cache
- reads a conflict with two or more addresses
- put any 1 or single cache block
- least recently used
- Allows any memory block to map to any location within the cache opposed to direct mapped

- Set-associative: a group of many addresses maps to one cache block
- can be placed in any unused block of the cache
- reads a conflict with two or more addresses
- put any 1 or single cache block
- least recently used
- Allows any memory block to map to any location within the cache opposed to direct mapped

- write-back: queues writes to main memory and permits the CPU to continue.
- especially bad when two cores run at different speeds
- If a processor temporarily goes into bursts of data too quickly for a consumer to handle, the worst data is stored in a buffer until the producer can continue on the other tasks, without waiting for the consumer
- write-back cache: the memory is not updated until the cache block needs to be replaced

- dirty bits: to indicate whether or not it must be saved to main memory before being replaced - otherwise, perhaps unnecessary write backs
- write-through: try to write to an address that is not already contained in the cache
- write-around policy: the write operation goes directly to main memory without affecting the cache.
- Allocate or write: load the newly written data into the cache.

- Performance:
 - include here: the memory, or split it into chunks that can be accessed independently

- Virtual memory:
 - Translational Look-aside Buffer (TLB): part of the cache that stores recent translations of virtual to physical addresses

- Cache as PLB's within the CPU

- Question 5:
 - For a byte addressable machine with a 16-bit address space, how many bits of data does the cache hold? (Do not count read-data bits)

- Cache parameters: Direct mapped, block size is one byte, Cache index is the four least significant bits of the address

- 16 blocks, 8 bits = 128 bits

- For a byte addressable machine with 16-bit address space, how many bits of meta-data does the following cache hold?

- Cache parameters: block size is one byte, Direct mapped, Cache index is the four least significant bits of the address

- 16 blocks, (12 + 16 bits) = 208 bits

- Tag = 12 bits, Valid bit = 1 bit
- Index: 4 more bits are not stored in the cache

- fully-associative vs. set-associative cache: for the following memory reference, what are the hits/miss results for a 4-sets fully-associative cache and a 2-way, 2-sets set-associative cache?

- Assume the caches are initially empty. Least recently used policy and each cache block is one byte.

	M[0]	M[1]	M[2]	M[3]	M[4]	M[5]	M[6]	M[7]	M[8]	M[9]	M[10]	M[11]	M[12]	M[13]	M[14]	M[15]
fully	M	M	M	M	M	M	M	M	M	M	M	M	M	M	M	M
4-sets	M	M	M	M	M	M	M	M	M	M	M	M	M	M	M	M
2-sets	M	M	M	M	M	M	M	M	M	M	M	M	M	M	M	M

- Question 6:
 - What is the purpose of TLB?
 - It is used to speed up the virtual to physical page translation. It holds the PTEs.
 - Calculate the memory size of a page table for 64-bit virtual address space that uses 4KB pages and has 60 bits per PTE entry (PTEs)?

- $2^{64} / 2^{12} = 2^{52}$ entries
- $2^{52} \cdot 2^4 = 2^{56}$ bytes in total = 64 PB
- What would the minimum page size need to be for the following virtually indexed, physically tagged cache not to have a binary?

- Assume byte-addressable memory space.
- Block size: 32B
- Total data size: 64KB
- Associativity: 8-way
- 56 bits for block address
- # of sets = $2^{16} / (2^5 \cdot 2^3) = 2^6$ sets
- Page offset = 5 + 8 = 13 bits
- = 213
- = 8 KB minimum page size

- This memory type is used to store frequently used data or program memory blocks to speed up the program. Cache accesses are much faster than DRAM accesses.

- Consider two cache structures, both initially empty.

Index	Tag	Valid	Tag	Valid
0				
1				
2				
3				
4				
5				
6				
7				

- For the following memory accesses, what are the hit and miss results for two different caches?

- mem[2], mem[5], mem[6], mem[10]
- mem[5], mem[2], mem[1], mem[11]
- mem[5], mem[6]

- Small cache → Miss, miss, miss, hit, hit, miss, miss, miss, miss, miss
- 1/10

- Large cache → miss, miss, hit, miss, miss, hit, miss, miss, miss, hit, hit