

CSE 560

Computer Systems Architecture

Cache

1

Why Caches?

2

Programs 101

C Code

```
int sum(int x, int y)
{
    int t = x+y;
    return t;
}
```

Generated IA32 Assembly

```
sum:
    pushl %ebp
    movl %esp,%ebp
    movl 12(%ebp),%eax
    addl 8(%ebp),%eax
    popl %ebp
    ret
```

High-level behavior: Instructions that read from/write to memory...

- Read data from memory (put in **registers**)
- Manipulate it
- Store it back to memory

3

The Need for Speed

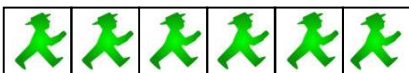
CPU Pipeline



4

The Need for Speed

CPU Pipeline



Instruction speeds:

- **add, sub, shift**: 1 cycle
- **mult**: 3 cycles
- **load/store**: **100 cycles**
off-chip 50(-70)ns
2(-3) GHz processor → 0.5 ns clock

5

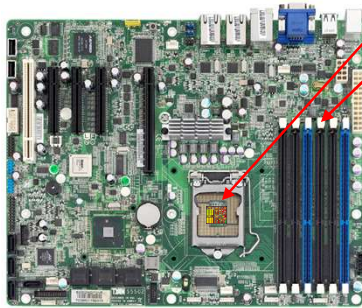
The Need for Speed

CPU Pipeline



6

What's the problem?



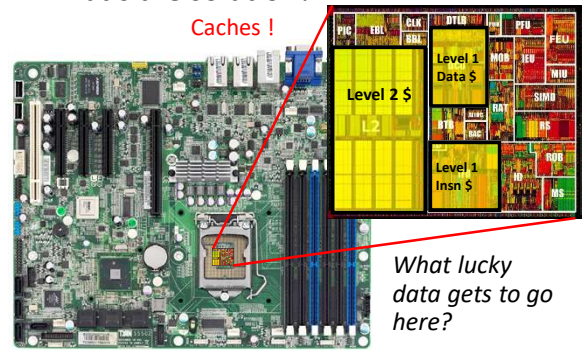
SandyBridge Motherboard, 2011
<http://news.softpedia.com>

- Processor
- Main Memory
- too slow
 - too far away

7

What's the solution?

Caches !



Intel Pentium 3, 1999

8

Locality Locality Locality

If you ask for something, you're likely to ask for:

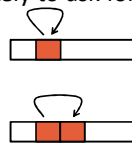
- the same thing again soon

→ Temporal Locality

- something near that thing, soon

→ Spatial Locality

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



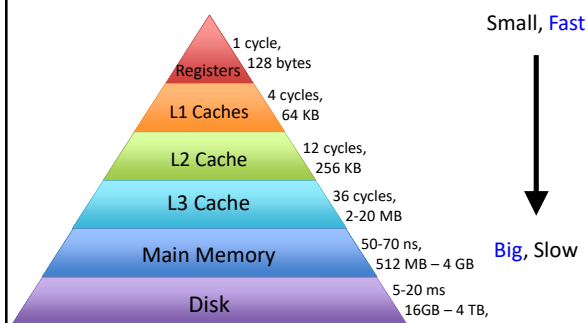
9

Your life is full of Locality



10

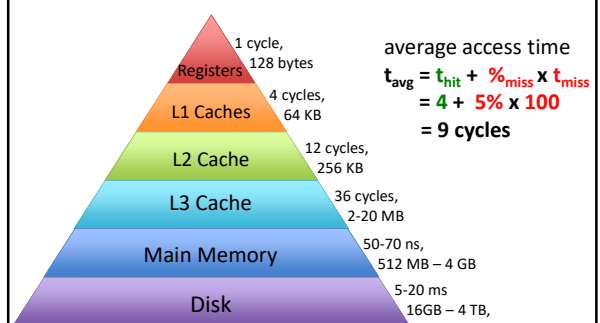
The Memory Hierarchy



Intel Haswell Processor, 2013

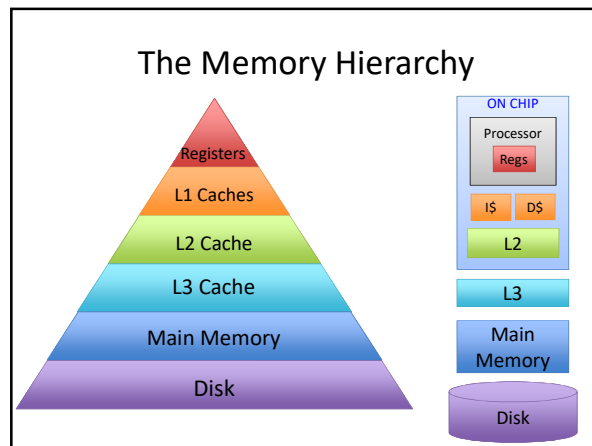
11

The Memory Hierarchy

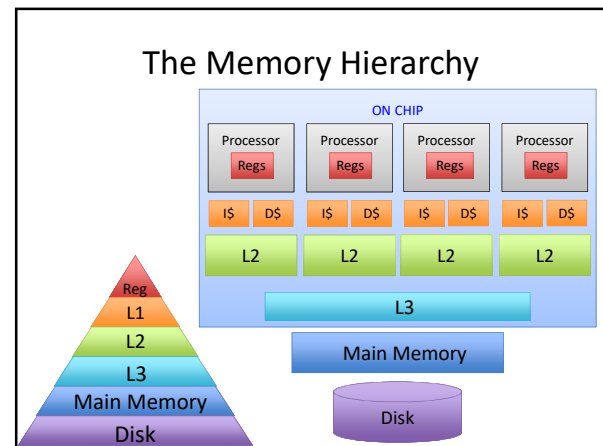


Intel Haswell Processor, 2013

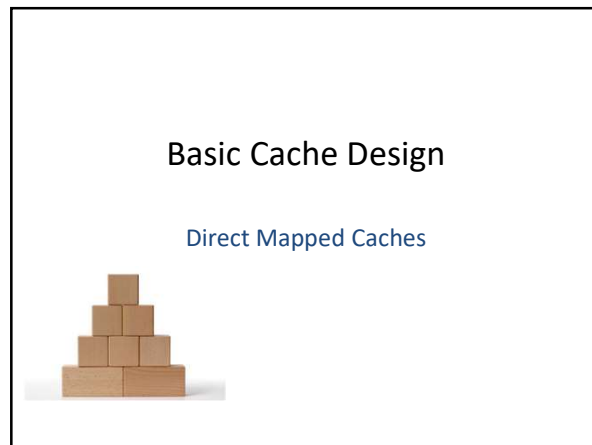
12



13



14



15

16 Byte Memory

load 0x1100 → r1

- Byte-addressable memory
- 4 address bits → 16 bytes total
- b addr bits → 2^b bytes in memory

MEMORY	
addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

16

4-Byte, Direct Mapped Cache

CACHE		
index	addr	data
00	xxxx	X
01	xxxx	X
10	xxxx	X
11	xxxx	X

- entry = row = **cache line** = **cache block**
- Block Size:** 1 byte
- Direct mapped:**
 - Each address mapped to specific cache block
 - 4 entries → 2 index bits ($2^n \rightarrow n$ bits)

MEMORY	
addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

17

Least Significant Bits as Index

index
XXXX

CACHE		
index	addr	data
00	0000	A
01	0001	B
10	0010	C
11	0011	D

- Supports spatial locality

MEMORY	
addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

18

Analogy to a Spice Rack

Spice Rack (Cache)

index spice

A B C D E F ... Z

Spice Wall (Memory)

- Compared to your spice wall
 - Smaller
 - Faster
 - More costly (per oz.)

http://www.bedbathandbeyond.com

19

Analogy to a Spice Rack

Spice Rack (Cache)

index tag spice

A B C D E F ... Z

innamon

Spice Wall (Memory)

- How do you know what's in the jar?
- Need labels

Tag = Ultra-minimalist label

20

4-Byte, Direct Mapped Cache

tag | index
XXXX

CACHE

index	tag	data
00	00	A
01	00	B
10	00	C
11	00	D

Tag: minimalist label/address
address = tag + index

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

21

Simulation #1 of a 4-byte, DM Cache

tag | index
XXXX

CACHE

index	tag	data
00	00	0
01	00	0
10	00	0
11	00	0

Cache starts empty

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x0000 Hit?

Lookup:
 ➔ Index into \$
 ➔ Check tag

22

4-Byte, Direct Mapped Cache

CACHE

index	V	tag	data
00	0	xx	X
01	0	xx	X
10	0	xx	X
11	0	xx	X

One last tweak: **valid bit**

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

23

Simulation #1 of a 4-byte, DM Cache

tag | index
XXXX

CACHE

index	V	tag	data
00	0	xx	X
01	0	xx	X
10	0	xx	X
11	0	xx	X

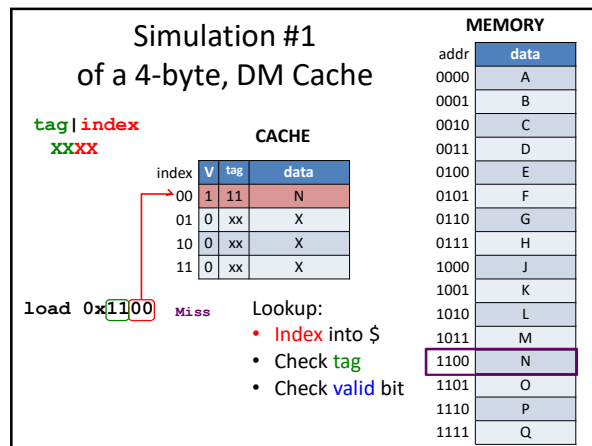
MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

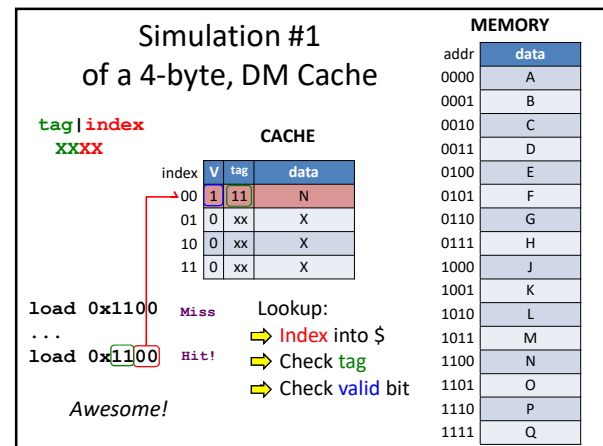
load 0x1100 Miss

Lookup:
 ➔ Index into \$
 ➔ Check tag
 ➔ Check valid bit

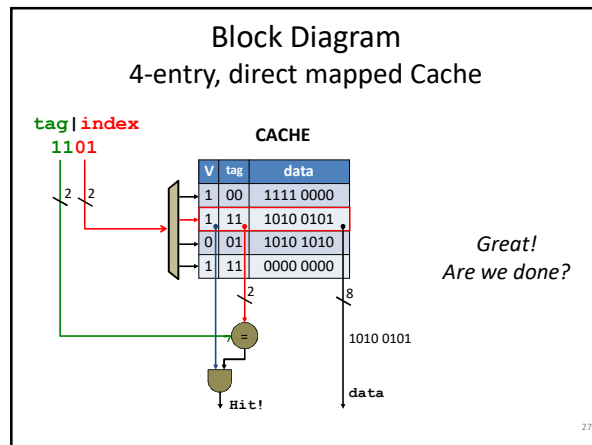
24



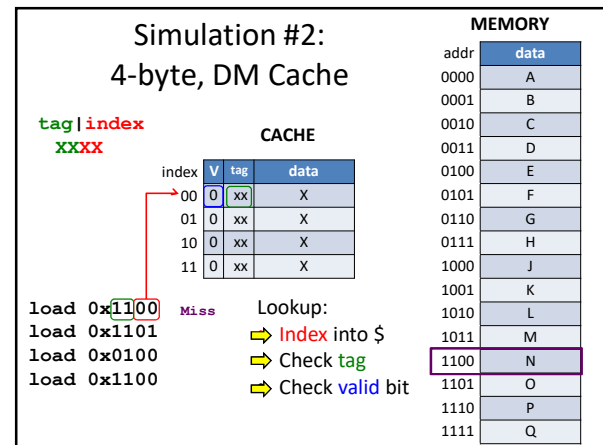
25



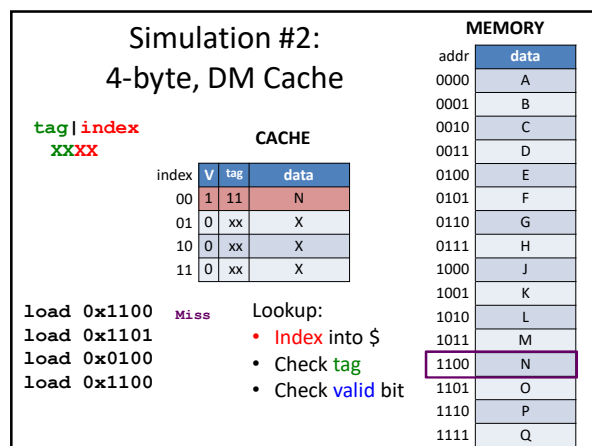
26



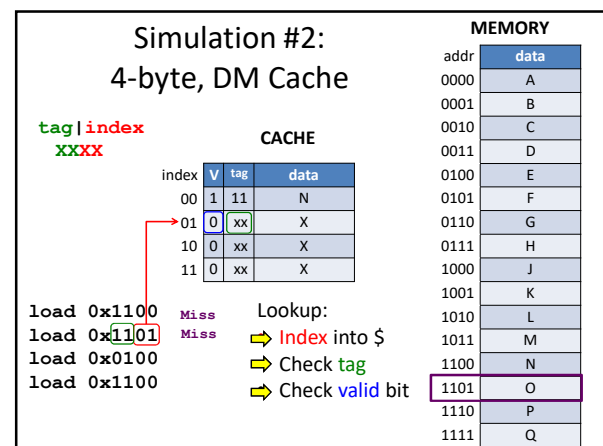
27



28



29



30

Simulation #2:
4-byte, DM Cache

tag|index
XXXX

index	V	tag	data
00	1	11	N
01	1	11	O
10	0	xx	X
11	0	xx	X

load 0x1100 Miss
load 0x1101 Miss
load 0x0100
load 0x1100

Lookup:
• Index into \$
• Check tag
• Check valid bit

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

31

Simulation #2:
4-byte, DM Cache

tag|index
XXXX

index	V	tag	data
00	1	11	N
01	1	11	O
10	0	xx	X
11	0	xx	X

load 0x1100 Miss
load 0x1101 Miss
load 0x0100 Miss
load 0x1100

Lookup:
⇒ Index into \$
⇒ Check tag
⇒ Check valid bit

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

32

Simulation #2:
4-byte, DM Cache

tag|index
XXXX

index	V	tag	data
00	1	01	E
01	1	11	O
10	0	xx	X
11	0	xx	X

load 0x1100 Miss
load 0x1101 Miss
load 0x0100 Miss
load 0x1100

Lookup:
• Index into \$
• Check tag
• Check valid bit

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

33

Simulation #2:
4-byte, DM Cache

tag|index
XXXX

index	V	tag	data
00	1	01	E
01	1	11	O
10	0	xx	X
11	0	xx	X

load 0x1100 Miss
load 0x1101 Miss
load 0x0100 Miss
load 0x1100

Lookup:
⇒ Index into \$
⇒ Check tag
⇒ Check valid bit

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

34

Simulation #2:
4-byte, DM Cache

tag|index
XXXX

index	V	tag	data
00	1	11	N
01	1	11	O
10	0	xx	X
11	0	xx	X

load 0x1100 Miss cold
load 0x1101 Miss cold
load 0x0100 Miss cold
load 0x1100 Miss

Disappointed! ☹️

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

35

Reducing Cold Misses
by Increasing Block Size

Leveraging Spatial Locality

36

Increasing Block Size

CACHE

offset	index	V	tag	data
xxxx	00	0	x	A B
	01	0	x	C D
	10	0	x	E F
	11	0	x	G H

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

- Block Size:** 2 bytes
- Block Offset:** least significant bits indicate where you live in the block
- Which bits are the index? tag?

37

Simulation #3: 8-byte, DM Cache

CACHE

offset	index	V	tag	data
xxxx	00	0	x	X X
	01	0	x	X X
	10	0	x	X X
	11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
 load 0x1101 ➔ Index into \$
 load 0x0100 ➔ Check tag
 load 0x1100 ➔ Check valid bit

38

Simulation #3: 8-byte, DM Cache

CACHE

offset	index	V	tag	data
xxxx	00	0	x	X X
	01	0	x	X X
	10	1	1	N O
	11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
 load 0x1101 • Index into \$
 load 0x0100 • Check tag
 load 0x1100 • Check valid bit

39

Simulation #3: 8-byte, DM Cache

CACHE

offset	index	V	tag	data
xxxx	00	0	x	X X
	01	0	x	X X
	10	1	1	N O
	11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
 load 0x1101 Hit! ➔ Index into \$
 load 0x0100 ➔ Check tag
 load 0x1100 ➔ Check valid bit

40

Simulation #3: 8-byte, DM Cache

CACHE

offset	index	V	tag	data
xxxx	00	0	x	X X
	01	0	x	X X
	10	1	1	N O
	11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
 load 0x1101 Hit! ➔ Index into \$
 load 0x0100 Miss ➔ Check tag
 load 0x1100 ➔ Check valid bit

41

Simulation #3: 8-byte, DM Cache

CACHE

offset	index	V	tag	data
xxxx	00	0	x	X X
	01	0	x	X X
	10	1	0	E F
	11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
 load 0x1101 Hit! • Index into \$
 load 0x0100 Miss • Check tag
 load 0x1100 • Check valid bit

42

Simulation #3:
8-byte, DM Cache

tag | index | offset
xxxx

CACHE

index	V	tag	data
00	0	x	X X
01	0	x	X X
10	1	0	E F
11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
load 0x1101 Hit! → Index into \$
load 0x0100 Miss → Check tag
load 0x1100 Miss → Check valid bit

43

Simulation #3:
8-byte, DM Cache

index | V | tag | data

index	V	tag	data
00	0	x	X X
01	0	x	X X
10	1	0	E F
11	0	x	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss cold 1 hit, 3 misses
load 0x1101 Hit! 3 bytes don't fit in
load 0x0100 Miss cold an 8 byte cache?
load 0x1100 Miss conflict

44

**Removing Conflict Misses
with Fully-Associative Caches**



45

**8 byte, fully-associative
Cache**

xxxx

CACHE

V	tag	data	V	tag	data	V	tag	data	V	tag	data
0	xxx	X X	0	xxx	X X	0	xxx	X X	0	xxx	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

What should the **offset** be?
What should the **index** be?
What should the **tag** be?

46

Simulation #4:
8-byte, FA Cache

tag | offset
xxxx

CACHE

V	tag	data	V	tag	data	V	tag	data	V	tag	data
0	xxx	X X	0	xxx	X X	0	xxx	X X	0	xxx	X X

MEMORY

addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
load 0x1101 Hit! → Index into \$
load 0x0100 → Check tags
load 0x1100 → Check valid bits

47

Simulation #4:
8-byte, FA Cache

tag | offset
xxxx

CACHE

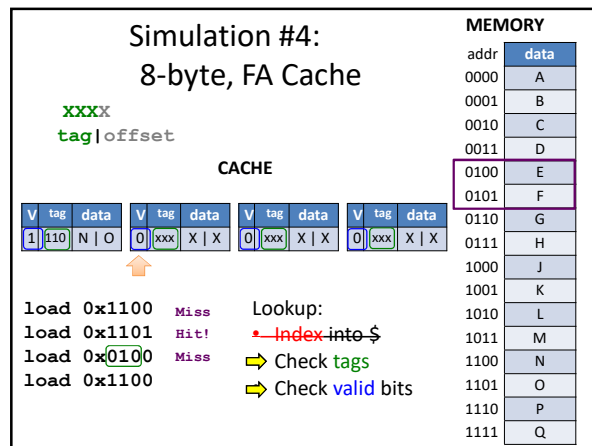
V	tag	data	V	tag	data	V	tag	data	V	tag	data
1	110	N O	0	xxx	X X	0	xxx	X X	0	xxx	X X

MEMORY

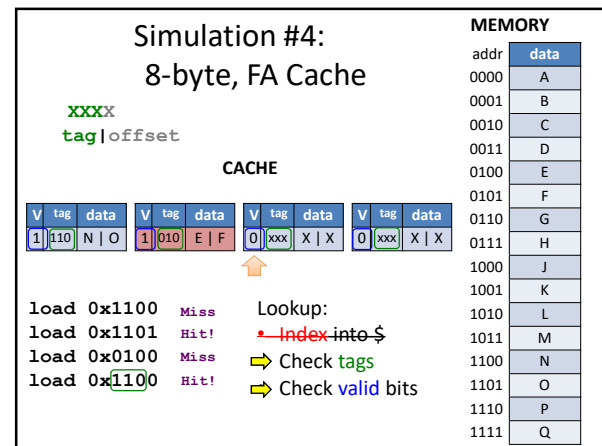
addr	data
0000	A
0001	B
0010	C
0011	D
0100	E
0101	F
0110	G
0111	H
1000	J
1001	K
1010	L
1011	M
1100	N
1101	O
1110	P
1111	Q

load 0x1100 Miss Lookup:
load 0x1101 Hit! → Index into \$
load 0x0100 → Check tags
load 0x1100 → Check valid bits

48



49




50

Pros and Cons of Full Associativity

+ No more conflicts!
+ Excellent utilization!
But...

Parallel Reads
• lots of reading!

Serial Reads
• lots of waiting



$$t_{avg} = t_{hit} + \%miss \times t_{miss}$$

= 4 + 5% x 100 = 6 + 3% x 100
= 9 cycles = 9 cycles

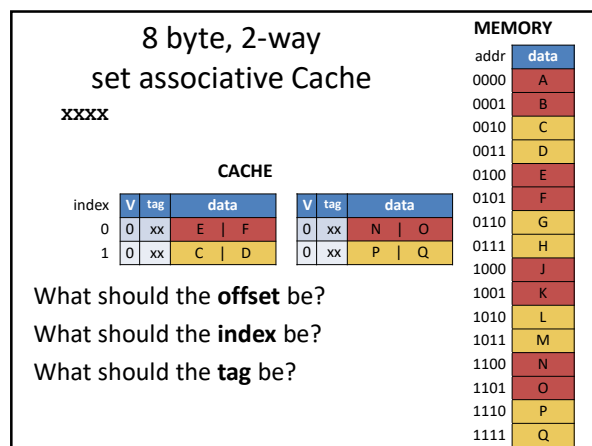
51

**Reducing Conflict Misses
with Set-Associative Caches**

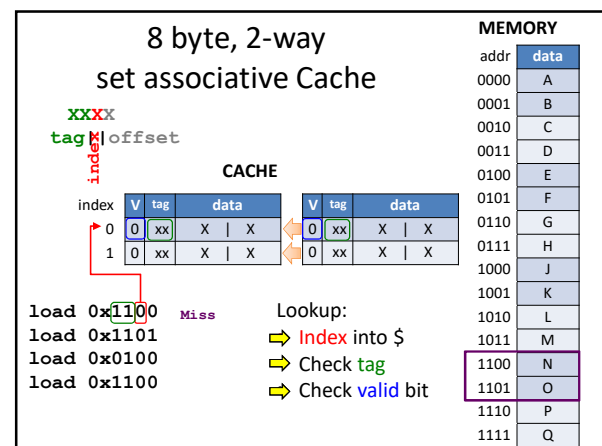
Not too conflict-y. Not too slow.
... Just Right!



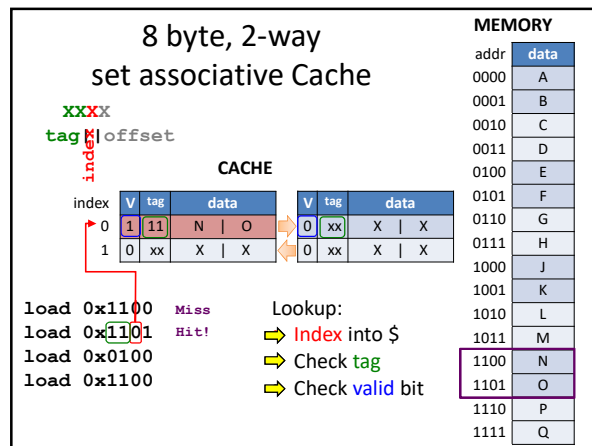
52



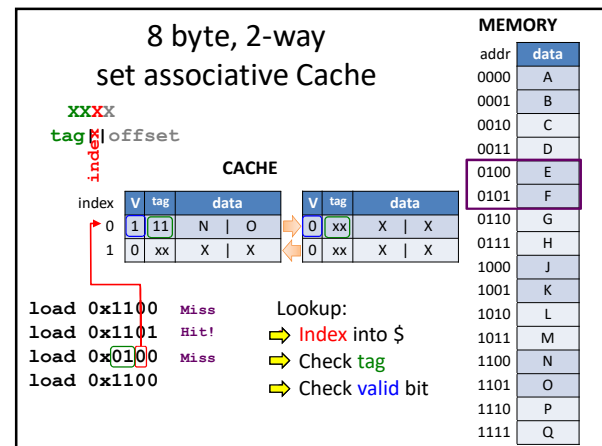
53



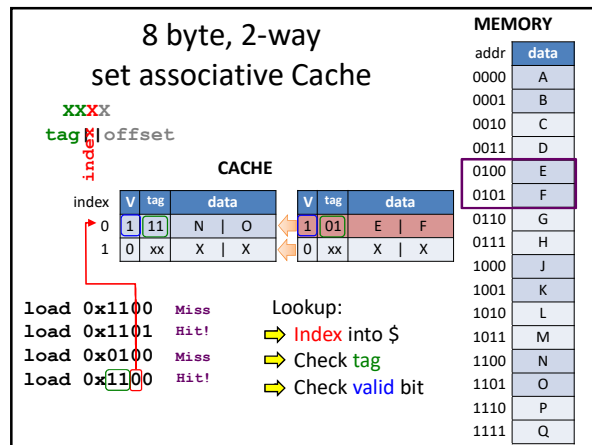
54



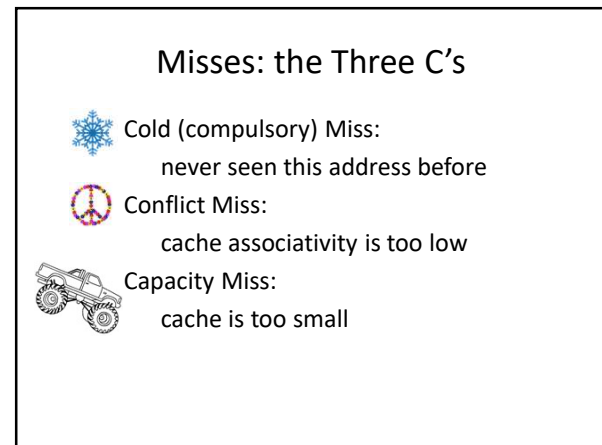
55



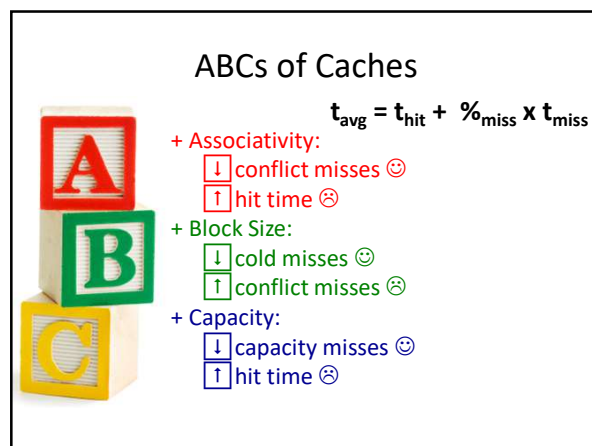
56



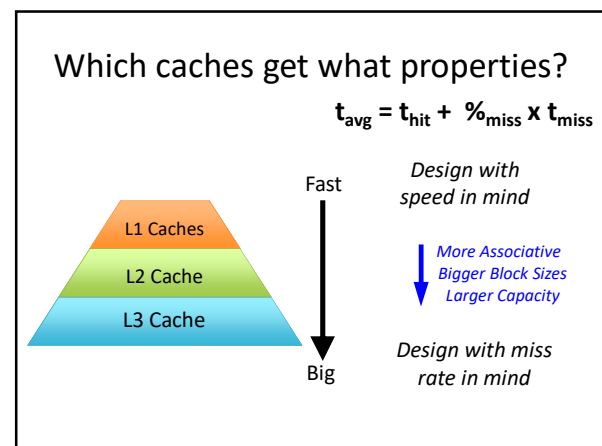
57



58



59



60

Summary so far

- Things we've covered:
 - The Need for Speed
 - Locality to the Rescue!
 - Calculating average memory access time
 - \$ Misses: Cold, Conflict, Capacity
 - \$ Characteristics: Associativity, Block Size, Capacity
- Things we skipped (and are about to cover):
 - Cache Overhead
 - Replacement Policies
 - Writes