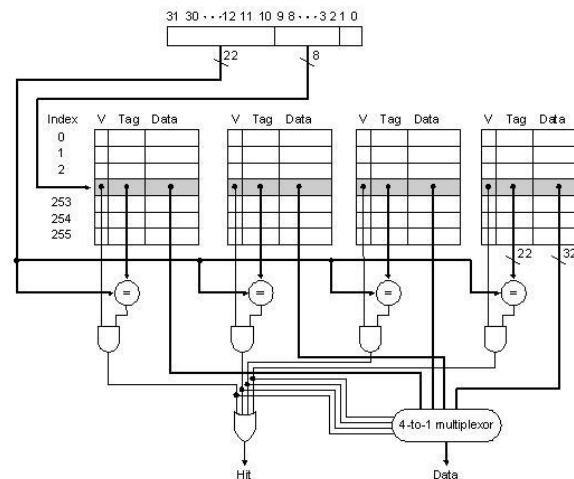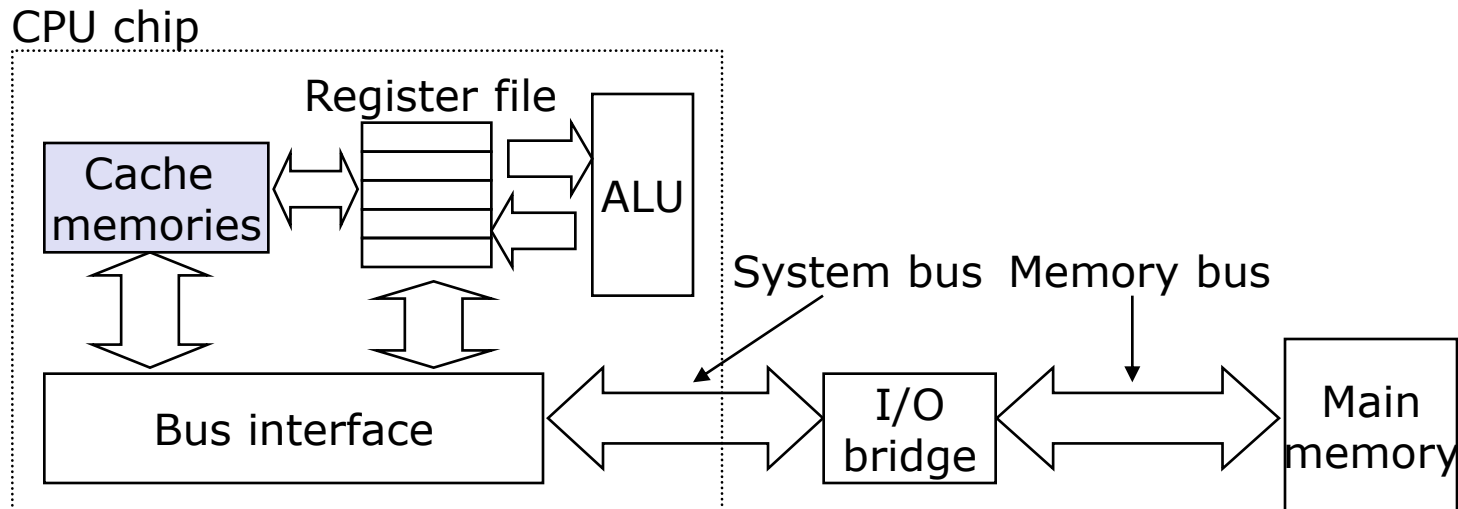# Accessing The Outside World

# Cache Memories

# Cache Memories
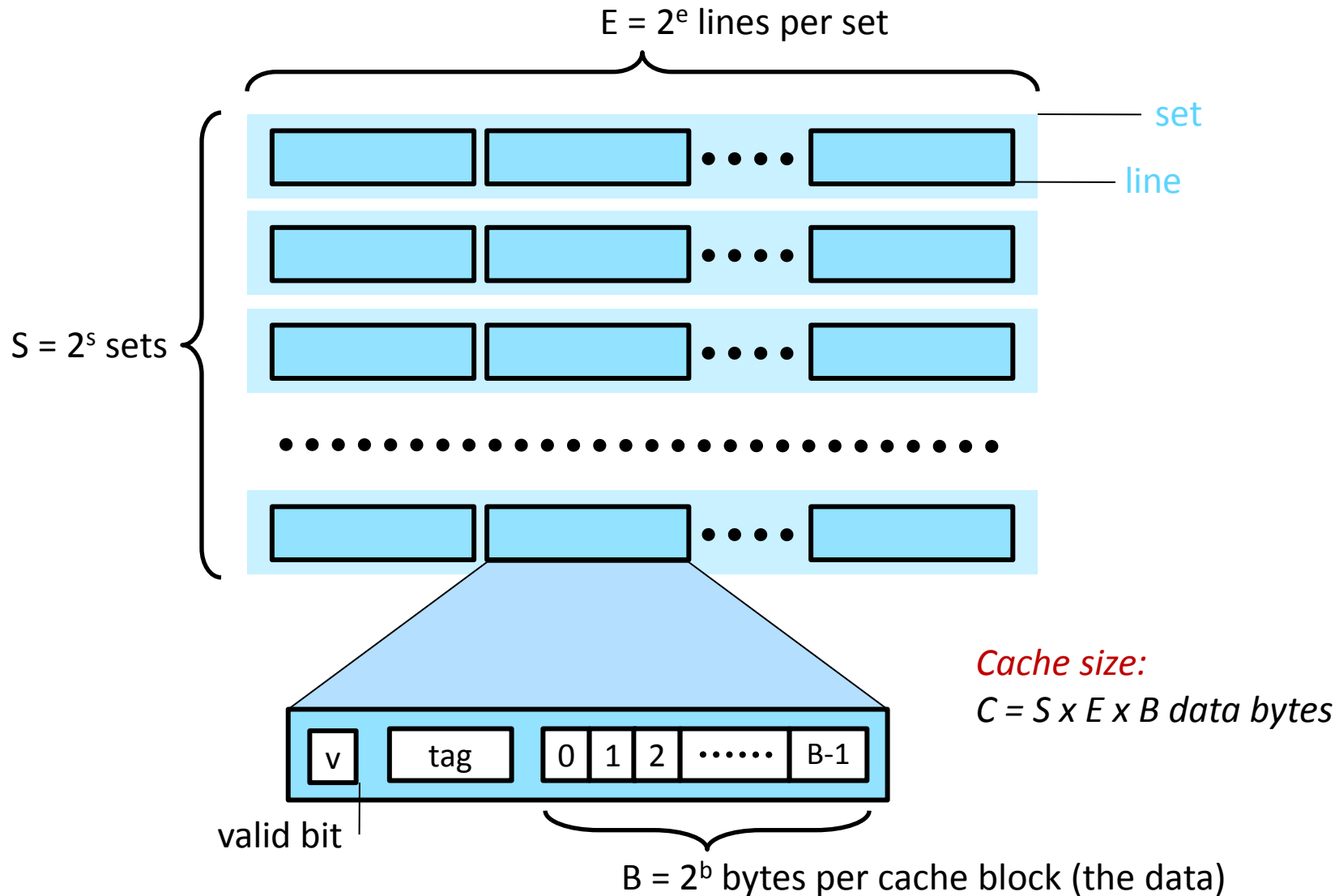
- Cache memory organization and operation
- Performance impact of caches
  - The memory mountain
  - Rearranging loops to improve spatial locality

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Cache Memories

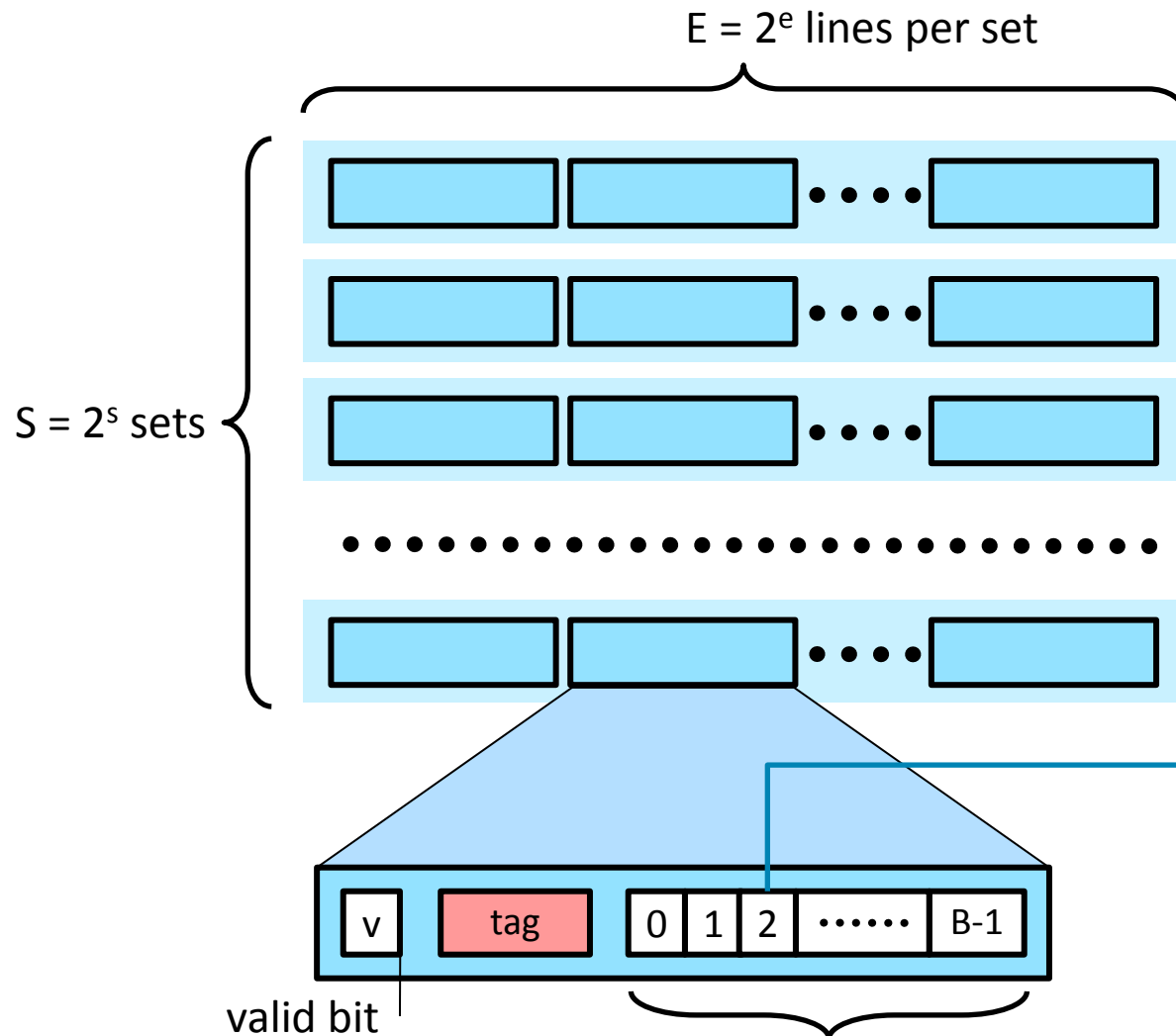- **Cache memories** are small, fast SRAM-based memories managed automatically in hardware.
    - Hold frequently accessed blocks of main memory
- CPU looks first for data in caches (e.g., L1, L2, and L3), then in main memory.
- Typical system structure:

CPU chip

Register file

Cache memories

ALU

System bus   Memory bus

Bus interface

I/O bridge

Main memory

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set

set

line

$S = 2^s$ sets

Cache size:

$C = S \times E \times B$ data bytes

| v | tag | 0 | 1 | 2 | ...... | B-1 |

valid bit

$B = 2^b$ bytes per cache block (the data)
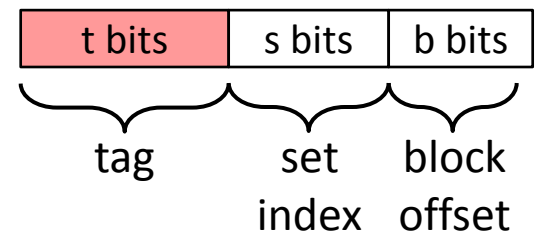
CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Cache Read

*• Locate set*
*• Check if any line in set has matching tag*
*• Yes + line valid: hit*
*• Locate data starting at offset*

$E = 2^e$ lines per set

$S = 2^s$ sets

Address of word:

| t bits | s bits | b bits |
|--------|--------|--------|
| tag | set index | block offset |

data begins at this offset

| v | tag | 0 | 1 | 2 | ⋯⋯ | B-1 |
|---|-----|---|---|---|-----|-----|

valid bit

$B = 2^b$ bytes per cache block (the data)

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size 8 bytes

$S = 2^s$ sets

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

• • • • • • • • • • • • • • • • • • • • • • •

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

Address of int:

| t bits | 0…01 | 100 |

find set

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set
Assume: cache block size 8 bytes



valid?  +  match: assume yes = hit

Address of int:

| t bits | 0...01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

# Example: Direct Mapped Cache (E = 1)

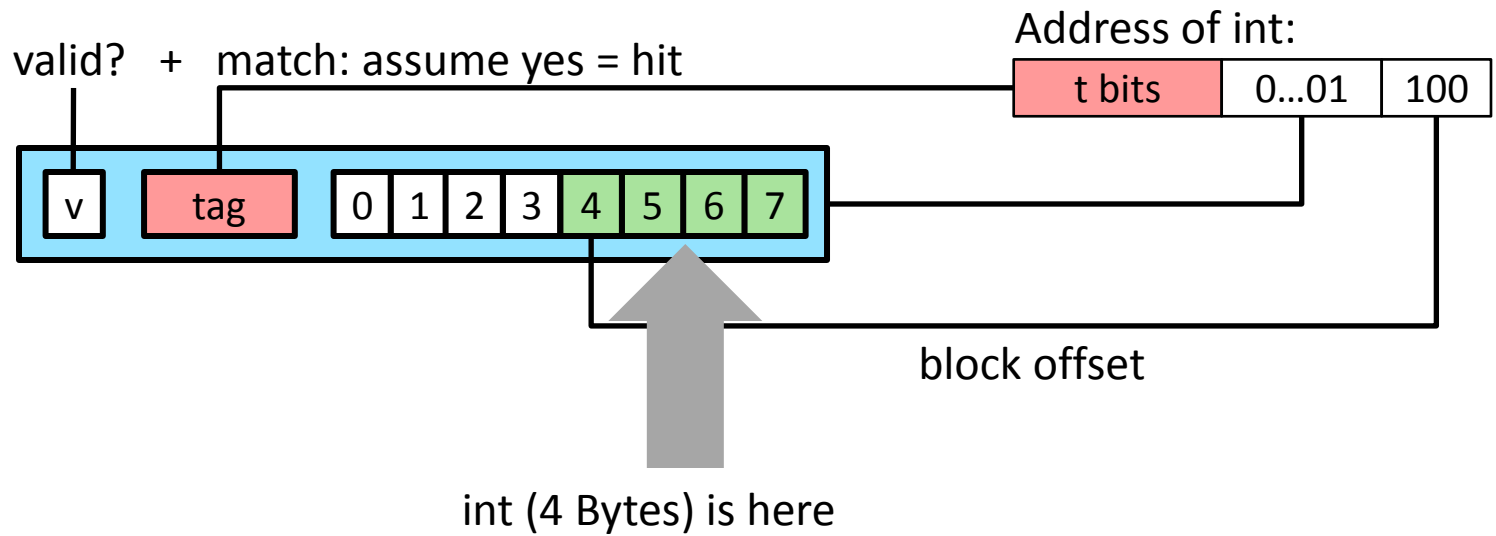Direct mapped: One line per set
Assume: cache block size 8 bytes

Address of int:

valid? + match: assume yes = hit

| t bits | 0…01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

int (4 Bytes) is here

No match: old line is evicted and replaced

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|-----|-----|-----|
| x   | xx  | x   |

M=16 byte addresses, B=2 bytes/block,
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | $[0\underline{00}0_2]$, | miss |
|---|------------------------|------|
| 1 | $[0\underline{00}1_2]$, | hit |
| 7 | $[0\underline{11}1_2]$, | miss |
| 8 | $[1\underline{00}0_2]$, | miss |
| 0 | $[0\underline{00}0_2]$ | miss |

|       | v | Tag | Block   |
|-------|---|-----|---------|
| Set 0 | 1 | 0   | M[0-1]  |
| Set 1 |   |     |         |
| Set 2 |   |     |         |
| Set 3 | 1 | 0   | M[6-7]  |

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# A Higher Level Example

assume: cold (empty) cache,
a[0][0] goes here

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

```
int sum_array_cols(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (j = 0; j < 16; j++)
        for (i = 0; i < 16; i++)
            sum += a[i][j];
    return sum;
}
```

32 B = 4 doubles

CSE 컴퓨터공학부
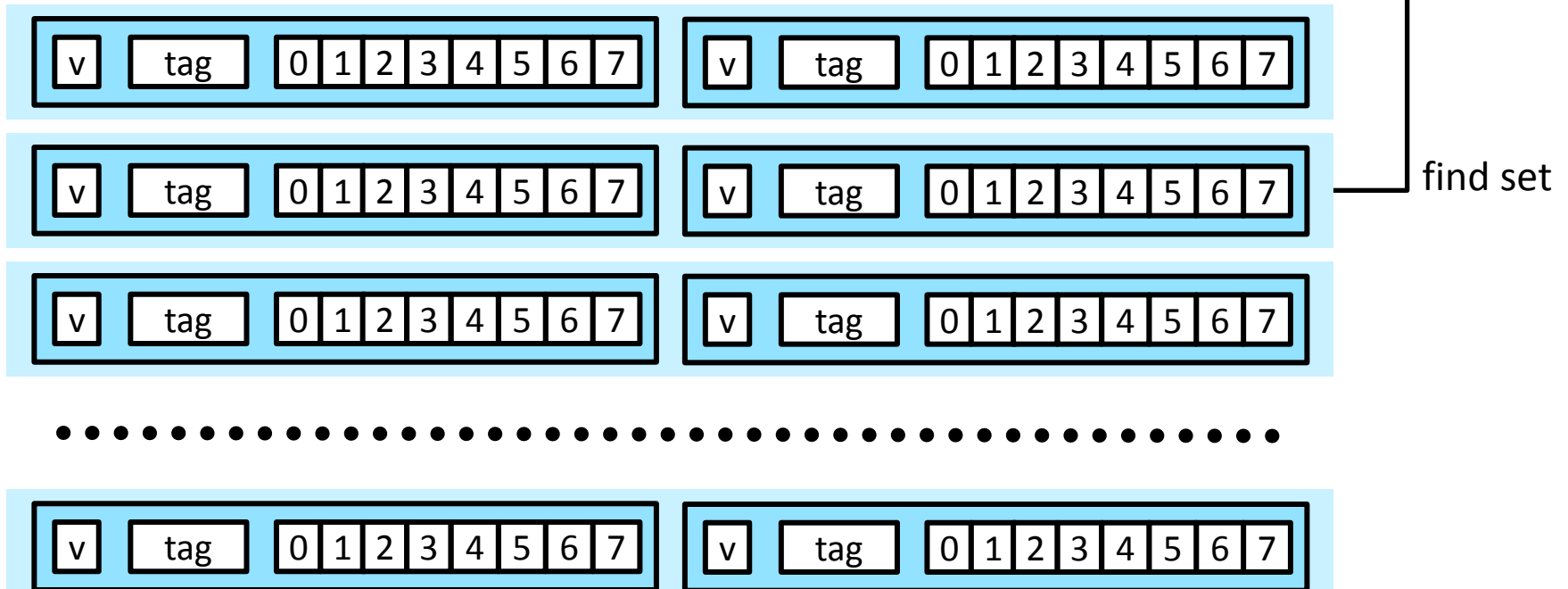Department of Computer Science & Engineering

# E-way Set Associative Cache

E = 2: Two lines per set
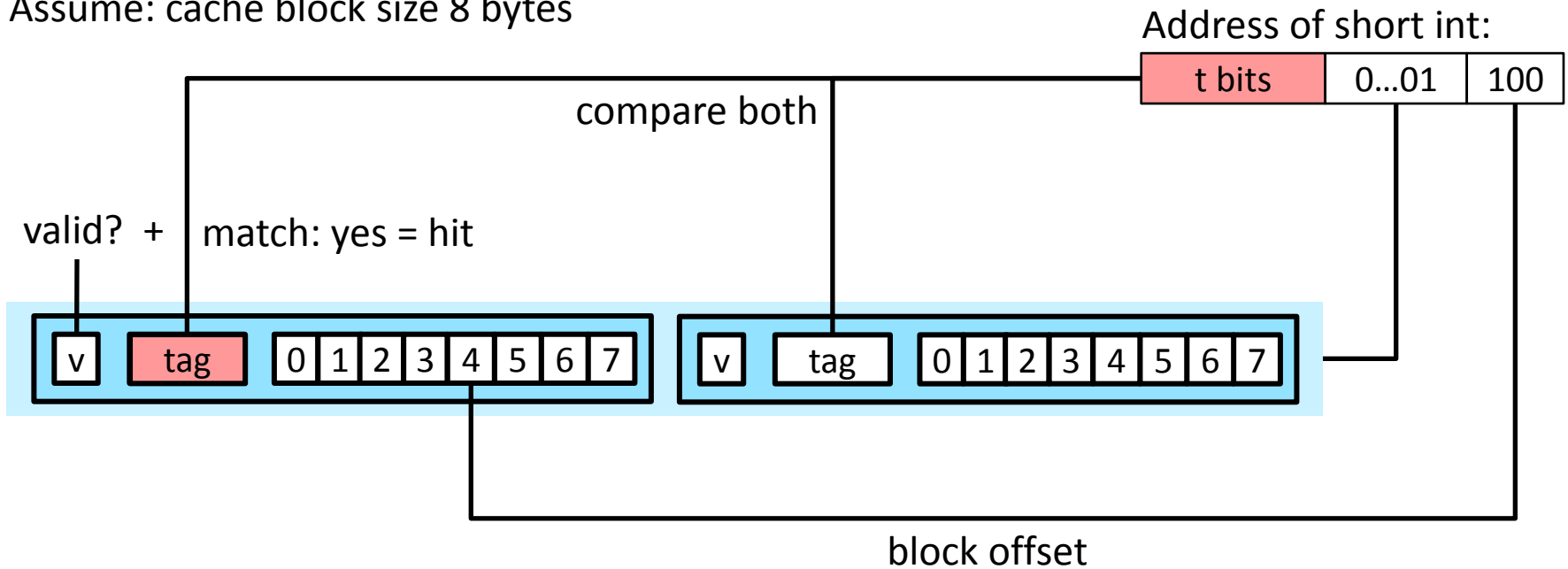Assume: cache block size 8 bytes

Address of short int:

| t bits | 0...01 | 100 |
|--------|--------|-----|

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |    | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |    | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |  find set

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |    | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

• • • • • • • • • • • • • • • • • • • • • • • • • • • • •

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |    | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# E-way Set Associative Cache

E = 2: Two lines per set
Assume: cache block size 8 bytes

Address of short int:

| t bits | 0...01 | 100 |
|---|---|---|

compare both

valid? + match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# E-way Set Associative Cache

E = 2: Two lines per set
Assume: cache block size 8 bytes

Address of short int:

| t bits | 0...01 | 100 |

compare both

valid? + match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

short int (2 Bytes) is here

No match:
- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# 2-Way Set Associative Cache Simulation

t=2   s=1   b=1

| xx | x | x |
|----|---|---|

M=16 byte addresses, B=2 bytes/block,
S=2 sets, E=2 blocks/set

Address trace (reads, one byte per read):

| 0 | [$0000_2$], | miss |
|---|------------|------|
| 1 | [$0001_2$], | hit |
| 7 | [$0111_2$], | miss |
| 8 | [$1000_2$], | miss |
| 0 | [$0000_2$] | hit |

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

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# A Higher Level Example

*Ignore the variables sum, i, j*

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

assume: cold (empty) cache,
a[0][0] goes here

32 B = 4 doubles

```
int sum_array_rows(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (j = 0; j < 16; j++)
        for (i = 0; i < 16; i++)
            sum += a[i][j];
    return sum;
}
```

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# What about writes?

- Multiple copies of data exist:
  - L1, L2, Main Memory, Disk
- What to do on a write-hit?
  - Write-through (write immediately to memory)
  - Write-back (defer write to memory until replacement of line)
    - ▸ Need a dirty bit (line different from memory or not)
- What to do on a write-miss?
  - Write-allocate (load into cache, update line in cache)
    - ▸ Good if more writes to the location follow
  - No-write-allocate (writes immediately to memory)
- Typical
  - Write-through + No-write-allocate
  - **Write-back + Write-allocate**

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# A Common Framework for Memory Hierarchies

- Question 1: Where can a Block be Placed?
  One place (direct-mapped), a few places (set associative),
  or any place (fully associative)

- Question 2: How is a Block Found?
  Indexing (direct-mapped), limited search (set associative),
  full search (fully associative)

- Question 3: Which Block is Replaced on a Miss?
  Typically LRU or random

- Question 4: How are Writes Handled?
  Write-through or write-back

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Intel Core i7 Cache Hierarchy

Processor package

Core 0

Regs

L1 d-cache

L1 i-cache

L2 unified cache

. . .

Core 3

Regs

L1 d-cache

L1 i-cache

L2 unified cache

L3 unified cache
(shared by all cores)

Main memory

L1 i-cache and d-cache:
32 KB, 8-way,
Access: 4 cycles

L2 unified cache:
256 KB, 8-way,
Access: 11 cycles

L3 unified cache:
8 MB, 16-way,
Access: 30-40
cycles

Block size: 64 bytes for
all caches.

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Cache Performance Metrics

- **Miss Rate**
  - Fraction of memory references not found in cache (misses / accesses) = 1 – hit rate
  - Typical numbers (in percentages):
    - 3-10% for L1
    - can be quite small (e.g., < 1%) for L2, depending on size, etc.
- **Hit Time**
  - Time to deliver a line in the cache to the processor
    - includes time to determine whether the line is in the cache
  - Typical numbers:
    - 1-2 clock cycle for L1
    - 5-20 clock cycles for L2
- **Miss Penalty**
  - Additional time required because of a miss
    - typically 50-200 cycles for main memory (Trend: increasing!)

# Lets think about those numbers

- Huge difference between a hit and a miss
  - Could be 100x, if just L1 and main memory

- Would you believe 99% hits is twice as good as 97%?
  - Consider:
    cache hit time of 1 cycle
    miss penalty of 100 cycles

  - Average access time:

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

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

- This is why "miss rate" is used instead of "hit rate"

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Writing Cache Friendly Code

- Make the common case go fast
  - Focus on the inner loops of the core functions

- Minimize the misses in the inner loops
  - Repeated references to variables are good (temporal locality)
  - Stride-1 reference patterns are good (spatial locality)

Key point: Our qualitative notion of locality is quantified through our understanding of cache memories.

# Cache Memories

- Cache organization and operation

- Performance impact of caches
    - The memory mountain
    - Rearranging loops to improve spatial locality

# The Memory Mountain

- **Read throughput** (read bandwidth)
    - Number of bytes read from memory per second (MB/s)

- **Memory mountain**:
  Measured read throughput as a function of spatial and temporal locality.
    - Compact way to characterize memory system performance.

# Memory Mountain Test Function

```c
/* The test function */
void test(int elems, int stride) {
    int i, result = 0;
    volatile int sink;

    for (i = 0; i < elems; i += stride)
        result += data[i];
    sink = result; /* So compiler doesn't optimize away the loop */
}

/* Run test(elems, stride) and return read throughput (MB/s) */
double run(int size, int stride, double Mhz)
{
    double cycles;
    int elems = size / sizeof(int);

    test(elems, stride);                       /* warm up the cache */
    cycles = fcyc2(test, elems, stride, 0);  /* call test(elems,stride) */
    return (size / stride) / (cycles / Mhz); /* convert cycles to MB/s */
}
```
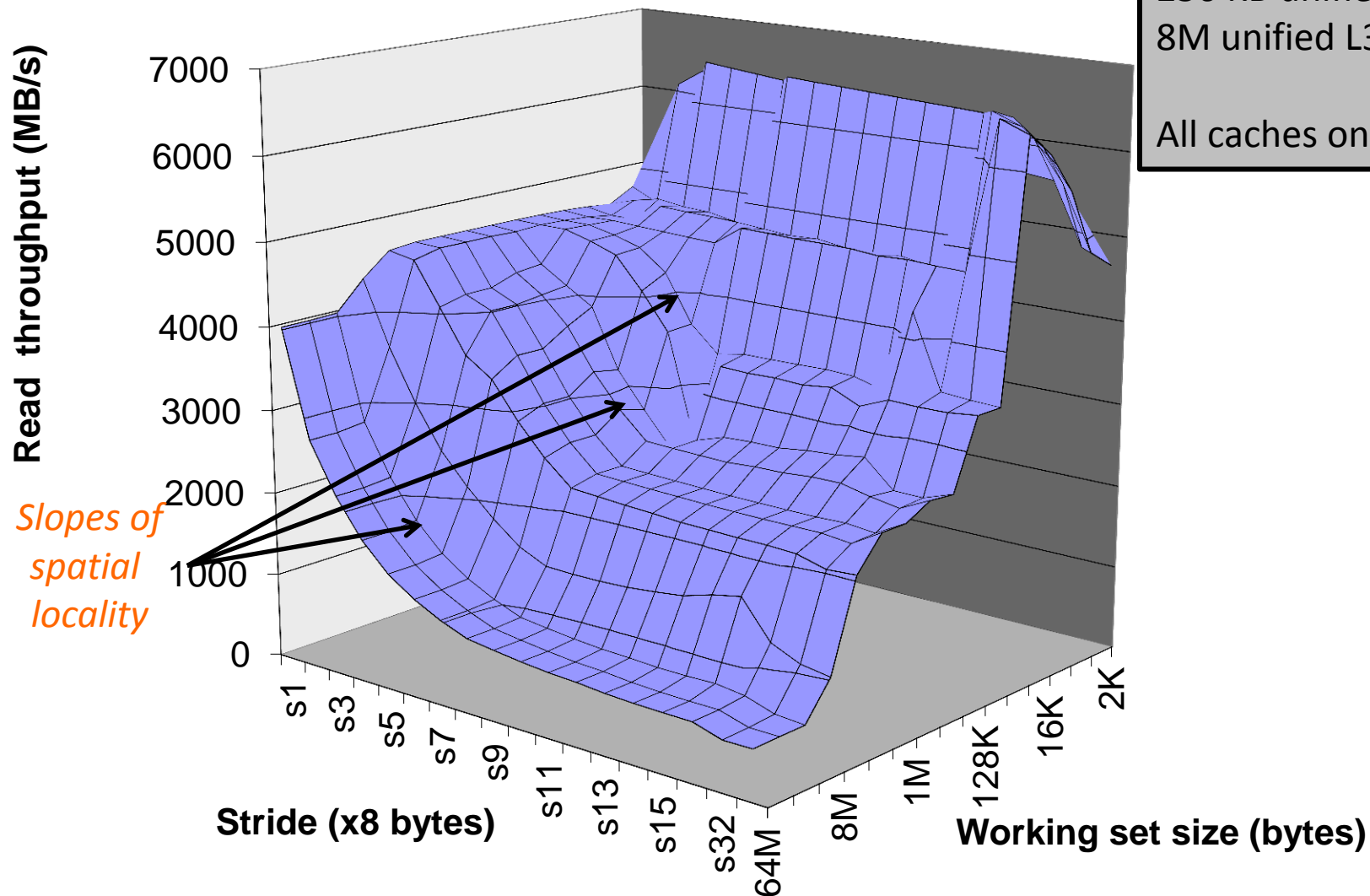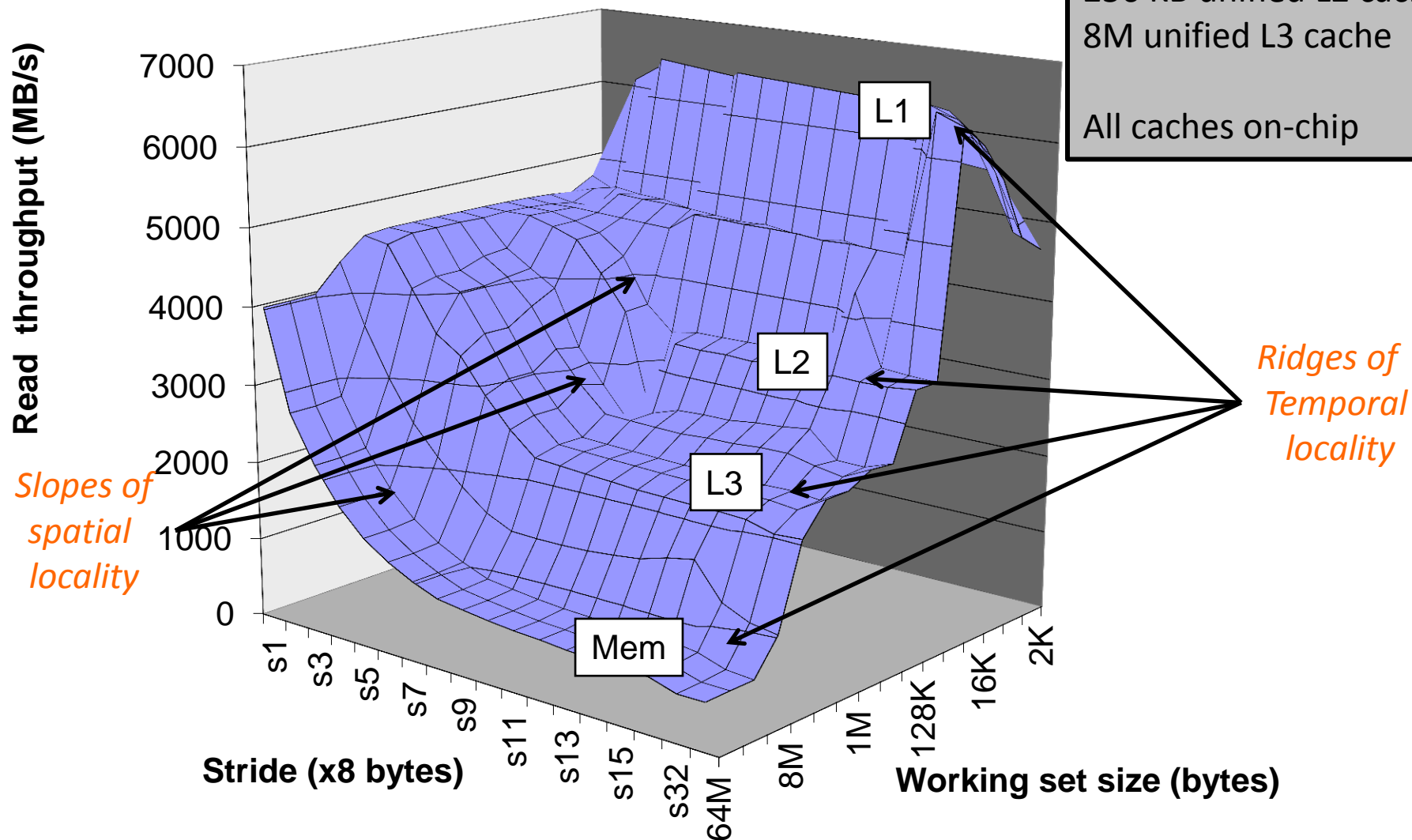
CSE 컴퓨터공학부
Department of Computer Science & Engineering

# 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



Read throughput (MB/s)

Stride (x8 bytes)

Working set size (bytes)

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# 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

**Read  throughput (MB/s)**

7000
6000
5000
4000
3000
2000
1000
0

*Slopes of spatial locality*

**Stride (x8 bytes)**

s1 s3 s5 s7 s9 s11 s13 s15 s32

**Working set size (bytes)**

64M 8M 1M 128K 16K 2K

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# 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

**Read  throughput (MB/s)**

7000

6000

5000

4000

3000

2000

1000

0

L1

L2

L3

Mem

*Ridges of Temporal locality*

*Slopes of spatial locality*

**Stride (x8 bytes)**

s1 s3 s5 s7 s9 s11 s13 s15 s32 64M

8M 1M 128K 16K 2K

**Working set size (bytes)**

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Cache Memories

- Cache organization and operation

- Performance impact of caches

  - The memory mountain

  - Rearranging loops to improve spatial locality

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Miss Rate Analysis for Matrix Multiply

- Assume:
    - Line size = 32B (big enough for four 64-bit words)
    - Matrix dimension (N) is very large
        - Approximate 1/N as 0.0
    - Cache is not even big enough to hold multiple rows
- Analysis Method:
    - Look at access pattern of inner loop

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Matrix Multiplication Example

- Description:
  - Multiply N x N matrices
  - $O(N^3)$ total operations
  - N reads per source element
  - N values summed per destination
    - but may be able to hold in register

```
/* 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;
  }
}
```

*Variable `sum` held in register*

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Layout of C Arrays in Memory (review)

- C arrays allocated in row-major order
  - each row in contiguous memory locations
- Stepping through columns in one row:
  - `for (i = 0; i < N; i++)`

    `sum += a[0][i];`
  - accesses successive elements
  - if block size (B) > 4 bytes, exploit spatial locality
    - compulsory miss rate = 4 bytes / B
- Stepping through rows in one column:
  - `for (i = 0; i < n; i++)`

    `sum += a[i][0];`
  - accesses distant elements
  - no spatial locality!
    - compulsory miss rate = 1 (i.e. 100%)

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Matrix Multiplication (ijk)

```
/* 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;
  }
}
```
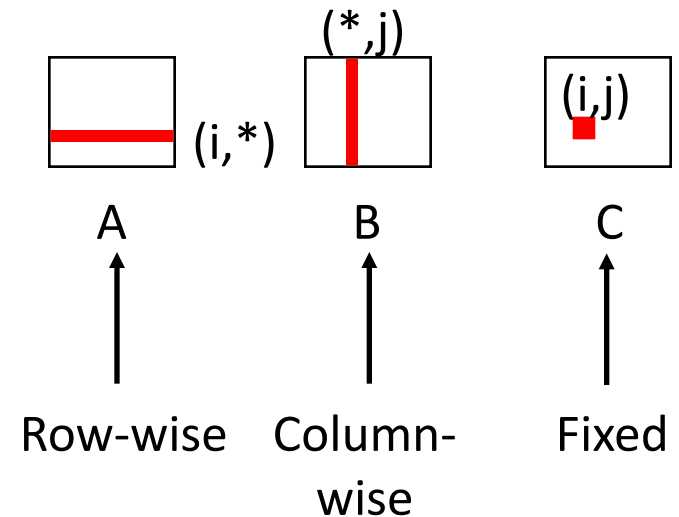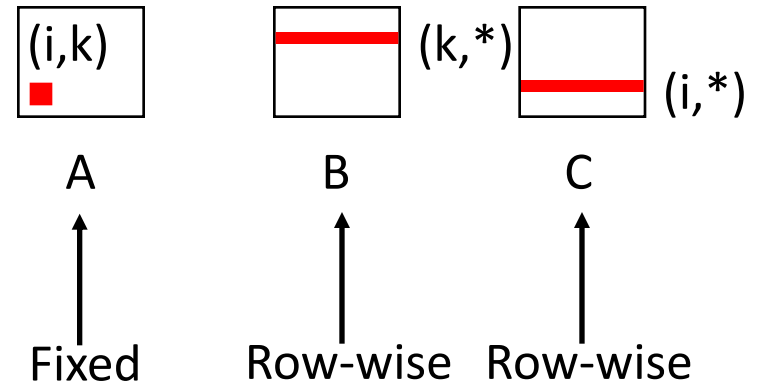
Inner loop:



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

Misses per inner loop iteration:

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

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# 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:



Misses per inner loop iteration:

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

# Matrix Multiplication (kij)

```
/* 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];
  }
}
```
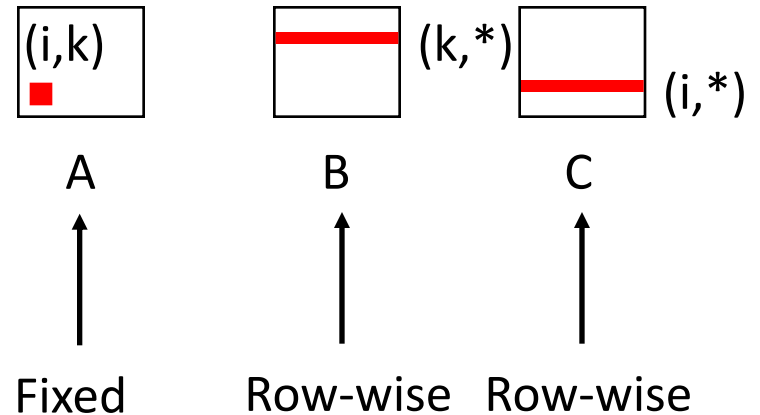
Inner loop:



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

Misses per inner loop iteration:

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

# Matrix Multiplication (ikj)

```
/* 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];
  }
}
```
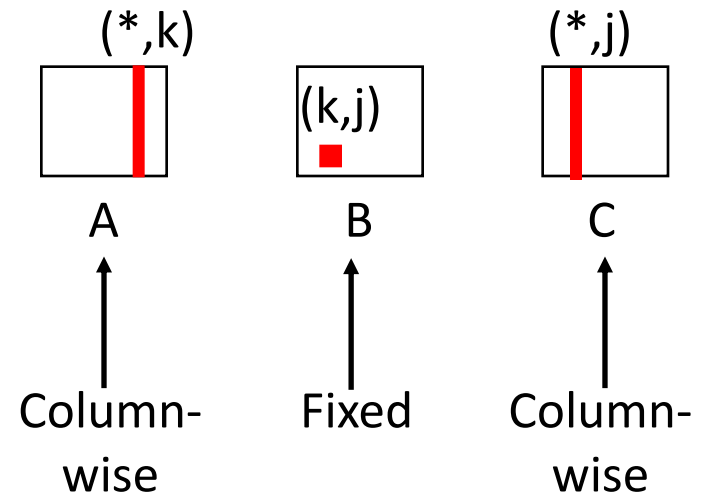
Inner loop:



| A | B | C |
|---|---|---|
| Fixed | Row-wise | Row-wise |

Misses per inner loop iteration:
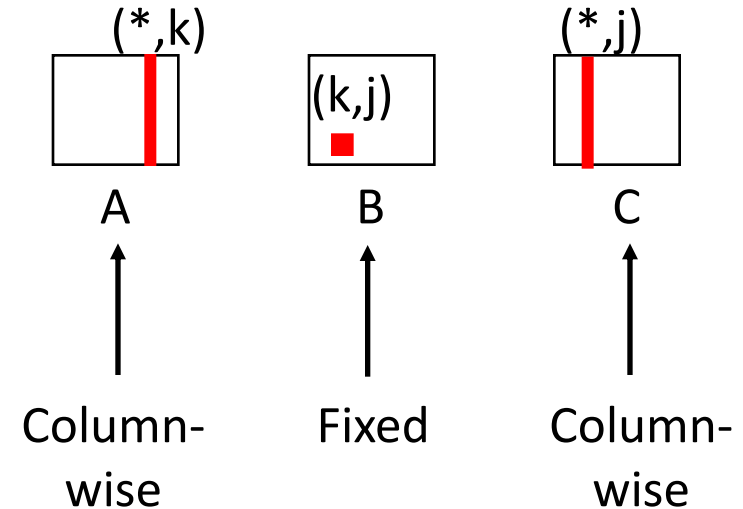
| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

# 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-      Fixed      Column-
wise                    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:



Misses per inner loop iteration:

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

# Summary of Matrix Multiplication

```
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;
  }
}
```

ijk (& jik):
- 2 loads, 0 stores
- misses/iter = 1.25

```
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];
 }
}
```
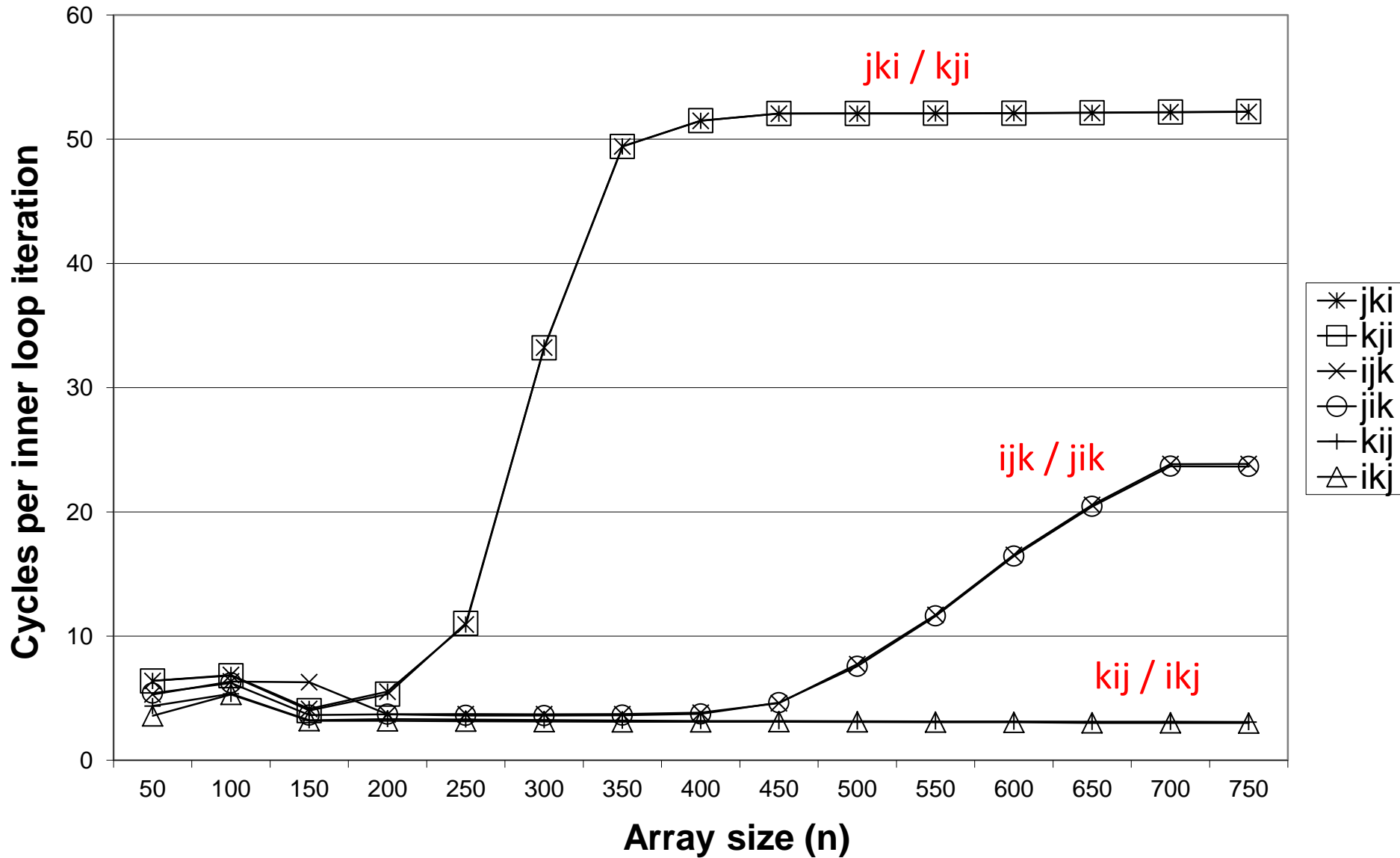
kij (& ikj):
- 2 loads, 1 store
- misses/iter = 0.5

```
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;
 }
}
```

jki (& kji):
- 2 loads, 1 store
- misses/iter = 2.0

CSE 컴퓨터공학부
Department of Computer Science & Engineering

# Core i7 Matrix Multiply Performance

# Summary

- Memory hierarchies are an optimization resulting from a perfect match between memory technology and two types of program locality
  - Temporal locality
  - Spatial locality

- The goal is to provide a "virtual" memory technology (an illusion) that has an access time of the highest-level memory with the size and cost of the lowest-level memory

- Cache memory is an instance of a memory hierarchy
  - exploits both temporal and spatial localities
  - direct-mapped caches are simple and fast but have higher miss rates
  - set-associative caches have lower miss rates but are complex and slow
  - multilevel caches are becoming increasingly popular

- Programmer can optimize for cache performance
  - How data structure are organized
  - How data are accessed (nested loop structure)