



Universidad
de Madrid
uc3m.es

Cartagena99

COMPUTER ARCHITECTURE

Cache Memory

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

--

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

roduction.

he performance.

de-offs in cache design.

ic cache optimizations.

--

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

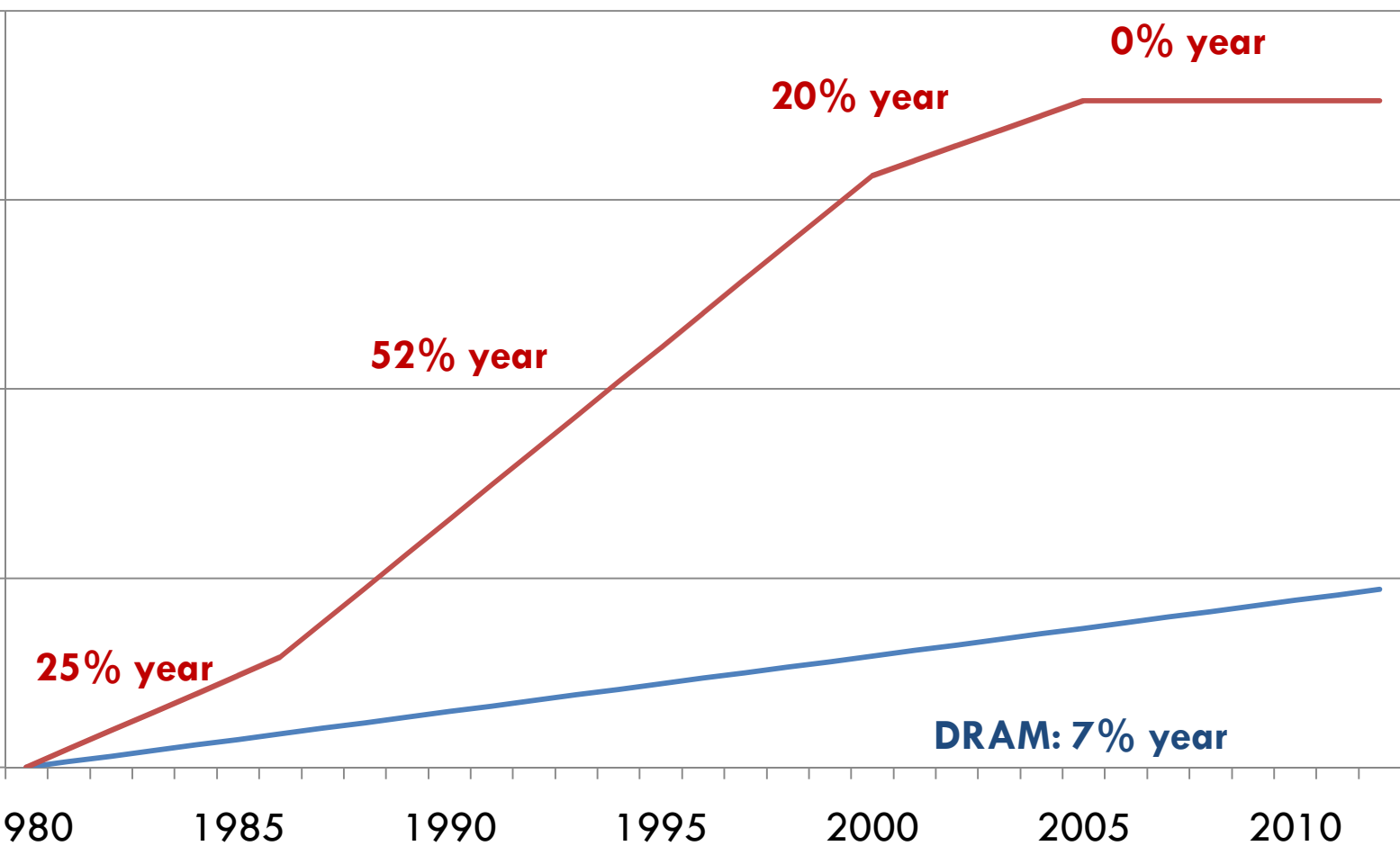
ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70



Performance evolution (1/latency)

Universidad
Complutense de Madrid
arc3m.es

— Memory — Processor



Computer Architecture - 2014 - ARCOS@uc3m



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70



Core i7

data access (64 bits) per cycle.
 cores, 3.2 GHz \rightarrow 25.6×10^9 access/sec
 instruction demand: 12.8×10^9 of 128 bits.
 peak bandwidth: 409.6 GB/sec

DRAM Memory

DDR2 (2003): 3.2 GB/sec – 8.5 GB/sec
 DDR3 (2007): 6.4 GB/sec – 17.06 GB/sec
 DDR4 (2014?): 17.05 GB/sec – 25.6 GB/sec

Techniques:

multi-gate memory, pipelined caches, multi-level caches, per-core caches, instruction/data separation.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70



Locality principle:

Program property exploited in hardware design.

Programs access to a relatively small portion of address space.

Types of locality:

Temporal locality: Recently accessed elements tend to be accessed again.

Examples: loops, variable reuse, ...

Spatial locality: Elements next to a recently accessed element tend to be accessed in near future.

Examples: sequential instruction execution, arrays, ...

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70



RAM – Static RAM

Access time: 0.5 ns – 2.5 ns

Cost per GB: 2000\$ - 5000\$

RAM – Dynamic RAM

Access time: 50ns – 70 ns

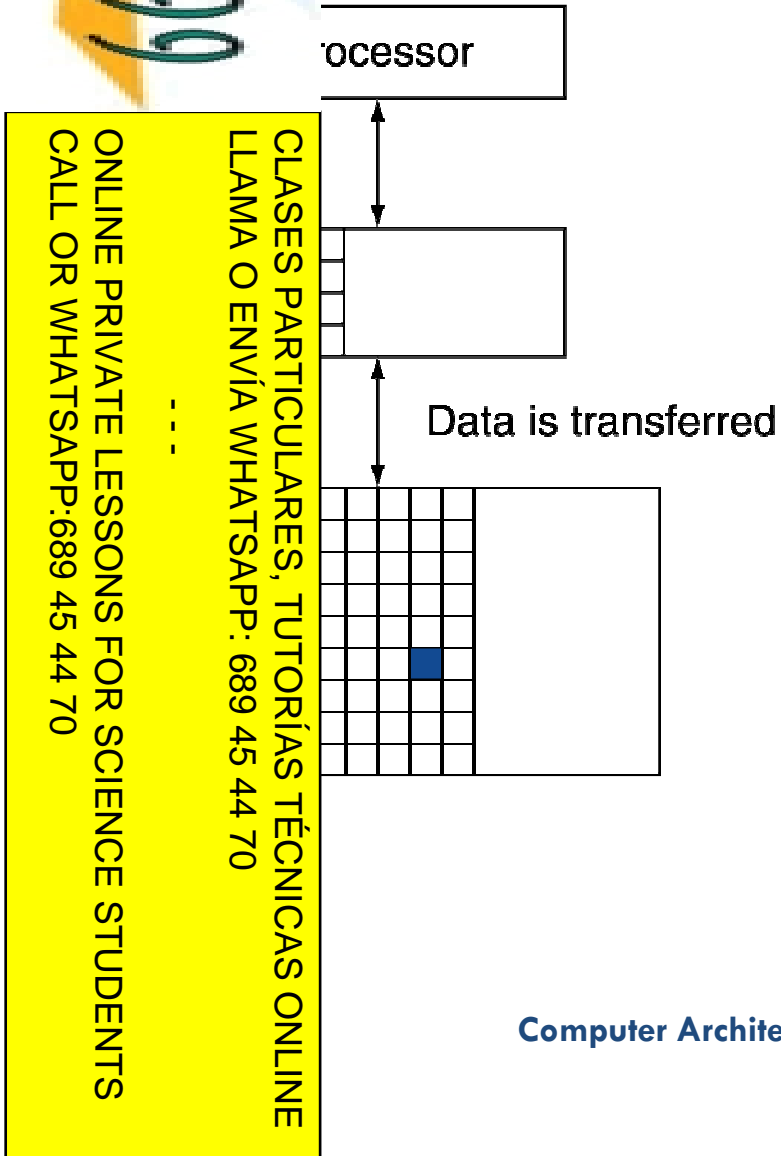
Cost per GB: 20\$ - 75\$

Magnetic disk

Access time: 5,000,000 ns – 20,000,000 ns

Cost per GB: 0.20 \$ - 2\$

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70
 - - -
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

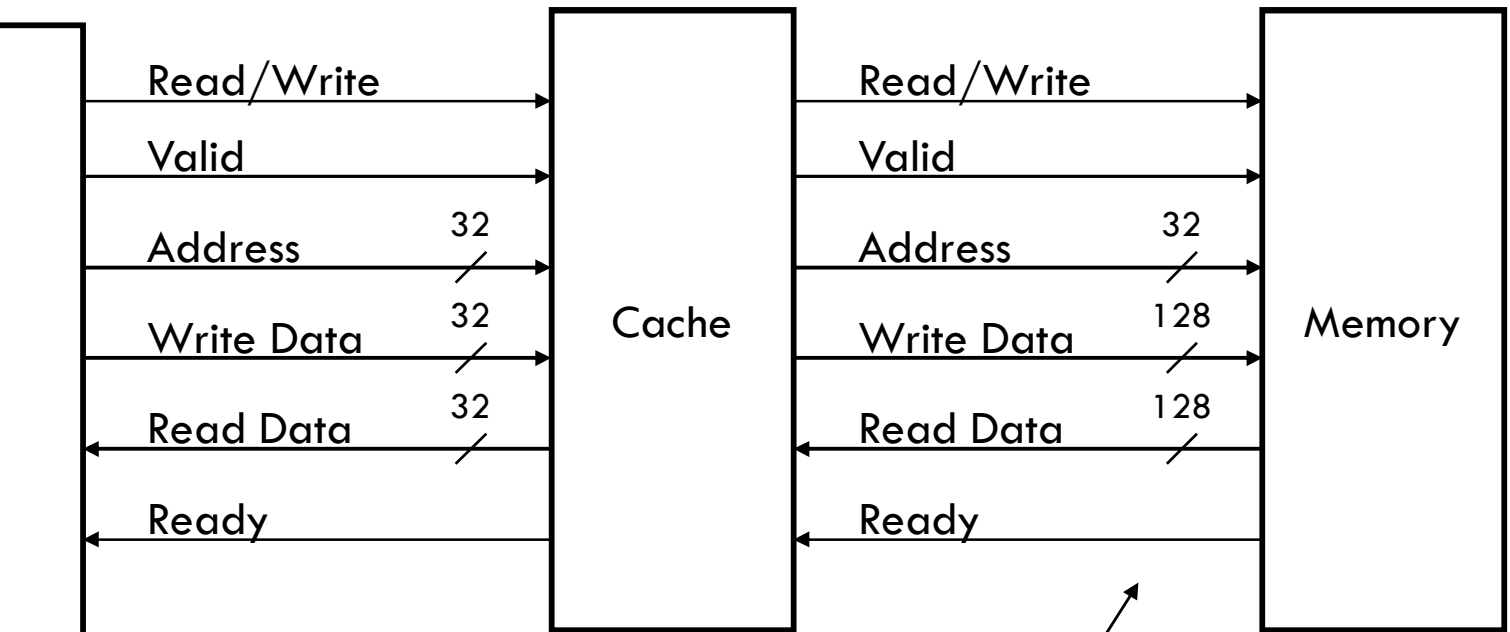


- **Block or line:** Unit of copy.
 - ▣ Usually several words.

- If accessed data present in higher level:
 - ▣ **Hit:** Delivered by higher level.
 - Hit rate = Hits / accesses.

- If accessed data is missing.
 - ▣ **Miss:** Block copied from lower level.
 - Needed time → Miss penalty.
 - Miss rate = Misses / accesses = 1 - Hit rate

Cache interconnection



Several cycles per access

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Example: Haswell Core i7



	Registers	L1I	L1D	L2	L3	DRAM	Disk
Capacity		32KB	32KB	4x256KB	8MB	32GB	1.8TB
Access time		4-5 cycles	4-5 cycles	12 cycles	36 cycles	36 cycles + 57 ns	2.8 ms

Goal: Give the illusion of large, fast, and cheap memory

Allow programs to address space scalable to disk size at register file speed.



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70



roduction.

he performance.

de-offs in cache design.

ic cache optimizations.

--

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70



average memory access time:

$$= t_H + (1-h) \cdot t_M$$

miss penalty:

Time to replace a block and deliver to CPU.

hit access time:

Time to get from lower level.

Dependent on lower level latency.

miss transfer time:

Time to transfer a block.

Depending on bandwidth across levels.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70
 --
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70



execution time

PU cycles + Memory stall cycles) x Cycle time

cycles

x CPI

Memory stall cycles

Misses number x Miss penalty

x Misses per instruction x Miss penalty

x Instruction Memory access x Miss rate x Miss penalty

where:

→ Instruction Count.

→ Cycles per instruction.

Beware: May be different for reads and writes.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

www.cartagena99.com no se hace responsable de la información contenida en el presente documento en virtud al Artículo 17.1 de la Ley de Servicios de la Sociedad de la Información y de Comercio Electrónico, de 11 de julio de 2002. Si la información contenida en el documento es ilícita o lesiona bienes o derechos de un tercero háganoslo saber y será retirada.



roduction.

he performance.

Trade-offs in cache design.

ic cache optimizations.

--

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70



Where can a block be placed in the upper level?

Block placement.

How is a block found in the upper level?

Block identification.

Which block should be replaced on a miss?

Block replacement.

What happens on a write?

Write strategy.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Q1: Block Placement



Direct mapping.

Block placement → block MOD blocks_in_cache

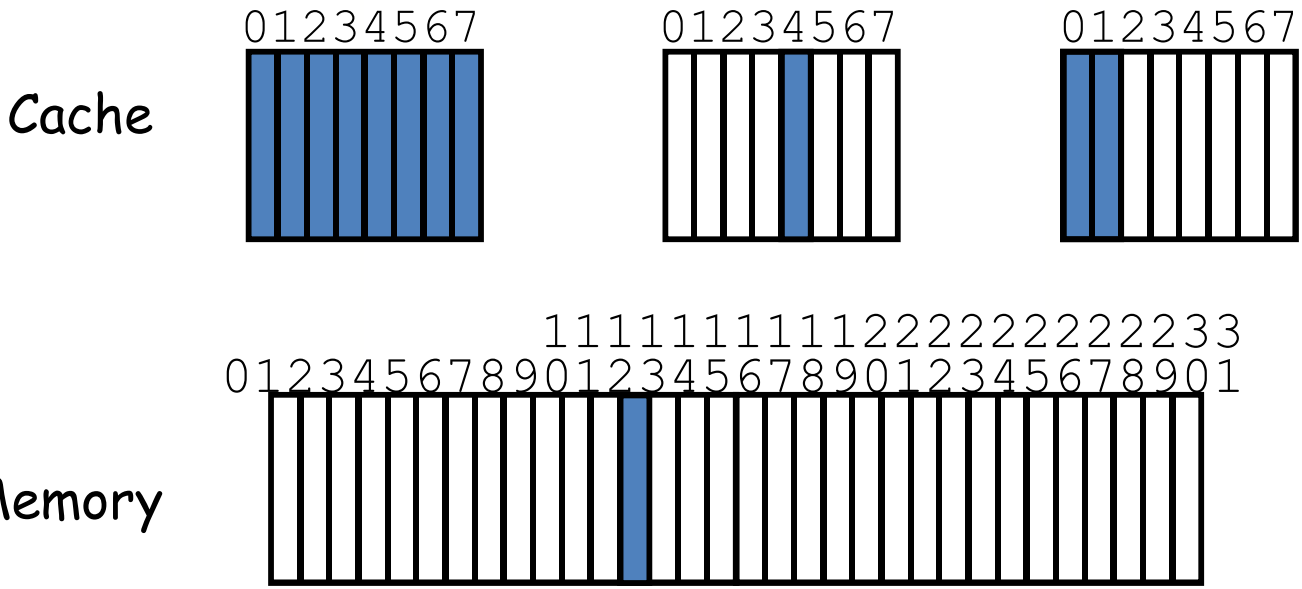
Set associative mapping.

Block placement → Anywhere.

Cache associative mapping.

Block placement → block MOD number_of_sets.

Block placement within set → Anywhere.



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Q2: Block Identification



Block Address:

Tag: Identifies the entry address.

Validity bit in every entry to flag if content is valid.

Index: Selects address.

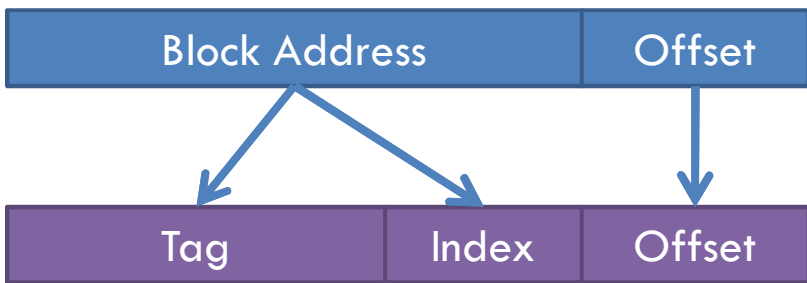
Block offset:

Select data within a block.

Number associativity means:

More bits for Index.

More bits for Tag.



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Q3: Block Replacement



relevant for associative and set associative mappings:
Rand.

Easy to implement.

LRU: Less Recently Used.

Increasing complexity as associativity increases.

FIFO: First In First Out.

Approximates LRU with lower complexity.

Misses per 1000 instr., SPEC 2000

	2 ways		4 ways			8 ways		
	Rand	FIFO	LRU	Rand	FIFO	LRU	Rand	FIFO
U	117.3	115.5	111.7	115.1	113.3	109.0	111.8	110.4
03.4	104.3	103.9	102.4	102.3	103.1	99.7	100.5	100.3
02.2	92.1	92.5	92.1	92.1	92.5	92.1	92.1	92.5

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

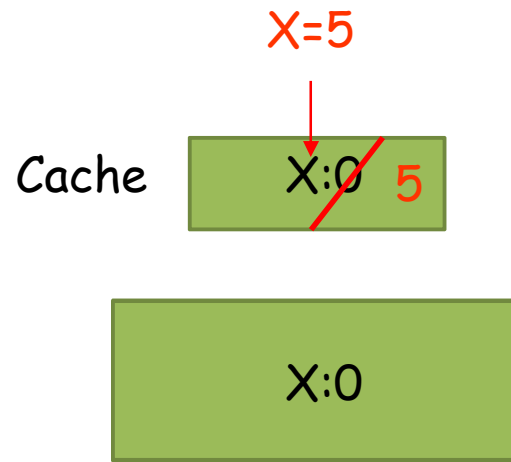
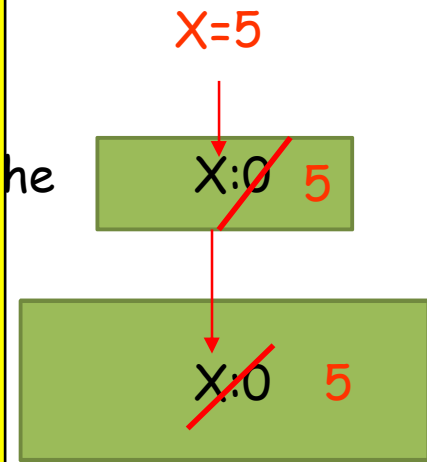
Q4: Write strategy

ough

writes sent to bus and
 copy.
 to implement.
 performance issues in SMPs .

Write back

- ❑ Many writes are a hit.
- ❑ Write hits **not** sent to bus and memory.
- ❑ Propagation and serialization issues.
- ❑ More complex.

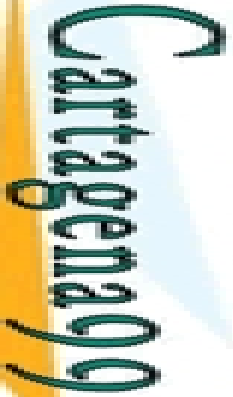


CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Q4: Write strategy

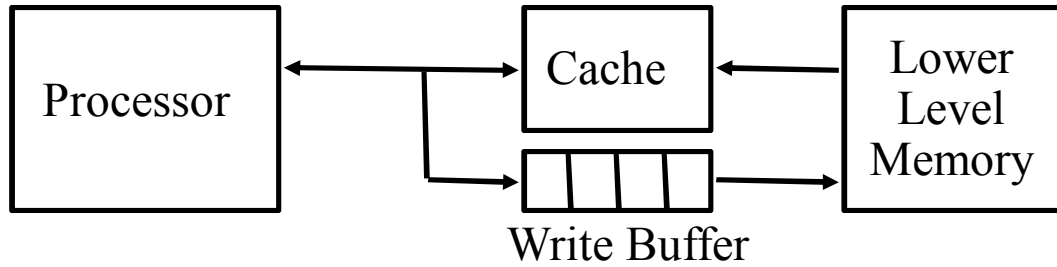
	Write-Through	Write-Back
Policy	Data written to cache blocks. Also written to next level in memory.	Data only written to cache. Update next levels when block evicted from cache
Bugging	Easy	Difficult
Write on miss?	No	Yes
Repeated writes sent to next level?	Yes	No



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Write Buffers in caches with write through



Why a buffer?

to avoid CPU stalls.

Why a buffer instead of a register?

Write bursts are common.

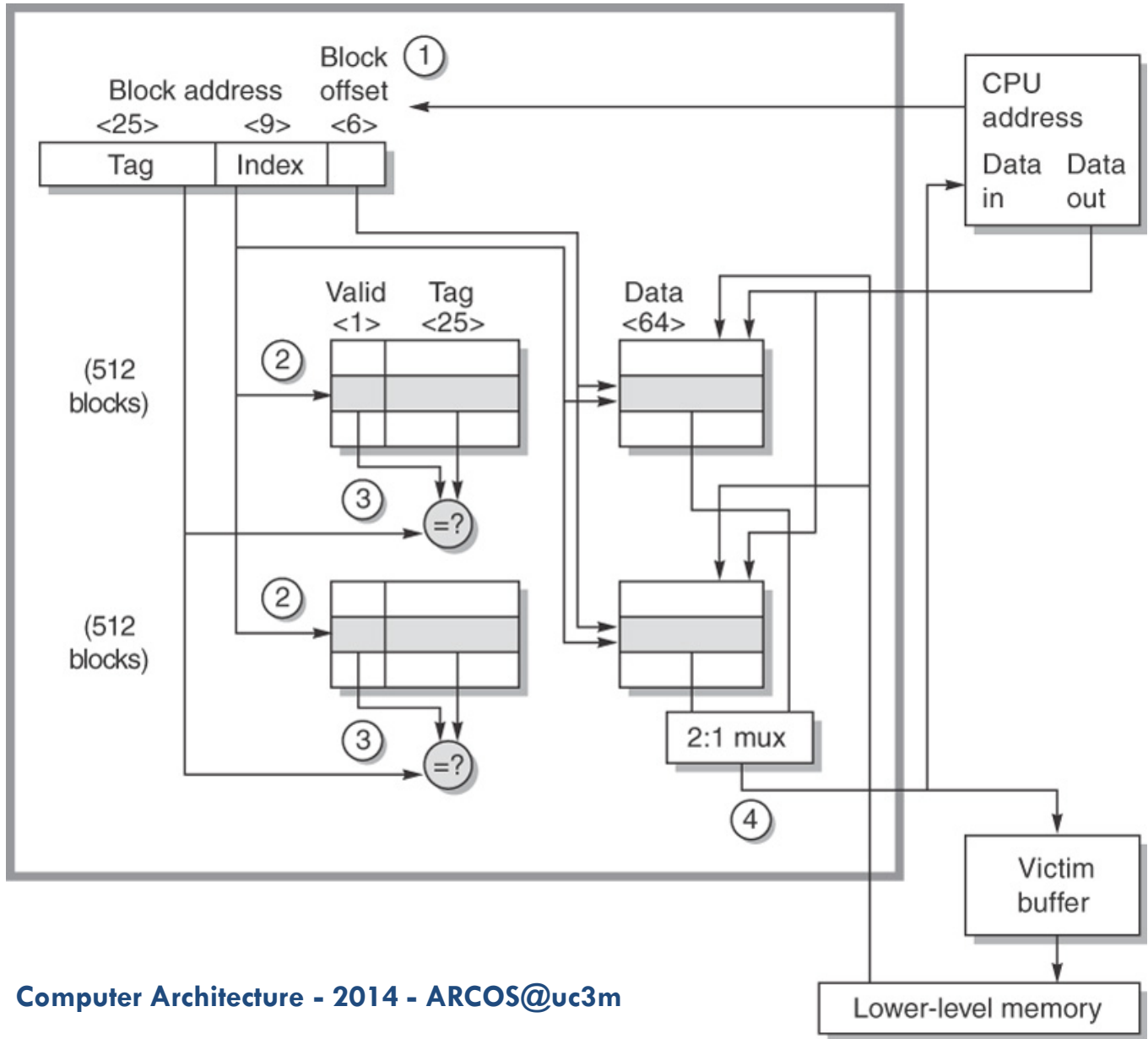
Are RAW hazards possible?

Yes.

Alternatives:

- Flush buffer before a read.
- Check buffer before a read.

Example: Opteron



Computer Architecture - 2014 - ARCOS@uc3m



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Miss penalty and Out-of-Order execution



Miss penalty definition:

Miss total latency.

Exposed latency (generating CPU stall).

Miss penalty:

Memory stalls / Instructions

$(\text{Misses/Instructions}) \times (\text{Total latency} - \text{overlapped latency})$

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70



roduction.

he performance.

de-offs in cache design.

ic cache optimizations.

--

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Reduce miss rate.

larger block size.

larger cache size.

higher associativity.

Reduce miss penalty.

multi-level caches.

prioritize read over writes.

Reduce hit time.

avoid address translation in cache indexing.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

1: Larger block size



ii: Reduce miss rate.

improve spatial locality exploitation.

Increases miss penalty.

Upon a miss, larger blocks need to be transferred.

More misses due to cache with less blocks.

Space needed:

High latency and high bandwidth memory:

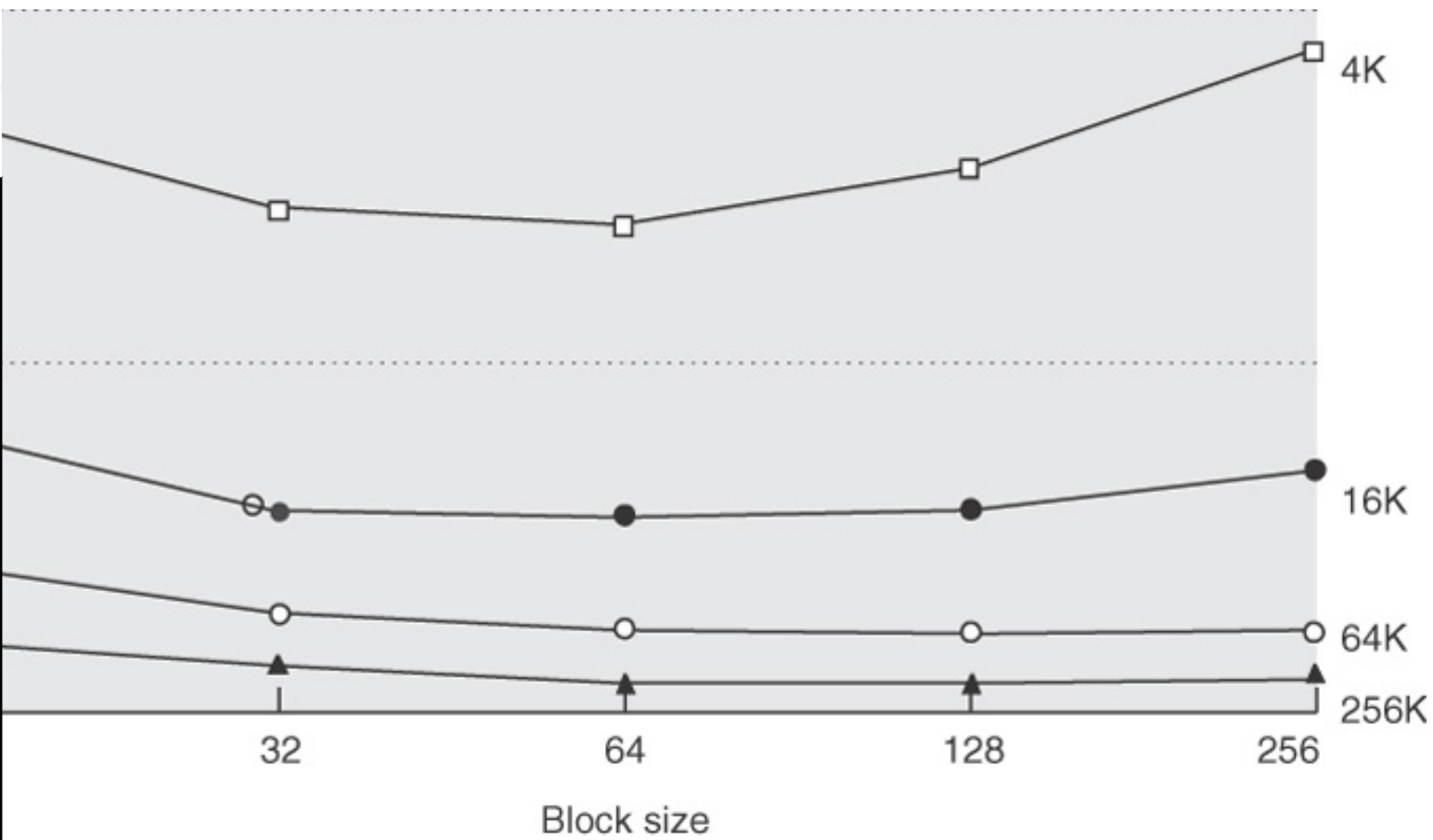
Increase block size.

Low latency and low bandwidth memory:

Reduced block size.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70
 --
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Miss rate and block size



© 2007 Elsevier, Inc. All rights reserved.



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

2: Larger cache size



1: Reduce miss rate.

More data fit in cache.

2: May increase hit time.

More time needed to find block.

3: Higher cost.

4: Higher energy consumption.

5: Harder to find balance:

Especially in on-chip caches.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

3: Higher associativity



Goal: Reduce miss rate.

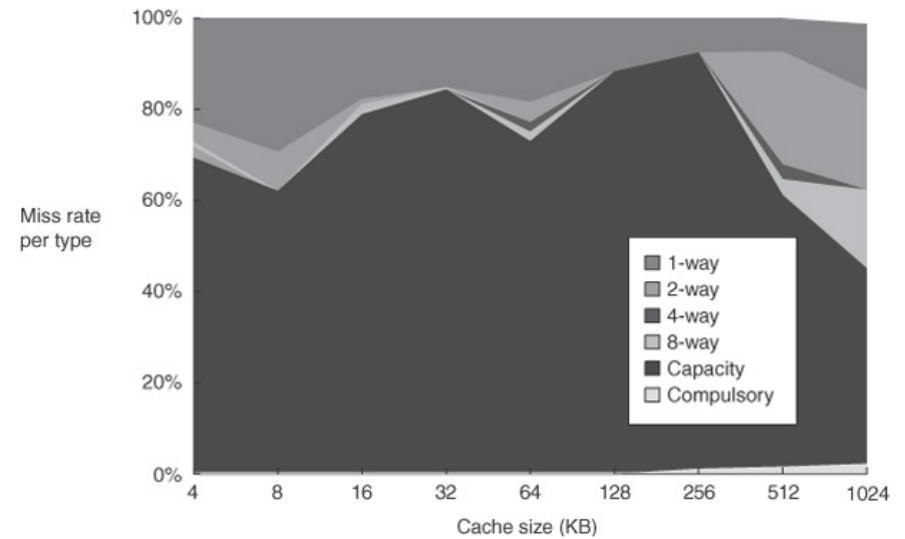
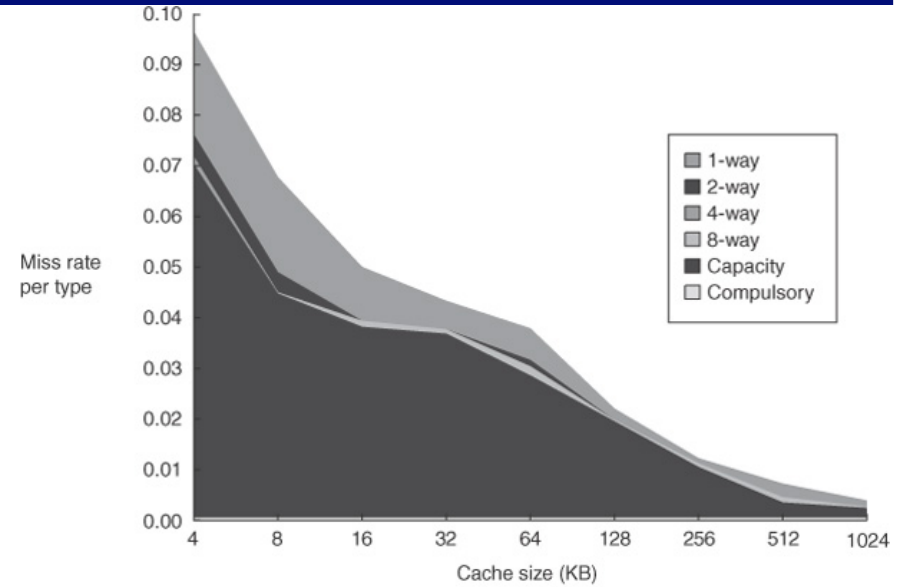
Less conflicts as more ways
in a set can be used.

Goal: Increase hit time.

More time needed to find
a block.

sequence:

ways \approx Fully associative



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

4: Multilevel caches



Goal: Reduce miss penalty.

Problem:

Increasing performance gap between DRAM and CPU.
Miss penalty cost increased over time.

Alternatives:

- Smaller caches.
- Larger caches.

Solution:

Both of them!
Several cache levels.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

average access time:

$$t_{\text{time}}_{L1} + \text{Miss rate}_{L1} \times \text{Miss penalty}_{L1}$$

$$t_{\text{time}}_{L1} + \text{Miss rate}_{L1} \times$$

$$(t_{\text{hit}}_{L2} + \text{Miss rate}_{L2} \times \text{Miss penalty}_{L2})$$

Local miss rate:

Misses at a cache level over accesses to that cache level.

$$\text{Local miss rate}_{L1} : \text{Miss rate}_{L1}$$

$$\text{Local miss rate}_{L2} : \text{Miss rate}_{L2}$$

Global miss rate:

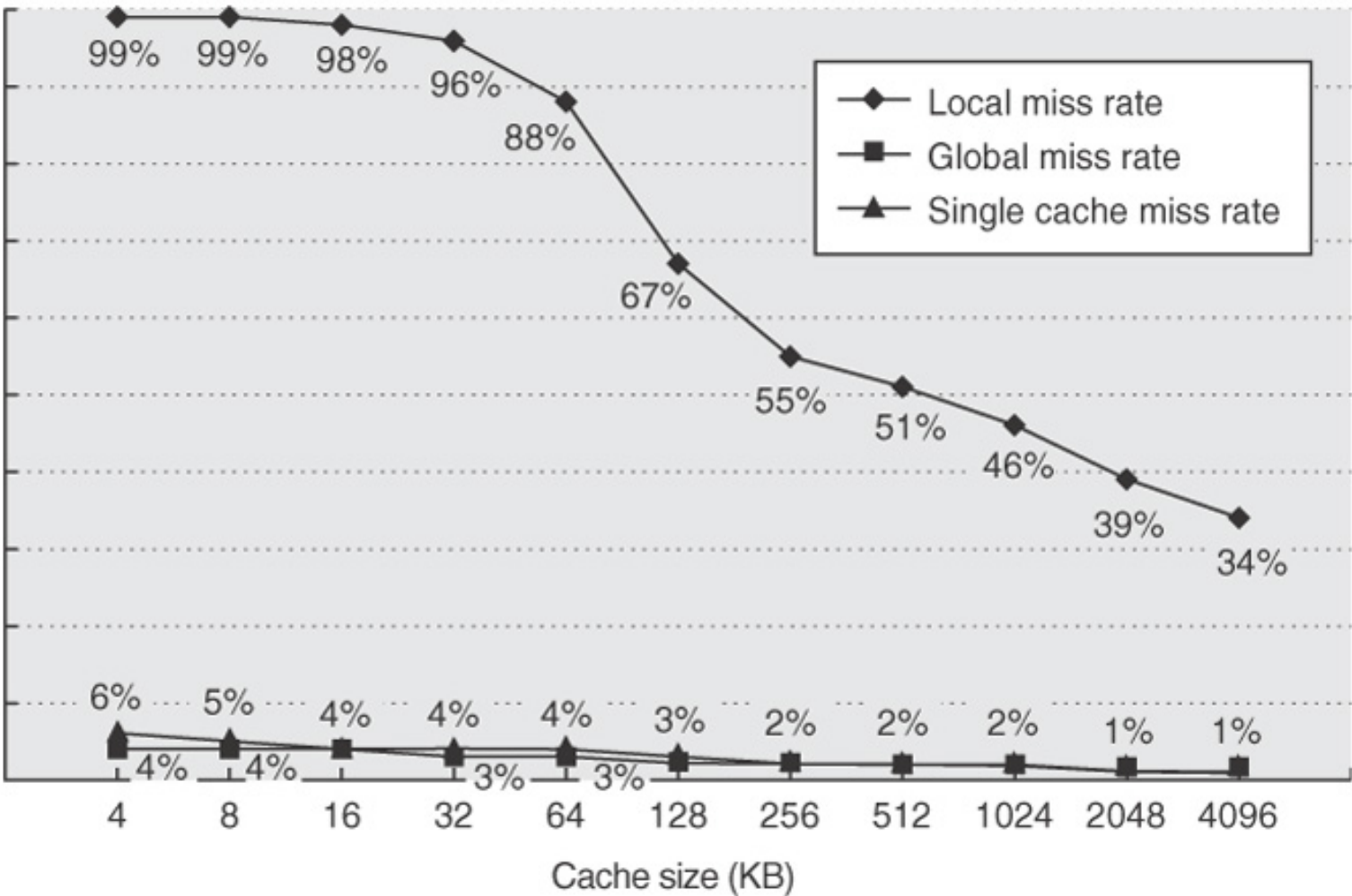
Misses at a cache level over **all** memory accesses.

$$\text{Global miss rate}_{L1} : \text{Miss rate}_{L1}$$

$$\text{Global miss rate}_{L2} : \text{Miss rate}_{L1} \times \text{Miss rate}_{L2}$$

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70
 --
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Miss rate versus cache size for multilevel caches



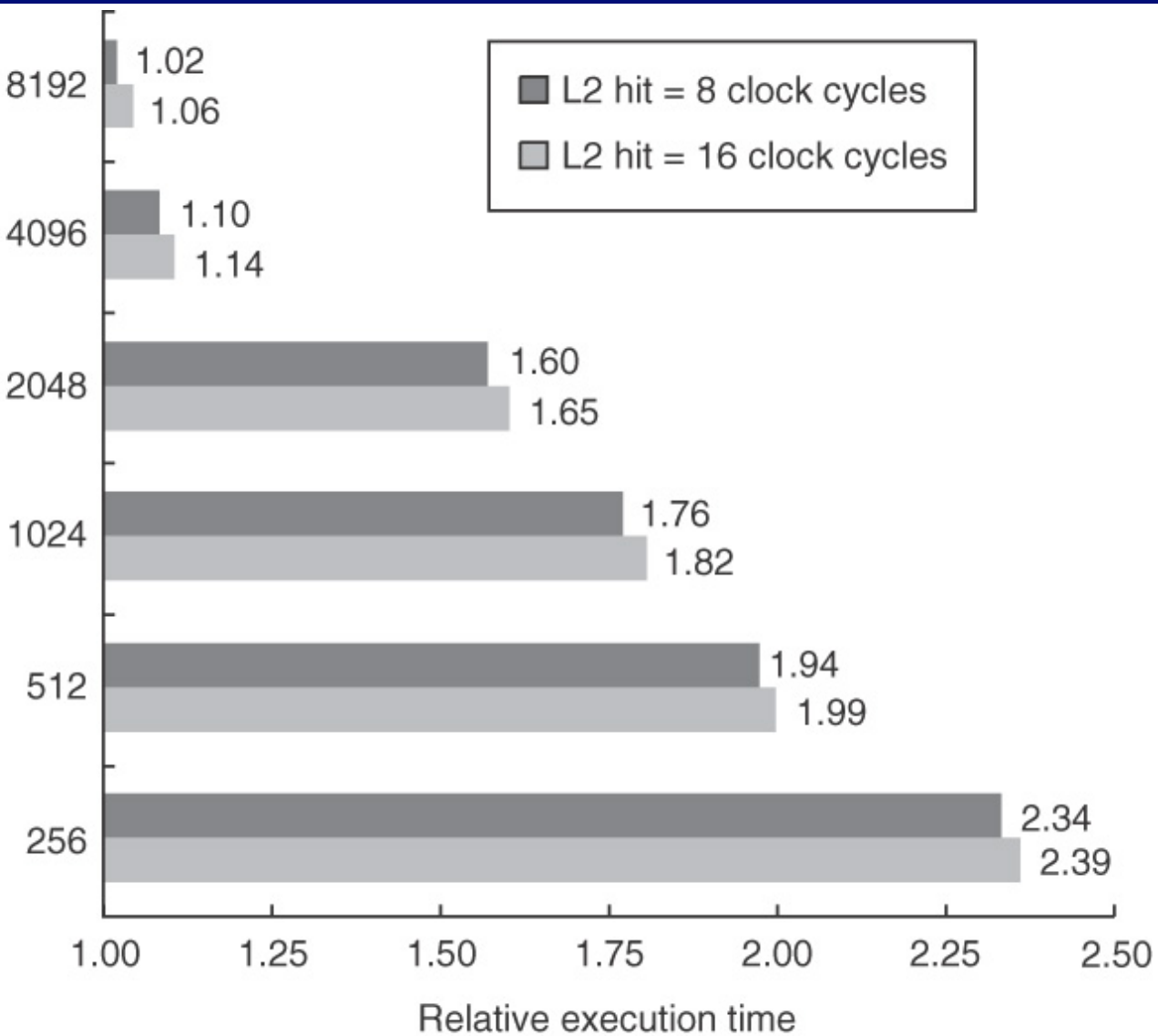
© 2007 Elsevier, Inc. All rights reserved.



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Relative execution time by second-level cache size



© 2007 Elsevier, Inc. All rights reserved.

Computer Architecture - 2014 - ARCOS@uc3m



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

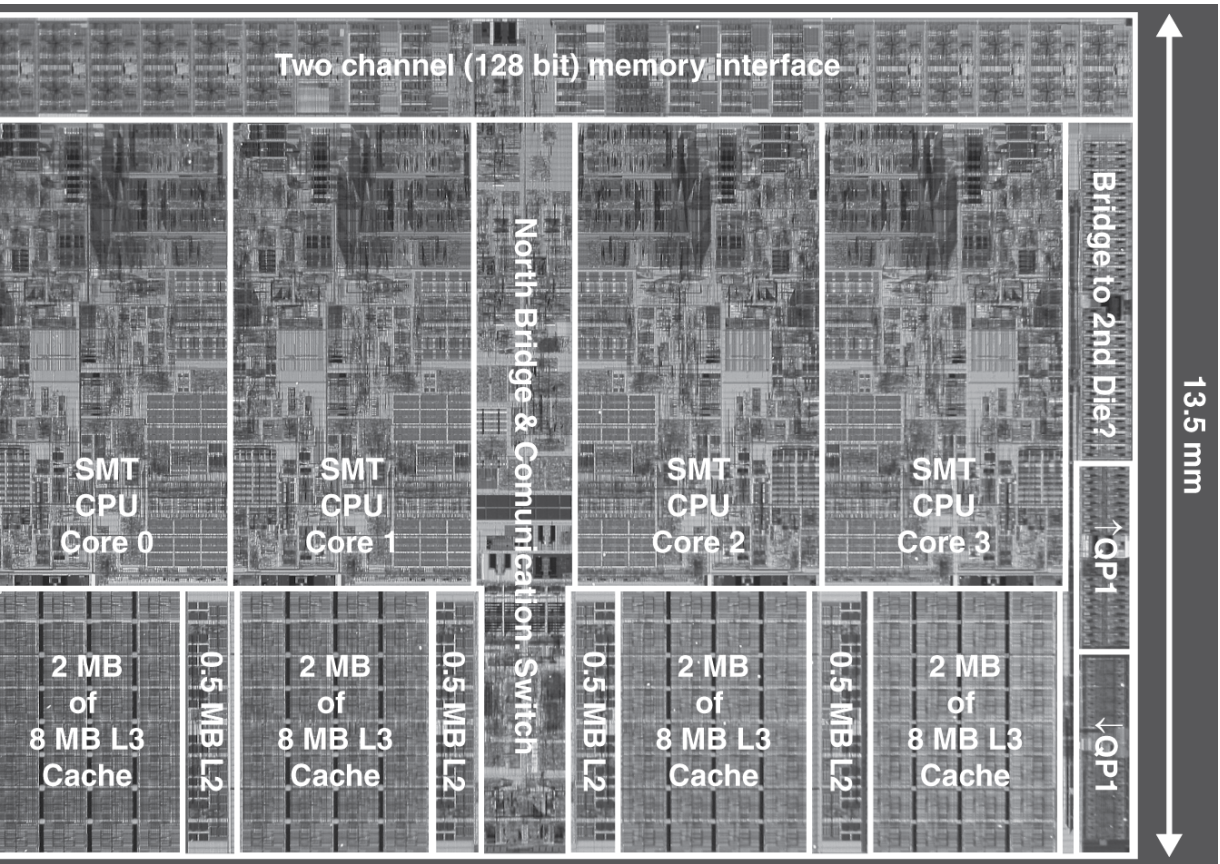
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

On-chip multilevel caches

Universidad
Complutense de Madrid
arc3m.es



Nehalem 4-core processor



Core: 32KB L1 I-cache, 32KB L1 D-cache, 512KB L2 cache

Computer Architecture - 2014 - ARCOS@uc3m

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

5: Prioritize read misses over writes

I: Reduce miss penalties.

Avoid that a read miss has to wait until writes are completed.

Write-through caches:

Write buffer could contain an updated value for the read address.

- A) Wait until write buffer is empty.
- B) Check contents of write buffer.
 - Continue with read miss if no conflict with buffer.

Write-back caches:

Read miss could replace a modified block.

Copy modified block to buffer, read, and dump block to memory.

Apply options A or B to buffer.

6: Avoid Address Translation during Indexing

ii: Reduce hit time.

Translation process:

Virtual address → Physical address.

May require additional accesses to memory.

Or at least to TLB.

ii: Optimize the most common case (hits).

Use virtual addresses for the cache.

Keys:

Indexing the cache.

Comparing tags.

Detection:

Page-level protection checked during virtual-to-physical translation.

Solution: Copy protection info from TLB on misses.

Process switch.

Virtual addresses refer to different physical addresses.

Old process virtual addresses.

Solutions:

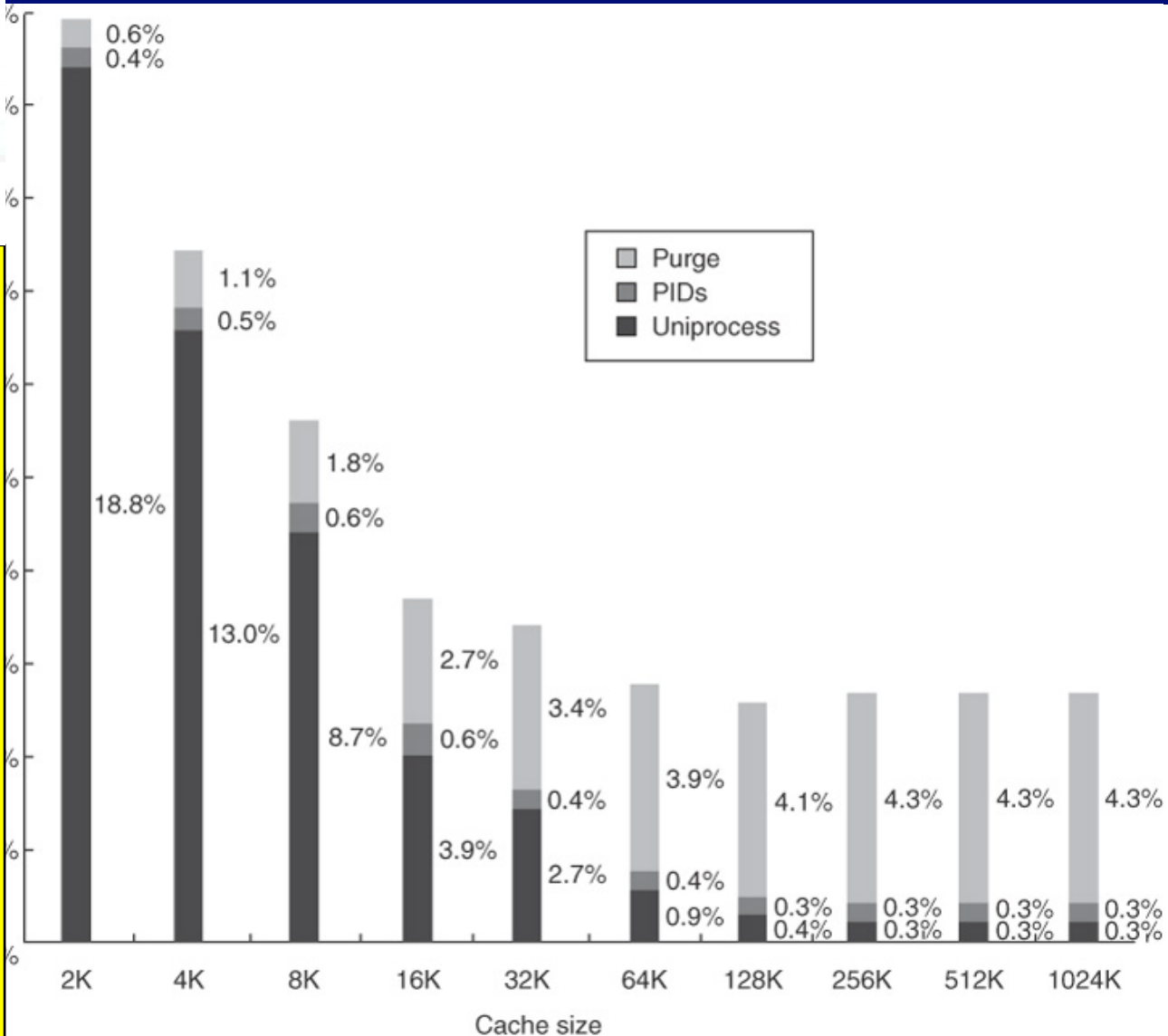
Flush the cache.

Add to cache address a PID tag.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Miss rate in virtual caches



© 2007 Elsevier, Inc. All rights reserved.
Computer Architecture - 2014 - ARCOS@uc3m



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Using:

Two different virtual addresses for the same physical address.

Anti-aliasing hardware: to guarantee that every cache lock corresponds to a unique physical address.

Check multiple addresses and invalidate.

Page coloring: Force all aliases to have identical their n last bits.

Makes impossible two alias to be at the same time in cache.

addresses:

One typically uses physical addresses.

Mapping to virtual addresses to interact with virtual caches.



Question:

Virtual indexing and physically tagging.

Steps:

Indexing the cache → Use page offset.

This part is identical in physical and virtual addresses

Comparing tags → Use translated physical address.

Tag matching uses physical address.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

hing as a **solution** to mitigate the memory wall.
 he performance due to **locality principle** (spatial
 temporal).

s **penalty** dependent on **access time** and **transfer**

key dimensions in **cache design**:

lock placement, **block identification**, **block replacement** and
 rite strategy.

basic cache optimizations:

reduce miss rate: larger block size, larger cache size,
 gher associativity.

reduce miss penalty: multilevel caches, prioritize reads
 ver writes.

reduce hit time: Avoid address translation when indexing.

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70
 --
 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70



Computer Architecture. A Quantitative Approach. 3rd Edition.

Benjamin P. Loh, John L. Hennessy y Patterson.

Exercises: B.1, B.2, B.3

Exercises:

B.2, B.3, B.4, B.5, B.6, B.7, B.8, B.9, B.10, B.11

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70