



CS24: INTRODUCTION TO COMPUTING SYSTEMS

Spring 2016

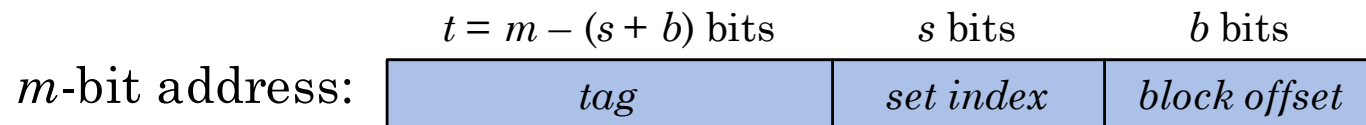
Lecture 15

LAST TIME: CACHE ORGANIZATION

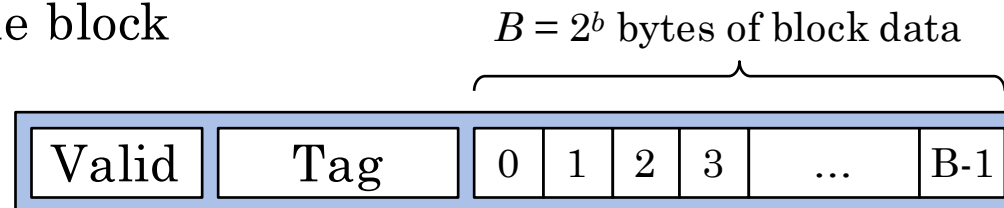
- Caches have several important parameters
 - $B = 2^b$ bytes to store the block in each cache line
 - $S = 2^s$ cache sets
 - E cache lines per set
 - Both S and B are powers of 2
- Cache is designed to work against a main memory with M bytes
 - $m = \log_2(M)$ bits in addresses to main memory
- Cache uses the m -bit address to determine:
 - Which cache set the memory block may appear in
 - The tag for the memory block, to see if the block is actually in the cache set

LAST TIME: MAPPING CACHE BLOCKS

- Relevant parameters for main memory and cache:
 - $m = \log_2(M)$ bits in addresses to main memory
 - $B = 2^b$ bytes to store the block in each cache line
 - $S = 2^s$ cache sets in the cache
- When CPU accesses a data value:



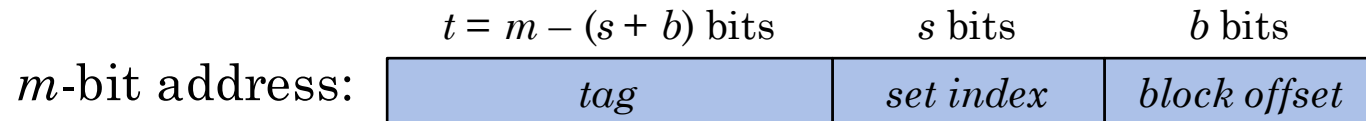
- Bottom-most b bits of the address specify offset of data value within the block



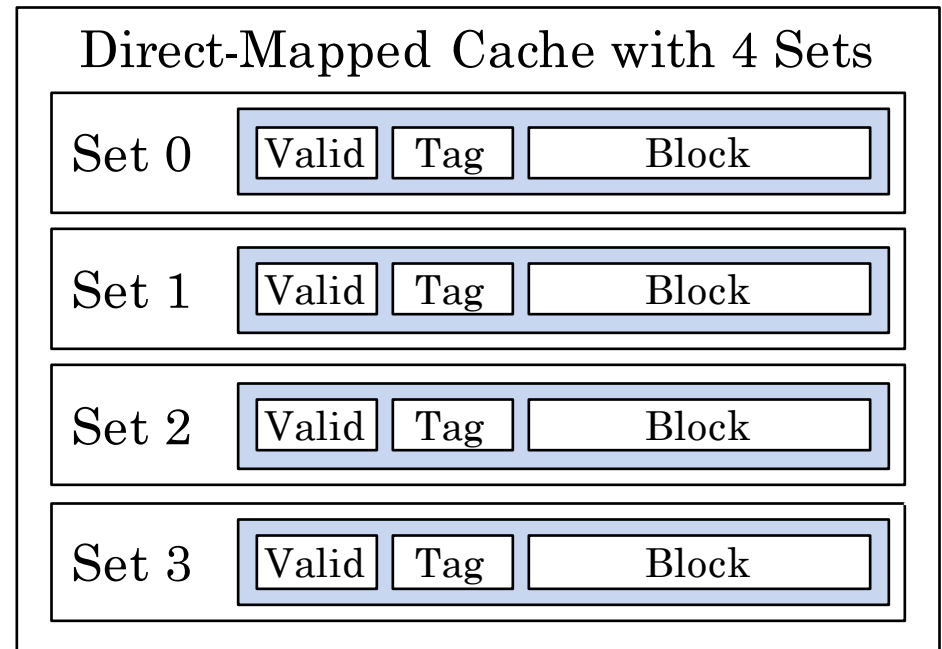
- Middle s bits specify the cache set where the block resides
- Remaining topmost bits constitute the block's tag
 - (Must be able to uniquely identify *all* blocks that can be cached)

DIRECT-MAPPED CACHES

- Direct-mapped caches have S sets, but each set contains only one cache line ($E = 1$)

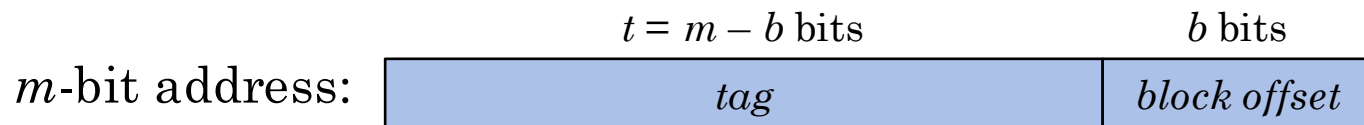


- Example: direct-mapped cache with 4 sets
 - 2 bits in set index
- Pros: Fast, very simple to implement in logic
- Cons: Can easily yield conflict-misses in code with otherwise good locality



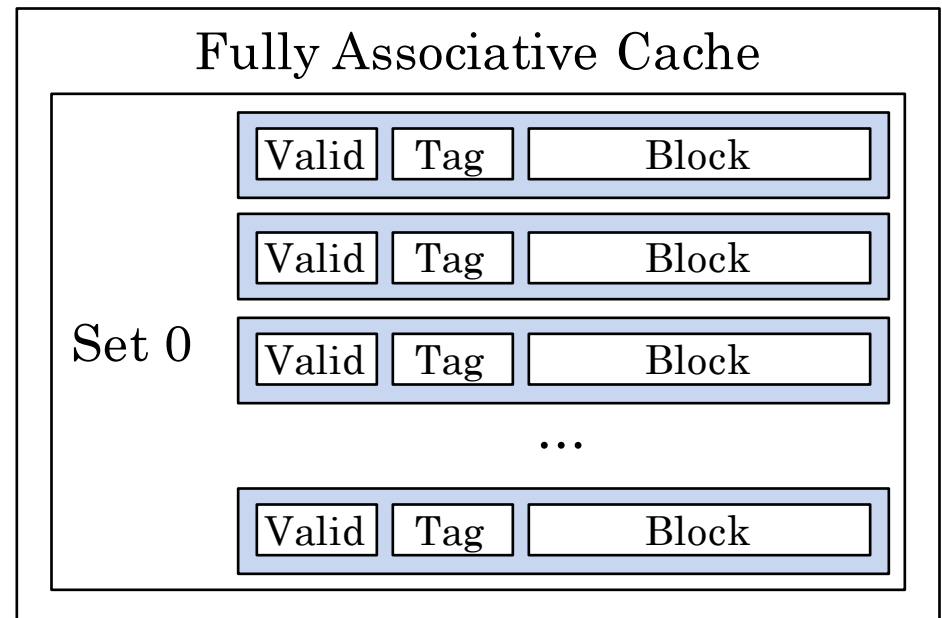
FULLY ASSOCIATIVE CACHES

- Fully associative caches have only one set, which contains all cache lines
 - $S = 2^s = 1 \Rightarrow s = 0$. No bits used for set index!



- Example: fully-associative cache

- Pros: No conflict-misses
- Cons: Slower; difficult to build in hardware!
 - Must examine the tags of all cache lines to see if a block is in the cache



SET-ASSOCIATIVE CACHES

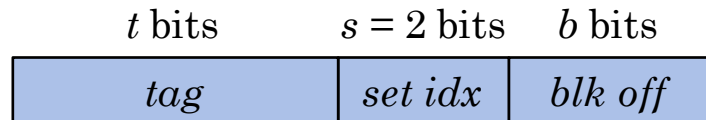
- Set-associative caches combine capabilities of both approaches into a single cache
 - Employ S cache sets, where $S > 1$
 - Each cache set contains E cache lines, where $E > 1$
- Achieves benefits of both techniques:
 - Significantly reduces potential for conflict misses
 - Limits the complexity of the logic that has to match block tags

SET-ASSOCIATIVE CACHES (2)

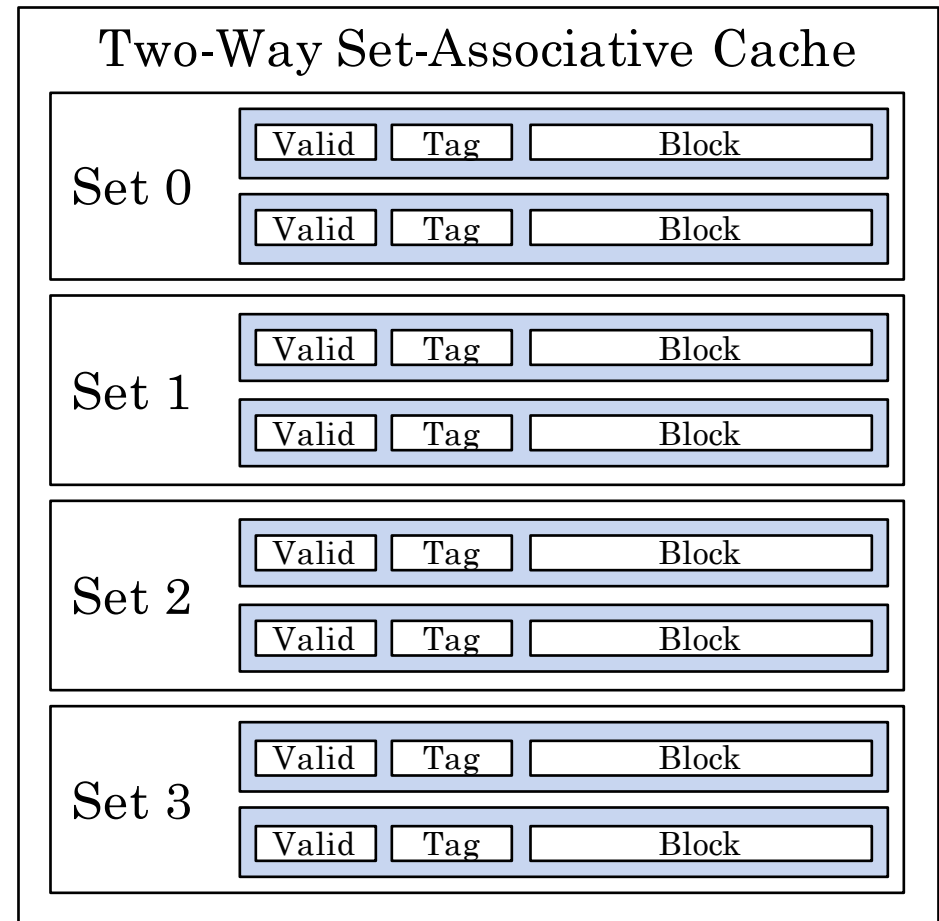
- Example: a two-way set-associative cache

- E is number of cache lines in each set...
- Cache is called “ E -way set-associative cache” when $S > 1$
- $E = 2$ for this example

- For an m -bit address:

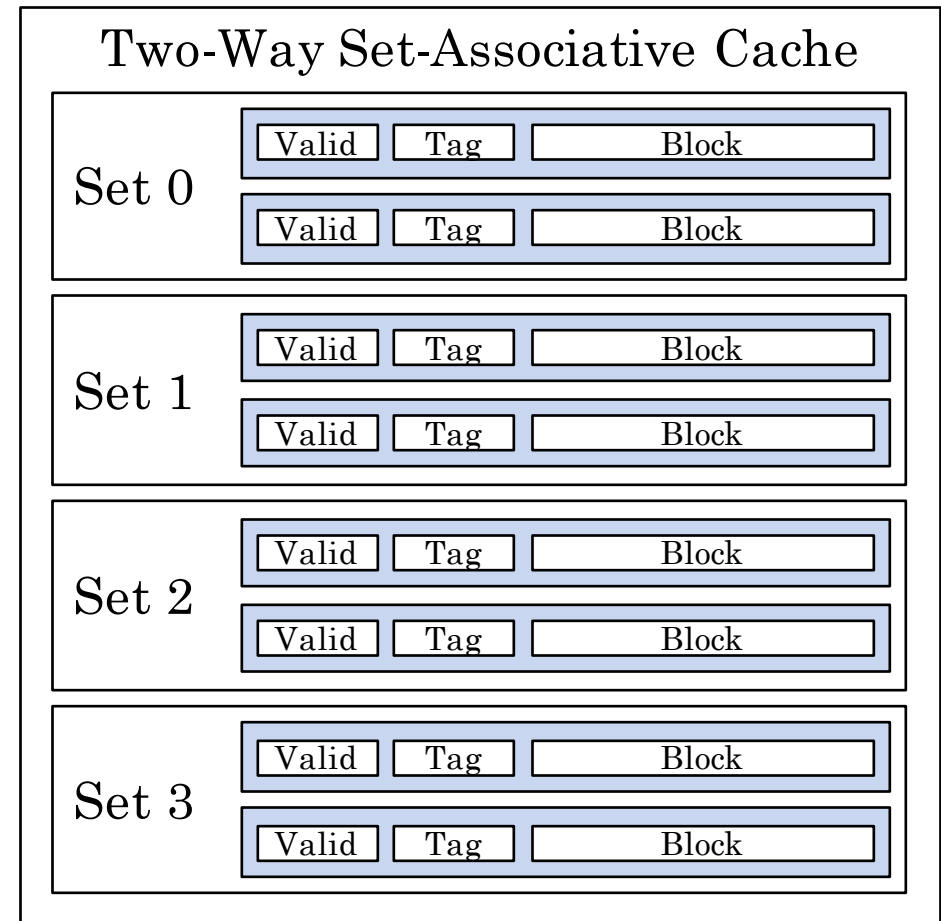


- Use set index to find cache set to examine
- Only have to check tags on small number of lines



LINE REPLACEMENT POLICIES

- When a cache miss occurs, must load a new block into the cache
 - Store in some cache line
- For direct-mapped caches, this is easy
 - Only one line per set
- For set-associative and fully associative caches, we get to choose a line!
- Replacement policy controls which cache line to evict when new block is loaded into cache

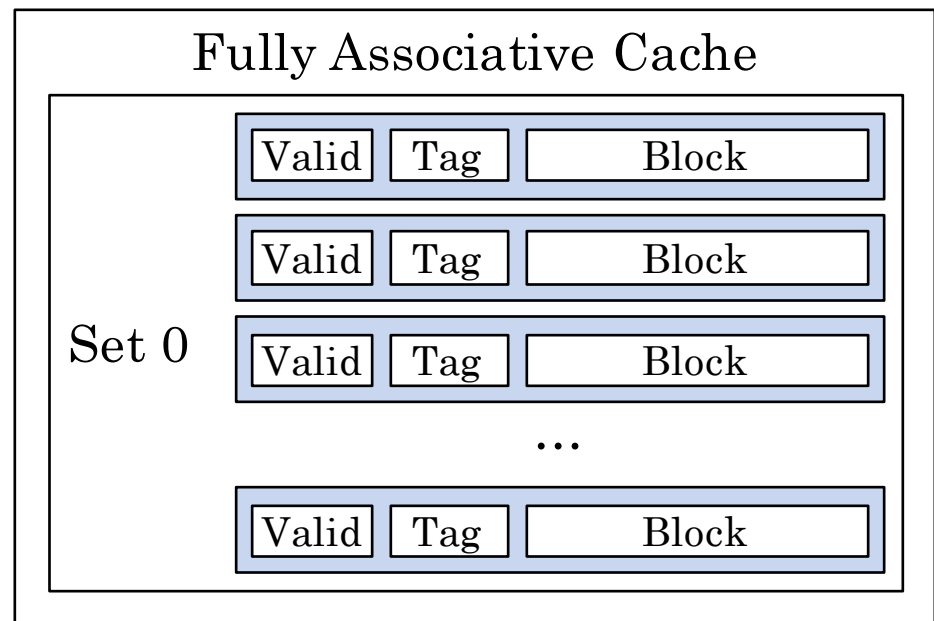


LINE REPLACEMENT POLICIES (2)

- Least Recently Used (LRU) policy
 - Evict cache line that was accessed furthest in past
- Least Frequently Used (LFU) policy
 - Evict cache line that was accessed the least frequently, over some time window
- These policies take extra time and hardware to implement
 - Not used as much in caches close to the CPU, where performance is critical
 - Used very often in caches further from the CPU, where cache misses are extremely costly
- For example, disk-block caches benefit *greatly* from more sophisticated replacement policies
 - ...when a cache miss costs 20 million clocks, spend a few thousand clocks to figure out what to keep in the cache...

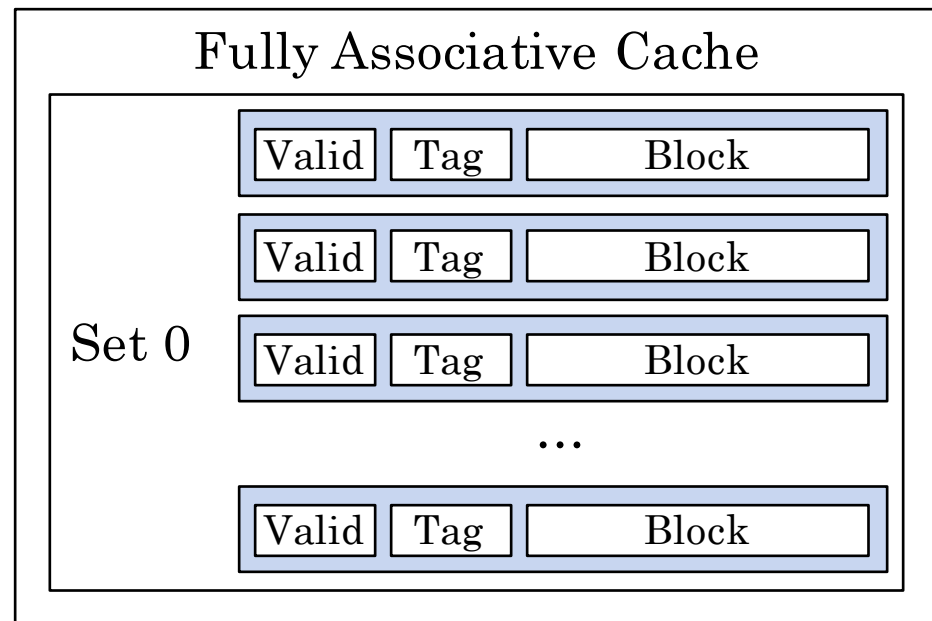
ASSOCIATIVE CACHES

- Where does the “associative” come from in set-associative caches and fully-associated caches?
- Each cache set has E cache lines in it...
 - Need to look up cache line using only the block's tag
 - The cache set is an *associative memory*
- Associative memory:
 - Not accessed with an address, like normal memories!
 - Associative memory stores (key, value) pairs
 - Key is the input to the associative memory
 - Memory returns value



ASSOCIATIVE CACHES (2)

- Associative caches must effectively implement associative memories for their cache sets
 - Keys are a concatenation of the tag, *plus* the valid flag
 - *No reason to look at the cache line if it isn't valid...*
 - Value is the block of data in the cache
- Set-associative caches:
 - Each cache set is an associative memory
 - Number of cache lines in each set is small, so logic is easier to implement
- Fully associative caches:
 - Need to examine *many* cache lines in parallel
 - *Much* more expensive...



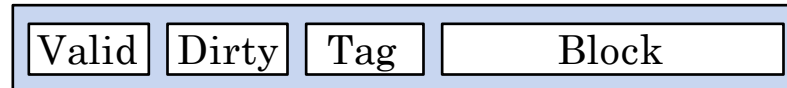
WRITING TO CACHES

- So far, only discussed reading from the cache...
- Most programs also write to memory...
 - Need to handle writes to the cache as well!
 - Ideally, want to minimize performance penalties on writes, just like on reads
- If CPU writes a word that is already in the cache:
 - Called a write hit
 - Several choices of how to handle this situation
- Option 1: employ a *write-through* strategy
 - Every write to cache causes cache to write the entire block back to memory
 - Problem: *every single write* to the cache causes a write to main memory! Can't exploit data locality!

WRITING TO CACHES (2)

- Option 2: use a *write-back* strategy
 - Add a “dirty flag” to each cache line
 - When a cached block is written to, set the dirty flag
 - When a cache line is evicted, write the block back to main memory

Cache line:



- Benefit:
 - Can now exploit locality of writes as well as reads
 - Writes to adjacent values likely to hit cached blocks
 - When dirty blocks are finally evicted, have a much smaller number of writes to main memory
- Drawback: more complex than write-through
 - These days, most caches use a write-back strategy

WRITING TO CACHES (3)

- If address being written to isn't already in the cache:
 - Called a write miss
 - Again we have several options
- Option 1: use a *write-allocate* strategy
 - When a write miss occurs, load the block into cache and then perform the write against the cache
 - Assumes the program will perform subsequent writes
- Option 2: use a *no-write-allocate* strategy
 - When a write miss occurs, just perform the write against main memory
- Write-back caches usually use write-allocate strategy
- Write-through caches typically use no-write-allocate strategy

CACHE PERFORMANCE ANALYSIS

- Cache performance modeling can be extremely complicated...
 - Usually easiest to measure actual system behaviors
- Basic ideas can be captured with only a few simple parameters
- Miss rate: the fraction of memory references that result in cache misses
 - Miss rate = # misses / # references ($0 \leq \text{miss rate} \leq 1$)
- Hit rate: the fraction of memory references that hit the cache
 - Hit rate = $1 - \text{miss rate}$

CACHE PERFORMANCE ANALYSIS (2)

- Hit time: time to deliver a value from the cache to the CPU
 - Includes all necessary steps for delivering the value!
 - Time to select the appropriate cache set
 - Time to find the cache line using the block's tag
 - Time to retrieve the specified word from block data
- Miss penalty: time required to handle cache miss
 - *(includes cache-set identification, checking tags, etc.)*
 - Need to fetch a block from main memory
 - May need to evict another cache line from the cache
 - (Evicted line may be dirty and need written back...)
 - Other associated bookkeeping for storing cache line

CACHE PERFORMANCE ANALYSIS (3)

- Simple example to see benefit of caching:
 - Hit time = 1 clock (typical goal for L1 hits)
 - Miss penalty = 100 clocks (main memory usu. 25-100)
- If all reads were from cache, each read is 1 clock
- Hit rate of 80%
 - $0.8 \times 1 \text{ clock} + 0.2 \times 100 \text{ clocks} = 20.8 \text{ clocks/access}$
- Hit rate of 90%
 - $0.9 \times 1 \text{ clock} + 0.1 \times 100 \text{ clocks} = 10.9 \text{ clocks/access}$
- Hit rate of 95%
 - $0.95 \times 1 \text{ clock} + 0.05 \times 100 \text{ clocks} = 5.95 \text{ clocks/access}$
- Hit rate is very important!
 - For programs with low miss-rate (good data locality), get a large memory at (nearly) the cost of a small one!

CACHE DESIGN TRADE-OFFS

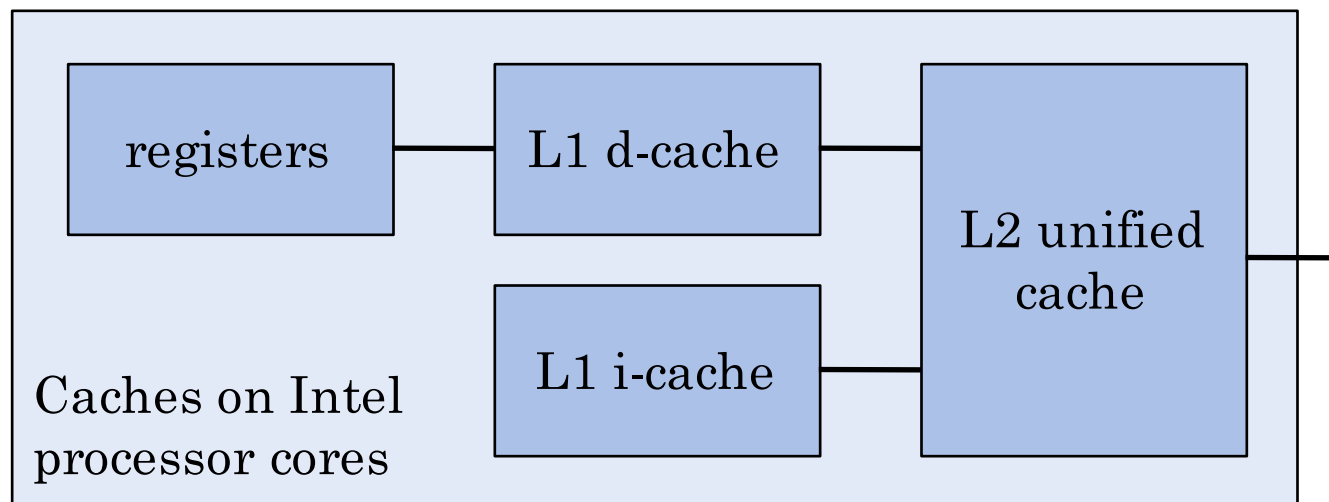
- Previous example shows importance of designing programs to maximize cache usage
 - ...and compilers that know how to optimize well...
- Cache designers also have important trade-offs to consider
- Try to maximize hit rate, without increasing hit time
- Can increase overall size of cache, but this will probably increase the hit time
 - Idea is to increase number of lines, or size of each line
 - Physics: larger memories are slower than smaller ones
- Can increase block size, keeping cache size constant
 - (Don't want to incur penalties due to physics)
 - Miss penalty will increase due to larger transfer times

CACHE DESIGN TRADE-OFFS (2)

- For associative caches, what about increasing number of cache lines per set (E) ?
- As number of cache lines per set increases:
 - Complexity of line-matching goes up, since more lines to compare. Will probably increase hit times.
 - Must choose a line to evict when a cache miss occurs; now there are more choices.
 - Will probably increase miss penalty
 - Complexity of replacement-policy logic also increases!
- Direct mapped caches avoid above penalties with *ultra-simple* mapping of blocks to cache lines...
 - ...but *dramatically* increase likelihood of thrashing due to conflict misses!

INTEL PROCESSOR CACHES

- So far have only discussed caching data...
 - Reading instructions also frequently has good locality
- Processor can manage separate instruction caches (i-caches) and data caches (d-caches)
 - Allows parallel data-access paths within the CPU
 - Also, instruction caches usually only support reads ☺



INTEL CORE 2 CACHES

- Intel Core 2 Duo/Quad cache information (typical):

Cache	Associativity (E)	Block Size (B)	Sets (S)	Cache Size (C)
L1 i-cache	8	64 bytes (16 dwords)	64	32KB
L1 d-cache	8	64 bytes (16 dwords)	64	32KB
Unified L2 cache	8	64 bytes (16 dwords)	2048 - 16384	2 MB – 6MB+

- For Core-2 Quad, L2 is divided into two segments, each of which is shared by two cores
- Small number of cache lines per set, but large number of cache sets
 - Fast to determine if a block is in a cache set, and many cache sets to minimize thrashing issues

INTEL CORE I7 CACHES

- Intel Core i7 cache information:

Cache	Associativity (E)	Block Size (B)	Sets (S)	Cache Size (C)
L1 i-cache	8	64 bytes (16 dwords)	64	32KB
L1 d-cache	8	64 bytes (16 dwords)	64	32KB
Unified L2 cache	8	64 bytes (16 dwords)	512	256KB
Unified L3 cache	16	64 bytes (16 dwords)	8192	8MB

- Each core has its own L2 cache
- All cores share the L3 cache
- With another level of caching, expect to see a smoother degradation in memory performance

MEASURING CACHE BEHAVIOR

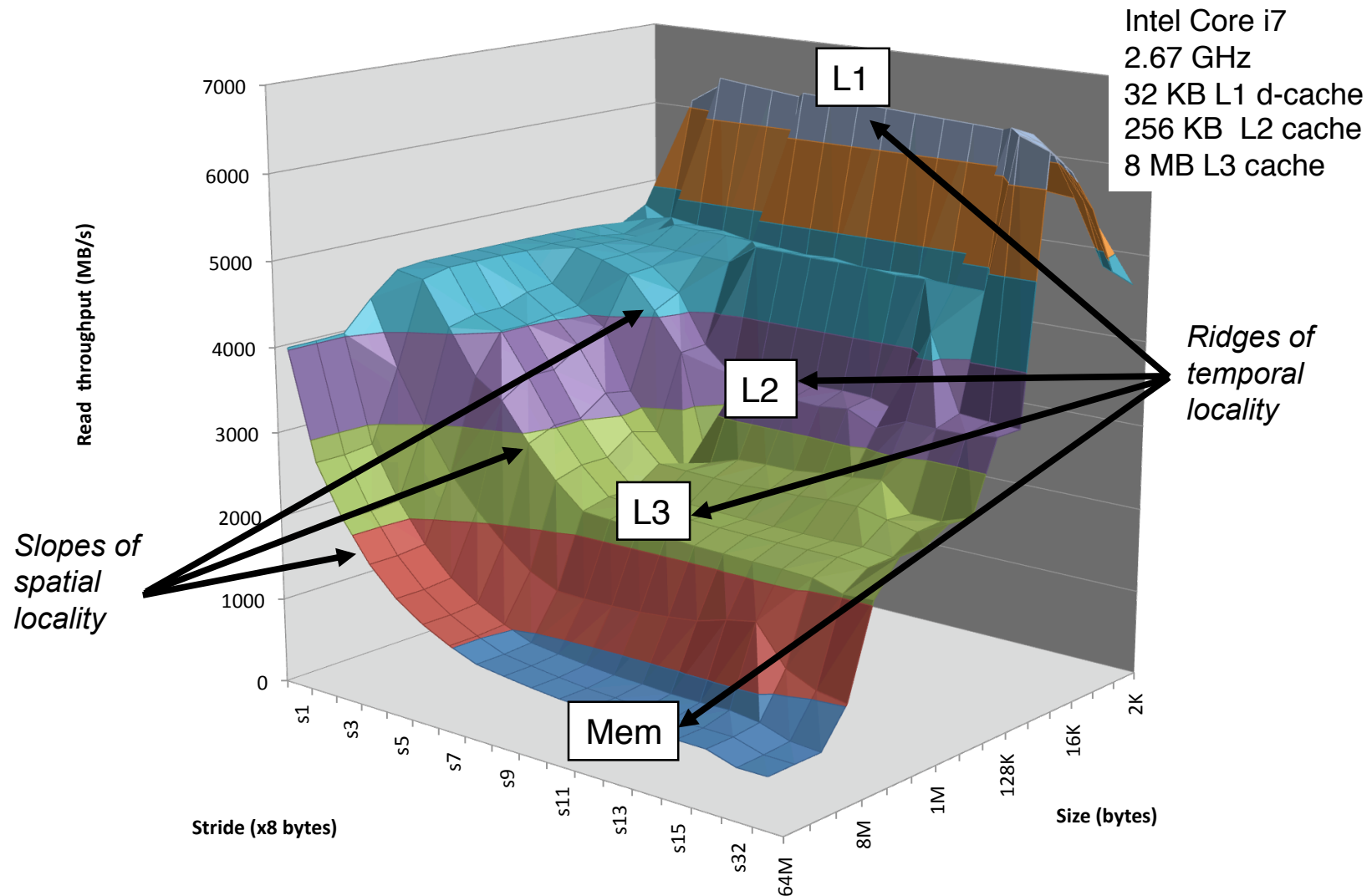
- Can measure the performance of various caches in our computer by constructing a special program:
 - Allocate an array of *size* elements to scan through
 - Access the array with *k*-stride reference patterns
 - Vary *k* from 1 to some large value
 - Scan through the memory for each stride value
- Measure memory access throughput as we vary both of these parameters
- Remember, two different kinds of locality!
- Spatial locality:
 - Program accesses data items that are close to other data items that have been recently accessed
- Temporal locality:
 - Program accesses the same data item multiple times

MEASURING CACHE BEHAVIOR (2)

- Can measure the performance of various caches in our computer by constructing a special program:
 - Allocate an array of *size* elements to scan through
 - Access the array with *k*-stride reference patterns
 - Vary *k* from 1 to some large value
 - Scan through the memory for each stride value
- When our working set and stride are small:
 - Should fit entirely into L1 cache, giving *fast* access!
- When working set doesn't fit completely within L1 cache, but fits in L2 cache, and stride is increased:
 - Should see performance taper off, as ratio of L1 misses to L1 hits increases
 - When *all* accesses miss L1, will see L2 cache performance
- Should see similar behavior from L2-L3, L3-DRAM

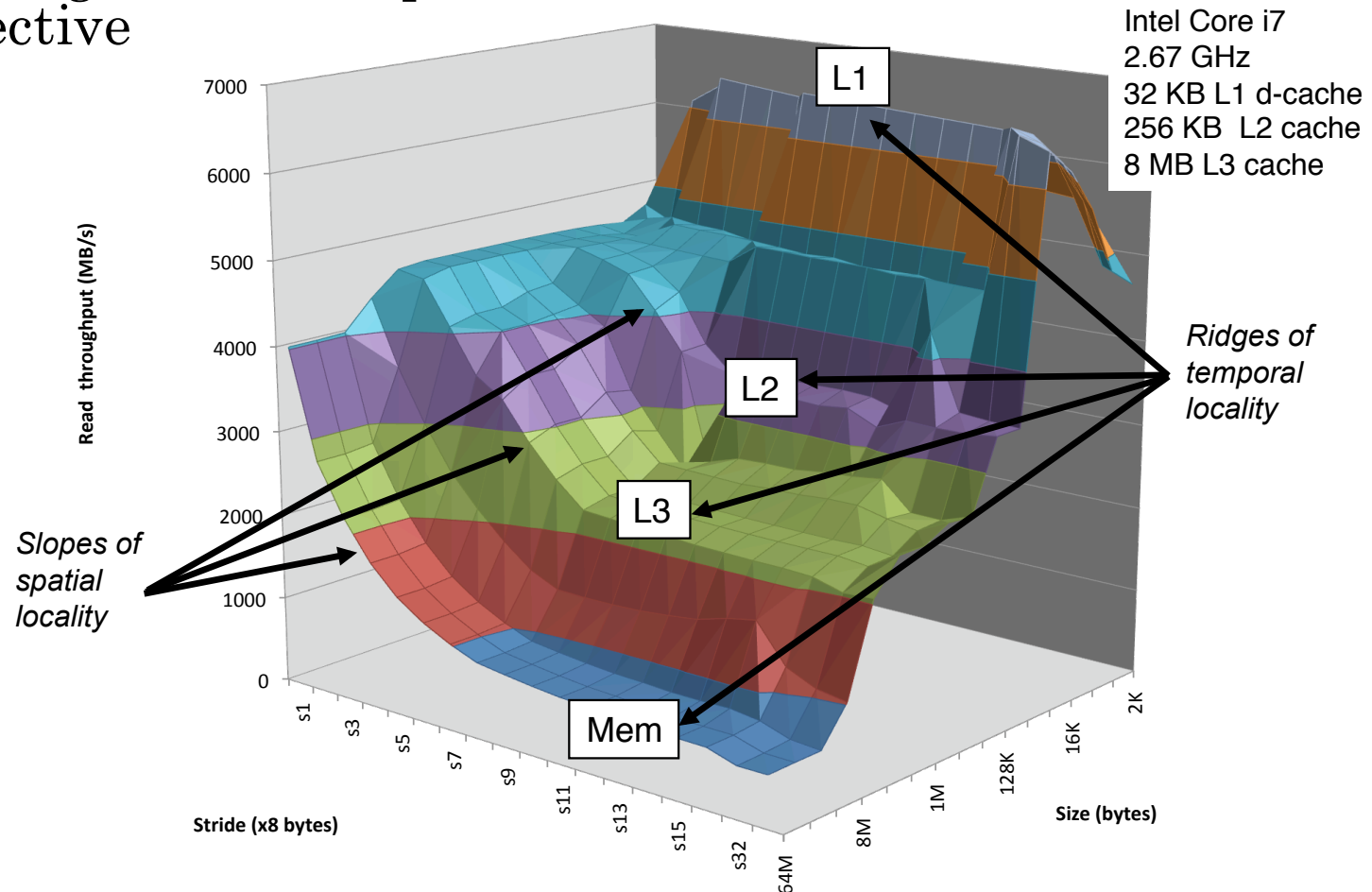
INTEL CORE I7 MEMORY MOUNTAIN

- Produces a 3D surface that shows our caches!



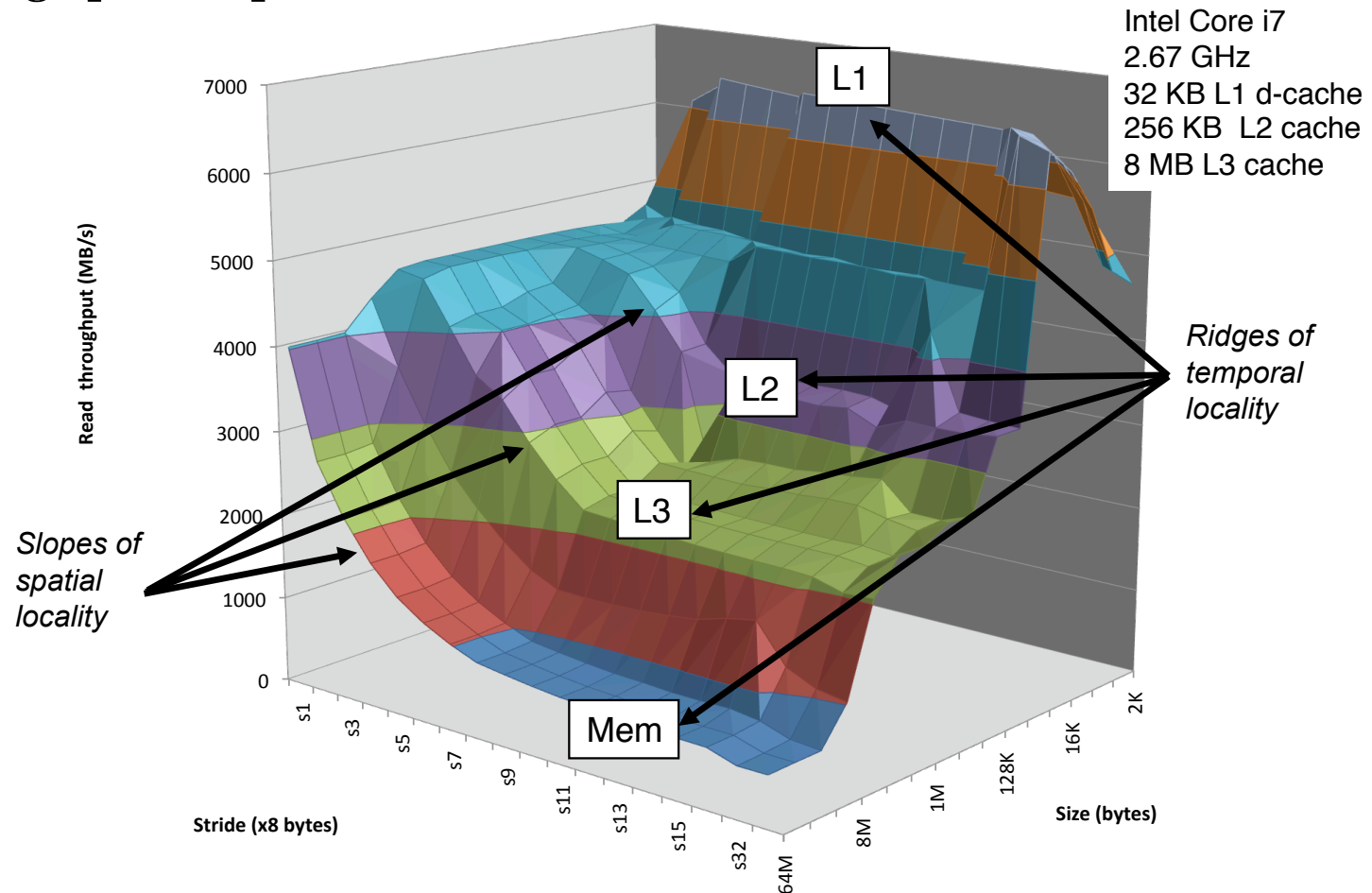
INTEL CORE I7 MEMORY MOUNTAIN (2)

- As total size increases, temporal locality decreases
 - Working set can no longer fit in a given cache...
 - See a significant dropoff as each cache level becomes ineffective



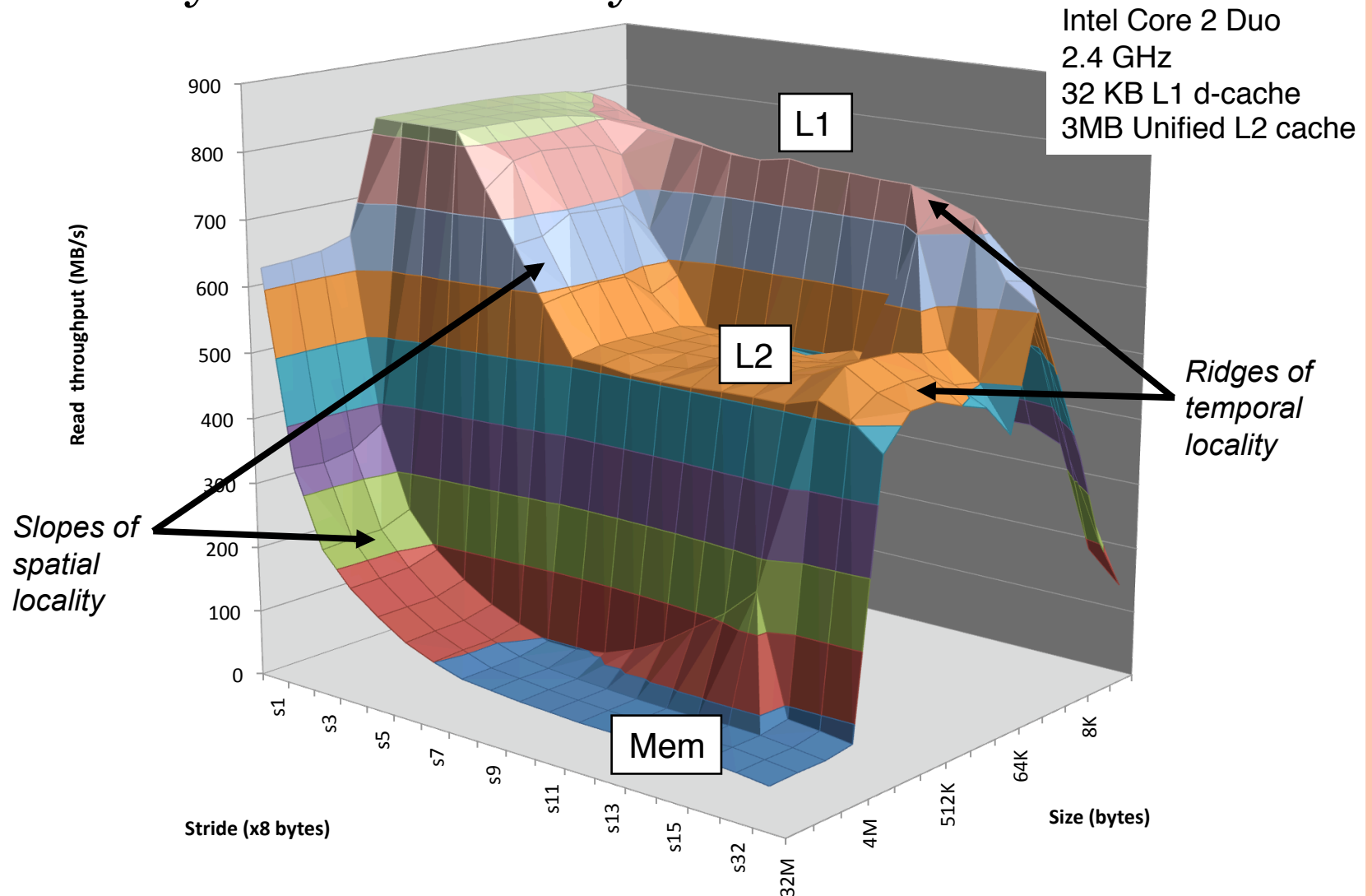
INTEL CORE I7 MEMORY MOUNTAIN (3)

- As stride increases, spatial locality decreases
 - Cache miss rate increases as stride increases
 - Throughput tapers off to a minimum for each level



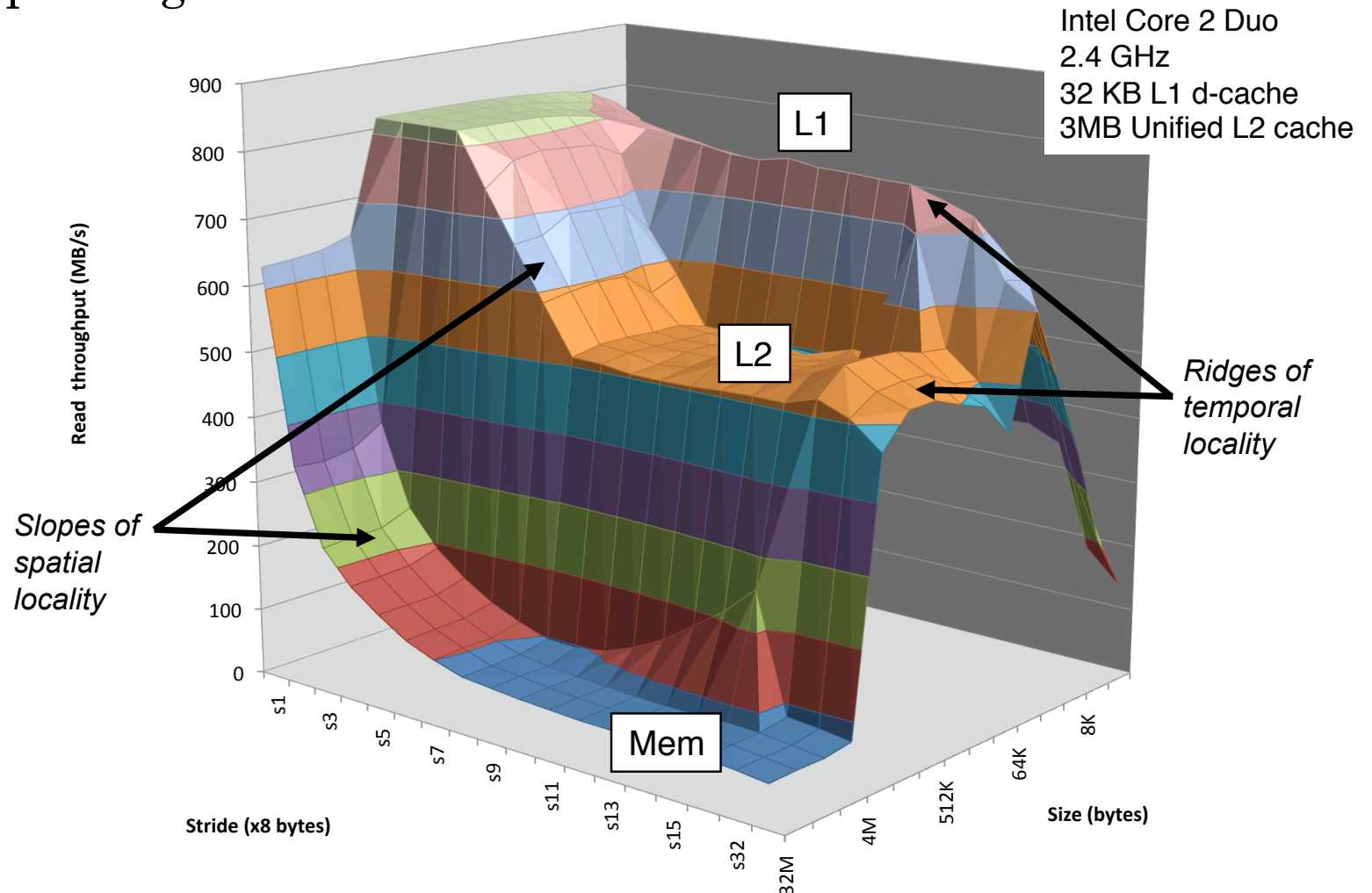
INTEL CORE 2 DUO – MACBOOK

- Different systems have very different statistics



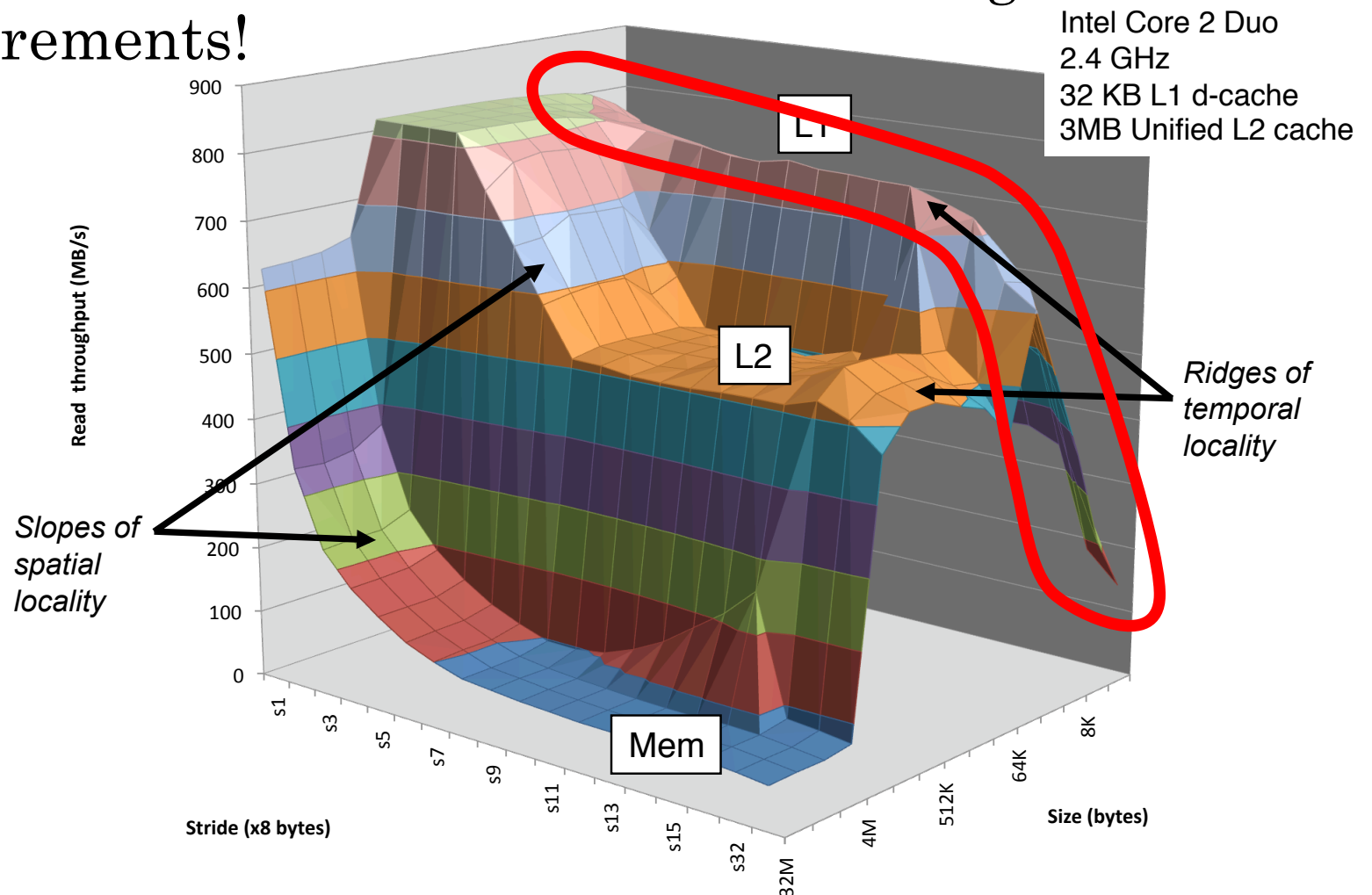
INTEL CORE 2 DUO – MACBOOK (2)

- Core 2 only has two caches before main memory
 - Sharper degradation on Core 2 than on Core i7



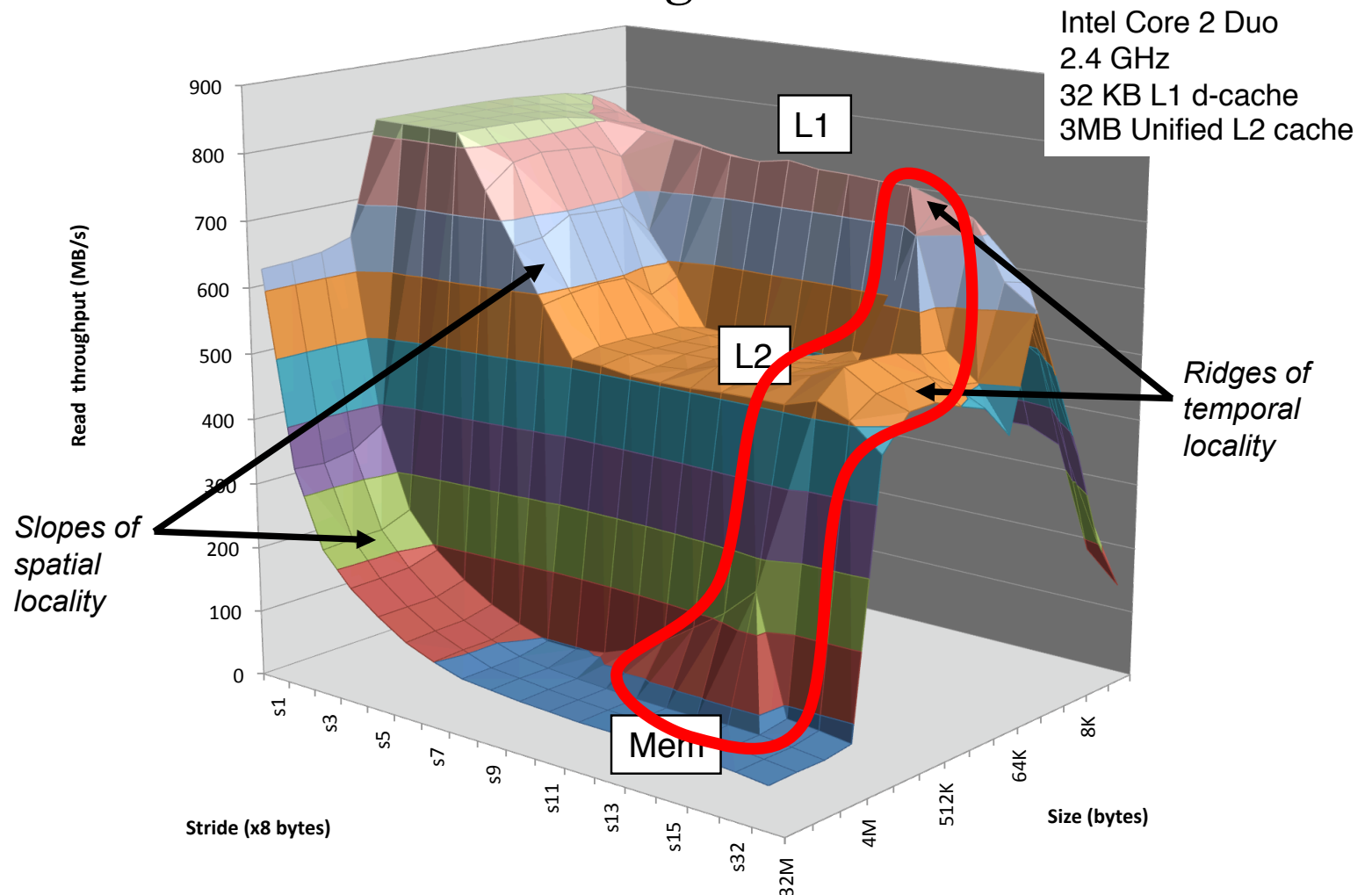
INTEL CORE 2 DUO – MACBOOK (3)

- Dip for very small working sets – at these sizes, cost of function invocation is overwhelming the measurements!



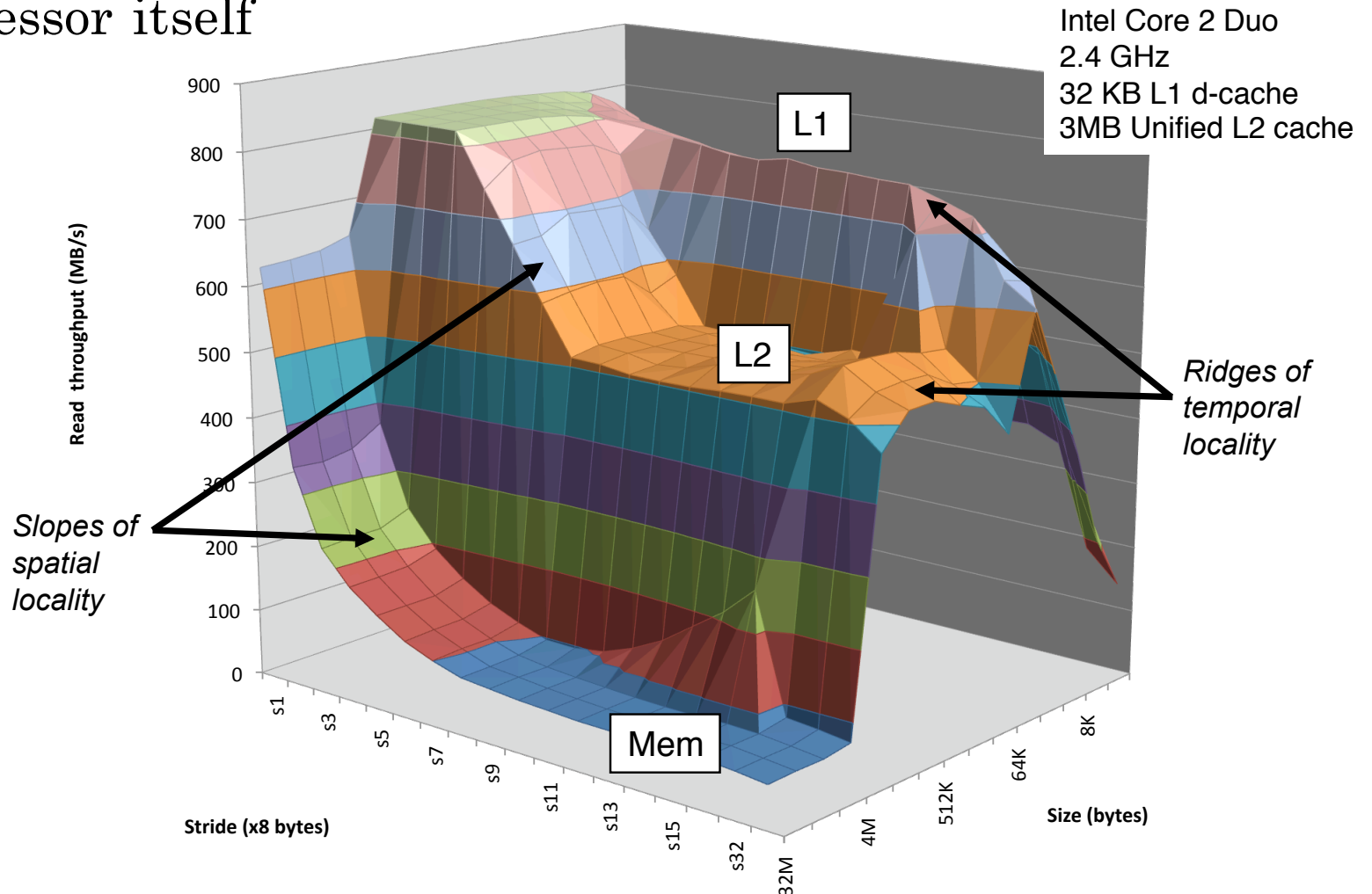
INTEL CORE 2 DUO – MACBOOK (4)

- Small bump where stride allows previously loaded cache lines to be used again...



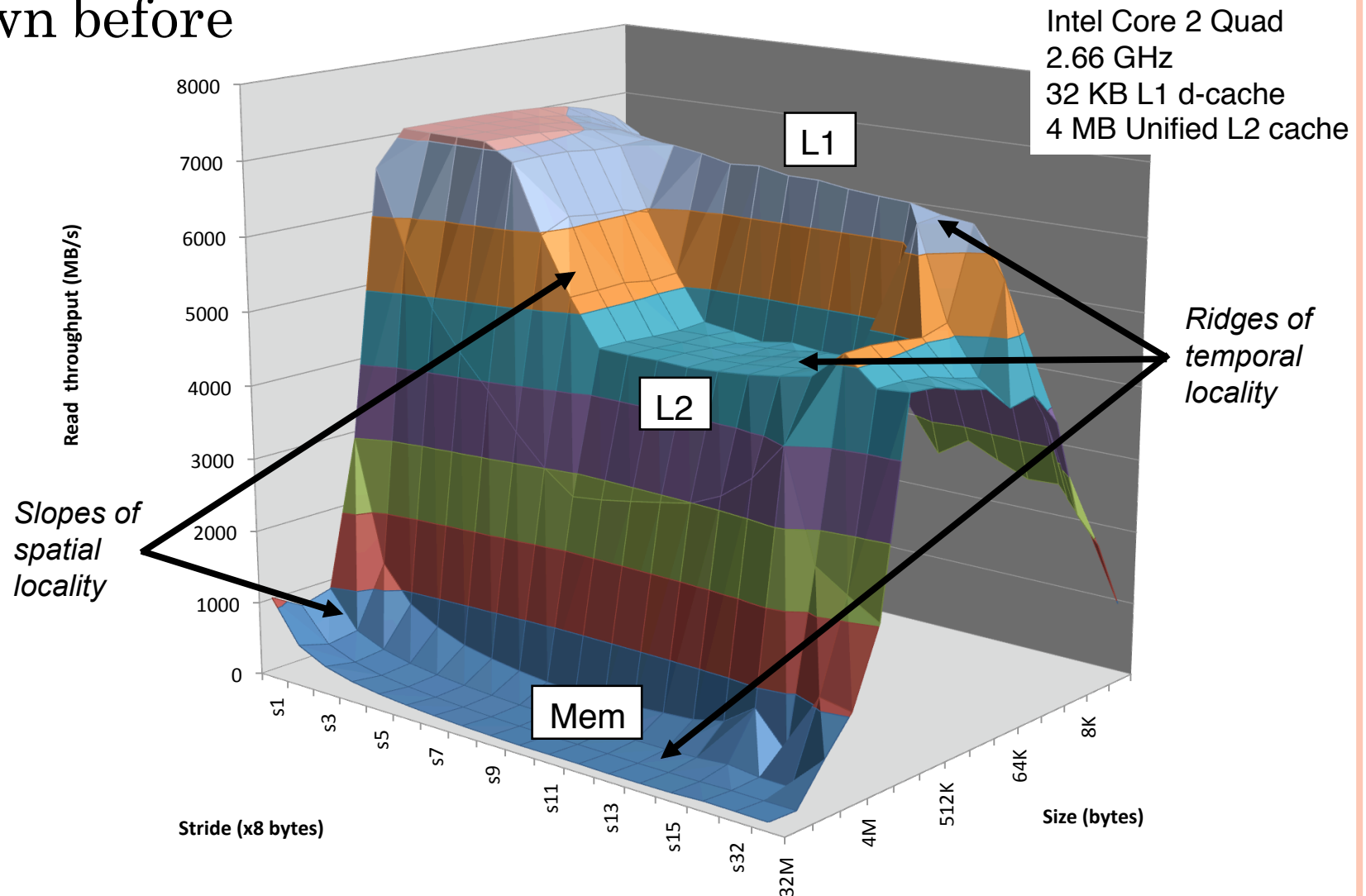
INTEL CORE 2 DUO – MACBOOK (5)

- Peak throughput is *much* lower than before...
 - Affected as much by motherboard chipset as by the processor itself



INTEL CORE 2 QUAD – DESKTOP

- Peak throughput is much more similar to Core i7 shown before



SUMMARY: “CACHE IS KING”

- Caches are essential components in computers
 - Performance gap between CPU and memory is large and increasing...
- Caches give us an opportunity to have large memories that are nearly as fast as small ones
 - As long as we keep our cache hit rates up, program performance will improve dramatically
- As programmers, it is extremely important to be aware of data locality and caching considerations
 - Sometimes, simple changes in a program can produce dramatic changes in memory throughput
 - Other times, must carefully design program to take advantage of processor caches

PROCESSOR PERFORMANCE

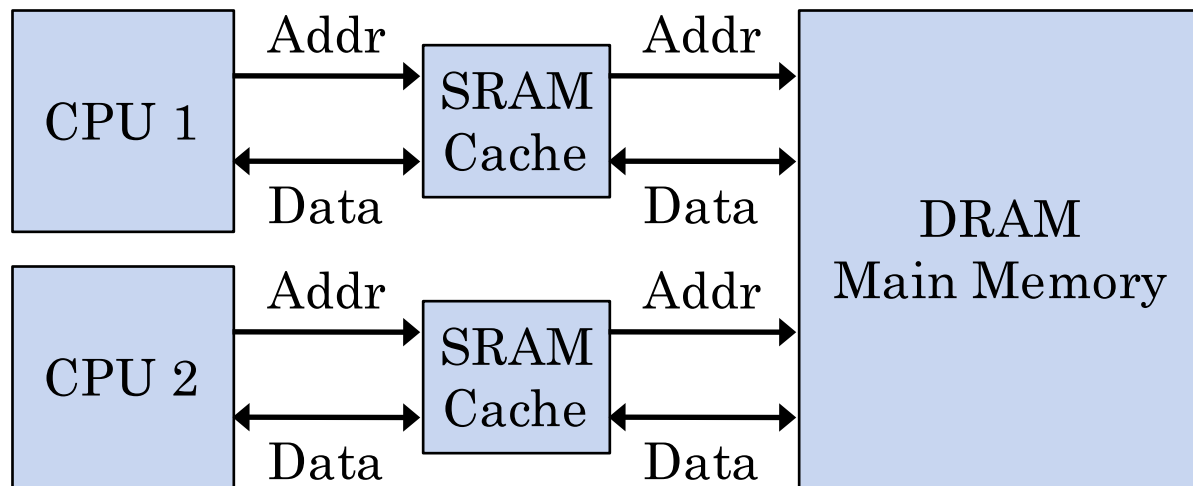
- Increases in CPU clock speeds have tapered off...
 - Issues: clock signal propagation, heat dissipation
- New focus: design CPUs that do more with less
 - Parallelize instruction execution
 - Pipelined instruction execution
 - Superscalar architectures – multiple ALUs, etc.
 - Out-of-order instruction execution
 - Disabling unused CPU components to reduce power
 - (See CS:APP chapters 4, 5 – very good discussion!)
- Although CPU speeds have remained in 1-3GHz ballpark, CPUs are still becoming more powerful

PROCESSOR PERFORMANCE (2)

- Another approach: multi-core processors
 - Put multiple independent CPUs onto a single chip
 - Intel Core 2/Core i7, IBM/Sony/Toshiba Cell, NVIDIA GeForce 9/GeForce 200 series GPUs, etc.
- Still suffers from the same processor-memory performance gap as before...
 - *(And, the solution is still caching...)*
- Explore the challenges and opportunities that multi-core introduces into cached memory access
 - Very rich topic, so our discussion will be a pretty high-level overview

MULTICORE AND CACHING

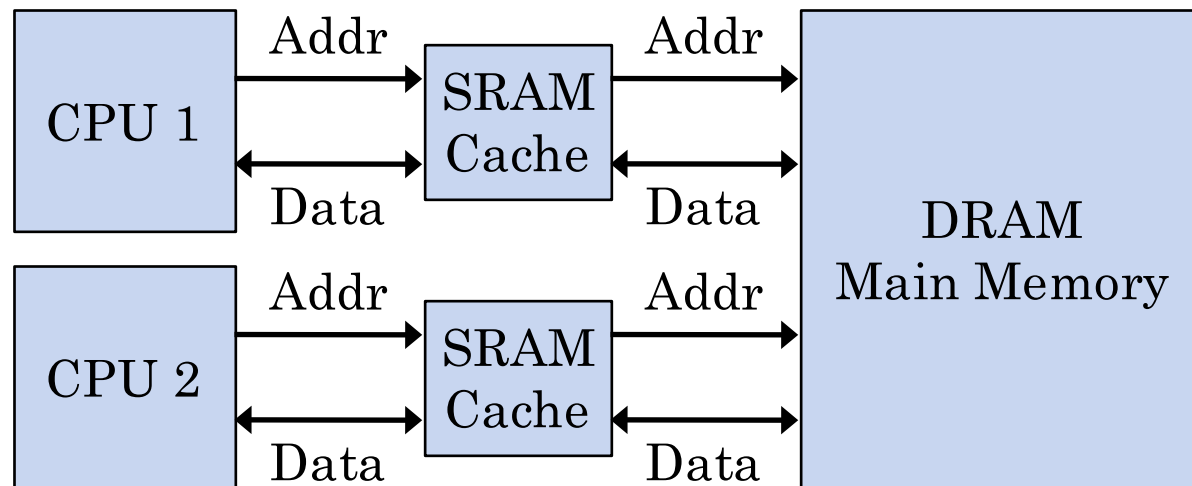
- Go ahead and stick multiple processors on a chip
- What happens to caching in these systems?
- For example:



- Any problems with this approach?
 - More accurately, “What *new* problems do we have?”

INDEPENDENT CACHES

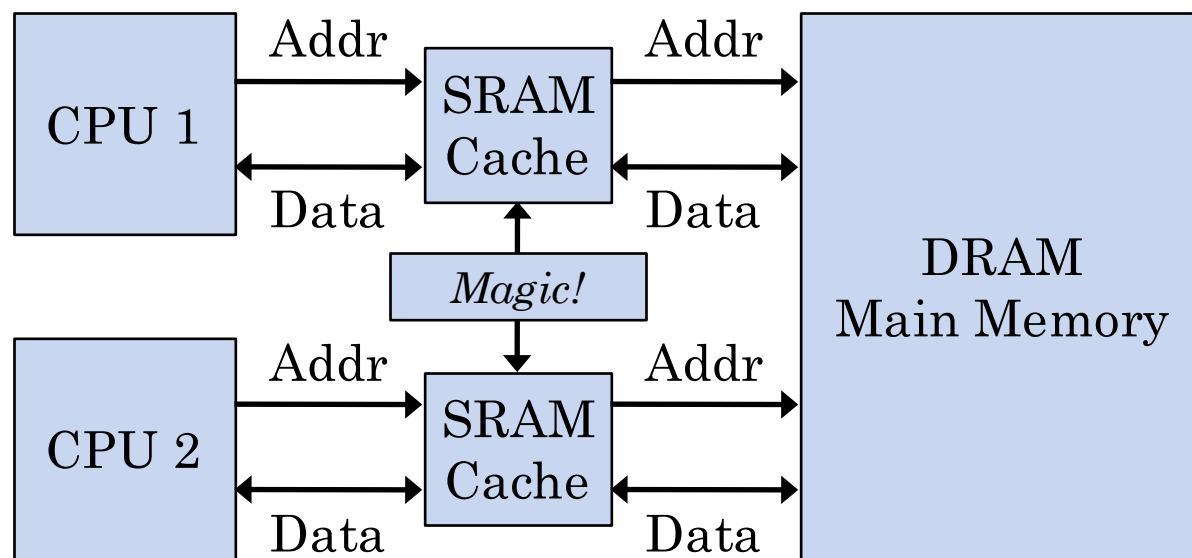
- Two independent caches of a common resource:



- CPU 1 reads data in block 23, starts using it...
 - Block 23 is loaded into CPU1's cache
- CPU 2 issues a write on block 23...
 - Block 23 also loaded into CPU2's cache
 - CPU 1 cache no longer consistent with memory state

COORDINATED CACHES

- Must coordinate state between separate caches



- CPU 1 reads data in block 23, starts using it...
 - Block 23 is loaded into CPU1's cache
- CPU 2 issues a write on block 23...
 - Block 23 also loaded into CPU2's cache
 - *Somehow*, tell CPU 1's cache that block 23 has changed...

CACHE COHERENCE

- Cache coherence constrains the behavior of reads and writes to caches of a shared resource
- Need to define how coherent memory should behave:
- Processor P reads a location X, then writes to X
 - No other processors modify X between read and write
 - Subsequent reads of X must return the value that P wrote
- Processor P1 writes to a location X
 - Subsequently, another processor P2 reads X
 - P2 must read the value that P1 wrote
 - Specifically, P2 must not read the old value of X
- Processor P1 writes value A to a location X
 - Then, processor P2 writes value B to the same location X
 - Reads of X must see A and then B, but never B and then A