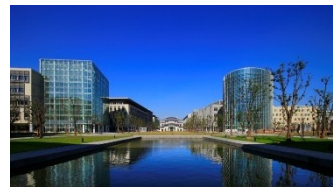




上海科技大学  
ShanghaiTech University



# CS110 Computer Architecture

## Lecture 17:

### *Caches Part III*

Chundong Wang & Siting Liu  
SIST, ShanghaiTech

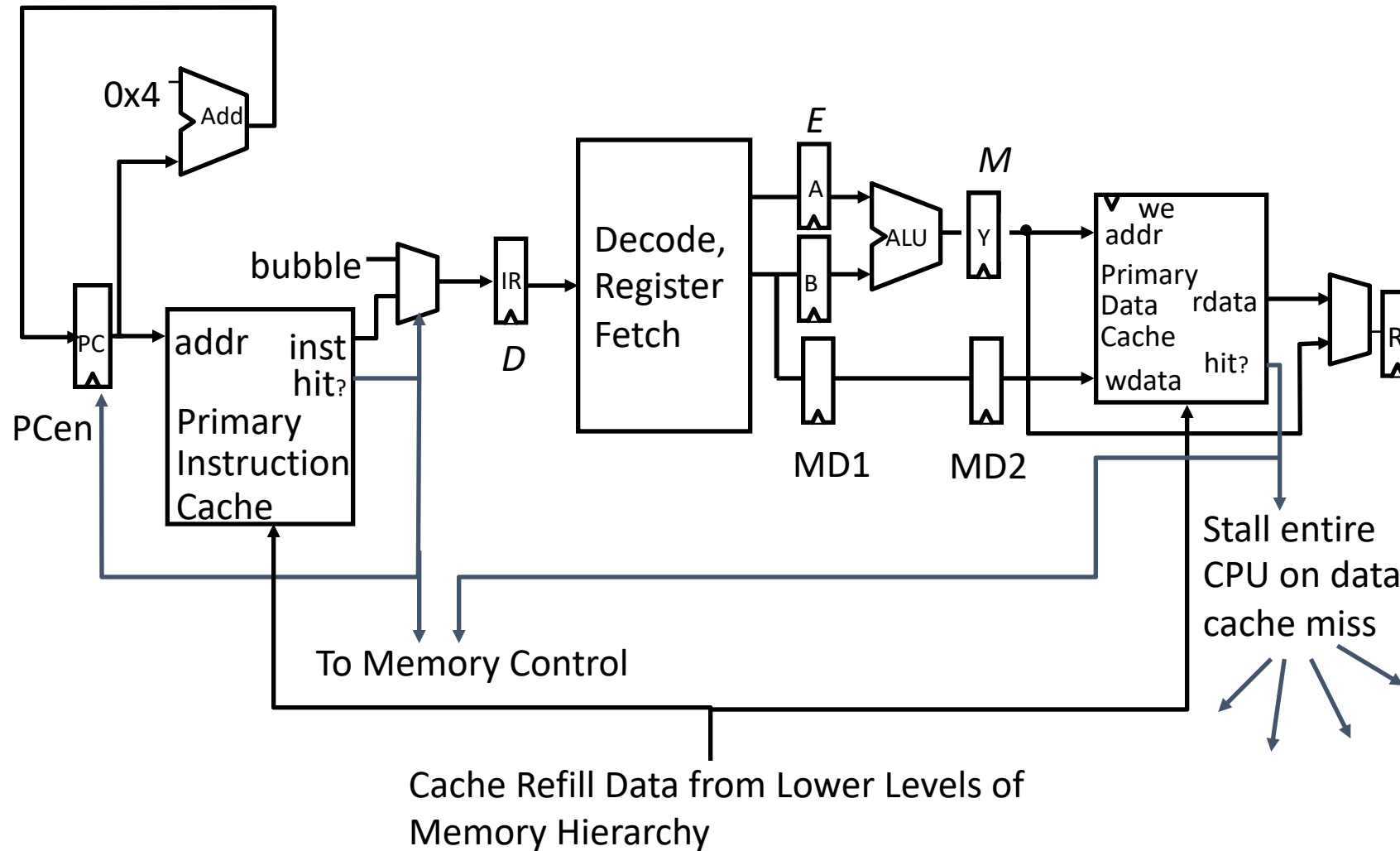


# Cache Terms I

- **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 **spatial locality**: The tendency to use data at nearby addresses
- 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 **line flush**: Invalidate and flush one entry
  - e.g., cflush of x86, but newer clwb may not invalidate the entry
- Cache **flush**: Invalidate all entries
  - e.g., wbinvd of x86



# CPU-Cache Interaction (5-stage pipeline)





# Cache Terms II

- 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





# Cache Terms III: Associativity

- **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 ( $N=1$ )
  - **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
- Total # of cache lines == capacity of cache/line size
- Total # of lines in a set == # ways ==  $N$  == associativity
- Total # of sets == # of cache lines / associativity



# Victim cache

- Conflict misses are a pain, but...
  - Perhaps a little associativity can help without having to be a fully associative cache
- In addition to the main cache...
  - Optionally have a very small (16-64 entry) **fully associative** "victim" cache
- Whenever we evict a cache entry
  - Don't just get rid of it, put it in the victim cache
- Now on cache misses...
  - Check the victim cache first, if it is in the victim cache you can just reload it from there



# Cache Terms IV: 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  $\log_2(\text{line/block size})$
  - 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 64-bit addresses it is **51-bit...**



# Cache Terms V: Writing

- **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





# Replacement Policy

- In an associative cache, which line from a set should be evicted when the set becomes full?
- Random
- Least-Recently Used (LRU)
  - LRU cache state must be updated on every access
  - True implementation only feasible for small sets (2-way)
  - Pseudo-LRU binary tree often used for 4-8 way
- First-In, First-Out (FIFO) a.k.a. Round-Robin
  - Used in highly associative caches
- Not-Most-Recently Used (NMRU)
  - FIFO with exception for most-recently used line or lines
- *This is a second-order effect. Why?*

*Replacement only happens on misses*



# Cache Terms VI: 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



# Understanding Cache Misses: The 3Cs

- **Compulsory** (cold start or process migration, 1<sup>st</sup> reference):
  - First access to block impossible to avoid; small effect for long running programs
  - Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)
- **Capacity**:
  - Cache cannot contain all blocks accessed by the program
  - Solution: increase cache size (may increase access time)
- **Conflict** (*collision*):
  - *Multiple memory locations mapped to the same cache location*
  - *Solution 1: increase cache size*
  - *Solution 2: increase associativity (may increase access time)*



# 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 x0 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



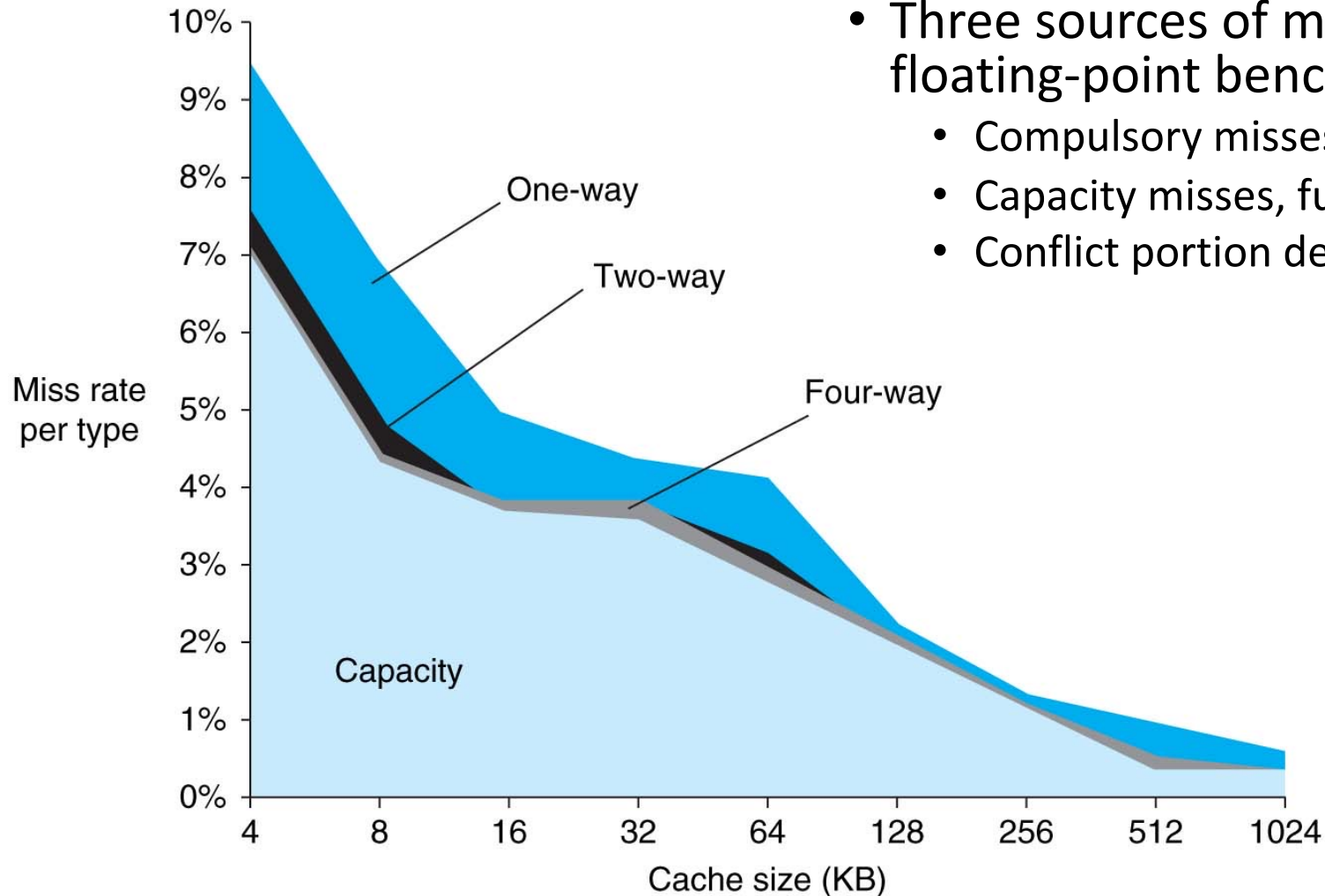


# How to Calculate 3C's using Cache Simulator

1. *Compulsory*: set cache size to infinity and fully associative, and count number of misses
2. *Capacity*: Change cache size from infinity, usually in powers of 2, and count misses for each reduction in size
  - 16 MB, 8 MB, 4 MB, ... 128 KB, 64 KB, 16 KB
3. *Conflict*: Change from fully associative to n-way set associative while counting misses
  - Fully associative, 16-way, 8-way, 4-way, 2-way, 1-way



# 3Cs Analysis



- Three sources of misses (SPEC2000 integer and floating-point benchmarks)
  - Compulsory misses 0.006%; not visible
  - Capacity misses, function of cache size
  - Conflict portion depends on associativity and cache size



# Improving Cache Performance

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

- Note: miss penalty is **additional** time for 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
- Hit and Miss, 3C



# Increasing Associativity?

- Hit time as associativity increases?
  - Increases, with large step from direct-mapped to  $\geq 2$  ways, as now need to mux correct way to processor
  - Smaller increases in hit time for further increases in associativity
- Miss rate as associativity increases?
  - Goes down due to reduced conflict misses, but most gain is from 1- $\rightarrow$ 2- $\rightarrow$ 4-way with limited benefit from higher associativities
- Miss penalty as associativity increases?
  - 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  $\sim 2x$  for every  $\sim 4x$  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**



# 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

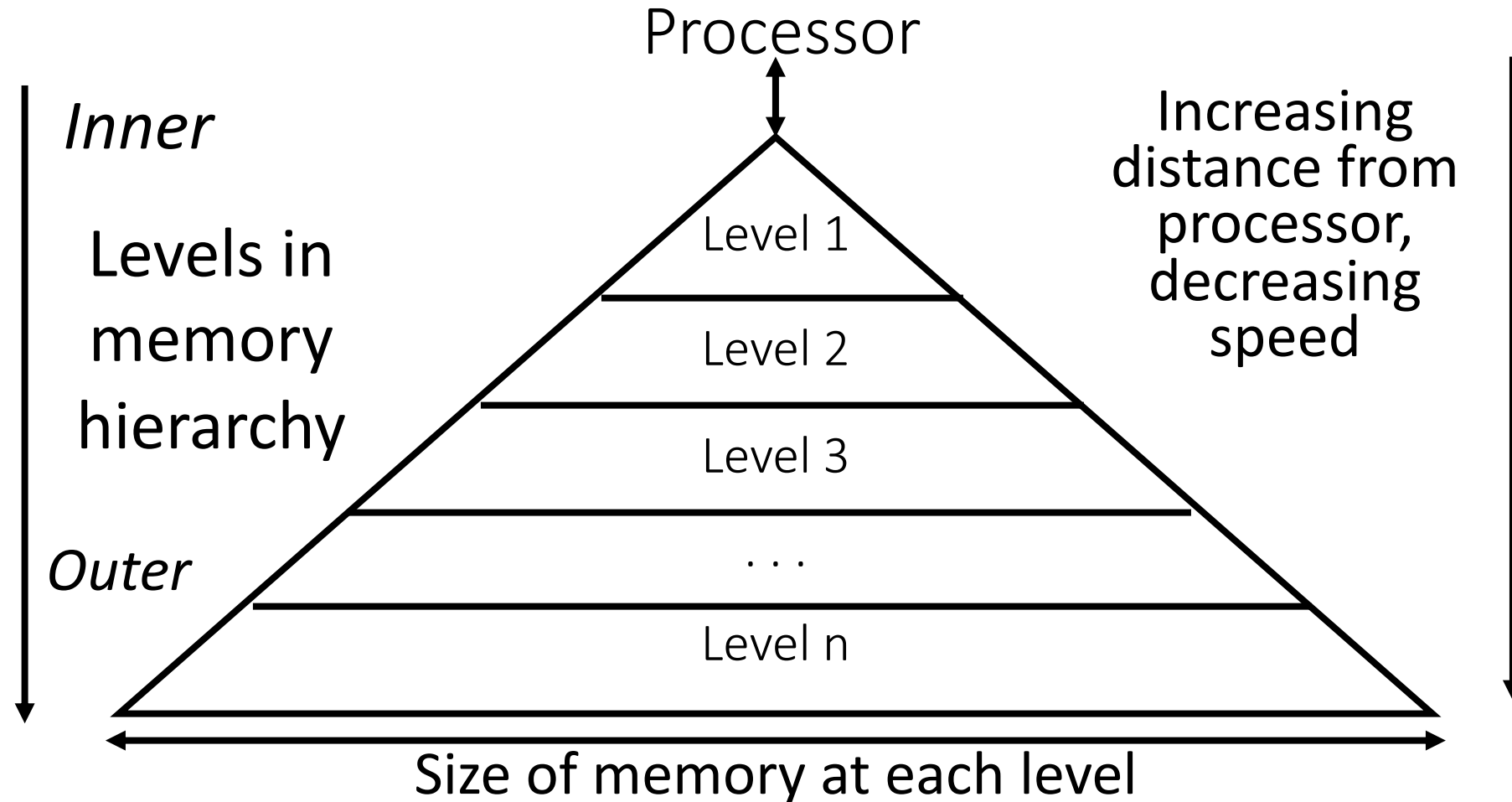


# How to Reduce Miss Penalty?

- 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 in some cases even an L3 cache!
- IBM mainframes have ~1GB L4 cache off-chip.



# Review: Memory Hierarchy

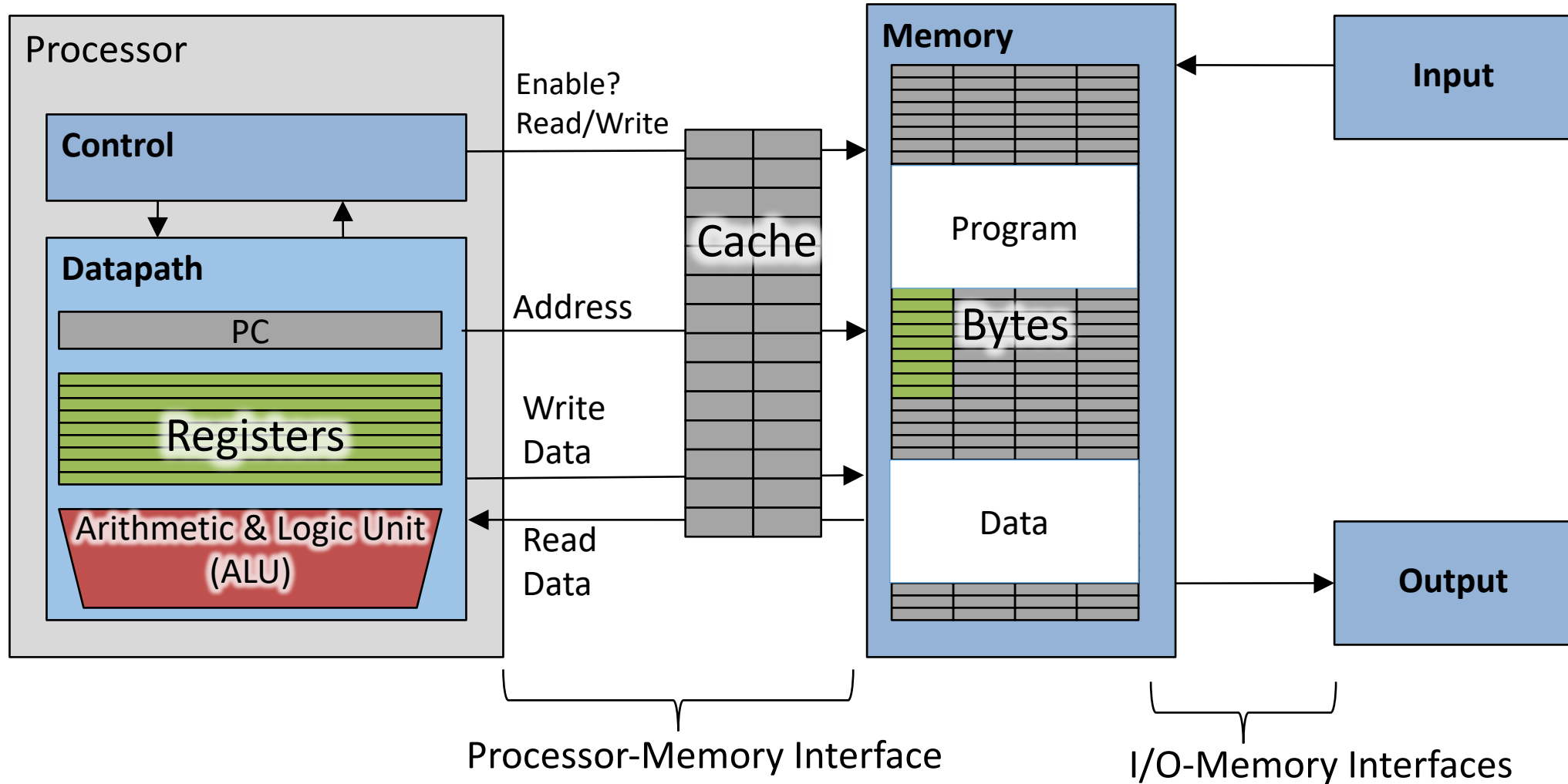


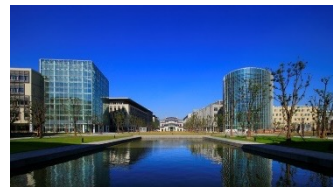
*As we move to outer levels the latency goes up and price per bit goes down.* 20



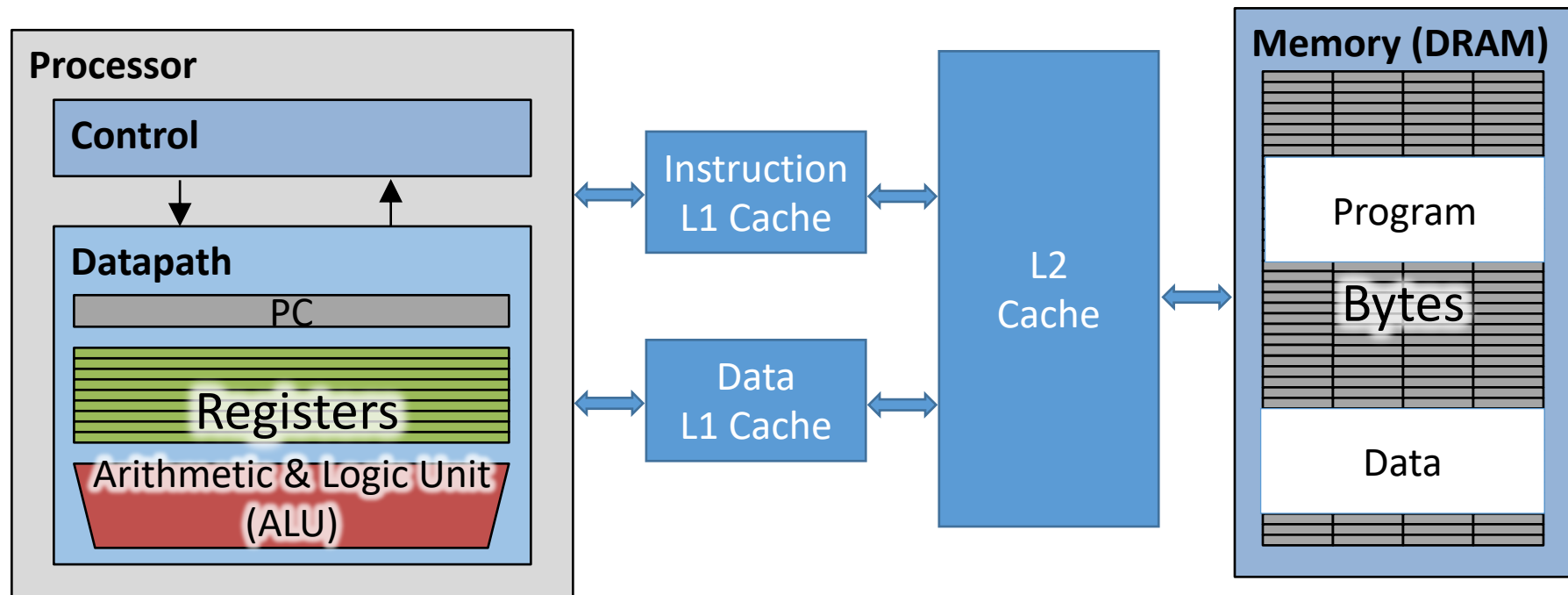


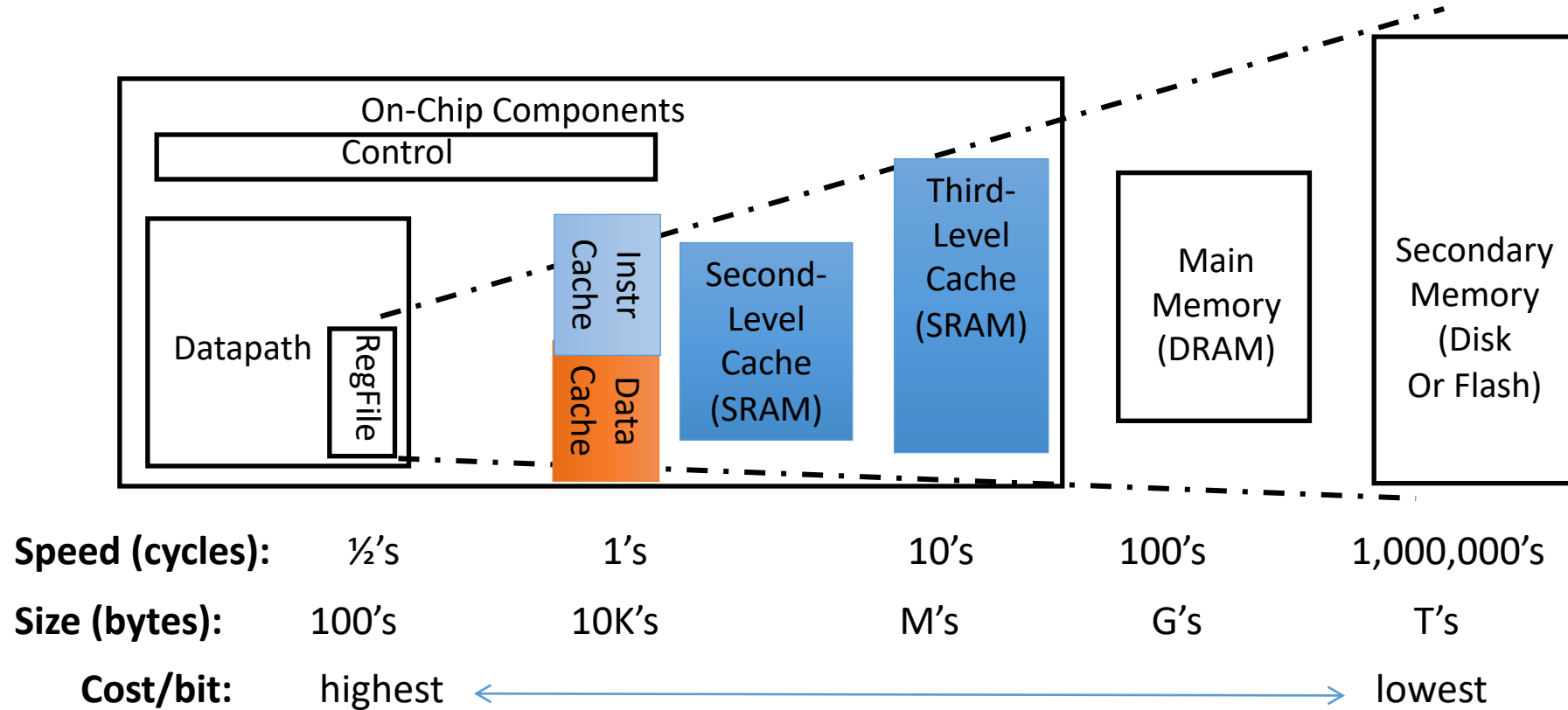
# Adding Cache to Computer





# L1 and L2 Caches





## Typical Memory Hierarchy

- **Principle of locality + memory hierarchy** presents programmer with  $\approx$  as much memory as is available in the *cheapest* technology at the  $\approx$  speed offered by the *fastest* technology



# How is the Hierarchy Managed?

- registers  $\leftrightarrow$  memory
  - By compiler (or assembly level programmer)
- cache  $\leftrightarrow$  main memory
  - By the cache controller hardware
- main memory  $\leftrightarrow$  disks (secondary storage)
  - By the operating system (virtual memory)
  - Virtual to physical address mapping assisted by the hardware ('translation lookaside buffer' or TLB)
  - By the programmer (files)

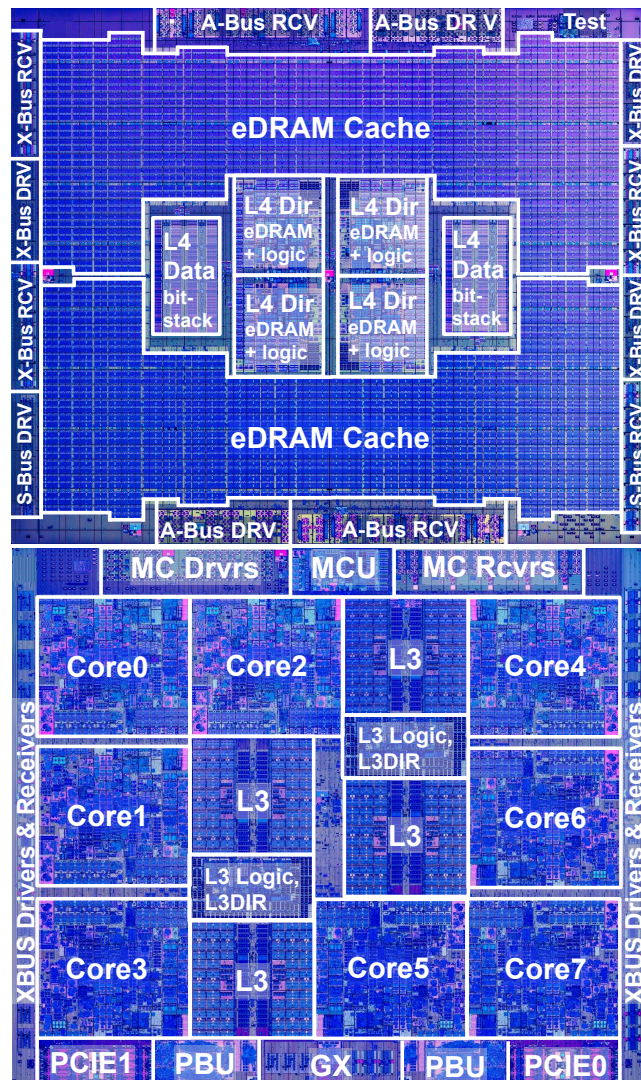
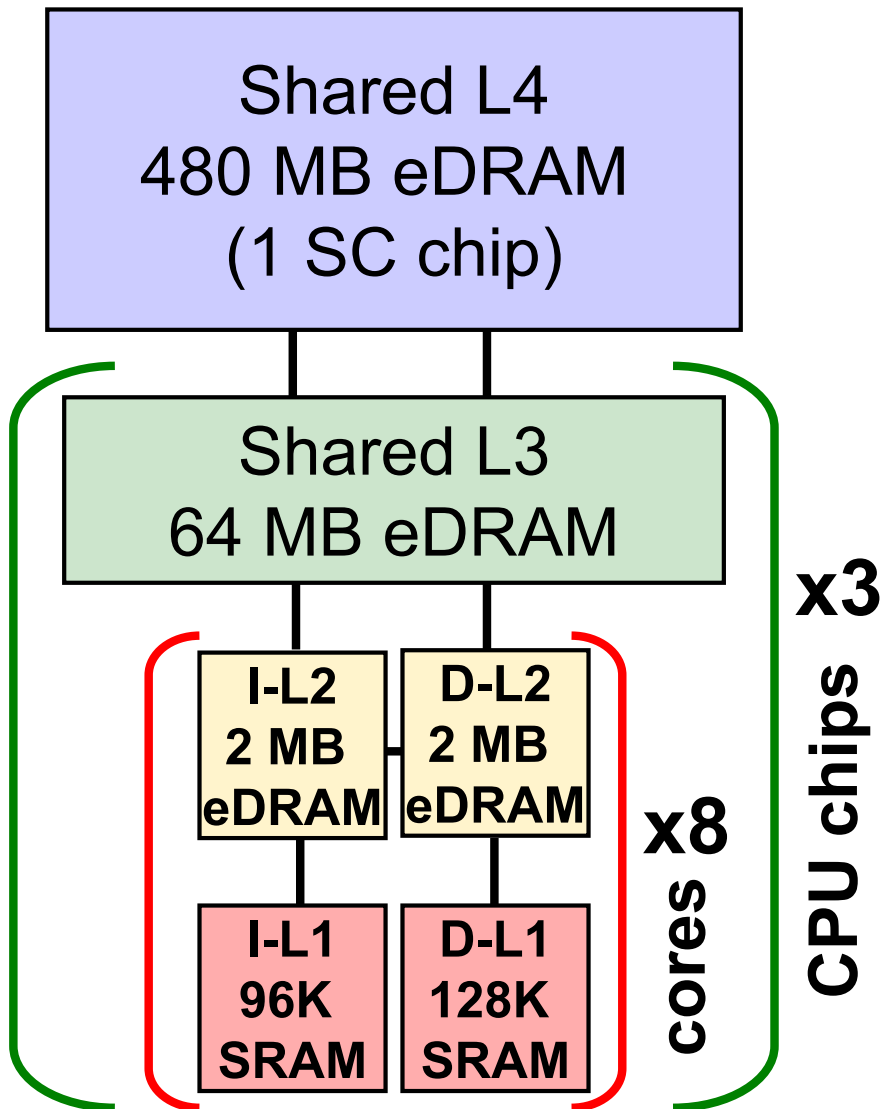
**Also a type of cache**





## 2015 IBM CPU

- z13 designed in 22nm SOI technology with **seventeen** metal layers, 4 billion transistors/chip
- 8 cores/chip, with 2MB L2 cache, 64MB L3 cache, and 480MB L4 off-chip cache.
- 5GHz clock rate, 6 instructions per cycle, 2 threads/core
- Up to 24 processor chips in shared memory node



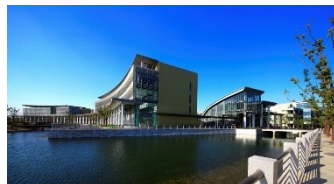
## IBM z13 Memory Hierarchy



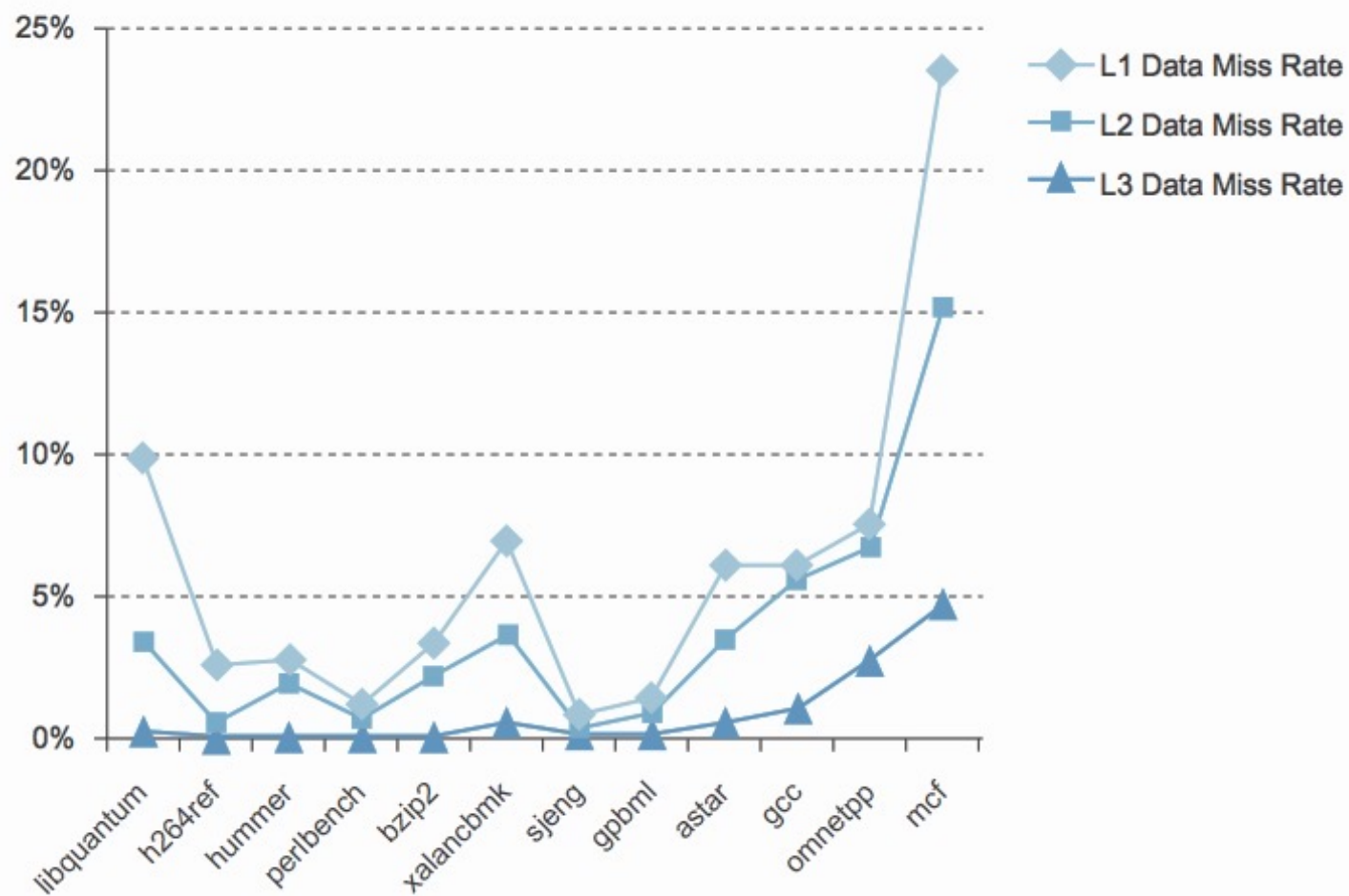
# Local vs. Global Miss Rates

- *Local miss rate* – the fraction of references to one level of a cache that miss
- Local Miss rate  $L2\$ = L2\$ \text{ Misses} / L1\$ \text{ Misses}$   
 $= L2\$ \text{ Misses} / \text{total\_L2\_accesses}$
- *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





# Miss rates



**FIGURE 5.47** The L1, L2, and L3 data cache miss rates for the Intel Core i7 920 running the full integer SPEC CPU2006 benchmarks.



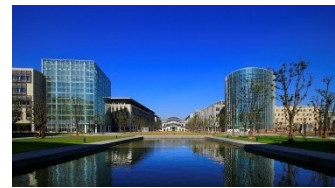
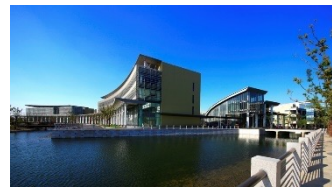


# Local vs. Global Miss Rates

- *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
- AMAT = Time for a L1\$ hit + (local) Miss rate L1\$  $\times$  (Time for a L2\$ hit + (local) Miss rate L2\$  $\times$  L2\$ Miss penalty)



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



# CPI/Miss Rates/DRAM Access (SpecInt2006)

Data Only

Data Only

Instructions and Data

Name	CPI	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

Cache Size

