

The Memory Hierarchy

15-213/14-513/15-513: Introduction to Computer Systems
9th Lecture, Sept 24, 2024

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

1

Writing & Reading Memory

■ Write

- Transfer data from CPU to memory
`movq %rax, 8(%rsp)`
- “Store” operation

■ Read

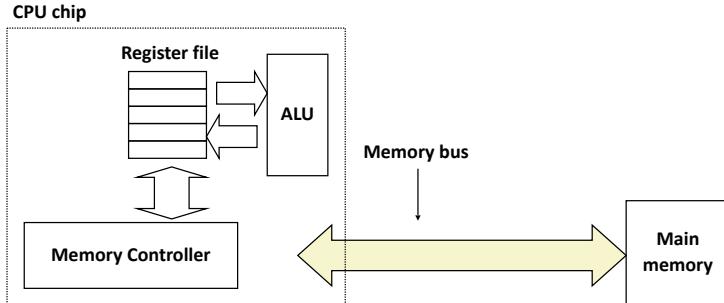
- Transfer data from memory to CPU
`movq 8(%rsp), %rax`
- “Load” operation

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

2

Modern Connection between CPU and Memory

- A **bus** is a collection of parallel wires that carry address, data, and control signals.
- Buses are typically shared by multiple devices.

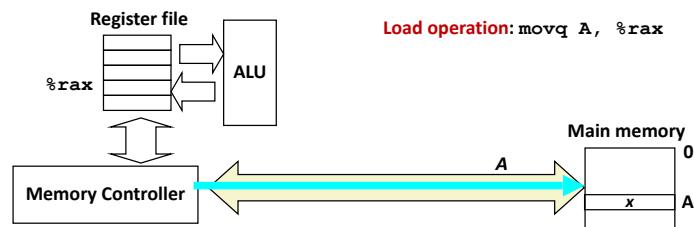


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

3

Memory Read Transaction (1)

- CPU places address A on the memory bus.

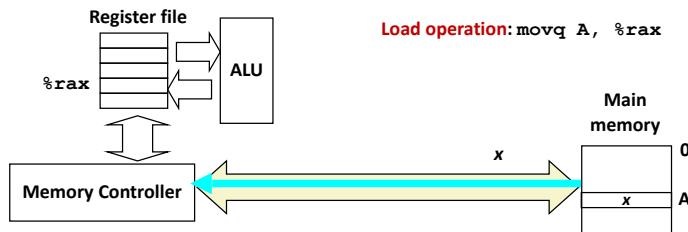


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

4

Memory Read Transaction (2)

- Main memory reads A from the memory bus, retrieves word x, and places it on the bus.

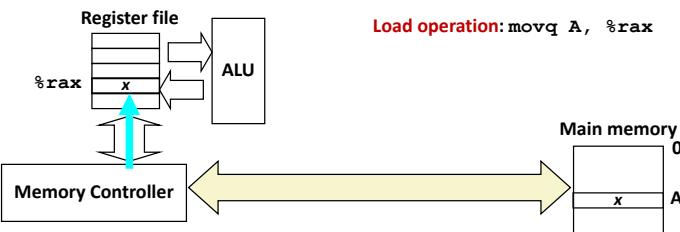


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

5

Memory Read Transaction (3)

- CPU reads word x from the bus and copies it into register $\%rax$.



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

6

Memory Write Transaction (1)

- CPU places address A on bus. Main memory reads it and waits for the corresponding data word to arrive.

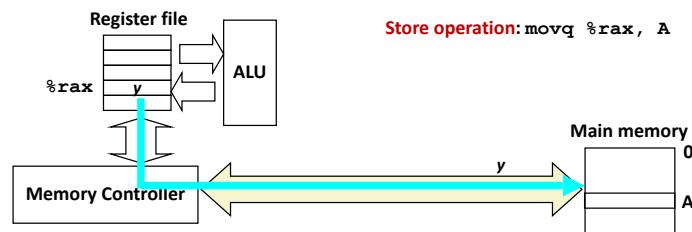


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

7

Memory Write Transaction (2)

- CPU places data word y on the bus.

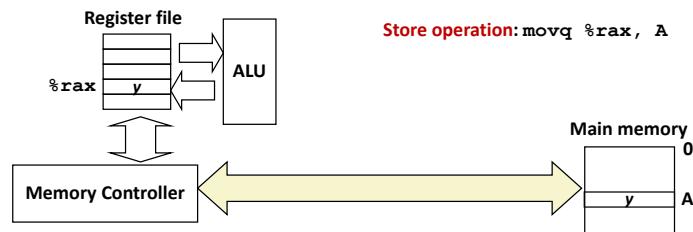


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

8

Memory Write Transaction (3)

- Main memory reads data word y from the bus and stores it at address A.



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

9

Random-Access Memory (RAM)

■ Key features

- RAM is traditionally packaged as a chip.
 - or embedded as part of processor chip
- Basic storage unit is normally a **cell** (one bit per cell).
- Multiple RAM chips form a memory.

■ RAM comes in two varieties:

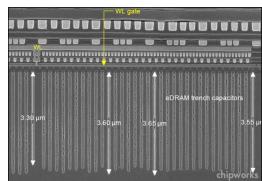
- SRAM (Static RAM)
- DRAM (Dynamic RAM)

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

10

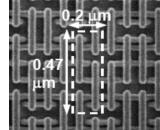
RAM Technologies

■ DRAM



- 1 Transistor + 1 capacitor / bit
 - Capacitor oriented vertically
- Must refresh state periodically

■ SRAM



- 6 transistors / bit
- Holds state indefinitely (but will still lose data on power loss)

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

11

SRAM vs DRAM Summary

	Trans. per bit	Access time	Needs refresh?	Needs EDC?	Cost	Applications
SRAM	6 or 8	1x	No	Maybe	100x	Cache memories
DRAM	1	10x	Yes	Yes	1x	Main memories, frame buffers

EDC: Error detection and correction

■ Trends

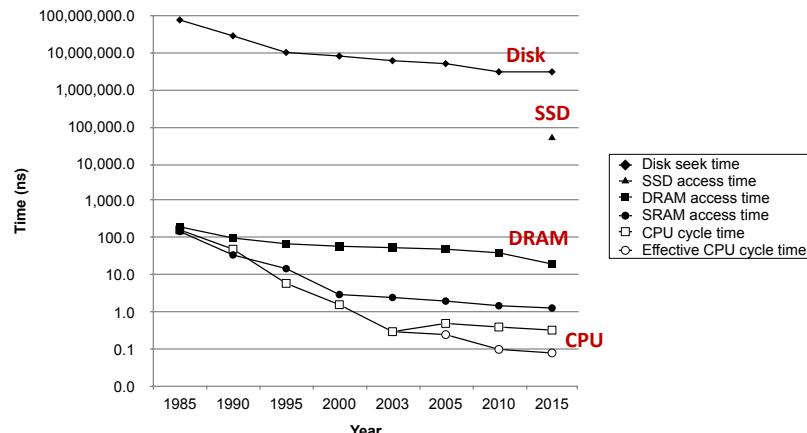
- SRAM scales with semiconductor technology
 - Reaching its limits
- DRAM scaling limited by need for minimum capacitance
 - Aspect ratio limits how deep can make capacitor
 - Also reaching its limits

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

12

The CPU-Memory Gap

The gap *widens* between DRAM, disk, and CPU speeds.



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

13

Locality to the Rescue!

The key to bridging this CPU-Memory gap is an important property of computer programs known as **locality**.

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

14

Locality

■ **Principle of Locality:** Many Programs tend to use data and instructions with addresses near or equal to those they have used recently.

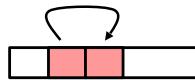
■ **Temporal locality:**

- Recently referenced items are likely to be referenced again in the near future



■ **Spatial locality:**

- Items with nearby addresses tend to be referenced close together in time



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

15

Locality Example

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

Spatial or Temporal
Locality?

spatial

temporal

spatial

temporal

■ **Data references**

- Reference array elements in succession (stride-1 reference pattern).
- Reference variable **sum** each iteration.

■ **Instruction references**

- Reference instructions in sequence.
- Cycle through loop repeatedly.

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

16

Qualitative Estimates of Locality

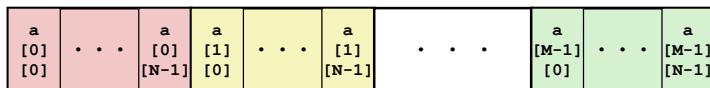
- **Claim:** Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.
- **Question:** Does this function have good locality with respect to array a?

Hint: array layout is row-major order

Answer: yes

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}
```



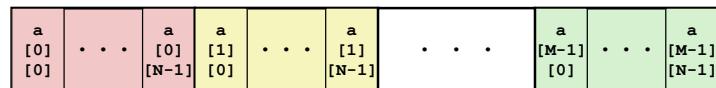
Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

17

Answer: no

Stride N reference pattern

Note: If M is very small then good locality. Why?



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

Locality Example

- **Question:** Can you permute the loops so that the function scans the 3-d array a with a stride-1 reference pattern (and thus has good spatial locality)?

```
int sum_array_3d(int a[M][N][N])
{
    int i, j, k, sum = 0;

    for (i = 0; i < N; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < M; k++)
                sum += a[k][i][j];
    return sum;
}
```

Answer: make j the inner loop

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

19

Memory Hierarchies

- Some fundamental and enduring properties of hardware and software:

- Fast storage technologies cost more per byte, have less capacity, and require more power (heat!).
- The gap between CPU and main memory speed is widening.
- Well-written programs tend to exhibit good locality.

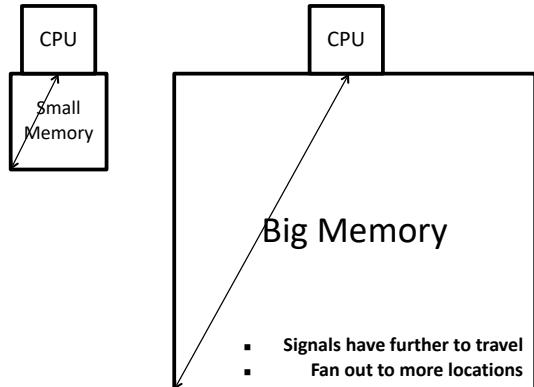
- These properties complement each other well for many types of programs.

- They suggest an approach for organizing memory and storage systems known as a **memory hierarchy**.

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

20

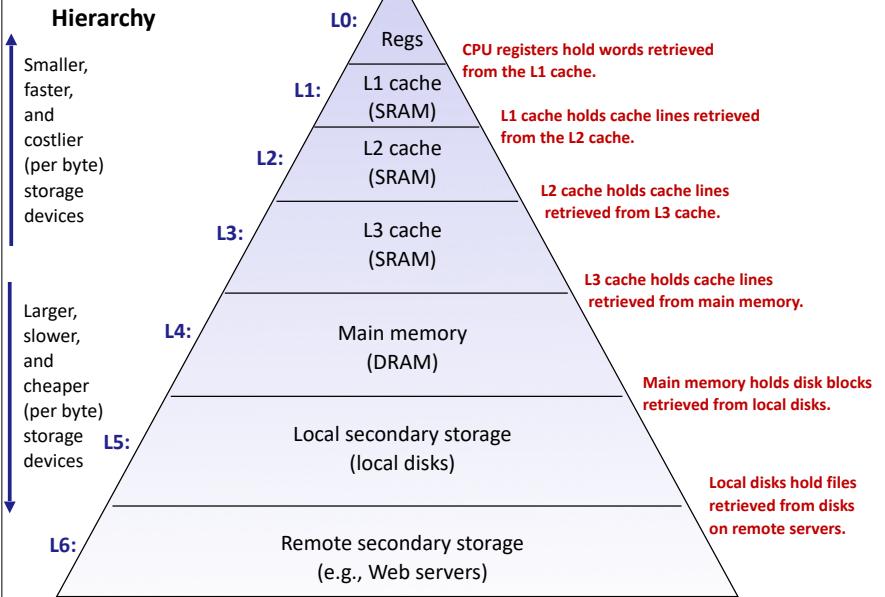
Memory size affects latency & energy



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

21

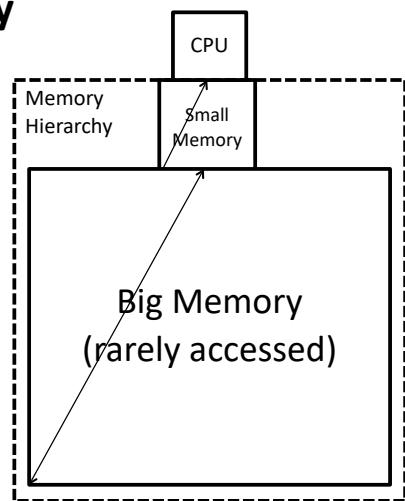
Example Memory Hierarchy



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

22

Hierarchy provides the illusion of large & fast memory

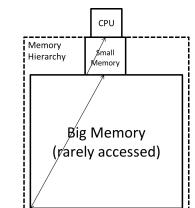


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

23

Caches

- **Cache:** A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.
- **Fundamental idea of a memory hierarchy:**
 - For each k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1.
- **Why do memory hierarchies work?**
 - Because of locality: programs tend to access the data at level k more often than they access the data at level k+1.
 - Thus, the storage at level k+1 can be slower, and thus larger and cheaper per bit.
- **Big Idea (Ideal):** The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

24

Cache vs Memory

Caches are invisible (transparent) to software

- Managed by hardware in response to loads & stores
- Performance & energy improve without software changes
- Caveat: Recent CPUs have some instructions to manage cache (e.g., prefetch, invalidate, partition...)

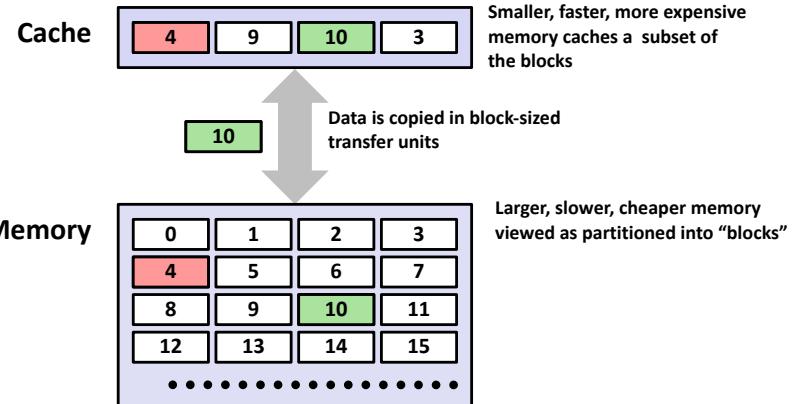
Memory is visible to software

- I.e., addressable directly by instructions (memory address, registers)
- Some optimization opportunities, but only w/ software changes

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

27

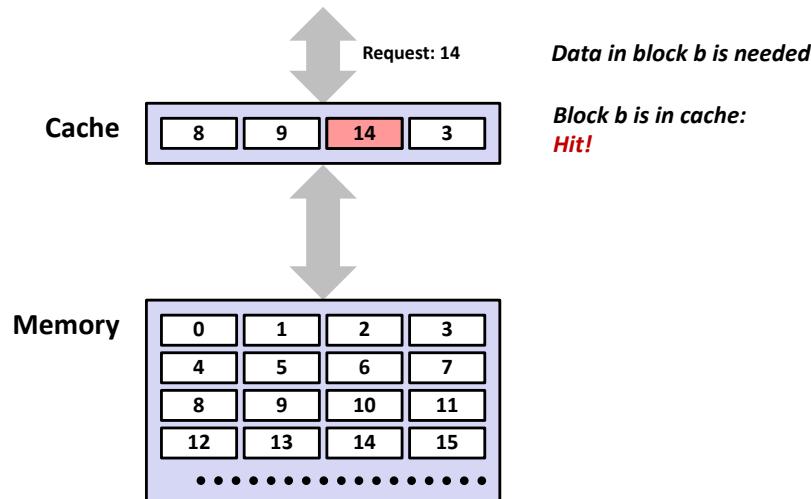
General Cache Concepts



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

26

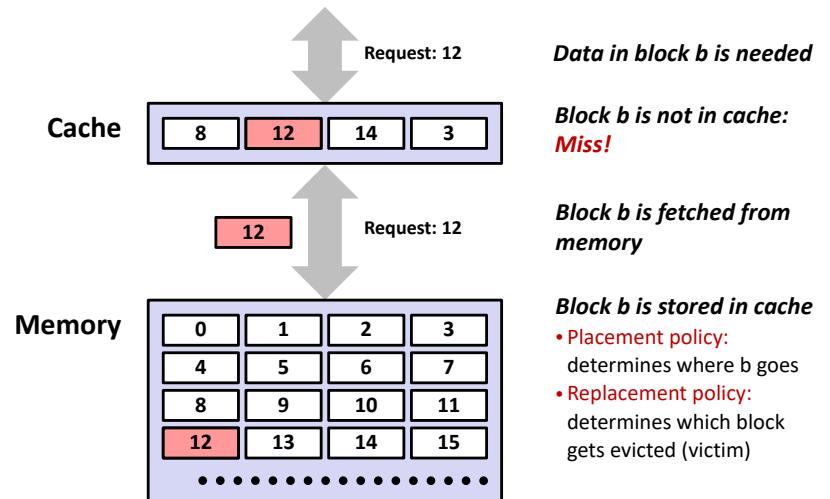
General Cache Concepts: Hit



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

27

General Cache Concepts: Miss



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

28

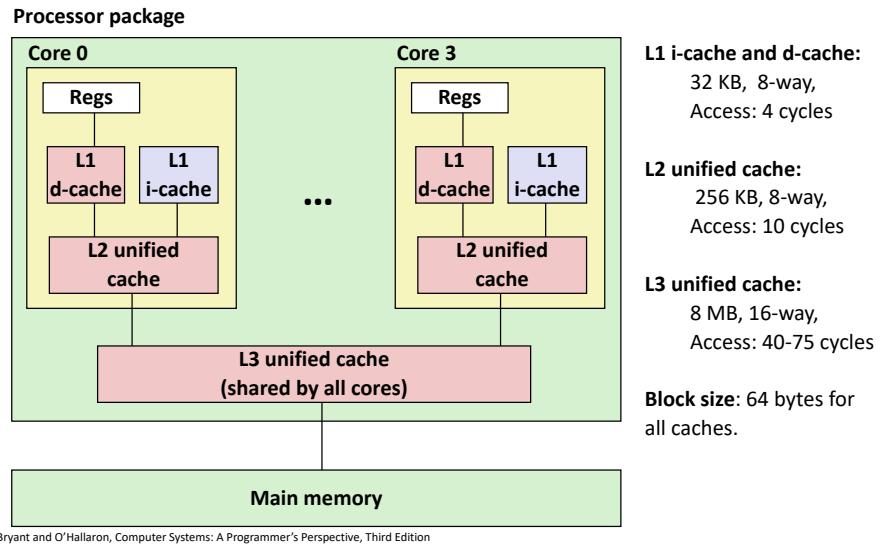
Examples of Caching in the Mem. Hierarchy

Cache Type	What is Cached?	Where is it Cached?	Latency (cycles)	Managed By
Registers	4-8 byte words	CPU core	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware MMU
L1 cache	64-byte blocks	On-Chip L1	4	Hardware
L2 cache	64-byte blocks	On-Chip L2	10	Hardware
Virtual Memory	4-KB pages	Main memory	100	Hardware + OS
Buffer cache	Parts of files	Main memory	100	OS
Disk cache	Disk sectors	Disk controller	100,000	Disk firmware
Network buffer cache	Parts of files	Local disk	10,000,000	NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

29

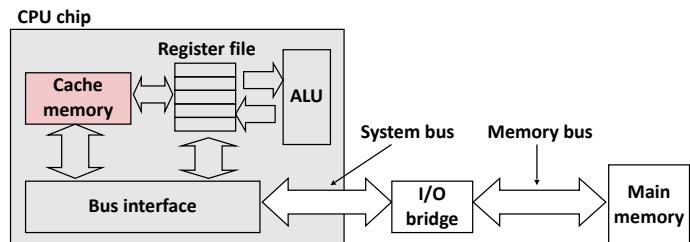
Intel Core i7 Cache Hierarchy



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

CPU Cache Memories

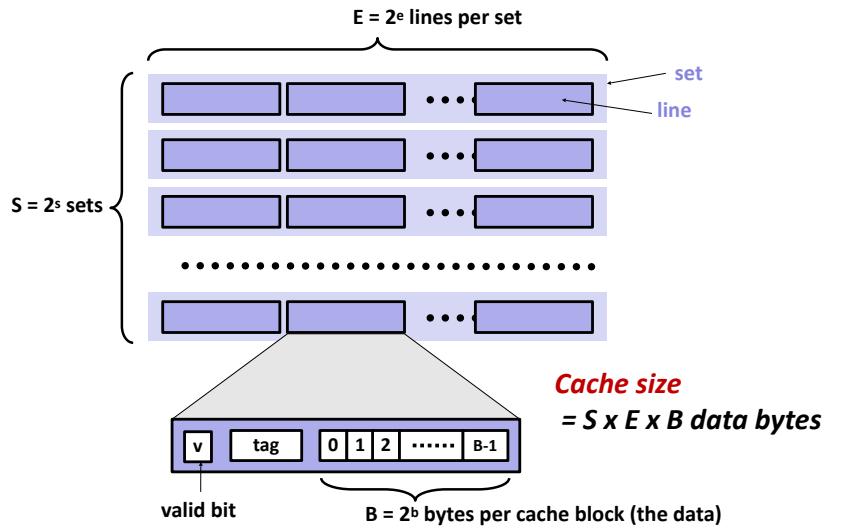
- **CPU cache memories** are small, fast SRAM-based memories managed automatically in hardware
 - Hold frequently accessed blocks of main memory
- CPU looks first for data in cache
- Typical system structure:



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

31

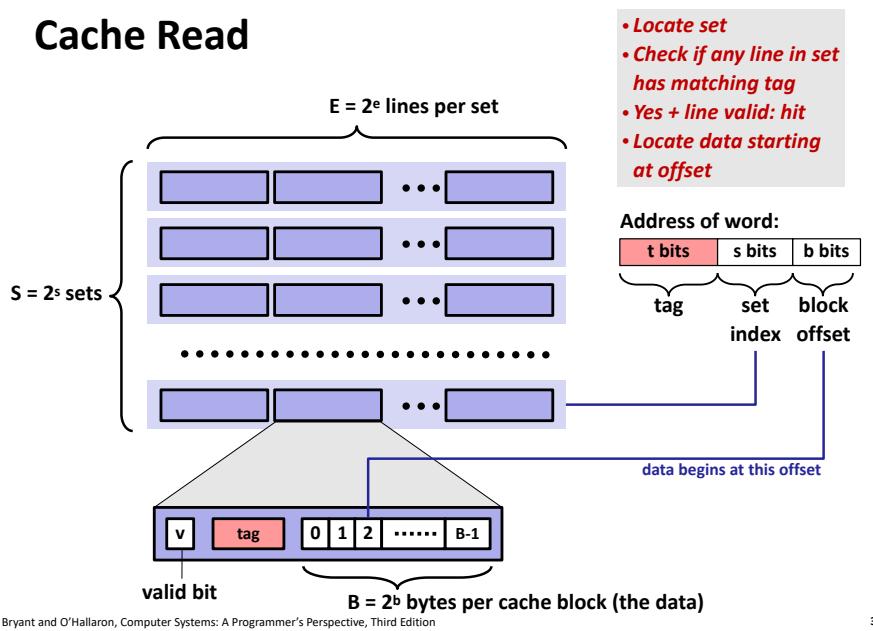
General Cache Organization (S, E, B)



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

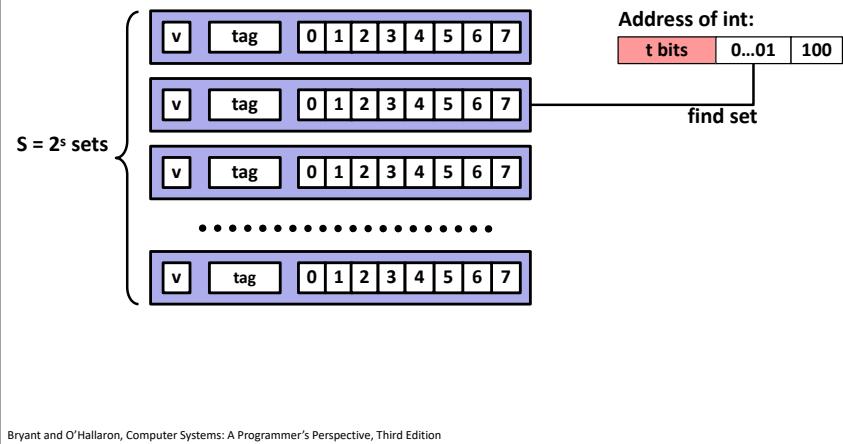
32

Cache Read



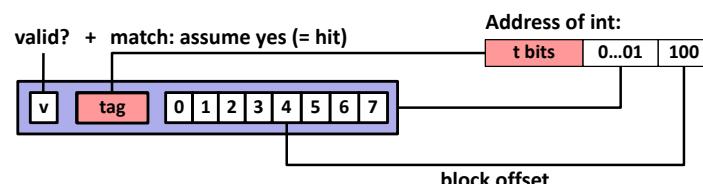
Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size $B=8$ bytes



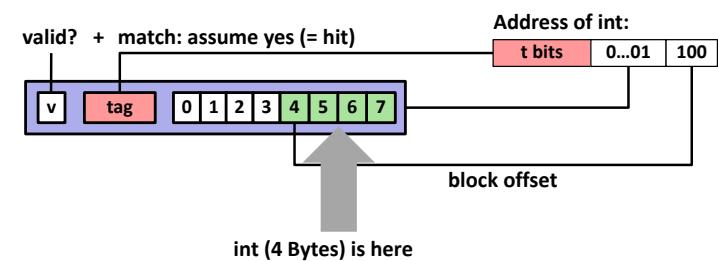
Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size $B=8$ bytes



Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size $B=8$ bytes



If tag doesn't match (= miss): old line is evicted and replaced

Direct-Mapped Cache Simulation

$t=1 \quad s=2 \quad b=1$

x	xx	x
---	----	---

4-bit addresses (address space size $M=16$ bytes)
 $S=4$ sets, $E=1$ blocks/set, $B=2$ bytes/block

Address trace (reads, one byte per read):

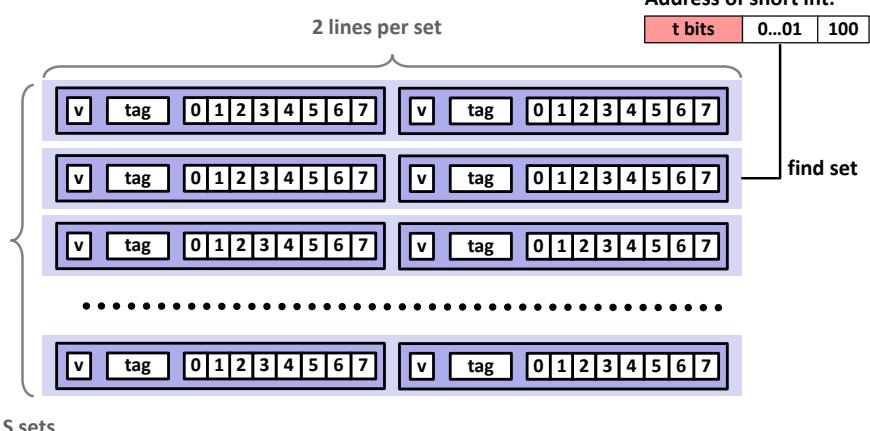
0	[0000] ₂	miss	(cold)
1	[0001] ₂	hit	
7	[0111] ₂	miss	(cold)
8	[1000] ₂	miss	(cold)
0	[0000] ₂	miss	(conflict)

	v	Tag	Block
Set 0	1	0	M[0-1]
Set 1	0		
Set 2	0		
Set 3	1	0	M[6-7]

E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set

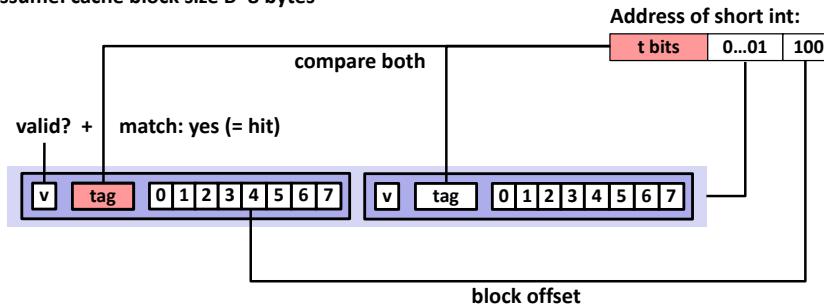
Assume: cache block size $B=8$ bytes



E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set

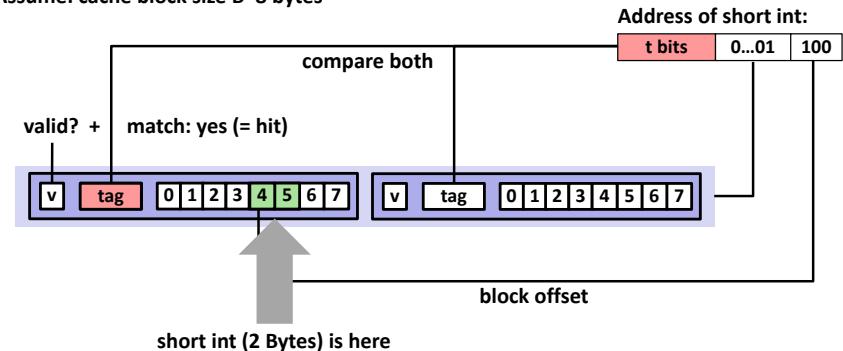
Assume: cache block size $B=8$ bytes



E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size $B=8$ bytes



No match or not valid (= miss):

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

2-Way Set Associative Cache Simulation

t=2 s=1 b=1
 xx x x

4-bit addresses (M=16 bytes)
 S=2 sets, E=2 blocks/set, B=2 bytes/block

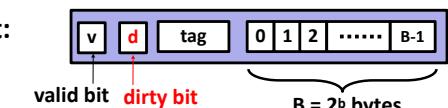
Address trace (reads, one byte per read):

0	[0000] ₂	miss	(cold)
1	[0001] ₂	hit	
7	[0111] ₂	miss	(cold)
8	[1000] ₂	miss	(cold)
0	[0000] ₂	hit	

	v	Tag	Block
Set 0	1	00	M[0-1]
	1	10	M[8-9]
Set 1	1	01	M[6-7]
	0		

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

41



What about writes?

■ Multiple copies of data exist:

- L1, L2, L3, Main Memory

■ What to do on a write-hit?

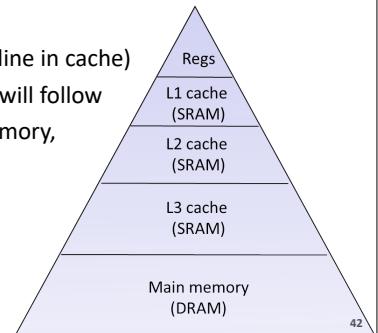
- Write-through (write immediately to memory)
- Write-back (defer write to memory until replacement of line)
 - Each cache line needs a dirty bit (set if data has been written to)

■ What to do on a write-miss?

- Write-allocate (load into cache, update line in cache)
 - Good if more writes to the location will follow
- No-write-allocate (writes straight to memory, does not load into cache)

■ Typical

- Write-through + No-write-allocate
- Write-back + Write-allocate

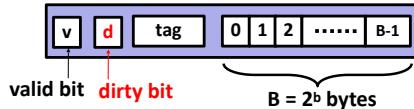


Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

42

Practical Write-back Write-allocate

■ A write to address X is issued



■ If it is a hit

- Update the contents of block
- Set dirty bit to 1 (bit is sticky and only cleared on eviction)

■ If it is a miss

- Fetch block from memory (per a read miss)
- Then perform the write operations (per a write hit)

■ If a line is evicted and dirty bit is set to 1

- The entire block of 2^b bytes are written back to memory
- Dirty bit is cleared (set to 0)
- Line is replaced by new contents

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

43

Cache Performance Metrics

■ Miss Rate

- Fraction of memory references not found in cache (misses / accesses) = 1 - hit rate
- Typical numbers (in percentages):
 - 3-10% for L1
 - can be quite small (e.g., < 1%) for L2, depending on size, etc.

■ Hit Time

- Time to deliver a line in the cache to the processor
 - includes time to determine whether the line is in the cache
- Typical numbers:
 - 4 clock cycle for L1
 - 10 clock cycles for L2

■ Miss Penalty

- Additional time required because of a miss
 - typically 50-200 cycles for main memory (Trend: increasing!)

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

44

Let's think about those numbers

■ Huge difference between a hit and a miss

- Could be 100x, if just L1 and main memory

■ Would you believe 99% hits is twice as good as 97%?

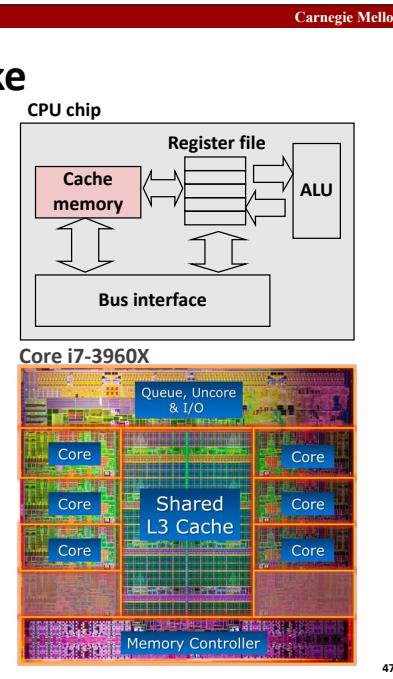
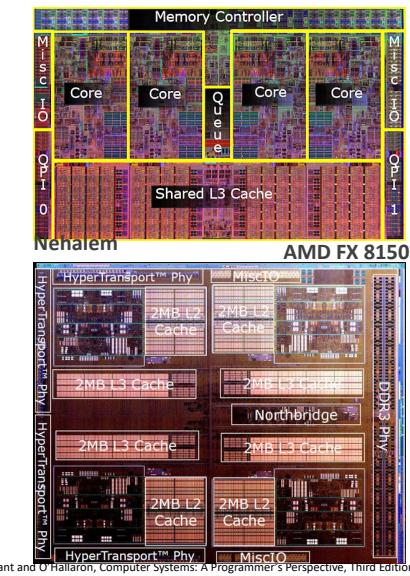
- Consider this simplified example:
cache hit time of 1 cycle
miss penalty of 100 cycles
- Average access time:
97% hits: 1 cycle + 0.03×100 cycles = **4 cycles**
99% hits: 1 cycle + 0.01×100 cycles = **2 cycles**

■ This is why “miss rate” is used instead of “hit rate”

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

45

What it Really Looks Like



47

Writing Cache Friendly Code

■ Make the common case go fast

- Focus on the inner loops of the core functions

■ Minimize the misses in the inner loops

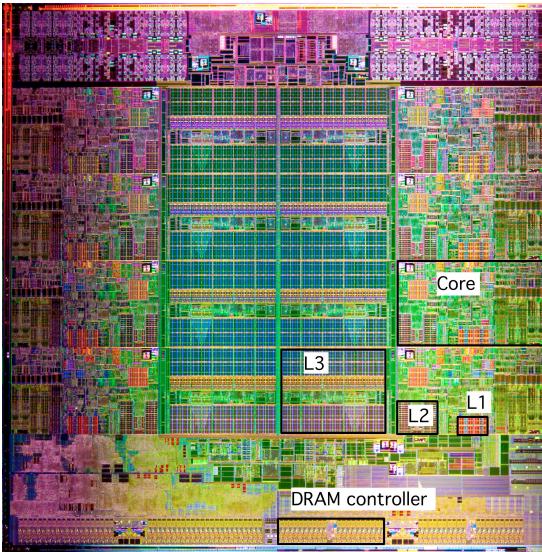
- Repeated references to variables are good (**temporal locality**)
- Stride-1 reference patterns are good (**spatial locality**)

Key idea: Our qualitative notion of locality is quantified through our understanding of cache memories

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

46

What it Really Looks Like (Cont.)



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

48

Example: Core i7 L1 Data Cache

32 kB 8-way set associative

64 bytes/block

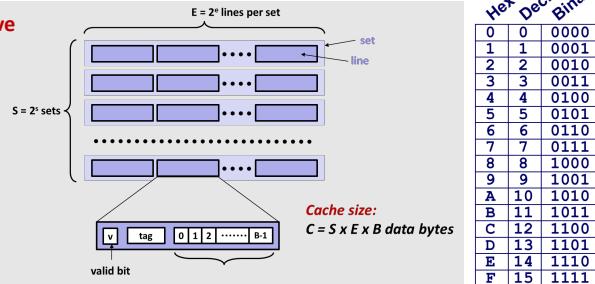
47 bit address range

B =

S = , **s** =

E = , **e** =

C =



Hex	Decimal	Binary
0	0	0000
1	1	0001
2	2	0010
3	3	0011
4	4	0100
5	5	0101
6	6	0110
7	7	0111
8	8	1000
9	9	1001
A	10	1010
B	11	1011
C	12	1100
D	13	1101
E	14	1110
F	15	1111

Address of word:
 t bits s bits b bits
 tag set index block offset
 Block offset: . bits
 Set index: . bits
 Tag: . bits

Stack Address:
 0x00007f7262a1e010

Block offset: 0x??
 Set index: 0x??
 Tag: 0x??

49

Example: Core i7 L1 Data Cache

32 kB 8-way set associative

64 bytes/block

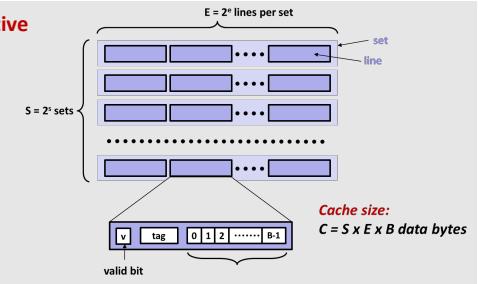
47 bit address range

B = 64

S = 64, **s** = 6

E = 8, **e** = 3

C = 64 x 64 x 8 = 32,768



Hex	Decimal	Binary
0	0	0000
1	1	0001
2	2	0010
3	3	0011
4	4	0100
5	5	0101
6	6	0110
7	7	0111
8	8	1000
9	9	1001
A	10	1010
B	11	1011
C	12	1100
D	13	1101
E	14	1110
F	15	1111

Address of word:
 t bits s bits b bits
 tag set index block offset
 Block offset: 6 bits
 Set index: 6 bits
 Tag: 35 bits

Stack Address:
 0x00007f7262a1e010
 Block offset: 0x10
 Set index: 0x0
 Tag: 0x7f7262a1e

50