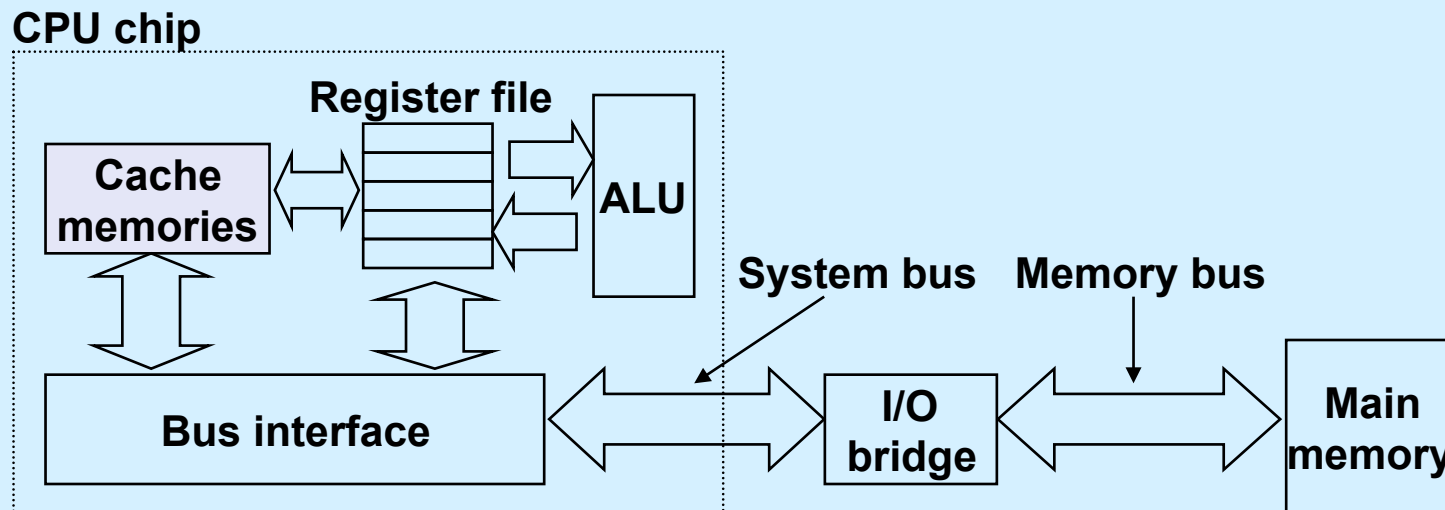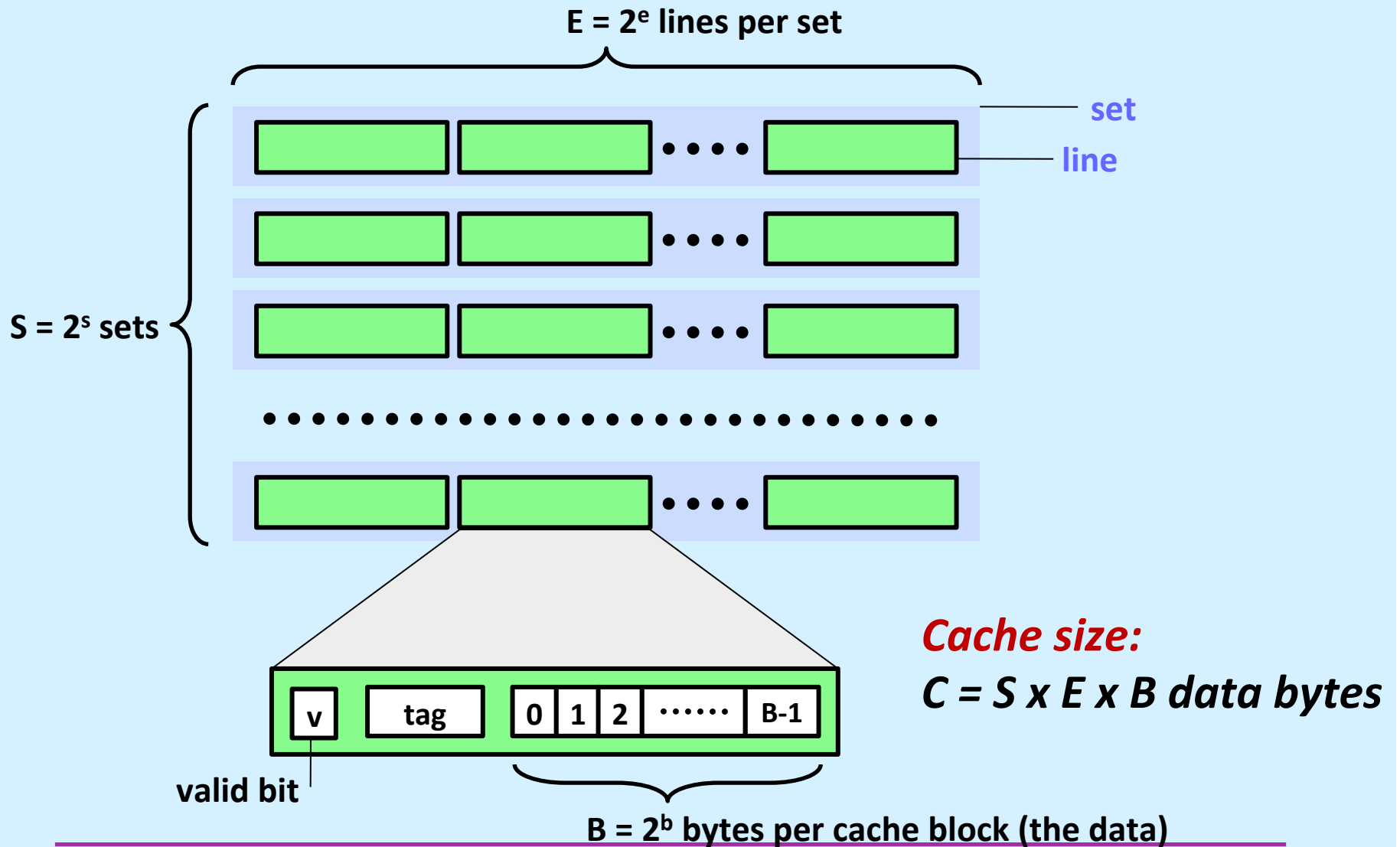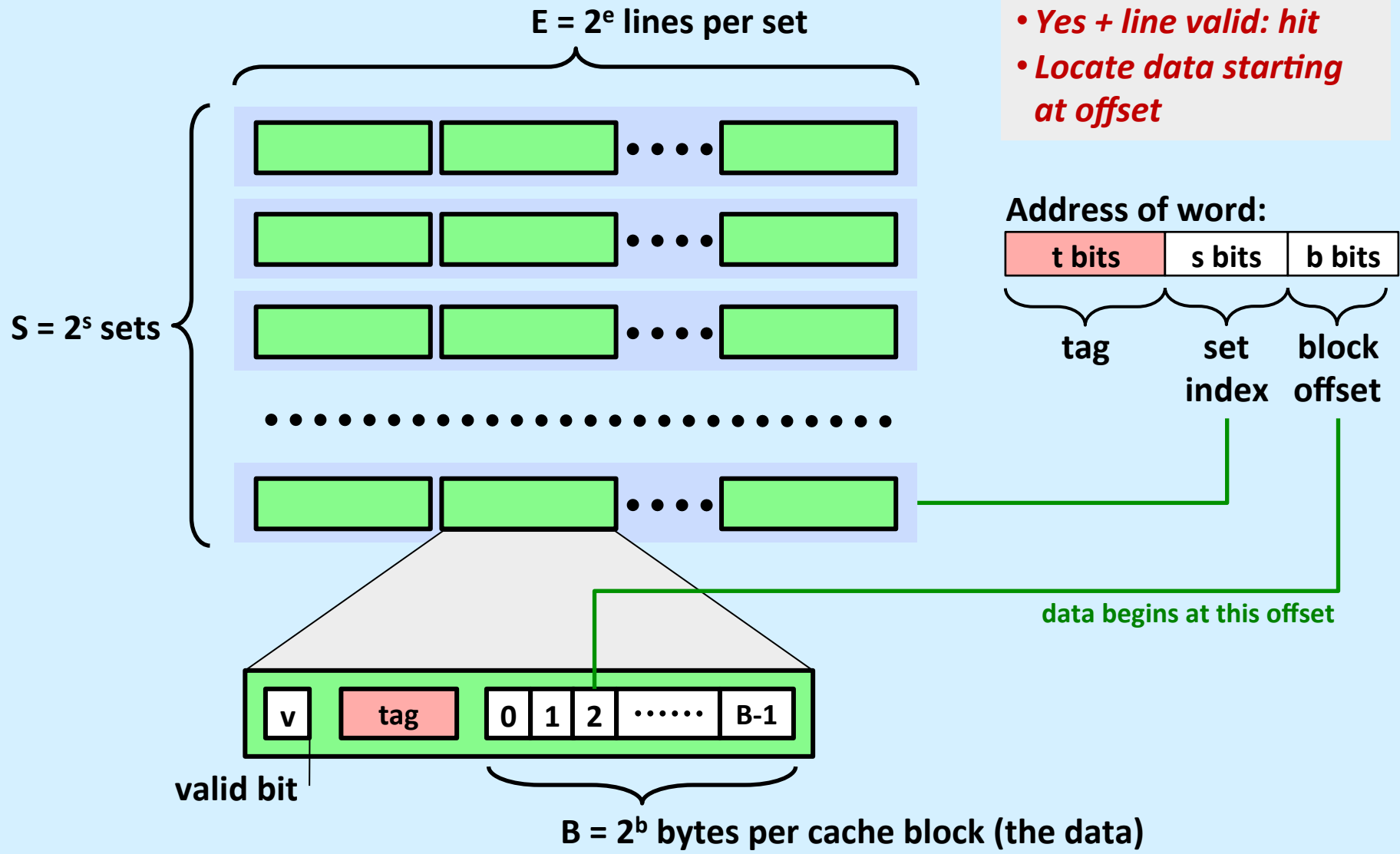# CS 33

## Caches

# 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** |
| **Bus interface** | | |

**System bus**    **Memory bus**

**I/O bridge**    **Main memory**

# 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)

# Cache Read

**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)**

# Example: Direct Mapped Cache (E = 1)

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



$S = 2^s$ sets

**Address of int:**

| t bits | 0…01 | 100 |
|--------|------|-----|

**find set**

# 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**

# 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**

**int (4 Bytes) is here**

**No match: old line is evicted and replaced**

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

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

# 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;

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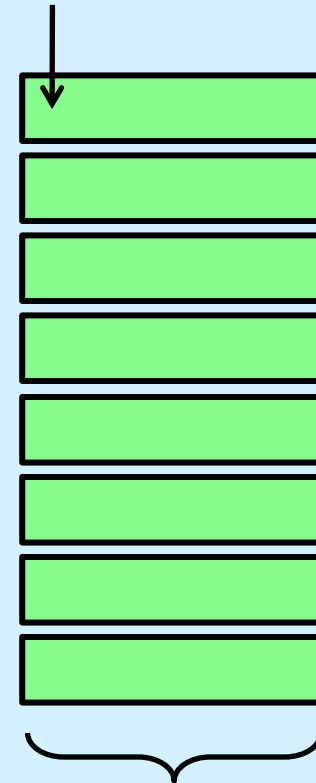
```
int sum_array_cols(double a[16][16])
{
    int i, j;
    double sum = 0;

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

}
```

**32 B = 4 doubles**

# A Higher-Level Example

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

| $a_{0,0}$ | $a_{0,1}$ | $a_{0,2}$ | $a_{0,3}$ |
|---|---|---|---|
| $a_{0,4}$ | $a_{0,5}$ | $a_{0,6}$ | $a_{0,7}$ |
| $a_{0,8}$ | $a_{0,9}$ | $a_{0,10}$ | $a_{0,11}$ |
| $a_{0,12}$ | $a_{0,13}$ | $a_{0,14}$ | $a_{0,15}$ |
| $a_{1,0}$ | $a_{1,1}$ | $a_{1,2}$ | $a_{1,3}$ |
| $a_{1,4}$ | $a_{1,5}$ | $a_{1,6}$ | $a_{1,7}$ |
| $a_{1,8}$ | $a_{1,9}$ | $a_{1,10}$ | $a_{1,11}$ |
| $a_{1,12}$ | $a_{1,13}$ | $a_{1,14}$ | $a_{1,15}$ |

**32 B = 4 doubles**

# A Higher-Level Example

```
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; i < 16; i++)
        for (i = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```
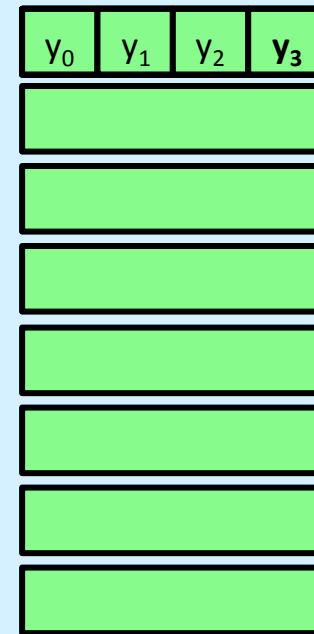
| $a_{2,0}$ | $a_{2,1}$ | $a_{2,2}$ | $a_{2,3}$ |
|-----------|-----------|-----------|-----------|

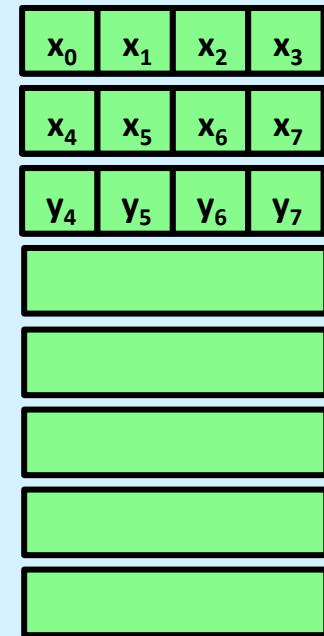| $a_{3,0}$ | $a_{3,1}$ | $a_{3,2}$ | $a_{3,3}$ |
|-----------|-----------|-----------|-----------|

**32 B = 4 doubles**

# Conflict Misses: Aligned

```
double dotprod(double x[8], double y[8]) {
  double sum = 0.0;
  int i;

  for (i=0; i<8; i++)
    sum += x[i] * y[i];

  return sum;
}
```

| $y_0$ | $y_1$ | $y_2$ | $y_3$ |
|---|---|---|---|

**32 B = 4 doubles**

# Different Alignments

```
double dotprod(double x[8], double y[8]) {
  double sum = 0.0;
  int i;

  for (i=0; i<8; i++)
    sum += x[i] * y[i];

  return sum;
}
```
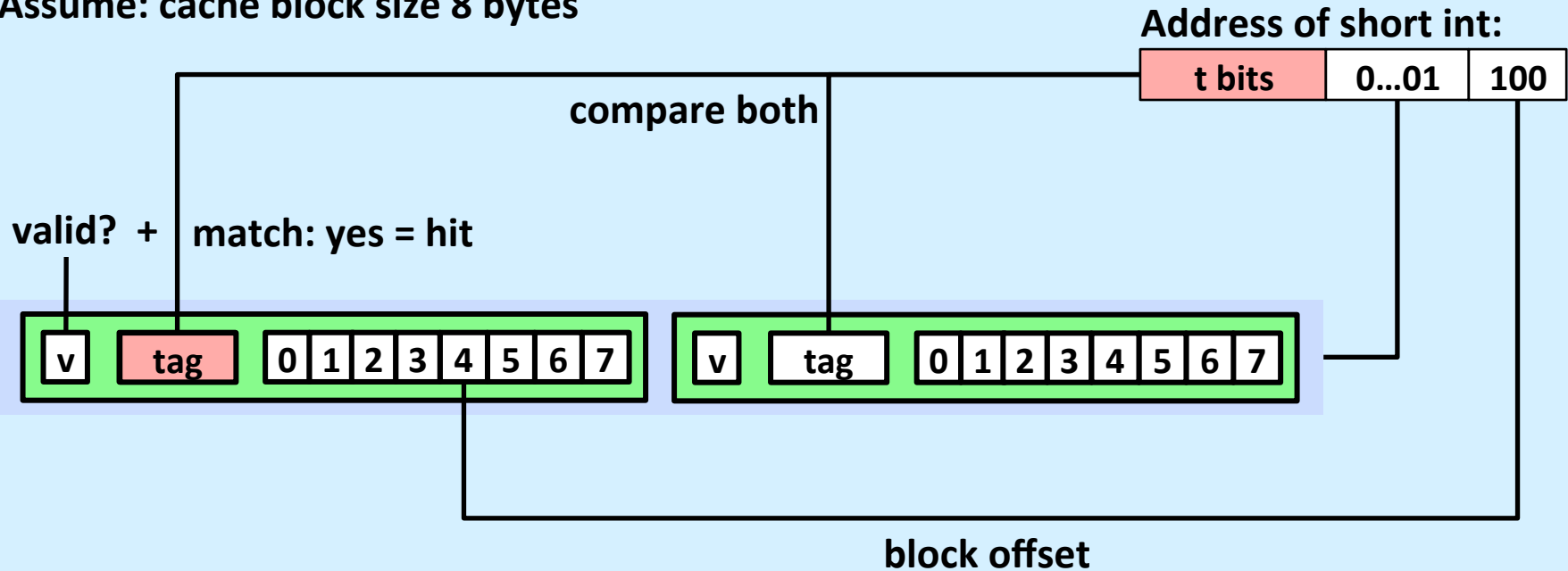
| $x_0$ | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $x_4$ | $x_5$ | $x_6$ | $x_7$ |
| $y_4$ | $y_5$ | $y_6$ | $y_7$ |

**32 B = 4 doubles**

# E-way Set-Associative Cache (Here: E = 2)

**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**

# E-way Set-Associative Cache (Here: E = 2)

**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), ...**

# Quiz 1

| 0 | v | tag=0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |   | v | tag=2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 |
|---|---|-------|---|---|---|---|---|---|---|---|---|---|-------|---|---|---|---|---|---|---|---|
| 1 | v | tag=0 | 4 | 4 | 4 | 4 | 5 | 5 | 5 | 5 |   | v | tag=4 | 6 | 6 | 6 | 6 | 7 | 7 | 7 | 7 |
| 2 | v | tag=2 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 9 |   | v | tag=3 | a | a | a | a | b | b | b | b |
| 3 | v | tag=4 | c | c | c | c | d | d | d | d |   | v | tag=a | e | e | e | e | f | f | f | f |

**Given the address above and the cache contents as shown, what is the value of the *int* at the given address?**

    a) 1111

    b) 3333

    c) 4444

    d) 7777

# 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 | [00$\underline{0}$0$_2$], | miss |
| 1 | [00$\underline{0}$1$_2$], | hit |
| 7 | [01$\underline{1}$1$_2$], | miss |
| 8 | [10$\underline{0}$0$_2$], | miss |
| 0 | [00$\underline{0}$0$_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 |    |         |

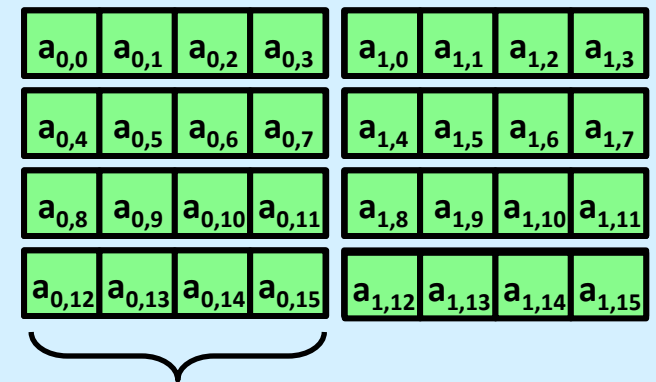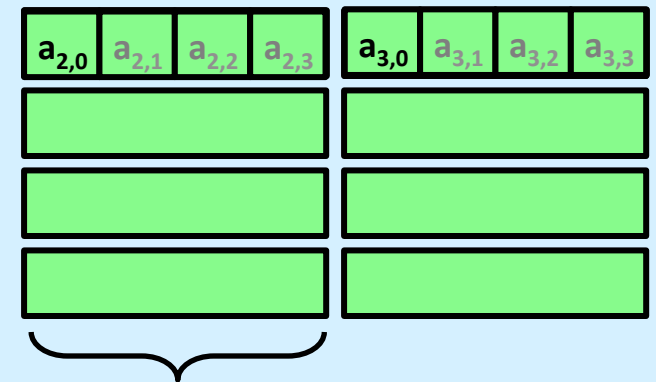# A Higher-Level Example

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

# A Higher-Level Example

*Ignore the variables sum, i, j*

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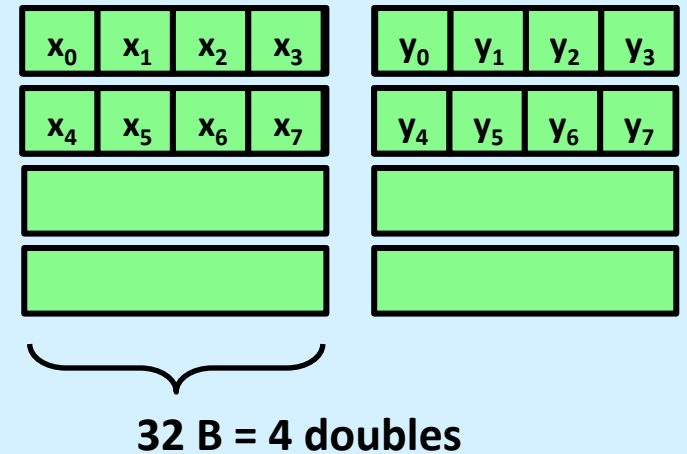
```c
int sum_array_cols(double a[16][16])
{
    int i, j;
    double sum = 0;

    for (j = 0; i < 16; i++)
        for (i = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```
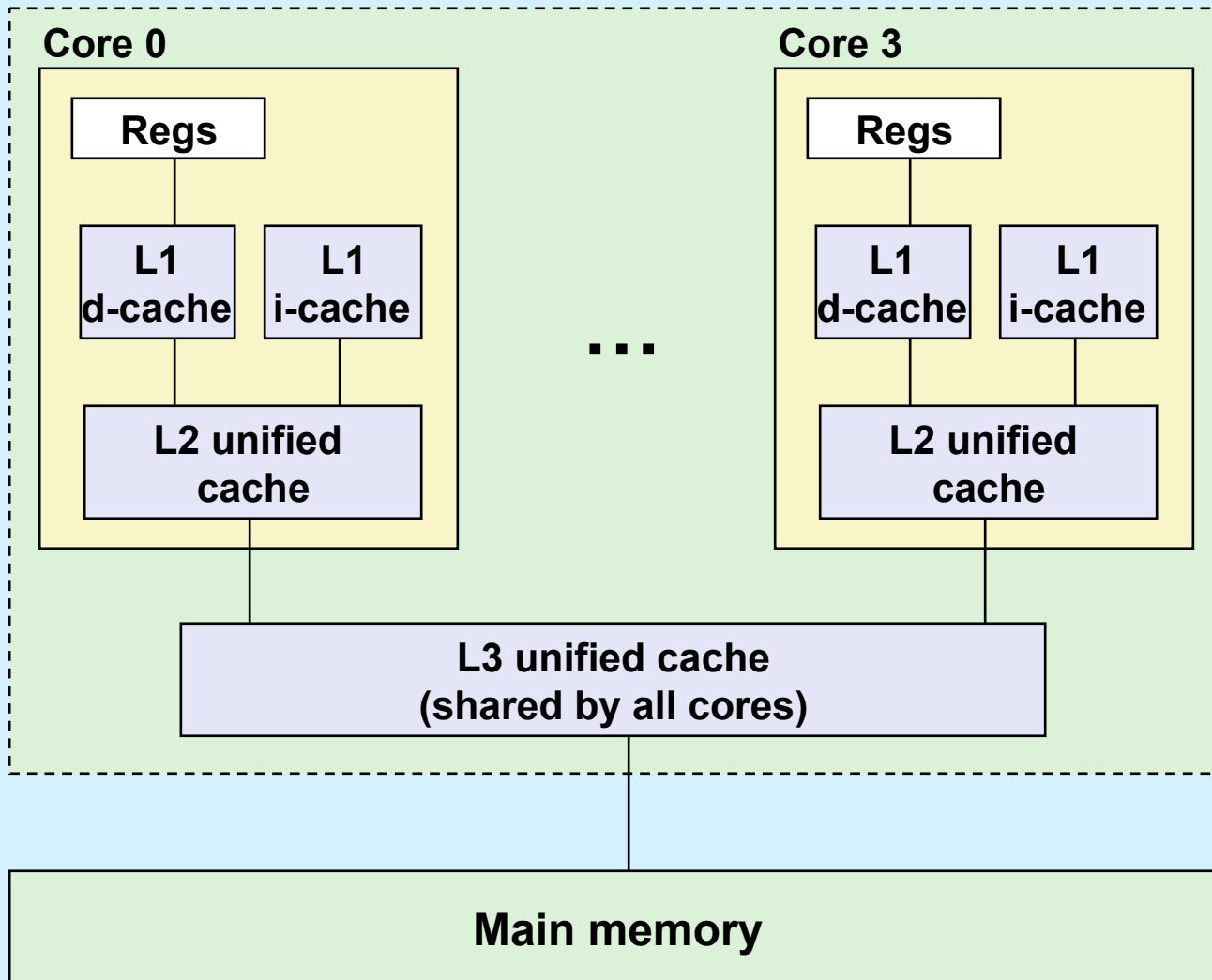
| $a_{0,0}$ | $a_{0,1}$ | $a_{0,2}$ | $a_{0,3}$ | $a_{1,0}$ | $a_{1,1}$ | $a_{1,2}$ | $a_{1,3}$ |
|---|---|---|---|---|---|---|---|
| $a_{0,4}$ | $a_{0,5}$ | $a_{0,6}$ | $a_{0,7}$ | $a_{1,4}$ | $a_{1,5}$ | $a_{1,6}$ | $a_{1,7}$ |
| $a_{0,8}$ | $a_{0,9}$ | $a_{0,10}$ | $a_{0,11}$ | $a_{1,8}$ | $a_{1,9}$ | $a_{1,10}$ | $a_{1,11}$ |
| $a_{0,12}$ | $a_{0,13}$ | $a_{0,14}$ | $a_{0,15}$ | $a_{1,12}$ | $a_{1,13}$ | $a_{1,14}$ | $a_{1,15}$ |

**32 B = 4 doubles**

# 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;
}
```

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

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

| $a_{2,0}$ | $a_{2,1}$ | $a_{2,2}$ | $a_{2,3}$ | $a_{3,0}$ | $a_{3,1}$ | $a_{3,2}$ | $a_{3,3}$ |
|---|---|---|---|---|---|---|---|

**32 B = 4 doubles**

# Conflict Misses

```
double dotprod(double x[8], double y[8]) {
  double sum = 0.0;
  int i;

  for (i=0; i<8; i++)
    sum += x[i] * y[i];

  return sum;
}
```

| $x_0$ | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|
| $x_4$ | $x_5$ | $x_6$ | $x_7$ |
| | | | |
| | | | |

| $y_0$ | $y_1$ | $y_2$ | $y_3$ |
|---|---|---|---|
| $y_4$ | $y_5$ | $y_6$ | $y_7$ |
| | | | |
| | | | |

**32 B = 4 doubles**

# 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

# 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**

# 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 cycles 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!)

# Let's 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% hit rate is twice as good as 97%?**
  - consider:
    cache hit time of 1 cycle
    miss penalty of 100 cycles
  - average access time:

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

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

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

# 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 idea: our qualitative notion of locality is quantified through our understanding of cache memories**

# Miss-Rate Analysis for Matrix Multiply

- ## Assume:
  - Block size = 32B (big enough for four 64-bit words)
  - matrix dimension (N) is very large
    - » approximate 1/N as 0.0
  - cache is not big enough to hold multiple rows

- ## Analysis method:
  - look at access pattern of inner loop



C   =   A   *   B

# Matrix Multiplication Example

- **Description:**
  - **multiply N x N matrices**
  - **O(N³) 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;
  }
}
```
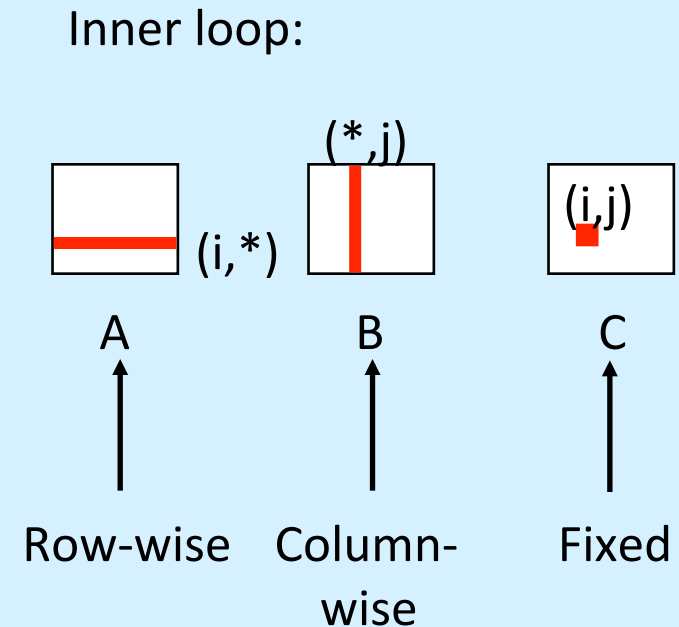
*Variable sum held in register*

# 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%)**

# 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;
  }
}
```
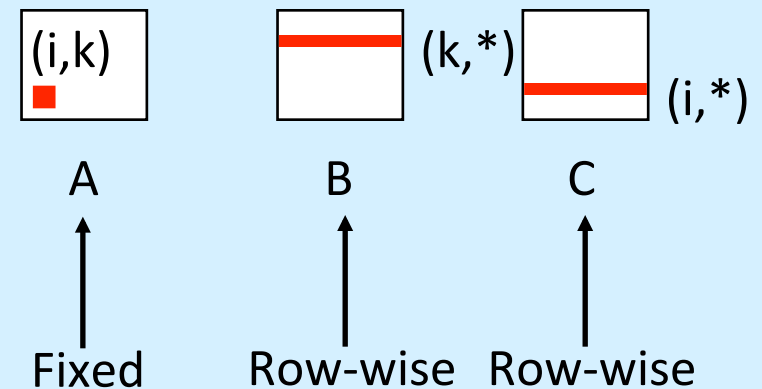
Inner loop:



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

Misses per inner loop iteration:

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

# 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:

(*,j)

(i,*)

(i,j)

A          B          C

Row-wise   Column-   Fixed
           wise

Misses per inner loop iteration:
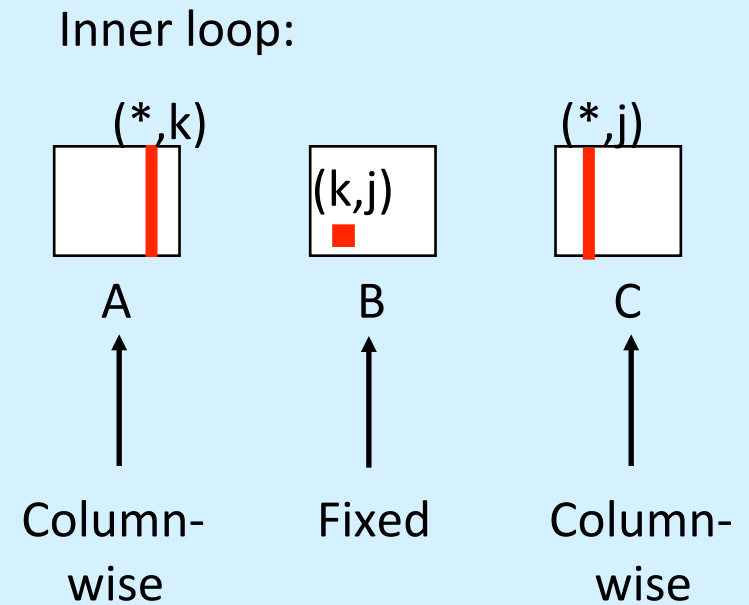
| A | B | C |
|---|---|---|
| 0.25 | 1.0 | 0.0 |

# 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:



(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 (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:

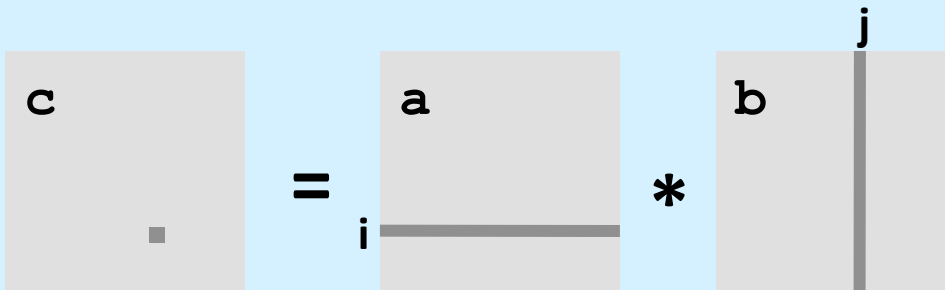(*,k)          (k,j)          (*,j)

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:

(*,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 |

# 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**

# Core i7 Matrix Multiply Performance



Cycles per inner loop iteration vs. Array size (n)

Legend:
- jki (*)
- kji (□)
- ijk (×)
- jik (○)
- kij (+)
- ikj (△)

jki / kji

ijk / jik

kij / ikj

# Matrix Multiplication: More Analysis

```
/* Multiply n x n matrices a and b  */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i++)
        for (j = 0; j < n; j++)
            for (k = 0; k < n; k++)
                c[i*n+j] += a[i*n + k]*b[k*n + j];
}
```
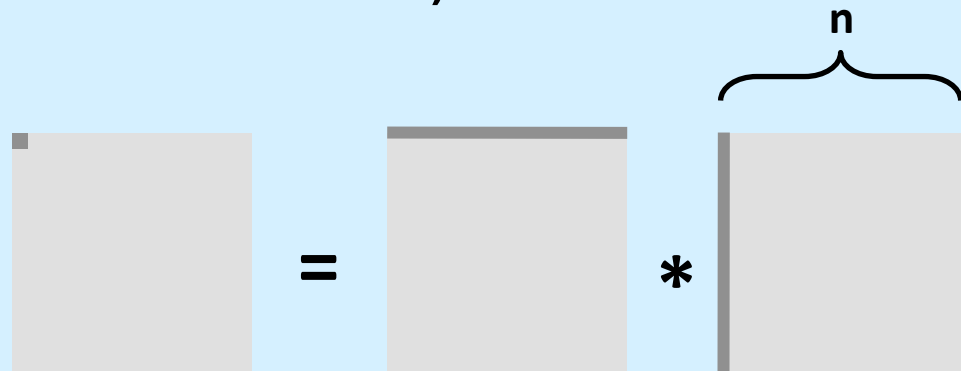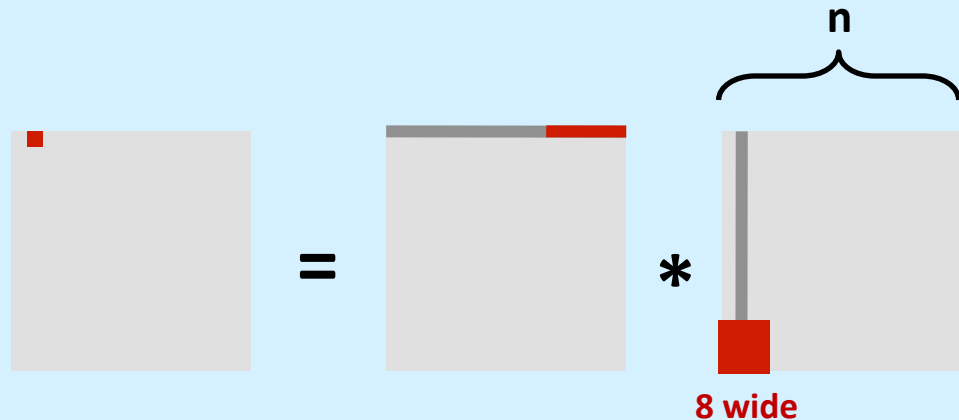
# Cache-Miss Analysis

- ## Assume:
    - matrix elements are doubles
    - cache block = 8 doubles
    - cache size C << n (much smaller than n)

- ## First iteration:
    - n/8 + n = 9n/8 misses

    - afterwards in cache: (schematic)



8 wide

# Cache-Miss Analysis

- ## Assume:
  - matrix elements are doubles
  - cache block = 8 doubles
  - cache size C << n (much smaller than n)

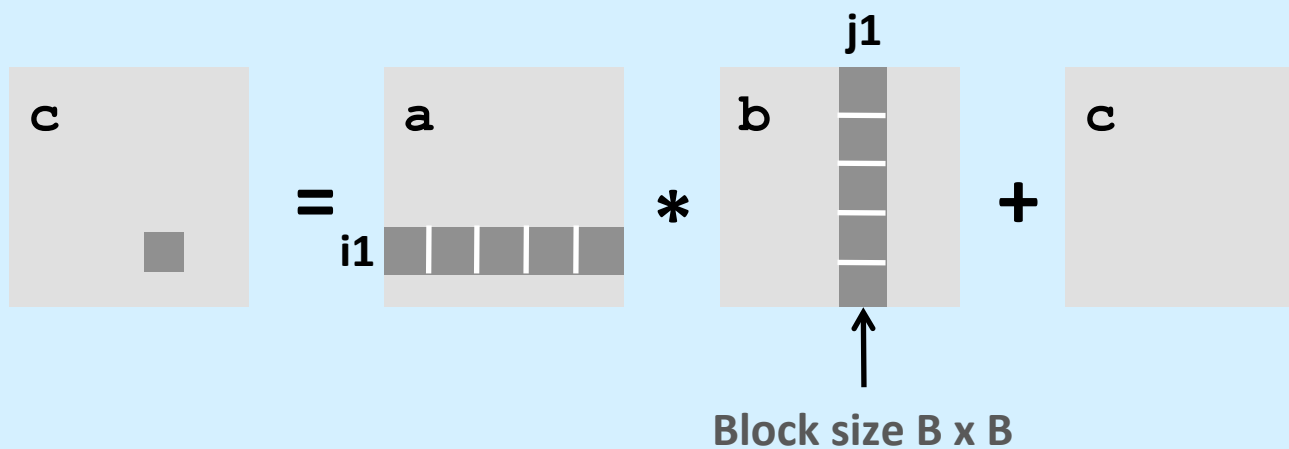- ## Second iteration:
  - again:
    n/8 + n = 9n/8 misses



n

= * 

8 wide

- ## Total misses:
  - $9n/8 * n^2 = (9/8) * n^3$

# Blocked Matrix Multiplication

```
/* Multiply n x n matrices a and b  */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=B)
        for (j = 0; j < n; j+=B)
            for (k = 0; k < n; k+=B)
                /* B x B mini matrix multiplications */
                for (i1 = i; i1 < i+B; i++)
                    for (j1 = j; j1 < j+B; j++)
                        for (k1 = k; k1 < k+B; k++)
                            c[i1*n+j1] += a[i1*n + k1]*b[k1*n + j1];
}
```
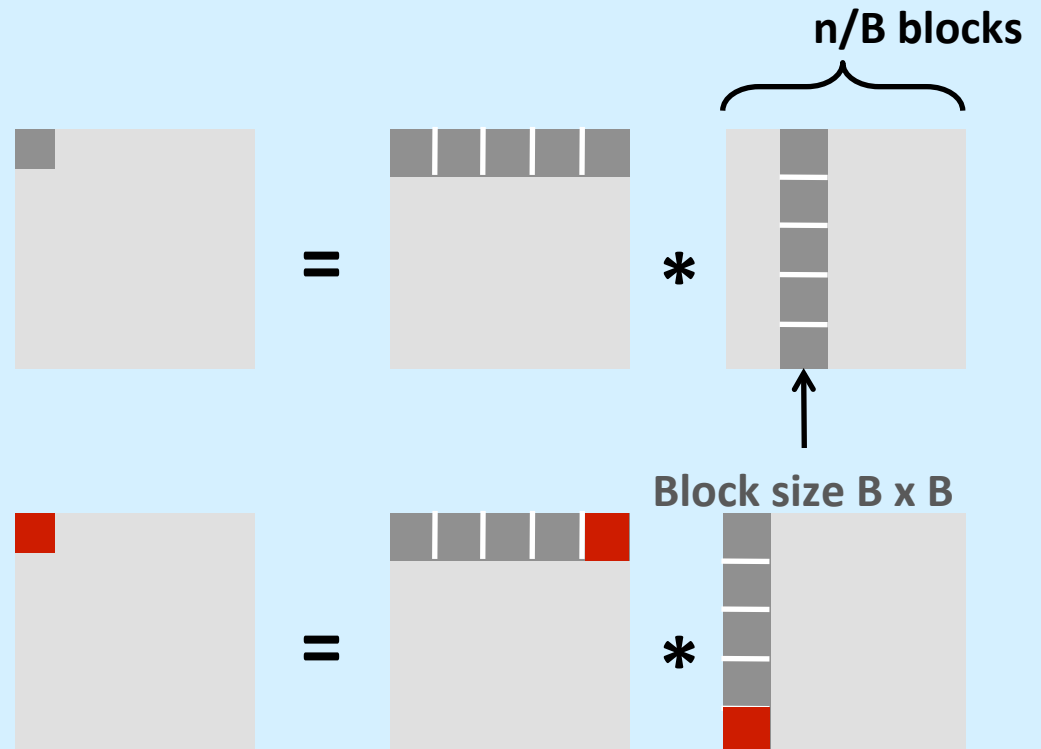


**j1**

c     =     a     *     b     +     c

i1

Block size B x B

# Cache-Miss Analysis

- ## Assume:
  - cache block = 8 doubles
  - cache size $C \ll n$ (much smaller than n)
  - three blocks ⬛ fit into cache: $3B^2 < C$

- ## First (block) iteration:
  - $B^2/8$ misses for each block
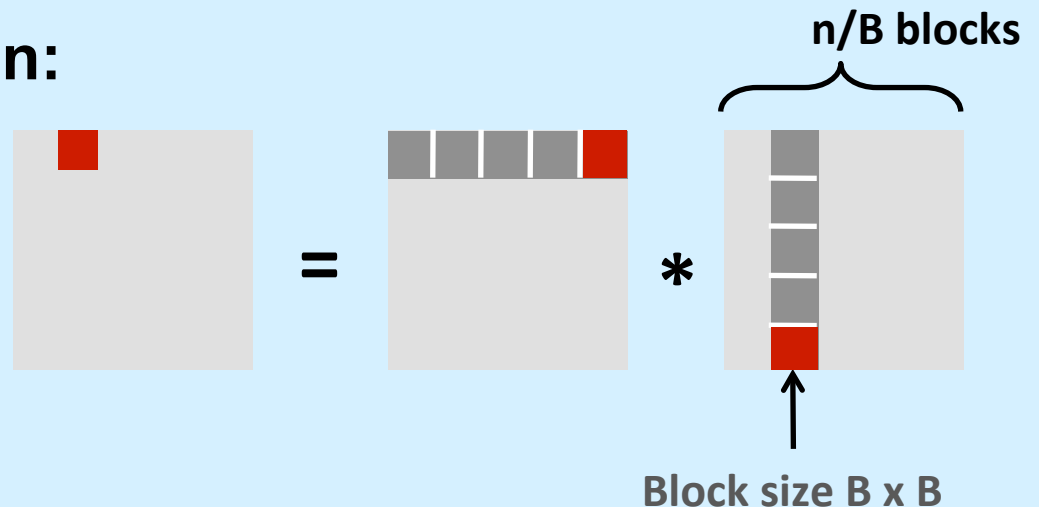  - $2n/B * B^2/8 = nB/4$ (omitting matrix c)

  n/B blocks

  ⬛ = ⬛⬛⬛⬛⬛ * ⬛

  Block size B x B

  - afterwards in cache (schematic)

  ⬛ = ⬛⬛⬛⬛⬛ * ⬛

# Cache-Miss Analysis

- ## Assume:
  - cache block = 8 doubles
  - cache size C << n (much smaller than n)
  - three blocks ▪ fit into cache: $3B^2 < C$

- ## Second (block) iteration:
  - same as first iteration
  - $2n/B * B^2/8 = nB/4$

**n/B blocks**

= * 

**Block size B x B**

- ## Total misses:
  - $nB/4 * (n/B)^2 = n^3/(4B)$

# Summary

- **No blocking: $(9/8) * n^3$**
- **Blocking: $1/(4B) * n^3$**

- **Suggest largest possible block size B, but limit $3B^2 < C$!**

- **Reason for dramatic difference:**
  - matrix multiplication has inherent temporal locality:
    » input data: $3n^2$, computation $2n^3$
    » every array element used $O(n)$ times!
  - but program has to be written properly

# Quiz 2

What is the smallest value of B (in 8-byte doubles) for which the cache-miss analysis works?

a) 1

b) 2

c) 4

d) 8

# Concluding Observations

- **Programmer can optimize for cache performance**
  - how data structures are organized
  - how data are accessed
    - » nested loop structure
    - » blocking is a general technique

- **All systems favor "cache-friendly code"**
  - getting absolute optimum performance is very platform specific
    - » cache sizes, line sizes, associativities, etc.
  - can get most of the advantage with generic code
    - » keep working set reasonably small (temporal locality)
    - » use small strides (spatial locality)