

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?

Processor
Main Memory

- too slow
- too far away

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

7

What's the solution?

Caches !

Level 1 Data \$
Level 2 \$
Level 1 Insn \$

What lucky data gets to go here?

Intel Pentium 3, 1999

8

Locality Locality Locality

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

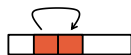
- the same thing again soon

→ Temporal Locality



- something near that thing, soon

→ Spatial Locality



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

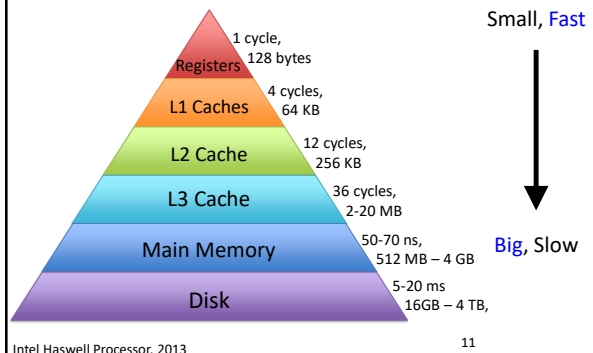
9

9

Your life is full of Locality

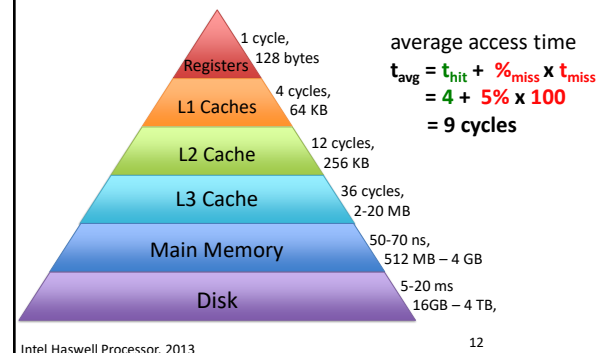
10

The Memory Hierarchy

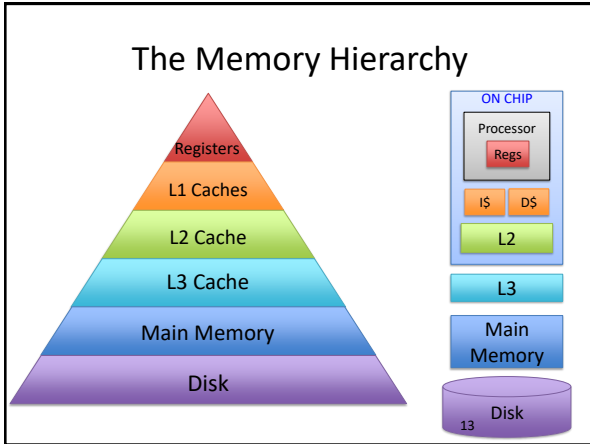


11

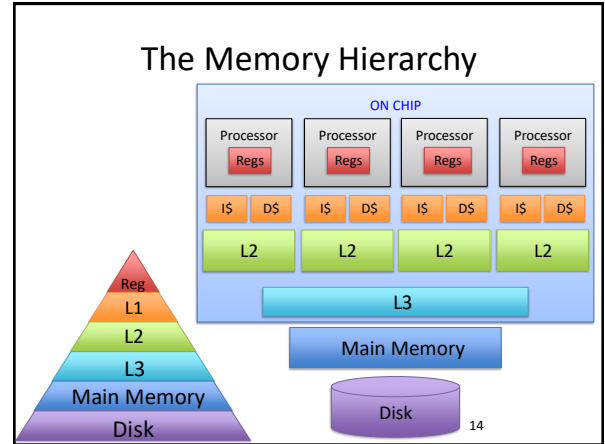
The Memory Hierarchy



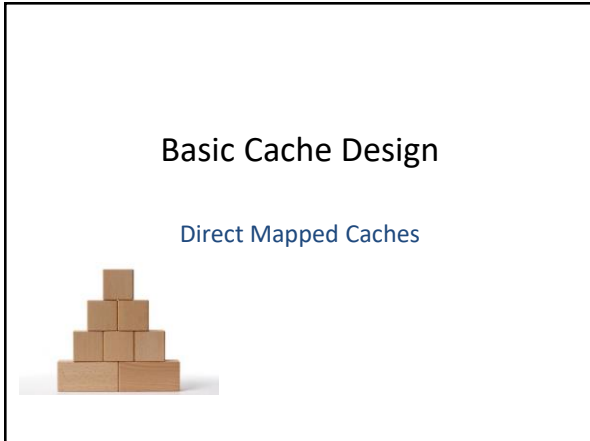
12



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 (2ⁿ → 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

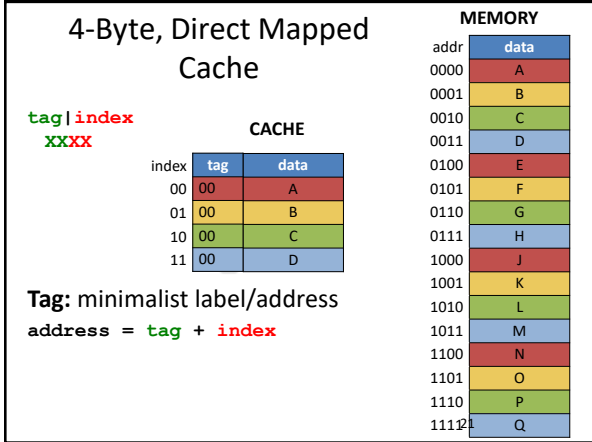
index
XXXX

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

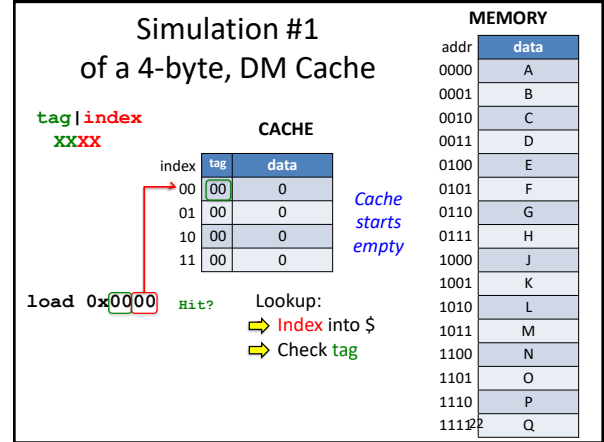
- Supports spatial locality

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

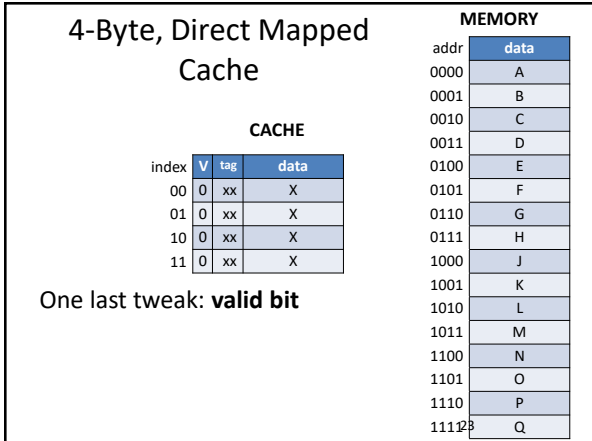
18



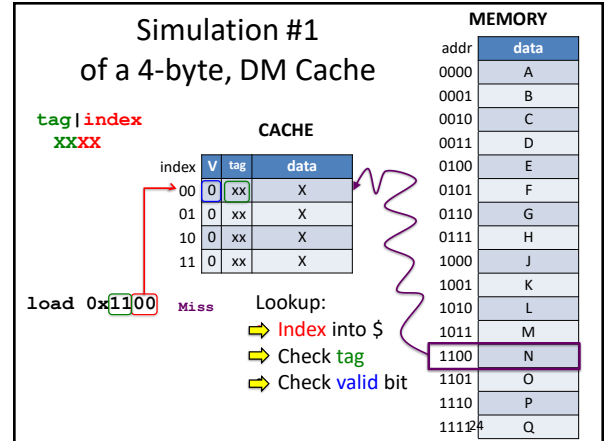
21



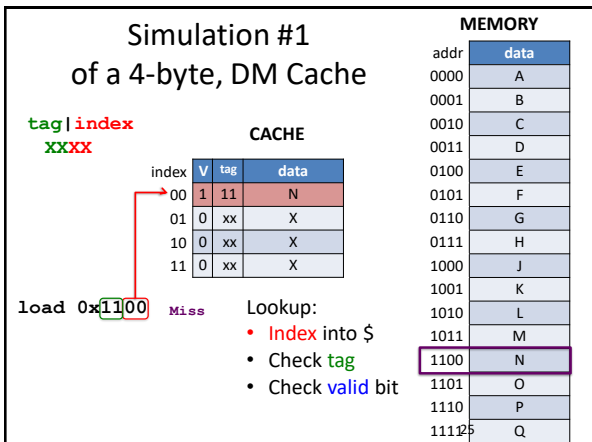
22



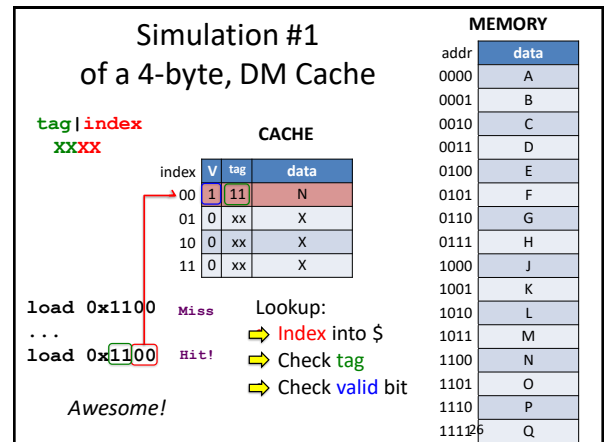
23



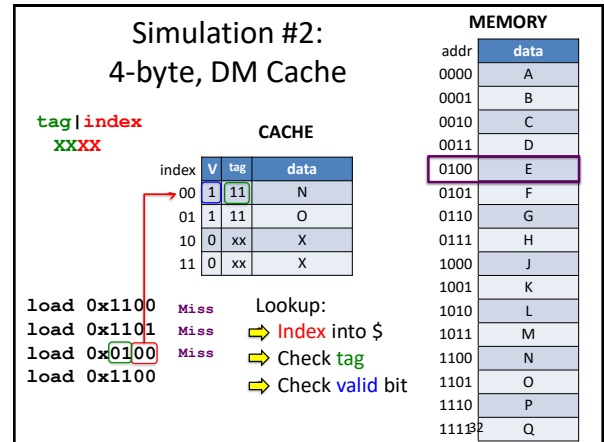
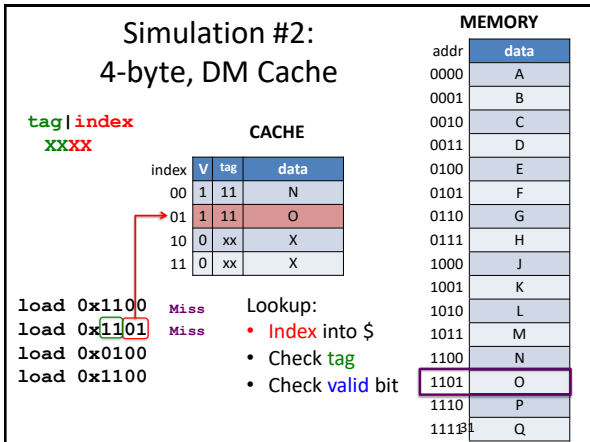
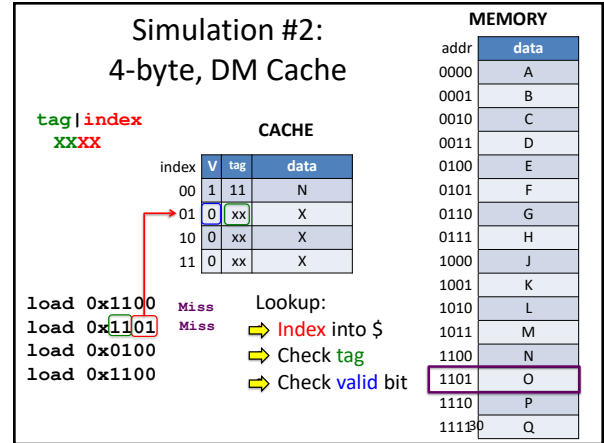
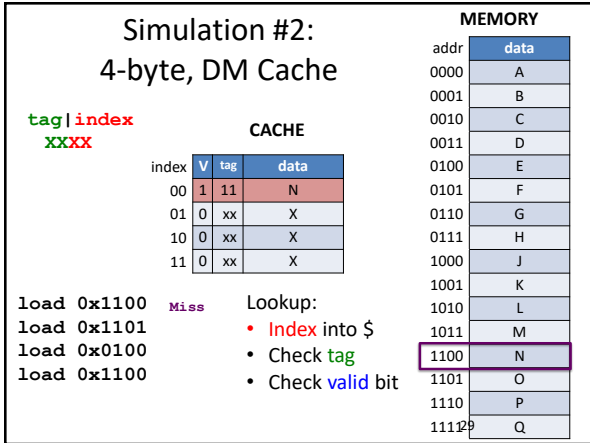
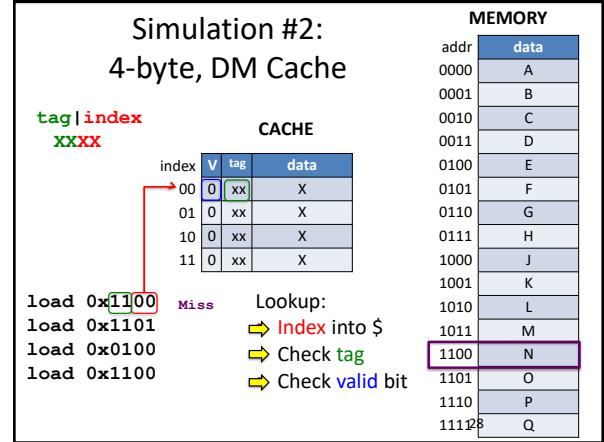
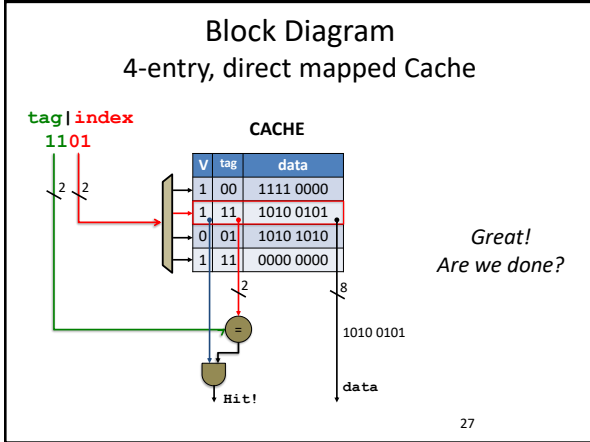
24



25



26



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

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

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

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

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! 😞

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

35

Reducing Cold Misses by Increasing Block Size

Leveraging Spatial Locality



36

Increasing Block Size

offset | index
XXXX

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

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

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

37

Simulation #3: 8-byte, DM Cache

tag | index | offset
XXXX

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

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

Lookup:

- Index into \$
- Check tag
- Check 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

38

Simulation #3: 8-byte, DM Cache

CACHE

index	V	tag	data
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

index	V	tag	data
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

index	V	tag	data
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 Miss • Check valid bit

41

Simulation #3: 8-byte, DM Cache

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

42

Simulation #3: 8-byte, DM Cache

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

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 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

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

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

46

Simulation #4: 8-byte, FA Cache

XXXX
tag|offset

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

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

47

Simulation #4: 8-byte, FA Cache

XXXX
tag|offset

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

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

48

Simulation #4: 8-byte, FA Cache

XXXX
tag|offset

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

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

49

Simulation #4: 8-byte, FA Cache

XXXX
tag|offset

CACHE

V	tag	data	V	tag	data	V	tag	data	V	tag	data
1	110	N O	1	010	E F	0	xxx	X X	0	xxx	X X

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

50

Pros and Cons of Full Associativity

- + No more conflicts!
- + Excellent utilization!
- But...
- Parallel Reads
 - lots of reading!
- Serial Reads
 - lots of waiting



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

$= 4 + 5\% \times 100 = 6 + 3\% \times 100 = 9 \text{ cycles}$

51

Reducing Conflict Misses with Set-Associative Caches

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



52

8 byte, 2-way set associative Cache

XXXX

CACHE		CACHE	
index	V tag data	V tag data	V tag data
0	0 xx E F	0 xx N O	
1	0 xx C D	0 xx P Q	

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

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
51111	Q

53

8 byte, 2-way set associative Cache

XXXX
tag | offset

CACHE		CACHE	
index	V tag data	V tag data	V tag data
0	0 xx X X	0 xx X X	
1	0 xx X X	0 xx X X	

- load 0x1100 Miss
 - load 0x1101 Miss
 - load 0x0100
 - load 0x1100
- Lookup:
- Index into \$
 - Check tag
 - Check 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
51111	Q

54

8 byte, 2-way set associative Cache

XXXX
tag | offset

CACHE		CACHE	
index	V tag data	V tag data	V tag data
0	1 11 N O	0 xx X X	
1	0 xx X X	0 xx X X	

- load 0x1100 Miss
 - load 0x1101 Hit!
 - load 0x0100
 - load 0x1100
- Lookup:
- Index into \$
 - Check tag
 - Check 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
51111	Q

55

8 byte, 2-way set associative Cache

XXXX
tag | offset

CACHE		CACHE	
index	V tag data	V tag data	V tag data
0	1 11 N O	0 xx X X	
1	0 xx X X	0 xx X X	

- load 0x1100 Miss
 - load 0x1101 Hit!
 - load 0x0100 Miss
 - load 0x1100
- Lookup:
- Index into \$
 - Check tag
 - Check 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
51111	Q

56

8 byte, 2-way set associative Cache

xxxX
tag: offset

CACHE

index	V	tag	data	
0	1	11	N	O
1	0	xx	X	X

index: 0, 1

load 0x1100 Miss
load 0x1101 Hit!
load 0x0100 Miss
load 0x1100 Hit!

Lookup:
 Index into \$
 Check tag
 Check 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

57

Misses: the Three C's

- Cold (compulsory) Miss:**
never seen this address before
- Conflict Miss:**
cache associativity is too low
- Capacity Miss:**
cache is too small

58

ABCs of Caches

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

- + Associativity:**
 - ↓ conflict misses ☹️
 - ↑ hit time ☹️
- + Block Size:**
 - ↓ cold misses ☹️
 - ↑ conflict misses ☹️
- + Capacity:**
 - ↓ capacity misses ☹️
 - ↑ hit time ☹️

59

Which caches get what properties?

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

Fast
↓
Big

Design with speed in mind

More Associative
Bigger Block Sizes
Larger Capacity

Design with miss rate in mind

L1 Caches
L2 Cache
L3 Cache

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

61