

Lec 17: Caches

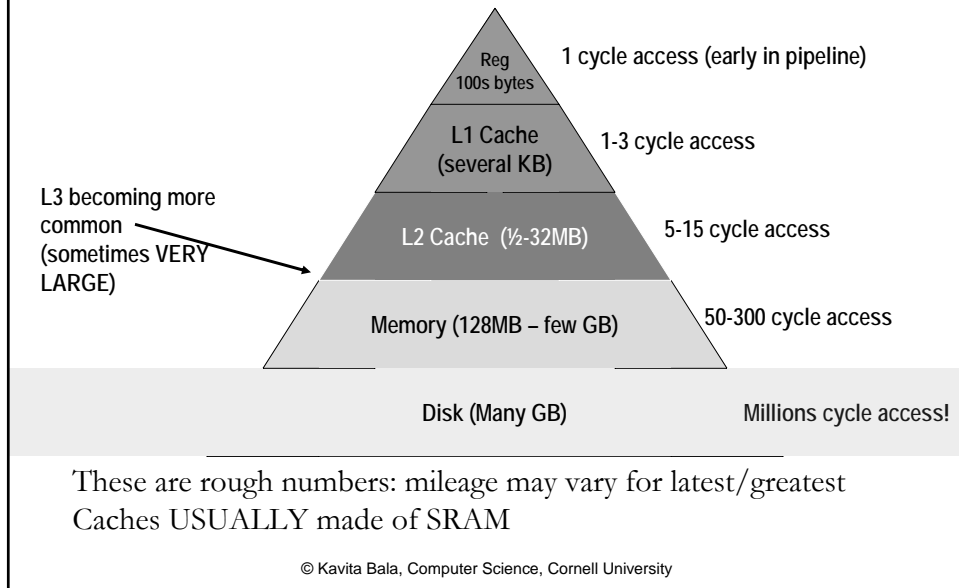
Kavita Bala
CS 3410, Fall 2008
Computer Science
Cornell University

Announcements

- Prelim: graded
- PA 2: graded
- HW 2: graded

- HW 3 out tonight: cache simulation
 - Recitations this week on C/Unix/etc.

Cache Design 101



Insight of Caches

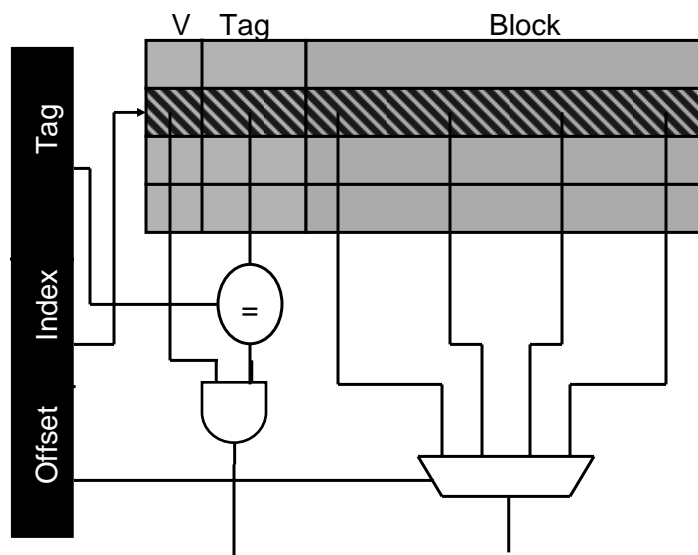
- Exploit locality
 - Two types: temporal and spatial
- Temporal locality
 - If memory location X is accessed, then it is more likely to be accessed again in the near future than some random location Y
 - Caches exploit temporal locality by placing a memory element that has been referenced into the cache
- Spatial locality
 - If memory location X is accessed, then locations near X are more likely to be accessed in the near future than some random location Y
 - Caches exploit spatial locality by allocating a cache line of data (including data near the referenced location)

Cache Lookups (Read)

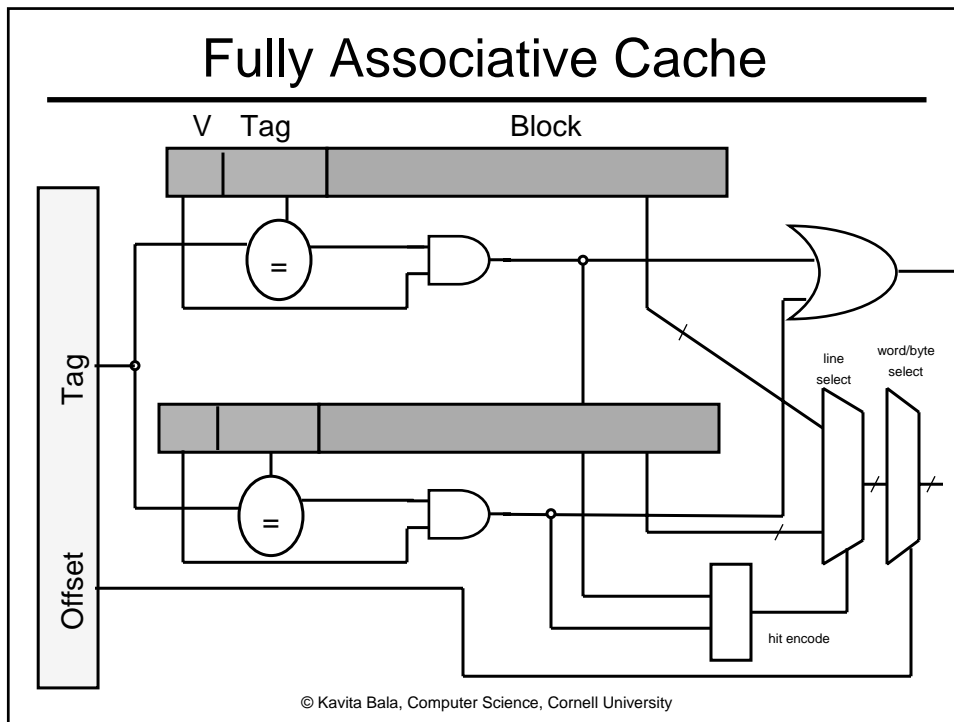
- Look at address issued by processor
- Search cache to see if that block is in the cache
 - Hit: Block is in the cache
 - return requested data
 - Miss: Block is not in the cache
 - read line from memory
 - evict an existing line from the cache
 - place new line in cache
 - return requested data

© Kavita Bala, Computer Science, Cornell University

Direct Mapped Cache



© Kavita Bala, Computer Science, Cornell University



- ## Cache Organization
- Three common designs
 - Fully associative: Block can be anywhere in the cache
 - Direct mapped: Block can only be in one line in the cache
 - Set-associative: Block can be in a few (2 to 8) places in the cache
- © Kavita Bala, Computer Science, Cornell University

Eviction

- Which cache line should be evicted from the cache to make room for a new line?
 - Direct-mapped
 - no choice, must evict line selected by index
 - Associative caches
 - random: select one of the lines at random
 - round-robin: similar to random
 - FIFO: replace oldest line
 - LRU: replace line that has not been used in the longest time

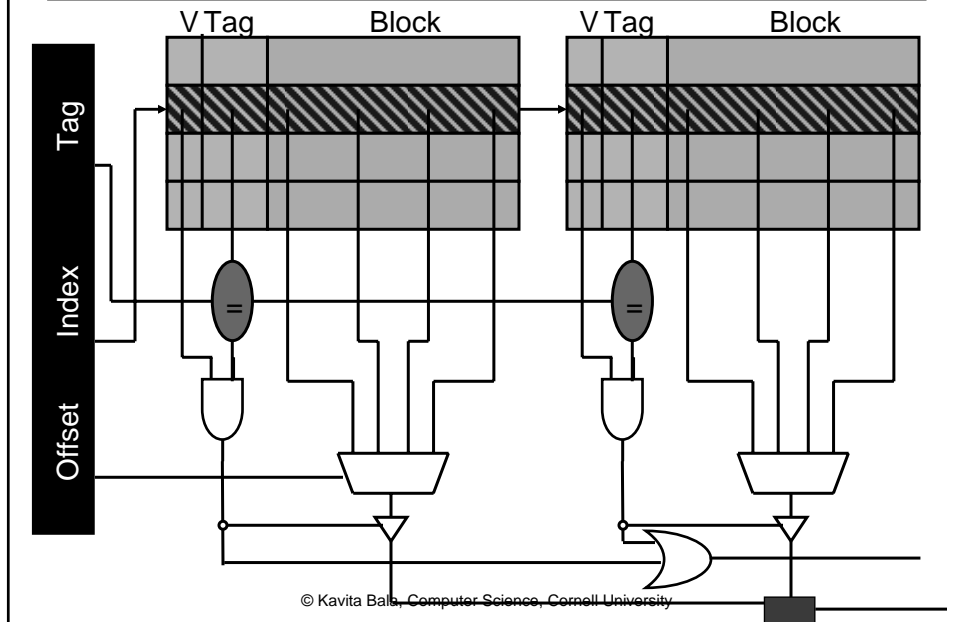
© Kavita Bala, Computer Science, Cornell University

Compromise

- Set-associative cache
- Like a direct-mapped cache
 - Index into a location
 - Fast
- Like a fully-associative cache
 - Can store multiple entries
 - decreases thrashing in cache
 - Search in each element

© Kavita Bala, Computer Science, Cornell University

2-Way Set-Associative Cache



Direct Mapped

Processor	Cache	Memory
Ld R1 ← M[1]	4 cache lines	0 100
Ld R2 ← M[5]	1 bit tag field	1 110
Ld R3 ← M[1]	2 word block	2 120
Ld R3 ← M[4]	tag data	3 130
Ld R2 ← M[0]	1 0 100	4 140
Ld R2 ← M[12]	0 110	5 150
Ld R2 ← M[5]	1 0 140	6 160
Ld R2 ← M[12]	0 150	7 170
➔ Ld R2 ← M[12]	Misses: 4+2+2	8 180
Ld R2 ← M[5]	Hits: 3	9 190
Ld R2 ← M[12]		10 200
Ld R2 ← M[5]		11 210
		12 220
		13 230
		14 240
		15 250

© Kavita Bala, Computer Science, Cornell University

Fully Assoc

Processor	Cache	Memory
	4 cache lines	0 100
	3 bit tag field	1 110
	2 word block	2 120
	tag data	3 130
Ld R1 ← M[1]	1 0 100	4 140
Ld R2 ← M[5]		5 150
Ld R3 ← M[1]		6 160
Ld R3 ← M[4]		7 170
Ld R2 ← M[0]	1 2 140	8 180
Ld R2 ← M[12]		9 190
Ld R2 ← M[5]		10 200
→ Ld R2 ← M[12]	1 6 220	11 210
Ld R2 ← M[5]		12 220
Ld R2 ← M[12]		13 230
Ld R2 ← M[5]		14 240
		15 250
	Misses: 3	
	Hits: 8	

© Kavita Bala, Computer Science, Cornell University

2 Way Set Assoc

Processor	Cache	Memory
	2 sets	0 100
	2 bit tag field	1 110
	2 word block	2 120
	tag data	3 130
Ld R1 ← M[1]	0 0	4 140
Ld R2 ← M[5]		5 150
Ld R3 ← M[1]		6 160
Ld R3 ← M[4]		7 170
Ld R2 ← M[0]	0 0	8 180
Ld R2 ← M[12]		9 190
Ld R2 ← M[5]		10 200
Ld R2 ← M[12]		11 210
Ld R2 ← M[5]		12 220
Ld R2 ← M[12]		13 230
Ld R2 ← M[5]		14 240
		15 250
	Misses: 4	
	Hits: 7	

© Kavita Bala, Computer Science, Cornell University

Cache Design

- Need to determine parameters
 - Block size
 - Number of ways of set-associativity
 - Eviction policy
 - Write policy
 - Separate I-cache from D-cache

© Kavita Bala, Computer Science, Cornell University

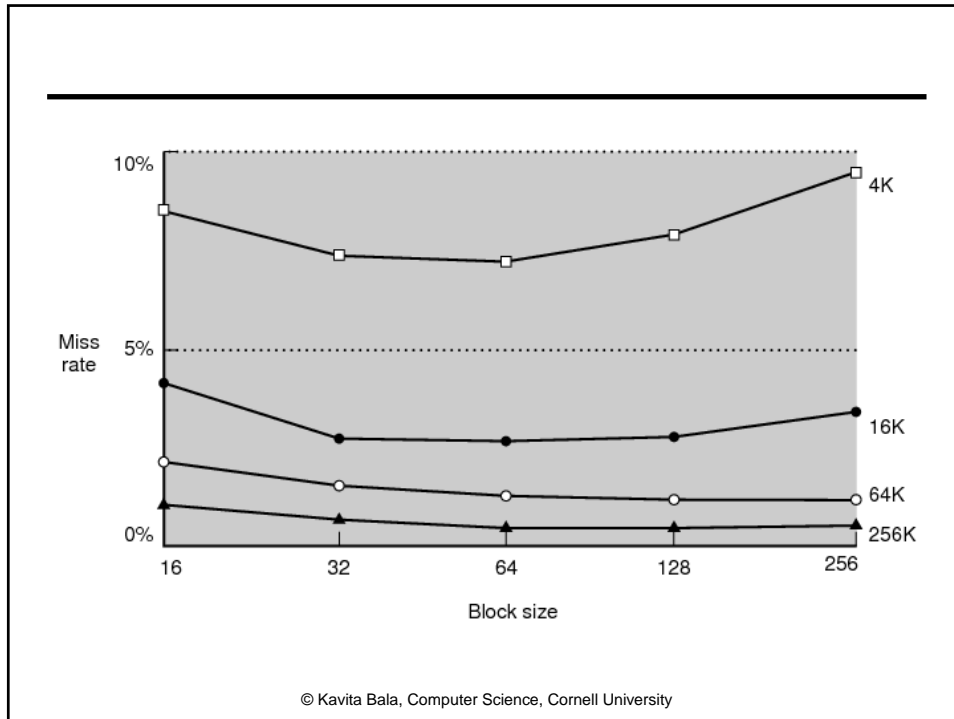
Basic Cache Organization

Decide on the block size

- How? Simulate lots of different block sizes and see which one gives the best performance
- Most systems use a block size between 32 bytes and 128 bytes



© Kavita Bala, Computer Science, Cornell University



Tradeoff

- Larger sizes reduce the overhead by
 - Reducing the number of tags
 - Reducing the size of each tag

- But
 - Have fewer blocks available
 - And the time to fetch the block on a miss is longer

© Kavita Bala, Computer Science, Cornell University

Valid Bits

- Valid bits indicate whether cache line contains an up-to-date copy of the values in memory
 - Must be 1 for a hit
 - Reset to 0 on power up
- An item can be removed from the cache by setting its valid bit to 0

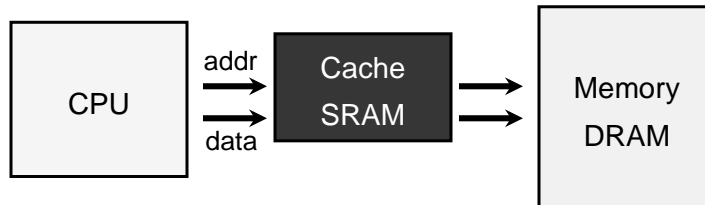
© Kavita Bala, Computer Science, Cornell University

Eviction

- Which cache line should be evicted from the cache to make room for a new line?
 - Direct-mapped
 - no choice, must evict line selected by index
 - Associative caches
 - random: select one of the lines at random
 - round-robin: similar to random
 - FIFO: replace oldest line
 - LRU: replace line that has not been used in the longest time

© Kavita Bala, Computer Science, Cornell University

Cache Writes



- No-Write
 - writes invalidate the cache and go to memory
- Write-Through
 - writes go to main memory and cache
- Write-Back
 - write cache, write main memory only when block is evicted

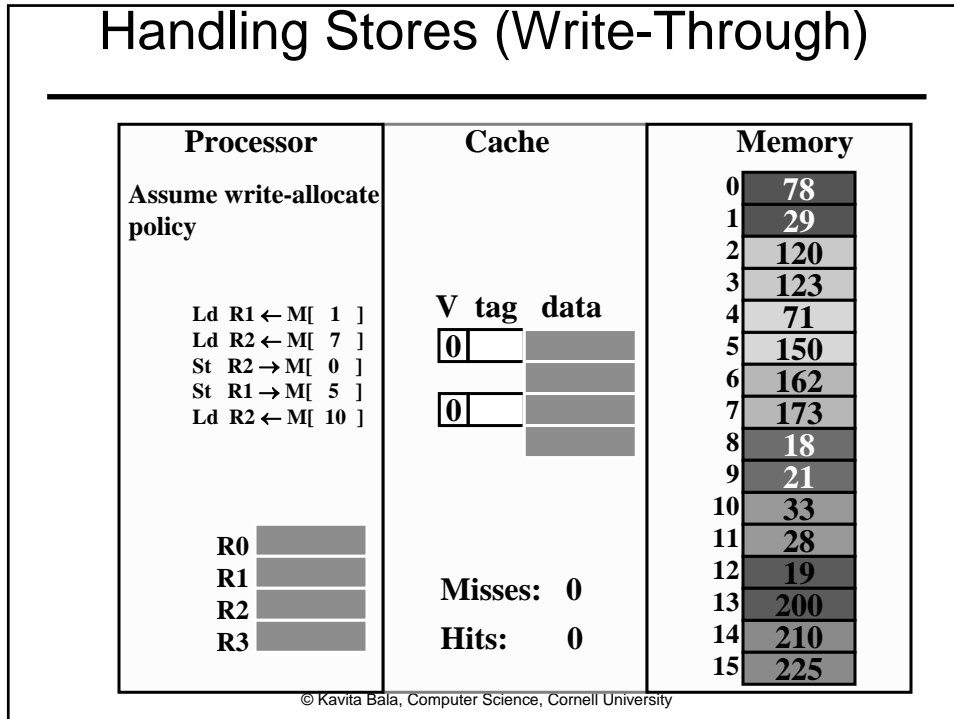
© Kavita Bala, Computer Science, Cornell University

What about Stores?

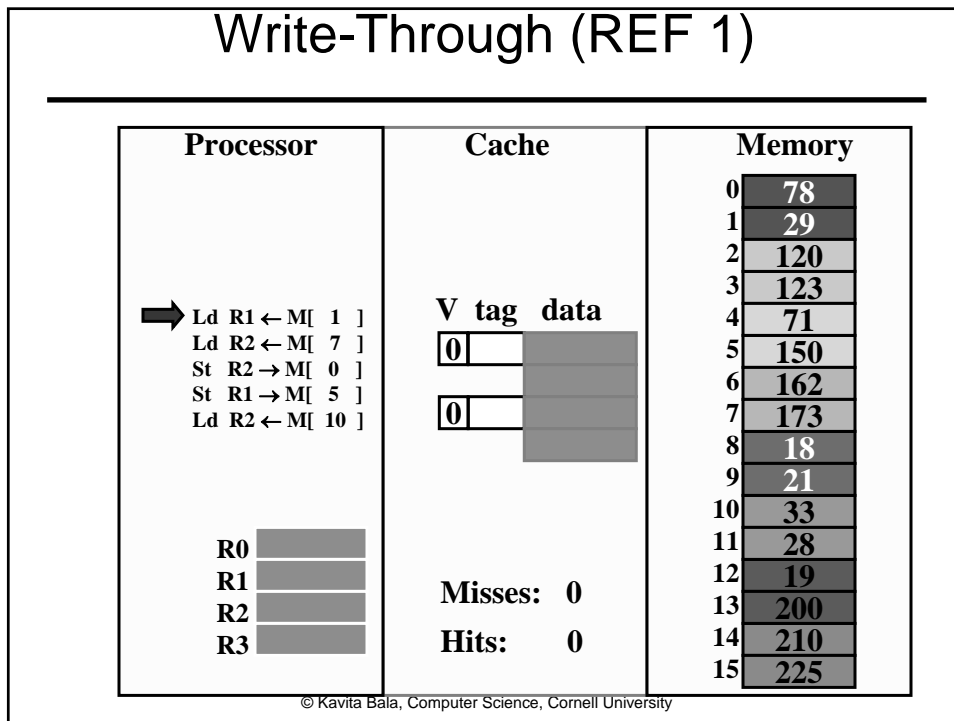
- Where should you write the result of a store?
 - If that memory location is in the cache?
 - Send it to the cache
 - Should we also send it to memory right away? (write-through policy)
 - Wait until we kick the block out (write-back policy)
 - If it is not in the cache?
 - Allocate the line (put it in the cache)? (write allocate policy)
 - Write it directly to memory without allocation? (no write allocate policy)

© Kavita Bala, Computer Science, Cornell University

Handling Stores (Write-Through)



Write-Through (REF 1)



Write-Through (REF 1)

Processor	Cache	Memory																																						
<p>➔ Ld R1 ← M[1] Ld R2 ← M[7] St R2 → M[0] St R1 → M[5] Ld R2 ← M[10]</p> <p>R0 R1 29 R2 R3 </p>	<p style="text-align: center;">V tag data</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td style="width: 10px; text-align: center;">1</td> <td style="width: 10px; text-align: center;">0</td> <td style="width: 20px; text-align: center;">78</td> </tr> <tr> <td></td> <td></td> <td style="text-align: center;">29</td> </tr> </table> <p>lru 0 </p> <p style="text-align: center;">Misses: 1 Hits: 0</p>	1	0	78			29	<table border="1" style="width: 100%; text-align: center;"> <tr><td>0</td><td>78</td></tr> <tr><td>1</td><td>29</td></tr> <tr><td>2</td><td>120</td></tr> <tr><td>3</td><td>123</td></tr> <tr><td>4</td><td>71</td></tr> <tr><td>5</td><td>150</td></tr> <tr><td>6</td><td>162</td></tr> <tr><td>7</td><td>173</td></tr> <tr><td>8</td><td>18</td></tr> <tr><td>9</td><td>21</td></tr> <tr><td>10</td><td>33</td></tr> <tr><td>11</td><td>28</td></tr> <tr><td>12</td><td>19</td></tr> <tr><td>13</td><td>200</td></tr> <tr><td>14</td><td>210</td></tr> <tr><td>15</td><td>225</td></tr> </table>	0	78	1	29	2	120	3	123	4	71	5	150	6	162	7	173	8	18	9	21	10	33	11	28	12	19	13	200	14	210	15	225
1	0	78																																						
		29																																						
0	78																																							
1	29																																							
2	120																																							
3	123																																							
4	71																																							
5	150																																							
6	162																																							
7	173																																							
8	18																																							
9	21																																							
10	33																																							
11	28																																							
12	19																																							
13	200																																							
14	210																																							
15	225																																							

© Kavita Bala, Computer Science, Cornell University

Write-Through (REF 2)

Processor	Cache	Memory																																						
<p>➔ Ld R1 ← M[1] Ld R2 ← M[7] St R2 → M[0] St R1 → M[5] Ld R2 ← M[10] St R1 → M[5] St R1 → M[10]</p> <p>R0 R1 29 R2 R3 </p>	<p style="text-align: center;">V tag data</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td style="width: 10px; text-align: center;">1</td> <td style="width: 10px; text-align: center;">0</td> <td style="width: 20px; text-align: center;">78</td> </tr> <tr> <td></td> <td></td> <td style="text-align: center;">29</td> </tr> </table> <p>lru 0 </p> <p style="text-align: center;">Misses: 1 Hits: 0</p>	1	0	78			29	<table border="1" style="width: 100%; text-align: center;"> <tr><td>0</td><td>78</td></tr> <tr><td>1</td><td>29</td></tr> <tr><td>2</td><td>120</td></tr> <tr><td>3</td><td>123</td></tr> <tr><td>4</td><td>71</td></tr> <tr><td>5</td><td>150</td></tr> <tr><td>6</td><td>162</td></tr> <tr><td>7</td><td>173</td></tr> <tr><td>8</td><td>18</td></tr> <tr><td>9</td><td>21</td></tr> <tr><td>10</td><td>33</td></tr> <tr><td>11</td><td>28</td></tr> <tr><td>12</td><td>19</td></tr> <tr><td>13</td><td>200</td></tr> <tr><td>14</td><td>210</td></tr> <tr><td>15</td><td>225</td></tr> </table>	0	78	1	29	2	120	3	123	4	71	5	150	6	162	7	173	8	18	9	21	10	33	11	28	12	19	13	200	14	210	15	225
1	0	78																																						
		29																																						
0	78																																							
1	29																																							
2	120																																							
3	123																																							
4	71																																							
5	150																																							
6	162																																							
7	173																																							
8	18																																							
9	21																																							
10	33																																							
11	28																																							
12	19																																							
13	200																																							
14	210																																							
15	225																																							

© Kavita Bala, Computer Science, Cornell University

Write-Through (REF 2)

Processor	Cache	Memory																																												
<p> Ld R1 ← M[1] Ld R2 ← M[7] St R2 → M[0] St R1 → M[5] Ld R2 ← M[10] St R1 → M[5] St R1 → M[10] </p> <p> R0 R1 29 R2 173 R3 </p>	<p style="text-align: center;">V tag data</p> <p>lru <table border="1" style="display: inline-table; border-collapse: collapse;"> <tr><td style="padding: 2px;">1</td><td style="padding: 2px;">0</td><td style="padding: 2px;">78</td></tr> <tr><td style="padding: 2px;"></td><td style="padding: 2px;"></td><td style="padding: 2px;">29</td></tr> <tr><td style="padding: 2px;">1</td><td style="padding: 2px;">3</td><td style="padding: 2px;">162</td></tr> <tr><td style="padding: 2px;"></td><td style="padding: 2px;"></td><td style="padding: 2px;">173</td></tr> </table></p> <p style="text-align: center;">Misses: 2 Hits: 0</p>	1	0	78			29	1	3	162			173	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td style="padding: 2px;">0</td><td style="padding: 2px;">78</td></tr> <tr><td style="padding: 2px;">1</td><td style="padding: 2px;">29</td></tr> <tr><td style="padding: 2px;">2</td><td style="padding: 2px;">120</td></tr> <tr><td style="padding: 2px;">3</td><td style="padding: 2px;">123</td></tr> <tr><td style="padding: 2px;">4</td><td style="padding: 2px;">71</td></tr> <tr><td style="padding: 2px;">5</td><td style="padding: 2px;">150</td></tr> <tr><td style="padding: 2px;">6</td><td style="padding: 2px;">162</td></tr> <tr><td style="padding: 2px;">7</td><td style="padding: 2px;">173</td></tr> <tr><td style="padding: 2px;">8</td><td style="padding: 2px;">18</td></tr> <tr><td style="padding: 2px;">9</td><td style="padding: 2px;">21</td></tr> <tr><td style="padding: 2px;">10</td><td style="padding: 2px;">33</td></tr> <tr><td style="padding: 2px;">11</td><td style="padding: 2px;">28</td></tr> <tr><td style="padding: 2px;">12</td><td style="padding: 2px;">19</td></tr> <tr><td style="padding: 2px;">13</td><td style="padding: 2px;">200</td></tr> <tr><td style="padding: 2px;">14</td><td style="padding: 2px;">210</td></tr> <tr><td style="padding: 2px;">15</td><td style="padding: 2px;">225</td></tr> </table>	0	78	1	29	2	120	3	123	4	71	5	150	6	162	7	173	8	18	9	21	10	33	11	28	12	19	13	200	14	210	15	225
1	0	78																																												
		29																																												
1	3	162																																												
		173																																												
0	78																																													
1	29																																													
2	120																																													
3	123																																													
4	71																																													
5	150																																													
6	162																																													
7	173																																													
8	18																																													
9	21																																													
10	33																																													
11	28																																													
12	19																																													
13	200																																													
14	210																																													
15	225																																													

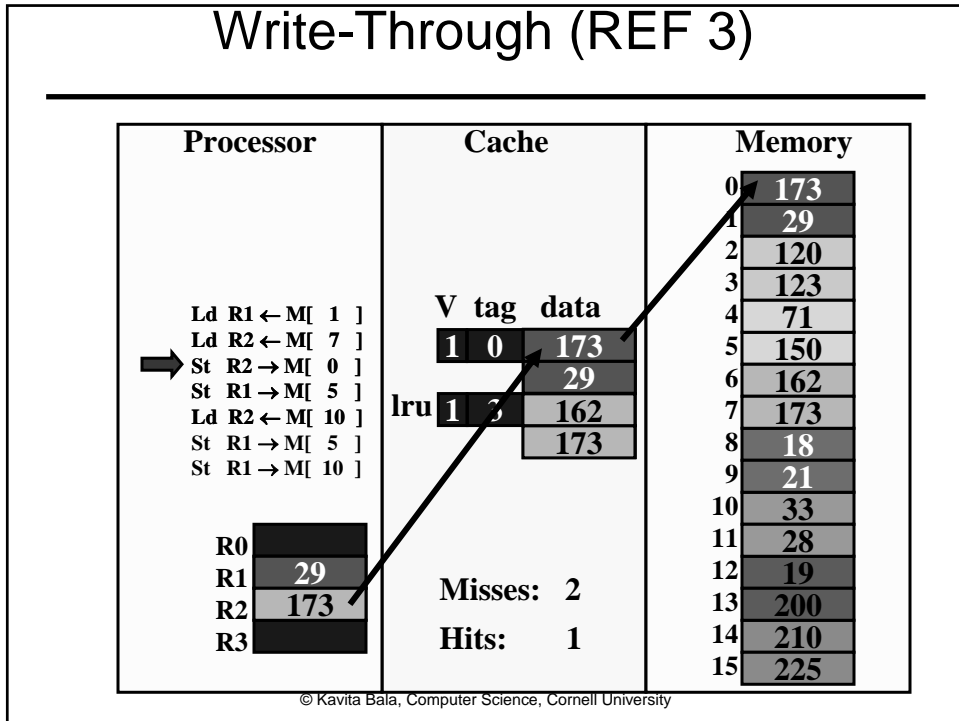
© Kavita Bala, Computer Science, Cornell University

Write-Through (REF 3)

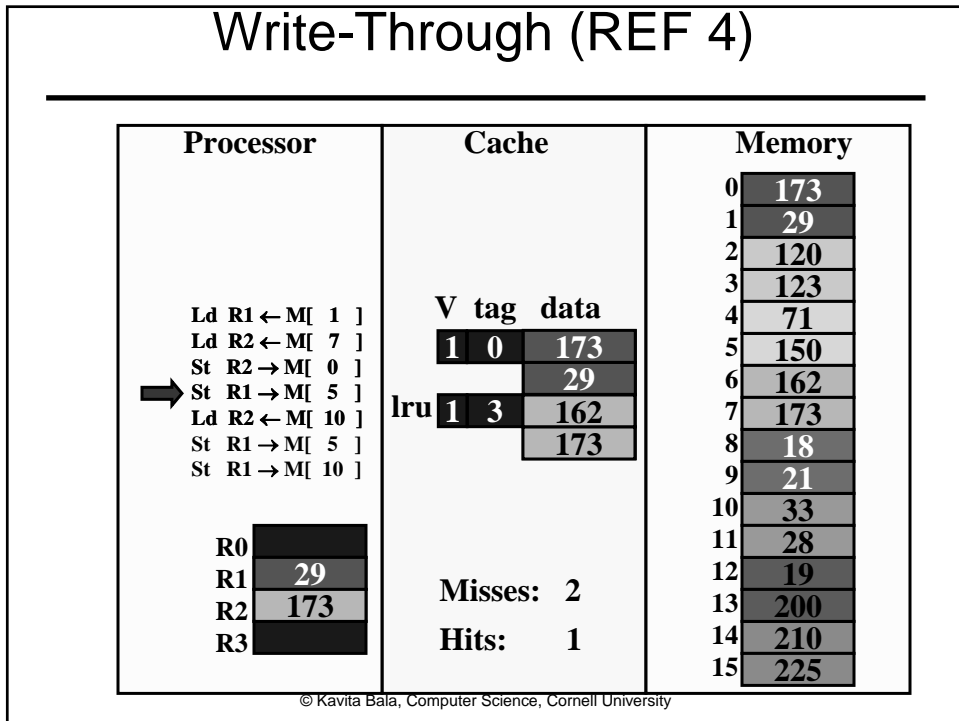
Processor	Cache	Memory																																												
<p> Ld R1 ← M[1] Ld R2 ← M[7] St R2 → M[0] St R1 → M[5] Ld R2 ← M[10] St R1 → M[5] St R1 → M[10] </p> <p> R0 R1 29 R2 173 R3 </p>	<p style="text-align: center;">V tag data</p> <p>lru <table border="1" style="display: inline-table; border-collapse: collapse;"> <tr><td style="padding: 2px;">1</td><td style="padding: 2px;">0</td><td style="padding: 2px;">78</td></tr> <tr><td style="padding: 2px;"></td><td style="padding: 2px;"></td><td style="padding: 2px;">29</td></tr> <tr><td style="padding: 2px;">1</td><td style="padding: 2px;">3</td><td style="padding: 2px;">162</td></tr> <tr><td style="padding: 2px;"></td><td style="padding: 2px;"></td><td style="padding: 2px;">173</td></tr> </table></p> <p style="text-align: center;">Misses: 2 Hits: 0</p>	1	0	78			29	1	3	162			173	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td style="padding: 2px;">0</td><td style="padding: 2px;">78</td></tr> <tr><td style="padding: 2px;">1</td><td style="padding: 2px;">29</td></tr> <tr><td style="padding: 2px;">2</td><td style="padding: 2px;">120</td></tr> <tr><td style="padding: 2px;">3</td><td style="padding: 2px;">123</td></tr> <tr><td style="padding: 2px;">4</td><td style="padding: 2px;">71</td></tr> <tr><td style="padding: 2px;">5</td><td style="padding: 2px;">150</td></tr> <tr><td style="padding: 2px;">6</td><td style="padding: 2px;">162</td></tr> <tr><td style="padding: 2px;">7</td><td style="padding: 2px;">173</td></tr> <tr><td style="padding: 2px;">8</td><td style="padding: 2px;">18</td></tr> <tr><td style="padding: 2px;">9</td><td style="padding: 2px;">21</td></tr> <tr><td style="padding: 2px;">10</td><td style="padding: 2px;">33</td></tr> <tr><td style="padding: 2px;">11</td><td style="padding: 2px;">28</td></tr> <tr><td style="padding: 2px;">12</td><td style="padding: 2px;">19</td></tr> <tr><td style="padding: 2px;">13</td><td style="padding: 2px;">200</td></tr> <tr><td style="padding: 2px;">14</td><td style="padding: 2px;">210</td></tr> <tr><td style="padding: 2px;">15</td><td style="padding: 2px;">225</td></tr> </table>	0	78	1	29	2	120	3	123	4	71	5	150	6	162	7	173	8	18	9	21	10	33	11	28	12	19	13	200	14	210	15	225
1	0	78																																												
		29																																												
1	3	162																																												
		173																																												
0	78																																													
1	29																																													
2	120																																													
3	123																																													
4	71																																													
5	150																																													
6	162																																													
7	173																																													
8	18																																													
9	21																																													
10	33																																													
11	28																																													
12	19																																													
13	200																																													
14	210																																													
15	225																																													

© Kavita Bala, Computer Science, Cornell University

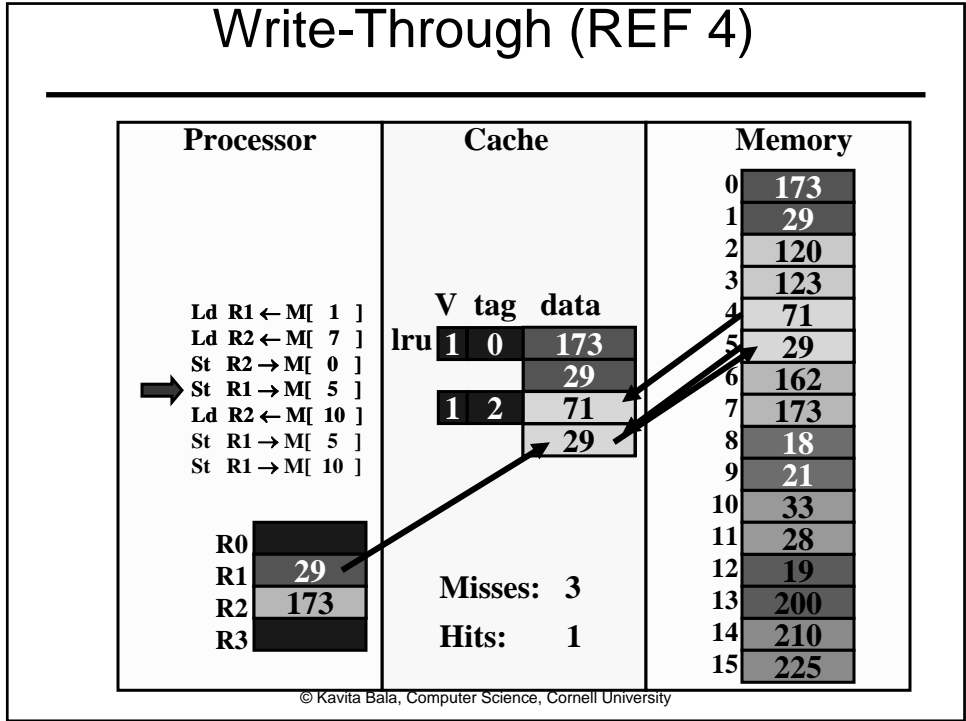
Write-Through (REF 3)



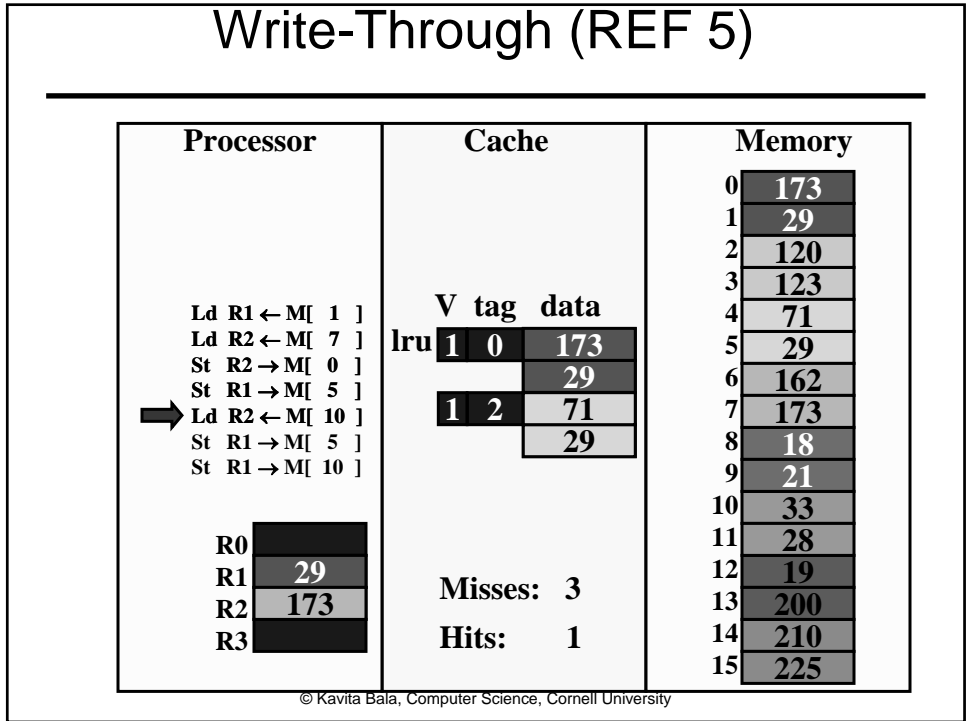
Write-Through (REF 4)



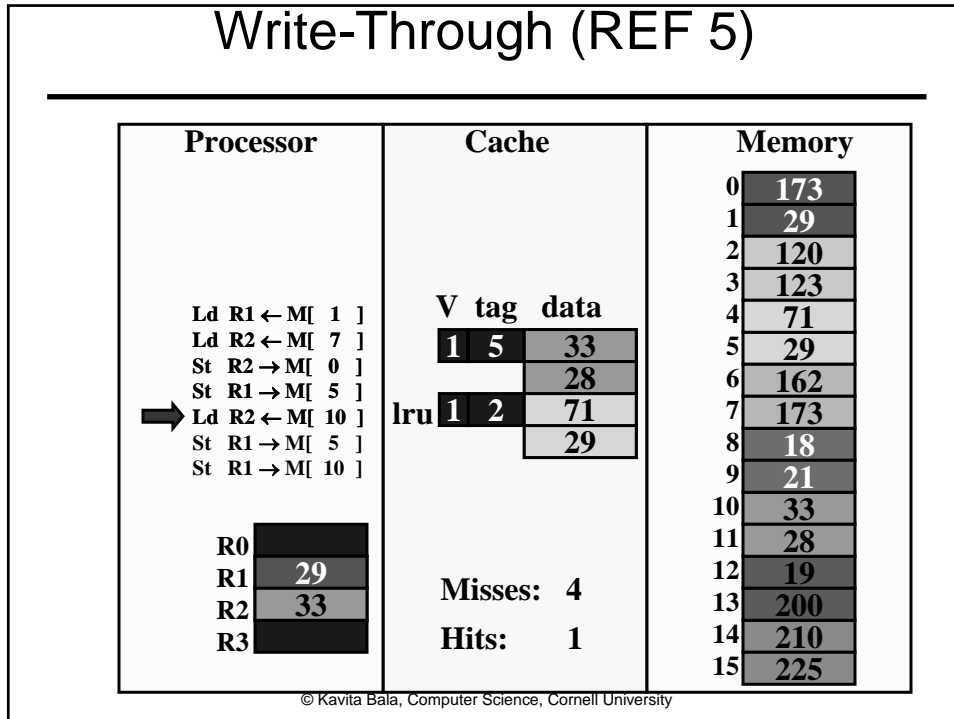
Write-Through (REF 4)



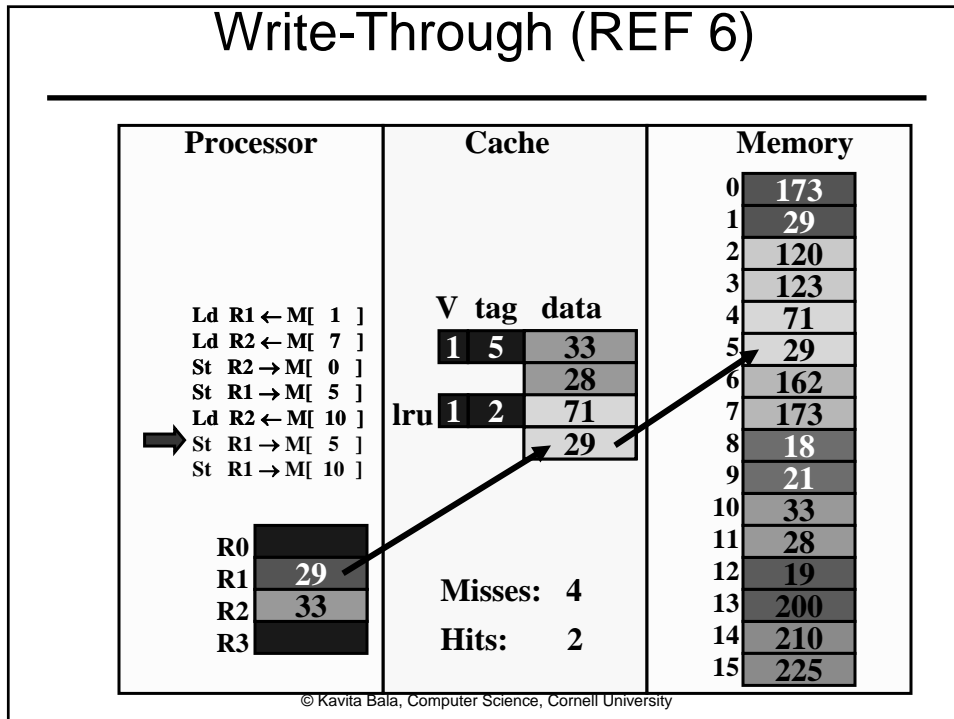
Write-Through (REF 5)



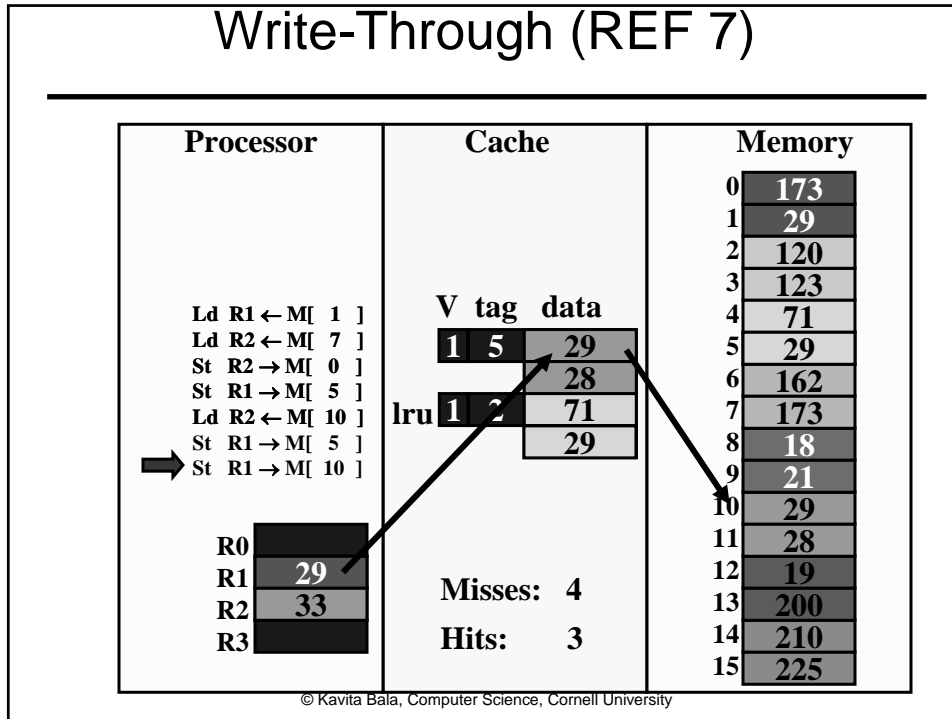
Write-Through (REF 5)



Write-Through (REF 6)



Write-Through (REF 7)

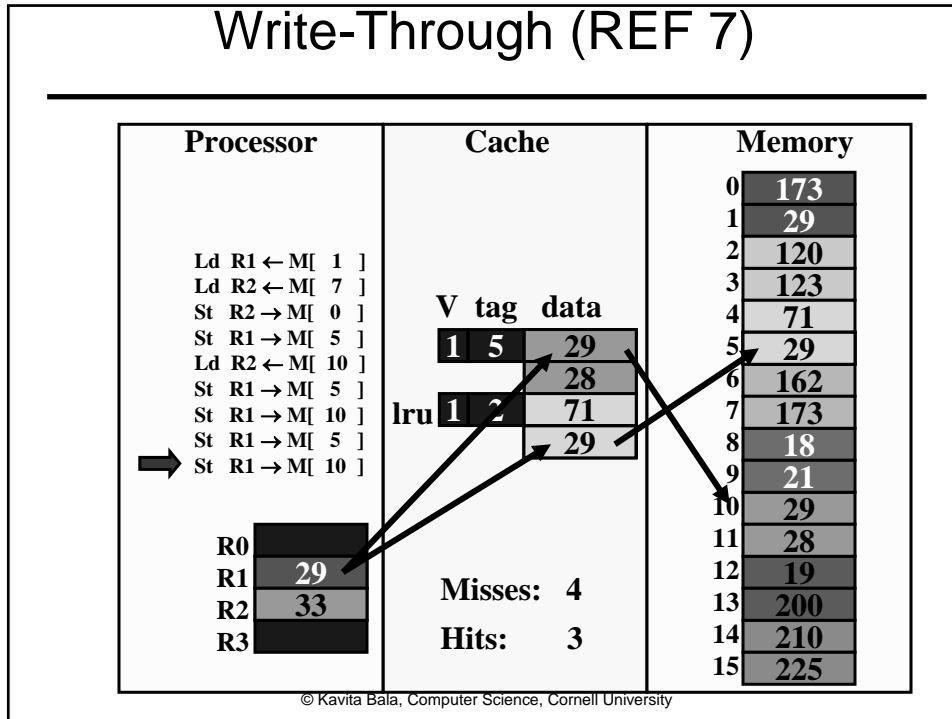


How Many Memory References?

- Each miss reads a block
(only two words in this cache)
- Each store writes a word
- Total reads: eight words
- Total writes: four words

but caches generally miss < 20%
usually much lower miss rates . . . but depends
on both cache and application!

Write-Through (REF 7)



How Many Memory References?

- Each miss reads a block
(only two words in this cache)
- Each store writes a word
- Total reads: eight words
- Total writes: six words, eight words, etc.

but caches generally miss < 20%
usually much lower miss rates . . . but depends
on both cache and application!

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)

© Kavita Bala, Computer Science, Cornell University

Dirty Bits and Write-Back Buffers

V	D	Tag	Data Byte 0, Byte 1 ... Byte N	Line
1	0			
1	1			
1	0			

- Dirty bits indicate which lines have been written
- Dirty bits enable the cache to handle multiple writes to the same cache line without having to go to memory
- Dirty bit reset when line is allocated
- Set when block is written
- Write-back buffer
 - A queue where dirty lines are placed
 - Items added to the end as dirty lines are evicted from the cache
 - Items removed from the front as memory writes are completed

© Kavita Bala, Computer Science, Cornell University

Handling Stores (Write-Back)

Processor	Cache	Memory
Ld R1 ← M[1]	V d tag data	0 78
Ld R2 ← M[7]	0	1 29
St R2 → M[0]		2 120
St R1 → M[5]	0	3 123
Ld R2 ← M[10]		4 71
St R1 → M[5]		5 150
St R1 → M[10]		6 162
		7 173
		8 18
		9 21
		10 33
		11 28
		12 19
		13 200
		14 210
		15 225
R0	Misses: 0	
R1	Hits: 0	
R2		
R3		

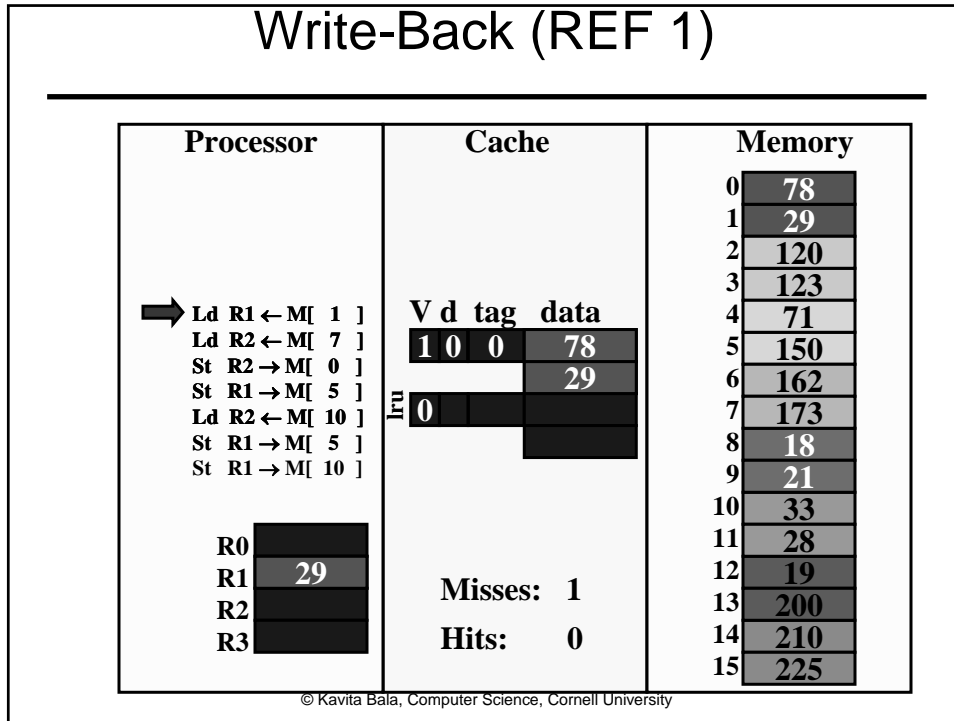
© Kavita Bala, Computer Science, Cornell University

Write-Back (REF 1)

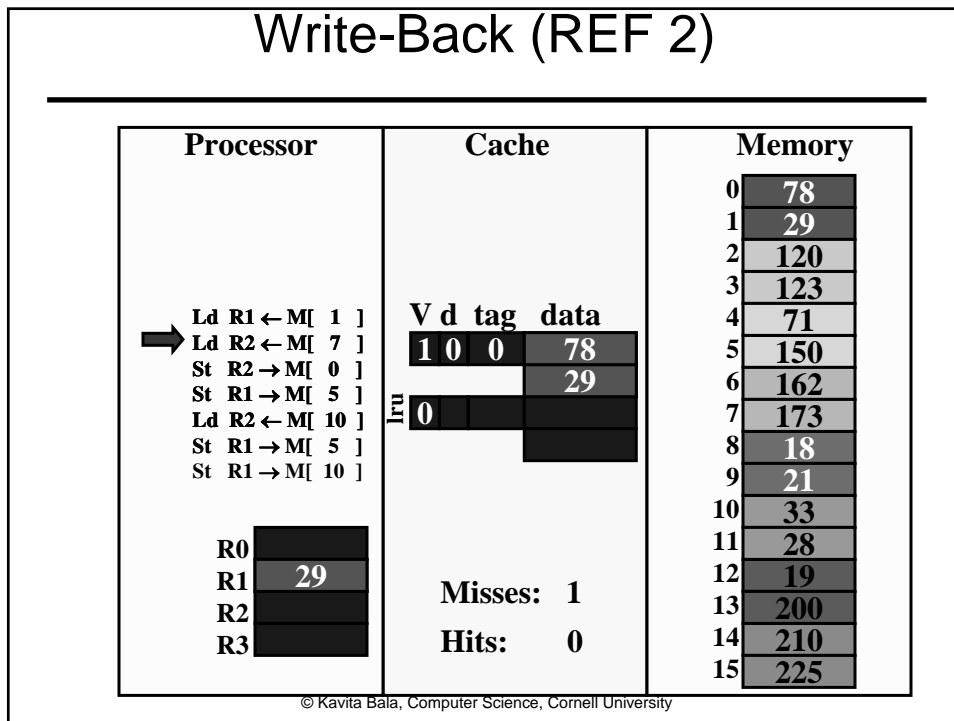
Processor	Cache	Memory
→ Ld R1 ← M[1]	V d tag data	0 78
Ld R2 ← M[7]	0	1 29
St R2 → M[0]		2 120
St R1 → M[5]	0	3 123
Ld R2 ← M[10]		4 71
St R1 → M[5]		5 150
St R1 → M[10]		6 162
		7 173
		8 18
		9 21
		10 33
		11 28
		12 19
		13 200
		14 210
		15 225
R0	Misses: 0	
R1	Hits: 0	
R2		
R3		

© Kavita Bala, Computer Science, Cornell University

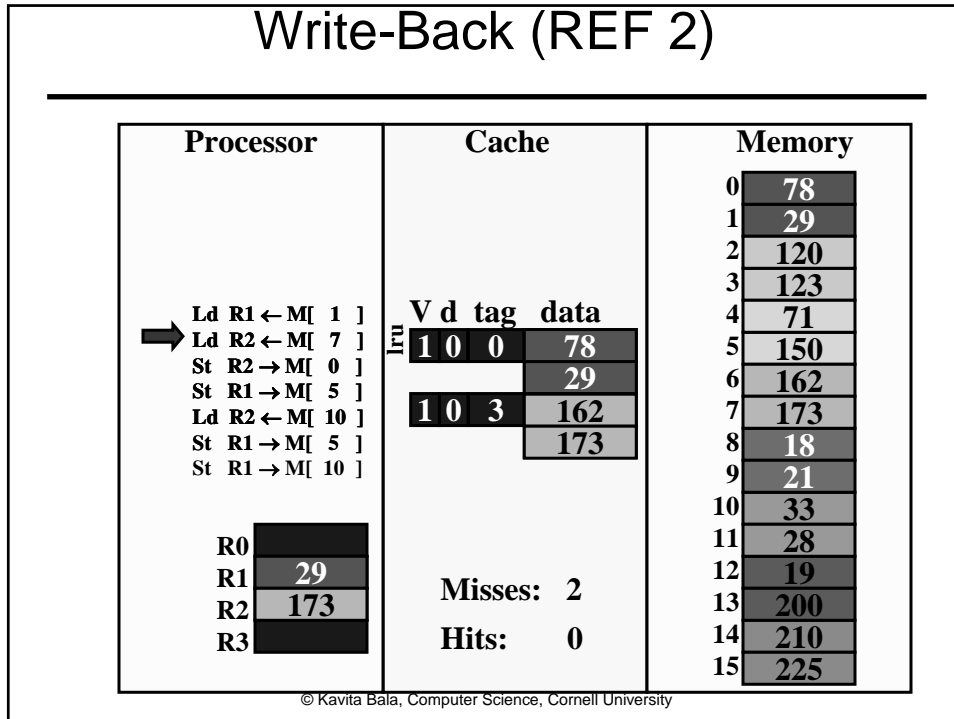
Write-Back (REF 1)



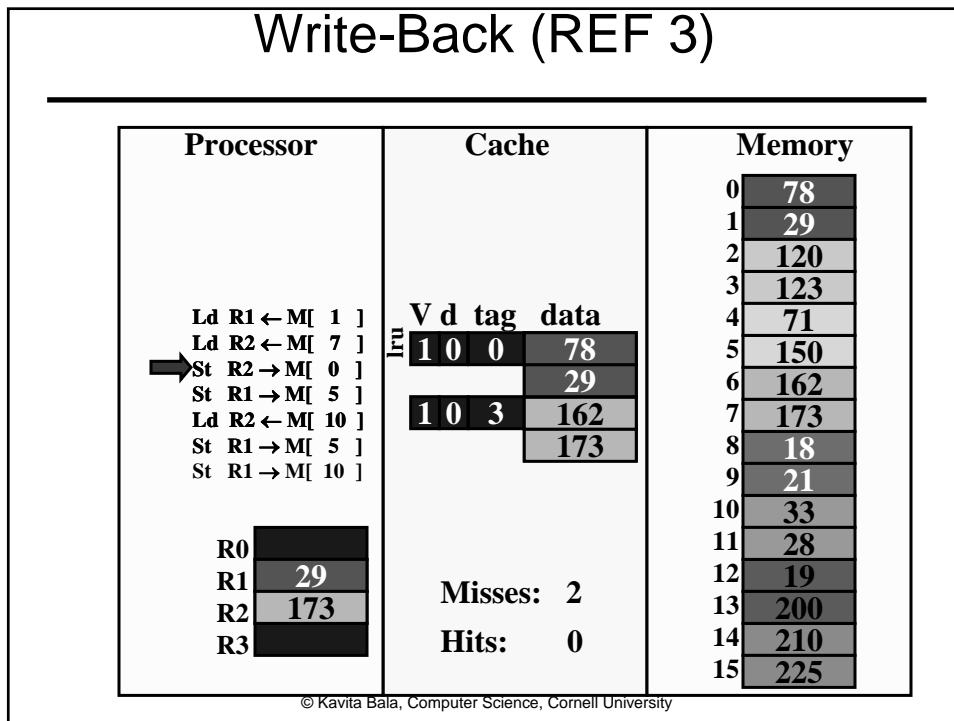
Write-Back (REF 2)



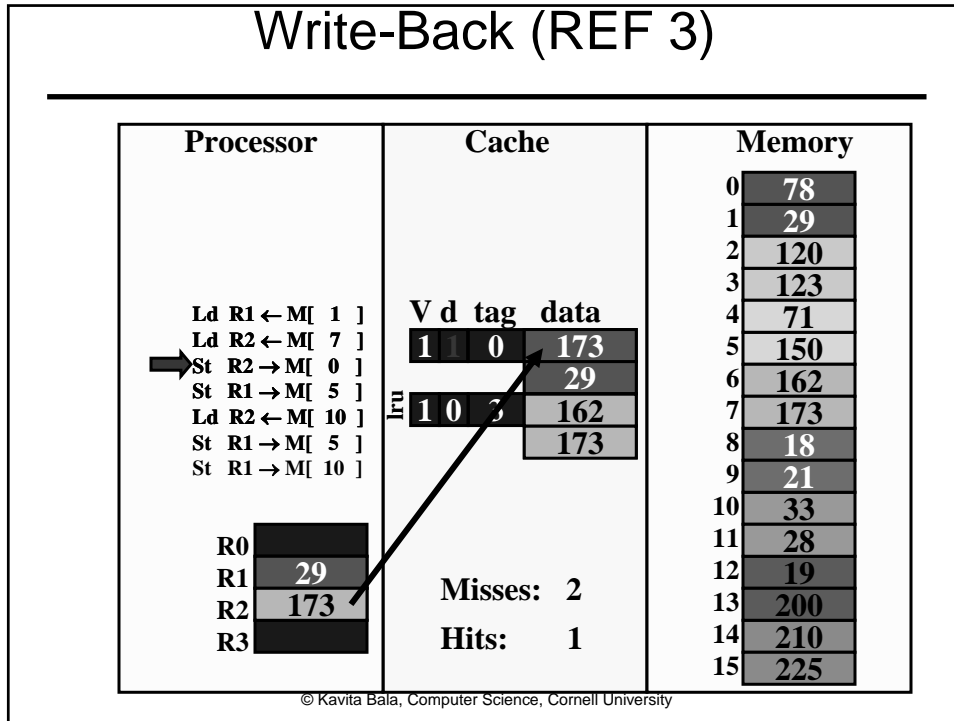
Write-Back (REF 2)



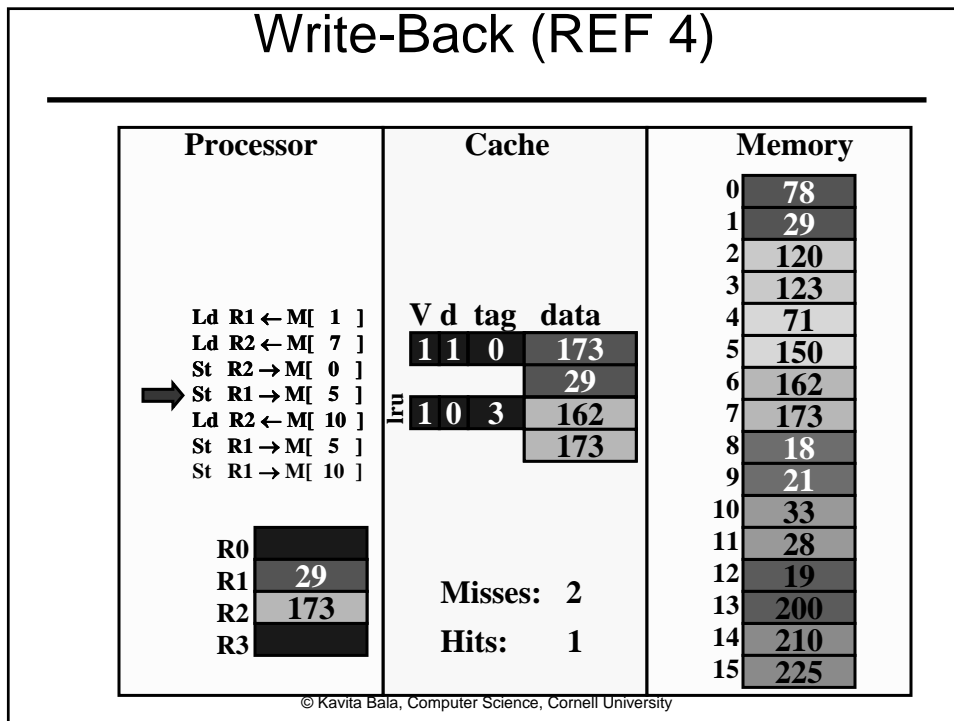
Write-Back (REF 3)



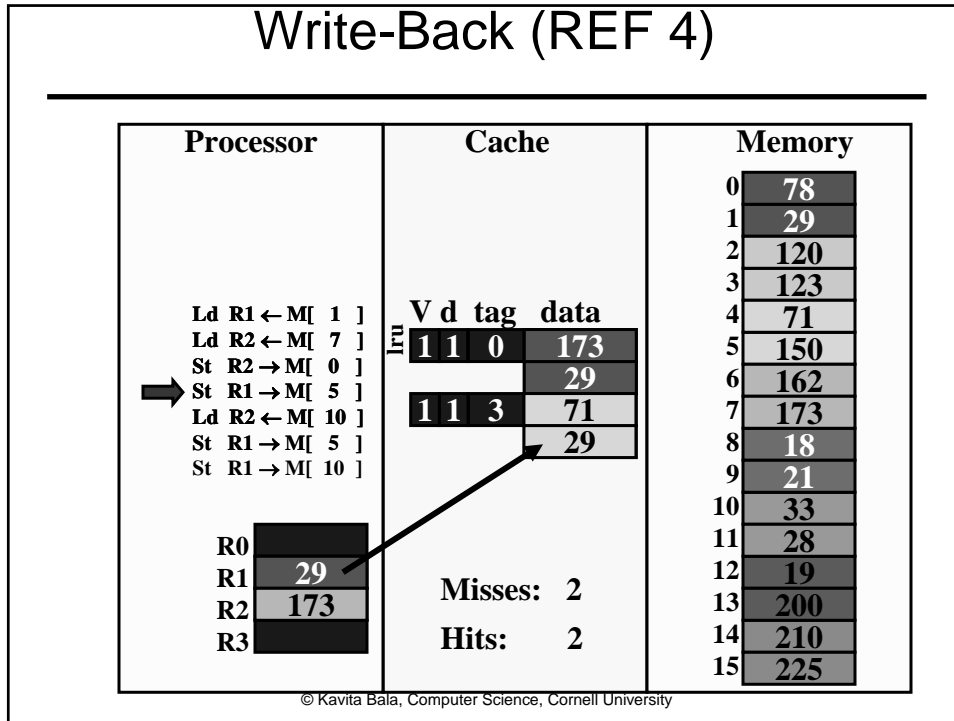
Write-Back (REF 3)



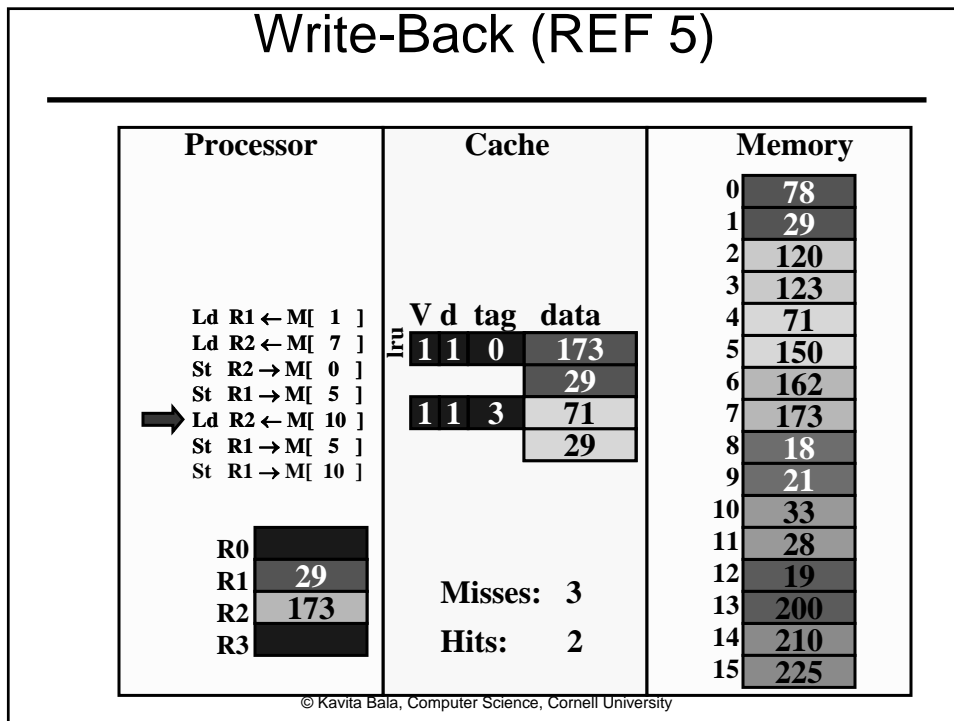
Write-Back (REF 4)



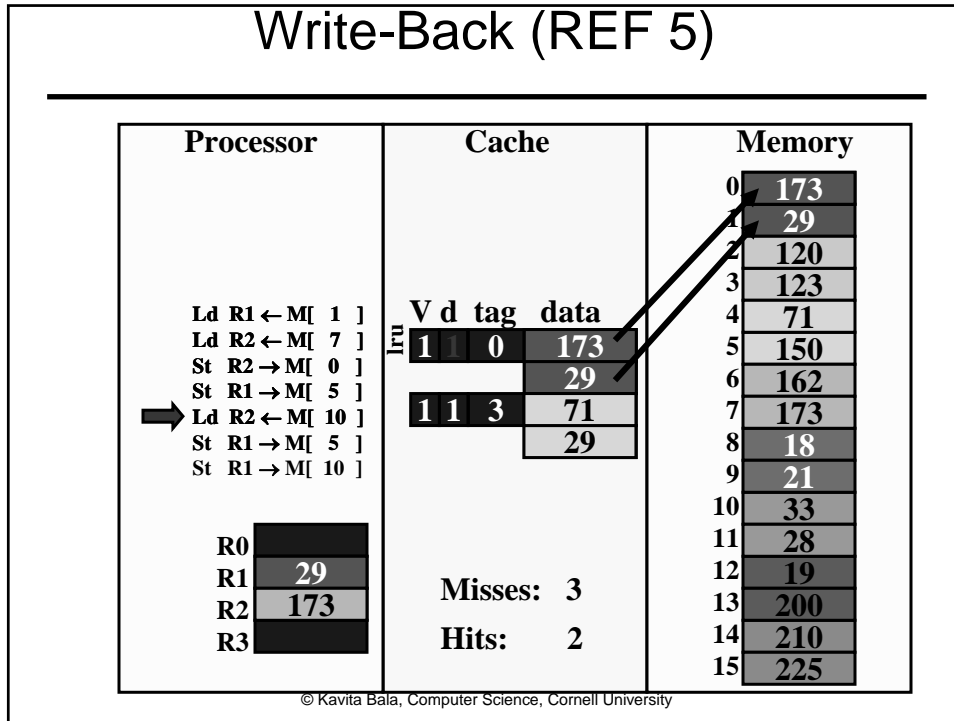
Write-Back (REF 4)



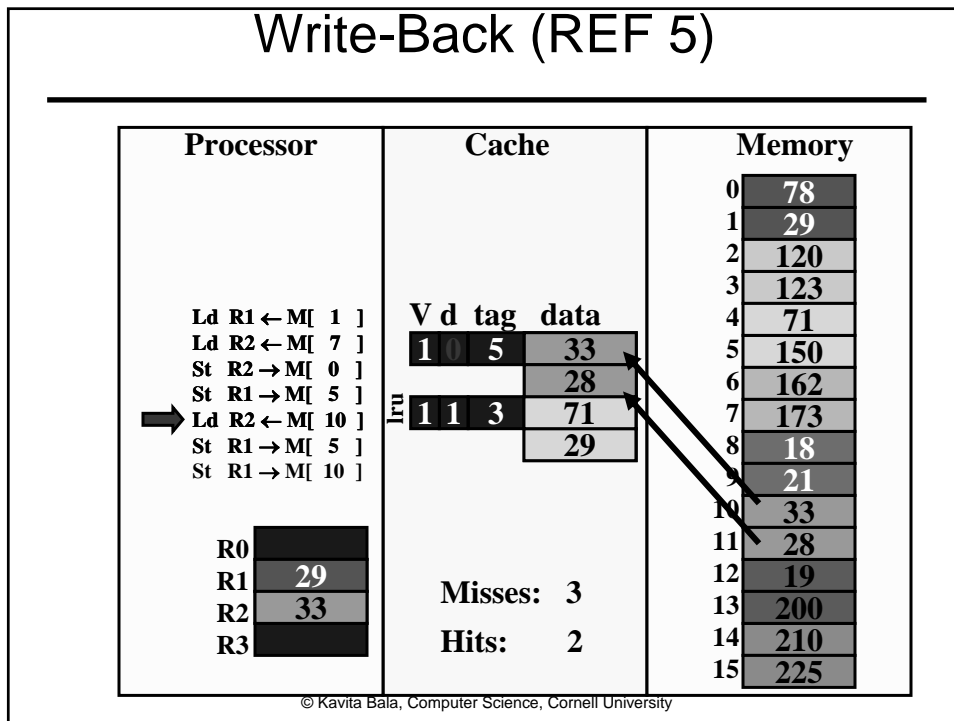
Write-Back (REF 5)



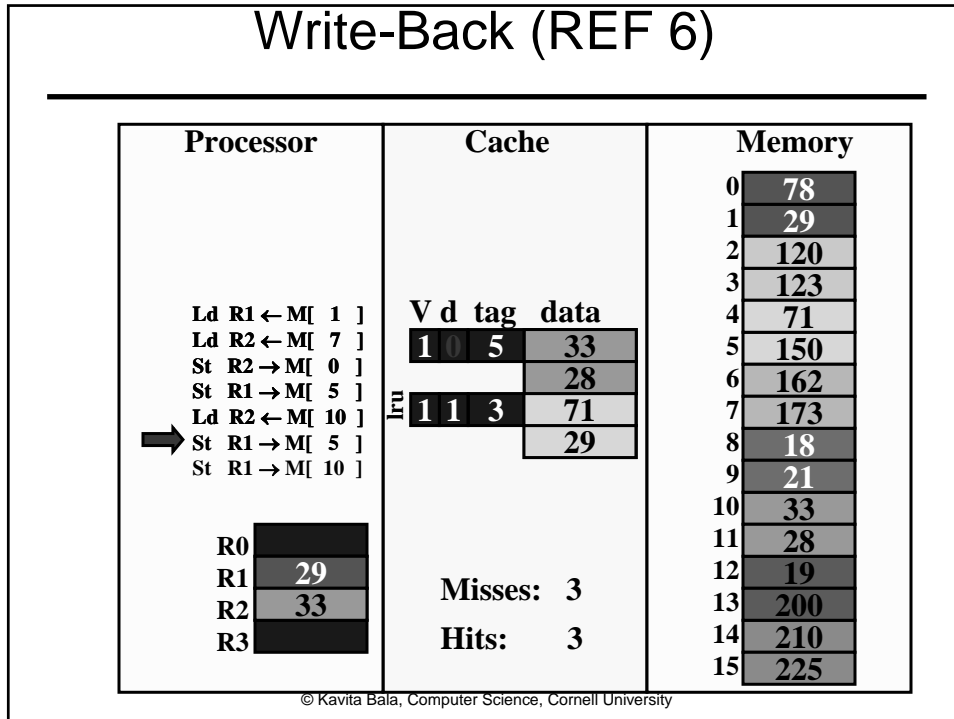
Write-Back (REF 5)



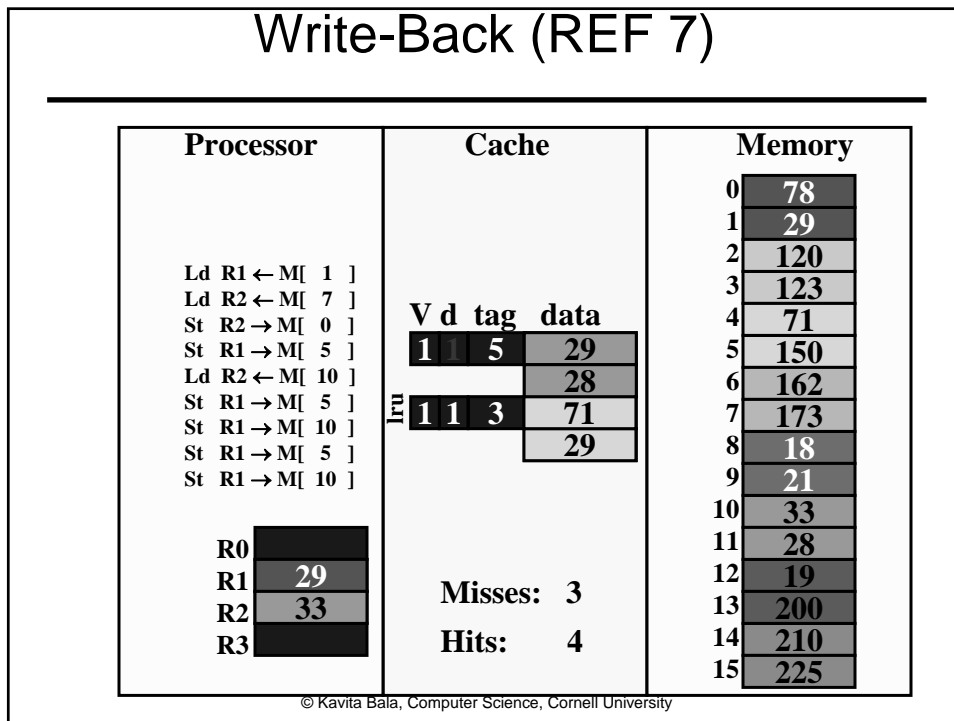
Write-Back (REF 5)



Write-Back (REF 6)



Write-Back (REF 7)



How many memory references?

- Each miss reads a block
 - Two words in this cache
- Each evicted dirty cache line writes a block
- Total reads: six words
- Total writes: 4/6 words (after final eviction)

© Kavita Bala, Computer Science, Cornell University

Write-Back (REF 8,9)

Processor	Cache	Memory																								
Ld R1 ← M[1]	<table border="1"> <thead> <tr> <th>V</th> <th>d</th> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1</td> <td></td> <td>5</td> <td>29</td> </tr> <tr> <td></td> <td></td> <td></td> <td>28</td> </tr> <tr> <td>lru</td> <td>1</td> <td>1</td> <td>3</td> </tr> <tr> <td></td> <td></td> <td></td> <td>71</td> </tr> <tr> <td></td> <td></td> <td></td> <td>29</td> </tr> </tbody> </table>	V	d	tag	data	1		5	29				28	lru	1	1	3				71				29	0 78
V		d	tag	data																						
1		5	29																							
			28																							
lru	1	1	3																							
			71																							
			29																							
Ld R2 ← M[7]		1 29																								
St R2 → M[0]		2 120																								
St R1 → M[5]		3 123																								
Ld R2 ← M[10]		4 71																								
St R1 → M[5]		5 150																								
St R1 → M[10]		6 162																								
St R1 → M[10]		7 173																								
St R1 → M[5]		8 18																								
St R1 → M[10]		9 21																								
		10 33																								
		11 28																								
		12 19																								
		13 200																								
		14 210																								
		15 225																								
<table border="1"> <tbody> <tr><td>R0</td><td></td></tr> <tr><td>R1</td><td>29</td></tr> <tr><td>R2</td><td>33</td></tr> <tr><td>R3</td><td></td></tr> </tbody> </table>	R0		R1	29	R2	33	R3		Misses: 3 Hits: 4																	
R0																										
R1	29																									
R2	33																									
R3																										

© Kavita Bala, Computer Science, Cornell University

How many memory references?

- Each miss reads a block
 - Two words in this cache
- Each evicted dirty cache line writes a block
- Total reads: six words
- Total writes: 4/6 words (after final eviction)
- By comparison write-through was
 - Reads: eight words
 - Writes: 6/8/10 etc words
- Write-through or Write-back?

© Kavita Bala, Computer Science, Cornell University

Write-through vs. Write-back

- Write-through is slower
 - But cleaner (memory always consistent)
- Write-back is faster
 - But complicated when multi cores sharing memory

© Kavita Bala, Computer Science, Cornell University