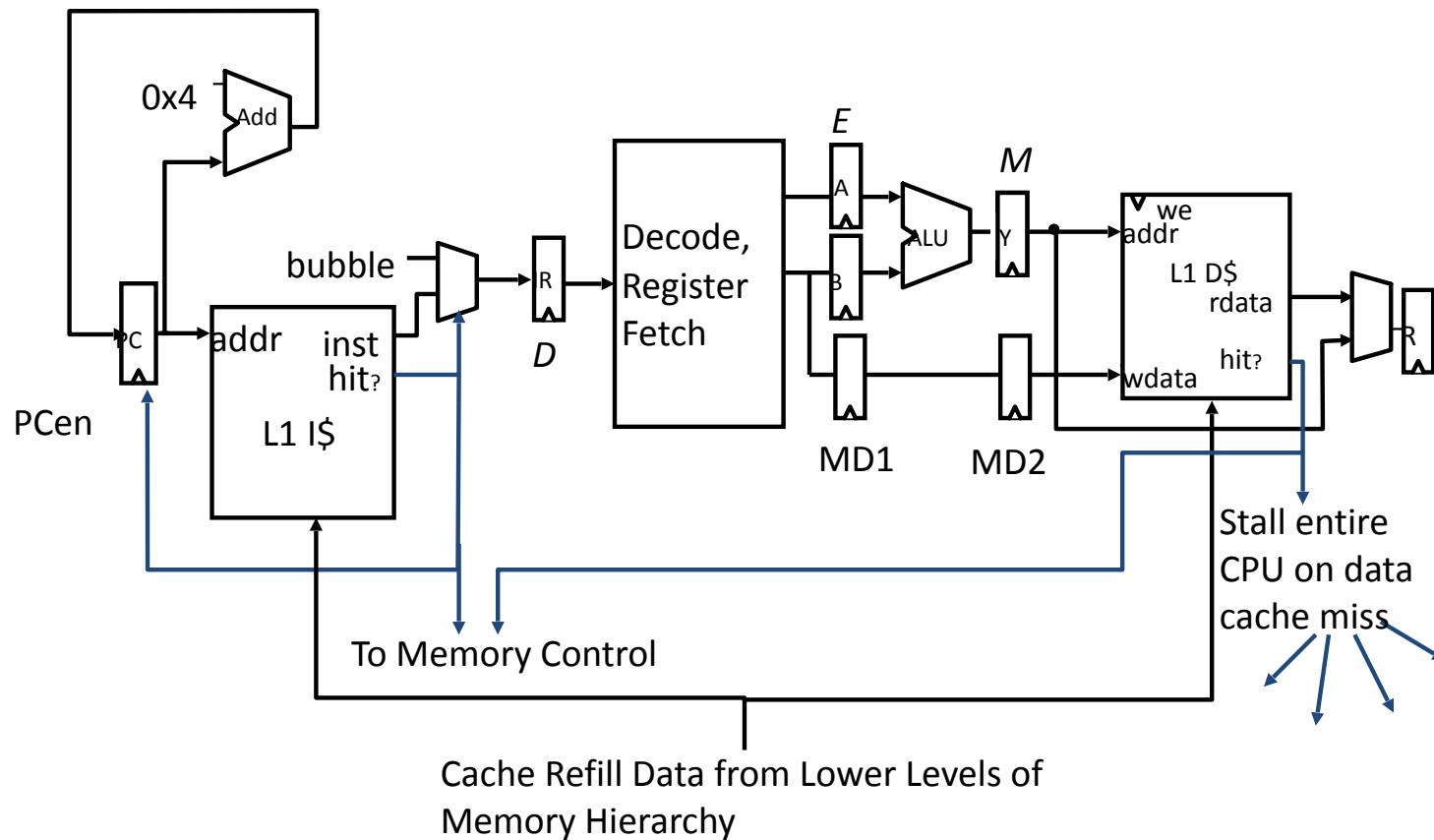


Mo Money, No Problems: Caches #2...

Reminder: Cache Terms...

- **Cache:**
 - A small and fast memory used to increase the performance of accessing a big and slow memory
 - Uses ***temporal locality***: The tendency to reuse data in the same space over time
 - Uses ***spacial locality***: The tendency to use data at addresses near
- Cache ***hit***: The address being fetched is in the cache
- Cache ***miss***: The address being fetched is not in the cache
- ***Valid bit***: Is a particular entry valid
- Cache ***flush***: Invalidate all entries

CPU-Cache Interaction (5-stage pipeline)



More Terms

- **Cache *level*:**
 - The order in the memory hierarchy: L1\$ is closest to the processor
 - L1 caches may only hold data (Data-cache, D\$) or instructions (Instruction Cache, I\$)
 - Most L2+ caches are "unified", can hold both instructions and data
- **Cache *capacity*:**
 - The total # of bytes in the cache
- **Cache *line* or cache *block*:**
 - A single entry in the cache
- **Cache *block size*:**
 - The number of bytes in each cache line

Even More Terms

- **Number of cache lines:**
 - Cache capacity / block size:
- **Cache associativity:**
 - The number of possible cache lines a given address may exist in.
 - Also the number of comparison operations needed to check for an element in the cache
 - **Direct mapped:** A data element can only be in one possible location
 - **N-way set associative:** A data element can be in one of N possible positions
 - **Fully associative:** A data element can be at any location in the cache.
 - Associativity == # of lines

Even More More terms: Parts of the address

- Address is divided into **|TAG|INDEX|OFFSET|**
- **Offset:**
 - The lowest bits of the memory address which say where data exists within the cache line.
 - It is $\log_2(\text{line size})$
 - So for a cache with 64B blocks it is 6 bits
- **Index:**
 - The portion of the address which says where in the cache an address may be stored
 - Takes $\log_2(\#\text{ of cache lines} / \text{associativity})$ bits
 - So for a 4 way associative cache with 512 lines it is 7 bits
- **Tag:** The portion of the address which must be stored in the cache to check if a location matches
 - # of bits of address - (# of bits for index + # of bits for offset)
 - So with 64b addresses it is 51b...

Even More More **More** terms...

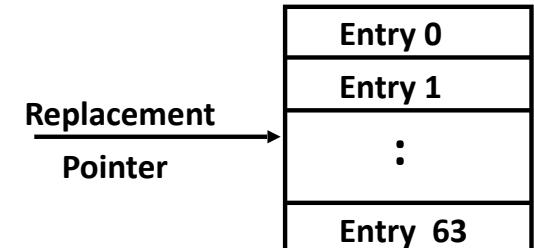
- **Eviction:**
 - The process of removing an entry from the cache
- **Write Back:**
 - A cache which only writes data up the hierarchy when a cache line is evicted
 - Instead set a **dirty bit** on cache entries
 - All i7 caches are **write back**
- **Write Through:**
 - A cache which always writes to memory
- **Write Allocate:**
 - If writing to memory **not in the cache** fetch it first
 - i7 L2 is Write Allocate
- **No Write Allocate:**
 - Just write to memory without a fetch
 - i7 L1 is no write allocate

Even Mostest Terms... Cache Performance

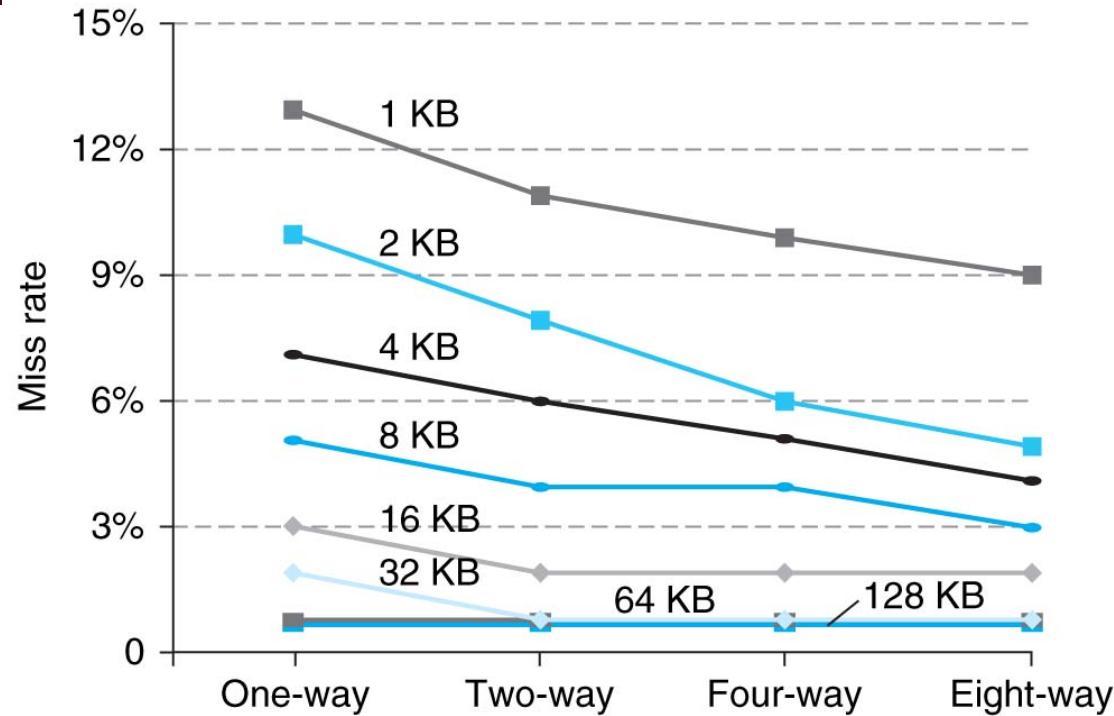
- ***Hit Time:***
 - Amount of time to return data in a given cache: depends on the cache
 - i7 L1 hit time: 4 clock cycles
- ***Miss Penalty:***
 - Amount of ***additional*** time to return an element if its not in the cache: depends on the cache
- ***Miss Rate:***
 - Fraction of a ***particular program's*** memory requests which miss in the cache
- **Average Memory Access Time (**AMAT**):**
 - Hit time + Miss Rate * Miss Penalty

Cache Replacement Policies

- Random Replacement
 - Hardware randomly selects a cache evict
- Least-Recently Used
 - Hardware keeps track of access history
 - Replace the entry that has not been used for the longest time
 - For 2-way set-associative cache, need one bit for LRU replacement
- Example of a Simple “Pseudo” LRU Implementation
 - Assume 64 Fully Associative entries
 - Hardware replacement pointer points to one cache entry
 - Whenever access is made to the entry the pointer points to:
 - Move the pointer to the next entry
 - Otherwise: do not move the pointer
 - (example of “not-most-recently used” replacement policy)
 - What Intel uses on the i7



Benefits of Set-Associative Caches



- Largest gains are in going from direct mapped to 2-way (20%+ reduction in miss rate)

Sources of Cache Misses (3 C's)

- *Compulsory* (cold start, first reference):
 - 1st access to a block, not a lot you can do about it.
 - If running billions of instructions, compulsory misses are insignificant
- *Capacity*:
 - Cache cannot contain all blocks accessed by the program
 - Misses that would not occur **with infinite cache**
- *Conflict* (collision):
 - Multiple memory locations mapped to same cache set
 - Misses that would not occur with ideal **fully associative cache of the same size**

Prefetching...

- Programmer/Compiler: I know that, later on, I will need this data... 
- Tell the computer to ***prefetch*** the data
 - Can be as an explicit prefetch instruction
 - Or an implicit instruction:
`lw 0 0($t0)`
 - Won't stall the pipeline on a cache miss: The processor control logic recognizes this situation
- Allows you to hide the cost of compulsory misses
 - You still need to fetch the data however

Improving Cache Performance

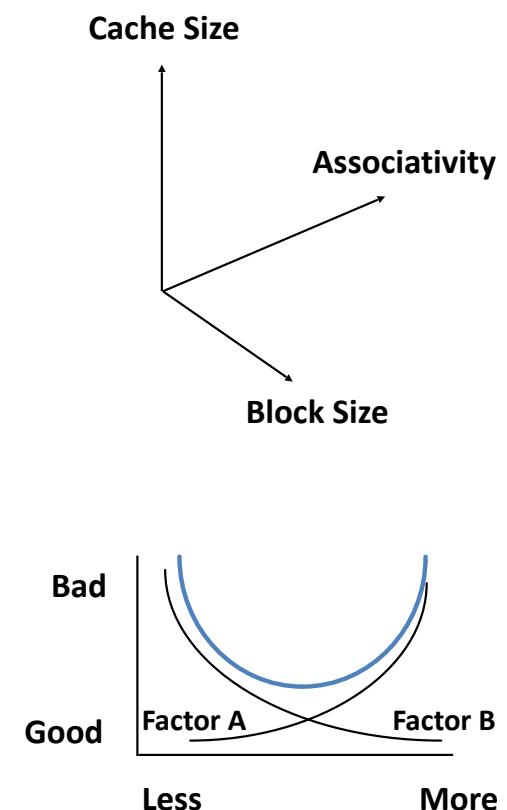
$$\text{AMAT} = \textit{Time for a hit} + \textit{Miss rate} \times \textit{Miss penalty}$$

- Note: miss penalty is the ***additional*** time required for a cache miss
- Reduce the time to hit in the cache
 - E.g., Smaller cache
- Reduce the miss rate
 - E.g., Bigger cache
Longer cache lines (somewhat: improves ability to exploit spatial locality at the cost of reducing the ability to exploit temporal locality)
 - E.g., Better programs!
- Reduce the miss penalty
 - E.g., Use multiple cache levels

Cache Design Space

Computer architects expend considerable effort optimizing organization of cache hierarchy – big impact on performance and power!

- Several interacting dimensions
 - Cache size
 - Block size
 - Associativity
 - Replacement policy
 - Write-through vs. write-back
 - Write allocation
- Optimal choice is a compromise
 - Depends on access characteristics
 - Workload
 - Use (I-cache, D-cache)
 - Depends on technology / cost
 - Simplicity often wins



Primary Cache Parameters

- Block size
 - how many bytes of data in each cache entry?
- Associativity
 - how many ways in each set?
 - Direct-mapped => Associativity = 1
 - Set-associative => $1 < \text{Associativity} < \#\text{Entries}$
 - Fully associative => Associativity = $\#\text{Entries}$
- ***Capacity (bytes) = Total #Entries * Block size***
- ***#Entries = #Sets * Associativity***

Clickers

- For fixed capacity and fixed block size, how does increasing associativity effect AMAT?
 - A: Increases hit time, decreases miss rate
 - B: Decreases hit time, decreases miss rate
 - C: Increases hit time, increases miss rate
 - D: Decreases hit time, increases miss rate

Administrivia

- HW3 Released tonight: Caches only (no floating point)
- Project 3 part 2 due the 23rd at 23:59PM PDT
- Guerilla section tonight 7-9 in 20 Barrows on pipelining/hazards
- MT2 on the 20th
 - Covering up through caches, not covering Floating Point

Increasing Associativity?

- Hit time as associativity increases?
 - Increases, with large step from direct-mapped to ≥ 2 ways, as now need to mux correct way to processor
 - Smaller increases in hit time for further increases in associativity
 - Able to build reasonably efficient wide muxes
- Miss rate as associativity increases?
 - Goes down due to reduced conflict misses, but most gain is from 1->2->4-way with limited benefit from higher associativities
- Miss penalty as associativity increases?
 - Mostly unchanged, replacement policy runs in parallel with fetching missing line from memory

Increasing #Entries?

- Hit time as #entries increases?
 - Increases, since reading tags and data from larger memory structures
- Miss rate as #entries increases?
 - Goes down due to reduced capacity and conflict misses
 - Architects rule of thumb: miss rate drops ~2x for every ~4x increase in capacity (only a gross approximation)
- Miss penalty as #entries increases?
 - Unchanged

At some point, increase in hit time for a larger cache may overcome the improvement in hit rate, yielding a decrease in performance

Clickers: Impact of larger blocks on AMAT

- For fixed total cache capacity and associativity, what is effect of larger **blocks** on each component of AMAT:
 - A: Decrease, B: Unchanged, C: Increase, D: $\nabla(\cup)\cup$
 - Hit Time?
 - Miss Penalty?
 - Miss Rate?

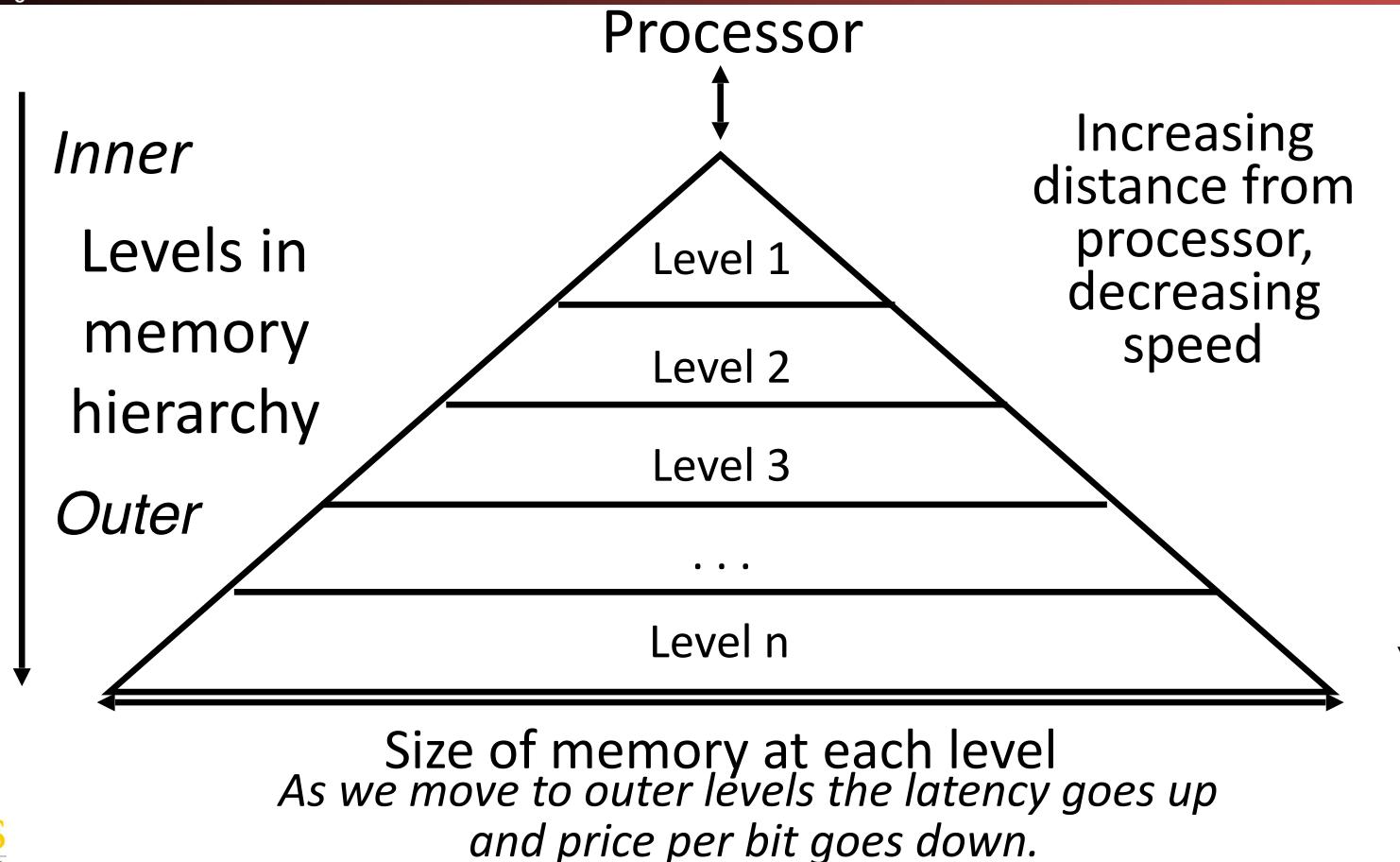
Increasing Block Size?

- Hit time as block size increases?
 - Hit time unchanged, but might be slight hit-time reduction as number of tags is reduced, so faster to access memory holding tags
- Miss rate as block size increases?
 - Goes down at first due to spatial locality, then increases due to increased conflict misses due to fewer blocks in cache
- Miss penalty as block size increases?
 - Rises with longer block size, but with fixed constant initial latency that is amortized over whole block

How to Reduce Miss Penalty?

- Could there be locality on misses from a cache?
- Use multiple cache levels!
- With Moore's Law, more room on die for bigger L1 caches and for second-level (L2) cache
- And now big L3 caches!
- IBM mainframes have ~1GB L4 cache off-chip.

Review: Memory Hierarchy



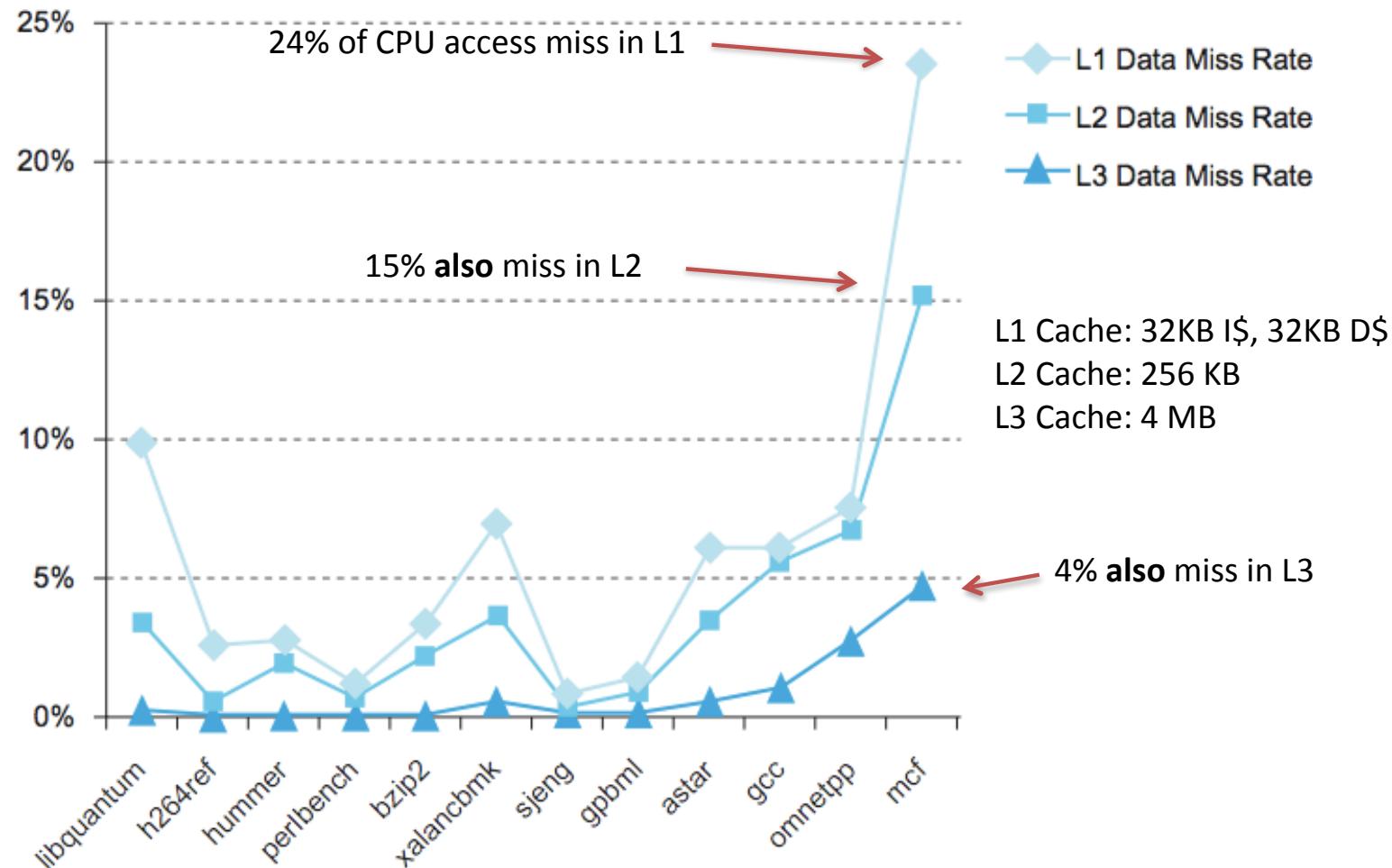


FIGURE 5.47 The L1, L2, and L3 data cache miss rates for the Intel Core i7 920 running the full integer SPECCPU2006 benchmarks.

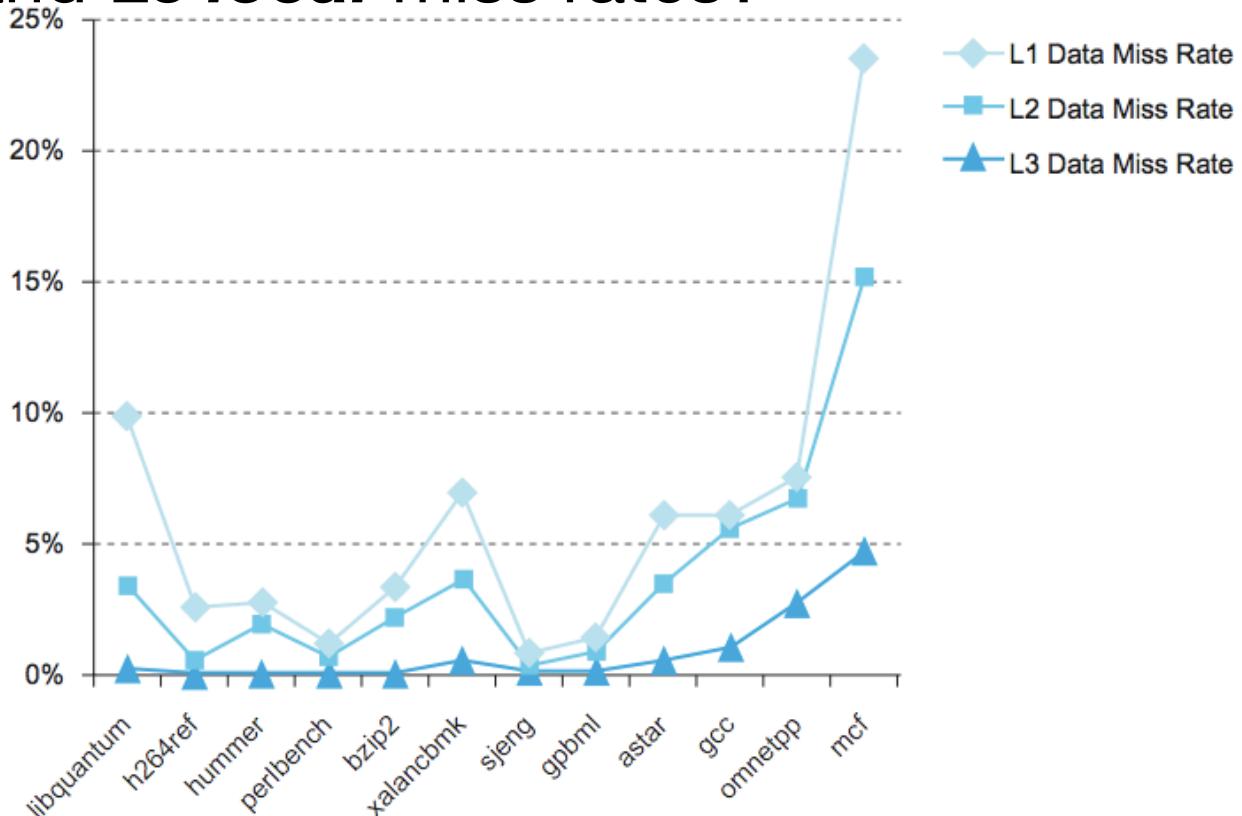
Local vs. Global Miss Rates

- ***Global*** miss rate – the fraction of references that miss some level of a multilevel cache
 - misses in this cache divided by the total number of memory accesses generated by the CPU
- ***Local*** miss rate – the fraction of references to one level of a cache that miss
 - Local Miss rate L2\$ = L2\$ Misses / L1\$ Misses
= L2\$ Misses / total_L2_accesses
 - L2\$ local miss rate \gg than the global miss rate

Clickers/Peer Instruction: Graph is for *global* miss rate

- Overall, what are L2 and L3 *local* miss rates?

- A: L2 > 50%, L3 > 50%
B: L2 ~ 50%, L3 < 50%
C: L2 ~ 50%, L3 ~ 50%
D: L2 < 50%, L3 < 50%
E: L2 > 50%, L3 ~50%



Local vs. Global Miss Rates

- Local miss rate – the fraction of references to one level of a cache that miss
 - Local Miss rate $L2\$ = L2\$ \text{ Misses} / L1\$ \text{ Misses}$
- Global miss rate – the fraction of references that miss in all levels of a multilevel cache
 - $L2\$$ local miss rate \gg than the global miss rate
 - Global Miss rate $= L2\$ \text{ Misses} / \text{Total Accesses}$
 $= (L2\$ \text{ Misses} / L1\$ \text{ Misses}) \times (L1\$ \text{ Misses} / \text{Total Accesses})$
 $= \text{Local Miss rate } L2\$ \times \text{Local Miss rate } L1\$$
- AMAT = Time for a hit + Miss rate \times Miss penalty
 - For 2-level cache system:
 $\text{AMAT} = \text{Time for a } L1\$ \text{ hit} + \text{Miss rate } L1\$ \times$
 $(\text{Time for a } L2\$ \text{ hit} + (\text{local}) \text{ Miss rate } L2\$ \times L2\$ \text{ Miss penalty})$

Real World Caches: Nehalem and Barcelona

Characteristic	Intel Nehalem	AMD Opteron X4 (Barcelona)
L1 cache organization	Split instruction and data caches	Split instruction and data caches
L1 cache size	32 KB each for instructions/data per core	64 KB each for instructions/data per core
L1 block size	64 bytes	64 bytes
L1 write policy	Write-back, Write-allocate	Write-back, Write-allocate
L1 hit time (load-use)	Not Available	3 clock cycles
L2 cache organization	Unified (instruction and data) per core	Unified (instruction and data) per core
L2 cache size	256 KB (0.25 MB)	512 KB (0.5 MB)
L2 block size	64 bytes	64 bytes
L2 write policy	Write-back, Write-allocate	Write-back, Write-allocate
L2 hit time	Not Available	9 clock cycles
L3 cache organization	Unified (instruction and data)	Unified (instruction and data)
L3 cache size	8192 KB (8 MB), shared	2048 KB (2 MB), shared
L3 block size	64 bytes	64 bytes
L3 write policy	Write-back, Write-allocate	Write-back, Write-allocate
L3 hit time	Not Available	38 (?)clock cycles

A Modern x86's Organization: Intel Xeon E7 v3 (Haswell EX)

- A 18 core processor!
 - And you can get more cores today... The v4 comes in a 24 core version!?!?
- Each core can run two separate threads
 - Two **separate** program counters
 - Very pipelined: 14-19 stages (depending on actual instruction)
 - Very superscalar: Issuing up to 7 μ OP per cycle between the two threads
 - Very out-of-order: 168 actual registers, 192 instruction window for reordering
- Addressing:
 - 64b addressing, 64B block size for all caches

The Caches

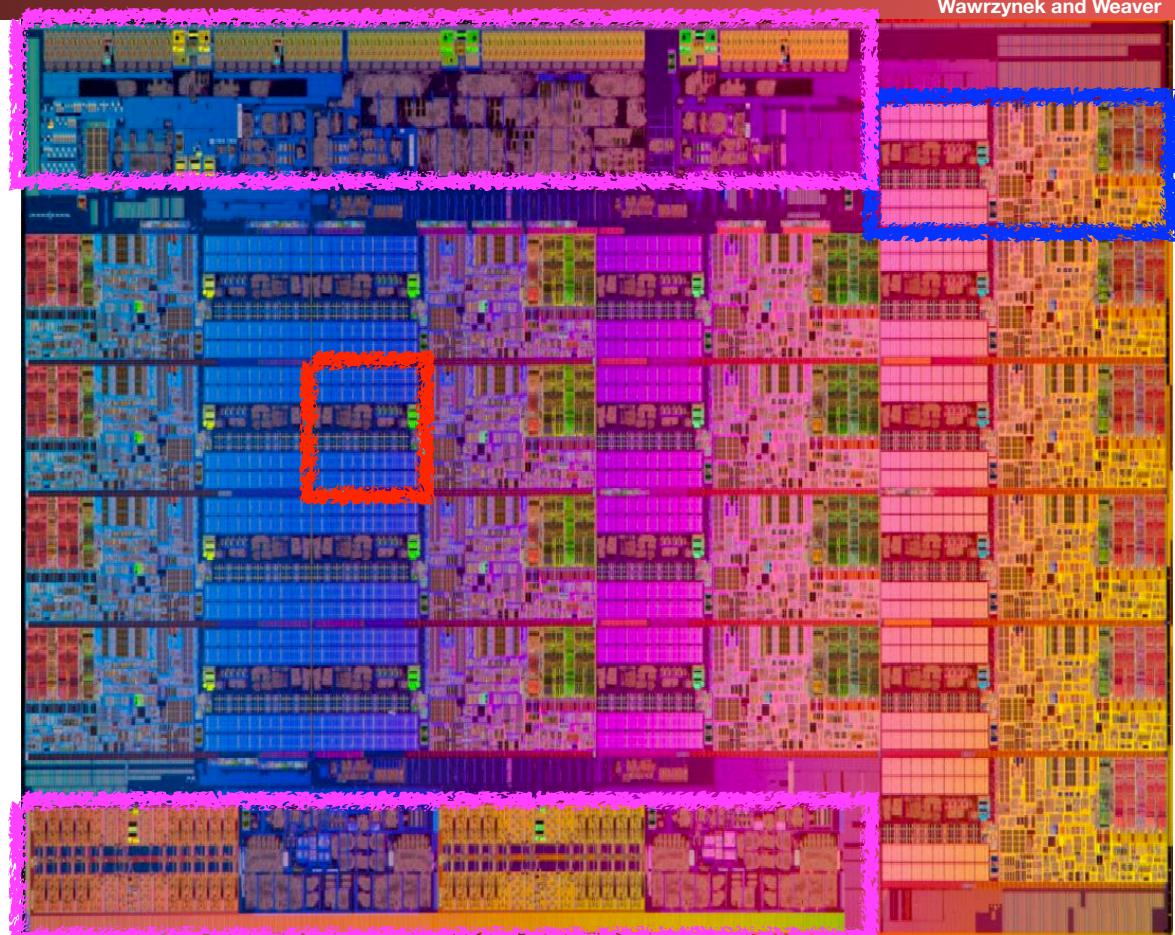
- Each core, 32kB 4-way associative L1 instruction cache, 64B block size
 - 4 cycles latency, pipelined
 - So the ***first 4 stages*** of the instruction pipeline!
 - 512 cache lines
 - Offset: $\lg(64) = 6b$, Index: $\lg((512/4)) = 7b$, Tag: 64-13 = 51b
- Each core, 32 kB L1, 8-way associative write-back data cache
 - 4 cycles latency, pipelined, write back but ***no write allocate!***
 - Pseudo-LRU replacement
- Each core, 256kB 8-way associative write-back L2 cache
 - 10 cycles latency, write back with write allocate
 - Pseudo-LRU
- Common cache, 45MB 16-way associative unified L3 cache (2.5MB per core)
 - Each core has its own section of cache, but all cores can read/write all entries
 - ***Almost***-random replacement

Die Photograph: 18 core Intel Haswell EX processor

Computer Science 61C Spring 2018

Wawrynek and Weaver

- Can see the patterns:
- One Core Block
- One L3 Cache Block
- I/O Control



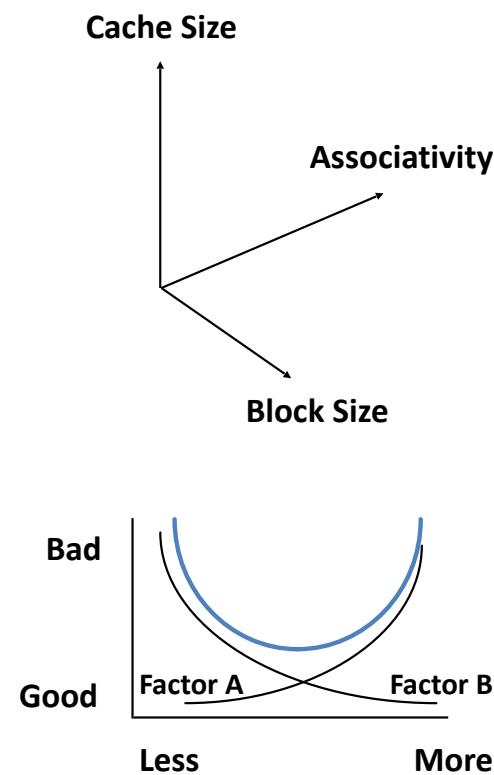
CPI/Miss Rates/DRAM Access

SpecInt2006

Name	CPI	Data Only	Data Only	Instructions and Data
		L1 D cache misses/1000 instr	L2 D cache misses/1000 instr	DRAM accesses/1000 instr
perl	0.75	3.5	1.1	1.3
bzip2	0.85	11.0	5.8	2.5
gcc	1.72	24.3	13.4	14.8
mcf	10.00	106.8	88.0	88.5
go	1.09	4.5	1.4	1.7
hmmer	0.80	4.4	2.5	0.6
sjeng	0.96	1.9	0.6	0.8
libquantum	1.61	33.0	33.1	47.7
h264avc	0.80	8.8	1.6	0.2
omnetpp	2.94	30.9	27.7	29.8
astar	1.79	16.3	9.2	8.2
xalancbmk	2.70	38.0	15.8	11.4
Median	1.35	13.6	7.5	5.4

In Conclusion, Cache Design Space

- Several interacting dimensions
 - Cache size
 - Block size
 - Associativity
 - Replacement policy
 - Write-through vs. write-back
 - Write-allocation
- Optimal choice is a compromise
 - Depends on access characteristics
 - Workload
 - Use (I-cache, D-cache)
 - Depends on technology / cost
 - Simplicity often wins



More Misses...

- We have **Compulsory**, **Capacity**, and **Conflict**...
- We also have **Coherence**
 - Two different processor may share memory...
 - They implement cache coherence so that both processors see the same shared memory
 - When one processor writes to memory, it invalidates the other processor's cache entry for that memory
 - Thus if both processors are working on the same data...
 - This causes Coherence misses
- A related problem can occur if one shared cache is working on two unrelated problems
 - You get additional **capacity** misses:
Can happen in "multithreaded" (aka 'Intel Hyperthreaded') processor cores

Fun Additional Stuff: Nick's Caches

- Note: These won't be on the exam, but they are interesting asides
 - Nick's research has used this material in multiple ways
 - Predictability and caches
 - Why its bad
 - Unpredictable caches: Permutation caches and location-associative permutation caches

Predictability and Caches

- Caches improve performance but...
 - The performance improvement depends on the input
 - E.g. conflict misses depend on input patterns
 - An attacker can take advantage of this
 - Timing of operations can tell something about the input
 - E.g. techniques to extract cryptographic keys!
 - Attacker selected inputs can degrade performance

Why Timing Matters

- Timing enables "side-channel" attacks on cryptography
 - The ability to know some detail of an encryption system based on how long operations take
 - Part of a larger class of side-channel attacks
- It is a fundamentally hard problem to build cryptographic systems that don't have sidechannels
 - Modern processors make this even harder
 - Take 161 and 261 for more on the deep voodoo that side channels can do

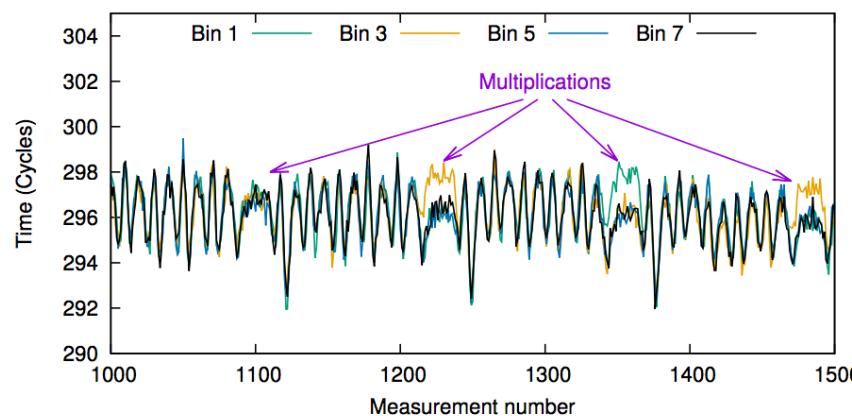


Fig. 5. Measurement trace of OpenSSL RSA decryption—detailed view

From "CacheBleed: A Timing Attack on OpenSSL Constant-Time RSA"

Attacker Selected Input

- Alternatively, if the attacker can select the input...
 - The attacker can select *hard* input:
E.G. Traffic that causes ping-ponging in a direct mapped cache
- Nick's problem:
 - He had to cache IP addresses (32 bit values)
 - This is a network application for security
 - He only wants to store a small amount of information
 - On chip storage expensive (in this case, on an FPGA)
 - He had plenty of room to pipeline
 - Each network packet is independent: No need to forward
 - Misses are expensive
 - Requires processing the packet in software

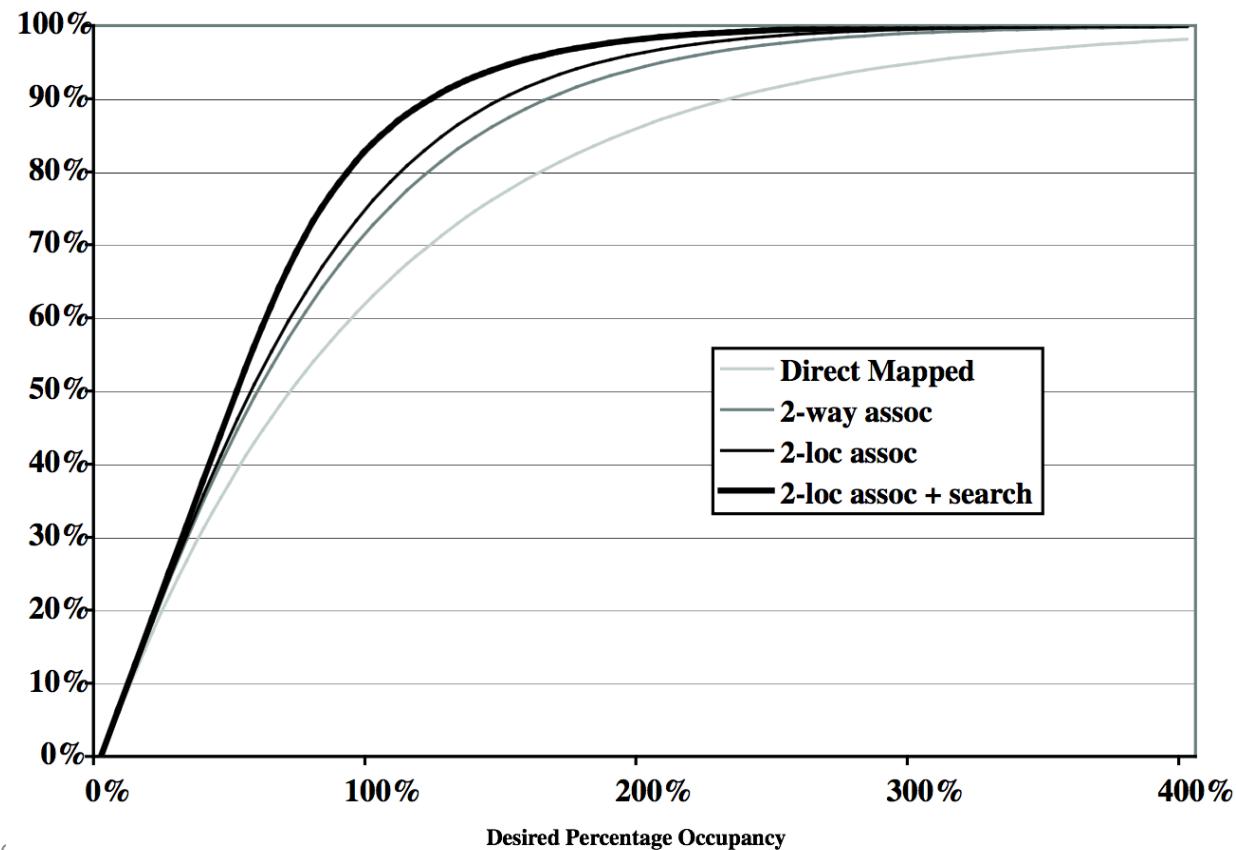
Nick's Cache #1: Permutation Cache

- Traditionally, you would hash the address
 - With a "salt" to randomize things
 - Attacker needs to break the hash function to predict whether two different addresses will match to the same location
 - But this requires storing the whole original IP for your tag
- So instead of a hash, use a 32b ***keyed permutation***
 - Aka a 32b block cypher
 - Now you can use a conventional tag/index approach
 - Requires only storing the tag -> space mattered in this application

Nick's Cache #2: Location Associativity

- The fabric Nick used allowed "dual-ported" memories
 - Like your register file on your processor design: two independent read ports that operate at the same time
- Rather than using set associativity...
 - Instead do two ***different*** permutations (keys) and have one of two possible locations
 - If X, Y, and Z map to the same location with one key...
 - They probably do not on the other key: fewer conflict misses
 - Even better, can probably move a value to further reduce conflict misses

Simulation... *Actual* Occupancy for Caching Random Elements



And Since the Permutation looks "random"...

- It is one of those cases where simulation matches reality...
 - Because I'm caching *random* values
 - Most cache studies are much more program dependent
 - But it does give a way of reducing **conflict** misses without having to increase associativity
 - Just need to calculate two permutations in parallel rather than just one