

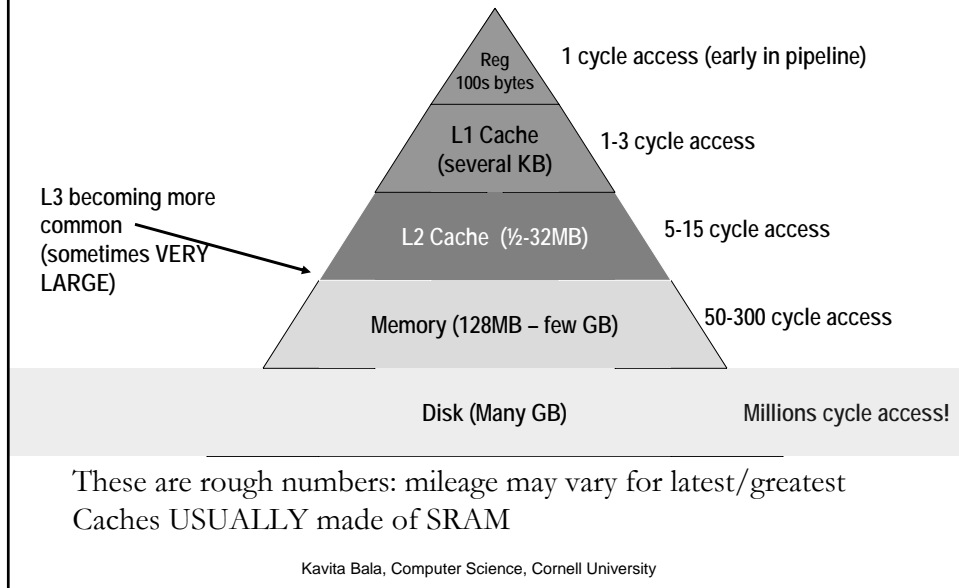
CS 316: Caches-II

Kavita Bala
Fall 2007
Computer Science
Cornell University

Announcements

- HW 2 is out
- Caches and Memory: Chapter 7 (H & P)

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

Cache Organization

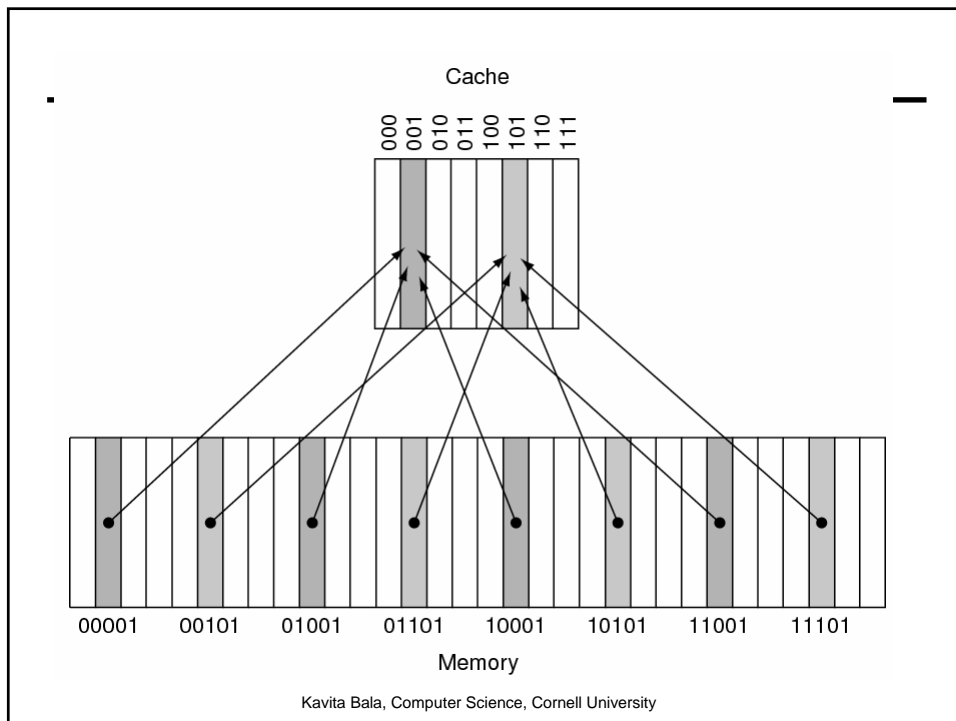
- Three common designs
 - Fully associative: Block can be anywhere in the cache
 - Direct mapped: Block can only be in one line in the cache
 - Set-associative: Block can be in a few (2 to 8) places in the cache

Kavita Bala, Computer Science, Cornell University

Direct Mapped Cache

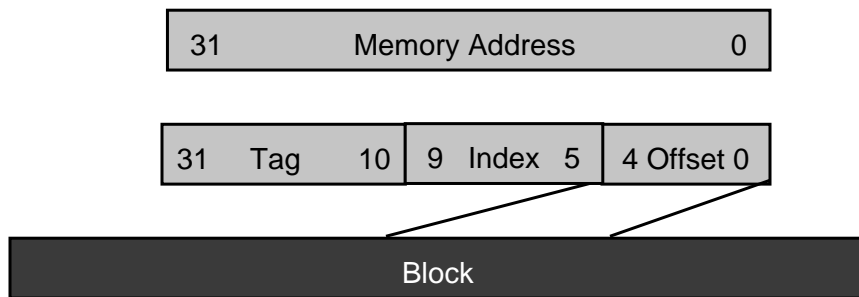
- Simplest
- Block can only be in one line in the cache
- How to determine this location?
 - Use modulo arithmetic
 - (Block address) modulo (# cache blocks)
 - For power of 2, use log (cache size in blocks)

Kavita Bala, Computer Science, Cornell University



Tags and Offsets

- Tag: matching
- Offset: within block
- Valid bit: is the data valid?



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

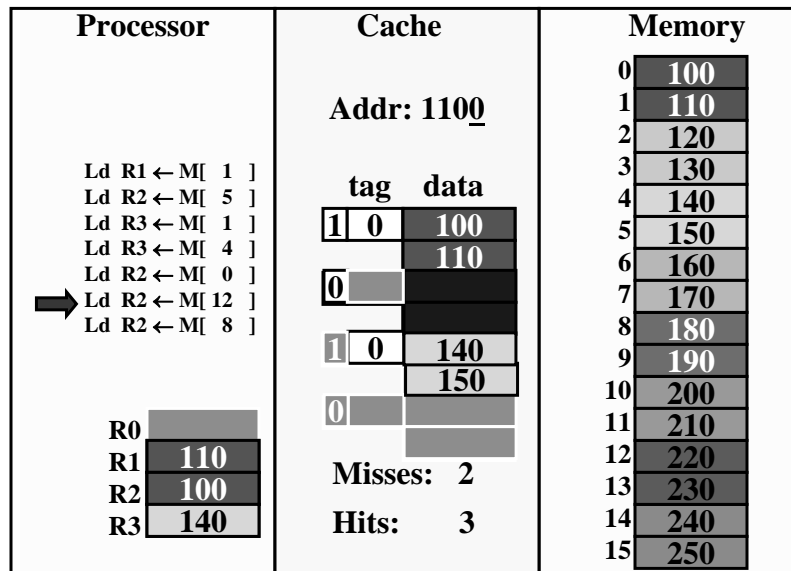
Cache Size

- Cache of size 2^n blocks (index: n bits)
- Block size of 2^m word (block index: $m+2$ bits)
- Tag field: $32 - (n + m + 2)$
- Valid bit: 1

- Bits in cache: $2^n \times (\text{block size} + \text{tag size} + \text{valid bit size})$
 $= 2^n (2^m \times 32 + (32 - n - m - 2) + 1)$

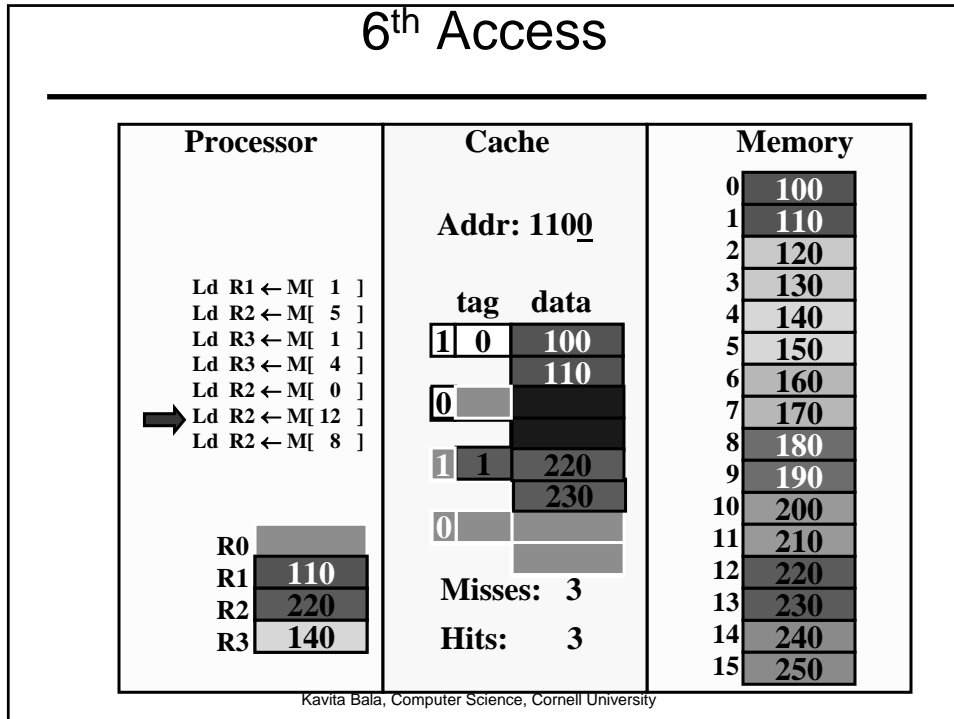
Kavita Bala, Computer Science, Cornell University

6th Access

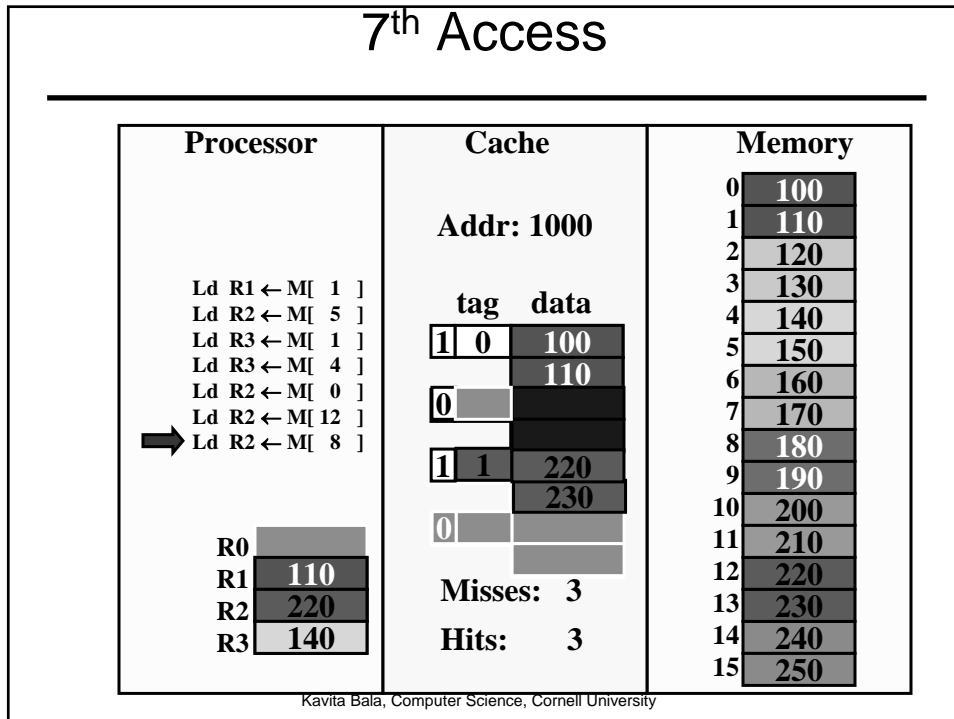


Kavita Bala, Computer Science, Cornell University

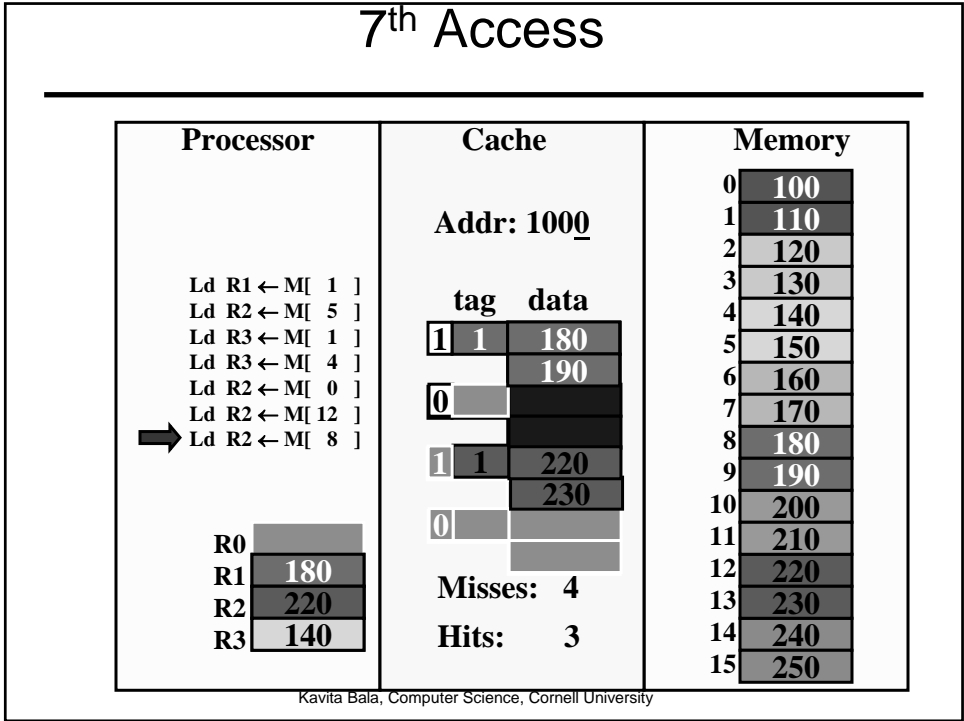
6th Access



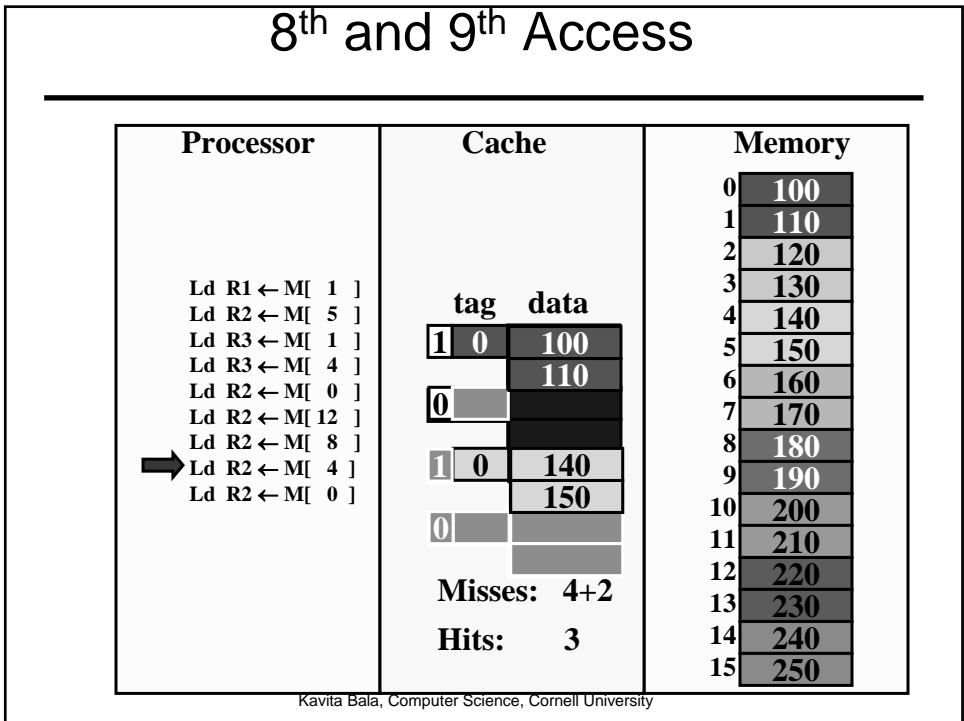
7th Access



7th Access



8th and 9th Access



DM: 10th and 11th Access

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

tag		data
1	1	180
0		190
1	1	220
0		230

Misses: 4+2+2
Hits: 3

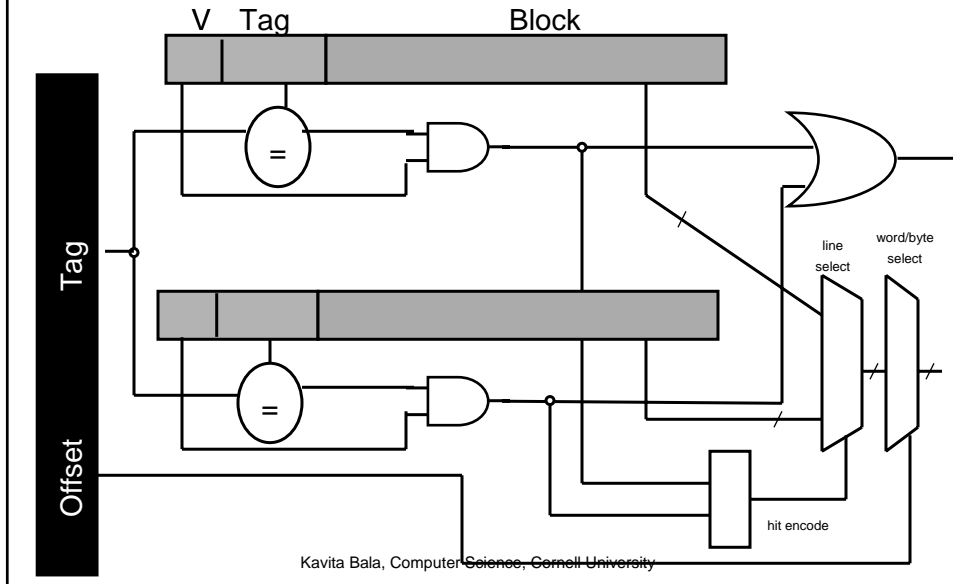
Kavita Bala, Computer Science, Cornell University

Misses

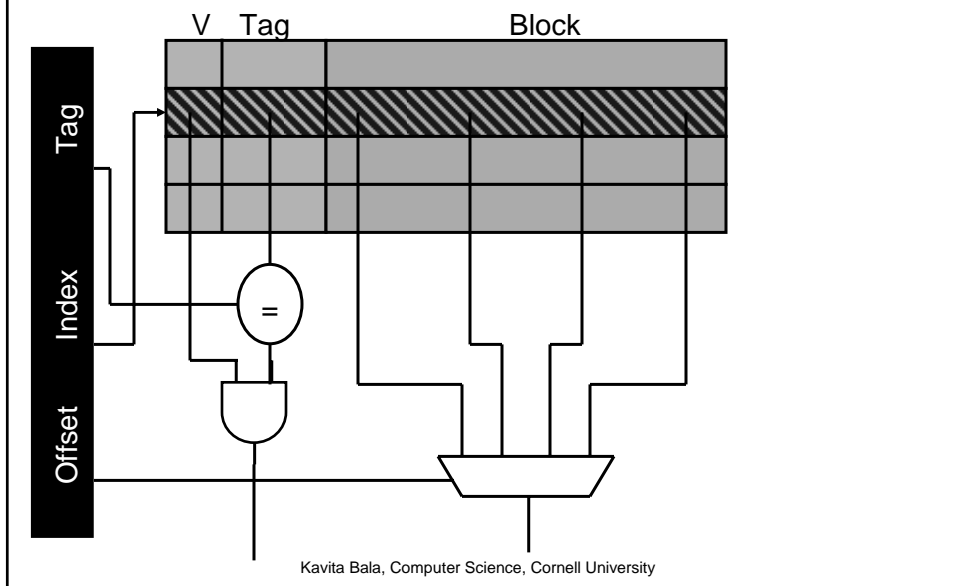
- Three types of misses
 - Cold
 - The line is being referenced for the first time
 - Capacity
 - The line was evicted because the cache was not large enough
 - Conflict
 - The line was evicted because of another access whose index conflicted

Kavita Bala, Computer Science, Cornell University

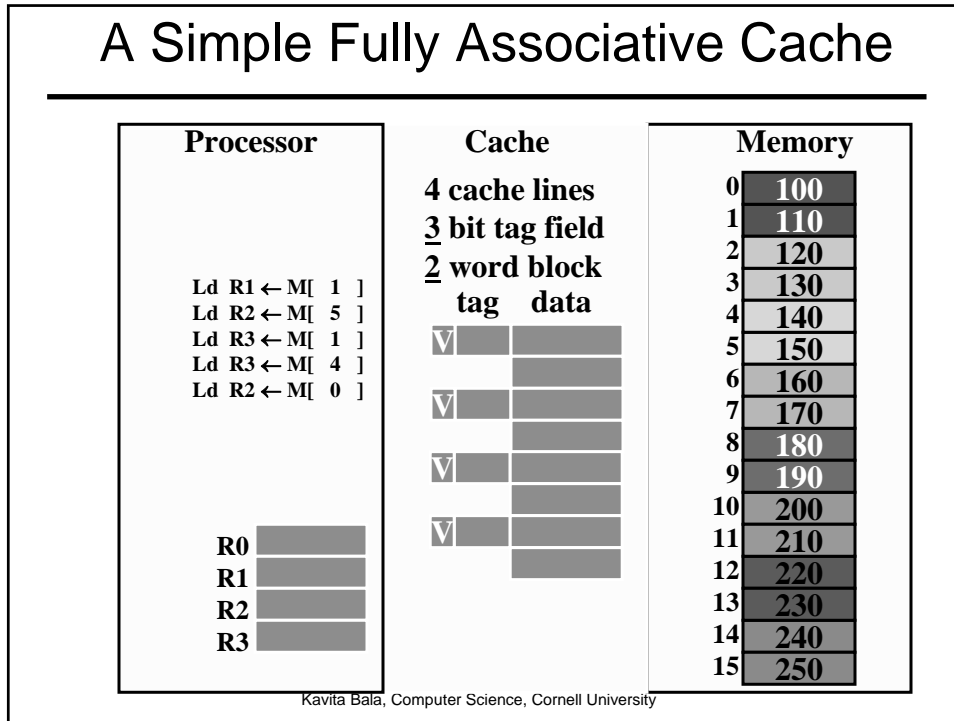
Fully Associative Cache



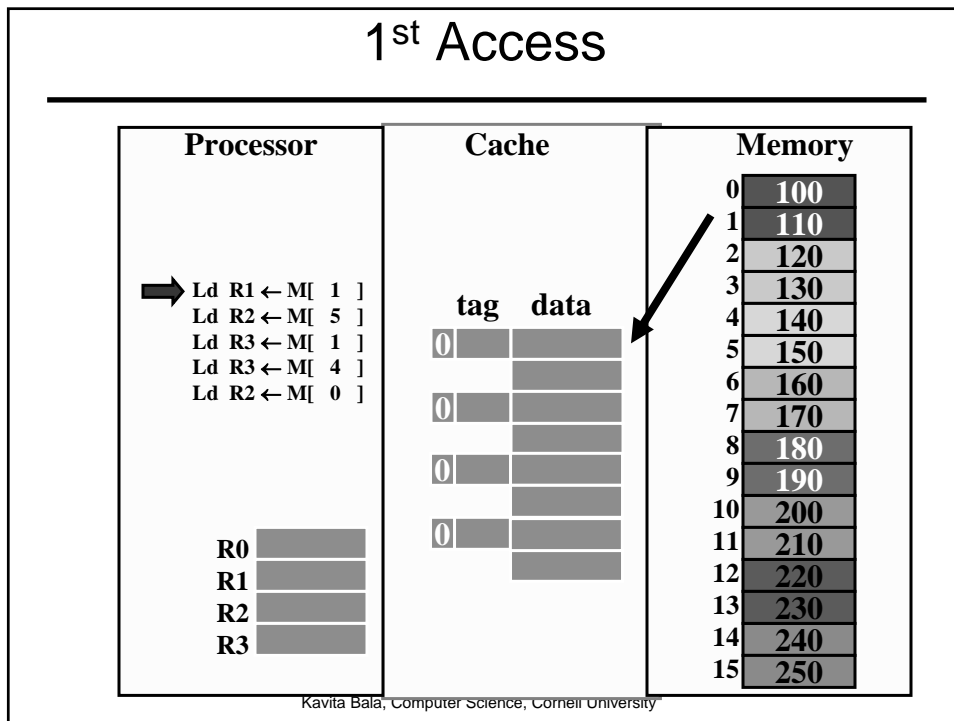
Comparison: Direct Mapped Cache



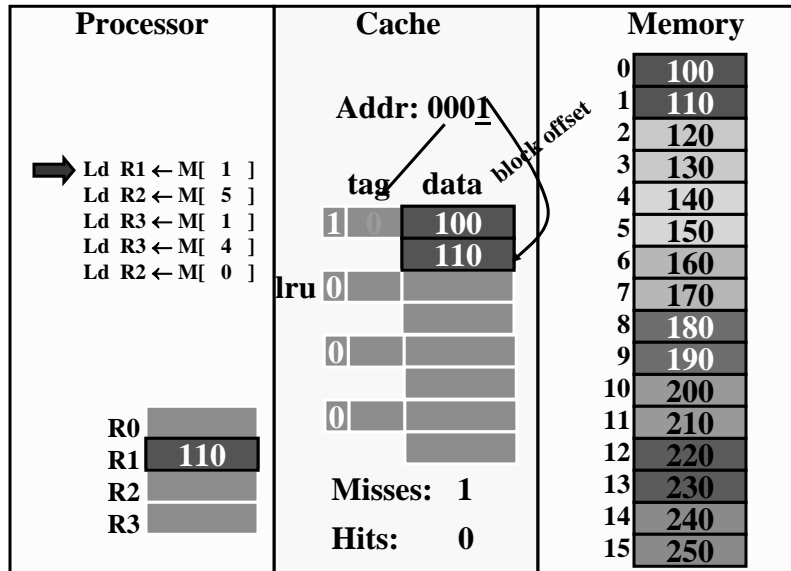
A Simple Fully Associative Cache



1st Access

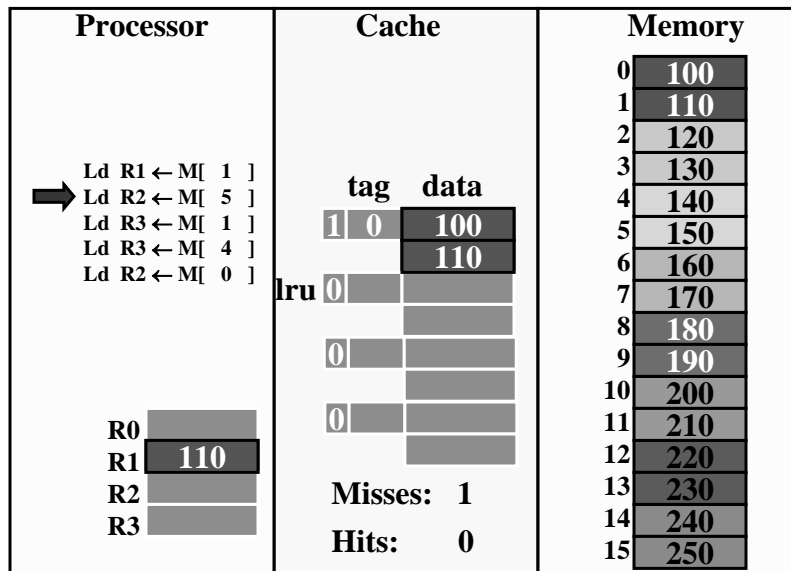


1st Access



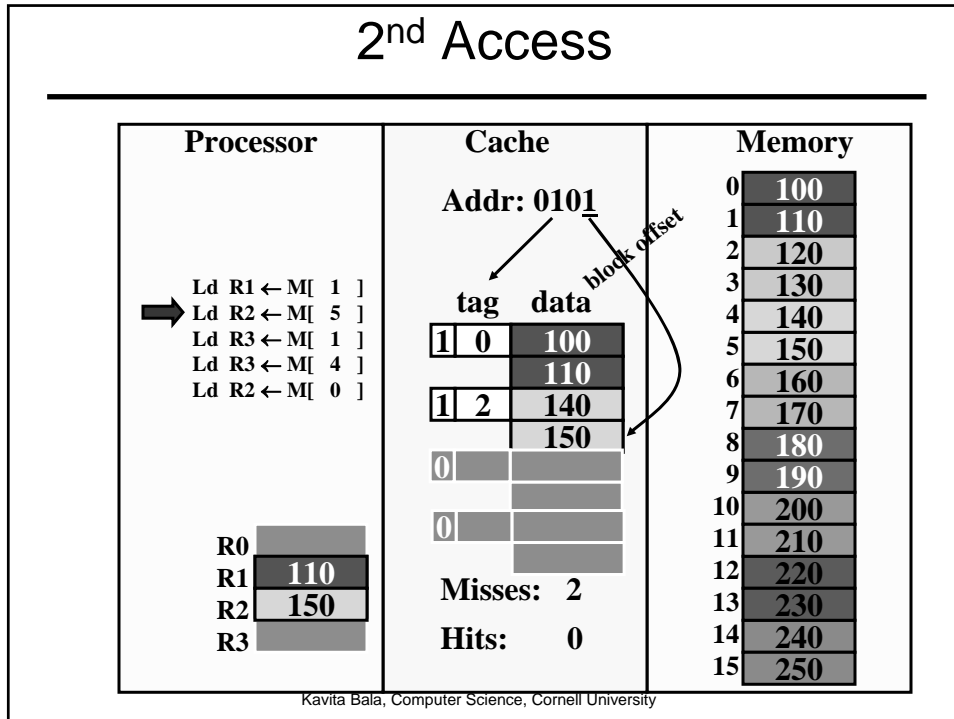
Kavita Bala, Computer Science, Cornell University

2nd Access

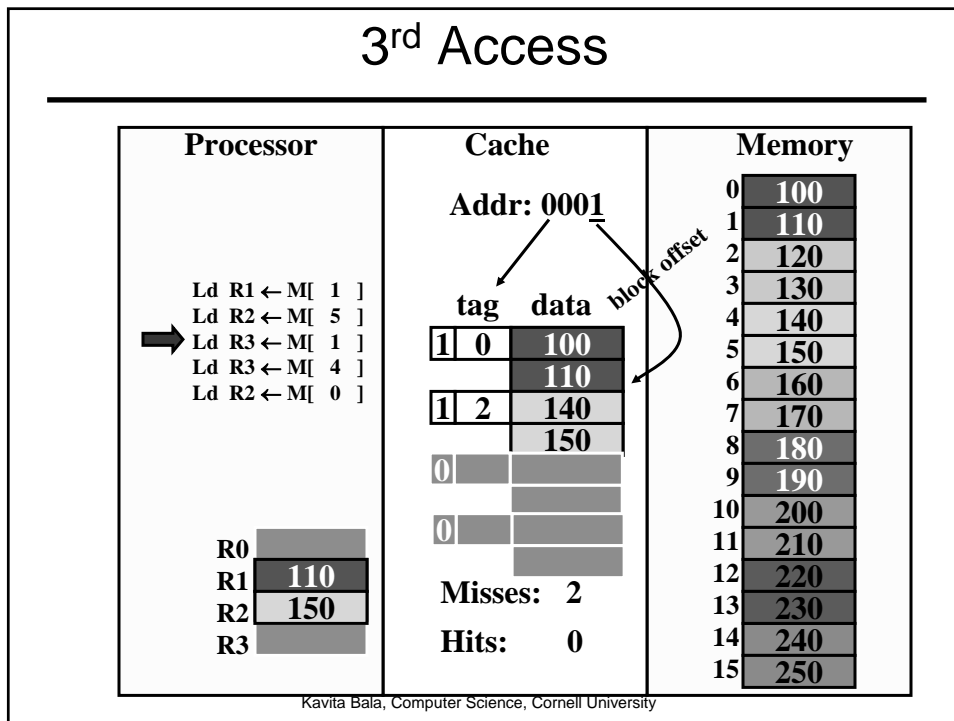


Kavita Bala, Computer Science, Cornell University

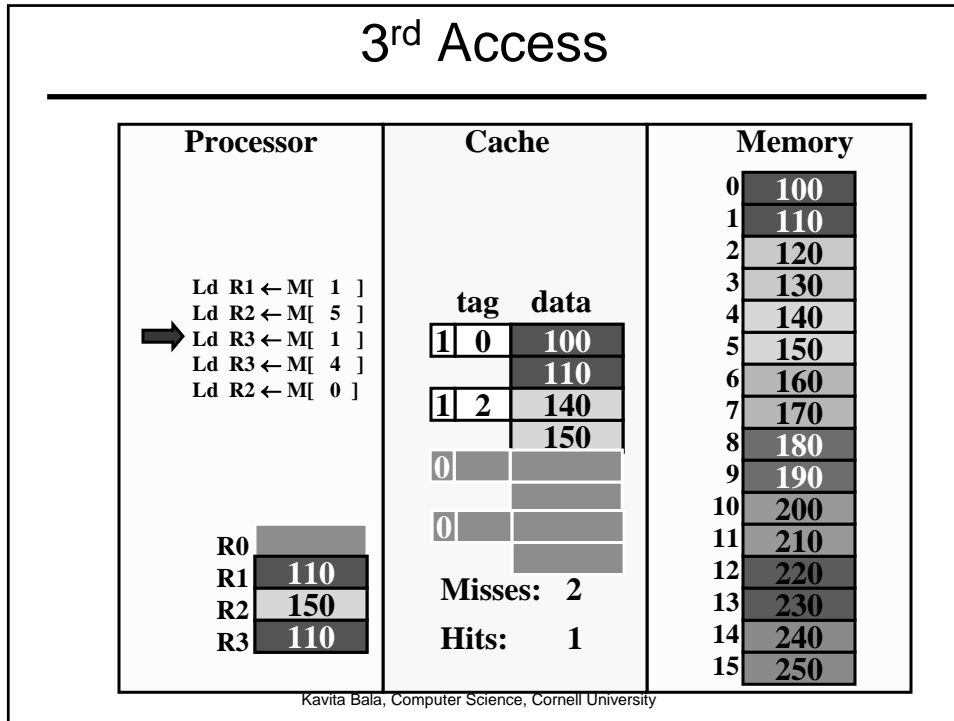
2nd Access



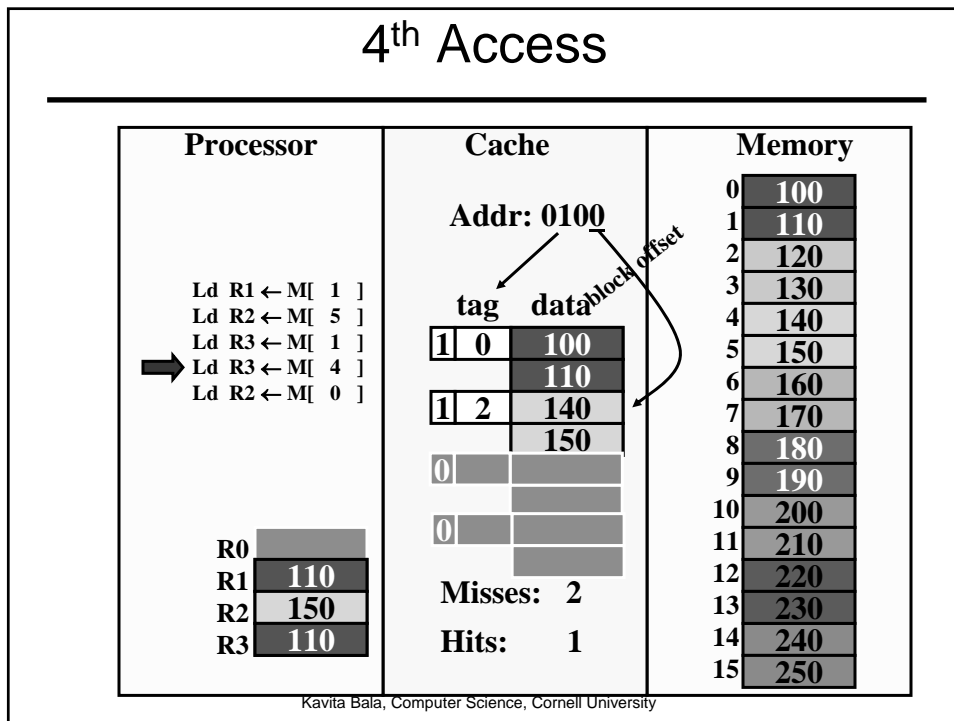
3rd Access



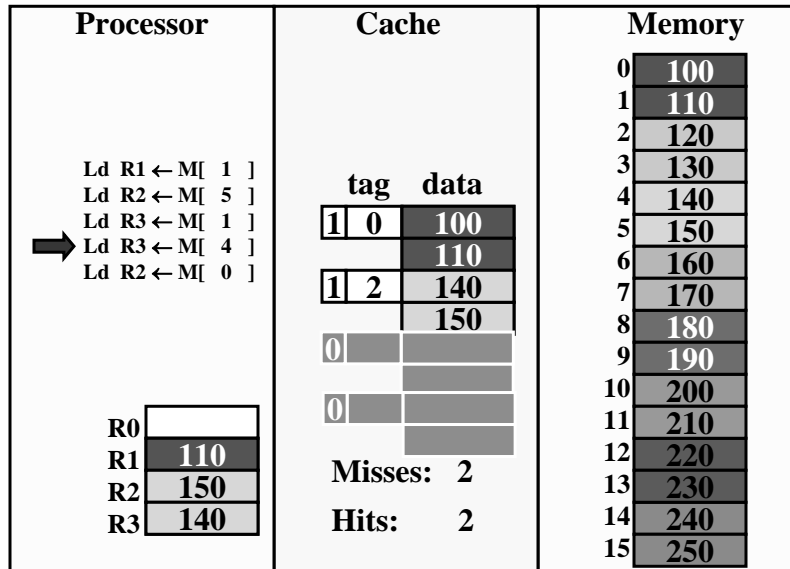
3rd Access



4th Access

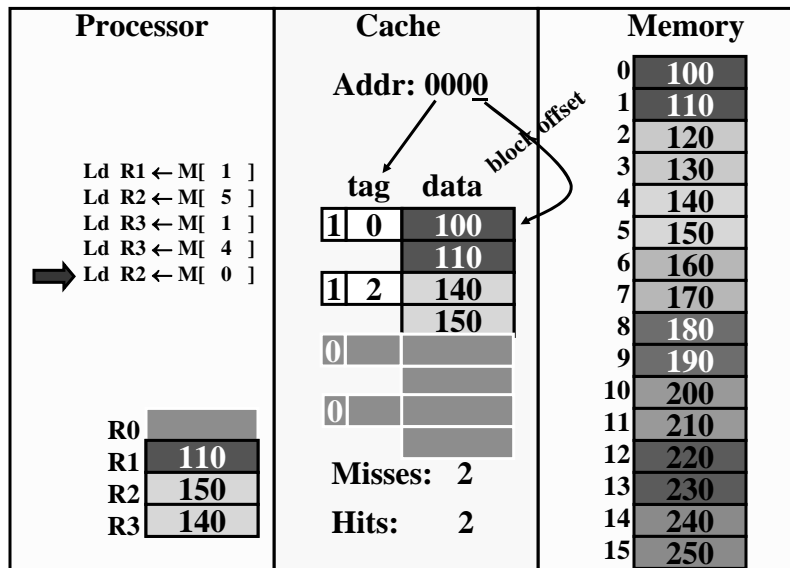


4th Access



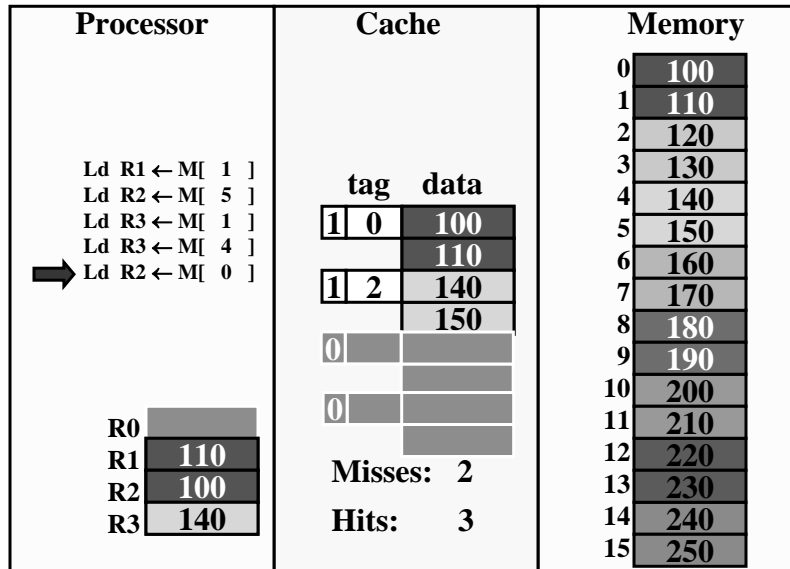
Kavita Bala, Computer Science, Cornell University

5th Access



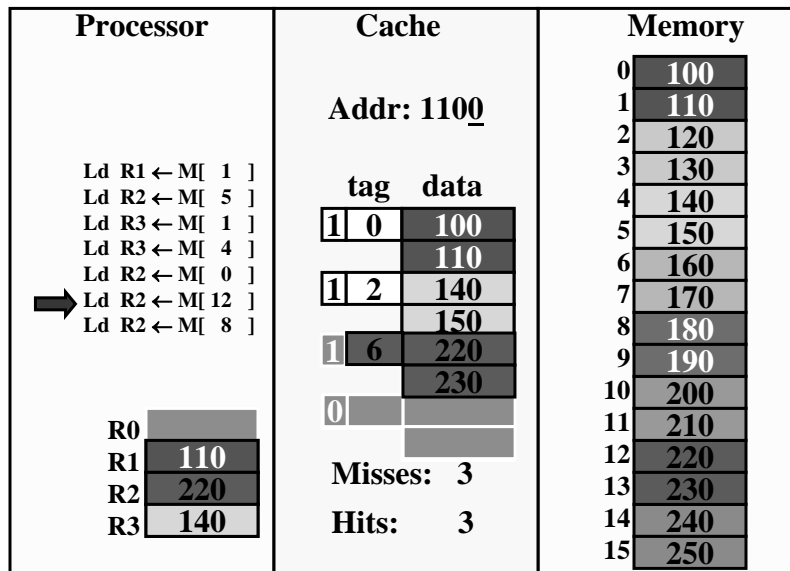
Kavita Bala, Computer Science, Cornell University

5th Access



Kavita Bala, Computer Science, Cornell University

6th Access



Kavita Bala, Computer Science, Cornell University

7th Access

Processor	Cache	Memory																																																								
<p>Ld R1 ← M[1] Ld R2 ← M[5] Ld R3 ← M[1] Ld R3 ← M[4] Ld R2 ← M[0] Ld R2 ← M[12] → Ld R2 ← M[8]</p> <p>R0 R1 110 R2 180 R3 140</p>	<p style="text-align: center;">Addr: 100<u>0</u></p> <p style="text-align: center;">tag data</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td style="width: 10%; text-align: center;">1</td><td style="width: 10%; text-align: center;">0</td><td style="width: 80%; text-align: center;">100</td></tr> <tr><td></td><td></td><td style="text-align: center;">110</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">2</td><td style="text-align: center;">140</td></tr> <tr><td></td><td></td><td style="text-align: center;">150</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">6</td><td style="text-align: center;">220</td></tr> <tr><td></td><td></td><td style="text-align: center;">230</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">4</td><td style="text-align: center;">180</td></tr> <tr><td></td><td></td><td style="text-align: center;">190</td></tr> </table> <p>Misses: 4 Hits: 3</p>	1	0	100			110	1	2	140			150	1	6	220			230	1	4	180			190	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td style="width: 5%; text-align: center;">0</td><td style="width: 95%; text-align: center;">100</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">110</td></tr> <tr><td style="text-align: center;">2</td><td style="text-align: center;">120</td></tr> <tr><td style="text-align: center;">3</td><td style="text-align: center;">130</td></tr> <tr><td style="text-align: center;">4</td><td style="text-align: center;">140</td></tr> <tr><td style="text-align: center;">5</td><td style="text-align: center;">150</td></tr> <tr><td style="text-align: center;">6</td><td style="text-align: center;">160</td></tr> <tr><td style="text-align: center;">7</td><td style="text-align: center;">170</td></tr> <tr><td style="text-align: center;">8</td><td style="text-align: center;">180</td></tr> <tr><td style="text-align: center;">9</td><td style="text-align: center;">190</td></tr> <tr><td style="text-align: center;">10</td><td style="text-align: center;">200</td></tr> <tr><td style="text-align: center;">11</td><td style="text-align: center;">210</td></tr> <tr><td style="text-align: center;">12</td><td style="text-align: center;">220</td></tr> <tr><td style="text-align: center;">13</td><td style="text-align: center;">230</td></tr> <tr><td style="text-align: center;">14</td><td style="text-align: center;">240</td></tr> <tr><td style="text-align: center;">15</td><td style="text-align: center;">250</td></tr> </table>	0	100	1	110	2	120	3	130	4	140	5	150	6	160	7	170	8	180	9	190	10	200	11	210	12	220	13	230	14	240	15	250
1	0	100																																																								
		110																																																								
1	2	140																																																								
		150																																																								
1	6	220																																																								
		230																																																								
1	4	180																																																								
		190																																																								
0	100																																																									
1	110																																																									
2	120																																																									
3	130																																																									
4	140																																																									
5	150																																																									
6	160																																																									
7	170																																																									
8	180																																																									
9	190																																																									
10	200																																																									
11	210																																																									
12	220																																																									
13	230																																																									
14	240																																																									
15	250																																																									

Kavita Bala, Computer Science, Cornell University

8th and 9th Access

Processor	Cache	Memory																																																								
<p>Ld R1 ← M[1] Ld R2 ← M[5] Ld R3 ← M[1] Ld R3 ← M[4] Ld R2 ← M[0] Ld R2 ← M[12] Ld R2 ← M[8] → Ld R2 ← M[4] Ld R2 ← M[0]</p>	<p style="text-align: center;">tag data</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td style="width: 10%; text-align: center;">1</td><td style="width: 10%; text-align: center;">0</td><td style="width: 80%; text-align: center;">100</td></tr> <tr><td></td><td></td><td style="text-align: center;">110</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">2</td><td style="text-align: center;">140</td></tr> <tr><td></td><td></td><td style="text-align: center;">150</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">6</td><td style="text-align: center;">220</td></tr> <tr><td></td><td></td><td style="text-align: center;">230</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">4</td><td style="text-align: center;">180</td></tr> <tr><td></td><td></td><td style="text-align: center;">190</td></tr> </table> <p>Misses: 4 Hits: 5</p>	1	0	100			110	1	2	140			150	1	6	220			230	1	4	180			190	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td style="width: 5%; text-align: center;">0</td><td style="width: 95%; text-align: center;">100</td></tr> <tr><td style="text-align: center;">1</td><td style="text-align: center;">110</td></tr> <tr><td style="text-align: center;">2</td><td style="text-align: center;">120</td></tr> <tr><td style="text-align: center;">3</td><td style="text-align: center;">130</td></tr> <tr><td style="text-align: center;">4</td><td style="text-align: center;">140</td></tr> <tr><td style="text-align: center;">5</td><td style="text-align: center;">150</td></tr> <tr><td style="text-align: center;">6</td><td style="text-align: center;">160</td></tr> <tr><td style="text-align: center;">7</td><td style="text-align: center;">170</td></tr> <tr><td style="text-align: center;">8</td><td style="text-align: center;">180</td></tr> <tr><td style="text-align: center;">9</td><td style="text-align: center;">190</td></tr> <tr><td style="text-align: center;">10</td><td style="text-align: center;">200</td></tr> <tr><td style="text-align: center;">11</td><td style="text-align: center;">210</td></tr> <tr><td style="text-align: center;">12</td><td style="text-align: center;">220</td></tr> <tr><td style="text-align: center;">13</td><td style="text-align: center;">230</td></tr> <tr><td style="text-align: center;">14</td><td style="text-align: center;">240</td></tr> <tr><td style="text-align: center;">15</td><td style="text-align: center;">250</td></tr> </table>	0	100	1	110	2	120	3	130	4	140	5	150	6	160	7	170	8	180	9	190	10	200	11	210	12	220	13	230	14	240	15	250
1	0	100																																																								
		110																																																								
1	2	140																																																								
		150																																																								
1	6	220																																																								
		230																																																								
1	4	180																																																								
		190																																																								
0	100																																																									
1	110																																																									
2	120																																																									
3	130																																																									
4	140																																																									
5	150																																																									
6	160																																																									
7	170																																																									
8	180																																																									
9	190																																																									
10	200																																																									
11	210																																																									
12	220																																																									
13	230																																																									
14	240																																																									
15	250																																																									

Kavita Bala, Computer Science, Cornell University

10th and 11th Access

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

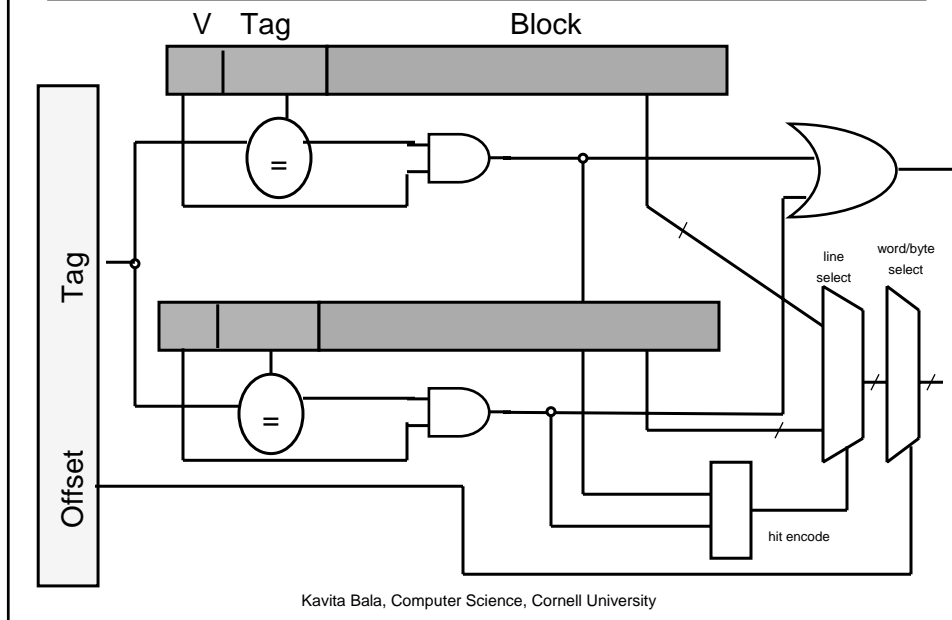
Kavita Bala, Computer Science, Cornell University

DM: 10th and 11th Access

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

Kavita Bala, Computer Science, Cornell University

Fully Associative Cache

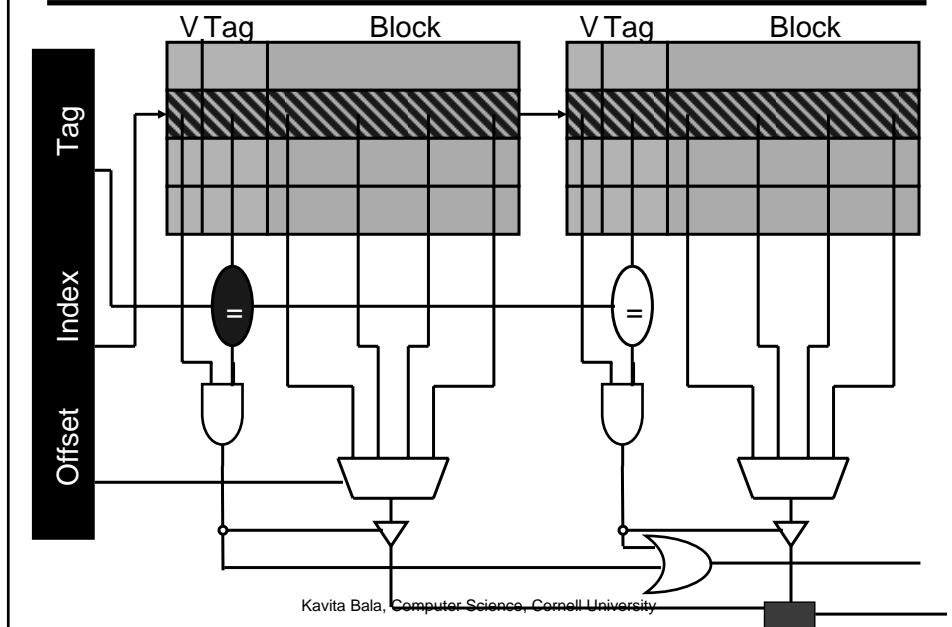


Cache Size

- Cache of size 2^n blocks
- Block size of 2^m word (block index: $m+2$ bits)
- Tag field: $32 - (m + 2)$
- Valid bit: 1
- Bits in cache: $2^n \times (\text{block size} + \text{tag size} + \text{valid bit size})$
 $= 2^n (2^m \times 32 + (32-m-2) + 1)$

Kavita Bala, Computer Science, Cornell University

2-Way Set-Associative Cache



2 Way Set Assoc: 10th and 11th Access

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

tag	data
0	
0	

Misses: 7
Hits: 4

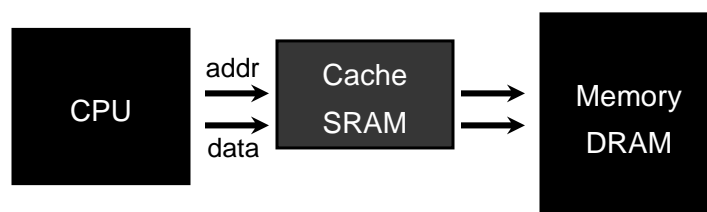
Kavita Bala, Computer Science, Cornell University

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 cache and to main memory
- Write-Back
 - write cache, write main memory only when block is evicted

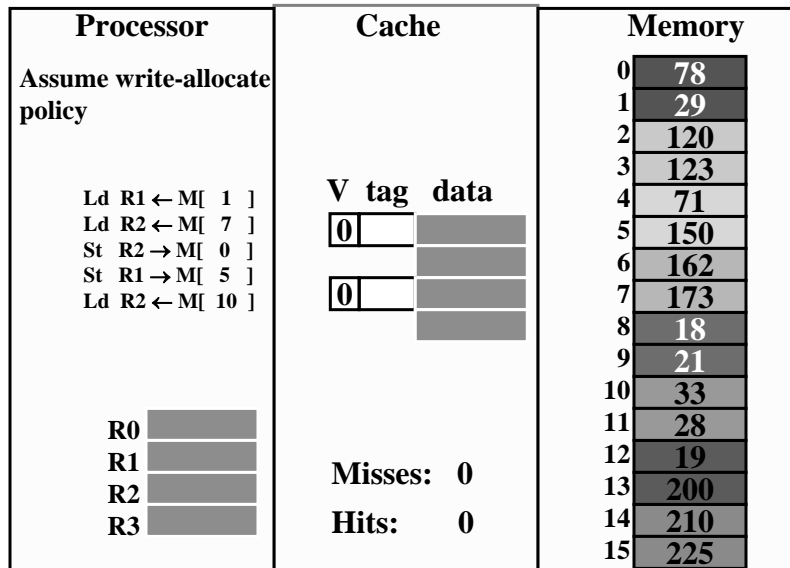
Kavita Bala, Computer Science, Cornell University

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)

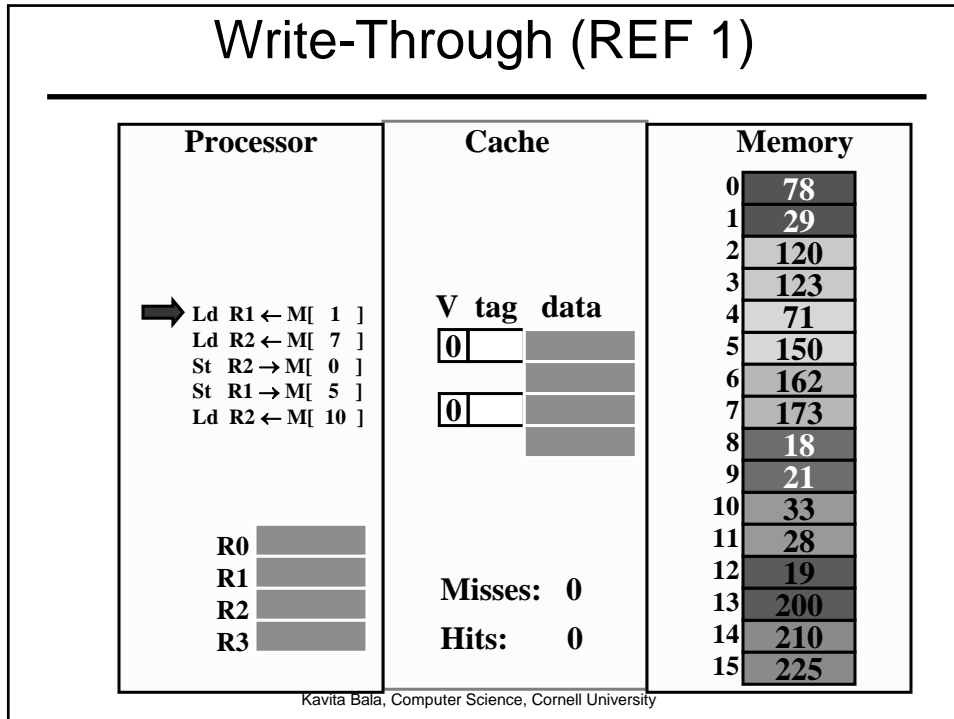
Kavita Bala, Computer Science, Cornell University

Handling Stores (Write-Through)

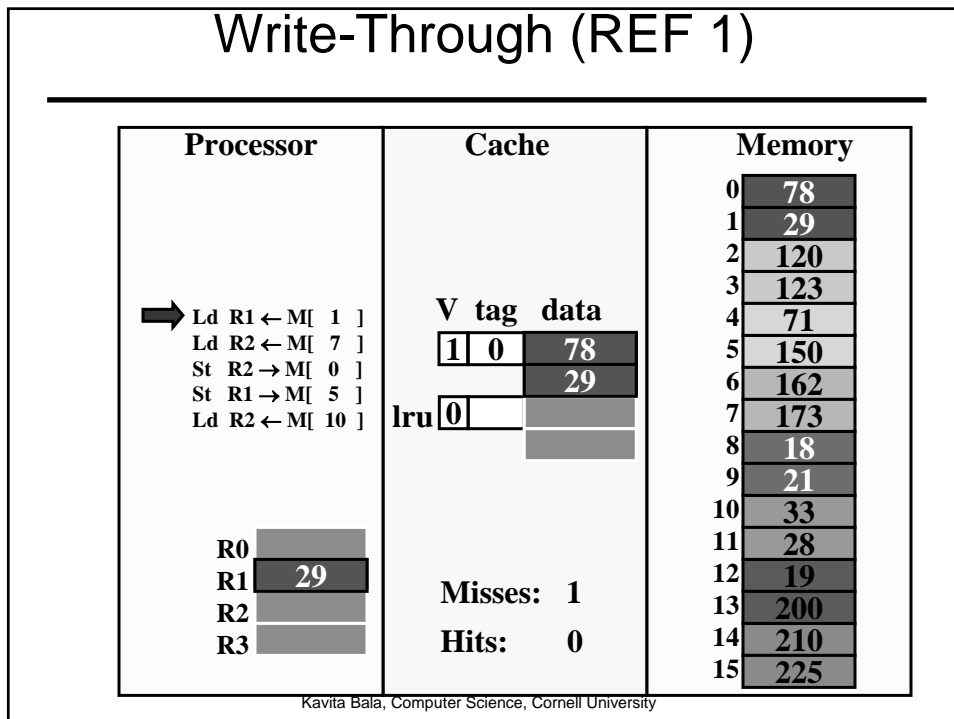


Kavita Bala, Computer Science, Cornell University

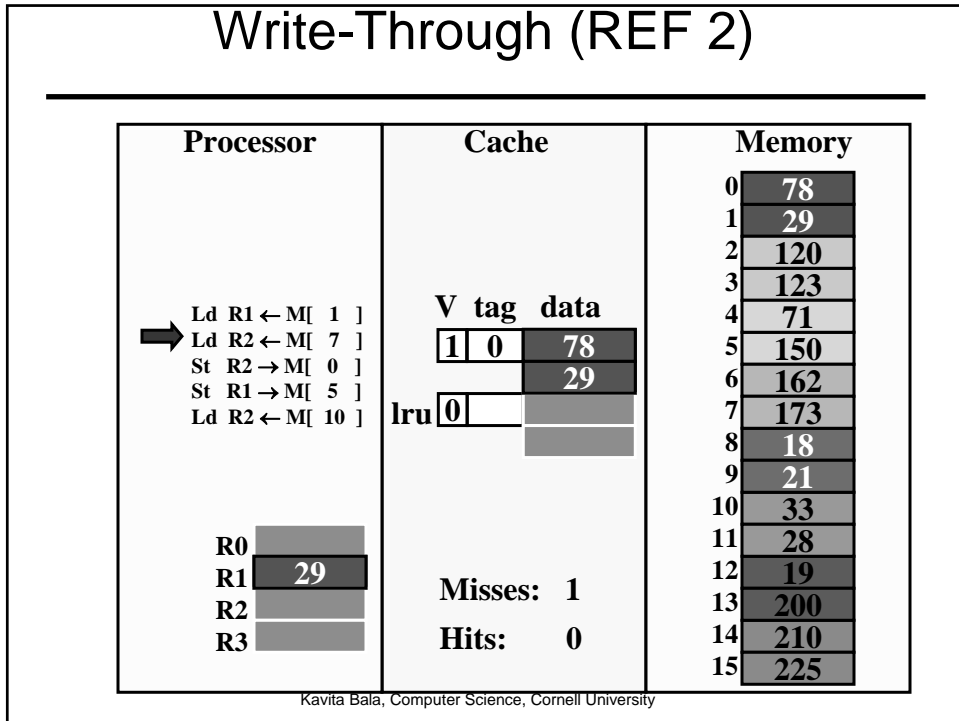
Write-Through (REF 1)



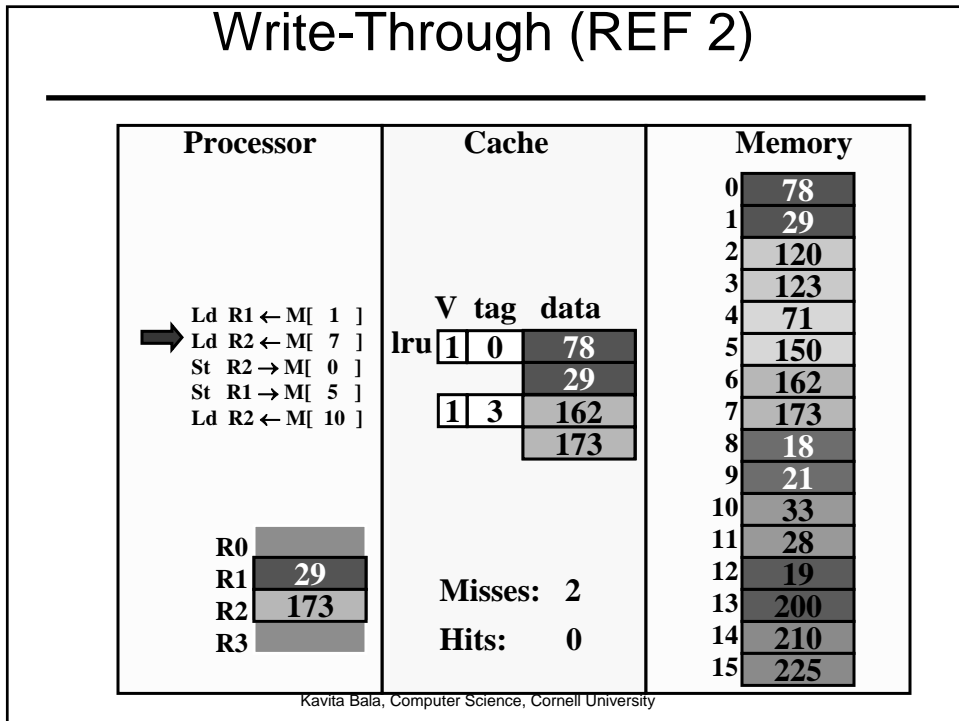
Write-Through (REF 1)



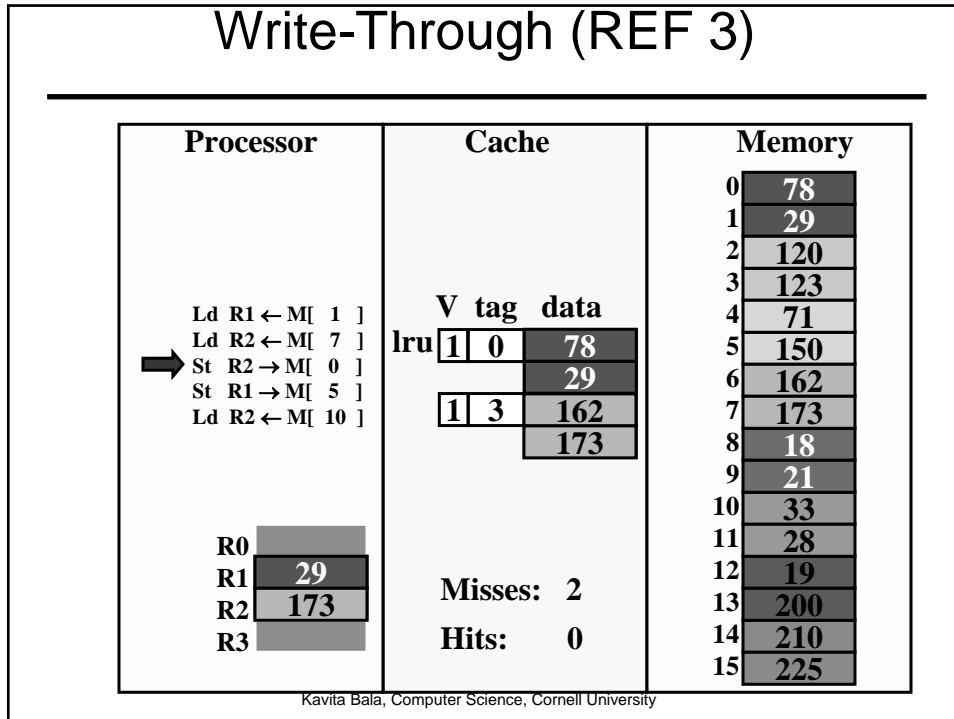
Write-Through (REF 2)



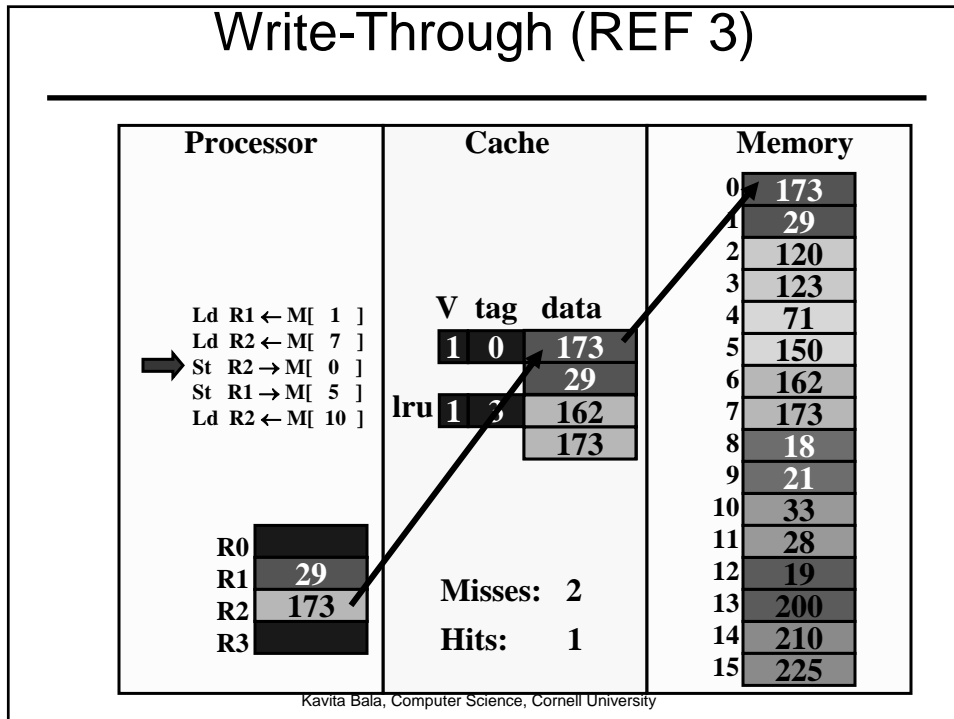
Write-Through (REF 2)



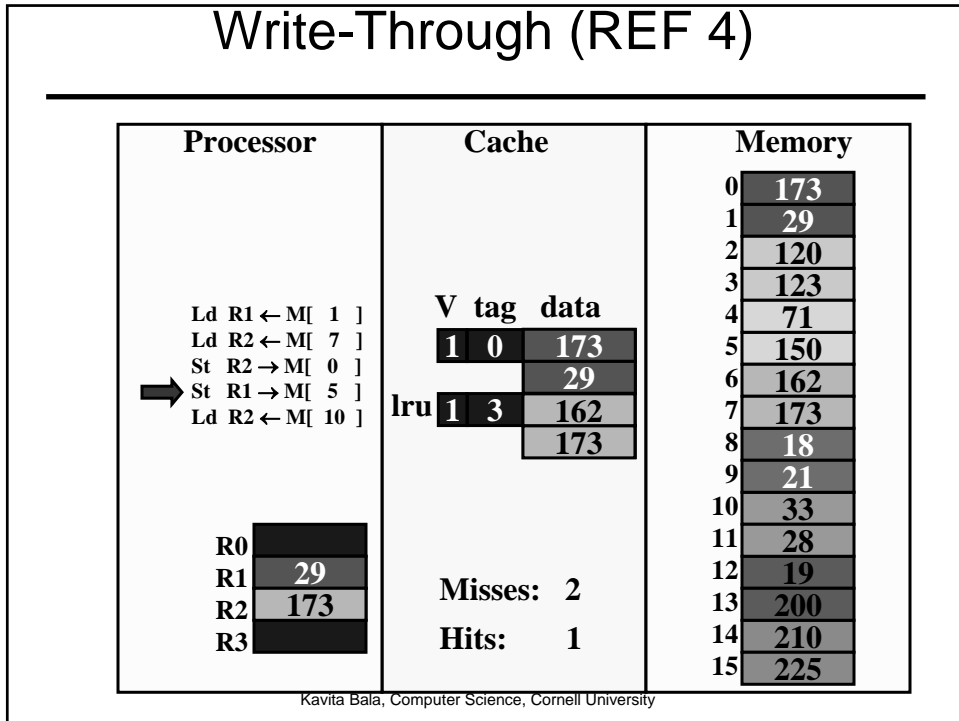
Write-Through (REF 3)



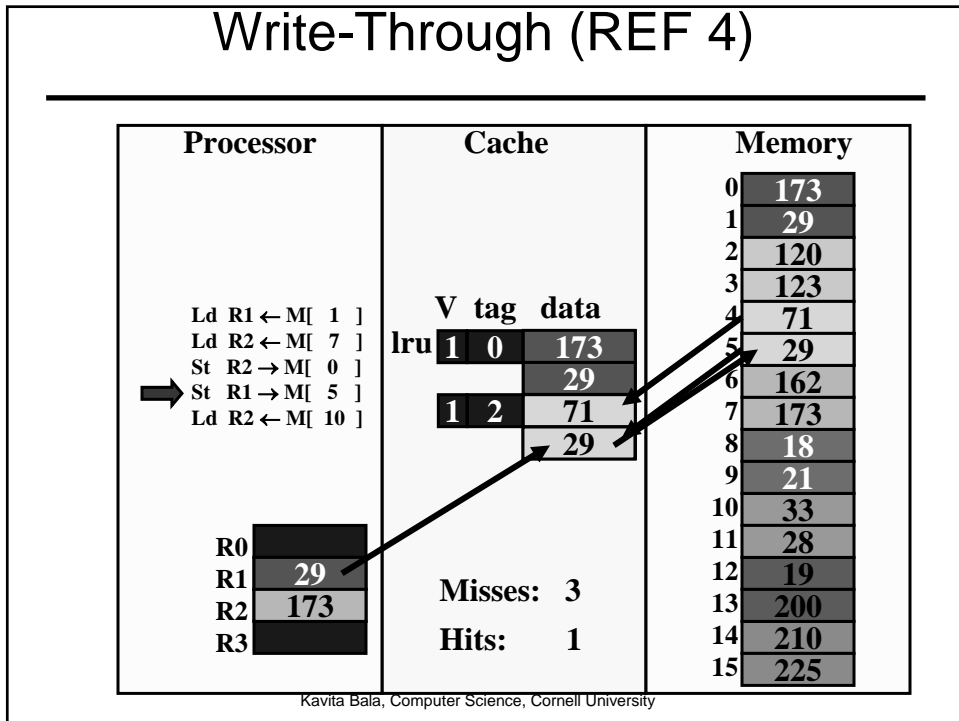
Write-Through (REF 3)



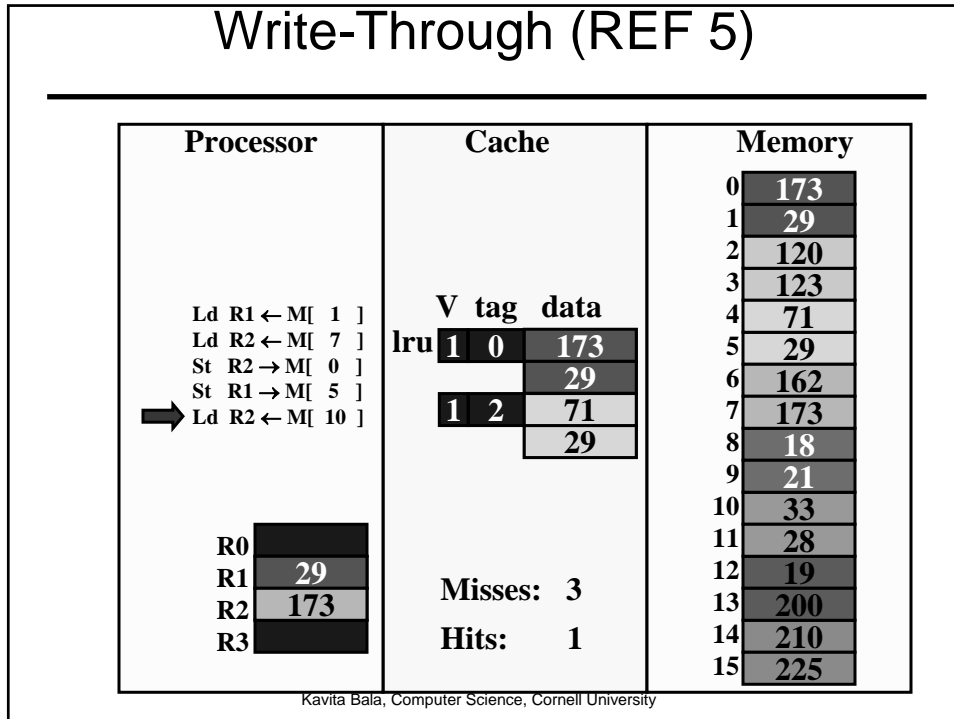
Write-Through (REF 4)



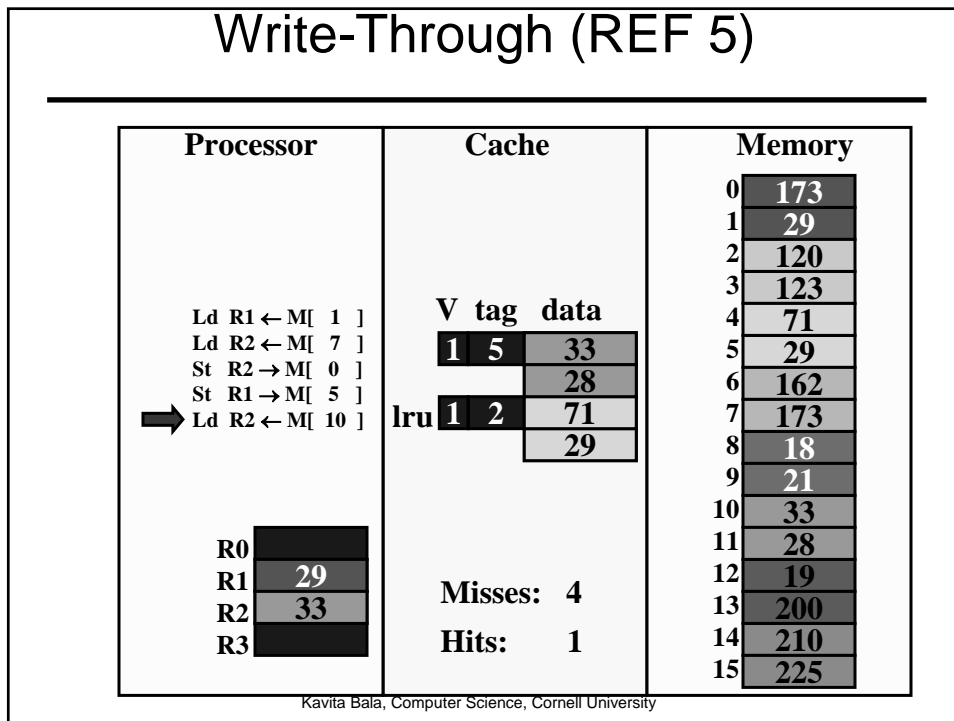
Write-Through (REF 4)



Write-Through (REF 5)



Write-Through (REF 5)



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!