











Cache

- Two questions to answer (in hardware):
 - Q1: How do we know if a data item is in the cache?
 - Q2: If it is, how do we find it?

Direct mapped

 For each item of data at the lower level, there is exactly one location in the cache where it might be - so lots of items at the lower level share locations in the upper level

Address mapping:

- (block address) modulo (# of blocks in the cache)
- First consider when block is one word











Handling Cache Hits

- Read hits (I\$ and D\$)
- this is what we want!
- Write hits (D\$ only)
 - allow cache and memory to be inconsistent
 - write the data only into the cache block (write-back the cache contents to the next level in the memory hierarchy when that cache block is "evicted")
 - need a dirty bit for each data cache block to tell if it needs to be written back to memory when it is evicted
 - require the cache and memory to be consistent
 - always write the data into both the cache block and the next level in the memory hierarchy (write-through) so don't need a dirty bit
 - writes run at the speed of the next level in the memory hierarchy so slow! – or use a write buffer, a buffer that holds the data that is waiting to be written, so would only have to stall if the write buffer is full









Sources of Cache Misses Compulsory (cold start or process migration, first reference): First access to a block, "cold" fact of life, not a whole lot you can do about it If you are going to run "billions" of instruction, compulsory misses are insignificant Conflict (collision): Multiple memory locations mapped to the same cache location

- Solution 1: increase cache size
- Solution 2: increase associativity

□ Capacity:

- Cache cannot contain all blocks accessed by the program
- Solution: increase cache size

Handling Cache Misses

Read misses (I\$ and D\$)

- stall the entire pipeline, fetch the block from the next level in the memory hierarchy, install it in the cache, then restart the pipeline and reissue the read request
- Write misses (D\$ only)
 - stall the pipeline, fetch the block from next level in the memory hierarchy, install it in the cache (which may involve having to evict a dirty block if a write-back cache), then restart the pipeline and reissue the write request
 - or (normally used in write-back caches)
 - Write allocate just write the word into the cache updating both the tag and data, no need to check for cache hit, no need to stall or (normally used in write-through caches with a write buffer)
 - No-write allocate skip the cache write and just write the word to the write buffer (and eventually to the next memory level), no need to stall if the write buffer isn't full; must invalidate the cache block since it will become inconsistent (on a hit)











Multiword Block Considerations

- Read misses (I\$ and D\$)
 - Processed the same as for single word blocks a miss returns the entire block from memory
 - Miss penalty grows as block size grows
 - Latency to first word in block + transfer time for remaining words
 Early restart resume execution as soon as the requested word of the block is returned
 - Requested word first requested word is transferred from the memory to the cache first
 - Nonblocking cache allows the processor to access the cache while the cache is handling an earlier miss

Write misses (D\$)

 Can't use write allocate or will end up with a "garbled" block in the cache (e.g., for 4 word blocks, a new tag, one word of data from the new block, and three words of data from the old block), so must fetch the block from memory first and pay the stall time

Cache Summary

• The Principle of Locality:

- Program likely to access a relatively small portion of the address space at any instant of time
 - Temporal Locality: Locality in TimeSpatial Locality: Locality in Space
 - opara zooany. zooany in opaso
- Three major categories of cache misses:
 - Compulsory misses: sad facts of life. Example: cold start misses
 Conflict misses: increase cache size and/or associativity
 - Nightmare Scenario: ping pong effect! • Capacity misses: increase cache size
- Cache design space

Cache design space

- total size, block size, associativity (replacement policy)
- write-hit policy (write-through, write-back)
- write-miss policy (write allocate, write buffers)