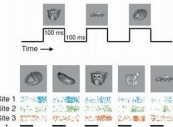


**Lecture #21**  
**Caches I**



**CPS Today!** 2005-11-14  
 There is **one** handout today at the front and back of the room!  
 Lecturer PSOE, new dad Dan Garcia  
[www.cs.berkeley.edu/~ddgarcia](http://www.cs.berkeley.edu/~ddgarcia)

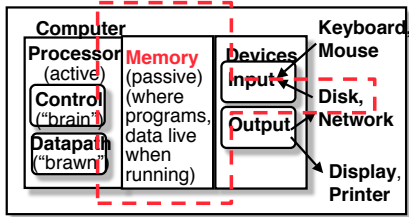
The matrix...reality? =>  
 MIT neuroscientists can now decipher part of the code involved in recognizing visual objects! A "classifier" was used on the monkey brain signals.  
[www.physorg.com/news7879.html](http://www.physorg.com/news7879.html)



**Review**

- Pipeline challenge is hazards
    - Forwarding helps w/many data hazards
    - Delayed branch helps with control hazard in 5 stage pipeline
  - More aggressive performance:
    - Superscalar
    - Out-of-order execution
  - You can be creative with your pipelines
    - Learn from our top 10 worst SW bugs...
      - Test, test, test. Expect the unexpected.
- Design w/failure as possibility! Redundancy!

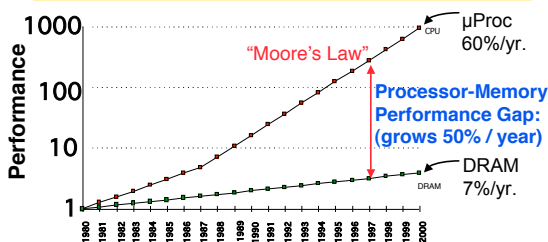
**The Big Picture**



**Memory Hierarchy (1/3)**

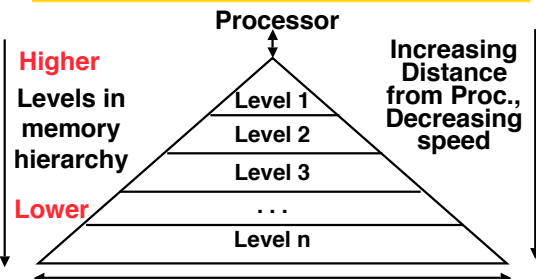
- Processor
  - executes instructions on order of nanoseconds to picoseconds
  - holds a small amount of code and data in registers
- Memory
  - More capacity than registers, still limited
  - Access time ~50-100 ns
- Disk
  - HUGE capacity (virtually limitless)
  - VERY slow: runs ~milliseconds

**Review: Why We Use Caches**



- 1989 first Intel CPU with cache on chip
- 1998 Pentium III has two levels of cache on chip

**Memory Hierarchy (2/3)**



Size of memory at each level  
 As we move to deeper levels the latency goes up and price per bit goes down.  
 Q: Can \$/bit go up as move deeper?



### Direct-Mapped Cache (1/2)

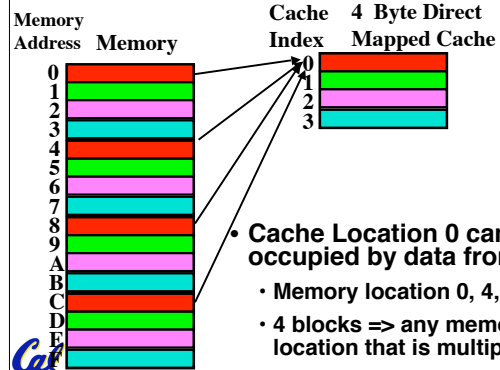
- In a **direct-mapped cache**, each memory address is associated with one possible **block** within the cache
- Therefore, we only need to look in a single location in the cache for the data if it exists in the cache
- Block is the unit of transfer between cache and memory



CS61C L21 Caches I (16)

Garcia, Fall 2005 © UCB

### Direct-Mapped Cache (2/2)



CS61C L21 Caches I (17)

Garcia, Fall 2005 © UCB

### Issues with Direct-Mapped

Tag Index Offset

- Since multiple memory addresses map to same cache index, how do we tell which one is in there?
- What if we have a block size > 1 byte?
- Answer: divide memory address into three fields



CS61C L21 Caches I (18)

Garcia, Fall 2005 © UCB

### Direct-Mapped Cache Terminology

- All fields are read as unsigned integers.
- Index**: specifies the cache index (which "row" of the cache we should look in)
- Offset**: once we've found correct block, specifies which byte within the block we want -- i.e., which "column"
- Tag**: the remaining bits after offset and index are determined; these are used to distinguish between all the memory addresses that map to the same location



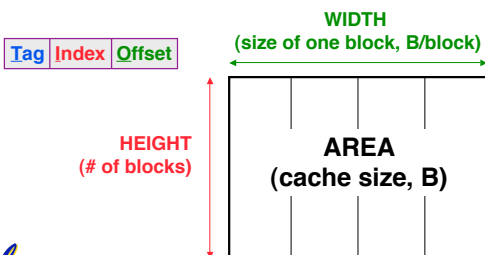
CS61C L21 Caches I (19)

Garcia, Fall 2005 © UCB

### TIO Dan's great cache mnemonic

$$\text{AREA (cache size, B)} = \text{HEIGHT (\# of blocks)} * \text{WIDTH (size of one block, B/block)}$$

$$2^{(H+W)} = 2^H * 2^W$$



CS61C L21 Caches I (20)

Garcia, Fall 2005 © UCB

### Direct-Mapped Cache Example (1/3)

- Suppose we have a 16KB of data in a direct-mapped cache with 4 word blocks
- Determine the size of the tag, index and offset fields if we're using a 32-bit architecture
- Offset**
  - need to specify correct byte within a block
  - block contains 4 words = 16 bytes =  $2^4$  bytes
  - need **4 bits** to specify correct byte



CS61C L21 Caches I (21)

Garcia, Fall 2005 © UCB

### Direct-Mapped Cache Example (2/3)

- **Index:** (~index into an “array of blocks”)
  - need to specify correct row in cache
  - cache contains 16 KB =  $2^{14}$  bytes
  - block contains  $2^4$  bytes (4 words)
  - # blocks/cache
    - =  $\frac{\text{bytes/cache}}{\text{bytes/block}}$
    - =  $\frac{2^{14} \text{ bytes/cache}}{2^4 \text{ bytes/block}}$
    - =  $2^{10}$  blocks/cache
  - need **10 bits** to specify this many rows



### Direct-Mapped Cache Example (3/3)

- **Tag:** use remaining bits as tag
  - tag length = addr length - offset - index  
= 32 - 4 - 10 bits  
= 18 bits
  - so tag is leftmost **18 bits** of memory address
- **Why not full 32 bit address as tag?**
  - All bytes within block need same address (4b)
  - Index must be same for every address within a block, so its redundant in tag check, thus can leave off to save memory (10 bits in this example)



### Caching Terminology

- When we try to read memory, 3 things can happen:
  1. **cache hit:** cache block is valid and contains proper address, so read desired word
  2. **cache miss:** nothing in cache in appropriate block, so fetch from memory
  3. **cache miss, block replacement:** wrong data is in cache at appropriate block, so discard it and fetch desired data from memory (cache always copy)



### Accessing data in a direct mapped cache

- **Ex.: 16KB of data, direct-mapped, 4 word blocks**

Address (hex)	Value of Word
...	...
00000010	a
00000014	b
00000018	c
0000001C	d
...	...
00000030	e
00000034	f
00000038	g
0000003C	h
...	...
00008010	i
00008014	j
00008018	k
0000801C	l
...	...
- **Read 4 addresses**
  1. 0x00000014
  2. 0x0000001C
  3. 0x00000034
  4. 0x00008014
- **Memory values on right:**
  - only cache/memory level of hierarchy



### Accessing data in a direct mapped cache

- **4 Addresses:**
  - 0x00000014, 0x0000001C, 0x00000034, 0x00008014
- **4 Addresses divided (for convenience) into Tag, Index, Byte Offset fields**

```

000000000000000000 0000000001 0100
000000000000000000 0000000001 1100
000000000000000000 0000000011 0100
000000000000000010 0000000001 0100
                Tag           Index   Offset
            
```



### 16 KB Direct Mapped Cache, 16B blocks

- **Valid bit:** determines whether anything is stored in that row (when computer initially turned on, all entries invalid)
- | Index | Tag | 0x0-3 | 0x4-7 | 0x8-b | 0xc-f |
|-------|-----|-------|-------|-------|-------|
| 0     | 0   |       |       |       |       |
| 1     | 0   |       |       |       |       |
| 2     | 0   |       |       |       |       |
| 3     | 0   |       |       |       |       |
| 4     | 0   |       |       |       |       |
| 5     | 0   |       |       |       |       |
| 6     | 0   |       |       |       |       |
| 7     | 0   |       |       |       |       |
| ...   | ... |       |       |       |       |
| 1022  | 0   |       |       |       |       |
| 1023  | 0   |       |       |       |       |



### 1. Read 0x00000014

• 00000000000000000000 0000000001 0100  
 Valid Tag field Index field Offset

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	0				
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				
...	...				
1022	0				
1023	0				



### So we read block 1 (000000001)

• 00000000000000000000 0000000001 0100  
 Valid Tag field Index field Offset

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1				
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				
...	...				
1022	0				
1023	0				



### No valid data

• 00000000000000000000 0000000001 0100  
 Valid Tag field Index field Offset

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	0				
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				
...	...				
1022	0				
1023	0				



### So load that data into cache, setting tag, valid

• 00000000000000000000 0000000001 0100  
 Valid Tag field Index field Offset

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	a	b	c	d
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				
...	...				
1022	0				
1023	0				



### Read from cache at offset, return word b

• 00000000000000000000 0000000001 0100  
 Valid Tag field Index field Offset

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	a	b	c	d
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				
...	...				
1022	0				
1023	0				



### 2. Read 0x0000001C = 0...00 0..001 1100

• 00000000000000000000 0000000001 1100  
 Valid Tag field Index field Offset

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	a	b	c	d
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				
...	...				
1022	0				
1023	0				



### Index is Valid

- 00000000000000000000 0000000001 1100

Index	Valid		Tag field		Index field		Offset
	Tag	0x0-3	0x4-7	0x8-b	0xc-f		
0	0						
1	1	0	a	b	c	d	
2	0						
3	0						
4	0						
5	0						
6	0						
7	0						
...							
1022	0						
1023	0						



CS61C L21 Caches I (34)

Garcia, Fall 2005 © UCB

### Index valid, Tag Matches

- 00000000000000000000 0000000001 1100

Index	Valid		Tag field		Index field		Offset
	Tag	0x0-3	0x4-7	0x8-b	0xc-f		
0	0						
1	1	0	a	b	c	d	
2	0						
3	0						
4	0						
5	0						
6	0						
7	0						
...							
1022	0						
1023	0						



CS61C L21 Caches I (35)

Garcia, Fall 2005 © UCB

### Index Valid, Tag Matches, return d

- 00000000000000000000 0000000001 1100

Index	Valid		Tag field		Index field		Offset
	Tag	0x0-3	0x4-7	0x8-b	0xc-f		
0	0						
1	1	0	a	b	c	d	
2	0						
3	0						
4	0						
5	0						
6	0						
7	0						
...							
1022	0						
1023	0						



CS61C L21 Caches I (36)

Garcia, Fall 2005 © UCB

### 3. Read 0x00000034 = 0...00 0..011 0100

- 00000000000000000000 0000000011 0100

Index	Valid		Tag field		Index field		Offset
	Tag	0x0-3	0x4-7	0x8-b	0xc-f		
0	0						
1	1	0	a	b	c	d	
2	0						
3	0						
4	0						
5	0						
6	0						
7	0						
...							
1022	0						
1023	0						



CS61C L21 Caches I (37)

Garcia, Fall 2005 © UCB

### So read block 3

- 00000000000000000000 0000000011 0100

Index	Valid		Tag field		Index field		Offset
	Tag	0x0-3	0x4-7	0x8-b	0xc-f		
0	0						
1	1	0	a	b	c	d	
2	0						
3	0						
4	0						
5	0						
6	0						
7	0						
...							
1022	0						
1023	0						



CS61C L21 Caches I (38)

Garcia, Fall 2005 © UCB

### No valid data

- 00000000000000000000 0000000011 0100

Index	Valid		Tag field		Index field		Offset
	Tag	0x0-3	0x4-7	0x8-b	0xc-f		
0	0						
1	1	0	a	b	c	d	
2	0						
3	0						
4	0						
5	0						
6	0						
7	0						
...							
1022	0						
1023	0						



CS61C L21 Caches I (39)

Garcia, Fall 2005 © UCB

### Load that cache block, return word f

- 00000000000000000000 0000000011 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	0	a	b	c
2	0				
3	1	0	e	f	g
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				



### 4. Read 0x00008014 = ...10 0..001 0100

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	0	a	b	c
2	0				
3	1	0	e	f	g
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				



### So read Cache Block 1, Data is Valid

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	0	a	b	c
2	0				
3	1	0	e	f	g
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				



### Cache Block 1 Tag does not match (0 != 2)

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	0	a	b	c
2	0				
3	1	0	e	f	g
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				



### Miss, so replace block 1 with new data & tag

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	2	i	j	k
2	0				
3	1	0	e	f	g
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				



### And return word j

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	1	2	i	j	k
2	0				
3	1	0	e	f	g
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				



### Do an example yourself. What happens?

- Chose from: Cache: Hit, Miss, Miss w. replace  
Values returned: a, b, c, d, e, ..., k, l
- Read address **0x00000030** ?  
000000000000000000 0000000011 0000
- Read address **0x0000001c** ?  
000000000000000000 0000000001 1100

Cache

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	2	i	j	k	l
2	0				
3	0	e	f	g	h
4	0				
5	0				
6	0				
7	0				

CS61C L21 Caches I (46) Garcia, Fall 2005 © UCB

### Answers

- **0x00000030** a **hit**  
Index = 3, Tag matches, Offset = 0, value = e
- **0x0000001c** a **miss**  
Index = 1, Tag mismatch, so replace from memory, Offset = 0xc, value = d
- Since reads, values must = memory values whether or not cached:
  - 0x00000030 = e
  - 0x0000001c = d

Memory

Address	Value of Word
...	...
00000010	a
00000014	b
00000018	c
0000001c	d
...	...
00000030	e
00000034	f
00000038	g
0000003c	h
...	...
00008010	i
00008014	j
00008018	k
0000801c	l
...	...

CS61C L21 Caches I (47) Garcia, Fall 2005 © UCB

### Peer Instruction

- Mem hierarchies were invented before 1950. (UNIVAC I wasn't delivered 'til 1951)
- If you know your computer's cache size, you can often make your code run faster.
- Memory hierarchies take advantage of spatial locality by keeping the most recent data items closer to the processor.

	ABC
1:	FFF
2:	FFT
3:	FTF
4:	FTT
5:	TFF
6:	TFT
7:	FTF
8:	FTT

CS61C L21 Caches I (48) Garcia, Fall 2005 © UCB

### Peer Instructions

- All caches take advantage of spatial locality.
- All caches take advantage of temporal locality.
- On a read, the return value will depend on what is in the cache.

	ABC
1:	FFF
2:	FFT
3:	FTF
4:	FTT
5:	TFF
6:	TFT
7:	FTF
8:	FTT

CS61C L21 Caches I (49) Garcia, Fall 2005 © UCB

### And in Conclusion (1/2)

- We would like to have the capacity of disk at the speed of the processor: unfortunately this is not feasible.
- So we create a memory hierarchy:
  - each successively lower level contains "most used" data from next higher level
  - exploits temporal locality
  - do the common case fast, worry less about the exceptions (design principle of MIPS)
- Locality of reference is a Big Idea

CS61C L21 Caches I (52) Garcia, Fall 2005 © UCB

### And in Conclusion (2/2)

- Mechanism for transparent movement of data among levels of a storage hierarchy
  - set of address/value bindings
  - address  $\Rightarrow$  index to set of candidates
  - compare desired address with tag
  - service hit or miss
    - load new block and binding on miss

address: tag index offset  
000000000000000000 0000000001 1100

Valid	Tag	0x0-3	0x4-7	0x8-b	0xc-f
0	0				
1	0	a	b	c	d
...	...	...	...	...	...

CS61C L21 Caches I (53) Garcia, Fall 2005 © UCB