Temporal locality → item itself

Spatial locality → nearby addresses

## Direct Mapped Cache Simulation (E=1)

$M = 16$ byte addresses $\quad 2^4 \Rightarrow 4$ bit

$B = 2$ bytes/block $\quad \Rightarrow 2^1 = 1 = |BO|$

$S = 4$ sets $= 2^2 \Rightarrow |SI| = 2$

$E = 1$ Blocks/set

PA

| | Tag | SI | BO |
|---|---|---|---|
| | | | |

| | V | Tag | Block | |
|---|---|---|---|---|
| 00 | 1 | 1 | M[8] | M[9] |
| 01 | 0 | | | |
| 10 | 0 | | | |
| 11 | 1 | 0 | M[6] | M[7] |

| 0 | 0000 | miss |
|---|---|---|
| 1 | 0001 | hit |
| 7 | 0111 | miss |
| 8 | 1000 | miss |
| 0 | 0000 | miss |

Cache size $= C = S \times E \times B$

For E-way set associative case, we need to compare all tags on that set with our tag

↳ no match: one line is selected for eviction and replacement

## 2-way set Associative Cache Simulation (E=2)

$M = 16$ byte addresses $= 2^4 \Rightarrow |PA| = 4$

$B = 2$ bytes/block $= 2^1 \Rightarrow |BO| = 1$

$S = 2$ sets $= 2^1 \Rightarrow |SI| = 1$

$|Tag| = 2$

$E = 2$ Blocks/set

| | Tag | SI | BO |
|---|---|---|---|
| | | | |

| | V | Tag | Block | | V | Tag | Block | |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 00 | M[0] | M[1] | 1 | 10 | M[8] | M[9] |
| 1 | 1 | 01 | M[6] | M[7] | | | | |

| 0 | 0000 | miss |
|---|---|---|
| 1 | 0001 | hit |
| 7 | 0111 | miss |
| 8 | 1000 | mis |
| 0 | 0000 | hit |

## Write Operations

**What to do on a write-hit?** → The line you want to write is on the cache.

    **Write-through** (write immediately to the memory)

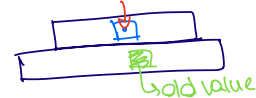        Write both in cache and to memory

    **Write-back** (defer write to mem until replacement of line)

        Need a dirty bit (line different from mem or not)

| D | V | Tag | Data |
|---|---|-----|------|

$D=0 \rightarrow$ no need to update the mem

$D=1 \rightarrow$

**What to do on a write-miss?** → You want to write into a mem address which is not in the cache.

old value

    **Write-allocate**

        Load into cache, update line in cache

    **No-write-allocate**

        Writes immediately to the memory

↳ when we have to replace the block, update data back to mem if the write bit is set.

↳ otherwise, don't.

**Typical**

    Write through + no-write-allocate

    Write-back + write allocate

---

<u>Örnek:</u>    4 i-cache, d-cache

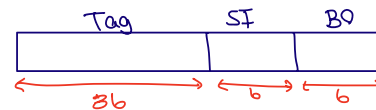        32 KB, 8 way, Mem Address = 48 bit

$C = 32 \text{ KB} = 2^5 \times 2^{10} = 2^{15}$

$E = 8 = 2^3 \rightarrow$ 8 lines in each set

$S = 2^6$                $|SI| = 6$

$C = 2^{15} = S \times E \times B \Rightarrow \dfrac{2^{15}}{2^3 \times 2^6} = 2^6 = B$   $|BO| = 6$

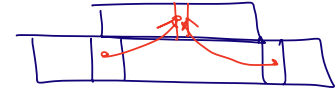| Tag | SI | BO |
|-----|----|----|
| 36 | 6 | 6 |

---

## TYPES OF CACHE MISSES

**Cold (compulsory) miss**

    The cache is empty

**Capacity miss**

    The set of active cache blocks (working set) is larger than the cache.

Conflict miss → Most likely to happen in direct-mapped access

Multiple data objects all map to the same level k block

## CACHE PERFORMANCE METRICS

### Miss rate

Fraction of memory references not found in cache (misses/accesses)

= 1 − hit rate

### Hit time

Time to deliver a line in the cache to the processor

↳ Includes time to determine whether the line is in the cache

### Miss Penalty

Additional time required because of a miss

Hit rate vs. Miss rate

Assume
{
cache hit time → 1 cycle
miss penalty → 100 cycles
}

97% hits : 1 cycle + 0.03 * 100 cycles = 4 cycles

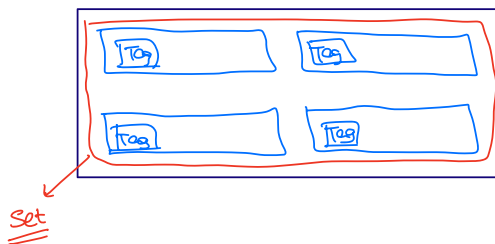99% hits : 1 cycle + 0.01 * 100 cycles = 2 cycles

Hit rate 97% is twice as good as 99%
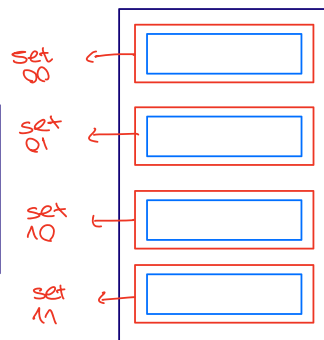
↳ use miss rate

---

Given  $C = S \times E \times B$

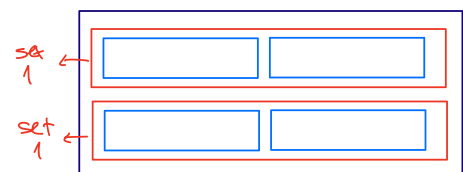Max number of Tag comparisons needed
Fully Associative (S=1)

Set

Address | Tag | BO |

1-tag comparison needed
Direct-mapped (E=1)

set 00
set 01
set 10
set 11

Tag | SI | BO |
← 2 →

(E=2)
2-way Set Associative
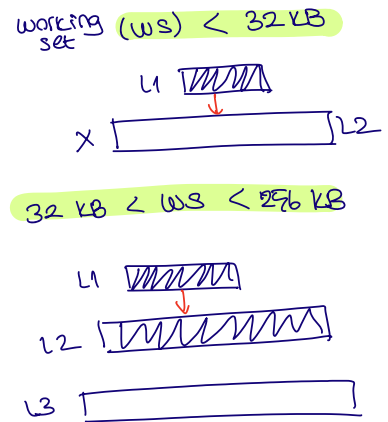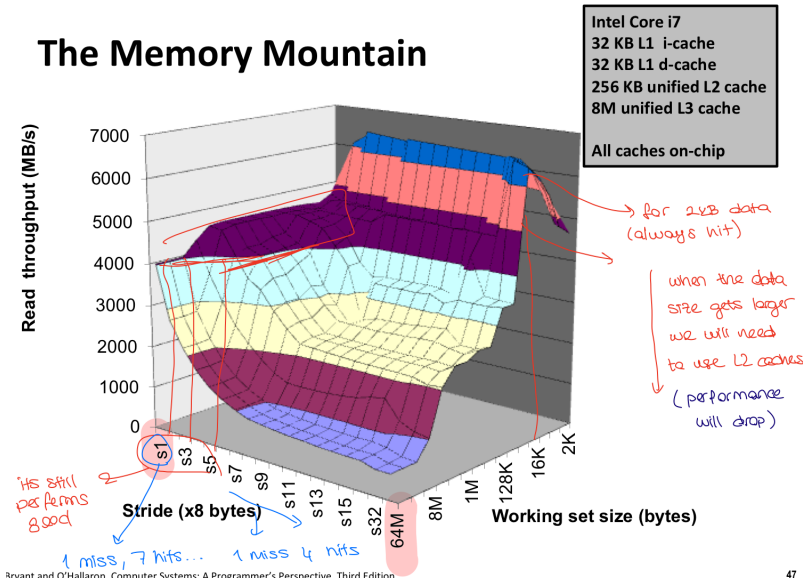
set 1
set 1

Tag | SI | BO |
← 1 →

*Read throughput*

Number of bytes read from mem per second (MB/s)

*Memory Mountain*

Measured read throughput as a function of spatial and temporal locality

## The Memory Mountain

Intel Core i7
32 KB L1 i-cache
32 KB L1 d-cache
256 KB unified L2 cache
8M unified L3 cache

All caches on-chip

*working set (WS) < 32 KB*

*32 KB < WS < 256 KB*

→ for 2KB data (always hit)

when the data size gets larger we will need to use L2 caches (performance will drop)

'its still performs good

1 miss, 7 hits... 1 miss 4 hits

Read throughput (MB/s): 7000, 6000, 5000, 4000, 3000, 2000, 1000, 0

Stride (x8 bytes): s1, s3, s5, s7, s9, s11, s13, s15, s32, s64

Working set size (bytes): 64M, 8M, 1M, 128K, 16K, 2K

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition
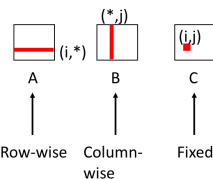
47

## Matrix Multiplication (ijk)

*2 loads, 0 stores*
*misses/iter = 1.25*

```
/* ijk */
for (i=0; i<n; i++)  {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

Inner loop:

(i,*) A — Row-wise
(*,j) B — Column-wise
(i,j) C — Fixed

|a_{ij}| = 8 bytes
Block size = 32 bytes

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.25 | 1.0 | 0.0 |

→ I'm accessing this 1 time (negligible)

## Matrix Multiplication (jik)

```
/* jik */
for (j=0; j<n; j++) {
  for (i=0; i<n; i++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum
  }
}
```
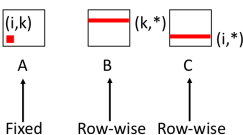
Inner loop:

(i,*) A — Row-wise
(*,j) B — Column-wise
(i,j) C — Fixed

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.25 | 1.0 | 0.0 |

## Matrix Multiplication (ikj)

*2 loads, 1 store*
*misses/iter = 0.5*

```
/* kij */
for (k=0; k<n; k++) {
  for (i=0; i<n; i++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
```

Inner loop:

(i,k) A — Fixed
(k,*) B — Row-wise
(i,*) C — Row-wise

Access only 1 time

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

```
/* ikj */
for (i=0; i<n; i++) {
  for (k=0; k<n; k++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
```

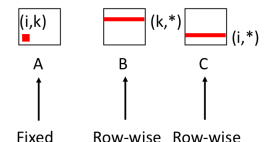Inner loop:

(i,k) A — Fixed
(k,*) B — Row-wise
(i,*) C — Row-wise

Misses per inner loop iteration:

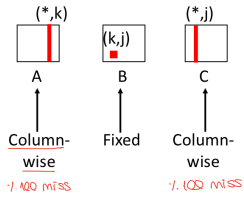| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

## Matrix Multiplication (jki)

```
/* jki */
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```

Inner loop:



| | (*,k) | (k,j) | (*,j) |
| | A | B | C |
| | Column-wise | Fixed | Column-wise |

Misses per inner loop iteration:

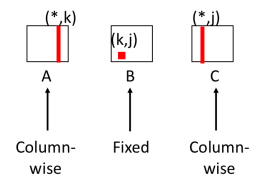| A | B | C |
|---|---|---|
| 1.0 | 0.0 | 1.0 |

## Matrix Multiplication (kji)

```
/* kji */
for (k=0; k<n; k++) {
  for (j=0; j<n; j++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```

Inner loop:



| | (*,k) | (k,j) | (*,j) |
| | A | B | C |
| | Column-wise | Fixed | Column-wise |

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 1.0 | 0.0 | 1.0 |