# caching 1

# last time

certificates and certificate authorities
> chains of certificates
> trust anchors
> the reality of website "verification"

cryptographic hash functions

Diffie-Hellman-style key agreement

TLS handshake authenticate + get symmetric keys for connection

denial of service how game is biased for attacker

(briefly) firewalls

# 2004 CPU



Floating Point Unit

Load/Store

Data Cache

Execution Units

Bus Unit

L2 Cache 1MB

Fetch Scan Align Micro-code

Instruction Cache

Memory Controller

Hyper Transport

DDR Memory Interface

Clock Generator

AMD

Image: approx 2004 AMD press image of Opteron die; approx register location via chip-architect.org (Hans de Vries)

# 2004 CPU



▲ Registers

3

# 2004 CPU



Registers
L1 cache

Floating Point Unit
Load/Store
Data Cache
Hyper Transport
Execution Units
Bus Unit
L2 Cache
1MB
DDR Memory Interface
Fetch Scan Align
Micro-code
Instruction Cache
Memory Controller
Clock Generator

3

# 2004 CPU



Floating Point Unit

Load/Store

Data Cache

Execution Units

Bus Unit

Fetch Scan Align Micro-code

Instruction Cache

L2 Cache 1MB

Memory Controller

Hyper Transport

DDR Memory Interface

Clock Generator

Registers
L1 cache
L2 cache

3

# 2004 CPU



Floating Point Unit
Load/Store
Data Cache
Execution Units
Bus Unit
Fetch Scan Align Micro-code
Instruction Cache
Memory Controller
Hyper Transport
DDR Memory Interface
L2 Cache 1MB
Clock Generator

Registers
L1 cache
L2 cache

3

# 2004 CPU



Image: approx 2004 AMD press image of Opteron die;
approx register location via chip-architect.org (Hans de Vries)

3

# 2004 CPU



Image: approx 2004 AMD press image of Opteron die; approx register location via chip-architect.org (Hans de Vries)

# the place of cache

# memory hierarchy goals

performance of the fastest (smallest) memory
    hide 100x latency difference? 99+% hit (= value found in cache) rate

capacity of the largest (slowest) memory

# memory hierarchy assumptions

temporal locality
"if a value is accessed now, it will be accessed again soon"
    caches should keep recently accessed values

spatial locality
"if a value is accessed now, adjacent values will be accessed soon"
    caches should store adjacent values at the same time

natural properties of programs — think about loops

## locality examples

```
double computeMean(int length, double *values) {
    double total = 0.0;
    for (int i = 0; i < length; ++i) {
        total += values[i];
    }
    return total / length;
}
```

temporal locality: machine code of the loop

spatial locality: machine code of most consecutive instructions

temporal locality: total, i, length accessed repeatedly

spatial locality: values[i+1] accessed after values[i]

# building a (direct-mapped) cache

Cache

| value |
|-------|
| 00 00 |
| 00 00 |
| 00 00 |
| 00 00 |

cache block: 2 bytes

Memory

| addresses | bytes |
|-----------|-------|
| 00000–00001 | 00 11 |
| 00010–00011 | 22 33 |
| 00100–00101 | 55 55 |
| 00110–00111 | 66 77 |
| 01000–01001 | 88 99 |
| 01010–01011 | AA BB |
| 01100–01101 | CC DD |
| 01110–01111 | EE FF |
| 10000–10001 | F0 F1 |
| ... | ... |

# building a (direct-mapped) cache

read byte at 01011?

Cache

| value |
|-------|
| 00 00 |
| 00 00 |
| 00 00 |
| 00 00 |

cache block: 2 bytes

Memory

| addresses | bytes |
|-----------|-------|
| 00000−00001 | 00 11 |
| 00010−00011 | 22 33 |
| 00100−00101 | 55 55 |
| 00110−00111 | 66 77 |
| 01000−01001 | 88 99 |
| 01010−01011 | AA BB |
| 01100−01101 | CC DD |
| 01110−01111 | EE FF |
| 10000−10001 | F0 F1 |
| ... | ... |

# building a (direct-mapped) cache

read byte at 01011?

exactly one place for each address
spread out what can go in a block



Cache

Memory

| index | value | addresses | bytes |
|-------|-------|-----------|-------|
| 00 | 00 00 | 00000-00001 | 00 11 |
| 01 | 00 00 | 00010-00011 | 22 33 |
| 10 | 00 00 | 00100-00101 | 55 55 |
| 11 | 00 00 | 00110-00111 | 66 77 |
| | | 01000-01001 | 88 99 |
| | | 01010-01011 | AA BB |
| | | 01100-01101 | CC DD |
| | | 01110-01111 | EE FF |
| | | 10000-10001 | F0 F1 |
| | | ... | ... |

cache block: 2 bytes
direct-mapped

# building a (direct-mapped) cache

read byte at 01011?

exactly one place for each address
spread out what can go in a block



Cache

| index | value |
|-------|-------|
| 00 | 00 00 |
| 01 | 00 00 |
| 10 | 00 00 |
| 11 | 00 00 |

cache block: 2 bytes
direct-mapped

Memory

| addresses | bytes |
|-----------|-------|
| 00000–00001 | 00 11 |
| 00010–00011 | 22 33 |
| 00100–00101 | 55 55 |
| 00110–00111 | 66 77 |
| 01000–01001 | 88 99 |
| 01010–01011 | AA BB |
| 01100–01101 | CC DD |
| 01110–01111 | EE FF |
| 10000–10001 | F0 F1 |
| ... | ... |

# building a (direct-mapped) cache

read byte at 01011?

exactly one place for each address
spread out what can go in a block



Cache

Memory

| index | value | addresses | bytes |
|-------|-------|-----------|-------|
| 00 | 00 00 | 00000–00001 | 00 11 |
| 01 | 00 00 | 00010–00011 | 22 33 |
| 10 | 00 00 | 00100–00101 | 55 55 |
| 11 | 00 00 | 00110–00111 | 66 77 |
| | | 01000–01001 | 88 99 |
| | | 01010–01011 | AA BB |
| | | 01100–01101 | CC DD |
| | | 01110–01111 | EE FF |
| | | 10000–10001 | F0 F1 |
| | | … | … |

cache block: 2 bytes
direct-mapped

# building a (direct-mapped) cache

read byte at 01011?



Cache

| index | valid | value |
|-------|-------|-------|
| 00 | 0 | 00 00 |
| 01 | 0 | 00 00 |
| 10 | 0 | 00 00 |
| 11 | 0 | 00 00 |

is this even a value?

need extra bit to know

cache block: 2 bytes
direct-mapped

Memory

| addresses | bytes |
|-----------|-------|
| | _1 |
| 00010–00011 | 22 33 |
| 0–00101 | 55 55 |
| 00110–00111 | 66 77 |
| 01000–01001 | 88 99 |
| 01010–01011 | AA BB |
| 01100–01101 | CC DD |
| 01110–01111 | EE FF |
| 10000–10001 | F0 F1 |
| ... | ... |

8

# building a (direct-mapped) cache

read byte at 01011?
invalid, fetch

## Cache

| index | valid | | value |
|-------|-------|--|-------|
| 00 | 0 | | 00 00 |
| 01 | 1 | | AA BB |
| 10 | 0 | | 00 00 |
| 11 | 0 | | 00 00 |

cache block: 2 bytes
direct-mapped

## Memory

| addresses | bytes |
|-----------|-------|
| 00000−00001 | 00 11 |
| 00010−00011 | 22 33 |
| 00100−00101 | 55 55 |
| 00110−00111 | 66 77 |
| 01000−01001 | 88 99 |
| 01010−01011 | AA BB |
| 01100−01101 | CC DD |
| 01110−01111 | EE FF |
| 10000−10001 | F0 F1 |
| ... | ... |

8

# building a (direct-mapped) cache

read byte at 01011?
invalid, fetch

Cache

Memory

value from `01010` or `00010`?

| index | valid | tag | value |
|-------|-------|-----|-------|
| 00 | 0 | 00 | 00 00 |
| 01 | 1 | 01 | AA BB |
| 10 | 0 | 00 | 00 00 |
| 11 | 0 | | |

need tag to know

cache block: 2 bytes
direct-mapped

| addresses | bytes |
|-----------|-------|
| 00000–00001 | 00 11 |
| 00010–00011 | 22 33 |
| 00100–00101 | 55 55 |
| 00110–00111 | 66 77 |
| 01000–01001 | 88 99 |
| 01010–01011 | AA BB |
| 01100–01101 | CC DD |
| 01110–01111 | EE FF |
| 10000–10001 | F0 F1 |
| … | … |

8

# building a (direct-mapped) cache

read byte at 01011?
invalid, fetch

Cache

| index | valid | tag | value |
|-------|-------|-----|-------|
| 00 | 0 | 00 | 00 00 |
| 01 | 1 | 01 | AA BB |
| 10 | 0 | 00 | 00 00 |
| 11 | 0 | 00 | 00 00 |

cache block: 2 bytes
direct-mapped

Memory

| addresses | bytes |
|-----------|-------|
| 00000–00001 | 00 11 |
| 00010–00011 | 22 33 |
| 00100–00101 | 55 55 |
| 00110–00111 | 66 77 |
| 01000–01001 | 88 99 |
| 01010–01011 | AA BB |
| 01100–01101 | CC DD |
| 01110–01111 | EE FF |
| 10000–10001 | F0 F1 |
| … | … |

# cache operation (read)

0b 11 100 10

| valid | tag | data |
|-------|-----|------|
| 1 | 10 | 00 11 22 33 |
| | | |
| | | |
| | | |
| 1 | 11 | B4 B5 B6 B7 |
| | | |
| | | |
| | | |

index

# cache operation (read)

0b11100010

# cache operation (read)



0b 11 100 10 ——— offset ———

valid **tag**   **data**

| valid | tag | data |
|---|---|---|
| 1 | 10 | 00 11 22 33 |
| | | |
| | | |
| **1** | **11** | **B4 B5 B6 B7** |
| | | |
| | | |

index

tag

= 

AND → is hit? (1)

→ data (B6)

# terminology

row = set
  preview: change how much is in a row

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|---|---|---|---|
| 2 byte blocks, 4 sets | | | |
| 2 byte blocks, 8 sets | | | |
| 4 byte blocks, 2 sets | | | |

### 2 byte blocks, 4 sets

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | -- | -- -- |
| 11 | 1 | 001 | EE FF |

### 4 byte blocks, 2 sets

| index | valid | tag | value |
|---|---|---|---|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

### 2 byte blocks, 8 sets

| index | valid | tag | value |
|---|---|---|---|
| 000 | 1 | 00 | 00 11 |
| 001 | 1 | 01 | F1 F2 |
| 010 | 0 | -- | -- -- |
| 011 | 0 | -- | -- -- |
| 100 | 0 | -- | -- -- |
| 101 | 1 | 00 | AA BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|---|---|---|---|
| 2 byte blocks, 4 sets | | | 1 |
| 2 byte blocks, 8 sets | | | 1 |
| 4 byte blocks, 2 sets | | | |

2 byte blocks, 4 sets

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | -- | -- -- |
| 11 | 1 | 001 | EE FF |

$2 = 2^1$ bytes in block
1 bit to say which byte

2 byte blocks, 8 sets

| index | valid | tag | value |
|---|---|---|---|
| | | | 00 11 |
| | | | F1 F2 |
| | | | -- -- |
| 011 | 0 | -- | -- -- |
| 100 | 0 | -- | -- -- |
| 101 | 1 | 00 | AA BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

4 byte blocks, 2 sets

| index | valid | tag | value |
|---|---|---|---|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|---|---|---|---|
| 2 byte blocks, 4 sets | | | 1 |
| 2 byte blocks, 8 sets | | | 1 |
| 4 byte blocks, 2 sets | | | 11 |

2 byte blocks, 4 sets

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | | |
| 11 | 1 | | |

2 byte blocks, 8 sets

| index | valid | tag | value |
|---|---|---|---|
| 000 | 1 | 00 | 00 11 |
| 001 | 1 | 01 | F1 F2 |
| | 0 | -- | -- -- |
| | 0 | -- | -- -- |
| | 0 | -- | -- -- |
| 101 | 1 | 00 | AA BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

4 byte

$4 = 2^2$ bytes in block
2 bits to say which byte

| index | valid | tag | value |
|---|---|---|---|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|---|---|---|---|
| 2 byte blocks, 4 sets | 11 | | 1 |
| 2 byte blocks, 8 sets | | | 1 |
| 4 byte blocks, 2 sets | 1 | | 11 |

2 byte blocks, 4 sets

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | -- | -- -- |
| 11 | 1 | 001 | EE FF |

$2^2 = 4$ sets
2 bits to index set

2 byte blocks, 8 sets

| index | valid | tag | value |
|---|---|---|---|
| 000 | 1 | 00 | 00 11 |
| 001 | | | F1 F2 |
| 010 | | | -- -- |
| 011 | | | -- -- |
| 100 | 0 | -- | -- -- |
| 101 | 1 | 00 | AA BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

4 byte blocks, 2 sets

| index | valid | tag | value |
|---|---|---|---|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

11

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|---|---|---|---|
| 2 byte blocks, 4 sets | 11 | 1 | |
| 2 byte blocks, 8 sets | 111 | 1 | |
| 4 byte blocks, 2 sets | 1 | 11 | |

2 byte blocks, 4 sets

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | -- | -- -- |
| 11 | 1 | | |

4

| index | valid | tag | value |
|---|---|---|---|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

$2^3 = 8$ sets
3 bits to index set

2 byte blocks, 8 sets

| index | valid | tag | value |
|---|---|---|---|
| 000 | 1 | 00 | 00 11 |
| 001 | 1 | 01 | F1 F2 |
| 010 | 0 | -- | -- -- |
| 011 | 0 | -- | -- -- |
| 100 | 0 | -- | -- -- |
| 101 | 1 | 00 | AA BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|-------|-----|-------|--------|
| 2 byte blocks, 4 sets | 11 | 1 | |
| 2 byte blocks, 8 sets | 111 | 1 | |
| 4 byte blocks, 2 sets | 1 | 11 | |

2 byte blocks, 4 sets

| index | valid | tag | value |
|-------|-------|-----|-------|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | -- | -- -- |
| 11 | 1 | 001 | EE FF |

4 byte blocks, 2 sets

| index | valid | tag | value |
|-------|-------|-----|-------|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

2 byte blocks, 8 sets

| index | valid | tag | value |
|-------|-------|-----|-------|
| 000 | 1 | 00 | 00 11 |
| 001 | 1 | 01 | F1 F2 |
| 010 | 0 | -- | -- -- |
| 011 | | | -- -- |
| 10 | | | -- -- |
| 10 | | | BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

$2^1 = 2$ sets

1 bit to index set

# Tag-Index-Offset (TIO)

address `001111` (stores value `0xFF`)

| cache | tag | index | offset |
|---|---|---|---|
| 2 byte blocks, 4 sets | 001 | 11 | 1 |
| 2 byte blocks, 8 sets | 00 | 111 | 1 |
| 4 byte blocks, 2 sets | 001 | 1 | 11 |

tag — whatever is left over

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 000 | 00 11 |
| 01 | 1 | 001 | AA BB |
| 10 | 0 | -- | -- -- |
| 11 | 1 | 001 | EE FF |

4 byte blocks, 2 sets

| index | valid | tag | value |
|---|---|---|---|
| 0 | 1 | 000 | 00 11 22 33 |
| 1 | 1 | 001 | CC DD EE FF |

2 byte blocks, 8 sets

| index | valid | tag | value |
|---|---|---|---|
| 000 | 1 | 00 | 00 11 |
| 001 | 1 | 01 | F1 F2 |
| 010 | 0 | -- | -- -- |
| 011 | 0 | -- | -- -- |
| 100 | 0 | -- | -- -- |
| 101 | 1 | 00 | AA BB |
| 110 | 0 | -- | -- -- |
| 111 | 1 | 00 | EE FF |

11

# TIO: exercise

64-byte blocks, 128 set cache

stores $64 \times 128 = 8192$ bytes (of data)

if addresses 32-bits, then how many tag/index/offset bits?

which bytes are stored in the same block as byte from 0x1037?

    A. byte from 0x1011
    B. byte from 0x1021
    C. byte from 0x1035
    D. byte from 0x1041

# Tag-Index-Offset formulas (direct-mapped)

(formulas derivable from prior slides)

$S = 2^s$          number of sets

$s$               (set) index bits

$B = 2^b$          block size

$b$               (block) offset bits

$m$             memory addreses bits

$t = m - (s + b)$    tag bits

$C = B \times S$       cache size (if direct-mapped)

# Tag-Index-Offset formulas (direct-mapped)

(formulas derivable from prior slides)

$S = 2^s$        number of sets

$s$        (set) index bits

$B = 2^b$        block size

$b$        (block) offset bits

$m$        memory addresses bits

$t = m - (s + b)$        tag bits

$C = B \times S$        cache size (if direct-mapped)

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | |
| 00000001 (01) | |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

| index | valid | tag | value |
|---|---|---|---|
| 00 | 0 | | |
| 01 | 0 | | |
| 10 | 0 | | |
| 11 | 0 | | |

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex)   | result |
|-----------------|--------|
| 00000000 (00)   |        |
| 00000001 (01)   |        |
| 01100011 (63)   |        |
| 01100001 (61)   |        |
| 01100010 (62)   |        |
| 00000000 (00)   |        |
| 01100100 (64)   |        |

| index | valid | tag | value |
|-------|-------|-----|-------|
| 00    | 0     |     |       |
| 01    | 0     |     |       |
| 10    | 0     |     |       |
| 11    | 0     |     |       |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

14

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| `00000000 (00)` | |
| `00000001 (01)` | |
| `01100011 (63)` | |
| `01100001 (61)` | |
| `01100010 (62)` | |
| `00000000 (00)` | |
| `01100100 (64)` | |

tag  index offset

| index | valid | tag | value |
|---|---|---|---|
| `00` | 0 | | |
| `01` | 0 | | |
| `10` | 0 | | |
| `11` | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| `00000000` (00) | miss |
| `00000001` (01) | |
| `01100011` (63) | |
| `01100001` (61) | |
| `01100010` (62) | |
| `00000000` (00) | |
| `01100100` (64) | |

tag  index  offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 00000 | `mem[0x00]` `mem[0x01]` |
| 01 | 0 | | |
| 10 | 0 | | |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag   index offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 00000 | mem[0x00] mem[0x01] |
| 01 | 0 | | |
| 10 | 0 | | |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| `00000000` (00) | miss |
| `00000001` (01) | hit |
| `01100011` (63) | miss |
| `01100001` (61) | |
| `01100010` (62) | |
| `00000000` (00) | |
| `01100100` (64) | |

tag  index  offset

| index | valid | tag | value |
|---|---|---|---|
| `00` | 1 | `00000` | `mem[0x00]`<br>`mem[0x01]` |
| `01` | 1 | `01100` | `mem[0x62]`<br>`mem[0x63]` |
| `10` | 0 | | |
| `11` | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

14

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag  index offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 01100 | mem[0x60] mem[0x61] |
| 01 | 1 | 01100 | mem[0x62] mem[0x63] |
| 10 | 0 | | |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| `00000000` (00) | miss |
| `00000001` (01) | hit |
| `01100011` (63) | miss |
| `01100001` (61) | miss |
| `01100010` (62) | hit |
| `00000000` (00) | |
| `01100100` (64) | |

tag  index  offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 01100 | mem[0x60] mem[0x61] |
| 01 | 1 | 01100 | mem[0x62] mem[0x63] |
| 10 | 0 | | |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

14

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | miss |
| 01100100 (64) | |

tag  index  offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 00000 | mem[0x00] mem[0x01] |
| 01 | 1 | 01100 | mem[0x62] mem[0x63] |
| 10 | 0 | | |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

14

# example access pattern (1)

$2$ byte blocks, $4$ sets

| address (hex) | result |
|---|---|
| `00000000` (00) | miss |
| `00000001` (01) | hit |
| `01100011` (63) | miss |
| `01100001` (61) | miss |
| `01100010` (62) | hit |
| `00000000` (00) | miss |
| `01100100` (64) | miss |

tag  index offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 00000 | mem[0x00] mem[0x01] |
| 01 | 1 | 01100 | mem[0x62] mem[0x63] |
| 10 | 1 | 01100 | mem[0x64] mem[0x65] |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | miss |
| 00000001 (01) | hit |
| 01100011 (63) | miss |
| 01100001 (61) | miss |
| 01100010 (62) | hit |
| 00000000 (00) | miss |
| 01100100 (64) | miss |

tag  index offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 00000 | mem[0x00] mem[0x01] |
| 01 | 1 | 01100 | mem[0x62] mem[0x63] |
| 10 | 1 | 01100 | mem[0x64] mem[0x65] |
| 11 | 0 | | |

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

14

# example access pattern (1)

2 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| `00000000` (00) | miss |
| `00000001` (01) | hit |
| `01100011` (63) | miss |
| `01100001` (61) | miss |
| `01100010` (62) | hit |
| `00000000` (00) | miss |
| `01100100` (64) | miss |

tag  index  offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | 1 | 00000 | `mem[0x00]` `mem[0x01]` |
| 01 | 1 | 01100 | `mem[0x62]` `mem[0x63]` |
| 10 | 1 | 01100 | `mem[0x64]` `mem[0x65]` |
| 11 | 0 | | |

miss caused by conflict

$B = 2 = 2^b$ byte block size
$b = 1$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 5$ tag bits

14

# exercise

4 byte blocks, 4 sets

| address (hex) | result |
|---------------|--------|
| 00000000 (00) | |
| 00000001 (01) | |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

| index | valid | tag | value |
|-------|-------|-----|-------|
| 00 | | | |
| 01 | | | |
| 10 | | | |
| 11 | | | |

# exercise

4 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | |
| 00000001 (01) | |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

| index | valid | tag | value |
|---|---|---|---|
| 00 | | | |
| 01 | | | |
| 10 | | | |
| 11 | | | |

how is the 8-bit address 61 (01100001) split up into tag/index/offset?

$b$ block offset bits;
$B = 2^b$ byte block size;
$s$ set index bits; $S = 2^s$ sets ;
$t = m - (s + b)$ tag bits (leftov

# exercise

4 byte blocks, 4 sets

| address (hex) | result |
|---------------|--------|
| 00000000 (00) |        |
| 00000001 (01) |        |
| 01100011 (63) |        |
| 01100001 (61) |        |
| 01100010 (62) |        |
| 00000000 (00) |        |
| 01100100 (64) |        |

| index | valid | tag | value |
|-------|-------|-----|-------|
| 00    |       |     |       |
| 01    |       |     |       |
| 10    |       |     |       |
| 11    |       |     |       |

$B = 4 = 2^b$ byte block size
$b = 2$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

$m = 8$ bit addresses
$t = m - (s + b) = 4$ tag bits

# exercise

4 byte blocks, 4 sets

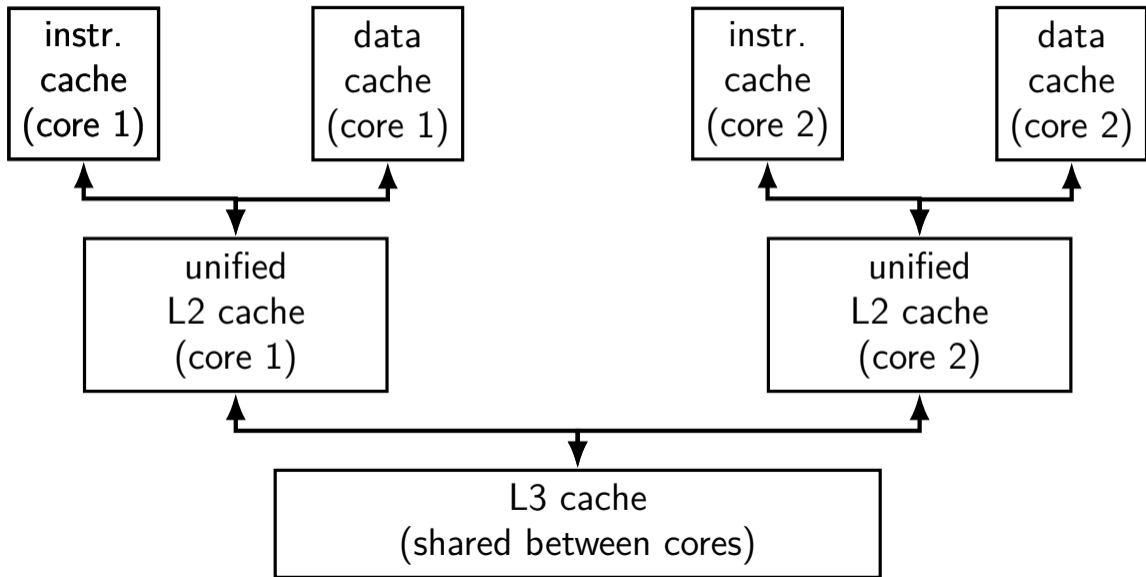| address (hex) | result |
|---------------|--------|
| `00000000` (00) | |
| `00000001` (01) | |
| `01100011` (63) | |
| `01100001` (61) | |
| `01100010` (62) | |
| `00000000` (00) | |
| `01100100` (64) | |

tag index offset

$B = 4 = 2^b$ byte block size
$b = 2$ (block) offset bits
$S = 4 = 2^s$ sets
$s = 2$ (set) index bits

| index | valid | tag | value |
|-------|-------|-----|-------|
| `00` | | | |
| `01` | | | |
| `10` | | | |
| `11` | | | |

$m = 8$ bit addresses
$t = m - (s + b) = 4$ tag bits

# exercise

4 byte blocks, 4 sets

| address (hex) | result |
|---|---|
| 00000000 (00) | |
| 00000001 (01) | |
| 01100011 (63) | |
| 01100001 (61) | |
| 01100010 (62) | |
| 00000000 (00) | |
| 01100100 (64) | |

tag index offset

| index | valid | tag | value |
|---|---|---|---|
| 00 | | | |
| 01 | | | |
| 10 | | | |
| 11 | | | |

exercise: which accesses are hits?

# split caches; multiple cores

# hierarchy and instruction/data caches

typically separate data and instruction caches for L1

(almost) never going to read instructions as data or vice-versa

avoids instructions evicting data and vice-versa

can optimize instruction cache for different access pattern

easier to build fast caches: that handles less accesses at a time

# cache accesses and C code (1)

```
int scaleFactor;

int scaleByFactor(int value) {
    return value * scaleFactor;
}
```

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

exericse: what data cache accesses does this function do?

# cache accesses and C code (1)

```
int scaleFactor;

int scaleByFactor(int value) {
    return value * scaleFactor;
}
```

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

exericse: what data cache accesses does this function do?

    4-byte read of scaleFactor
    8-byte read of return address

# possible scaleFactor use

```
for (int i = 0; i < size; ++i) {
    array[i] = scaleByFactor(array[i]);
}
```

# misses and code (2)

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

suppose each time this is called in the loop:

    return address located at address 0x7fffffffe43b8

    scaleFactor located at address 0x6bc3a0

with direct-mapped 32KB cache w/64 B blocks, what is their:

|        | return address | scaleFactor |
|--------|----------------|-------------|
| tag    |                |             |
| index  |                |             |
| offset |                |             |

# misses and code (2)

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

suppose each time this is called in the loop:
    return address located at address 0x7fffffe43b8
    scaleFactor located at address 0x6bc3a0

with direct-mapped 32KB cache w/64 B blocks, what is their:

|        | return address | scaleFactor |
|--------|----------------|-------------|
| tag    | 0xfffffffc     | 0xd7        |
| index  | 0x10e          | 0x10e       |
| offset | 0x38           | 0x20        |

# misses and code (2)

```
scaleByFactor:
    movl scaleFactor, %eax
    imull %edi, %eax
    ret
```

suppose each time this is called in the loop:
    return address located at address 0x7fffffffe43b8
    scaleFactor located at address 0x6bc3a0

with direct-mapped 32KB cache w/64 B blocks, what is their:

|        | return address | scaleFactor |
|--------|----------------|-------------|
| tag    | 0xffffffffc    | 0xd7        |
| index  | 0x10e          | 0x10e       |
| offset | 0x38           | 0x20        |

# conflict miss coincidences?

obviously I set that up to have the same index
> have to use exactly the right amount of stack space…

but one of the reasons we'll want something better than
direct-mapped cache

# C and cache misses (warmup 1)

```
int array[4];
...
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
odd_sum += array[1];
even_sum += array[2];
odd_sum += array[3];
```

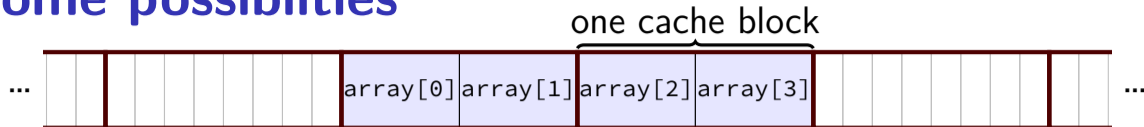Assume everything but array is kept in registers (and the compiler does not do anything funny).

How many data cache misses on a 1-set direct-mapped cache with 8B blocks?

# some possiblities



| … | | | | | | | | | array[0] | array[1] | array[2] | array[3] | | | | | | | | | … |

Q1: how do cache blocks correspond to array elements?
not enough information provided!

# some possiblities

… | | | | | | |array[0]|array[1]|array[2]|array[3]| | | | | | | | …

if array[0] starts at beginning of a cache block…
array split across two cache blocks

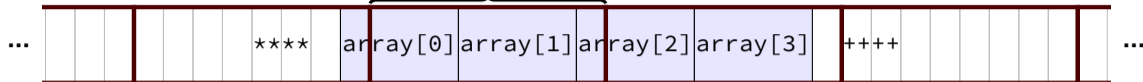| memory access | cache contents afterwards |
|---|---|
| — | (empty) |
| read array[0] (miss) | {array[0], array[1]} |
| read array[1] (hit) | {array[0], array[1]} |
| read array[2] (miss) | {array[2], array[3]} |
| read array[3] (hit) | {array[2], array[3]} |

# some possiblities

one cache block



if array[0] starts right in the middle of a cache block
array split across three cache blocks

| memory access | cache contents afterwards |
|---|---|
| — | (empty) |
| read array[0] (miss) | {****, array[0]} |
| read array[1] (miss) | {array[1], array[2]} |
| read array[2] (hit) | {array[1], array[2]} |
| read array[3] (miss) | {array[3], ++++} |

# some possiblities

... | | | | | | | **** | ar|ray[0]|array[1]|ar|ray[2]|array[3] | ++++ | | | | | | | ...

if array[0] starts at an odd place in a cache block,
need to read two cache blocks to get most array elements

| memory access | cache contents afterwards |
|---|---|
| — | (empty) |
| read `array[0]` byte 0 (miss) | { ****, array[0] byte 0 } |
| read `array[0]` byte 1-3 (miss) | { array[0] byte 1-3, array[2], array[3] byte 0 } |
| read `array[1]` (hit) | { array[0] byte 1-3, array[2], array[3] byte 0 } |
| read `array[2]` byte 0 (hit) | { array[0] byte 1-3, array[2], array[3] byte 0 } |
| read `array[2]` byte 1-3 (miss) | {part of array[2], array[3], ++++} |
| read `array[3]` (hit) | {part of array[2], array[3], ++++} |

# aside: alignment

compilers and malloc/new implementations usually try align values

align = make address be multiple of something

most important reason: don't cross cache block boundaries

# C and cache misses (warmup 2)

```
int array[4];
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
even_sum += array[2];
odd_sum += array[1];
odd_sum += array[3];
```

Assume everything but `array` is kept in registers (and the compiler does not do anything funny).

Assume array[0] at beginning of cache block.

How many data cache misses on a 1-set direct-mapped cache with 8B blocks?
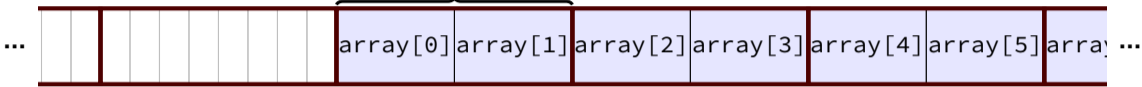
# C and cache misses (warmup 3)

```
int array[8];
...
int even_sum = 0, odd_sum = 0;
even_sum += array[0];
odd_sum += array[1];
even_sum += array[2];
odd_sum += array[3];
even_sum += array[4];
odd_sum += array[5];
even_sum += array[6];
odd_sum += array[7];
```

Assume everything but array is kept in registers (and the compiler does not do anything funny), and array[0] at beginning of cache block.
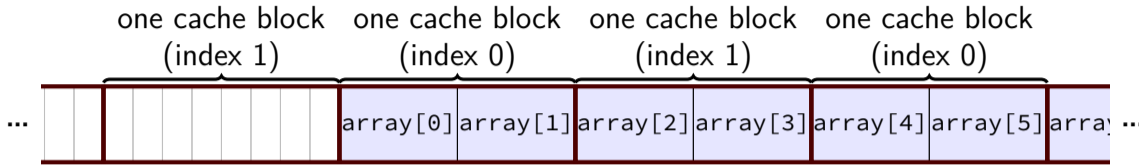
How many data cache misses on a **2**-set direct-mapped cache with 8B blocks?

# exercise solution

one cache block
(index 0)



... | | | | | | | | | | | array[0] | array[1] | array[2] | array[3] | array[4] | array[5] | arra... ...

# exercise solution



one cache block (index 1)   one cache block (index 0)   one cache block (index 1)   one cache block (index 0)

... | | | | | | | | | | | `array[0]`|`array[1]`|`array[2]`|`array[3]`|`array[4]`|`array[5]`|`array` ...

# exercise solution

one cache block (index 1)    one cache block (index 0)    one cache block (index 1)    one cache block (index 0)

… | | | | | | | | | `array[0]` `array[1]` `array[2]` `array[3]` `array[4]` `array[5]` `array…` …

| memory access | set 0 afterwards | set 1 afterwards |
|---|---|---|
| — | (empty) | (empty) |
| read array[0] (miss) | {array[0], array[1]} | (empty) |
| read array[1] (hit) | {array[0], array[1]} | (empty) |
| read array[2] (miss) | {array[0], array[1]} | {array[2], array[3]} |
| read array[3] (hit) | {array[0], array[1]} | {array[2], array[3]} |
| read array[4] (miss) | {array[4], array[5]} | {array[2], array[3]} |
| read array[5] (hit) | {array[4], array[5]} | {array[2], array[3]} |
| read array[6] (miss) | {array[4], array[5]} | {array[6], array[7]} |
| read array[7] (hit) | {array[4], array[5]} | {array[6], array[7]} |

# exercise solution

one cache block   one cache block   one cache block   one cache block

… | | | | | … | array | …

**observation:** what happens in set 0 doesn't affect set 1
when evaluating set 0 accesses,
can ignore non-set 0 accesses/content

| memory a... | (empty) | (empty) |
|---|---|---|
| — | (empty) | (empty) |
| read array[0] (miss) | {array[0], array[1]} | (empty) |
| read array[1] (hit) | {array[0], array[1]} | (empty) |
| read array[2] (miss) | {array[0], array[1]} | {array[2], array[3]} |
| read array[3] (hit) | {array[0], array[1]} | {array[2], array[3]} |
| read array[4] (miss) | {array[4], array[5]} | {array[2], array[3]} |
| read array[5] (hit) | {array[4], array[5]} | {array[2], array[3]} |
| read array[6] (miss) | {array[4], array[5]} | {array[6], array[7]} |
| read array[7] (hit) | {array[4], array[5]} | {array[6], array[7]} |

# exercise solution

one cache block   one cache block   one cache block   one cache block

... array ...

observation: what happens in set 0 doesn't affect set 1
when evaluating set 0 accesses,
can ignore non-set 0 accesses/content

| memory ac | | |
|---|---|---|
| — | (empty) | (empty) |
| read array[0] (miss) | {array[0], array[1]} | (empty) |
| read array[1] (hit) | {array[0], array[1]} | (empty) |
| read array[2] (miss) | {array[0], array[1]} | {array[2], array[3]} |
| read array[3] (hit) | {array[0], array[1]} | {array[2], array[3]} |
| read array[4] (miss) | {array[4], array[5]} | {array[2], array[3]} |
| read array[5] (hit) | {array[4], array[5]} | {array[2], array[3]} |
| read array[6] (miss) | {array[4], array[5]} | {array[6], array[7]} |
| read array[7] (hit) | {array[4], array[5]} | {array[6], array[7]} |

# exercise solution

| | one cache block (index 1) | one cache block (index 0) | one cache block (index 1) | one cache block (index 0) | |
|---|---|---|---|---|---|

... | | | | | array[0] | array[1] | array[2] | array[3] | array[4] | array[5] | array ...

| memory access | set 0 afterwards | set 1 afterwards |
|---|---|---|
| — | (empty) | (empty) |
| read array[0] (miss) | {array[0], array[1]} | (empty) |
| read array[1] (hit) | {array[0], array[1]} | (empty) |
| read array[2] (miss) | {array[0], array[1]} | {array[2], array[3]} |
| read array[3] (hit) | {array[0], array[1]} | {array[2], array[3]} |
| read array[4] (miss) | {array[4], array[5]} | {array[2], array[3]} |
| read array[5] (hit) | {array[4], array[5]} | {array[2], array[3]} |
| read array[6] (miss) | {array[4], array[5]} | {array[6], array[7]} |
| read array[7] (hit) | {array[4], array[5]} | {array[6], array[7]} |

# exercise solution



| | one cache block (index 1) | one cache block (index 0) | one cache block (index 1) | one cache block (index 0) |
|---|---|---|---|---|

... | | array[0] | array[1] | array[2] | array[3] | array[4] | array[5] | array ...

| memory access | set 0 afterwards | set 1 afterwards |
|---|---|---|
| — | (empty) | (empty) |
| read array[0] (miss) | {array[0], array[1]} | (empty) |
| read array[1] (hit) | {array[0], array[1]} | (empty) |
| read array[2] (miss) | {array[0], array[1]} | {array[2], array[3]} |
| read array[3] (hit) | {array[0], array[1]} | {array[2], array[3]} |
| read array[4] (miss) | {array[4], array[5]} | {array[2], array[3]} |
| read array[5] (hit) | {array[4], array[5]} | {array[2], array[3]} |
| read array[6] (miss) | {array[4], array[5]} | {array[6], array[7]} |
| read array[7] (hit) | {array[4], array[5]} | {array[6], array[7]} |

# backup slides
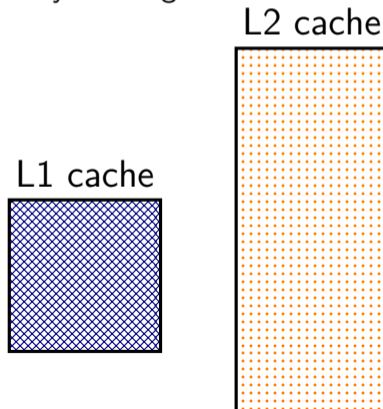
# inclusive versus exclusive

## L2 inclusive of L1

everything in L1 cache duplicated in L2
adding to L1 also adds to L2
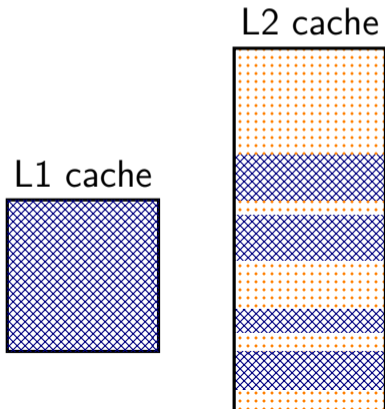
## L2 exclusive of L1

L2 contains different data than L1
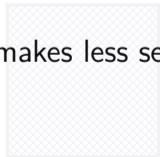adding to L1 must remove from L2
probably evicting from L1 adds to L2

# inclusive versus exclusive

L2 inclusive of L1

everything in L1 cache duplicated in L2
adding to L1 also adds to L2

L2 cache

L1 cache

L2 exclusive of L1

L2 contains different data than L1
adding to L1 must remove from L2
probably evicting from L1 adds to L2

L2 cache

inclusive policy:
no extra work on eviction
but duplicated data

L1 cache

easier to explain when
L$k$ shared by multiple L$(k-1)$ caches?
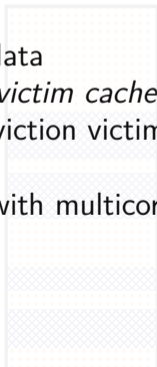
# inclusive versus exclusive

L2 inclusive of L1

everything in L1 cache duplicated in L2
adding to L1 also adds to L2

L2 cache

exclusive policy:
avoid duplicated data
sometimes called *victim cache*
(contains cache eviction victims)

makes less sense with multicore
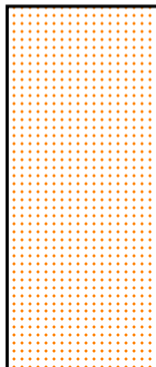
L2 exclusive of L1

L2 contains different data than L1
adding to L1 must remove from L2
probably evicting from L1 adds to L2

L2 cache

L1 cache