# Master Informatics Eng.

2019/20

*A.J.Proença*
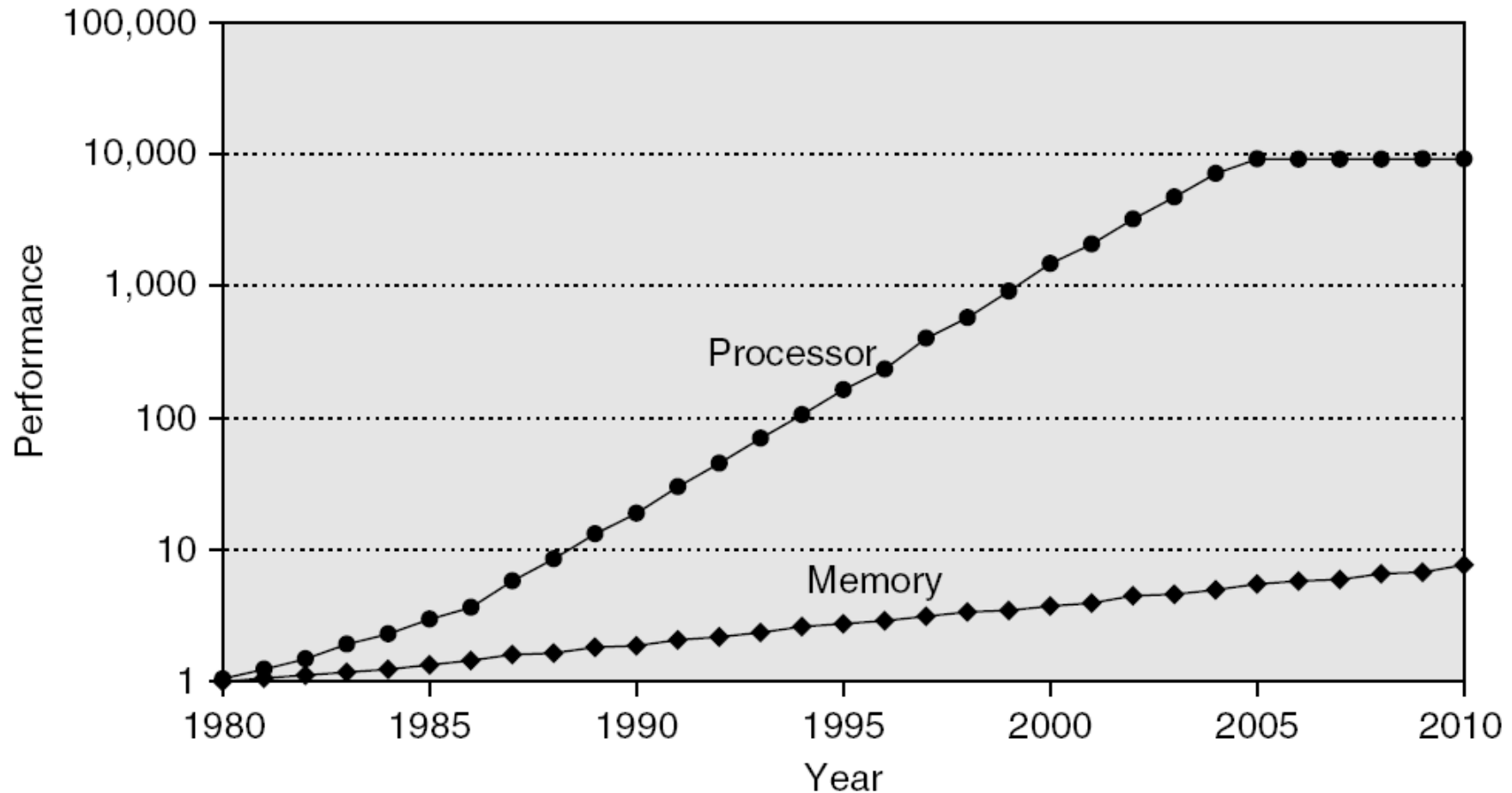
## Memory Hierarchy

**(some slides are borrowed, mod's in green)**

# Introduction

- Programmers want unlimited amounts of memory with low latency

- Fast memory technology is more expensive per bit than slower memory

- Solution: organize memory system into a hierarchy
    - Entire addressable memory space available in largest, slowest memory
    - Incrementally smaller and faster memories, each containing a subset of the memory below it, proceed in steps up toward the processor

- Temporal and spatial locality insures that nearly all references can be found in smaller memories
    - Gives the illusion of a large, fast memory being presented to the processor

# Memory Performance Gap

3

# Memory Hierarchy Design

- Memory hierarchy design becomes more crucial with recent multi-core processors:
  - Aggregate peak bandwidth grows with # cores:
    - Intel Core i7 can generate two references per core per clock
    - Four cores and 3.2 GHz clock
      - 25.6 billion* 64-bit data references/second +
      - 12.8 billion* 128-bit instruction references
      - = 409.6 GB/s!
  - DRAM bandwidth is only 6% of this (25 GB/s)
  - Requires:
    - Multi-port, pipelined caches
    - Two levels of cache per core
    - Shared third-level cache on chip
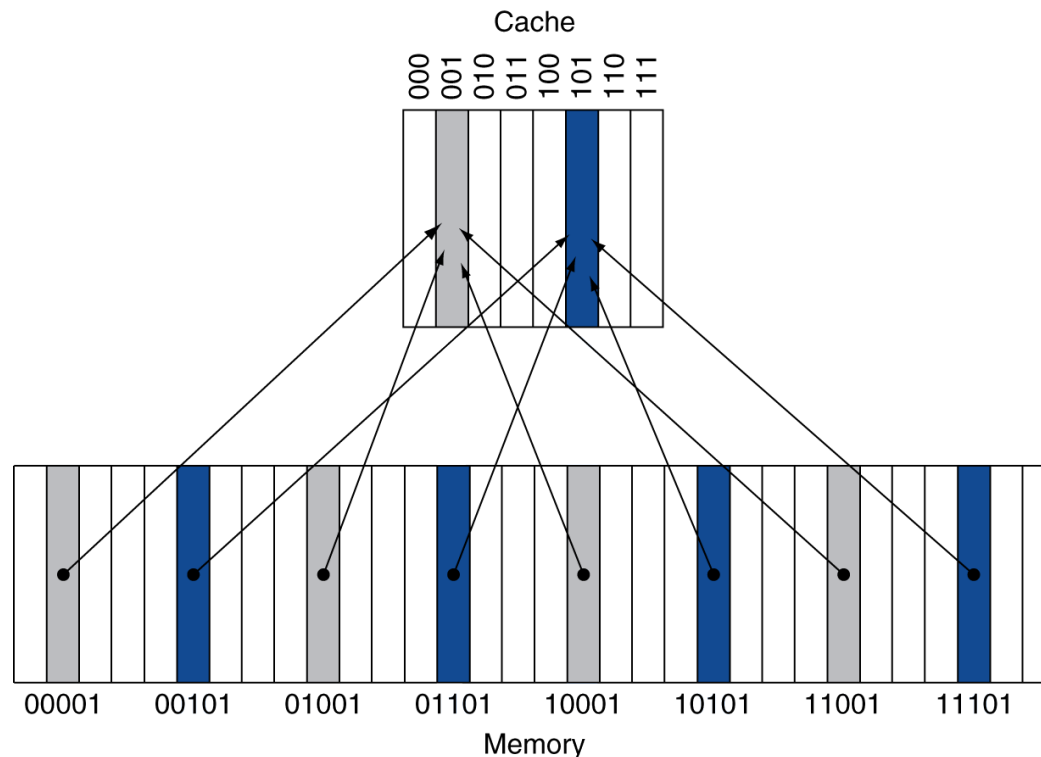
*US billion = $10^9$

# The Memory Hierarchy

## The BIG Picture

- Common principles apply at all levels of the memory hierarchy
  - Based on notions of caching
- At each level in the hierarchy
  - Block placement
  - Finding a block
  - Replacement on a miss
  - Write policy

# Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
  - (Block address) modulo (#Blocks in cache)



- #Blocks is a power of 2
- Use low-order address bits

# Associative Caches

- **Fully associative**
    - Allow a given block to go in any cache entry
    - Requires all entries to be searched at once
    - Comparator per entry (expensive)
- *n*-way set associative
    - Each set contains *n* entries
    - Block number determines which set
        - (Block number) modulo (#Sets in cache)
    - Search all entries in a given set at once
    - *n* comparators (less expensive)

# How Much Associativity

- Increased associativity decreases miss rate
    - But with diminishing returns
- Simulation of a system with 64KB D-cache, 16-word blocks, SPEC2000
    - 1-way: 10.3%
    - 2-way: 8.6%
    - 4-way: 8.3%
    - 8-way: 8.1%

# Block Placement

- Determined by associativity
  - Direct mapped (1-way associative)
    - One choice for placement
  - n-way set associative
    - n choices within a set
  - Fully associative
    - Any location

- Higher associativity reduces miss rate
  - Increases complexity, cost, and access time

# Replacement Policy

- Direct mapped: no choice
- Set associative
  - Prefer non-valid entry, if there is one
  - Otherwise, choose among entries in the set
- Least-recently used (LRU)
  - Choose the one unused for the longest time
    - Simple for 2-way, manageable for 4-way, too hard beyond that
- Random
  - Gives approximately the same performance as LRU for high associativity

# Write Policy

- ## Write-through
  - Update both upper and lower levels
  - Simplifies replacement, but may require write buffer

- ## Write-back
  - Update upper level only
  - Update lower level when block is replaced
  - Need to keep more state

- ## Virtual memory
  - Only write-back is feasible, given disk write latency

# Memory Hierarchy Basics

$$\text{CPU}_{\text{exec-time}} = (\text{CPU}_{\text{clock-cycles}} + \text{Mem}_{\text{stall-cycles}}) \times \text{Clock cycle time}$$

$$\text{CPU}_{\text{exec-time}} = (\text{IC} \times \text{CPI}_{\text{CPU}} + \text{Mem}_{\text{stall-cycles}}) \times \text{Clock cycle time}$$

$$\text{Mem}_{\text{stall-cycles}} = \text{IC} \times \text{ ... Miss rate ... Mem accesses ... Miss penalty...}$$

# Memory Hierarchy Basics

$$CPU_{\text{exec-time}} = (CPU_{\text{clock-cycles}} + Mem_{\text{stall-cycles}}) \times \text{Clock cycle time}$$

$$Mem_{\text{stall-cycles}} = IC \times \text{Misses}/\text{Instruction} \times \text{Miss Penalty}$$

$$\frac{\text{Misses}}{\text{Instruction}} = \frac{\text{Miss rate} \times \text{Memory accesses}}{\text{Instruction count}} = \text{Miss rate} \times \frac{\text{Memory accesses}}{\text{Instruction}}$$

$$\text{Average memory access time} = \text{Hit time} + \text{Miss rate} \times \text{Miss penalty}$$

- Note1: miss rate/penalty are often different for reads and writes

- Note2: speculative and multithreaded processors may execute other instructions during a miss
  - Reduces performance impact of misses

# Cache Performance Example

- ## Given
  - I-cache miss rate = 2%
  - D-cache miss rate = 4%
  - Miss penalty = 100 cycles
  - Base CPI (ideal cache) = 2
  - Load & stores are 36% of instructions
- ## Miss cycles per instruction
  - I-cache:
  - D-cache:
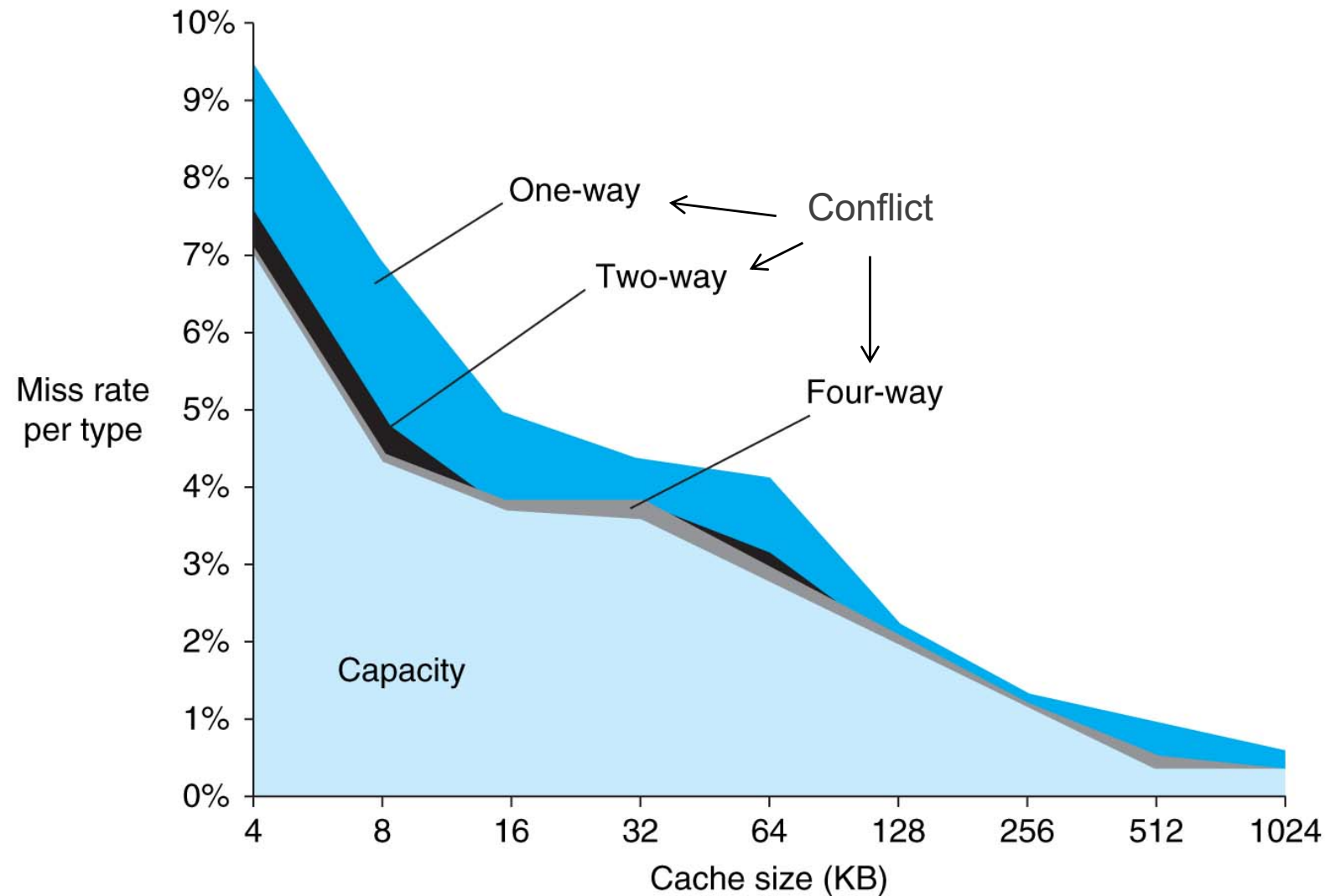- ## Actual CPI = 2 + ?? + ?? = ??

# Cache Performance Example

- Given
    - I-cache miss rate = 2%
    - D-cache miss rate = 4%
    - Miss penalty = 100 cycles
    - Base CPI (ideal cache) = 2
    - Load & stores are 36% of instructions
- Miss cycles per instruction
    - I-cache: $0.02 \times 100 = 2$
    - D-cache: $0.36 \times 0.04 \times 100 = 1.44$
- Actual CPI = 2 + 2 + 1.44 = 5.44

# Memory Hierarchy Basics

- ## Miss rate
  - ### Fraction of cache access that result in a miss

- ## Causes of misses (3C's +1)
  - ### Compulsory
    - First reference to a block
  - ### Capacity
    - Blocks discarded and later retrieved
  - ### Conflict
    - Program makes repeated references to multiple addresses from different blocks that map to the same location in the cache
  - ### Coherency
    - Different processors should see same value in same location

16

# The 3C's in diff cache sizes

# Cache Coherence

- ## Coherence
  - All reads by any processor must return the most recently written value
  - Writes to the same location by any two processors are seen in the same order by all processors

    *(Coherence defines the behaviour of reads & writes to the same memory location)*

- ## Consistency
  - When a written value will be returned by a read
  - If a processor writes location A followed by location B, any processor that sees the new value of B must also see the new value of A

    *(Consistency defines the behaviour of reads & writes with respect to accesses to other memory locations)*

18

# Enforcing Coherence

- Coherent caches provide:
  - *Migration*: movement of data
  - *Replication*: multiple copies of data

- Cache coherence protocols
  - Directory based
    - Sharing status of each block kept in one location
  - Snooping
    - Each core tracks sharing status of each block

# Memory Hierarchy Basics

- Six basic cache optimizations:
    - Larger block size
        - Reduces compulsory misses
        - Increases capacity and conflict misses, increases miss penalty
    - Larger total cache capacity to reduce miss rate
        - Increases hit time, increases power consumption
    - Higher associativity
        - Reduces conflict misses
        - Increases hit time, increases power consumption
    - Multilevel caches to reduce miss penalty
        - Reduces overall memory access time
    - Giving priority to read misses over writes
        - Reduces miss penalty
    - Avoiding address translation in cache indexing
        - Reduces hit time

# Multilevel Caches

- Primary cache attached to CPU
    - Small, but fast
- Level-2 cache services misses from primary cache
    - Larger, slower, but still faster than main memory
- Main memory services L-2 cache misses
- Some high-end systems include L-3 cache

# Multilevel Cache Example

- Given
  - CPU base CPI = 1, clock rate = 4GHz
  - Miss rate/instruction = 2%
  - Main memory access time = 100ns

- With just primary cache
  - Miss penalty = ??? = 400 cycles
  - Effective CPI = 1 + ??? = 9
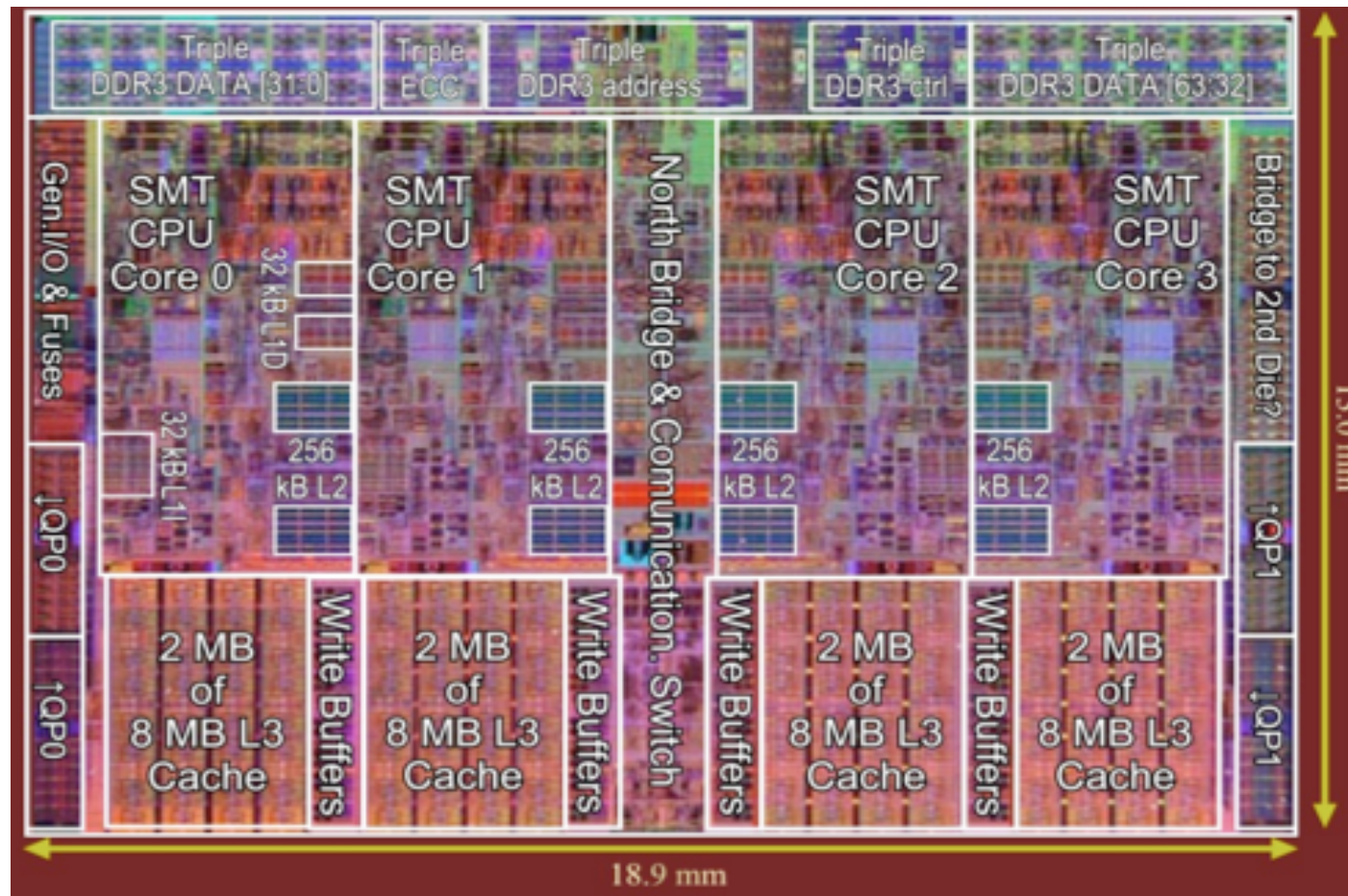
- Now add L-2 cache …

# Multilevel Cache Example

- ## Given
  - CPU base CPI = 1, clock rate = 4GHz
  - Miss rate/instruction = 2%
  - Main memory access time = 100ns
- ## With just primary cache
  - Miss penalty = 100ns/0.25ns = 400 cycles
  - Effective CPI = 1 + 0.02 × 400 = 9
- ## Now add L-2 cache …

# Example (cont.)

- Now add L-2 cache
  - Access time = 5ns
  - <u>Global</u> miss rate to main memory = 0.5%
- Primary miss with L-2 hit
  - Penalty = 5ns/0.25ns = 20 cycles
- Primary miss with L-2 miss
  - Extra penalty = 400 cycles
- CPI = 1 + 0.02 × 20 + 0.005 × 400 = 3.4
- Performance ratio = 9/3.4 = 2.6

# Multilevel On-Chip Caches

Intel Nehalem 4-core processor



Per core: 32KB L1 I-cache, 32KB L1 D-cache, 512KB L2 cache

# 3-Level Cache Organization

|  | Intel Nehalem | AMD Opteron X4 |
|---|---|---|
| L1 caches (per core) | L1 I-cache: 32KB, 64-byte blocks, 4-way, approx LRU replacement, hit time n/a<br><br>L1 D-cache: 32KB, 64-byte blocks, 8-way, approx LRU replacement, write-back/allocate, hit time n/a | L1 I-cache: 32KB, 64-byte blocks, 2-way, approx LRU replacement, hit time 3 cycles<br><br>L1 D-cache: 32KB, 64-byte blocks, 2-way, approx LRU replacement, write-back/allocate, hit time 9 cycles |
| L2 unified cache (per core) | 256KB, 64-byte blocks, 8-way, approx LRU replacement, write-back/allocate, hit time n/a | 512KB, 64-byte blocks, 16-way, approx LRU replacement, write-back/allocate, hit time n/a |
| L3 unified cache (shared) | 8MB, 64-byte blocks, 16-way, replacement n/a, write-back/allocate, hit time n/a | 2MB, 64-byte blocks, 32-way, replace block shared by fewest cores, write-back/allocate, hit time 32 cycles |

n/a: data not available