## Caching

*Caching is a general technique for exploiting the locality of data reference to make it appear as if there is a large amount of fast memory.*

| Memory Technology | Typical Access Time | $/MB ('97) |
|---|---|---|
| SRAM | 5 - 25ns | $100 - $250 |
| DRAM | 60 - 120ns | $5 - $10 |
| Magnetic Disk | $10^7$ - $2x10^7$ns | $0.10 - $0.20 |

## Terminology

*Locality* -- the property of memory references to cluster

*Temporal locality* -- the tendency of the time intervals between references to a given address to be small

*Spatial locality* -- the tendency of the distances between consecutive memory references to be small

*Memory hierarchy* -- a characteristic of computer design in which a series of storage technologies are used such that the access time is faster as the memory is closer to the processor and the capacity is larger as the memory is further from the processor

```
add    $4,$5,$6
lcw1   $f2,0($4)
lcw1   $f3,4($4)
sw     $7,0($29)
sw     $8,4($29)
bne    $9,$0,loop
```

## Terminology, continued

*Cache* -- memory closest to the processor in a memory hierarchy

*Caching* -- any storage management technique exploiting locality

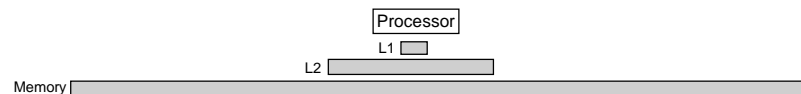*Upper/lower level* -- memory closer/further from the processor

*Block* -- unit of memory transfer between two levels in a memory hierarchy.  Also called a *cache line*

*Hit/Miss* -- accessing data present/not present in a hierarchy level

*Hit rate* -- ratio of hits to total references.  *Miss rate* = 1 - Hit rate

*Hit time* -- time to hit in the cache

*Miss penalty* -- time to move a block from a lower level in the hierarchy and satisfy the processor's request

Processor

L1

L2

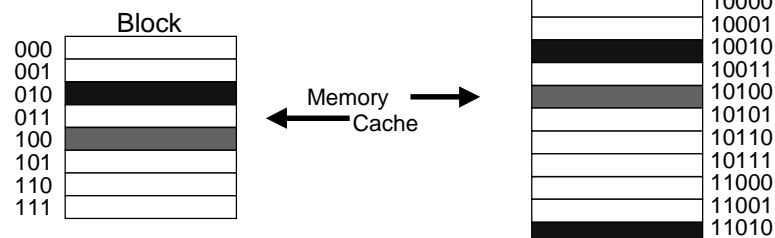Memory

---

## Direct-mapped Cache

Questions asked of a caching technique:

Where is a block stored?
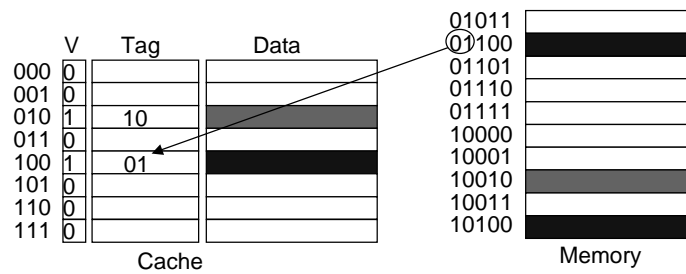
How is a block found?

What block is stored at a location?

A direct mapped cache of size $2^k$ uses the k lsb's of the (block) address.

Block

| 000 |
| 001 |
| 010 |
| 011 |
| 100 |
| 101 |
| 110 |
| 111 |

Memory ⟶
⟵ Cache

Block

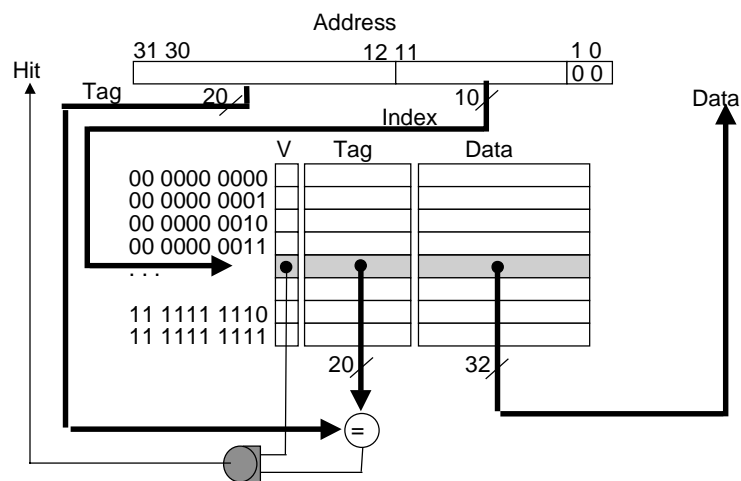| 00000 |
| 00001 |
| 00010 |
| 00011 |
| 00100 |
| 00101 |
| 00110 |
| 00111 |
| 01000 |
| 01001 |
| 01010 |
| 01011 |
| 01100 |
| 01101 |
| 01110 |
| 01111 |
| 10000 |
| 10001 |
| 10010 |
| 10011 |
| 10100 |
| 10101 |
| 10110 |
| 10111 |
| 11000 |
| 11001 |
| 11010 |

# Direct Mapped Cache Fields

- The Tag field stores the msbs of the address.
- The Valid Bit indicates whether the data in the cache block is correct and available.
- The Data field stores the contents of the block.

|  | V | Tag | Data |
|---|---|---|---|
| 000 | 0 | | |
| 001 | 0 | | |
| 010 | 1 | 10 | |
| 011 | 0 | | |
| 100 | 1 | 01 | |
| 101 | 0 | | |
| 110 | 0 | | |
| 111 | 0 | | |

Cache

Memory

01011
01100
01101
01110
01111
10000
10001
10010
10011
10100

# Direct-Mapped Cache Operation

Address

Hit

31 30      12 11      1 0

Tag    20        10    Data

Index

|  | V | Tag | Data |
|---|---|---|---|
| 00 0000 0000 | | | |
| 00 0000 0001 | | | |
| 00 0000 0010 | | | |
| 00 0000 0011 | | | |
| . . . | | | |
| 11 1111 1110 | | | |
| 11 1111 1111 | | | |

20      32

=

# Handling Misses

The processor has to stall the instruction that missed; instruction misses stall the pipeline at IF while data misses stall in MEM.
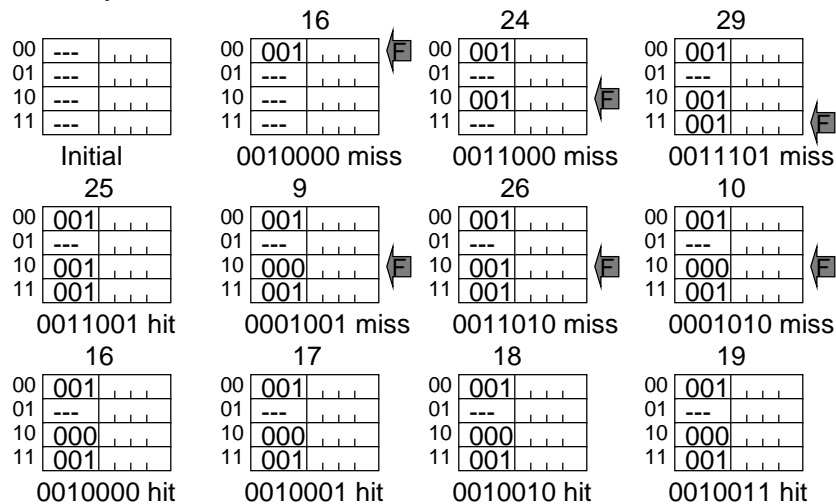
Operations by controller on a miss:

1. Compute PC-4

2. Access address in main memory and wait for completion.

3. Move data to cache, write tag bits, set Valid.

4. Restart execution pipeline at the fetch for instruction misses, or MEM for data misses.

---

# Reference Sequence

Block size = 4 bytes
Address = 7 bits
T = 3, CA = 2, BA = 2

Memory Refs: 16, 24, 29, 25, 9, 26, 10, 16, 17, 18, 19