## Improving on Caches CS448



### #5: Hardware Prefetch

- Get proactive!
- Modify our hardware to prefetch into the cache instructions and data we are likely to use
  - Alpha AXP 21064 fetches two blocks on a miss from the Icache
    - Requested block and the next consecutive block
    - Consecutive block catches 15-25% of misses on a 4K direct mapped cache, can improve with fetching multiple blocks

- Similar approach on data accesses not so good, however
- Works well if we have extra memory bandwidth that is unused
- Not so good if the prefetch slows down instructions trying to get to memory



![](_page_2_Figure_0.jpeg)

![](_page_2_Figure_1.jpeg)

![](_page_3_Figure_0.jpeg)

![](_page_3_Figure_1.jpeg)

![](_page_4_Figure_0.jpeg)

![](_page_4_Figure_1.jpeg)

### Problems with Write Buffers

- Consider this code sequence
  - SW 512(R0), R3  $\leftarrow$  Maps to cache index 0
  - LW R1, 1024(R0) ← Maps to cache index 0
  - LW R2, 512(R0)  $\leftarrow$  Maps to cache index 0
- There is a RAW data hazard
  - Store is put into write buffer
  - First load puts data from M[1024] into cache index 0
  - Second load results in a miss, if the write buffer isn't done writing, the read of M[512] could put the old value in the cache and then R2
- Solutions
  - Make the read wait for write to finish
  - Check the write buffer for contents first, associative memory<sup>11</sup>

### #2 Other Ways to Reduce Miss Penalties

- Sub-Block Placement
  - Large blocks reduces tag storage and increases spatial locality, but more collisions and a higher penalty in transferring big chunks of data
  - Compromise is Sub-Blocks
  - Add a "valid" bit to units smaller than the full block, called sub-blocks
    - Allow a single sub-block to be read on a miss to reduce transfer time
    - In other modes of operation, we fetch a regular-sized block to get the benefits of more spatial locality

### #3 Early Restart & Critical Word First

- CPU often needs just one word of a block at a time
  - Idea : Don't wait for full block to load, just pass on the requested word to the CPU and finish filling up the block while the CPU processes the data
- Early Start
  - As soon as the requested word of the block arrives, send it to the CPU
- Critical Word First
  - Request the missed word first from memory and send it to the CPU as soon as it arrives; let the CPU continue execution while filling in the rest of the block

13

# <section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item>

### #5 Second Level Caches

- Probably the best miss-penalty reduction technique, but does throw in a few extra complications on the analysis side...
- L1 = Level 1 cache, L2 = Level 2 cache

Average \_ Memory \_ Access \_ Time = Hit \_ Time(L1) + Miss \_ Rate(L1) × Miss \_ Penalty(L1)

 $Miss\_Penalty(L1) = Hit\_Time(L2) + Miss\_Rate(L2) \times Miss\_Penalty(L2)$ 

• Combining gives:

Average \_ Memory \_ Access \_ Time = Hit \_ Time(L1) + Miss \_ Rate(L1) × (Hit \_ Time(L2) + Miss \_ Rate(L2) × Miss \_ Penalty(L2))

- little to be done for compulsory misses and the penalty goes up
- capacity misses in L1 end up with a significant penalty reduction since they likely will get supplied from L2
- conflict misses in L1 will get supplied by L2 unless they also conflict in L2

![](_page_7_Figure_10.jpeg)

![](_page_8_Figure_0.jpeg)

![](_page_8_Figure_1.jpeg)

![](_page_9_Figure_0.jpeg)

![](_page_9_Figure_1.jpeg)

## **Reducing Hit Time**

- We've seen ways to reduce misses, and reduce the penalty.. next is reducing the hit time
- #1 Simplest technique: Small and Simple Cache
  - Small  $\rightarrow$  Faster, less to search
  - Must be small enough to fit on-chip
    - Some compromises to keep tags on chip, data off chip but not used today with the shrinking manufacturing process
  - Use direct-mapped cache
    - Choice if we want an aggressive cycle time
    - Trades off hit time for miss rate, since set-associative has a better miss rate

21

### #2 Virtual Caches • Virtual Memory - Map a virtual address to a physical address or to disk, allowing a virtual memory to be larger than physical memory More on virtual memory later Traditional caches or Physical caches - Take a physical address and look it up in the cache Virtual caches Same idea as physical caches, but start with the virtual address instead of the physical address - If data is in the cache, it avoids the costly lookup to map from a virtual address to a physical address • Actually, we still need to the do the translation to make sure there is no protection fault Too good to be true? 22

### Virtual Cache Problems

- Process Switching
  - When a process is switched, the same virtual address from a previous process can now refer to a different physical addresses
    - Cache must be flushed
    - Too expensive to safe the whole cache and re-load it
    - One solution: add PID's to the cache tag so we know what process goes with what cache entry
  - Comparison of results and the penalty on the next slide

![](_page_11_Figure_7.jpeg)

### More Virtual Cache Problems...

- Aliasing
  - Two processes might access different virtual addresses that are really the same physical address
  - Duplicate values in the virtual cache
  - Anti-aliasing hardware guarantees every cache block has a unique physical address
- Memory-Mapped I/O
  - Would also need to map memory-mapped I/O devices to a virtual address to interact with them
- Despite these issues...
  - Virtual caches used in some of today's processors
    - Alpha, HP...

25

## #3 Pipelining Writes for Fast Hits

- Write hits take longer than read hits
  - Need to check the tags first before writing data to avoid writing to the wrong address
  - To speed up the process we can pipeline the writes (Alpha)
    - First, split up the tags and the data to address each independently
    - On a write, cache compares the tag with the write address
    - Write to the data portion of the cache can occur in parallel with a comparison of some other tag
      - We just overlapped two stages
    - · Allows back-to-back writes to finish one per clock cycle
- Reads play no part in this pipeline, can already operate in parallel with the tag check

Technique	Miss rate	Miss penalty	Hit time	Hardware complexity	Comment
Larger block size	+	-		0	Trivial; RS/6000 550 uses 128
Higher associativity	+		-	1	e.g., MIPS R10000 is 4-way
Victim caches	+			2	Similar technique in HP 7200
Pseudo-associative caches	+			2	Used in L2 of MIPS R10000
Hardware prefetching of instructions and data	+			2	Data are harder to prefetch; tried in a few machines; Alpha 21064
Compiler-controlled prefetching	+			3	Needs nonblocking cache too; several machines support it
Compiler techniques to reduce cache misses	+			0	Software is challenge; some ma- chines give compiler option
Giving priority to read misses over writes		+		1	Trivial for uniprocessor, and widely used
Subblock placement		+		1	Used primarily to reduce tags
Early restart and critical word first		+		2	Used in MIPS R10000, IBM 620
Nonblocking caches		+		3	Used in Alpha 21064, R10000
Second-level caches		+		2	Costly hardware; harder if block size L1 ≠ L2; widely used
Small and simple caches	-		+	0	Trivial; widely used
Avoiding address translation during indexing of the cache			+	2	Trivial if small cache; used in Alpha 21064
Pipelining writes for fast write			+	1	Used in Alpha 21064