

Caches 3

Prof. Hakim Weatherspoon

CS 3410, Spring 2015

Computer Science

Cornell University

See P&H Chapter: 5.1-5.4, 5.8, 5.10, 5.15; Also, 5.13 & 5.17

Overview

Writing to caches: policies, performance

Cache tradeoffs and performance

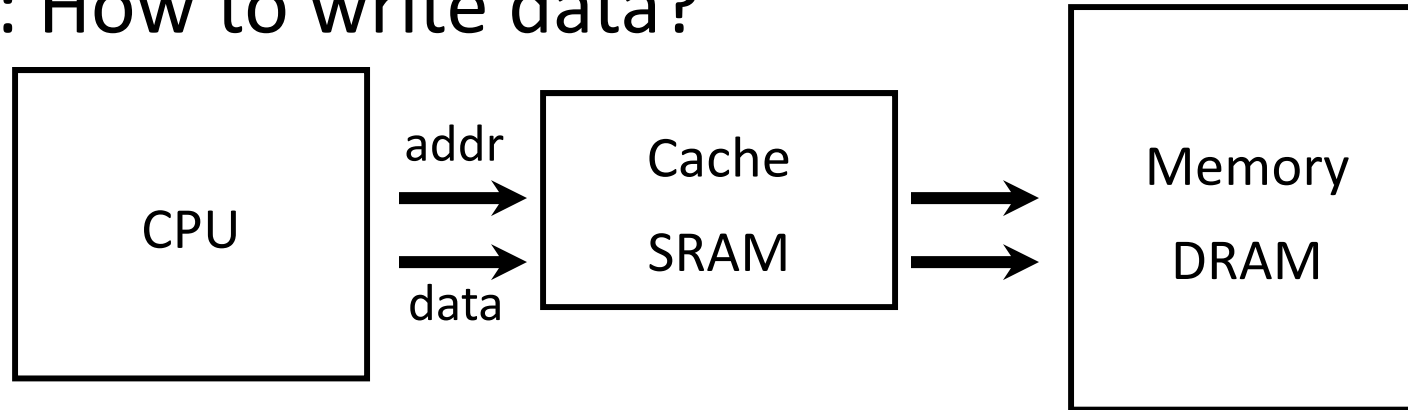
What about Stores?

Where should you write the result of a store?

- If that memory location is in the cache?
 - Send it to the cache
 - Should we also send it to memory right away?
(write-through policy)
 - Wait until we evict the block (write-back policy)
- If it is not in the cache?
 - Allocate the line (put it in the cache)?
(write allocate policy)
 - Write it directly to memory without allocation?
(no write allocate policy)

Cache Write Policies

Q: How to write data?



If data is already in the cache...

No-Write

writes invalidate the cache and go directly to memory

Write-Through

writes go to main memory and cache

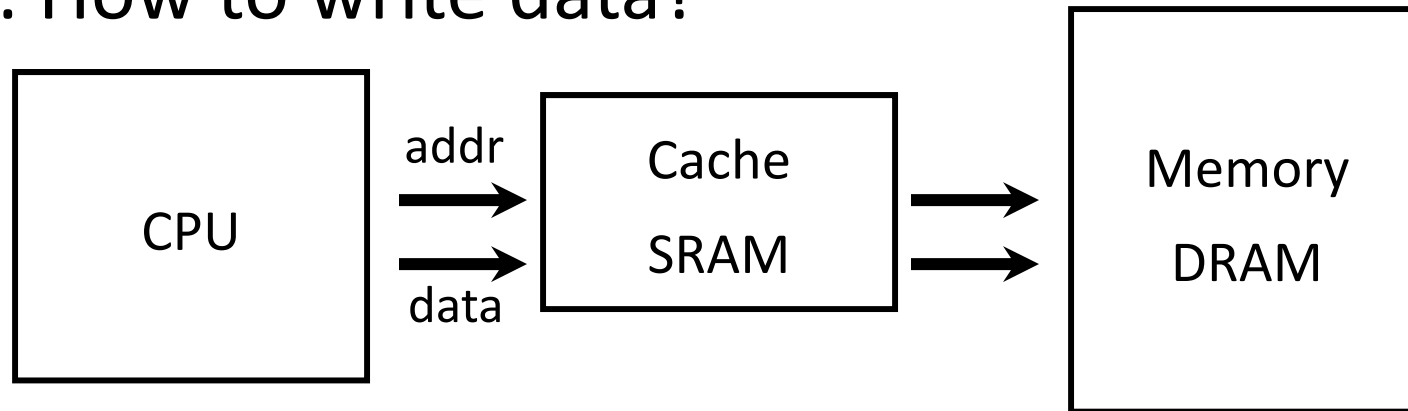
Write-Back

CPU writes only to cache

cache writes to main memory later (when block is evicted)

Write Allocation Policies

Q: How to write data?



If data is not in the cache...

Write-Allocate

allocate a cache line for new data (and maybe write-through)

No-Write-Allocate

ignore cache, just go to main memory

No write

- 2-way set associative cache
- 8 sets, use first letter to index set
- Use the initial to offset!

Speed Dial

2	ABC	Baker, J.	Baker, S.		
3	DEF	Dunn, A.	Foster, D.		
4	GHI	Gill, D.	Harris, F.	Henry, J.	Isaacs, M.
5	JKL				
6	MNO	Mann, B.	Moore, F.		
7	PQRS	Powell, C.	Sam, J.		
8	TUV	Taylor, B.	Taylor, O.		
9	WXYZ	Wright, T.	Zhang, W.		

Contacts

Baker, J.	111-111-1111
Baker, S.	222-222-2222
Dunn, A.	333-333-3333
Foster, D.	444-444-4444
Gill, D.	555-555-5555
Harris, F.	666-666-6666
Henry, J.	777-777-7777
Isaacs, M.	888-888-8888
Mann, B.	111-111-1119
Moore, F.	222-222-2229
Powell, C.	333-333-3339
Sam, J.	444-444-4449
Taylor, B.	555-555-5559
Taylor, O.	666-666-6669
Wright, T.	777-777-7779
Zhang, W.	888-888-8889

Block Size: 2 Contact Numbers

Write through

- 2-way set associative cache
- 8 sets, use first letter to index set
- Use the initial to offset!

Speed Dial

2	ABC	Baker, J.	Baker, S.		
3	DEF	Dunn, A.	Foster, D.		
4	GHI	Gill, D.	Harris, F.	Henry, J.	Isaacs, M.
5	JKL				
6	MNO	Mann, B.	Moore, F.		
7	PQRS	Powell, C.	Sam, J.		
8	TUV	Taylor, B.	Taylor, O.		
9	WXYZ	Wright, T.	Zhang, W.		

Contacts

Baker, J.	111-111-1111
Baker, S.	222-222-2222
Dunn, A.	333-333-3333
Foster, D.	444-444-4444
Gill, D.	555-555-5555
Harris, F.	666-666-6666
Henry, J.	777-777-7777
Isaacs, M.	888-888-8888
Mann, B.	111-111-1119
Moore, F.	222-222-2229
Powell, C.	333-333-3339
Sam, J.	444-444-4449
Taylor, B.	555-555-5559
Taylor, O.	666-666-6669
Wright, T.	777-777-7779
Zhang, W.	888-888-8889

Block Size: 2 Contact Numbers

Write back

- 2-way set associative cache
- 8 sets, use first letter to index set
- Use the initial to offset!

Speed Dial

2	ABC	Baker, J.	Baker, S.		
3	DEF	Dunn, A.	Foster, D.		
4	GHI	Gill, D.	Harris, F.	Henry, J.	Isaacs, M.
5	JKL				
6	MNO	Mann, B.	Moore, F.		
7	PQRS	Powell, C.	Sam, J.		
8	TUV	Taylor, B.	Taylor, O.		
9	WXYZ	Wright, T.	Zhang, W.		

Contacts

Baker, J.	111-111-1111
Baker, S.	222-222-2222
Dunn, A.	333-333-3333
Foster, D.	444-444-4444
Gill, D.	555-555-5555
Harris, F.	666-666-6666
Henry, J.	777-777-7777
Isaacs, M.	888-888-8888
Mann, B.	111-111-1119
Moore, F.	222-222-2229
Powell, C.	333-333-3339
Sam, J.	444-444-4449
Taylor, B.	555-555-5559
Taylor, O.	666-666-6669
Wright, T.	777-777-7779
Zhang, W.	888-888-8889

Block Size: 2 Contact Numbers

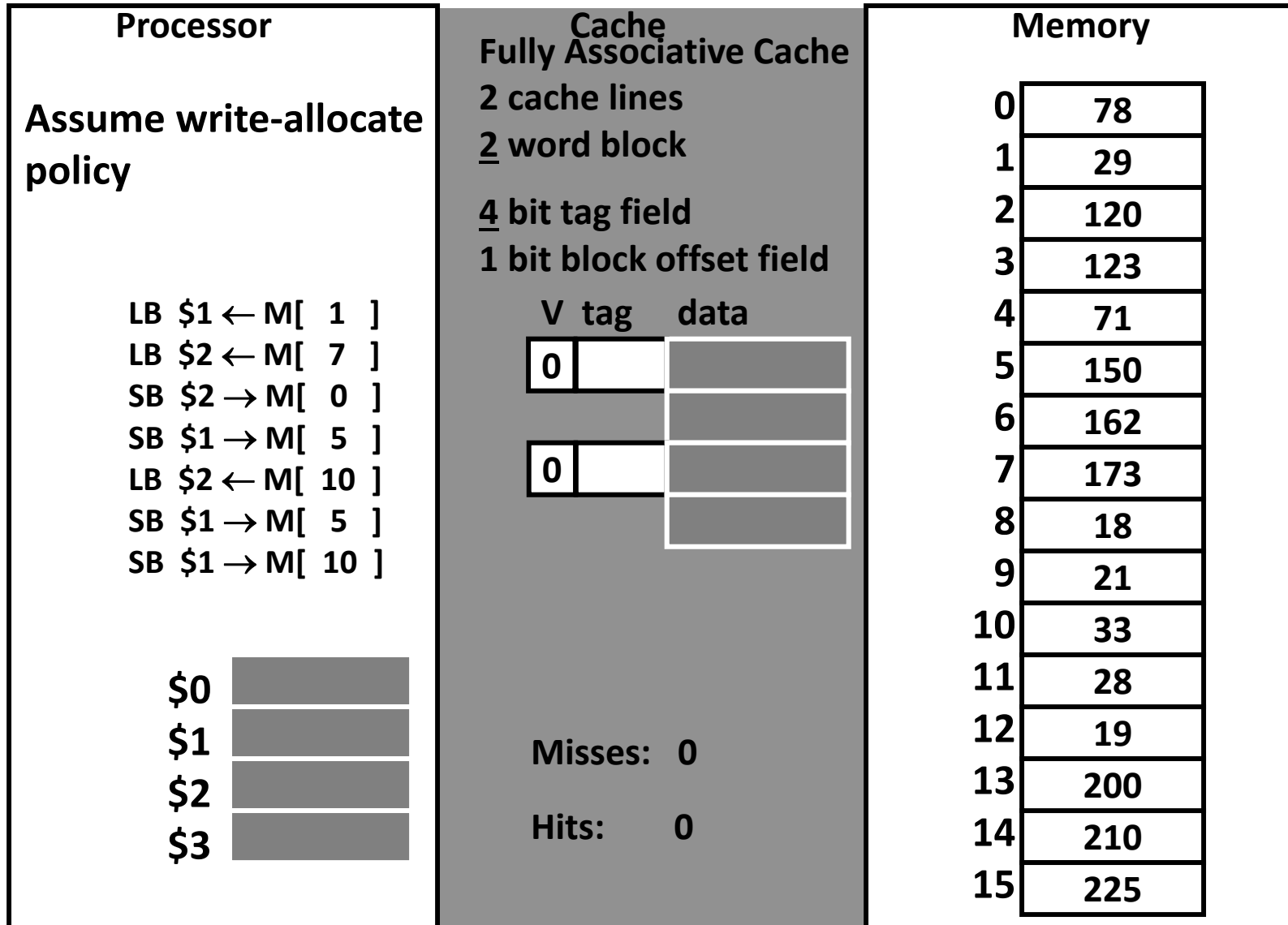
Next Goal

How does a write-through cache work?

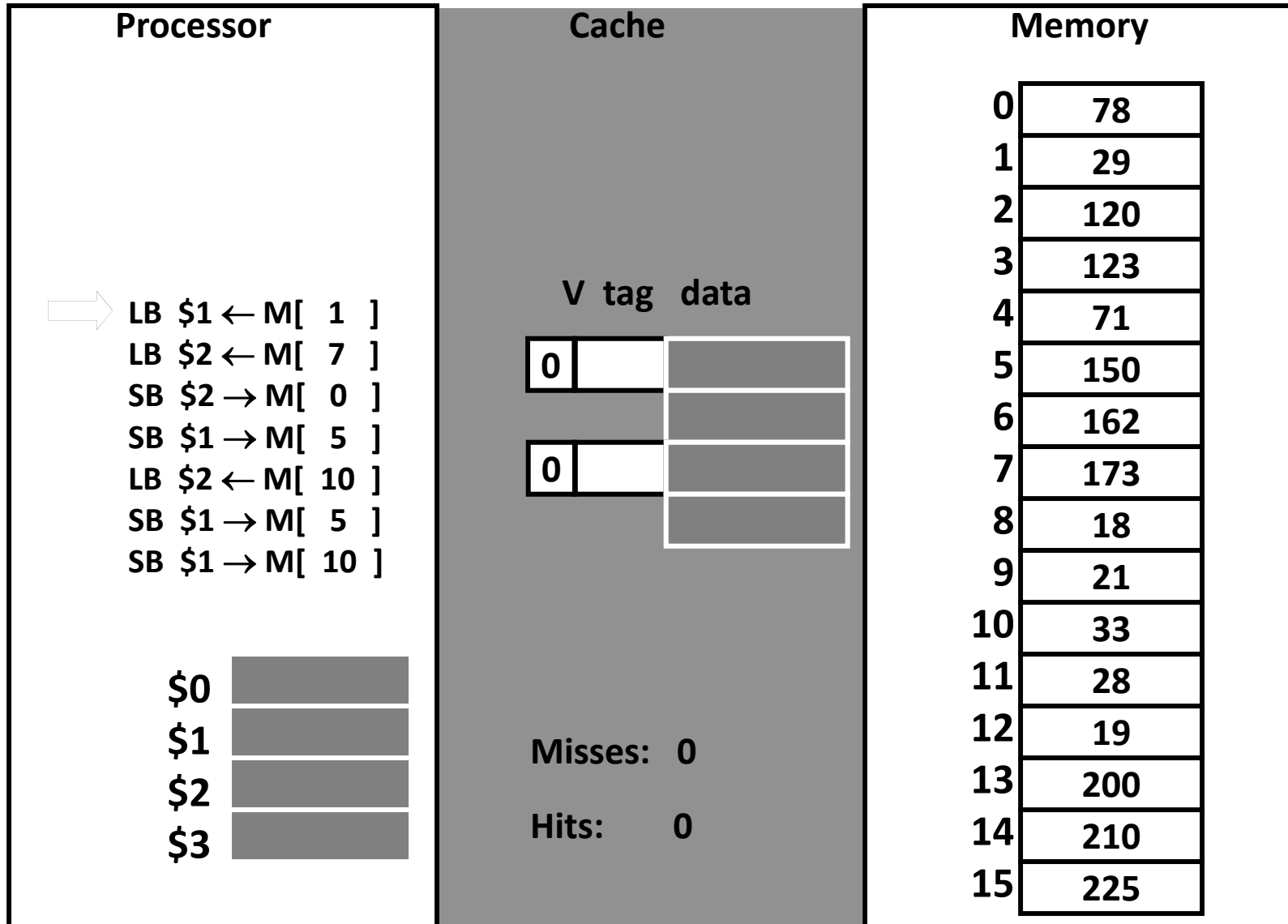
Assume write-allocate

Handling Stores (Write-Through)

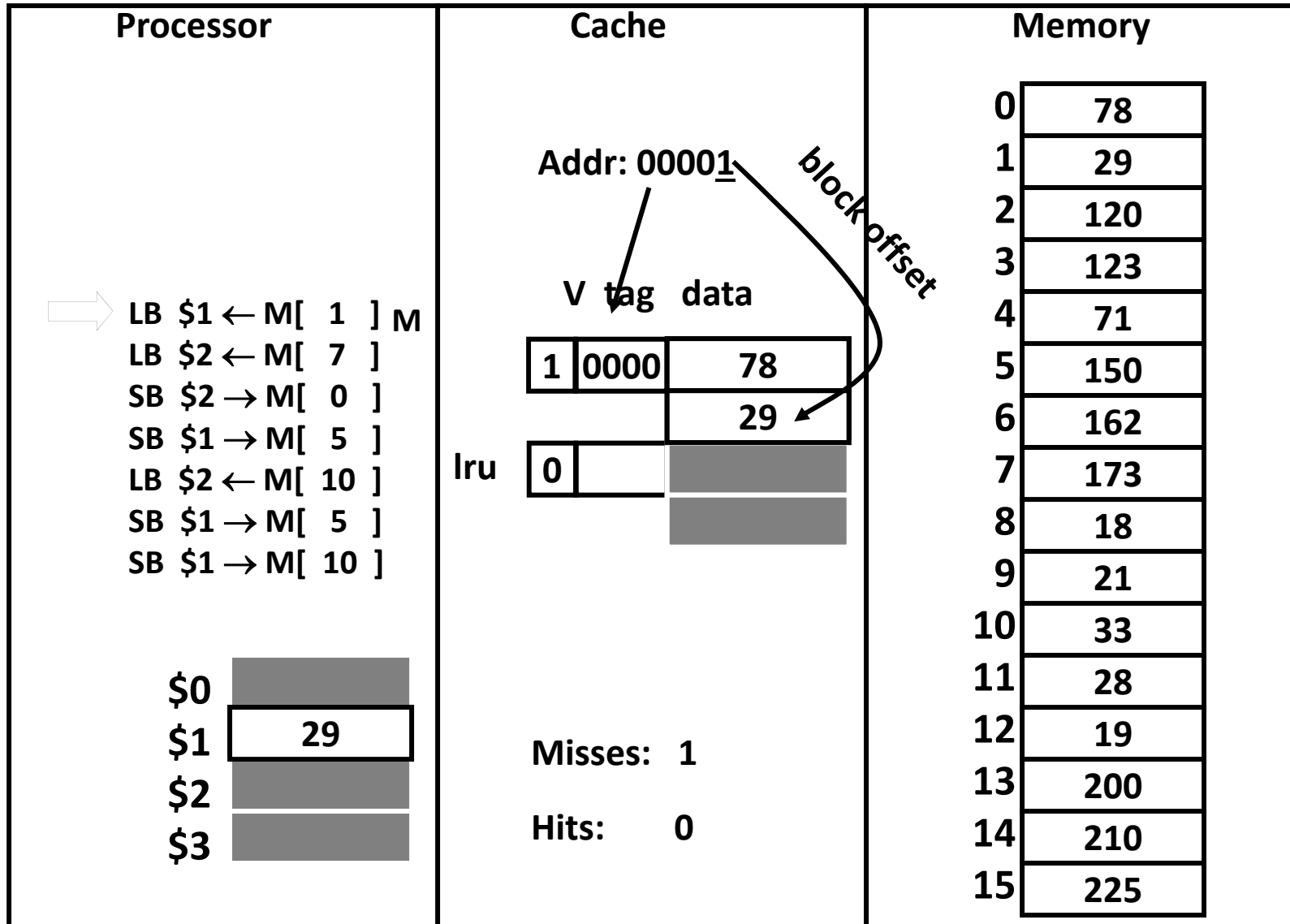
Using **byte addresses** in this example! Addr Bus = 5 bits



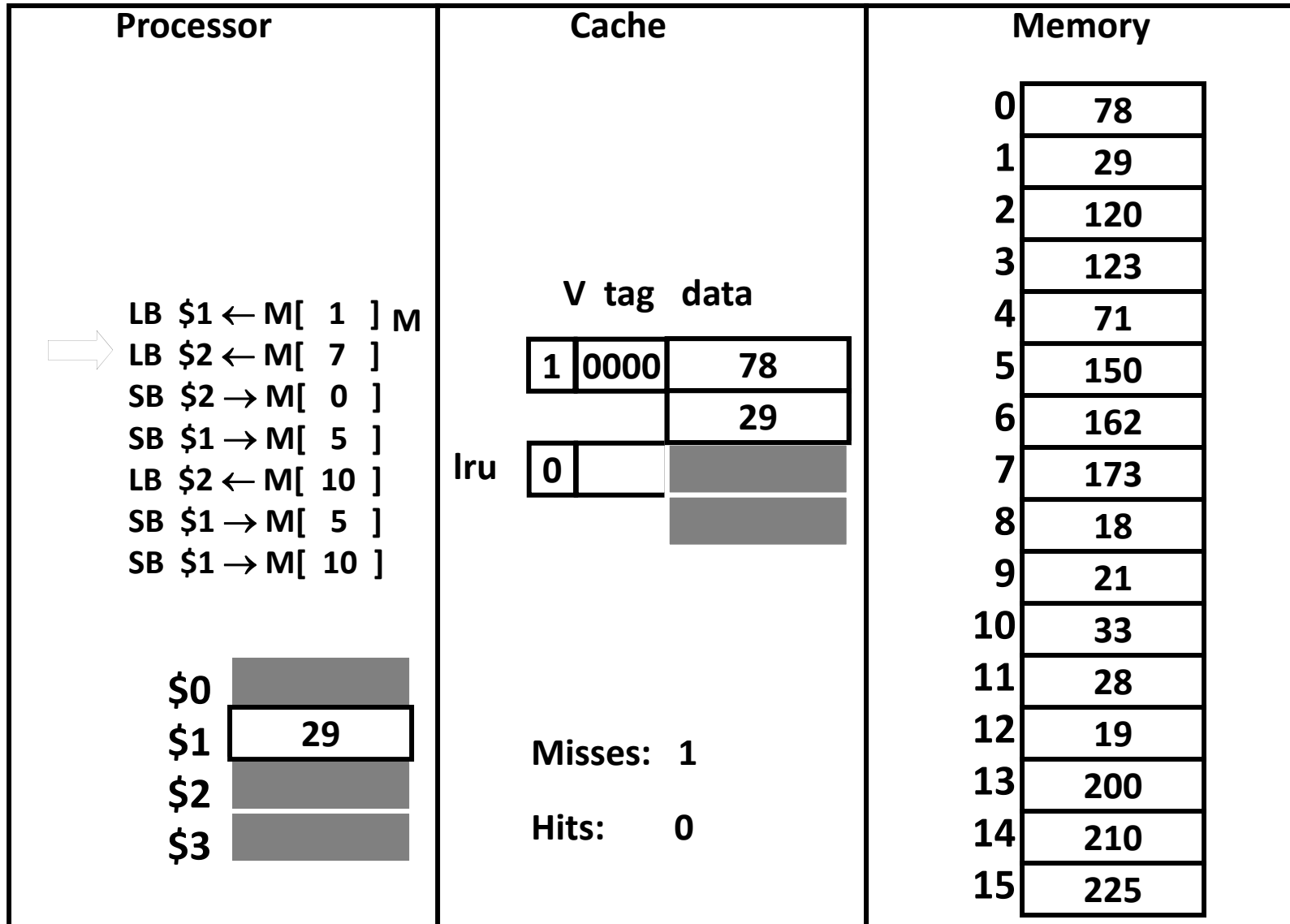
Write-Through (REF 1)



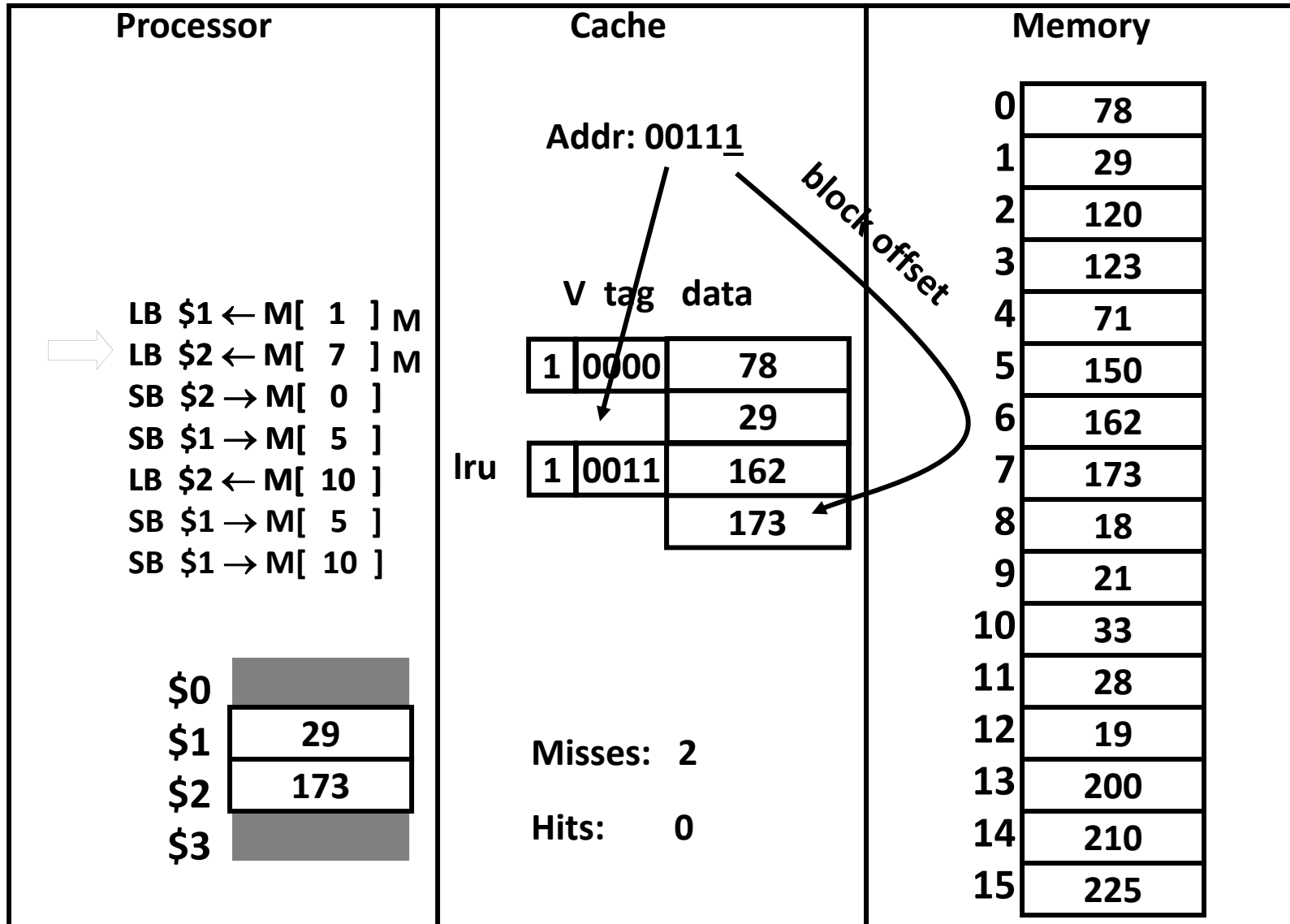
Write-Through (REF 1)



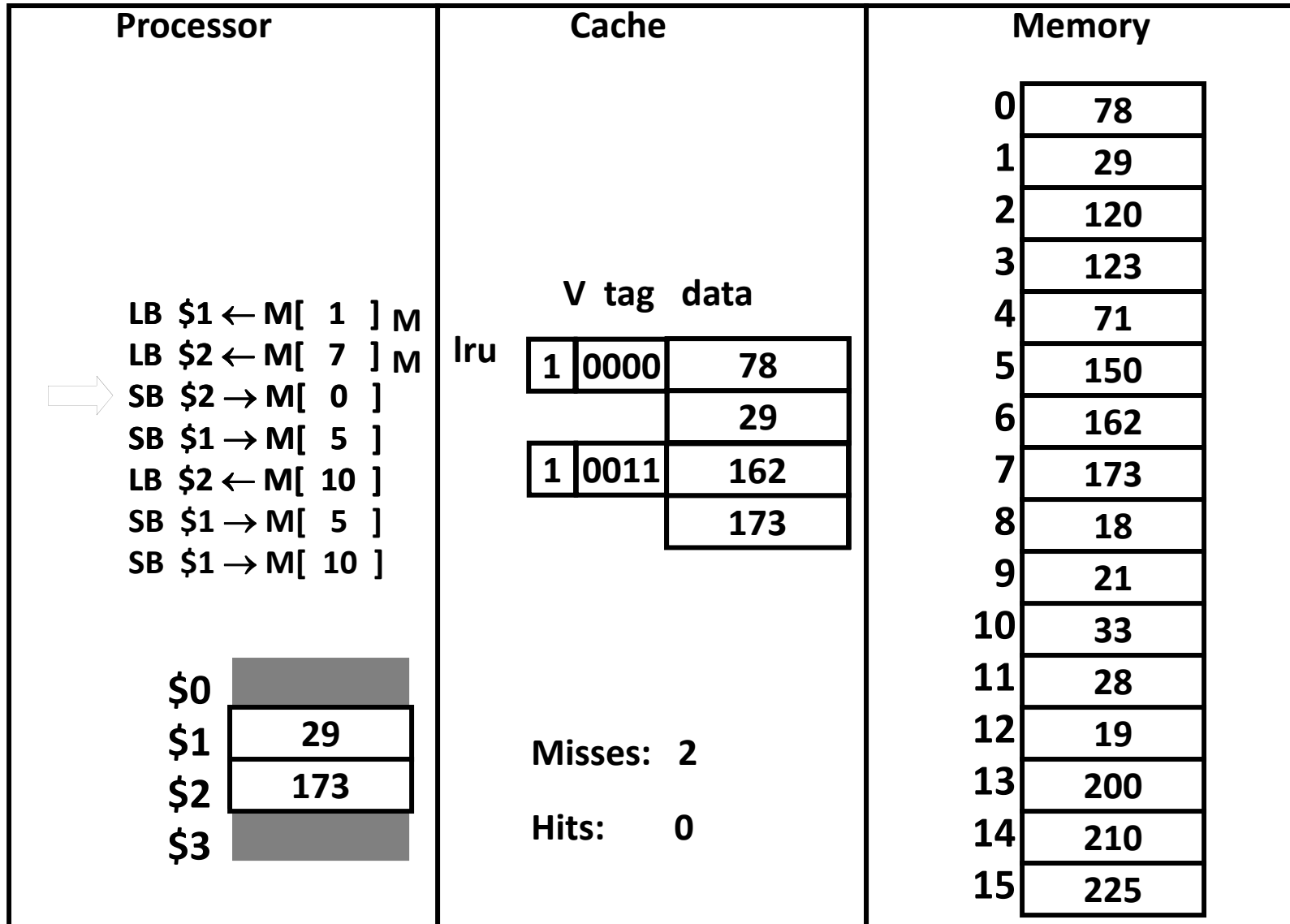
Write-Through (REF 2)



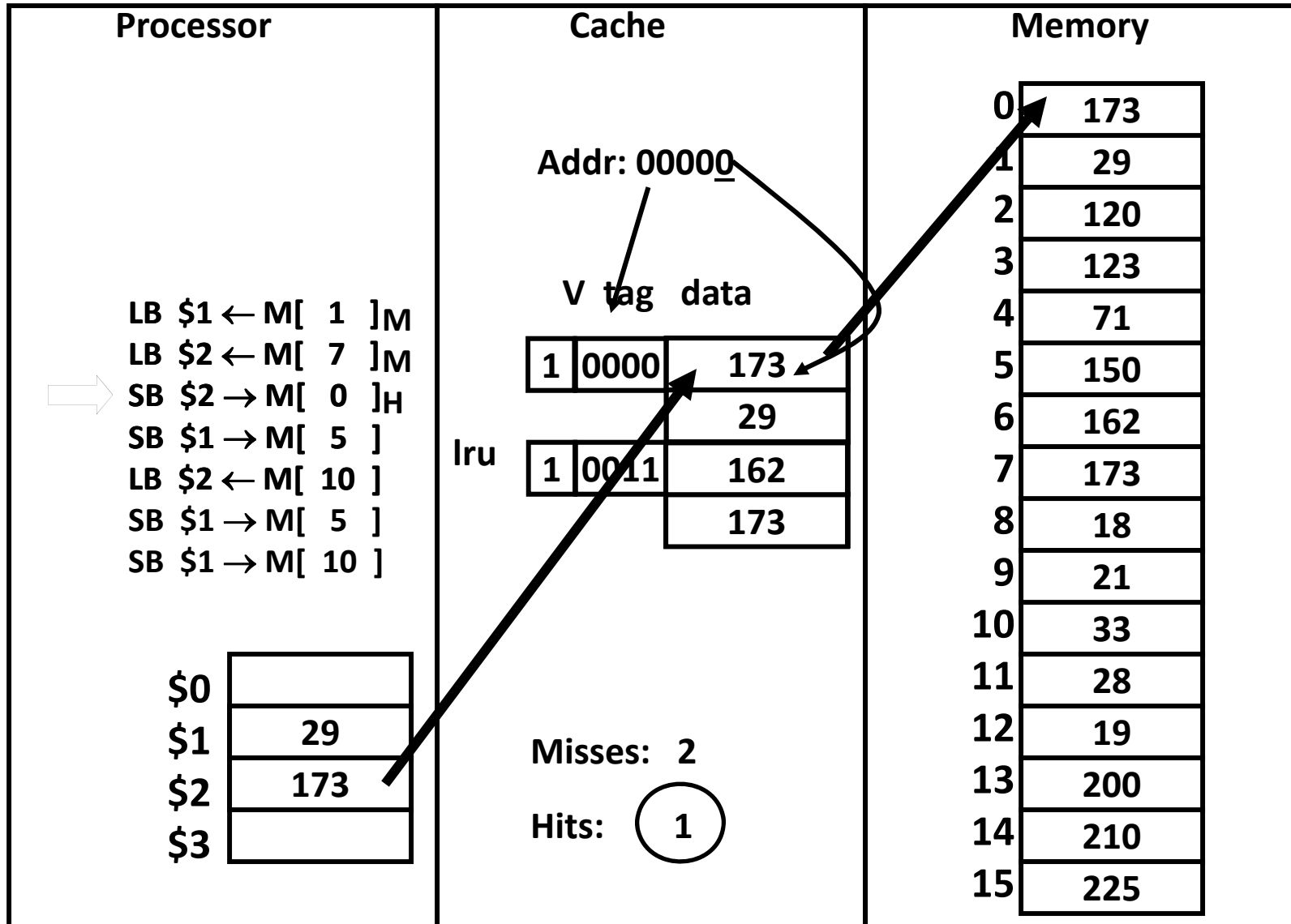
Write-Through (REF 2)



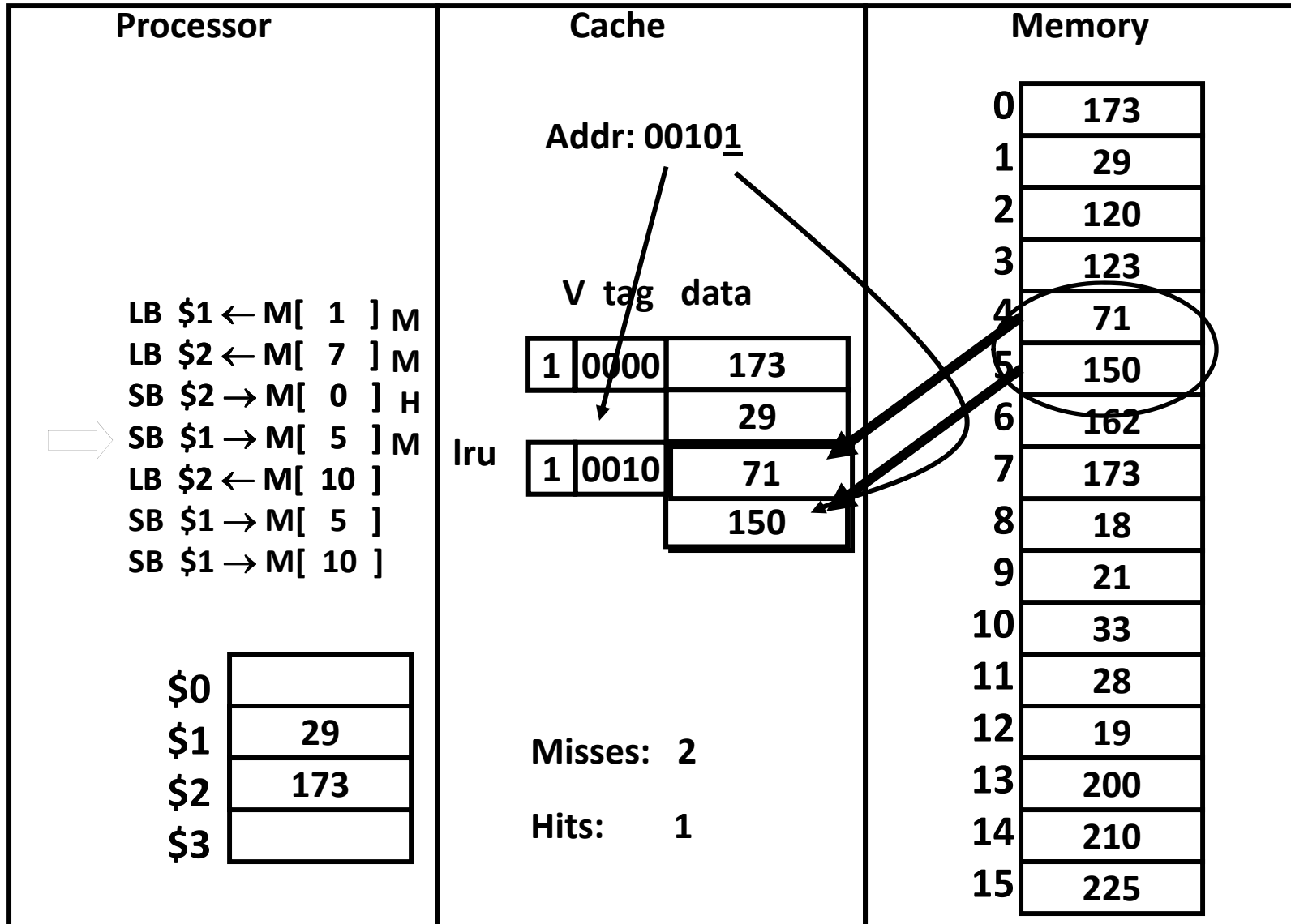
Write-Through (REF 3)



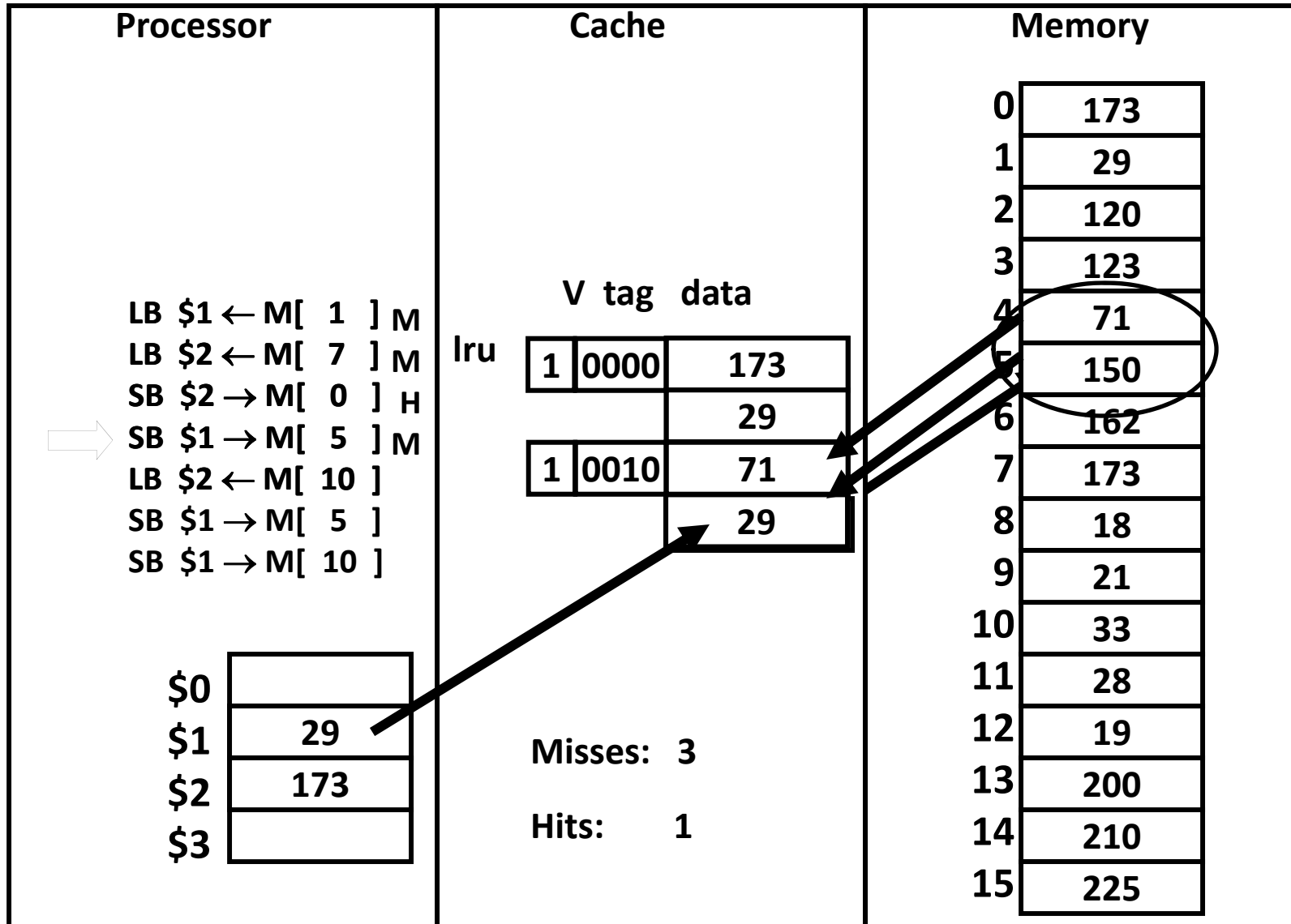
Write-Through (REF 3)



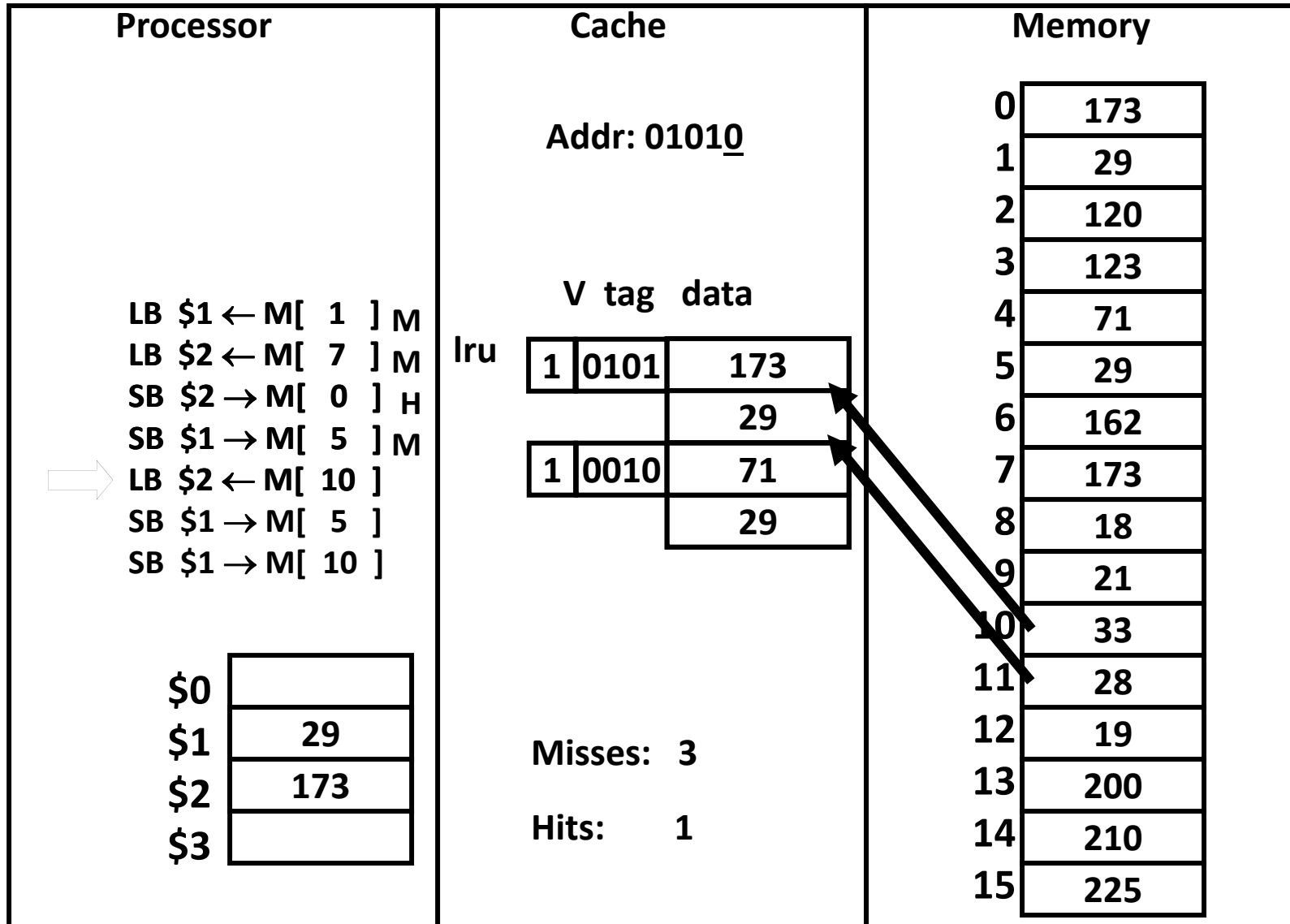
Write-Through (REF 4)



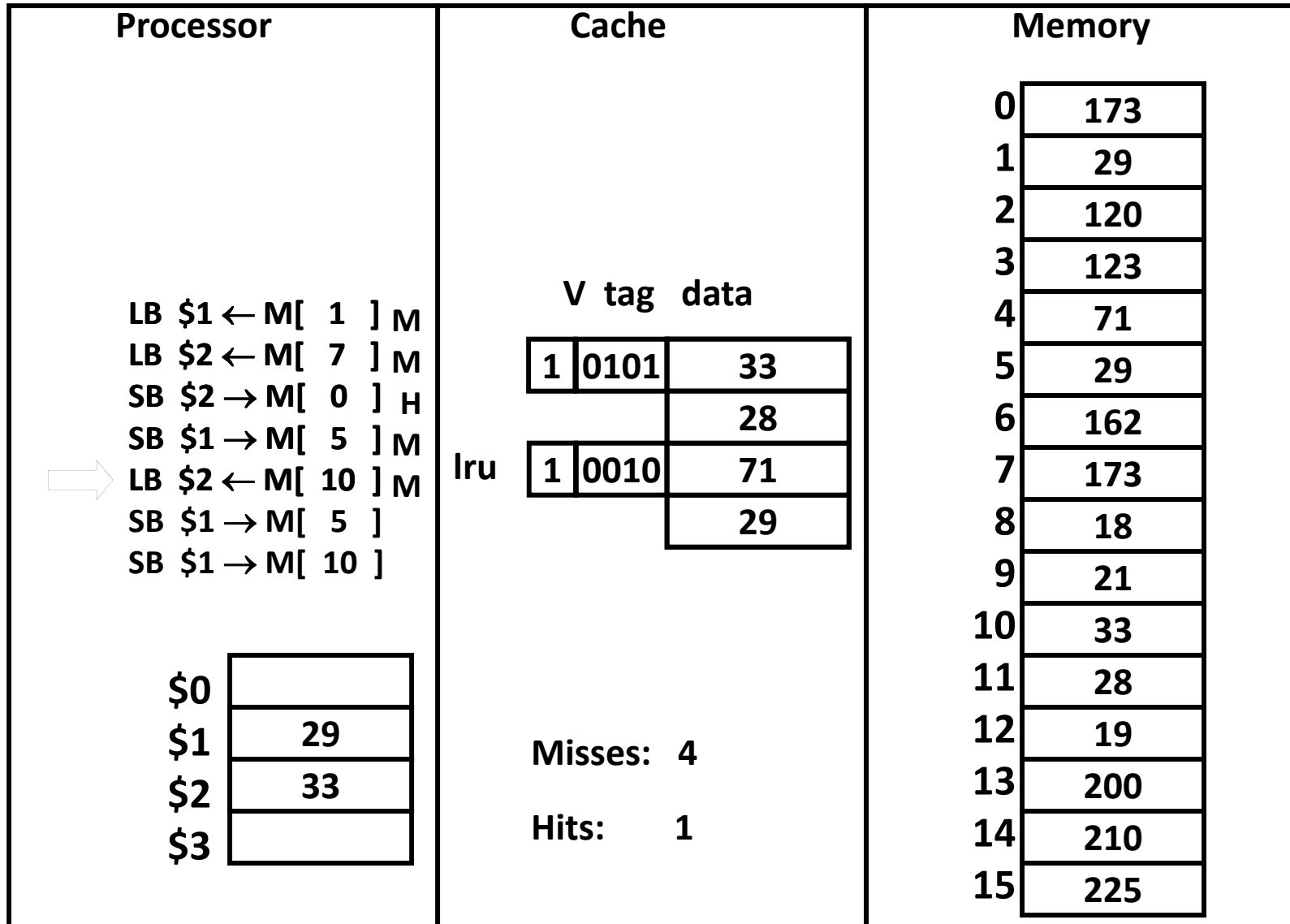
Write-Through (REF 4)



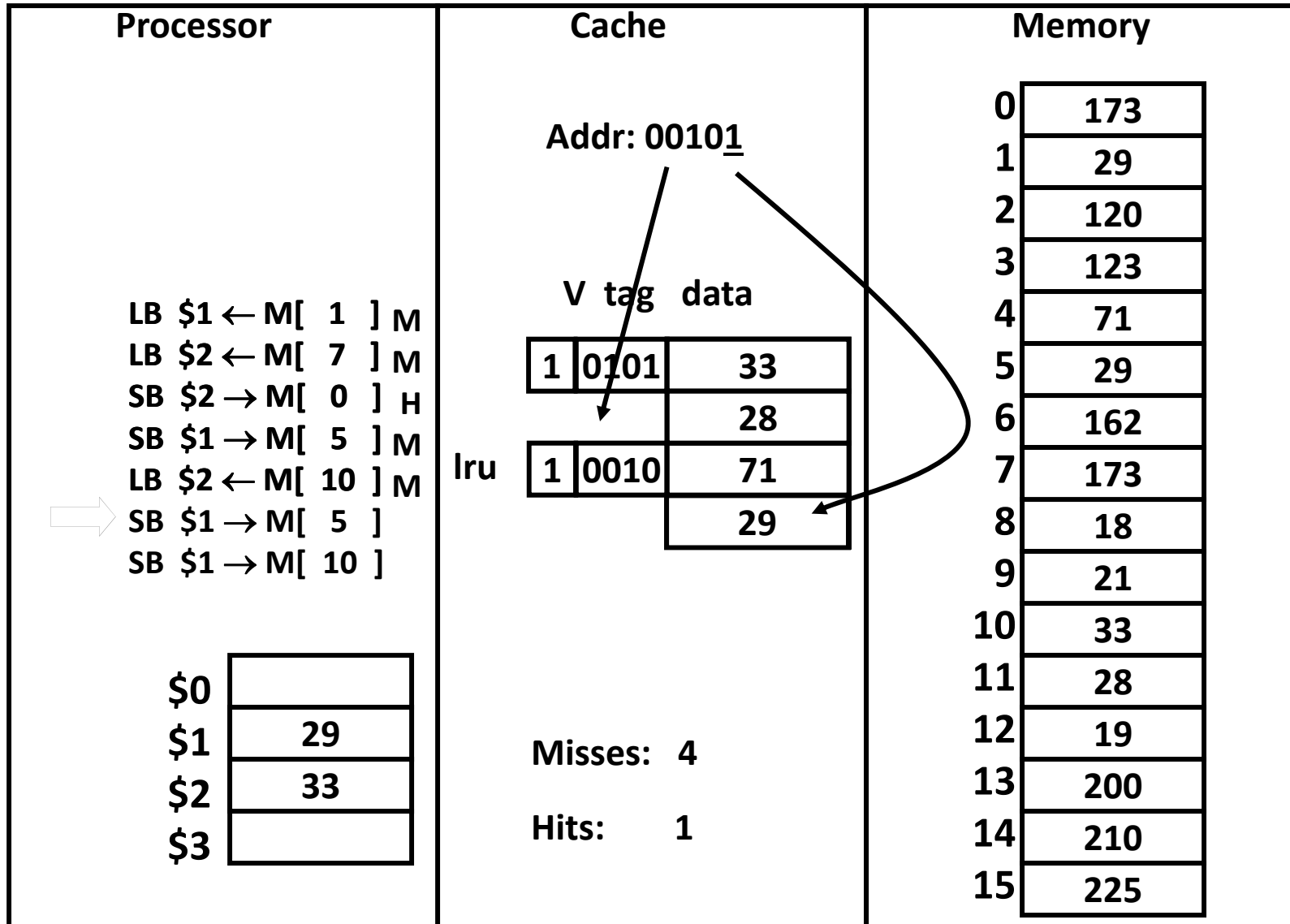
Write-Through (REF 5)



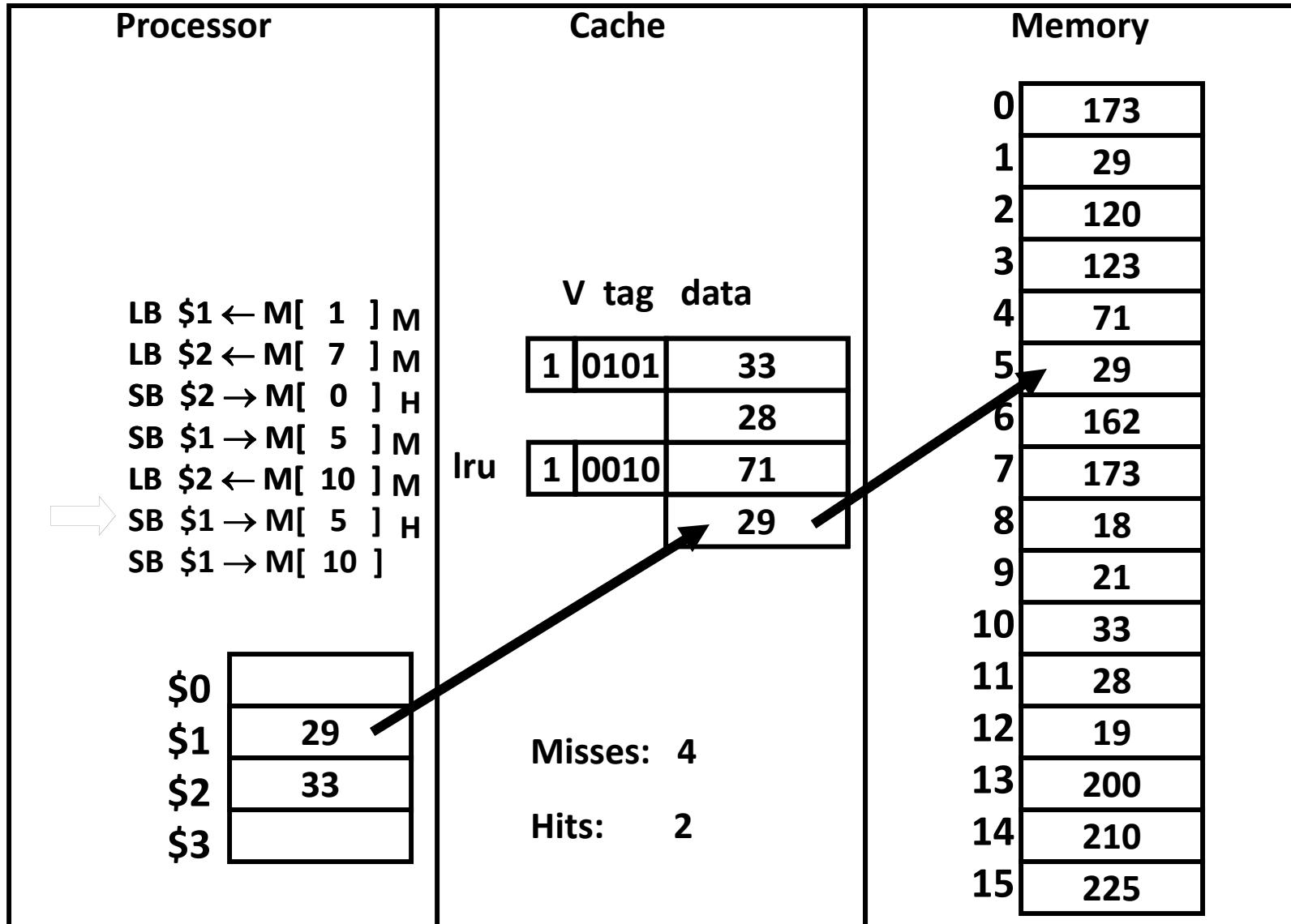
Write-Through (REF 5)



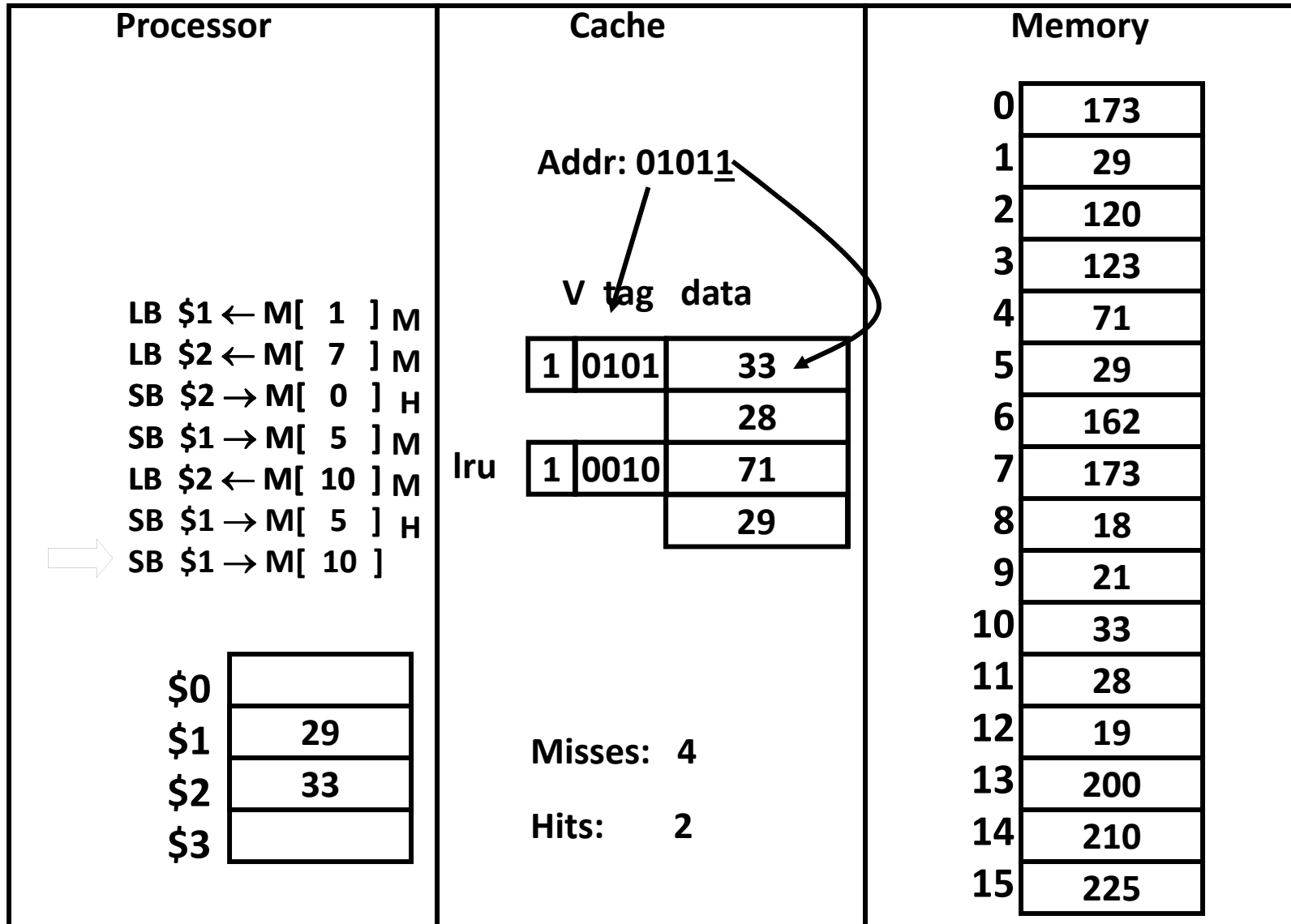
Write-Through (REF 6)



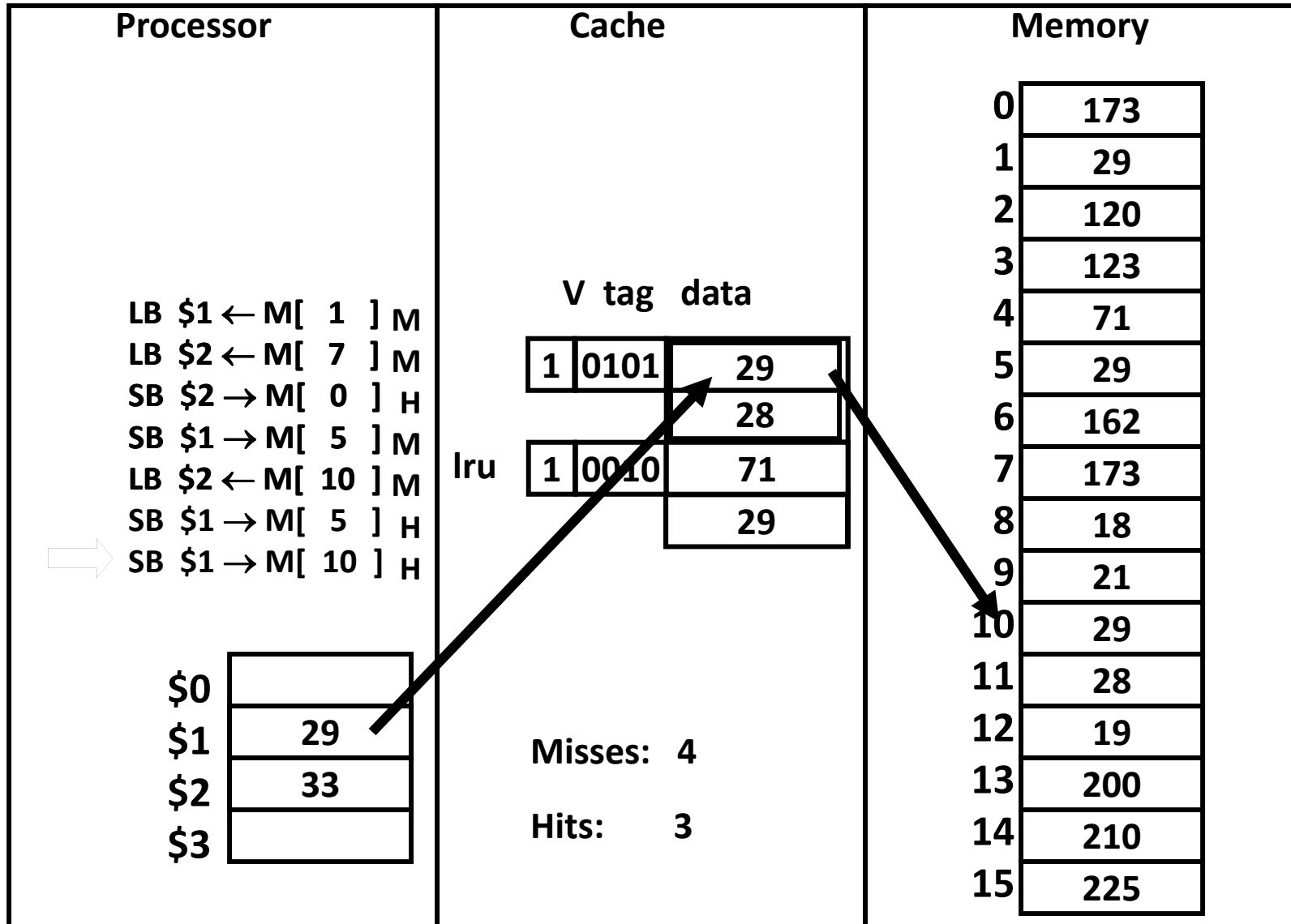
Write-Through (REF 6)



Write-Through (REF 7)



Write-Through (REF 7)



How Many Memory References?

Write-through performance

Each miss (read or write) reads a block from mem

- 4 misses → 8 mem reads

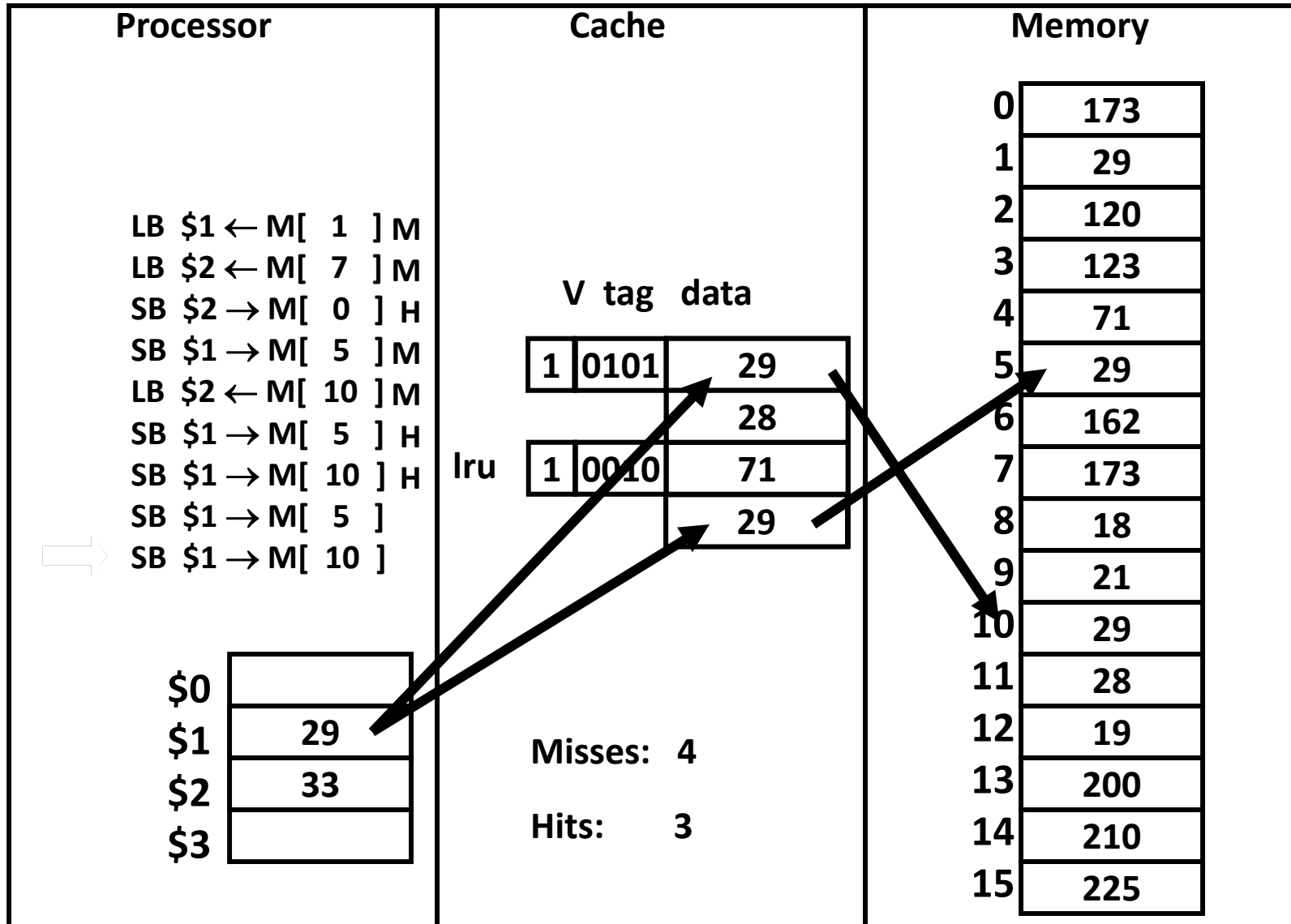
Each store writes an item to mem

- 4 mem writes

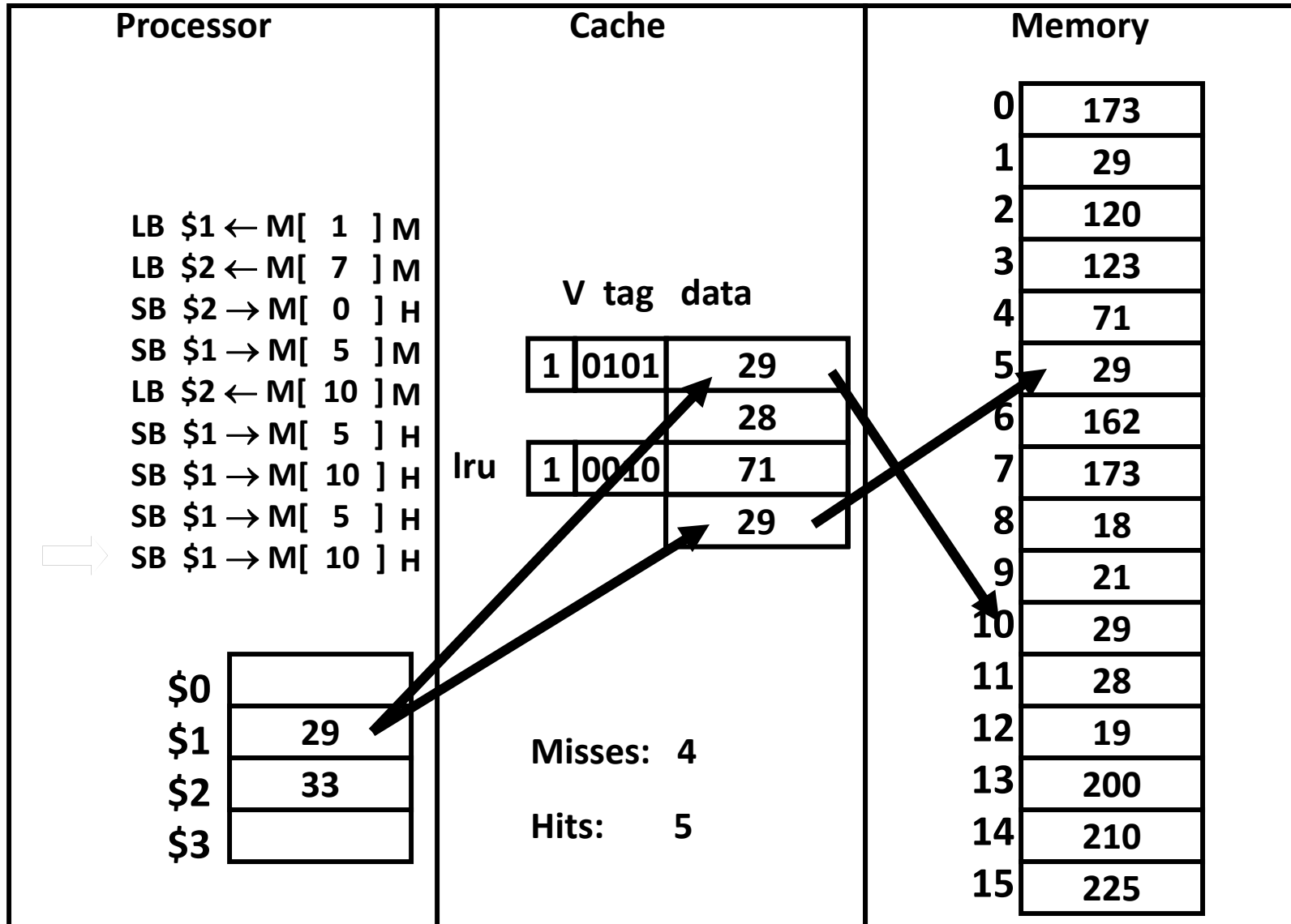
Evictions don't need to write to mem

- no need for dirty bit

Write-Through (REF 8,9)



Write-Through (REF 8,9)



Summary: Write Through

Write-through policy with write allocate

Cache miss: read entire block from memory

Write: write only updated item to memory

Eviction: no need to write to memory

Next Goal: Write-Through vs. Write-Back

Can we also design the cache NOT to write all stores immediately to memory?

- Keep the most current copy in cache, and update memory when that data is **evicted** (write-back policy)
- Do we need to write-back all evicted lines?
 - No, only blocks that have been stored into (written)

Write-Back Meta-Data

V	D	Tag	Byte 1	Byte 2	... Byte N

V = 1 means the line has valid data

D = 1 means the bytes are newer than main memory

When allocating line:

- Set V = 1, D = 0, fill in Tag and Data

When writing line:

- Set D = 1

When evicting line:

- If D = 0: just set V = 0
- If D = 1: write-back Data, then set D = 0, V = 0

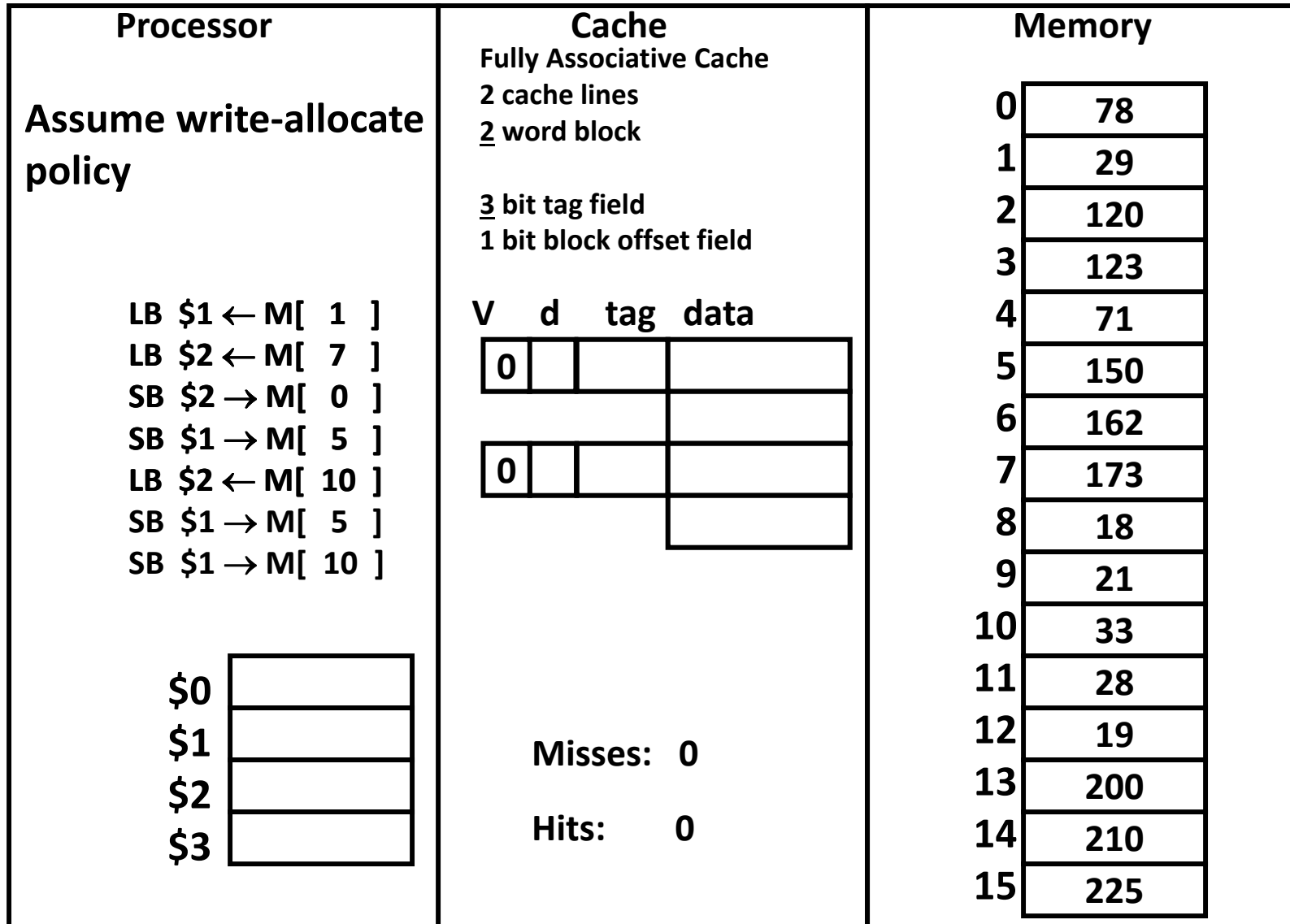
Write-back Example

Example: How does a write-back cache work?

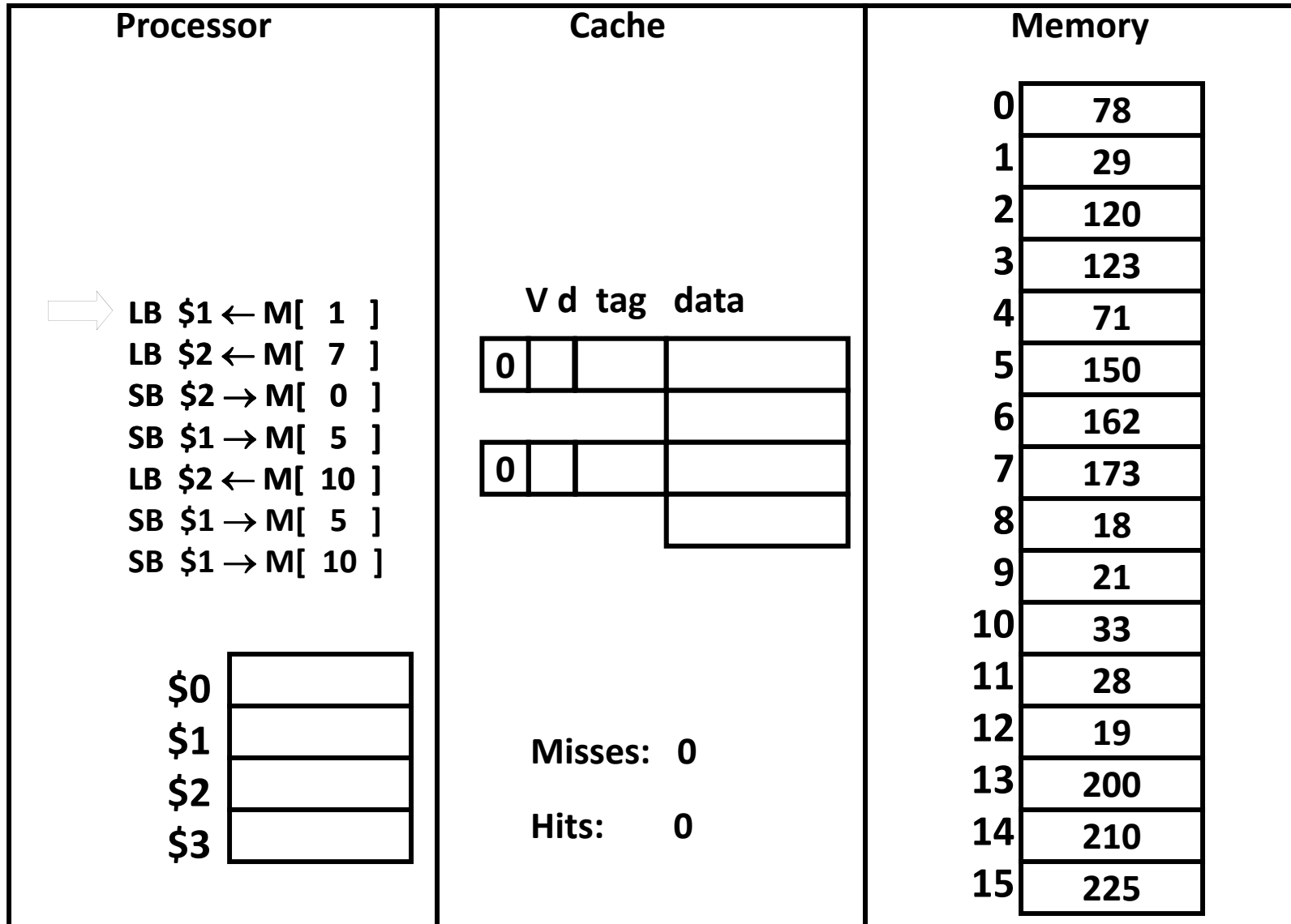
Assume write-allocate

Handling Stores (Write-Back)

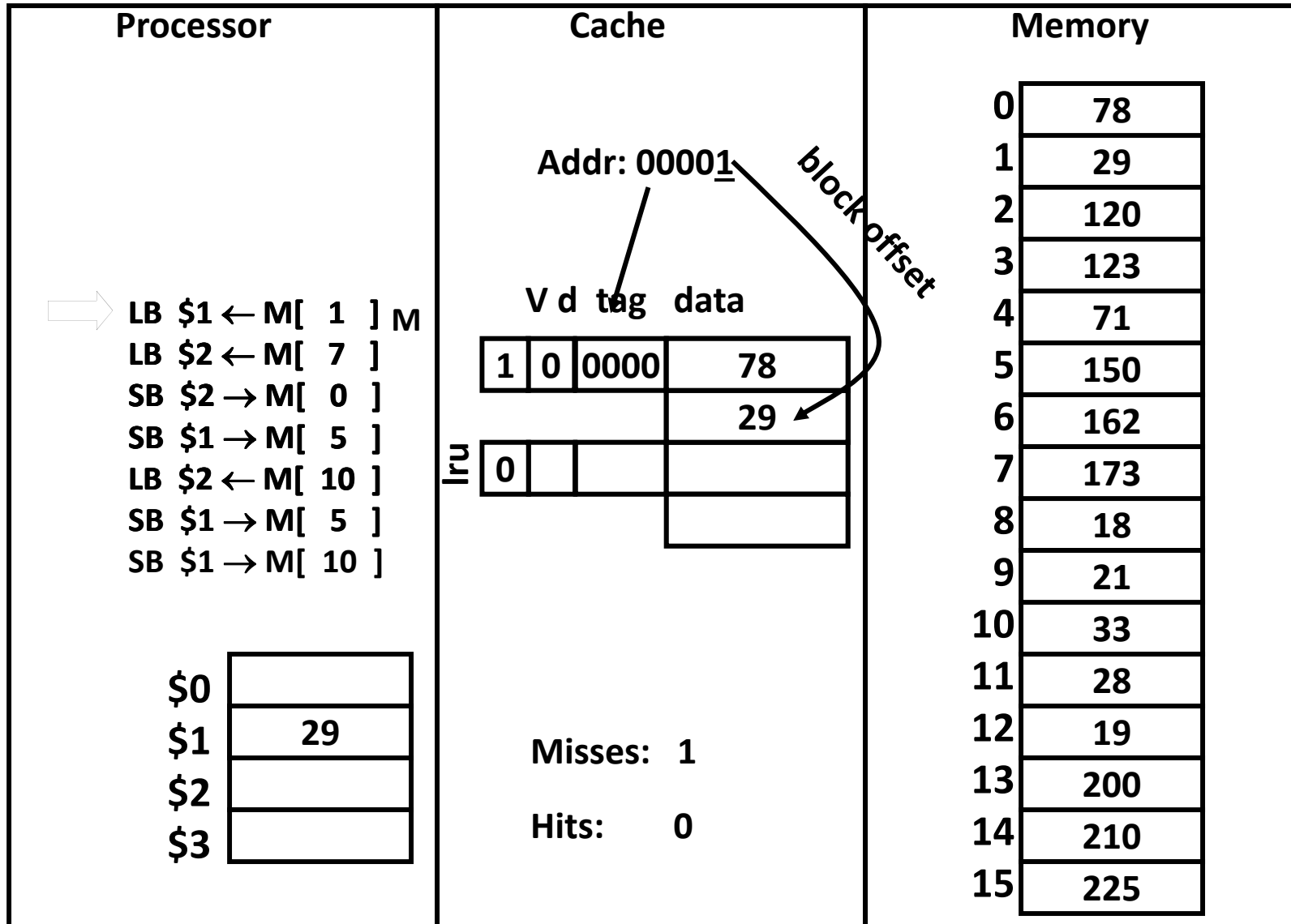
Using **byte addresses** in this example! Addr Bus = 5 bits



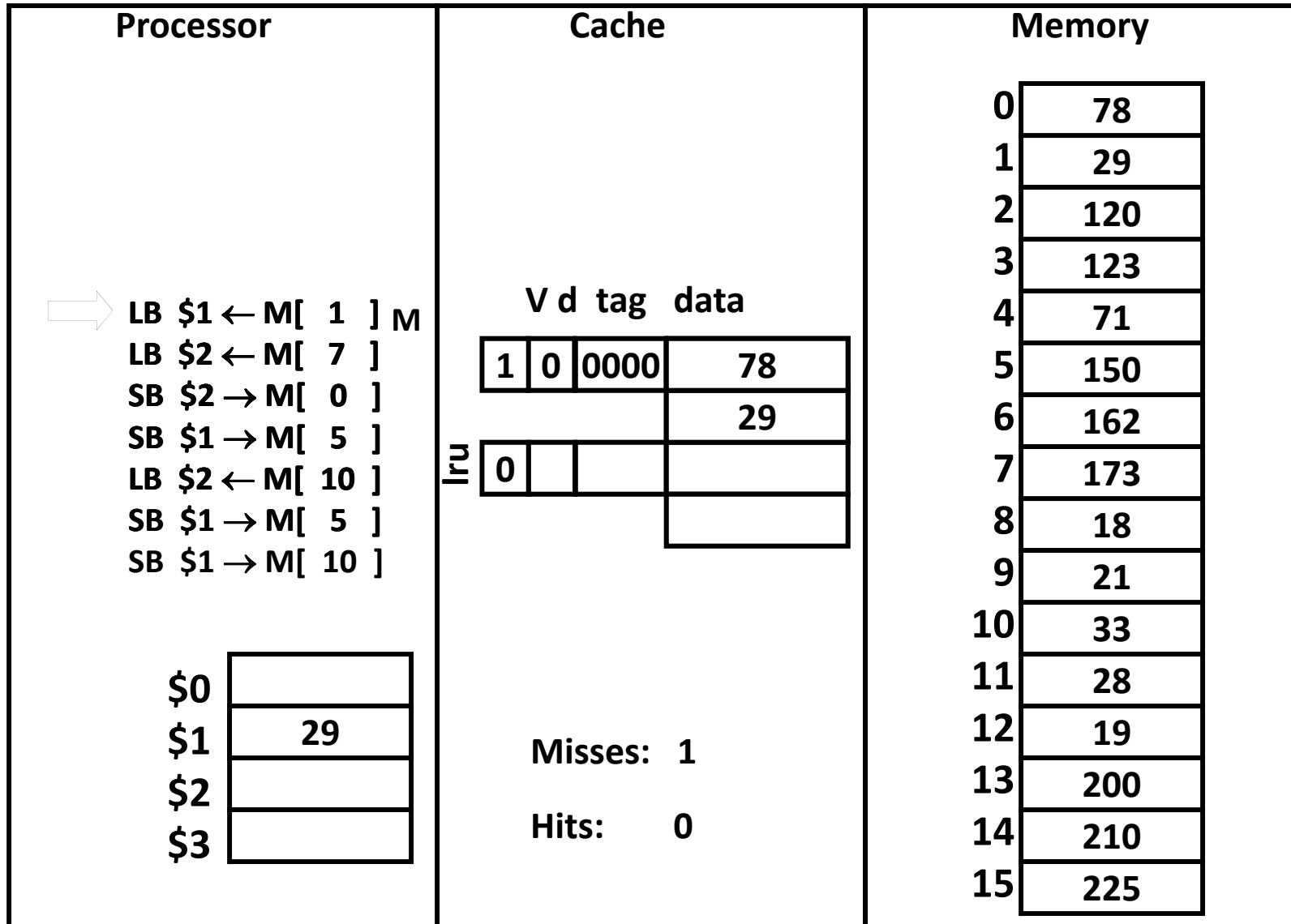
Write-Back (REF 1)



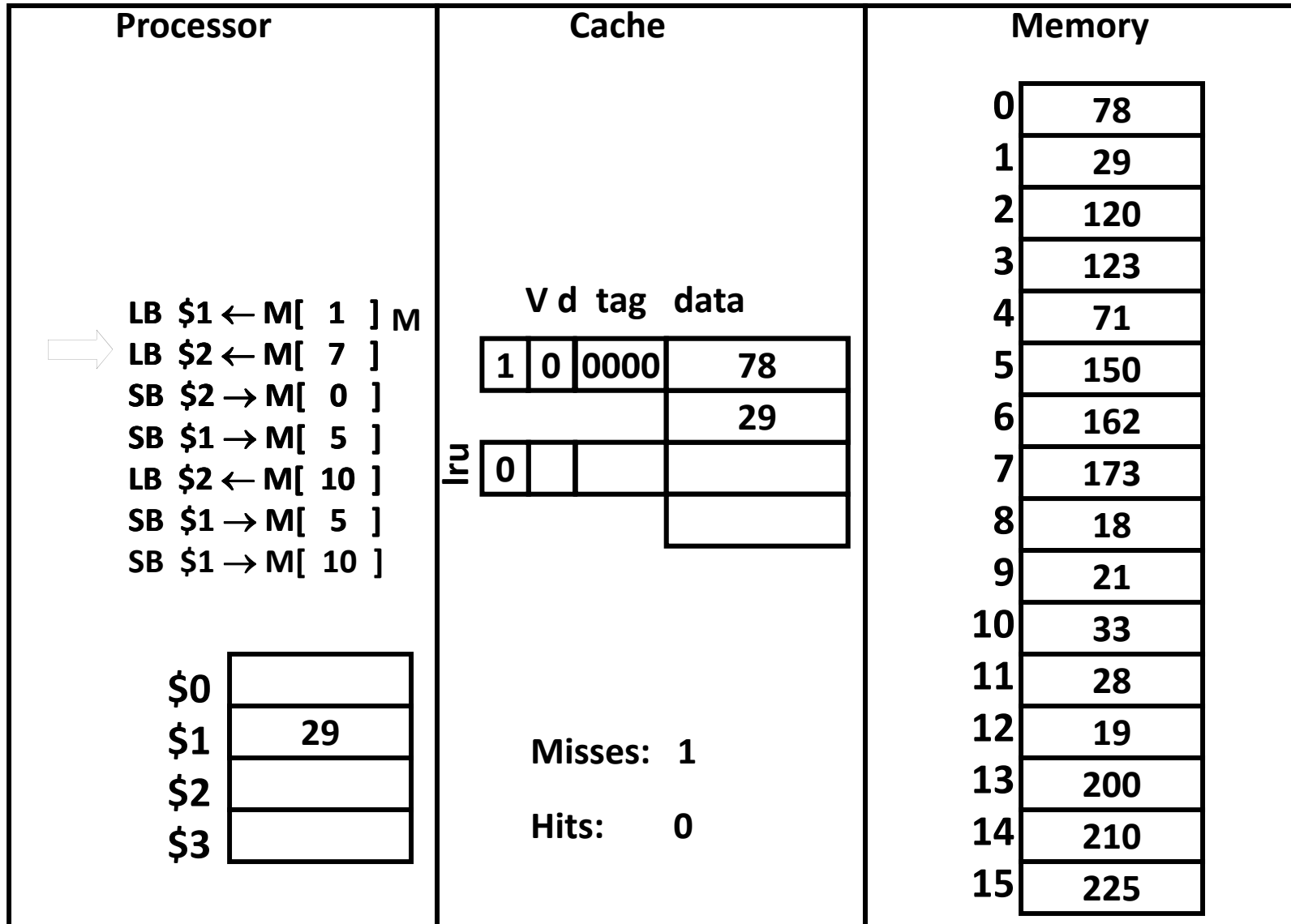
Write-Back (REF 1)



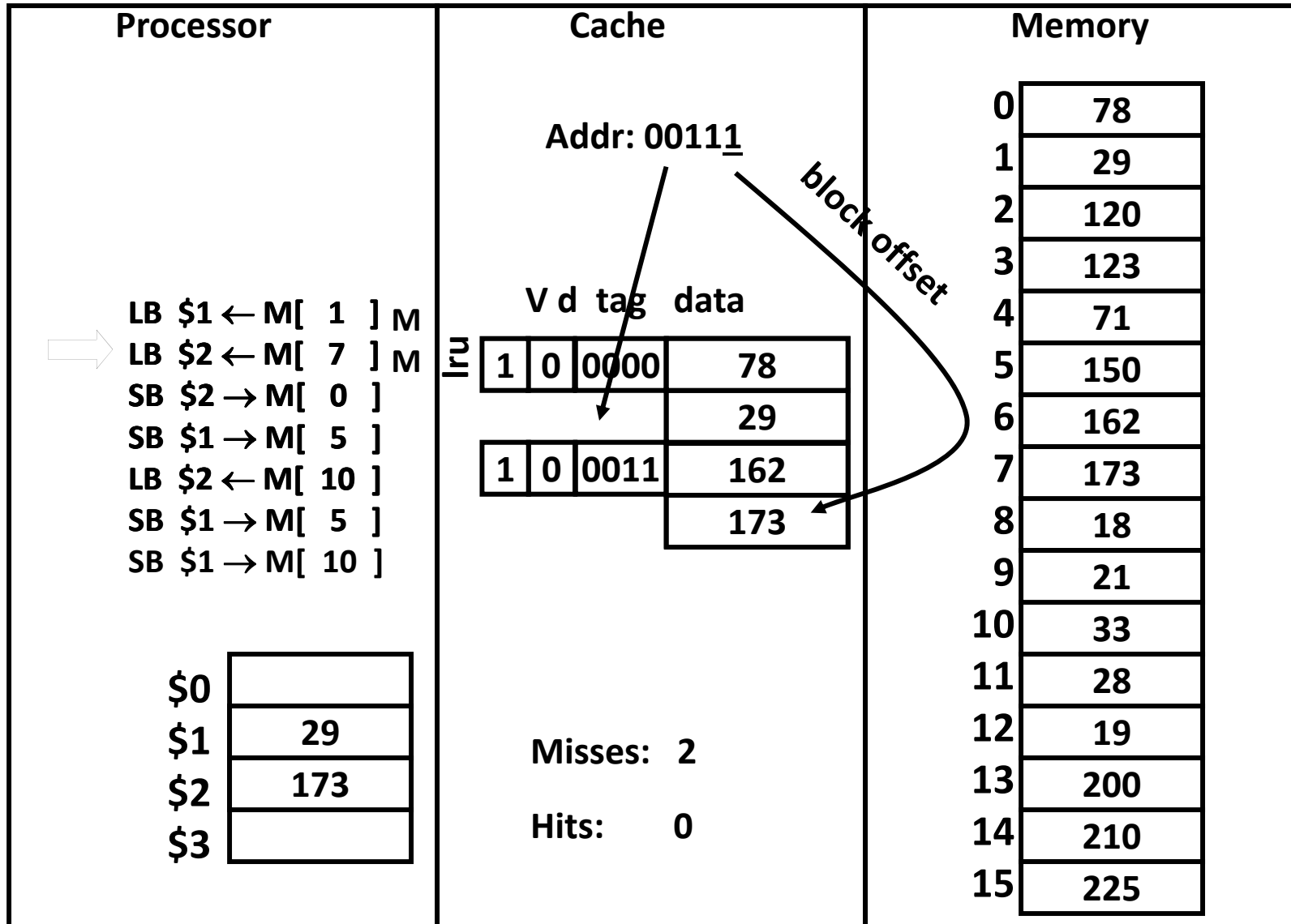
Write-Back (REF 1)



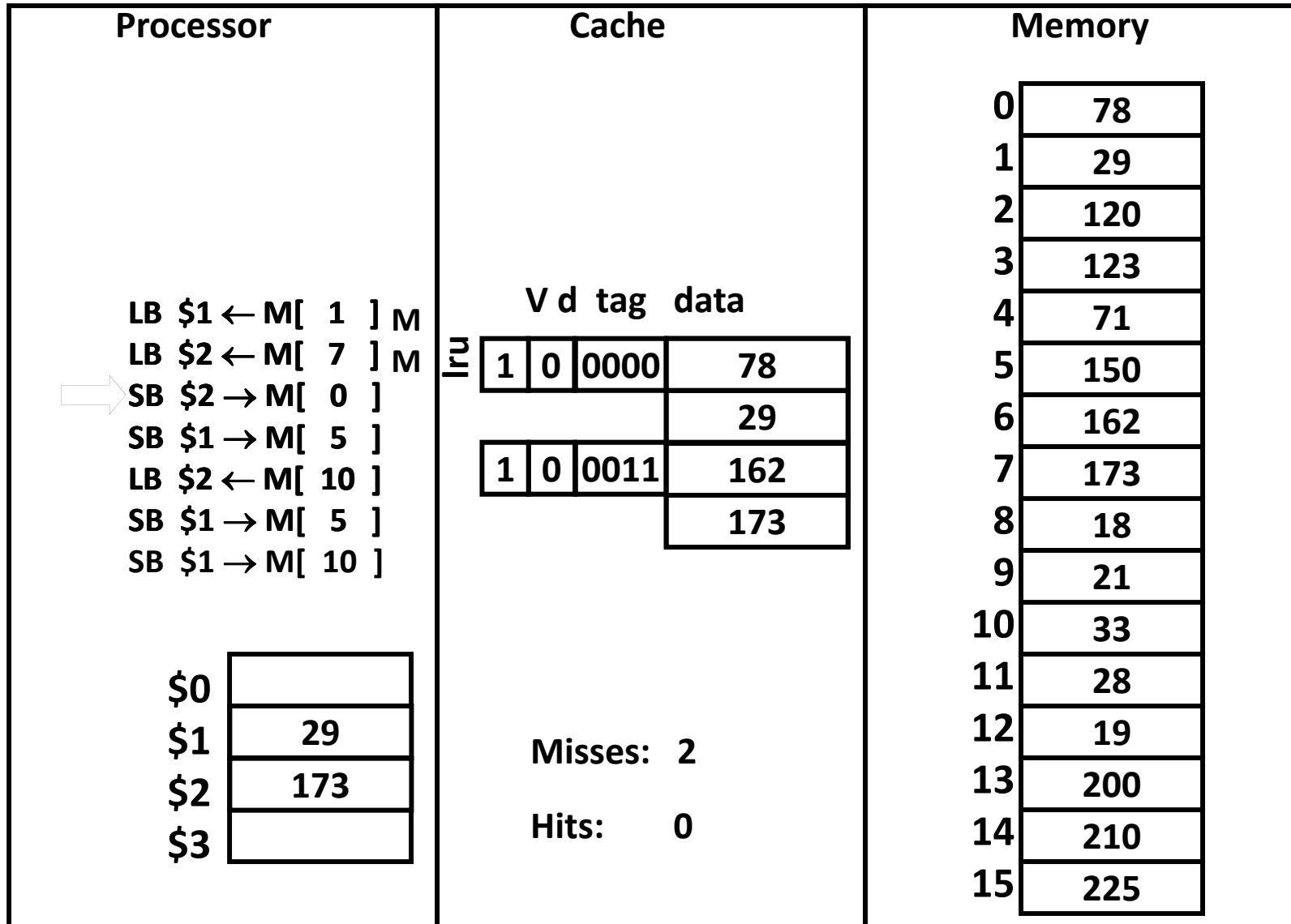
Write-Back (REF 2)



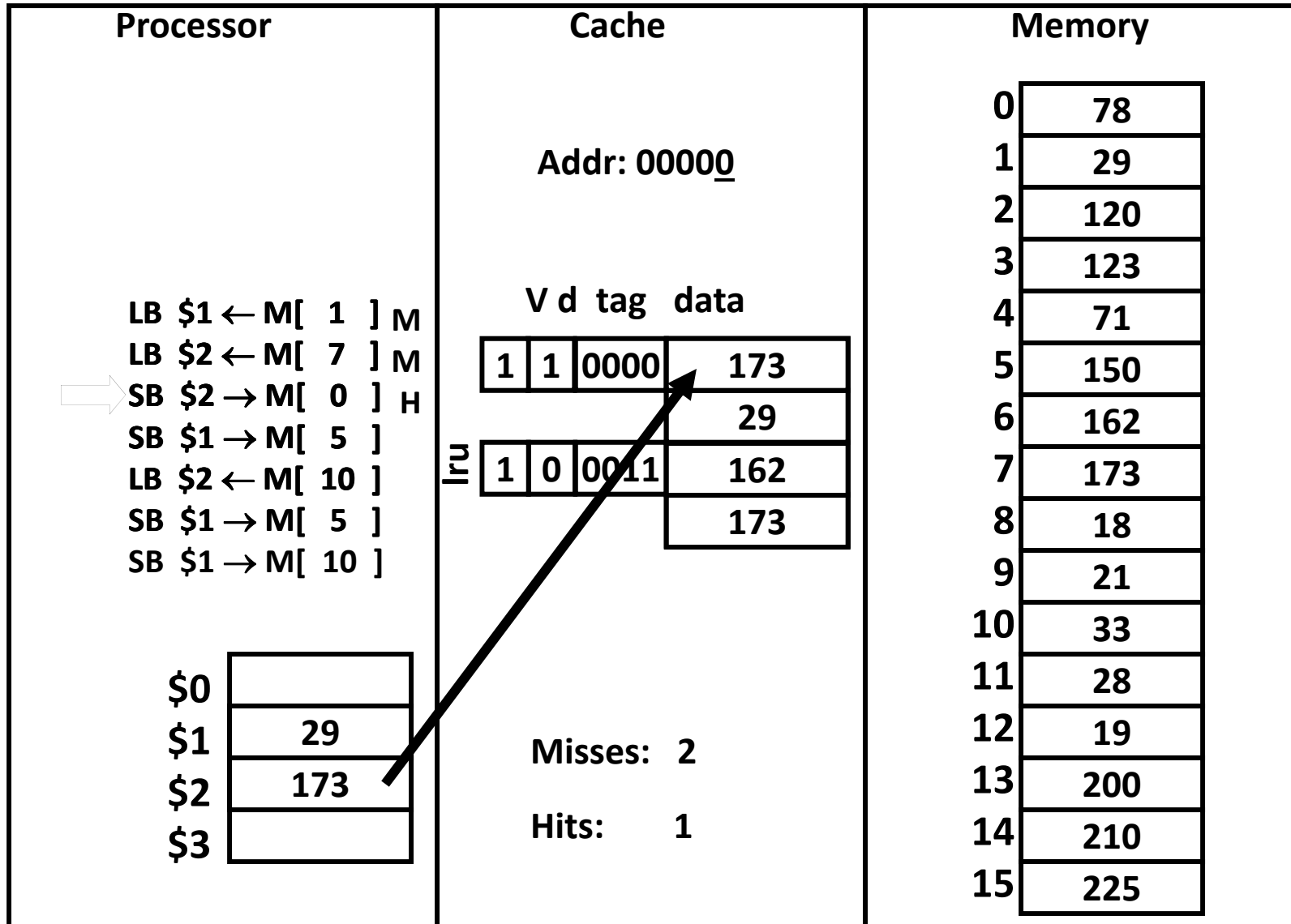
Write-Back (REF 2)



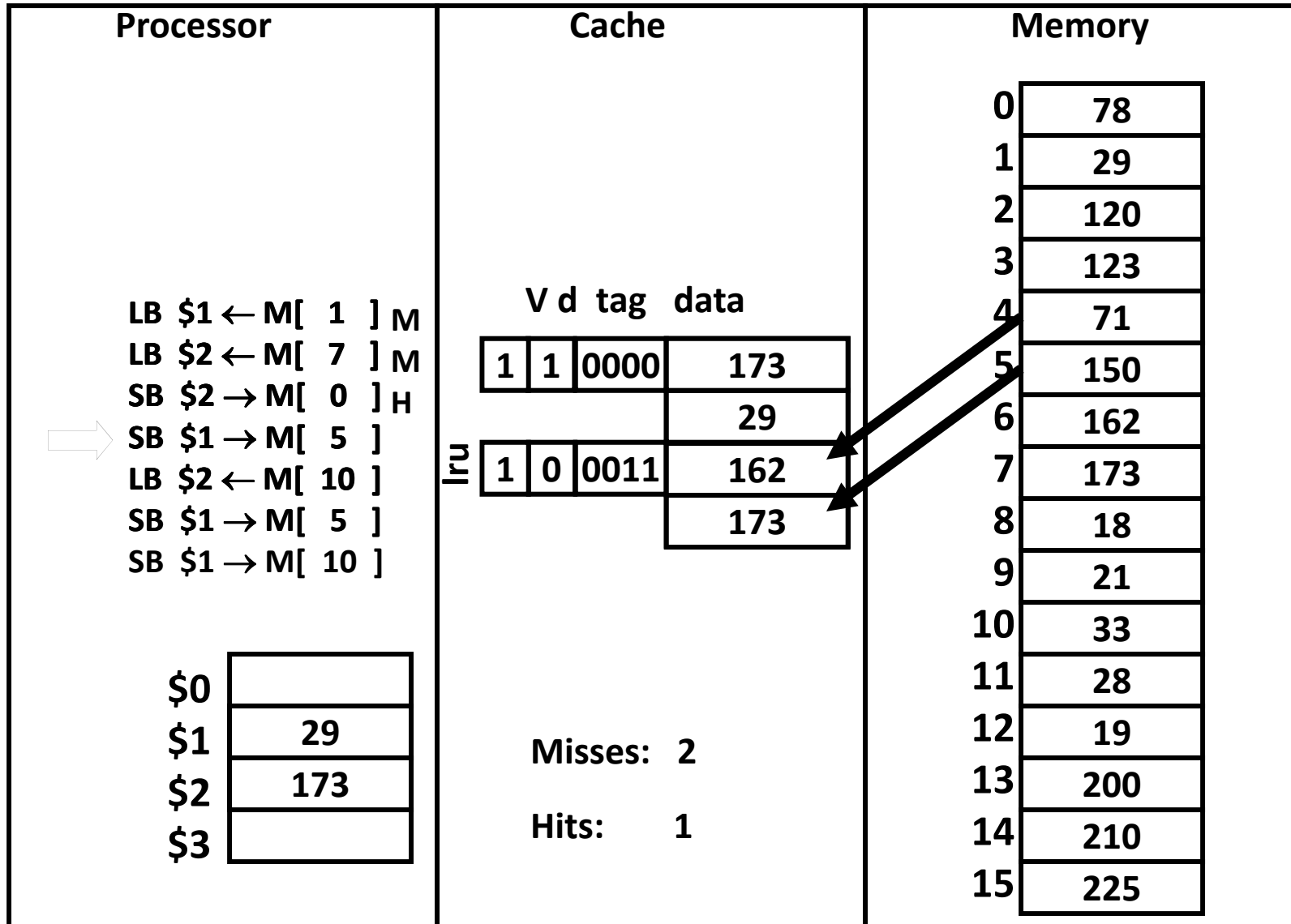
Write-Back (REF 3)



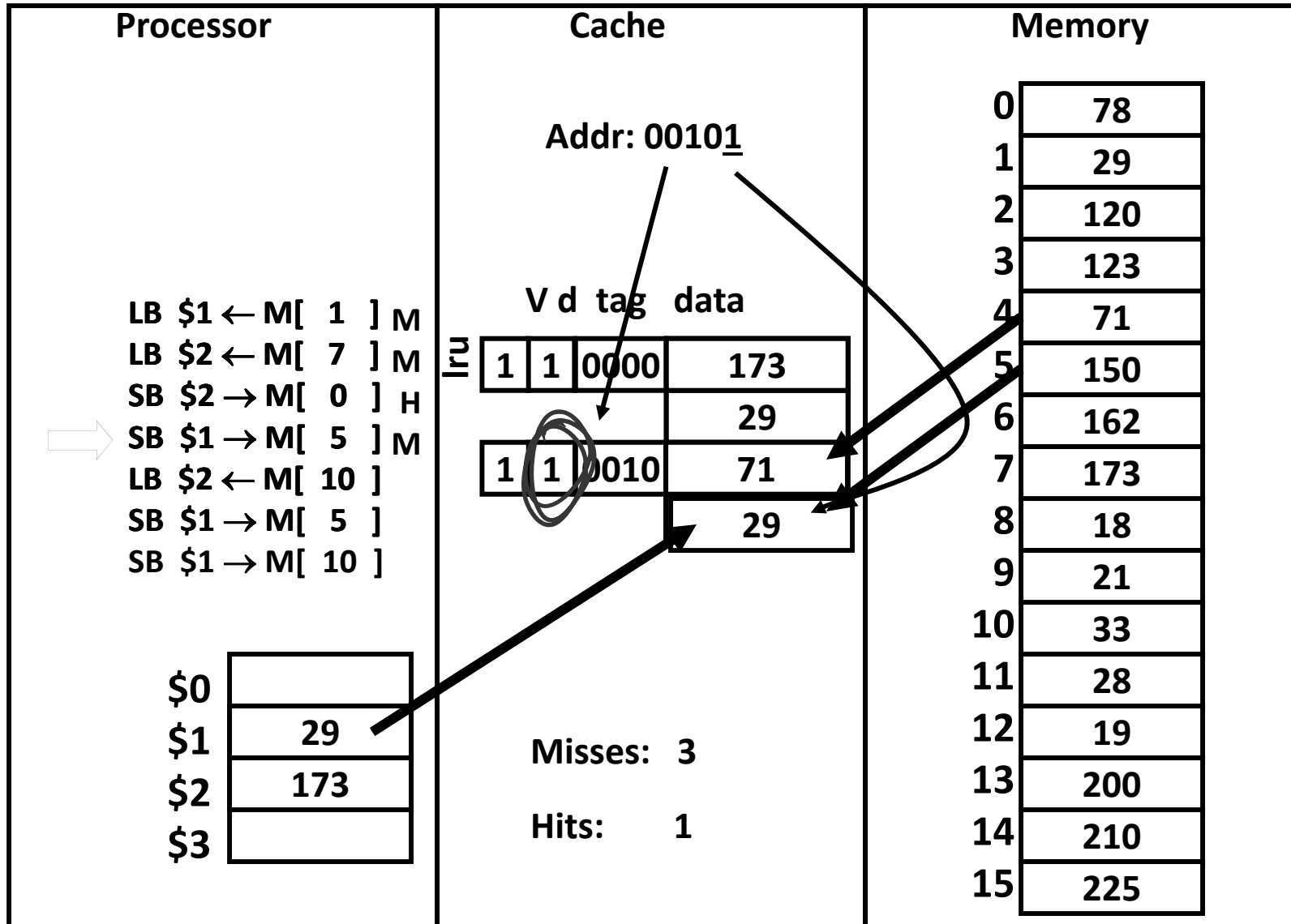
Write-Back (REF 3)



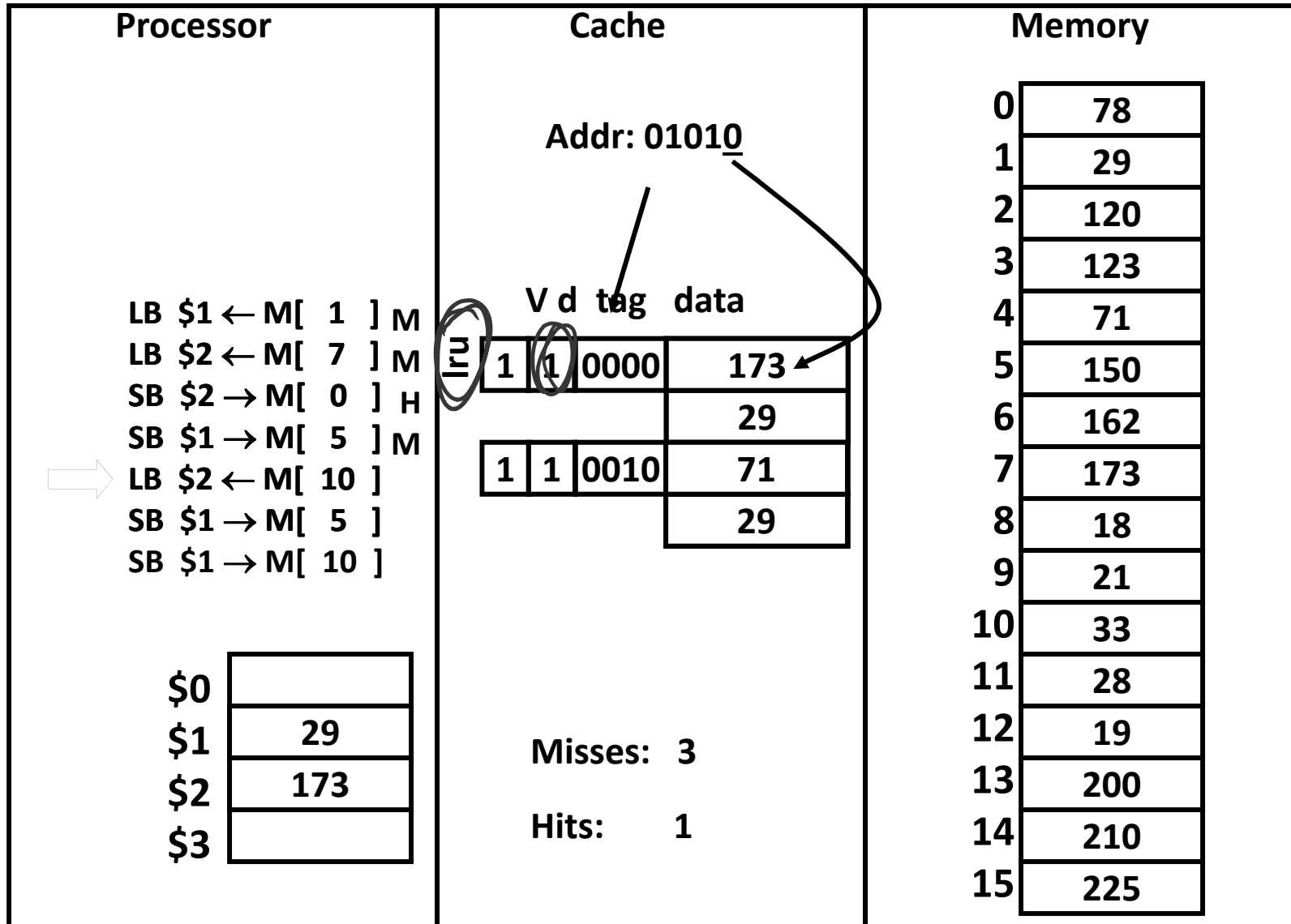
Write-Back (REF 4)



Write-Back (REF 4)

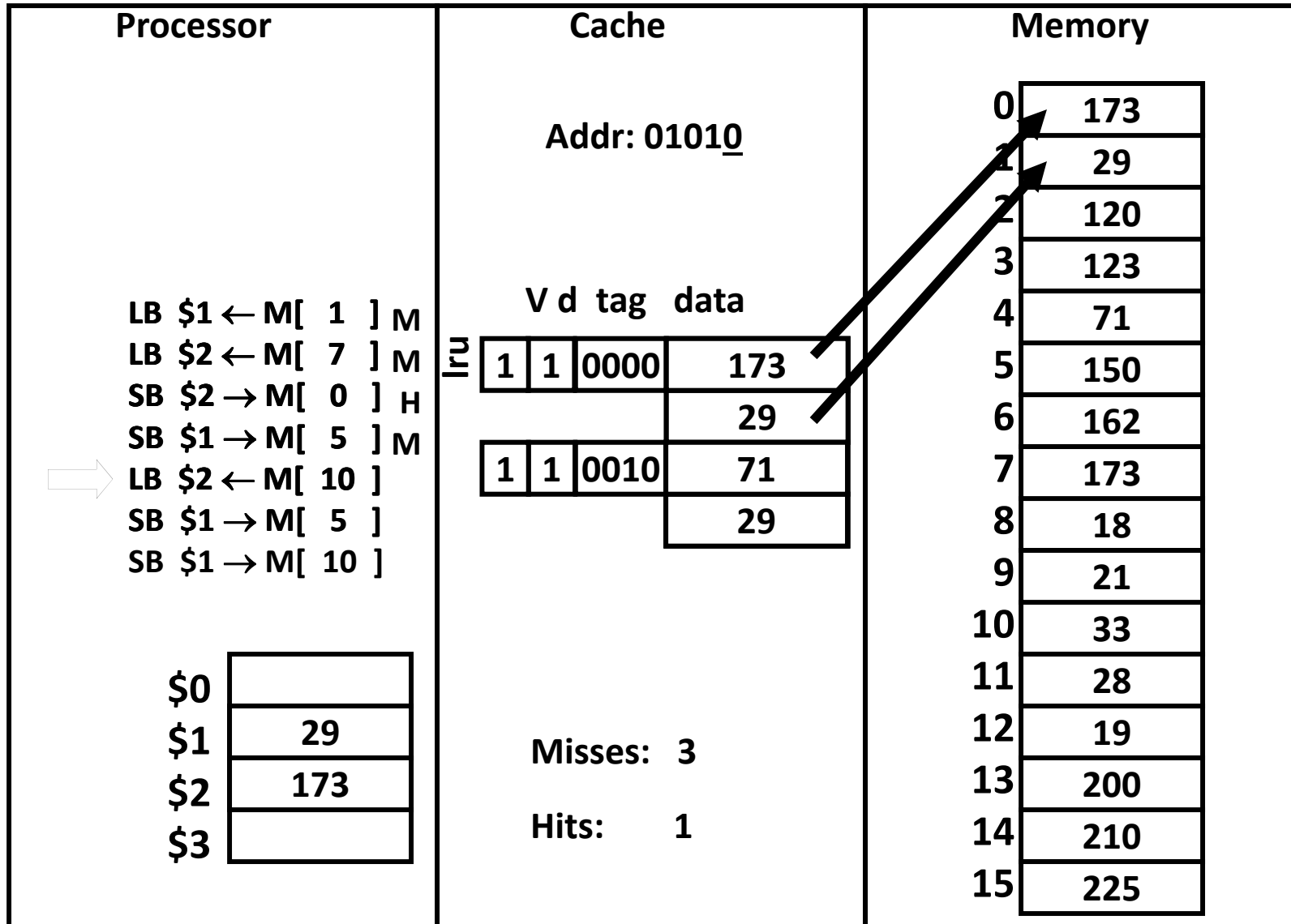


Write-Back (REF 5)

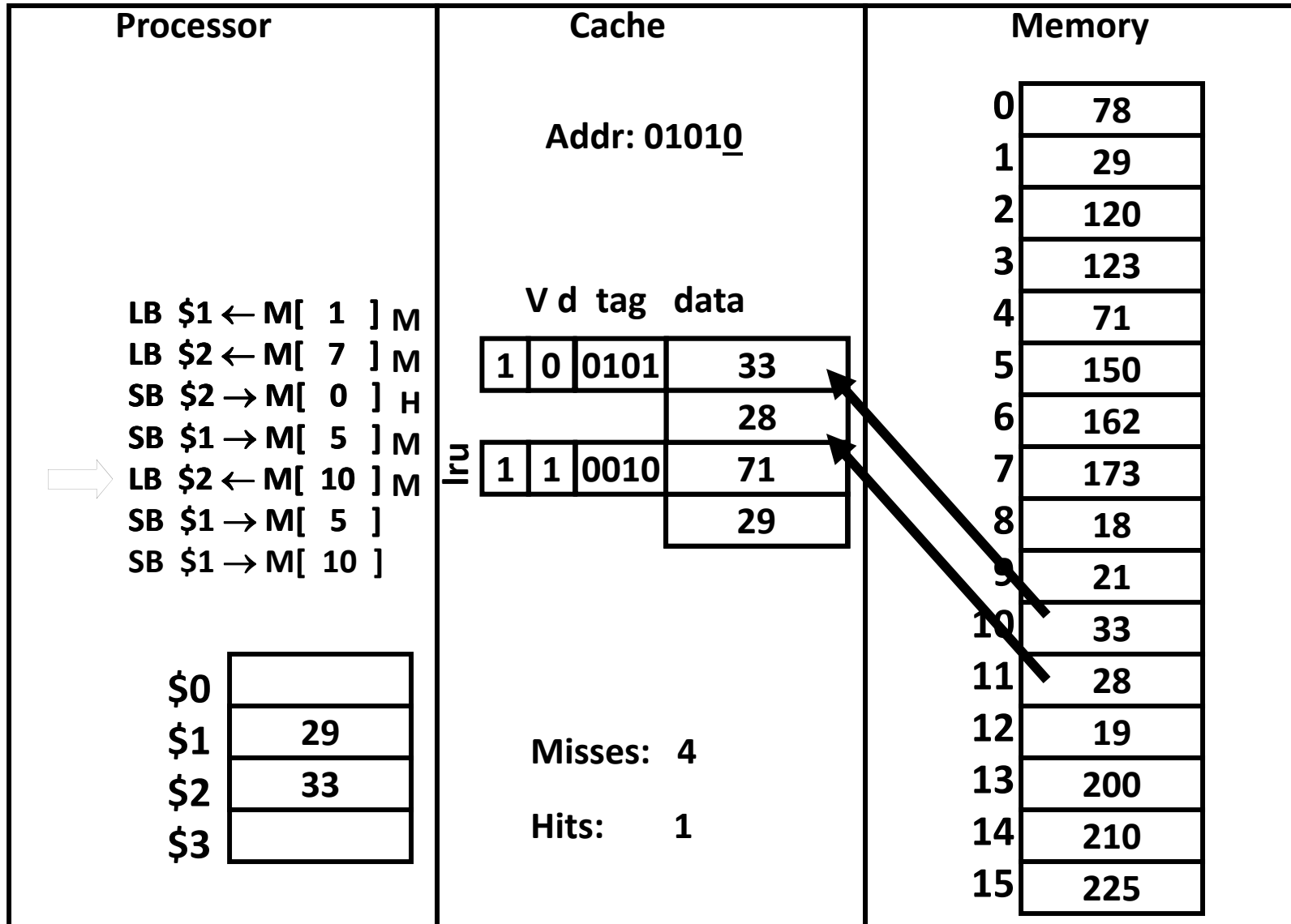


\$0	
\$1	29
\$2	173
\$3	

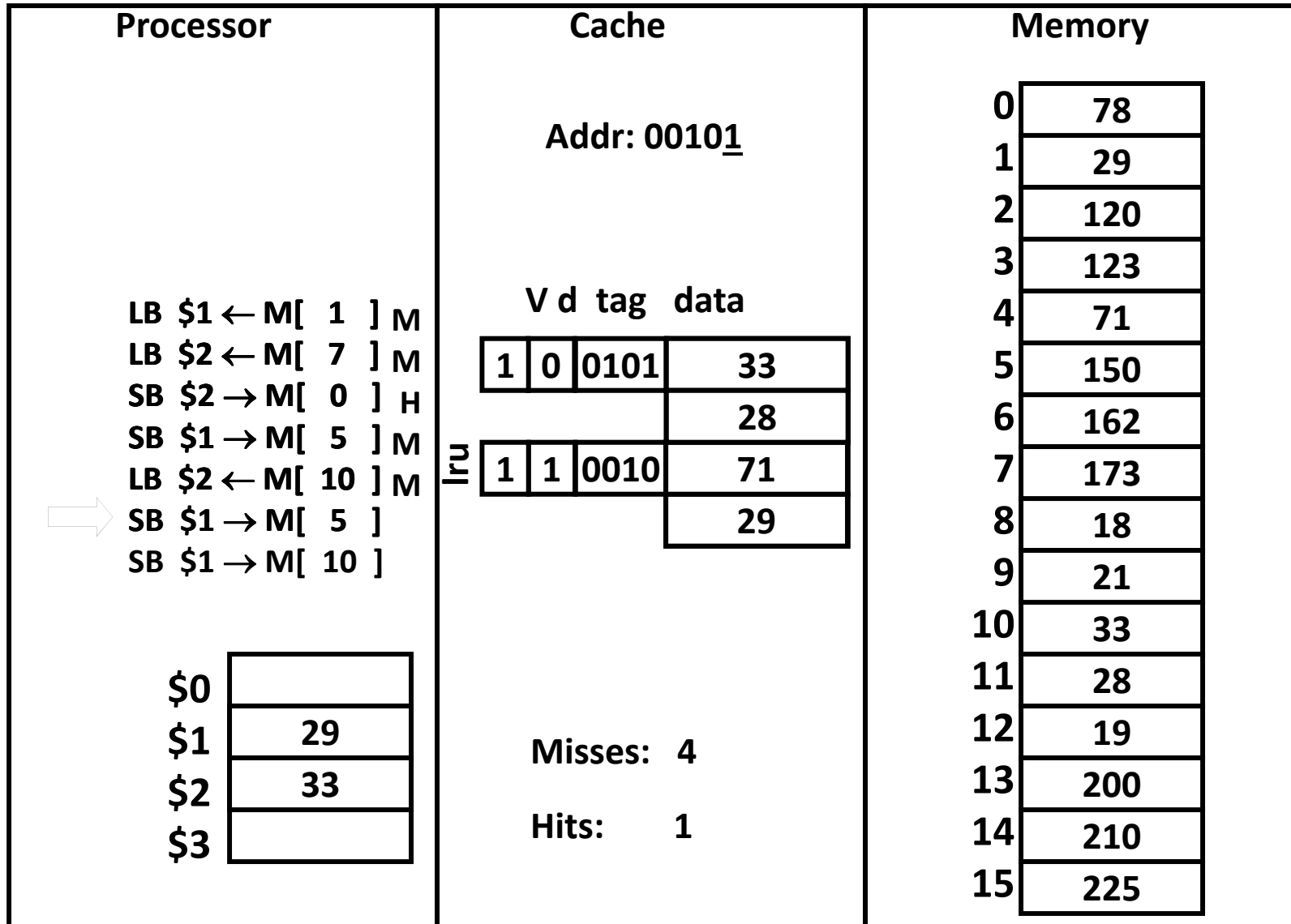
Write-Back (REF 5)



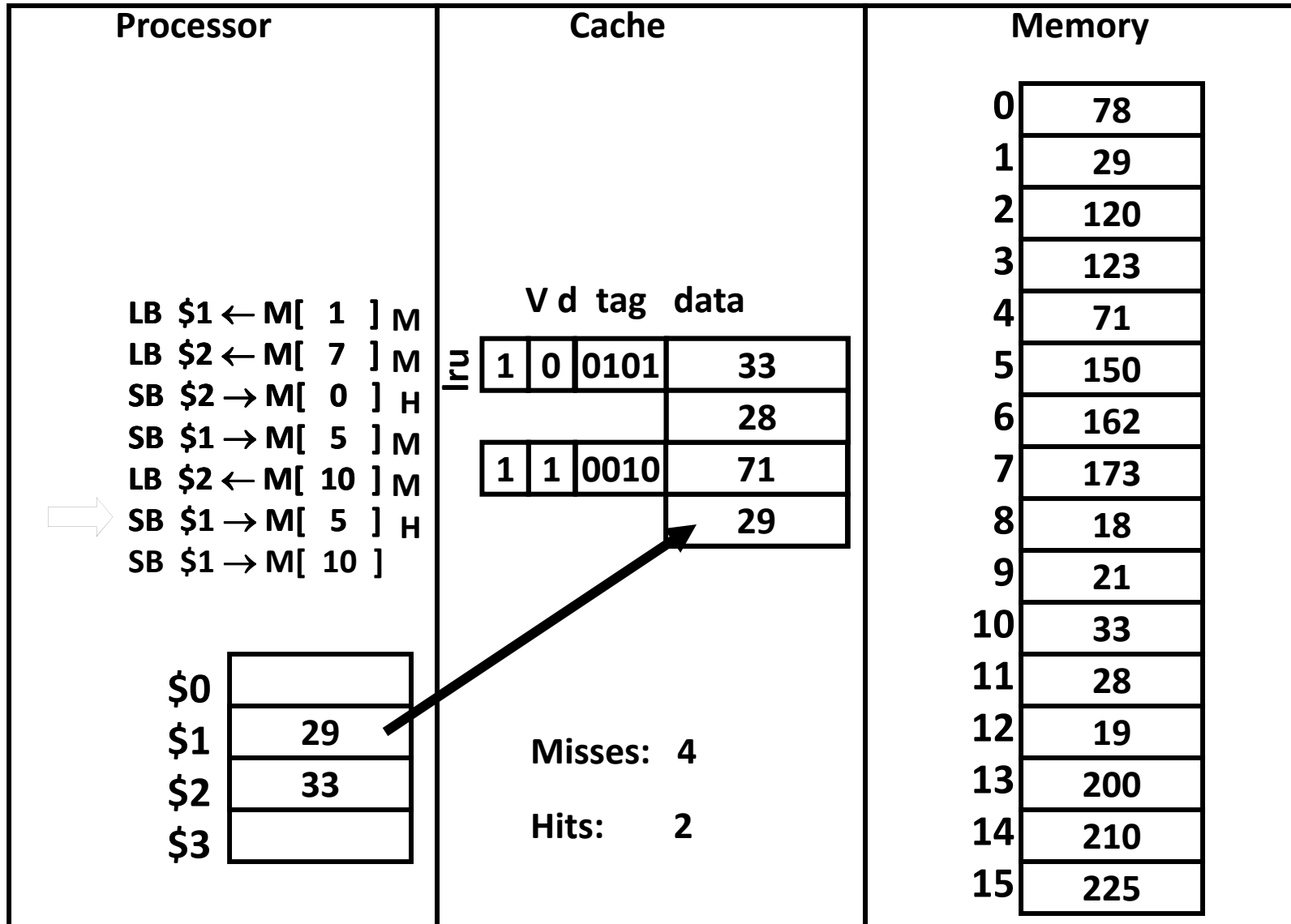
Write-Back (REF 5)



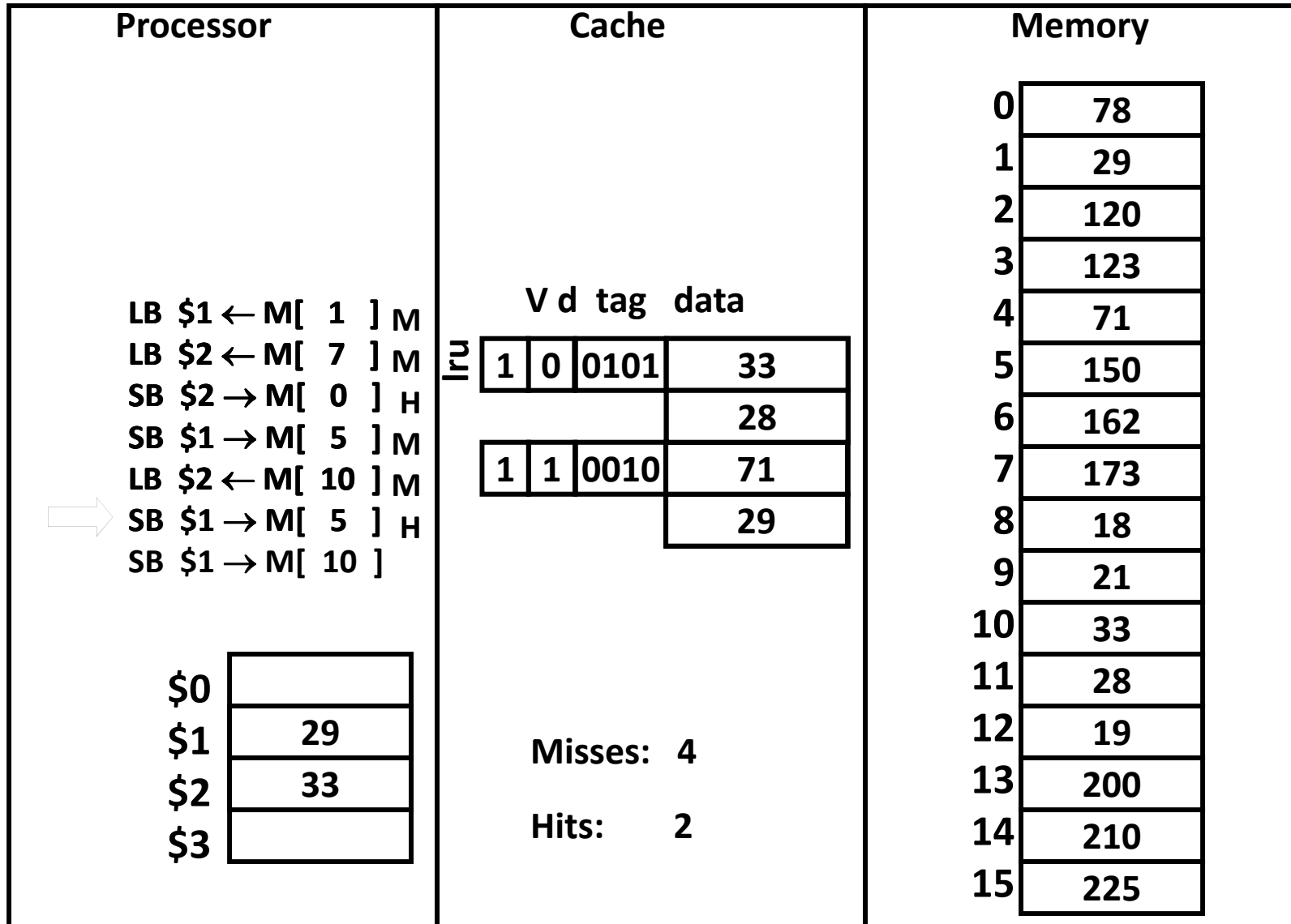
Write-Back (REF 6)



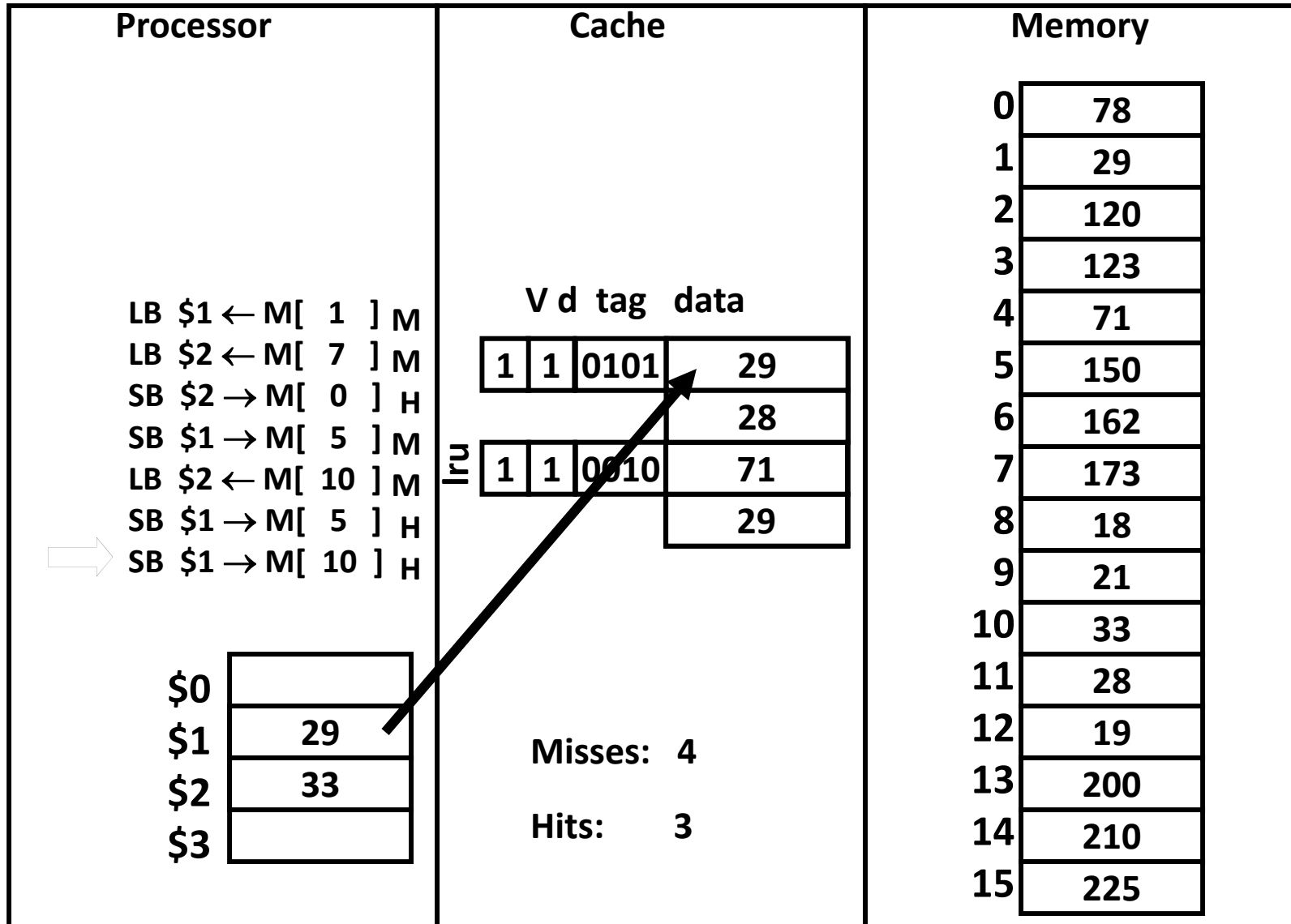
Write-Back (REF 6)



Write-Back (REF 7)



Write-Back (REF 7)



How Many Memory References?

Write-back performance

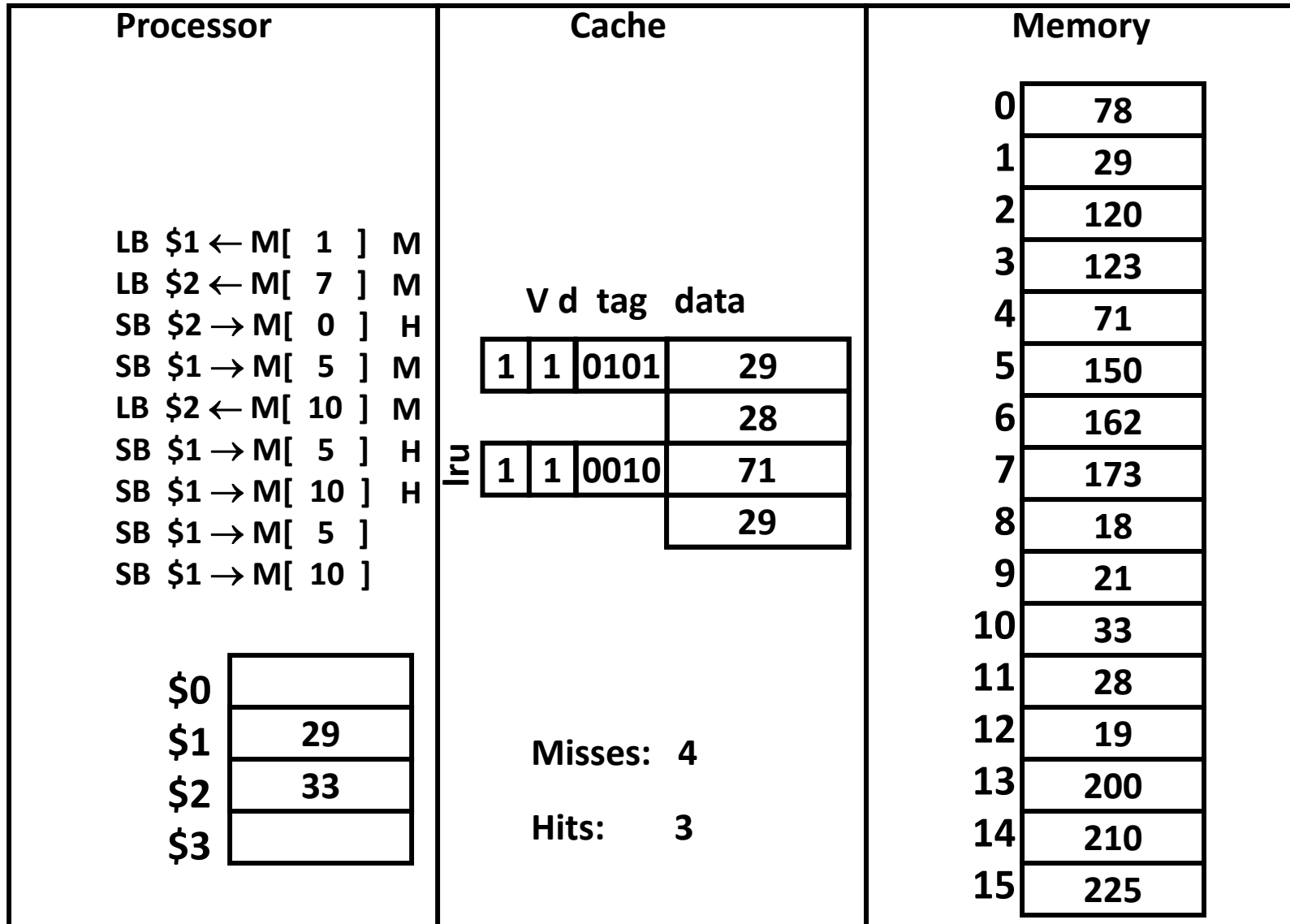
Each miss (read or write) reads a block from mem

- 4 misses \rightarrow 8 mem reads

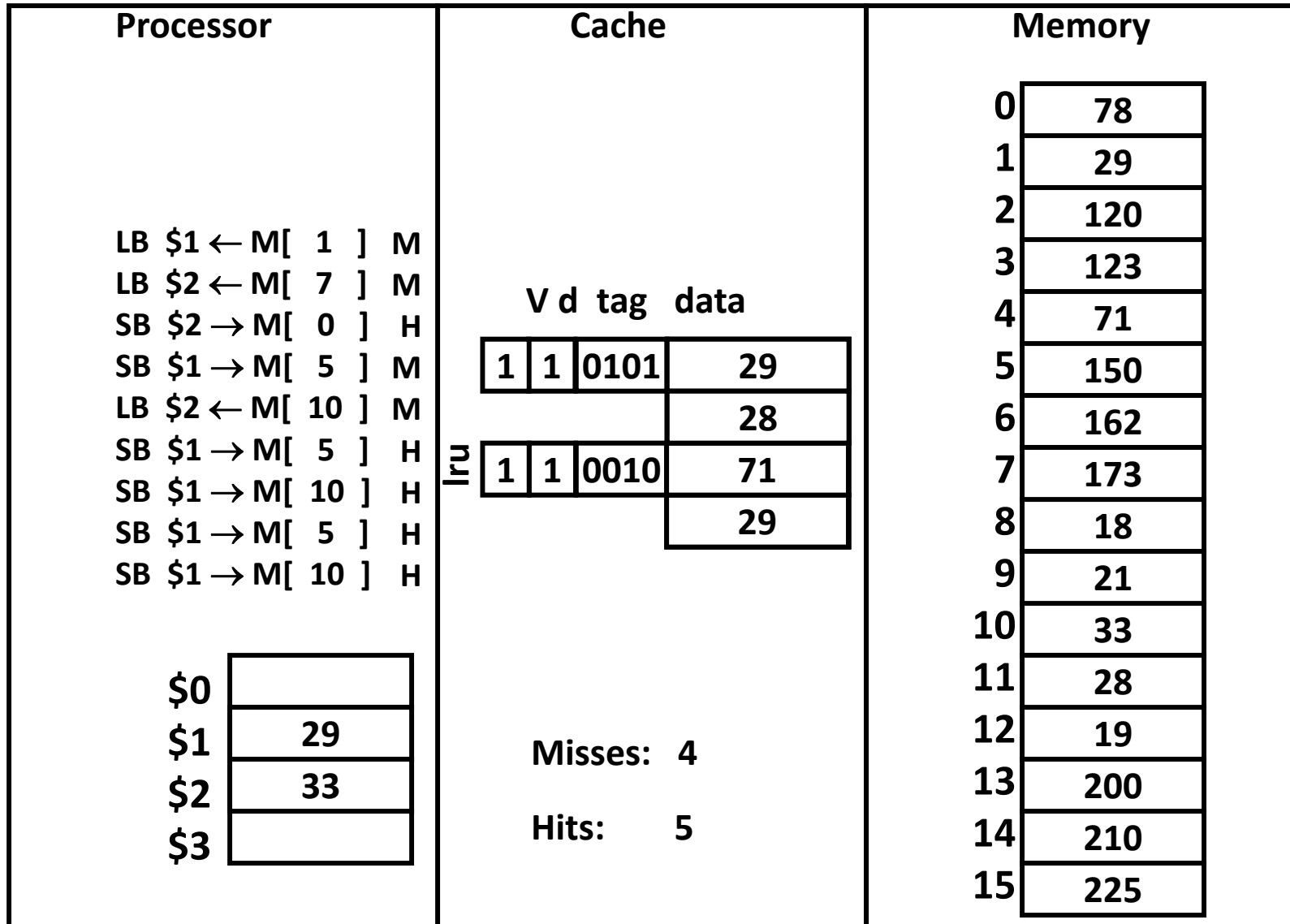
Some evictions write a block to mem

- 1 dirty eviction \rightarrow 2 mem writes
- (+ 2 dirty evictions later \rightarrow +4 mem writes)

Write-Back (REF 8,9)



Write-Back (REF 8,9)



How Many Memory References?

Write-back performance

Each miss (read or write) reads a block from mem

- 4 misses \rightarrow 8 mem reads

Some evictions write a block to mem

- 1 dirty eviction \rightarrow 2 mem writes
- (+ 2 dirty evictions later \rightarrow +4 mem writes)

By comparison write-through was

- Reads: eight words
- Writes: 4/6/8/10/12/... etc words

So is write back just better?

What are other performance tradeoffs between write-through and write-back?

How can we further reduce penalty for cost of writes to memory?

Performance Tradeoffs

Q: Hit time: write-through vs. write-back?

A: Write-through slower on writes

Q: Miss penalty: write-through vs. write-back?

A: Write-back slower on evictions

Performance: An Example

Performance: Write-back versus Write-through

Assume: large associative cache, 16-byte lines

```
for (i=1; i<n; i++)  
    A[0] += A[i];
```

N words

Write-through: $n/16$ reads
n writes

Write-back: $n/16$ reads
1 write

```
for (i=0; i<n; i++)  
    B[i] = A[i]
```

Write-through: $2 \times n/16$ reads
n writes

Write-back: $2 \times n/16$ reads
n write

Write Buffering

Q: Writes to main memory are **slow!**

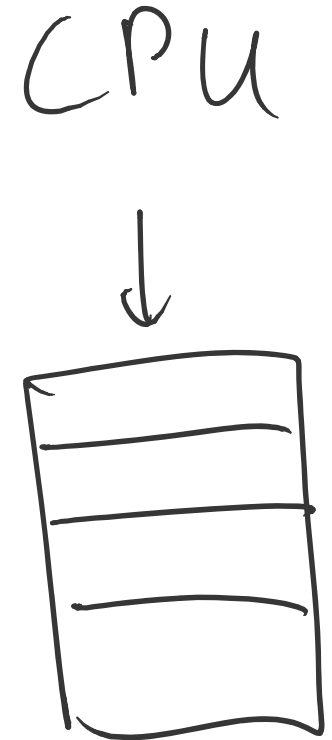
A: Use a write-back buffer

- A small queue holding dirty lines
- Add to end upon eviction
- Remove from front upon completion

Q: When does it help?

A: short bursts of writes (but not sustained writes)

A: fast eviction reduces miss penalty



Write-through vs. Write-back

Write-through is slower

- But simpler (memory always consistent)

Write-back is almost always faster

- write-back buffer hides large eviction cost
- But what about multiple cores with separate caches but sharing memory?

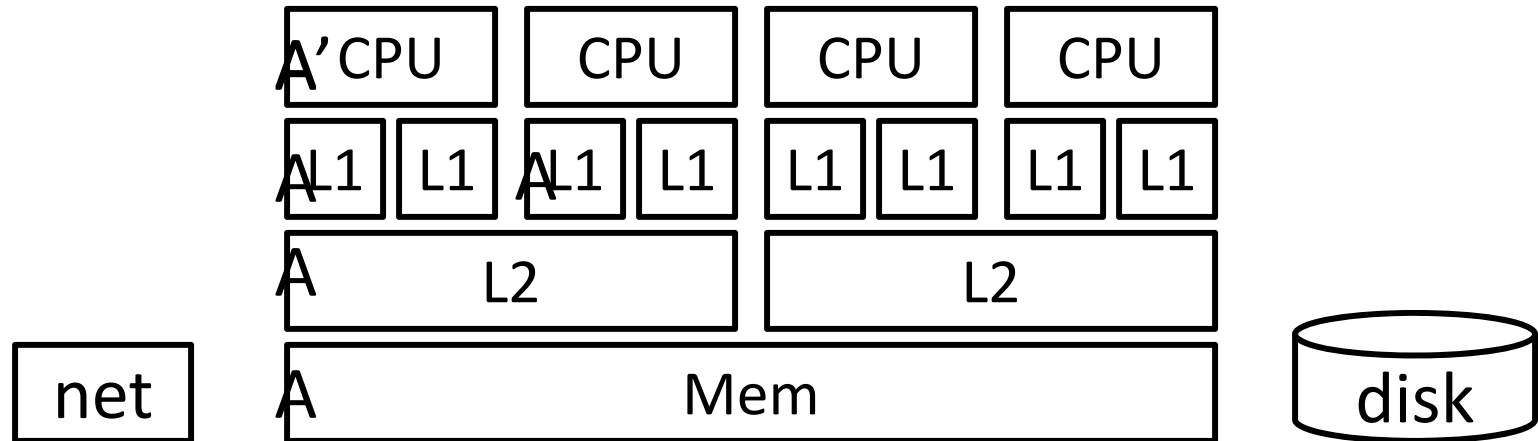
Write-back requires a cache coherency protocol

- Inconsistent views of memory
- Need to “snoop” in each other’s caches
- Extremely complex protocols, very hard to get right

Cache-coherency

Q: Multiple readers and writers?

A: Potentially inconsistent views of memory



Cache coherency protocol

- May need to snoop on other CPU's cache activity
- Invalidate cache line when other CPU writes
- Flush write-back caches before other CPU reads
- Or the reverse: Before writing/reading...
- Extremely complex protocols, very hard to get right

Summary: Write Through

Write-through policy with write allocate

- Cache miss: read entire block from memory
- Write: write only updated item to memory
- Eviction: no need to write to memory
- Slower, but cleaner

Write-back policy with write allocate

- Cache miss: read entire block from memory
 - **But may need to write dirty cacheline first**
- Write: nothing to memory
- Eviction: have to write to memory, entire cacheline because don't know what is dirty (only 1 dirty bit)
- Faster, but complicated with multicore

Next Goal

Performance: What is the average memory access time (AMAT) for a cache?

$$\text{AMAT} = \% \text{hit} \times \text{hit time} + \% \text{miss} \times \text{miss time}$$

Cache Performance Example

Average Memory Access Time (AMAT)

Cache Performance (very simplified):

L1 (SRAM): 512 x 64 byte cache lines, direct mapped

Data cost: 3 cycle per word access

Lookup cost: 2 cycle

16 words (i.e. $64 / 4 = 16$)

Mem (DRAM): 4GB

Data cost: 50 cycle for first word, plus 3 cycles per subsequent word

Cache Performance Example

Average Memory Access Time (AMAT)

Cache Performance (very simplified):

L1 (SRAM): 512 x 64 byte cache lines, direct mapped

Data cost: 3 cycle per word access

Lookup cost: 2 cycle

16 words (i.e. $64 / 4 = 16$)

Mem (DRAM): 4GB

Data cost: 50 cycle for first word, plus 3 cycles per subsequent word

AMAT = %hit x hit time + % miss x miss time

Hit time = 5 cycles

Miss time = hit time + 50 (first word) + 15 x 3 (words)
= 100 cycles

If %hit = 90%, then

AMAT = .9 x 5 + .1 x 100 = 14.5 cycles

Multi Level Caching

Cache Performance (very simplified):

L1 (SRAM): 512 x 64 byte cache lines, direct mapped

Hit time: 5 cycles

L2 cache: bigger

Hit time = 20 cycles

Mem (DRAM): 4GB

Hit rate: 90% in L1, 90% in L2

AMAT = %hit x hit time + % miss x miss time

AMAT = .9 x 5 + .1 (.9 x 20 + .1 x 120) = 4.5 + .1 (18 + 12) = 7.5

Often: L1 fast and direct mapped, L2 bigger and higher associativity

Next Goal

How to decide on Cache Organization

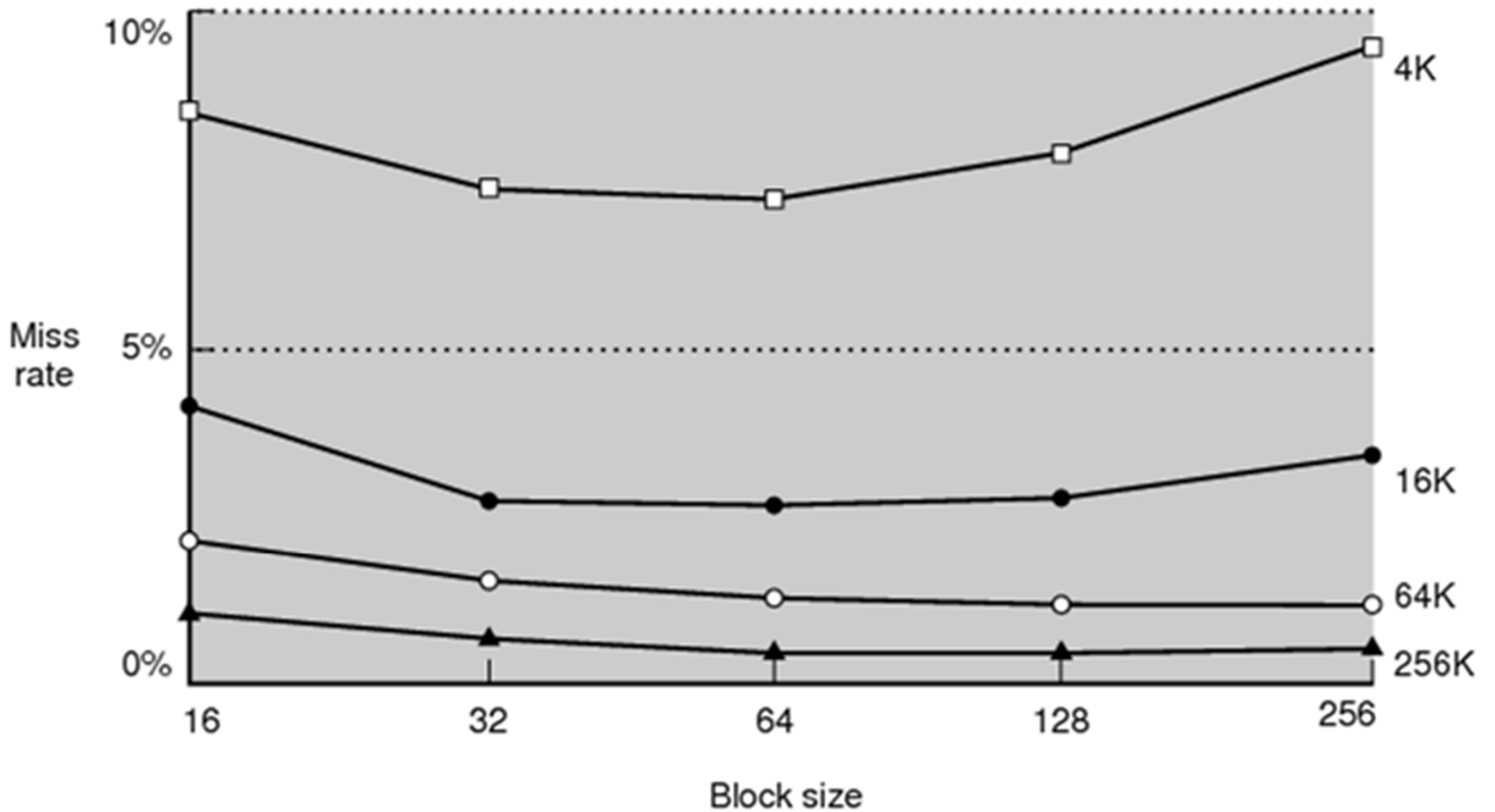
Q: How to decide block size?

A: Try it and see

But: depends on cache size, workload, associativity, ...

Experimental approach!

Experimental Results



Tradeoffs

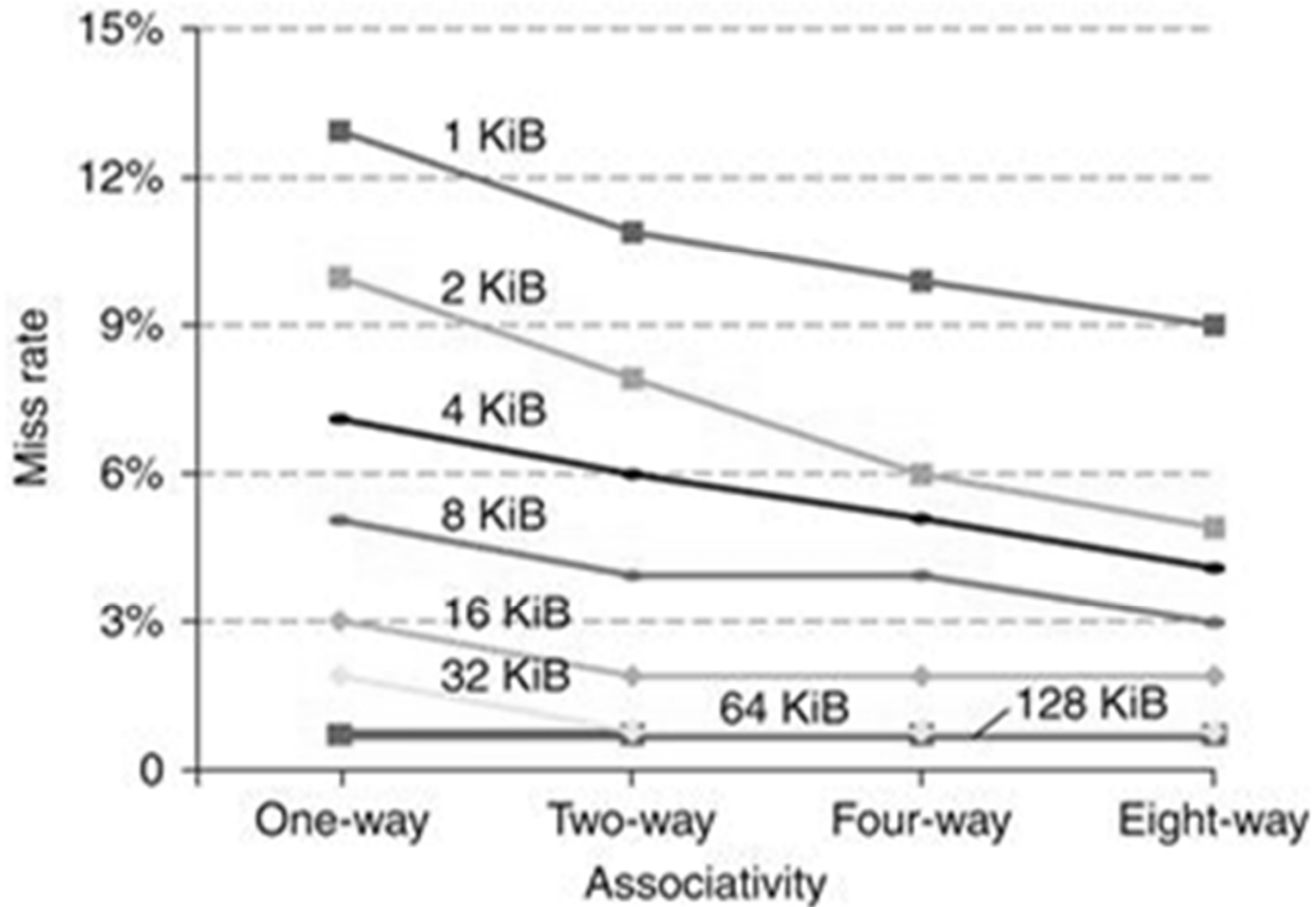
For a given total cache size,
larger block sizes mean....

- fewer lines
- so fewer tags, less overhead
- and fewer cold misses (within-block “prefetching”)

But also...

- fewer blocks available (for scattered accesses!)
- so more conflicts
- and larger miss penalty (time to fetch block)

Associativity



Other Designs

Multilevel caches

Performance Summary

Average memory access time (AMAT)

depends on cache architecture and size

access time for hit,

miss penalty, miss rate

Cache design a very complex problem:

- Cache size, block size (aka line size)
- Number of ways of set-associativity (1, N, ∞)
- Eviction policy
- Number of levels of caching, parameters for each
- Separate I-cache from D-cache, or Unified cache
- Prefetching policies / instructions
- Write policy

Cache Conscious Programming

```
// H = 12, W = 10
int A[H][W];

for(x=0; x < W; x++)
    for(y=0; y < H; y++)
        sum += A[y][x];
```

1									
2									
3									
4									
5									
6									
7									
8									
9									
10									
11									
12									

Every access is a cache miss!
(unless *entire* matrix can fit in cache)

A Real Example

```
> dmidecode -t cache
```

Cache Information

```
Socket Designation: L1 Cache
Configuration: Enabled, Not Socketed, Level 1
Operational Mode: Write Back
Location: Internal
Installed Size: 32 kB
Maximum Size: 32 kB
Supported SRAM Types:
    Synchronous
Installed SRAM Type: Synchronous
Speed: Unknown
Error Correction Type: Single-bit ECC
System Type: Instruction
Associativity: 8-way Set-associative
```

Cache Information

```
Socket Designation: L2 Cache
Configuration: Enabled, Not Socketed,
    Level 2
Operational Mode: Write Back
Location: Internal
Installed Size: 256 kB
Maximum Size: 256 kB
Supported SRAM Types:
    Synchronous
Installed SRAM Type: Synchronous
Speed: Unknown
Error Correction Type: Single-bit ECC
System Type: Unified
Associativity: 8-way Set-associative
```

Lenovo Yoga 2 Pro laptop

Dual core

Intel i7-4500 CPU @ 2.4 GHz
(purchased in 2014)

Cache Information

```
Socket Designation: L3 Cache
Configuration: Enabled, Not Socketed,
    Level 3
Operational Mode: Write Back
Location: Internal
Installed Size: 4096 kB
Maximum Size: 4096 kB
Supported SRAM Types:
    Synchronous
Installed SRAM Type: Synchronous
Speed: Unknown
Error Correction Type: Single-bit ECC
System Type: Unified
Associativity: 16-way Set-associative
```


A Real Example

Dual-core 3.16GHz Intel
(purchased in 2011)

```
> dmidecode -t cache
```

```
Cache Information
```

```
Configuration: Enabled, Not Socketed, Level 1  
Operational Mode: Write Back  
Installed Size: 128 KB  
Error Correction Type: None
```

```
Cache Information
```

```
Configuration: Enabled, Not Socketed, Level 2  
Operational Mode: Varies With Memory Address  
Installed Size: 6144 KB  
Error Correction Type: Single-bit ECC
```

```
> cd /sys/devices/system/cpu/cpu0; grep cache/*/*
```

```
cache/index0/level:1  
cache/index0/type:Data  
cache/index0/ways_of_associativity:8  
cache/index0/number_of_sets:64  
cache/index0/coherency_line_size:64  
cache/index0/size:32K  
cache/index1/level:1  
cache/index1/type:Instruction  
cache/index1/ways_of_associativity:8  
cache/index1/number_of_sets:64  
cache/index1/coherency_line_size:64  
cache/index1/size:32K  
cache/index2/level:2  
cache/index2/type:Unified  
cache/index2/shared_cpu_list:0-1  
cache/index2/ways_of_associativity:24  
cache/index2/number_of_sets:4096  
cache/index2/coherency_line_size:64  
cache/index2/size:6144K
```

A Real Example

Dual-core 3.16GHz Intel
(purchased in 2009)

Dual 32K L1 Instruction caches

- 8-way set associative
- 64 sets
- 64 byte line size

Dual 32K L1 Data caches

- Same as above

Single 6M L2 Unified cache

- 24-way set associative (!!!)
- 4096 sets
- 64 byte line size

4GB Main memory

1TB Disk

By the end of the cache lectures...

MacBook Pro

Retina, Mid 2012

Processor 2.7 GHz Intel Core i7

Memory 16 GB 1600 MHz DDR3

Graphics NVIDIA GeForce GT 650M 1024 MB

Serial Number C02J70TTDKQ5

Software OS X 10.9.2 (13C64)

Model Name:	MacBook Pro
Model Identifier:	MacBookPro10,1
Processor Name:	Intel Core i7
Processor Speed:	2.7 GHz
Number of Processors:	1
Total Number of Cores:	4
L2 Cache (per Core):	256 KB
L3 Cache:	8 MB
Memory:	16 GB
Boot ROM Version:	MBP101.00EE.B02
SMC Version (system):	2.3f36
Serial Number (system):	C02J70TTDKQ5
Hardware UUID:	F588E08C-60BF-5B35-A087-07714C2B2D11

- 32 KB data + 32 KB instruction L1 cache (3 clocks) and 256 KB L2 cache (8 clocks) per core.
- Shared L3 cache includes the processor graphics (LGA 1155).
- 64-byte cache line size.

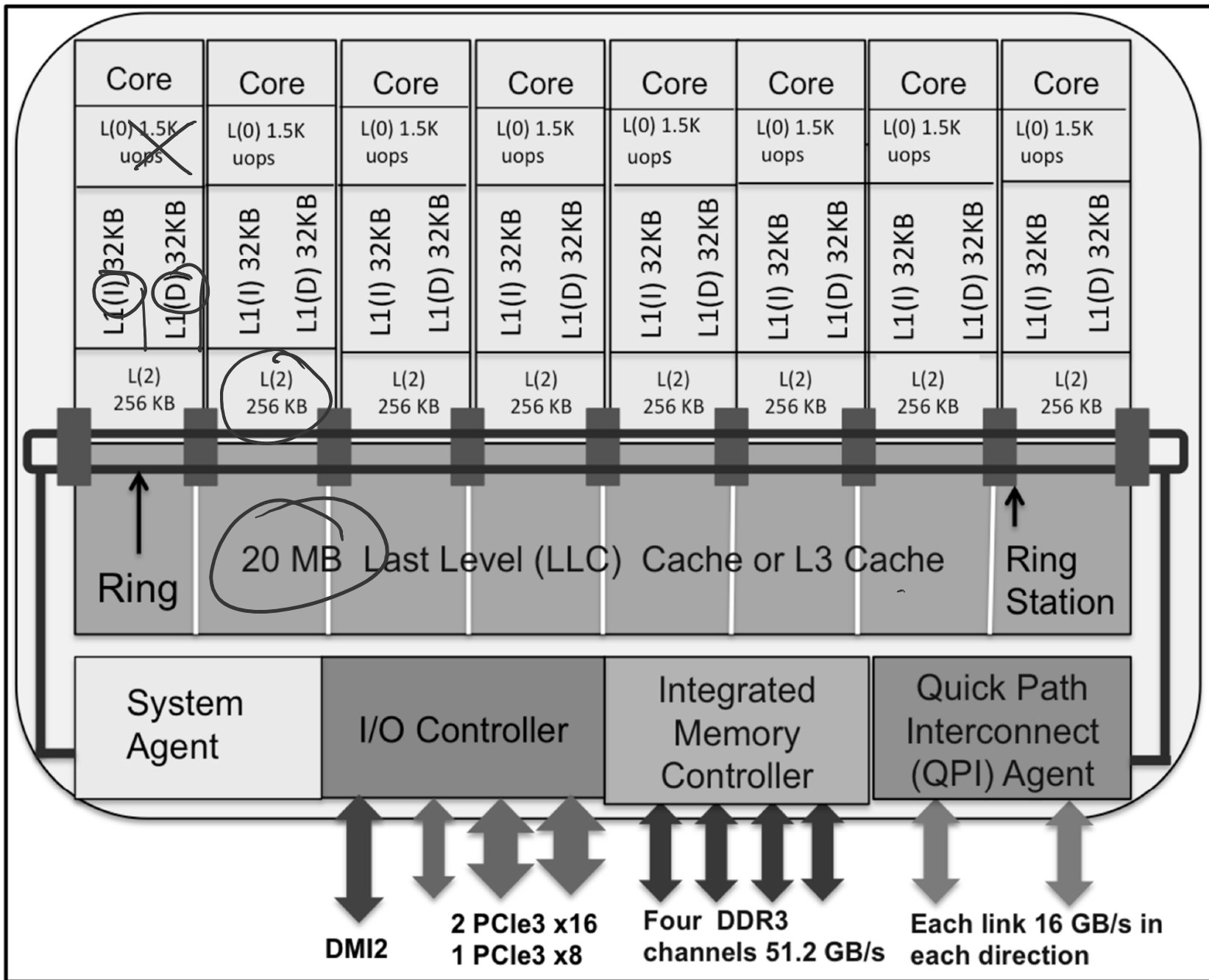


Figure 1. Schematic diagram of a Sandy Bridge processor.

Summary

Memory performance matters!

- often more than CPU performance
- ... because it is the bottleneck, and not improving much
- ... because most programs move a LOT of data

Design space is huge

- Gambling against program behavior
- Cuts across all layers:
users → programs → os → hardware

Multi-core / Multi-Processor is complicated

- Inconsistent views of memory
- Extremely complex protocols, very hard to get right

Have a great Spring Break!!