

CS 316: Memory and Caches

Kavita Bala
Fall 2007
Computer Science
Cornell University

Announcements

- Prelim graded
 - Average: 69.8
 - Standard Deviation: 16
 - Handed out after class
 - Discussed in recitation
- HW 2 is out
- Caches and Memory: Chapter 7 (H & P)

Memory So Far

Big array of storage with more complex indexing than registers

- Addressing modes help us access memory
- $A[i]$; use base (i) + displacement (A)

- Use less space in instruction than 32-bit immediate field

Kavita Bala, Computer Science, Cornell University

SRAM vs. DRAM

- SRAM (**static random access memory**)
 - Faster than DRAM
 - Each storage cell is larger (4-6 transistors)
 - So smaller capacity for same area
 - 2-10ns access time
- DRAM (**dynamic random access memory**)
 - Each storage cell tiny (capacitance on wire)
 - Can get 1-2GB chips
 - 50-150ns access time
 - Leaky—needs to refresh data periodically

Kavita Bala, Computer Science, Cornell University

Dynamic RAM: DRAM

- Dynamic-RAM
 - Data values require constant refresh
 - Internal circuitry keeps capacitor charges

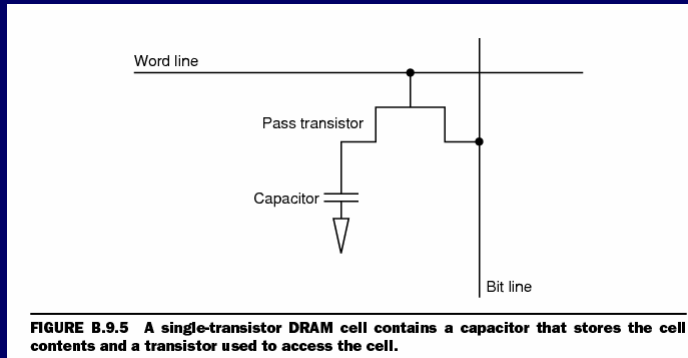


FIGURE B.9.5 A single-transistor DRAM cell contains a capacitor that stores the cell contents and a transistor used to access the cell.

Kavita Bala, Computer Science, Cornell University

Performance

- CPU clock rates $\sim 0.2\text{ns}$ - 2ns (5GHz-500MHz)

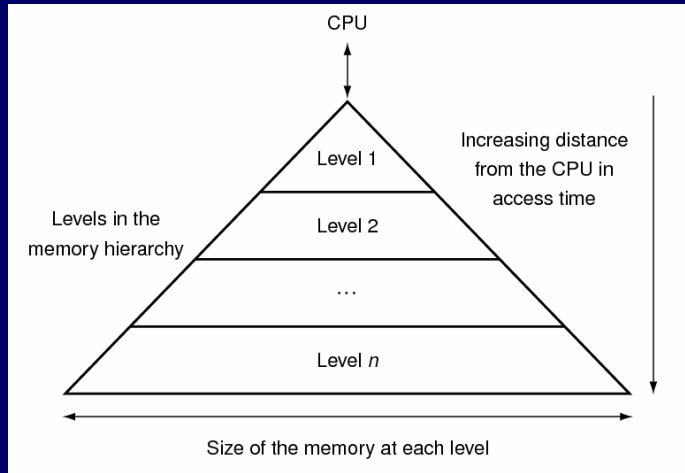
Technology	Capacity	Cost/GB	Latency
Tape	1 TB	\$.17	100s
Disk	300 GB	\$.34	Millions cycles (ms)
DRAM	4GB	\$100s	50-300 cycles (10s of ns)
SRAM off	512KB	\$4-10'sk	5-15 cycles (few ns)
SRAM on	16 KB	???	1-3 cycles (ns)

- Capacity and latency are closely coupled
- Cost is inversely proportional
- How do we create the illusion of large and fast memory?

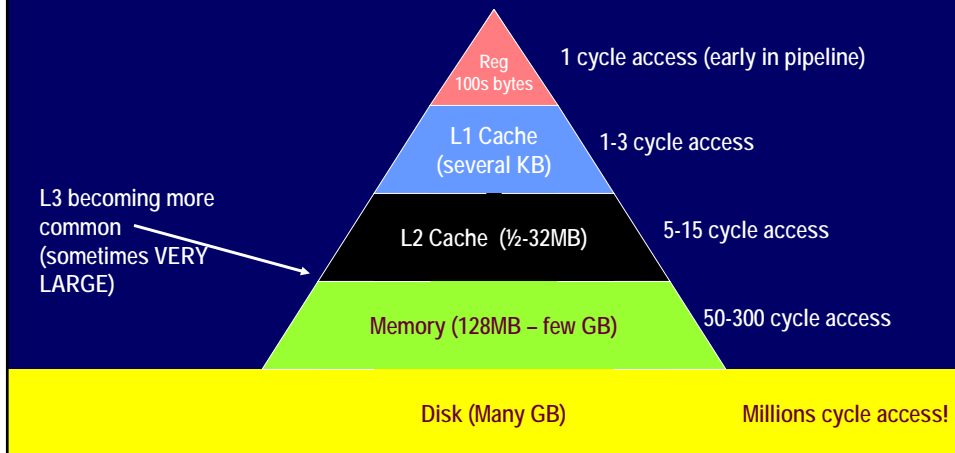
Kavita Bala, Computer Science, Cornell University

Memory Hierarchy

- Idea: Hide latency using small, fast memories called caches



Cache Design 101

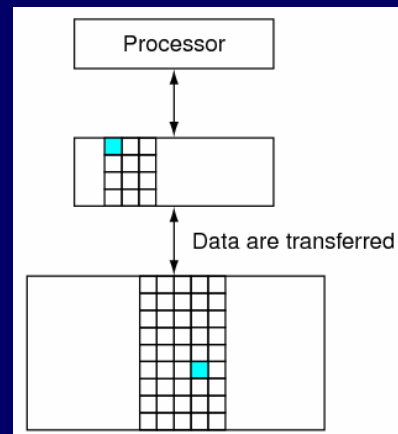


These are rough numbers: mileage may vary for latest/greatest Caches USUALLY made of SRAM

Kavita Bala, Computer Science, Cornell University

Memory Hierarchy

- Closer to processor
 - Subset of memory farther from processor
 - Faster and smaller
- Transfer an entire block



Kavita Bala, Computer Science, Cornell University

Insight of Caches

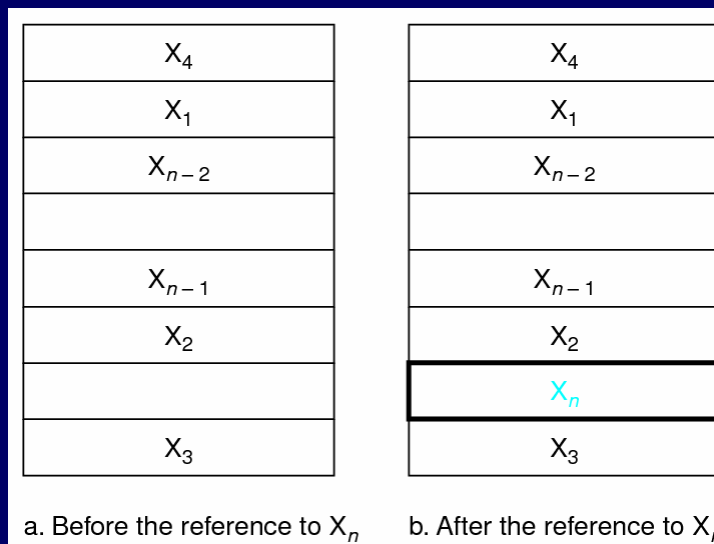
- Exploit locality
 - Two types: temporal and spatial
- Temporal locality
 - If memory location X is accessed, then it is more likely to be accessed again in the near future than some random location Y
 - Caches exploit temporal locality by placing a memory element that has been referenced into the cache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

Kavita Bala, Computer Science, Cornell University

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



Kavita Bala, Computer Science, Cornell University

Cache Organization

- Cache has to be fast and small
 - Gain speed by performing lookups in parallel, requires die real estate
 - Reduce hardware required by limiting where in the cache a block might be placed

Kavita Bala, Computer Science, Cornell University

Cache Organization

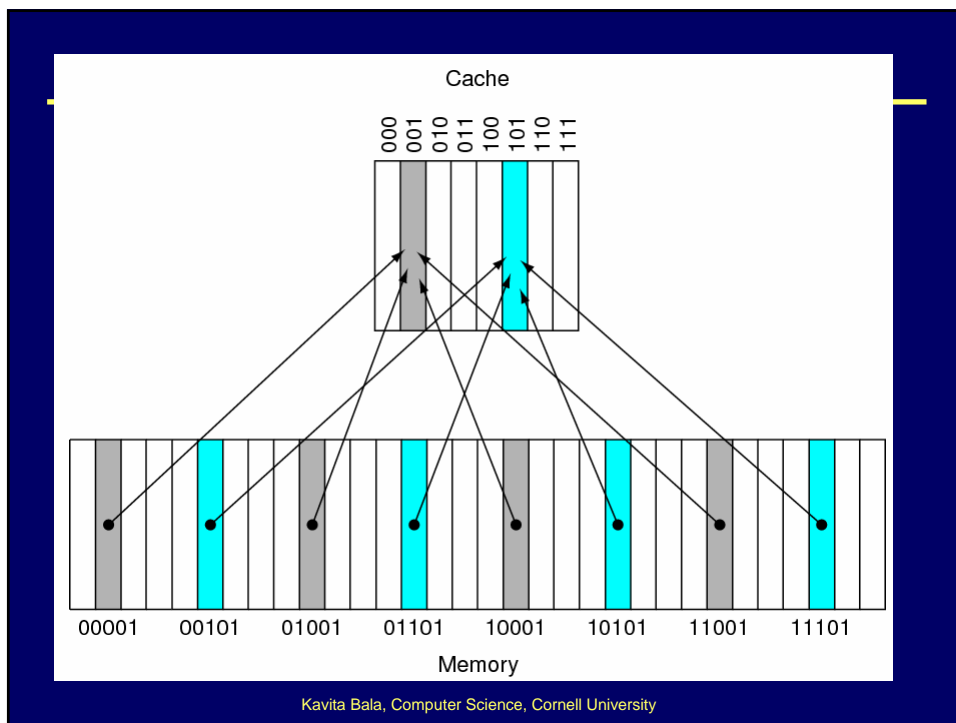
- Three common designs
 - Fully associative: Block can be anywhere in the cache
 - Direct mapped: Block can only be in one line in the cache
 - Set-associative: Block can be in a few (2 to 8) places in the cache

Kavita Bala, Computer Science, Cornell University

Direct Mapped Cache

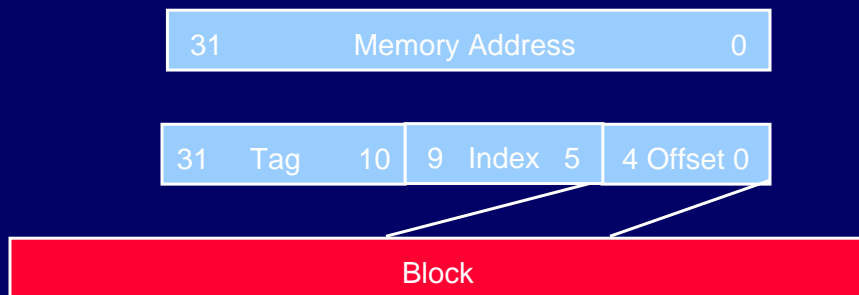
- Simplest
- Block can only be in one line in the cache
- How to determine this location?
 - Use modulo arithmetic
 - (Block address) modulo (# cache blocks)
 - For power of 2, use log (cache size in blocks)

Kavita Bala, Computer Science, Cornell University



Tags and Offsets

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

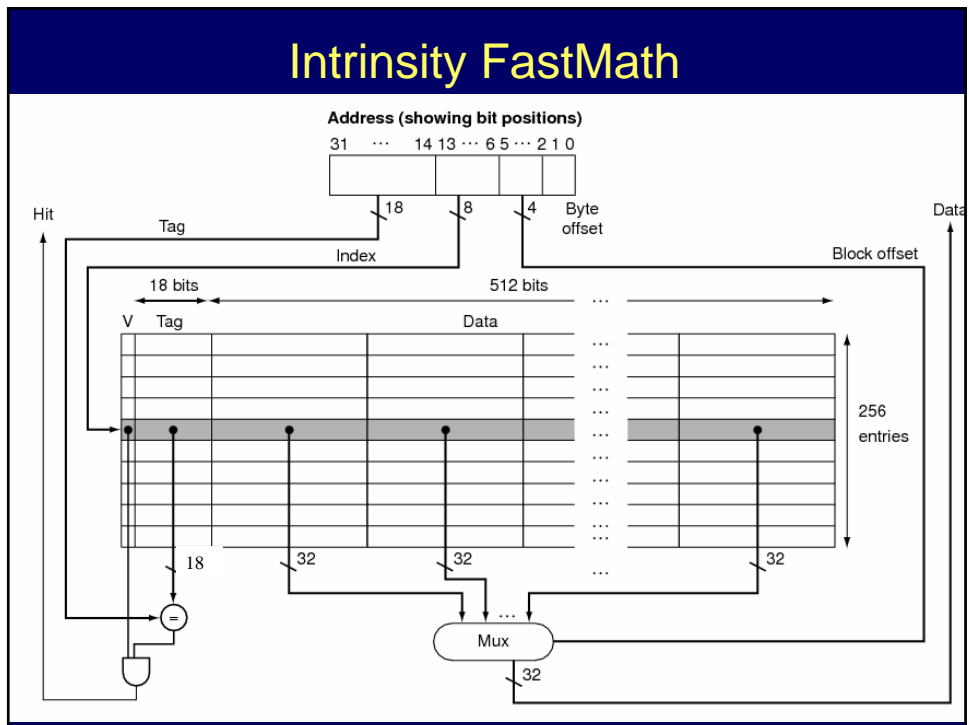
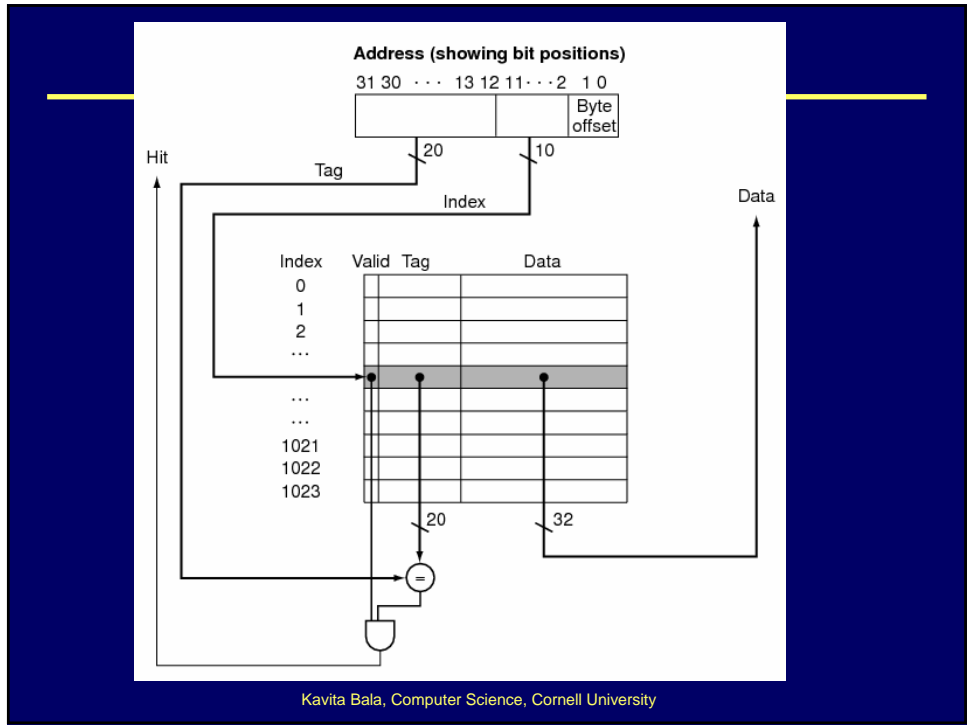


Kavita Bala, Computer Science, Cornell University

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

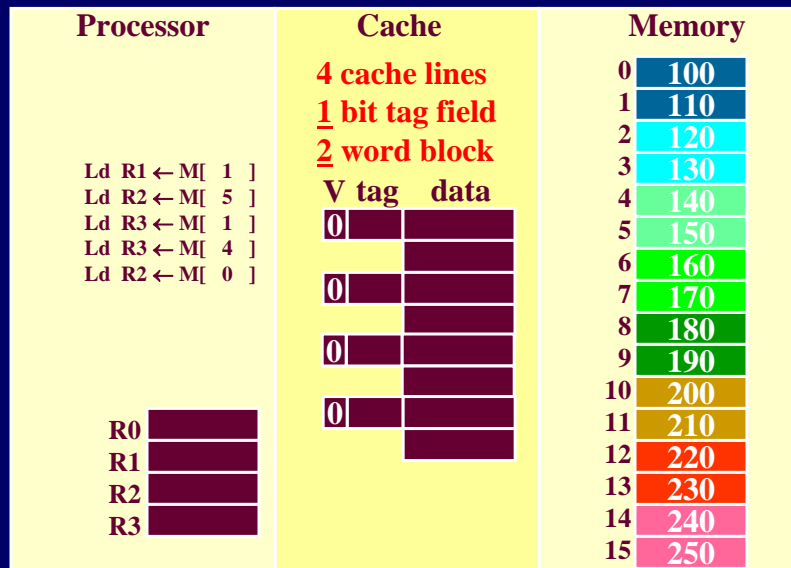


Cache Size

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

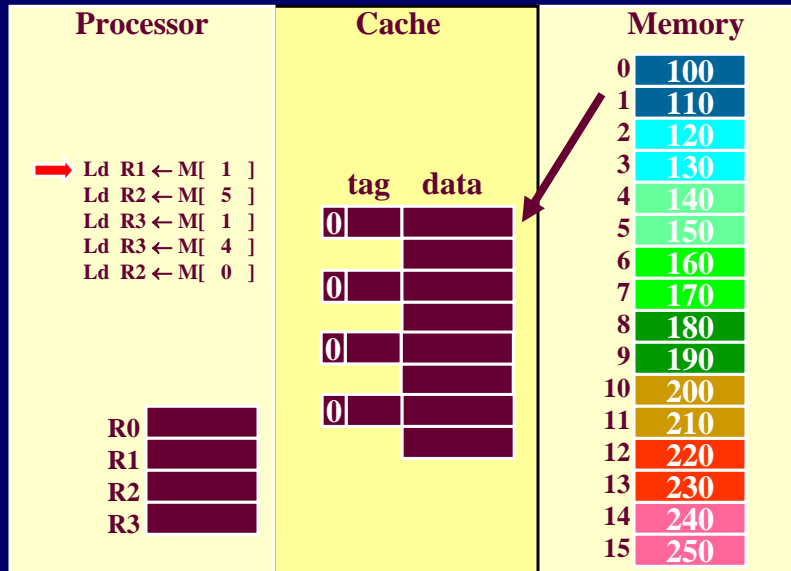
Kavita Bala, Computer Science, Cornell University

A Simple Direct Mapped Cache



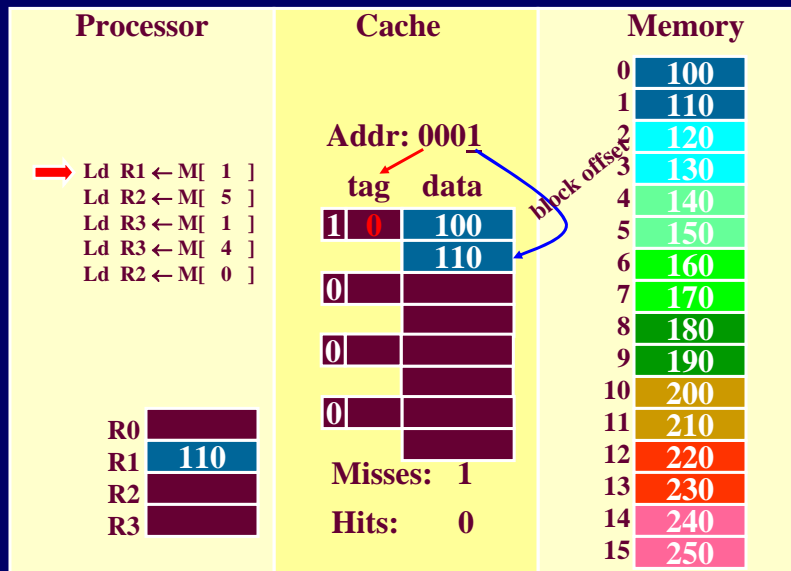
Kavita Bala, Computer Science, Cornell University

1st Access



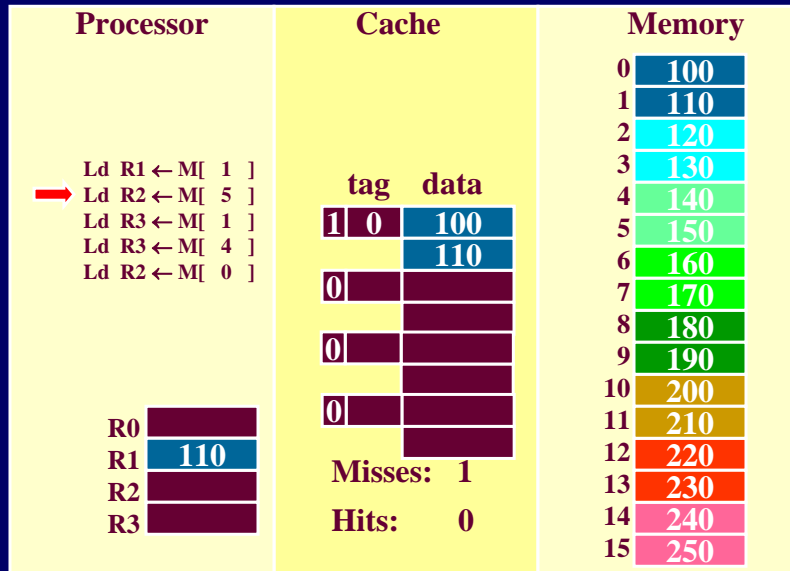
Kavita Bala, Computer Science, Cornell University

1st Access



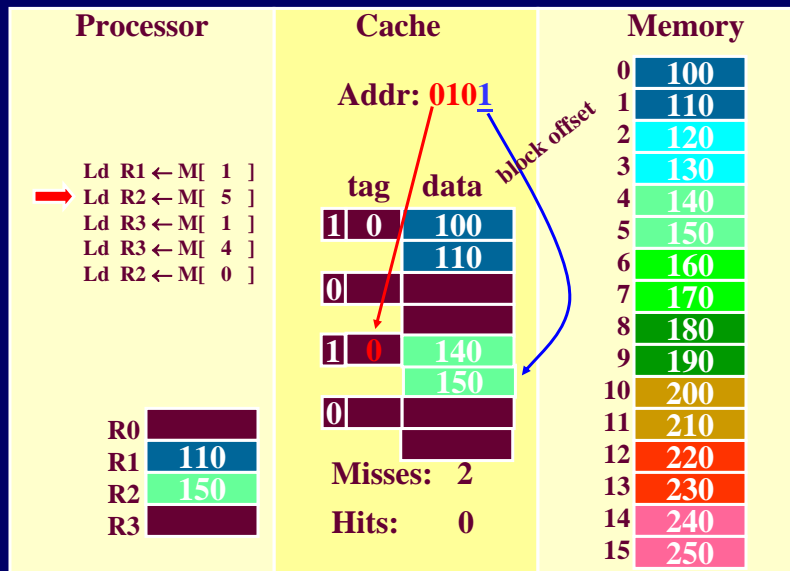
Kavita Bala, Computer Science, Cornell University

2nd Access



Kavita Bala, Computer Science, Cornell University

2nd Access



Kavita Bala, Computer Science, Cornell University

3rd Access

Processor	Cache	Memory
	Addr: 000<u>1</u>	
Ld R1 ← M[1]	tag data	0 100
Ld R2 ← M[5]	1 0 100	1 110
→ Ld R3 ← M[1]	0	2 120
Ld R3 ← M[4]	1 0 140	3 130
Ld R2 ← M[0]	0	4 140
	0	5 150
		6 160
R0		7 170
R1 110		8 180
R2 150		9 190
R3		10 200
	Misses: 2	11 210
	Hits: 0	12 220
		13 230
		14 240
		15 250

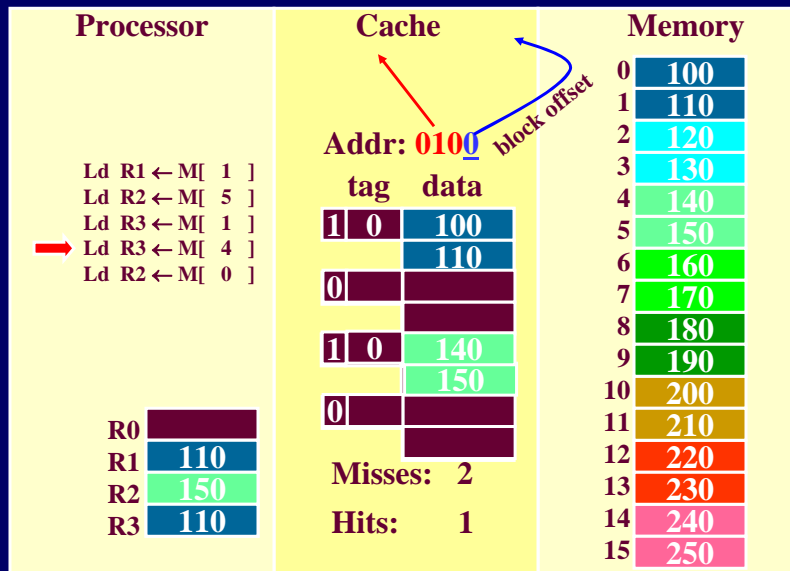
Kavita Bala, Computer Science, Cornell University

3rd Access

Processor	Cache	Memory
	tag data	
Ld R1 ← M[1]	1 0 100	0 100
Ld R2 ← M[5]	0	1 110
→ Ld R3 ← M[1]	1 0 140	2 120
Ld R3 ← M[4]	0	3 130
Ld R2 ← M[0]	0	4 140
		5 150
R0		6 160
R1 110		7 170
R2 150		8 180
R3 110		9 190
	Misses: 2	10 200
	Hits: 1	11 210
		12 220
		13 230
		14 240
		15 250

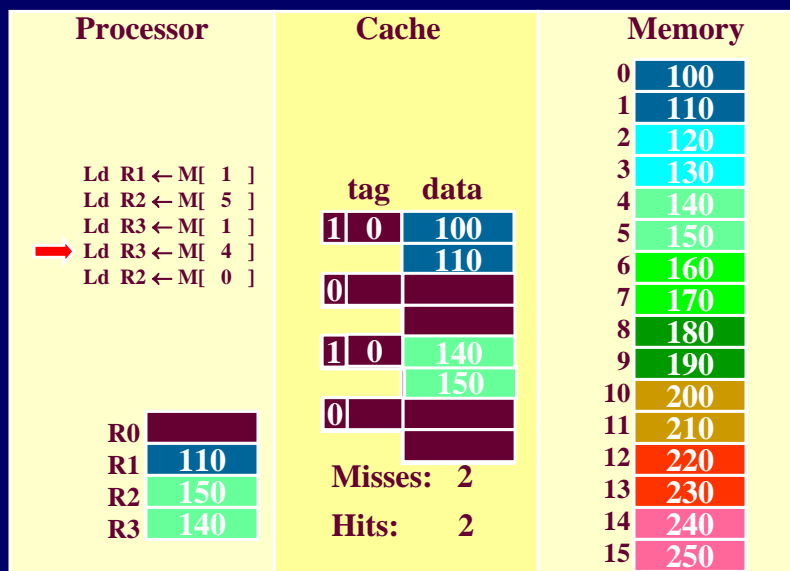
Kavita Bala, Computer Science, Cornell University

4th Access



Kavita Bala, Computer Science, Cornell University

4th Access



Kavita Bala, Computer Science, Cornell University

5th Access

Processor	Cache	Memory												
Ld R1 ← M[1]	Addr: 000<u>0</u> <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td>0</td> <td>110</td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>0</td> <td></td> </tr> </tbody> </table> Misses: 2 Hits: 2	tag	data	1 0	100	0	110	1 0	140		150	0		0 100
tag		data												
1 0		100												
0		110												
1 0		140												
	150													
0														
Ld R2 ← M[5]		1 110												
Ld R3 ← M[1]		2 120												
Ld R3 ← M[4]		3 130												
→ Ld R2 ← M[0]		4 140												
		5 150												
		6 160												
		7 170												
		8 180												
		9 190												
R0		10 200												
R1 110		11 210												
R2 150		12 220												
R3 140		13 230												
		14 240												
		15 250												

Kavita Bala, Computer Science, Cornell University

5th Access

Processor	Cache	Memory												
Ld R1 ← M[1]	Addr: 000<u>0</u> <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td>0</td> <td>110</td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>0</td> <td></td> </tr> </tbody> </table> Misses: 2 Hits: 3	tag	data	1 0	100	0	110	1 0	140		150	0		0 100
tag		data												
1 0		100												
0		110												
1 0		140												
	150													
0														
Ld R2 ← M[5]		1 110												
Ld R3 ← M[1]		2 120												
Ld R3 ← M[4]		3 130												
→ Ld R2 ← M[0]		4 140												
		5 150												
		6 160												
		7 170												
		8 180												
		9 190												
R0		10 200												
R1 110		11 210												
R2 100		12 220												
R3 140		13 230												
		14 240												
		15 250												

Kavita Bala, Computer Science, Cornell University

5th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1010 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td></td> <td>110</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> </tbody> </table> Misses: 2 Hits: 3	tag	data	1 0	100		110	0				1 0	140		150	0				0 100
tag		data																		
1 0		100																		
		110																		
0																				
1 0	140																			
	150																			
0																				
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[10]		5 150																		
Ld R2 ← M[15]		6 160																		
		7 170																		
		8 180																		
R0		9 190																		
R1 110		10 200																		
R2 100		11 210																		
R3 140		12 220																		
		13 230																		
		14 240																		
		15 250																		

Kavita Bala, Computer Science, Cornell University

5th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1010 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td></td> <td>110</td> </tr> <tr> <td>1 1</td> <td>200</td> </tr> <tr> <td></td> <td>210</td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> </tbody> </table> Misses: 3 Hits: 3	tag	data	1 0	100		110	1 1	200		210	1 0	140		150	0				0 100
tag		data																		
1 0		100																		
		110																		
1 1		200																		
		210																		
1 0	140																			
	150																			
0																				
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[10]		5 150																		
Ld R2 ← M[15]		6 160																		
		7 170																		
		8 180																		
R0		9 190																		
R1 110		10 200																		
R2 200		11 210																		
R3 140		12 220																		
		13 230																		
		14 240																		
		15 250																		

Kavita Bala, Computer Science, Cornell University

6th Access

Processor	Cache	Memory																
Ld R1 ← M[1]	Addr: 1111 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td></td> <td>110</td> </tr> <tr> <td>1 1</td> <td>200</td> </tr> <tr> <td></td> <td>210</td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>0</td> <td></td> </tr> </tbody> </table>	tag	data	1 0	100		110	1 1	200		210	1 0	140		150	0		0 100
tag		data																
1 0		100																
		110																
1 1		200																
		210																
1 0	140																	
	150																	
0																		
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[10]		5 150																
→ Ld R2 ← M[15]		6 160																
		7 170																
		8 180																
		9 190																
R0		10 200																
R1 110		11 210																
R2 200	Misses: 3	12 220																
R3 140	Hits: 3	13 230																
		14 240																
		15 250																

Kavita Bala, Computer Science, Cornell University

6th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1111 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td></td> <td>110</td> </tr> <tr> <td>1 1</td> <td>200</td> </tr> <tr> <td></td> <td>210</td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>1 1</td> <td>240</td> </tr> <tr> <td></td> <td>250</td> </tr> </tbody> </table>	tag	data	1 0	100		110	1 1	200		210	1 0	140		150	1 1	240		250	0 100
tag		data																		
1 0		100																		
		110																		
1 1		200																		
		210																		
1 0	140																			
	150																			
1 1	240																			
	250																			
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[10]		5 150																		
→ Ld R2 ← M[15]		6 160																		
		7 170																		
		8 180																		
		9 190																		
R0		10 200																		
R1 110		11 210																		
R2 250	Misses: 4	12 220																		
R3 140	Hits: 3	13 230																		
		14 240																		
		15 250																		

Kavita Bala, Computer Science, Cornell University

Cache

- Exploiting Spatial and Temporal Locality

Kavita Bala, Computer Science, Cornell University

7th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1000 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr><td>1 0</td><td>100</td></tr> <tr><td></td><td>110</td></tr> <tr><td>1 1</td><td>200</td></tr> <tr><td></td><td>210</td></tr> <tr><td>1 0</td><td>140</td></tr> <tr><td></td><td>150</td></tr> <tr><td>1 1</td><td>240</td></tr> <tr><td></td><td>250</td></tr> </tbody> </table>	tag	data	1 0	100		110	1 1	200		210	1 0	140		150	1 1	240		250	0 100
tag		data																		
1 0		100																		
		110																		
1 1		200																		
		210																		
1 0		140																		
		150																		
1 1		240																		
		250																		
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[10]		5 150																		
Ld R2 ← M[15]		6 160																		
→ Ld R2 ← M[8]		7 170																		
		8 180																		
		9 190																		
R0		10 200																		
R1 110		11 210																		
R2 250	Misses: 4	12 220																		
R3 140	Hits: 3	13 230																		
		14 240																		
		15 250																		

Kavita Bala, Computer Science, Cornell University

7th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1000 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr><td>1</td><td>180</td></tr> <tr><td></td><td>190</td></tr> <tr><td>1</td><td>200</td></tr> <tr><td></td><td>210</td></tr> <tr><td>1</td><td>140</td></tr> <tr><td></td><td>150</td></tr> <tr><td>1</td><td>240</td></tr> <tr><td></td><td>250</td></tr> </tbody> </table>	tag	data	1	180		190	1	200		210	1	140		150	1	240		250	0 100
tag		data																		
1		180																		
		190																		
1		200																		
		210																		
1		140																		
		150																		
1		240																		
		250																		
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[10]		5 150																		
Ld R2 ← M[15]		6 160																		
→ Ld R2 ← M[8]		7 170																		
		8 180																		
		9 190																		
R0		10 200																		
R1 110		11 210																		
R2 180	Misses: 5	12 220																		
R3 140	Hits: 3	13 230																		
		14 240																		
		15 250																		

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

6th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1100 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td></td> <td>110</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> <tr> <td>1 0</td> <td>140</td> </tr> <tr> <td></td> <td>150</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> </tbody> </table> Misses: 2 Hits: 3	tag	data	1 0	100		110	0				1 0	140		150	0				0 100
tag		data																		
1 0		100																		
		110																		
0																				
1 0	140																			
	150																			
0																				
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																		
		7 170																		
		8 180																		
		9 190																		
R0		10 200																		
R1 110		11 210																		
R2 100		12 220																		
R3 140		13 230																		
		14 240																		
		15 250																		

Kavita Bala, Computer Science, Cornell University

6th Access

Processor	Cache	Memory																		
Ld R1 ← M[1]	Addr: 1100 <table border="1"> <thead> <tr> <th>tag</th> <th>data</th> </tr> </thead> <tbody> <tr> <td>1 0</td> <td>100</td> </tr> <tr> <td></td> <td>110</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> <tr> <td>1 1</td> <td>220</td> </tr> <tr> <td></td> <td>230</td> </tr> <tr> <td>0</td> <td></td> </tr> <tr> <td></td> <td></td> </tr> </tbody> </table> Misses: 3 Hits: 3	tag	data	1 0	100		110	0				1 1	220		230	0				0 100
tag		data																		
1 0		100																		
		110																		
0																				
1 1	220																			
	230																			
0																				
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																		
		7 170																		
		8 180																		
		9 190																		
R0		10 200																		
R1 110		11 210																		
R2 220		12 220																		
R3 140		13 230																		
		14 240																		
		15 250																		

Kavita Bala, Computer Science, Cornell University

7th Access

Processor	Cache	Memory
Ld R1 ← M[1]	Addr: 1000	0 100
Ld R2 ← M[5]		1 110
Ld R3 ← M[1]	tag data	2 120
Ld R3 ← M[4]	1 0 100	3 130
Ld R2 ← M[0]	0	4 140
Ld R2 ← M[12]	1 1 220	5 150
→ Ld R2 ← M[8]	0	6 160
R0	1 1 230	7 170
R1 110	0	8 180
R2 220	Misses: 3	9 190
R3 140	Hits: 3	10 200
		11 210
		12 220
		13 230
		14 240
		15 250

Kavita Bala, Computer Science, Cornell University

7th Access

Processor	Cache	Memory
Ld R1 ← M[1]	Addr: 1000	0 100
Ld R2 ← M[5]		1 110
Ld R3 ← M[1]	tag data	2 120
Ld R3 ← M[4]	1 1 180	3 130
Ld R2 ← M[0]	0	4 140
Ld R2 ← M[12]	1 1 190	5 150
→ Ld R2 ← M[8]	0	6 160
R0	1 1 220	7 170
R1 180	0	8 180
R2 220	Misses: 4	9 190
R3 140	Hits: 3	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
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]	tag data	7 170
Ld R2 ← M[0]	1 0 100	8 180
	0	9 190
	1 0 140	10 200
	0	11 210
		12 220
		13 230
		14 240
		15 250

Misses: 4+2
Hits: 3

Kavita Bala, Computer Science, Cornell University

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]	tag data	9 190
Ld R2 ← M[8]	1 1 180	10 200
	0	11 210
	1 1 220	12 220
	0	13 230
		14 240
		15 250

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