# Fundamentals of Computer Systems
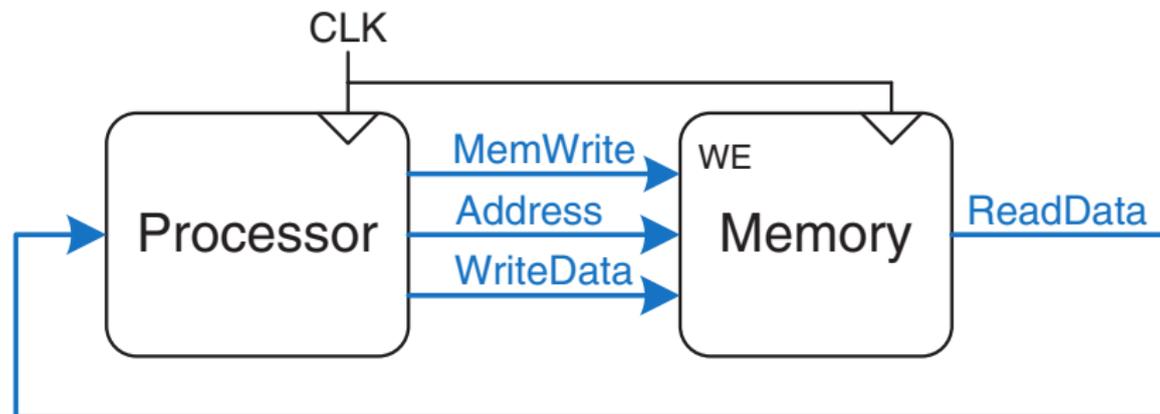## Caches

Stephen A. Edwards
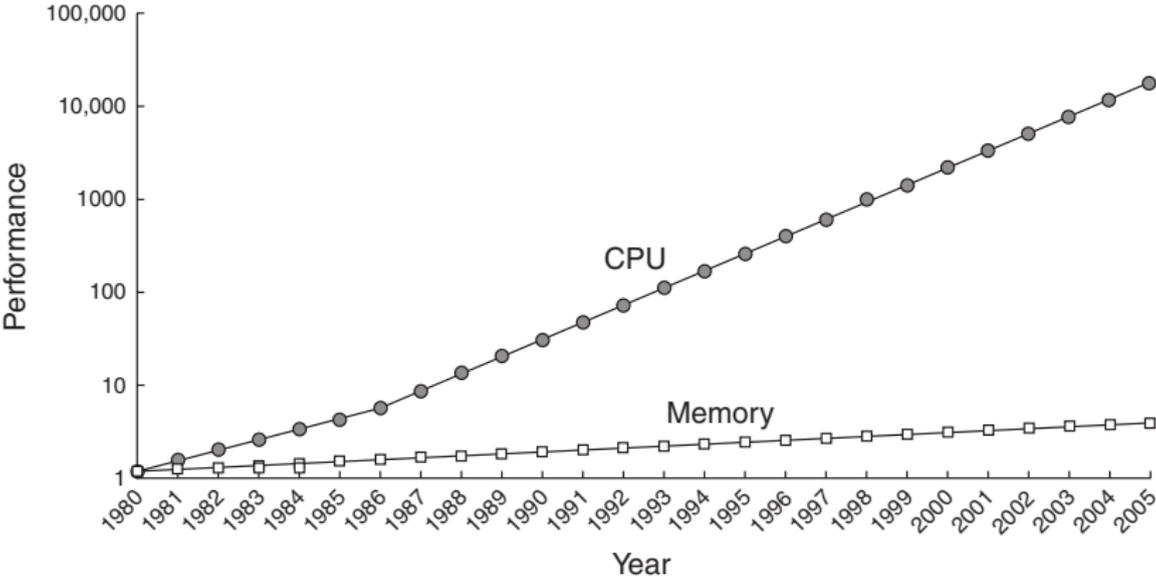
Columbia University

## Summer 2021

Illustrations Copyright © 2007 Elsevier

# Computer Systems

Performance depends on which is slowest:
the processor or the memory system
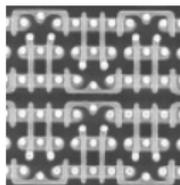
# Memory Speeds Haven't Kept Up



Our single-cycle memory assumption has been wrong since 1980.

Hennessy and Patterson. *Computer Architecture: A Quantitative Approach.* 3rd ed., Morgan Kaufmann, 2003.

# Your Choice of Memories

| | Fast | Cheap | Large |
|---|:---:|:---:|:---:|
|  On-Chip SRAM | ✔ | ✔ | |
|  Commodity DRAM | | ✔ | ✔ |
|  Supercomputer | ✔ | | ✔ |

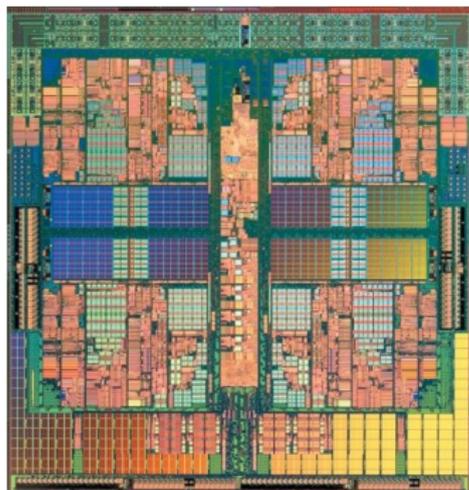# Memory Hierarchy

Fundmental trick to making a big memory appear fast

| Technology | Cost ($/Gb) | Access Time (ns) | Density (Gb/cm2) |
|---|---|---|---|
| SRAM | 30 000 | 0.5 | 0.00025 |
| DRAM | 10 | 100 | 1 – 16 |
| Flash | 2 | 300* | 8 – 32 |
| Hard Disk | 0.1 | 10 000 000 | 500 – 2000 |

*Read speed; writing much, much slower

# A Modern Memory Hierarchy



AMD Phenom 9600
Quad-core
2.3 GHz
1.1–1.25 V
95 W
65 nm

A desktop machine:

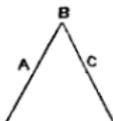| Level | Size | Tech. |
|---|---|---|
| L1 Instruction* | 64 K | SRAM |
| L1 Data* | 64 K | SRAM |
| L2* | 512 K | SRAM |
| L3 | 2 MB | SRAM |
| Memory | 4 GB | DRAM |
| Disk | 500 GB | Magnetic |

*per core

# Temporal Locality

## FIRST BOOK

### DEFINITIONS.

1. A point is that which is without parts.
2. A line is length without breadth.
3. The extremities of a line are points.
4. A right line, is that which lies evenly between its extremities.
5. A superficies, is that which has only length and breadth.
6. The boundings of a superficies are lines.
7. A plane superficies, is that which lies evenly between its extreme right lines.
8. A rectilineal angle, is the inclination of two right lines to each other, which touch, but do not form one straight line.

An angle is designated either by one letter at the vertex; or three, of which the middle one is at the vertex, the remaining two any place on the legs.

9. The legs of an angle, are the lines which make the angle.
10. The vertex of an angle is the point in which legs mutually touch each other.

What path do your eyes take when you read this?

Did you look at the drawings more than once?

Euclid's *Elements*

# Spatial Locality



If you need something, you may also need something nearby

# Memory Performance

Hit: Data is found in the level of memory hierarchy

Miss: Data not found; will look in next level



Hit Rate $= \frac{\text{Number of hits}}{\text{Number of accesses}}$

Miss Rate $= \frac{\text{Number of misses}}{\text{Number of accesses}}$

Hit Rate $+$ Miss Rate $= 1$

The expected access time $E_L$ for a memory level $L$ with latency $t_L$ and miss rate $M_L$:

$$E_L = t_L + M_L \cdot E_{L+1}$$

# Memory Performance Example

Two-level hierarchy: Cache and main memory

Program executes 1000 loads & stores

750 of these are found in the cache

*What's the cache hit and miss rate?*

# Memory Performance Example

Two-level hierarchy: Cache and main memory

Program executes 1000 loads & stores

750 of these are found in the cache

*What's the cache hit and miss rate?*

$$\text{Hit Rate} = \frac{750}{1000} = 75\%$$
$$\text{Miss Rate} = 1 - 0.75 = 25\%$$

---

If the cache takes 1 cycle and the main memory 100,

*What's the expected access time?*

# Memory Performance Example

Two-level hierarchy: Cache and main memory

Program executes 1000 loads & stores

750 of these are found in the cache

*What's the cache hit and miss rate?*

$$\text{Hit Rate} = \frac{750}{1000} = 75\%$$
$$\text{Miss Rate} = 1 - 0.75 = 25\%$$

---

If the cache takes 1 cycle and the main memory 100,

*What's the expected access time?*

Expected access time of main memory: $E_1 = 100$ cycles

Access time for the cache: $t_0 = 1$ cycle

Cache miss rate: $M_0 = 0.25$

$$E_0 = t_0 + M_0 \cdot E_1 = 1 + 0.25 \cdot 100 = 26 \text{ cycles}$$

# Cache

Highest levels of memory hierarchy

Fast: level 1 typically 1 cycle access time

With luck, supplies most data

Cache design questions:



What data does it hold?      *Recently accessed*

How is data found?      *Simple address hash*

What data is replaced?      *Often the oldest*

# What Data is Held in the Cache?

Ideal cache: always correctly guesses what you want before you want it.

Real cache: never that smart

**Caches Exploit**

**Temporal Locality**

Copy newly accessed data into cache, replacing oldest if necessary

**Spatial Locality**

Copy nearby data into the cache at the same time

Specifically, always read and write a block at a time (e.g., 64 bytes), never a single byte.
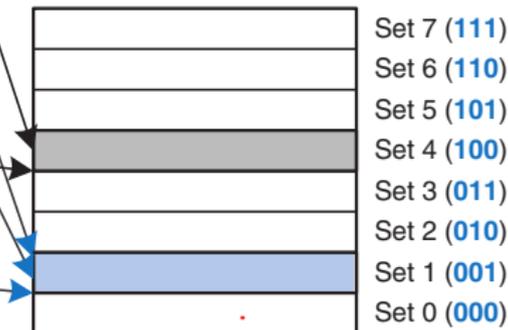
# A Direct-Mapped Cache

| Address | Data |
|---|---|
| 11...11111100 | mem[0xFFFFFFFC] |
| 11...11111000 | mem[0xFFFFFFF8] |
| 11...11110100 | mem[0xFFFFFFF4] |
| 11...11110000 | mem[0xFFFFFFF0] |
| 11...11101100 | mem[0xFFFFFFEC] |
| 11...11101000 | mem[0xFFFFFFE8] |
| 11...11100100 | mem[0xFFFFFFE4] |
| 11...11100000 | mem[0xFFFFFFE0] |
| ⋮ | ⋮ |
| 00...00100100 | mem[0x00000024] |
| 00...00100000 | mem[0x00000020] |
| 00...00011100 | mem[0x0000001C] |
| 00...00011000 | mem[0x00000018] |
| 00...00010100 | mem[0x00000014] |
| 00...00010000 | mem[0x00000010] |
| 00...00001100 | mem[0x0000000C] |
| 00...00001000 | mem[0x00000008] |
| 00...00000100 | mem[0x00000004] |
| 00...00000000 | mem[0x00000000] |

Set number 2^30-Word Main Memory

This simple cache has

- ▶ 8 *sets*
- ▶ 1 *block* per set
- ▶ 4 bytes per block

To simplify answering "is this memory in the cache?," each byte is mapped to exactly one set.

| | |
|---|---|
| | Set 7 (**111**) |
| | Set 6 (**110**) |
| | Set 5 (**101**) |
| | Set 4 (**100**) |
| | Set 3 (**011**) |
| | Set 2 (**010**) |
| | Set 1 (**001**) |
| | Set 0 (**000**) |

2^3-Word Cache

# Direct-Mapped Cache Hardware



Address bits:
0–1: byte within block
2–4: set number
5–31: block "tag"

8-entry x
(1+27+32)-bit
SRAM

Cache hit if
in the set of the address,

- ▶ block is valid (V=1)
- ▶ tag (address bits 5–31) matches

# Direct-Mapped Cache Behavior

A dumb loop:

repeat 5 times

load from 0x4;
load from 0xC;
load from 0x8.

```
      li    $t0, 5
l1:   beq   $t0, $0, done
      lw    $t1, 0x4($0)
      lw    $t2, 0xC($0)
      lw    $t3, 0x8($0)
      addiu $t0, $t0, -1
      j     l1
done:
```

Memory Address

| | Tag | Set | Byte Offset |
|---|---|---|---|
| | 00...00 | 001 | 00 |



| V | Tag | Data | |
|---|---|---|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 1 | 00...00 | mem[0x00...0C] | Set 3 (011) |
| 1 | 00...00 | mem[0x00...08] | Set 2 (010) |
| 1 | 00...00 | mem[0x00...04] | Set 1 (001) |
| 0 | | | Set 0 (000) |

Cache when reading 0x4 last time

Assuming the cache starts empty, what's the miss rate?

# Direct-Mapped Cache Behavior

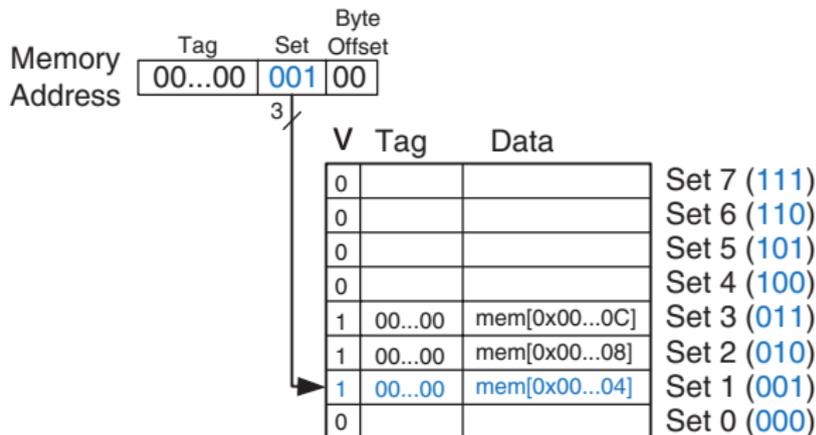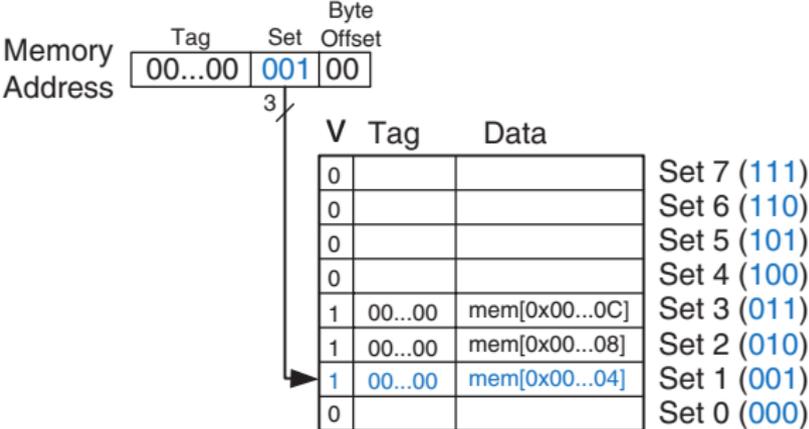A dumb loop:

repeat 5 times

load from 0x4;
load from 0xC;
load from 0x8.

```
    li    $t0, 5
l1: beq   $t0, $0, done
    lw    $t1, 0x4($0)
    lw    $t2, 0xC($0)
    lw    $t3, 0x8($0)
    addiu $t0, $t0, -1
    j     l1
done:
```



Memory Address

| Tag | Set | Byte Offset |
|-----|-----|-------------|
| 00...00 | 001 | 00 |

| V | Tag | Data | |
|---|-----|------|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 1 | 00...00 | mem[0x00...0C] | Set 3 (011) |
| 1 | 00...00 | mem[0x00...08] | Set 2 (010) |
| 1 | 00...00 | mem[0x00...04] | Set 1 (001) |
| 0 | | | Set 0 (000) |

Cache when reading 0x4 last time

Assuming the cache starts empty, what's the miss rate?

| 4 | C | 8 | 4 | C | 8 | 4 | C | 8 | 4 | C | 8 | 4 | C | 8 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| M | M | M | H | H | H | H | H | H | H | H | H | H | H | H |

3/15 = 0.2 = 20%

compulsory misses

# Direct-Mapped Cache: Conflict



Memory Address

| Tag | Set | Byte Offset |
|-----|-----|-------------|
| 00...01 | 001 | 00 |

A dumber loop:
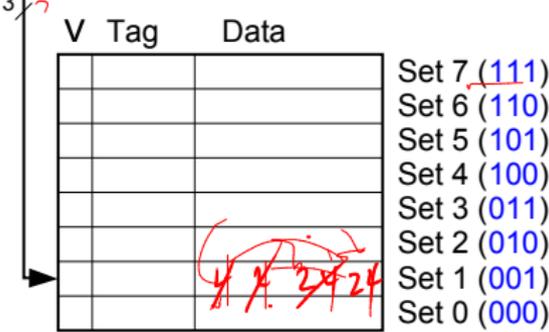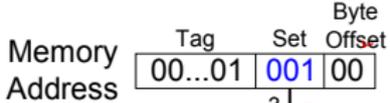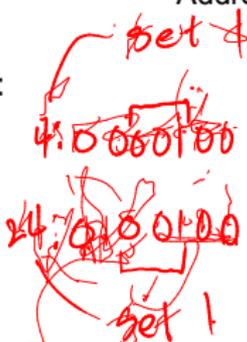
repeat 5 times

load from 0x4;
load from 0x24

```
     li    $t0, 5
l1:  beq   $t0, $0, done
     lw    $t1,  0x4($0)
     lw    $t2, 0x24($0)
     addiu $t0, $t0, -1
     j     l1
done:
```
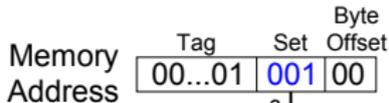
V  Tag  Data

Set 7 (111)
Set 6 (110)
Set 5 (101)
Set 4 (100)
Set 3 (011)
Set 2 (010)
Set 1 (001)
Set 0 (000)

Cache State

Assuming the cache starts empty, what's the miss rate?

These are *conflict misses*

# Direct-Mapped Cache: Conflict

Memory Address

| Tag | Set | Byte Offset |
|-----|-----|-------------|
| 00...01 | 001 | 00 |

k-way
blocks 34 byte.

A dumber loop:

repeat 5 times

load from 0x4;
load from 0x24

0000 0100  4
tag  set offset

0000 10000
tag  set offset

| V | Tag | Data | |
|---|-----|------|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 0 | | | Set 3 (011) |
| 0 | | | Set 2 (010) |
| 1 | 00...00 | mem[0x00...0x0] | Set 1 (001) |
| 0 | | | Set 0 (000) |

Cache State

```
        li    $t0, 5
l1:     beq   $t0, $0, done
        lw    $t1,  0x4($0)
        lw    $t2, 0x24($0)
        addiu $t0, $t0, -1
        j     l1
done:
```

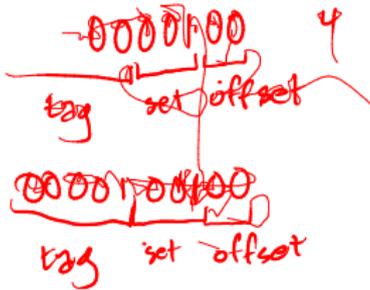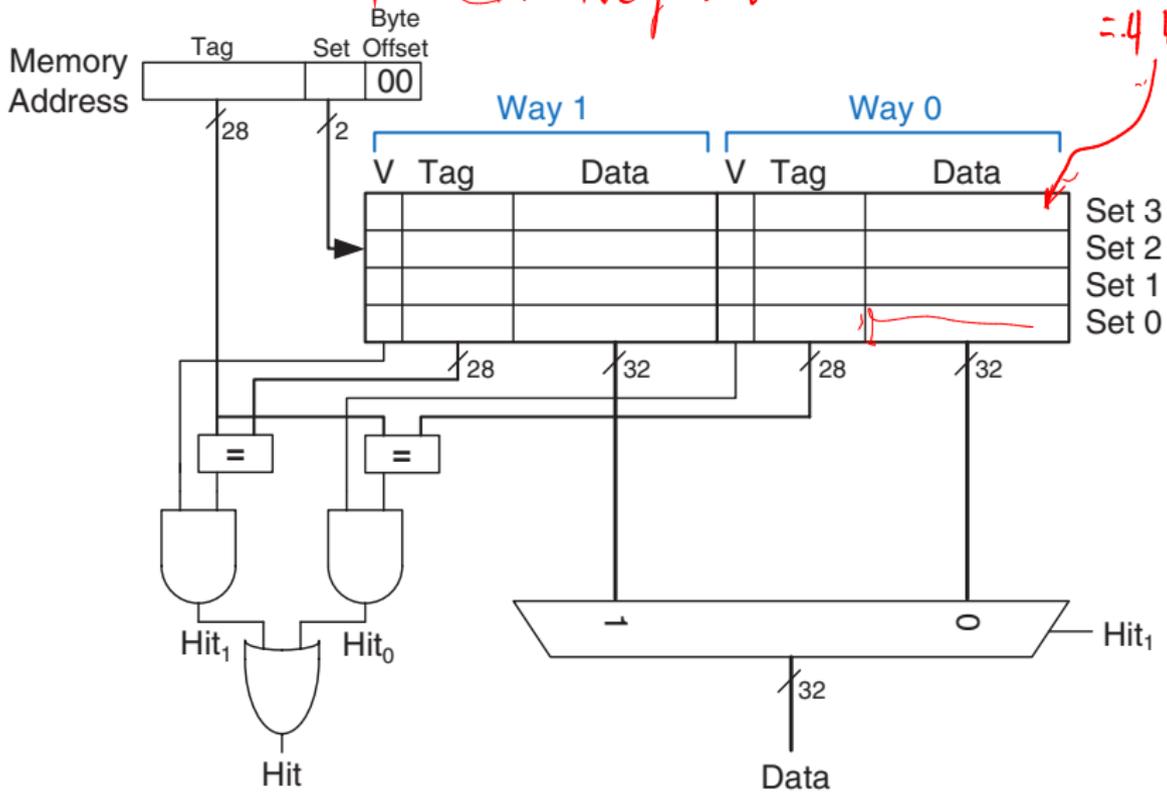Assuming the cache starts empty, what's the miss rate?

| 4 | 24 | 4 | 24 | 4 | 24 | 4 | 24 | 4 | 24 |
|---|----|---|----|---|----|---|----|---|----|
| M | M  | M | M  | M | M  | M | M  | M | M  |

10/10 = 1 = 100% Oops

These are _conflict misses_

# No Way! Yes Way! 2-Way Set Associative Cache



Direct-mapped "1-way set associative" block = 4 bytes

Memory Address

Tag | Set | Byte Offset

00

28 | 2

Way 1 | Way 0

V Tag Data | V Tag Data

Set 3
Set 2
Set 1
Set 0

28 | 32 | 28 | 32

= | =

Hit₁ | Hit₀

Hit

1 | 0 | Hit₁

32

Data

# 2-Way Set Associative Behavior

```
    li    $t0, 5
l1: beq   $t0, $0, done
    lw    $t1,  0x4($0)
    lw    $t2, 0x24($0)
    addiu $t0, $t0, -1
    j     l1
done:
```

*Associativity reduces conflict misses*

| | Way 1 | | | Way 0 | | |
|---|---|---|---|---|---|---|
| V | Tag | Data | V | Tag | Data | |
| 0 | | | 0 | | | Set 3 |
| 0 | | | 0 | | | Set 2 |
| 1 | 00...00 | mem[0x00...24] | 1 | 00...10 | mem[0x00...04] | Set 1 |
| 0 | | | 0 | | | Set 0 |

# An Eight-way Fully Associative Cache

*set # 0 bits*  *byte/ = 1 set only*

*tag*  *block offset!*

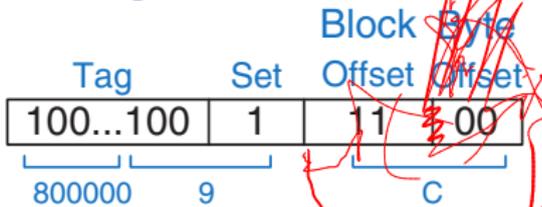| | Way 7 | | | Way 6 | | | Way 5 | | | Way 4 | | | Way 3 | | | Way 2 | | | Way 1 | | | Way 0 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data |
| | | | | | | | | | | | | | | | | | | | | | | | |

*Set 0*

No conflict misses: only compulsory or capacity misses

*Run out of space*

Either very expensive or slow because of all the associativity

# Exploiting Spatial Locality: Larger Blocks



**Handwritten annotations (top):** tag | Set size | block size | 64 block | 32-5-c = 23 bits · 4 sets 3 bits | 2c = 24 bits

**Memory Address:** 0x8000 0009C:

| Tag | Set | Block Offset | Byte Offset |
|---|---|---|---|
| 100...100 | 1 | 11 | 00 |
| 800000 | 9 | | C |

**Handwritten annotations (right):** byte offset with in the block · 16-byte block · 64-byte

**Diagram labels:**
Memory Address — Tag | Set | Block Offset | Byte Offset | 00

27 / 2

V  Tag      Data      Set 1 / Set 0

27  32  32  32  32

11  10  01  00

32

Hit

Data

2 sets
1 block per set (Direct Mapped)
4 words per block

# Direct-Mapped Cache Behavior w/ 4-word block



The dumb loop:

repeat 5 times

load from 0x4;
load from 0xC;
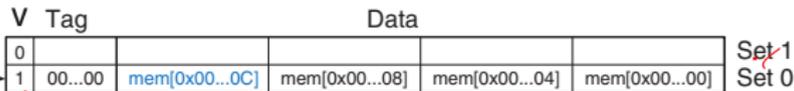load from 0x8.

```
      li    $t0, 5
l1:   beq   $t0, $0, done
      lw    $t1, 0x4($0)
      lw    $t2, 0xC($0)
      lw    $t3, 0x8($0)
      addiu $t0, $t0, -1
      j     l1
done:
```

Cache when reading 0xC

Assuming the cache starts empty, what's the miss rate?

# Direct-Mapped Cache Behavior w/ 4-word block

Memory Address

| | | Block Offset | Byte Offset |
|---|---|---|---|
| Tag | Set | | |
| 00...00 | 0 | 11 | 00 |

| V | Tag | Data | | | | |
|---|---|---|---|---|---|---|
| 0 | | | | | | Set 1 |
| 1 | 00...00 | mem[0x00...0C] | mem[0x00...08] | mem[0x00...04] | mem[0x00...00] | Set 0 |

Cache when reading 0xC

The dumb loop:

repeat 5 times

load from 0x4;
load from 0xC;
load from 0x8.

```
    li    $t0, 5
l1: beq   $t0, $0, done
    lw    $t1, 0x4($0)
    lw    $t2, 0xC($0)
    lw    $t3, 0x8($0)
    addiu $t0, $t0, -1
    j     l1
done:
```
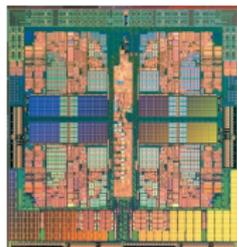
Assuming the cache starts empty, what's the miss rate?

$$\frac{4 \ C \ 8 \ 4 \ C \ 8 \ 4 \ C \ 8 \ 4 \ C \ 8 \ 4 \ C \ 8}{M \ H \ H \ H \ H \ H \ H \ H \ H \ H \ H \ H \ H \ H \ H}$$

$1/15 = 0.0666 = 6.7\%$

*Larger blocks reduce compulsory misses by exploting spatial locality*

# The Desktop Machine Revisited

cache = blocks × ways × sets.
size size



AMD Phenom
9600
Quad-core
2.3 GHz
1.1–1.25 V
95 W
65 nm

On-chip caches:

| Cache | Size | Sets | Ways | Block |
|-------|------|------|------|-------|
| L1I* | 64 K | 512 | 2-way | 64-byte |
| L1D* | 64 K | 512 | 2-way | 64-byte |
| L2* | 512 K | 512 | 16-way | 64-byte |
| L3 | 2 MB | 1024 | 32-way | 64-byte |

*per core

# Intel On-Chip Caches

| Chip | Year | Freq. (MHz) | L1 | | L2 |
|------|------|-------------|-----|-------|-----|
| | | | Data | Instr | |
| 80386 | 1985 | 16–25 | off-chip | | none |
| 80486 | 1989 | 25–100 | 8K unified | | off-chip |
| Pentium | 1993 | 60–300 | 8K | 8K | off-chip |
| Pentium Pro | 1995 | 150–200 | 8K | 8K | 256K–1M (MCM) |
| Pentium II | 1997 | 233–450 | 16K | 16K | 256K–512K (Cartridge) |
| Pentium III | 1999 | 450–1400 | 16K | 16K | 256K–512K |
| Pentium 4 | 2001 | 1400–3730 | 8–16K | 12k op trace cache | 256K–2M |
| Pentium M | 2003 | 900–2130 | 32K | 32K | 1M–2M |
| Core 2 Duo | 2005 | 1500–3000 | 32K per core | 32K per core | 2M–6M |