# Cache memories

**Slides by:** Randal E. Bryant and David R. O'Hallaron (CMU)

Presented by Xin Jin and David Hovemeyer for CSF

October 15, 2021

601.229 Computer Systems Fundamentals

# Cache writes and performance

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, 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 straight to memory, does not load into cache)
- **Typical**
  - Write-through + No-write-allocate
  - **Write-back + Write-allocate**

# Zoom poll #1!

Consider the following code:

```
for (int i = 0; i < 8; i++) {
  a[i] = i * 2;
}
```

Assume that

► address of a[0] is a multiple of 16
► cache is cold initially
► i is a register
► sizeof(int)=4
► 16 bytes per block
► cache is direct-mapped
► loads and stores always access exactly 4 bytes

If there are 8 stores to *memory*, what cache configuration is likely?

A. write-allocate + write-through
B. no-write-allocate + write-through
C. write-allocate + write-back
D. no-write-allocate + write-back

# Zoom poll #2!

Consider the following code:

```
for (int i = 0; i < 8; i++) {
  a[i] = i * 2;
}
```

Assume that

► address of a[0] is a multiple of 16
► cache is cold initially
► i is a register
► sizeof(int)=4
► 16 bytes per block
► cache is direct-mapped
► loads and stores always access exactly 4 bytes

If the cache is configured for write-allocate + write-back, how many loads from *memory* are there?

A. 0
B. 2
C. 8
D. 10
E. 16

# Zoom poll #3!

Consider the following code:

```
for (int i = 0; i < 8; i++) {
  a[i] = i * 2;
}
```
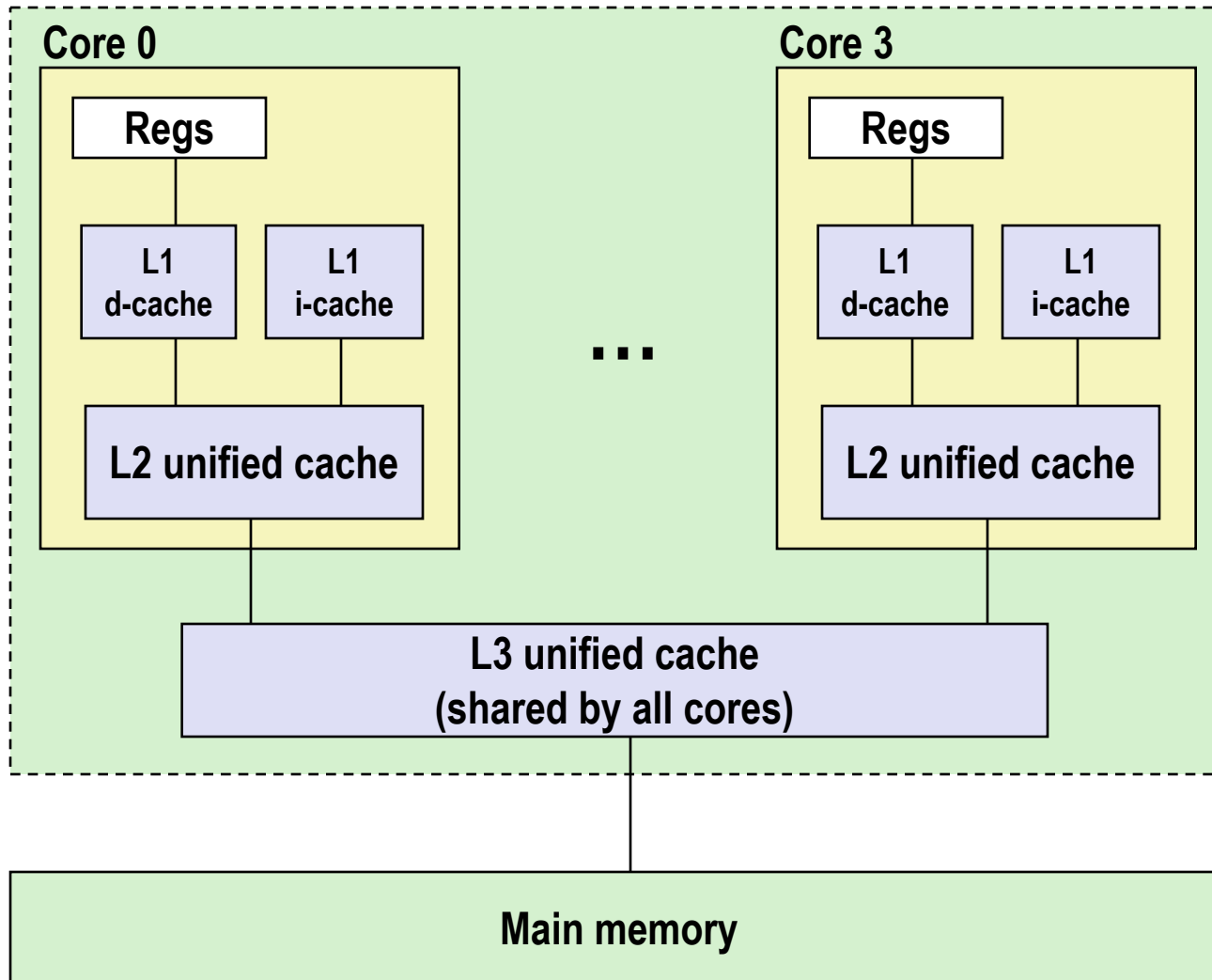
Assume that

- ▶ address of a[0] is a multiple of 16
- ▶ cache is cold initially
- ▶ i is a register
- ▶ sizeof(int)=4
- ▶ 16 bytes per block
- ▶ cache is direct-mapped
- ▶ loads and stores always access exactly 4 bytes

If the cache is configured for write-allocate + write-back, how many stores to *memory* are there?

A. 0
B. 2
C. 8
D. 10
E. 16

# 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: 10 cycles

**L3 unified cache:**
8 MB, 16-way,
Access: 40-75 cycles

**Block size**: 64 bytes for all caches.

# 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:
    - 4 clock cycle for L1
    - 10 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% 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"**

# Writing cache-friendly code

# 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 (<span style="color:red">temporal locality</span>)

  ▶ Stride-1 reference patterns are good (<span style="color:red">spatial locality</span>)

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

# Matrix Multiplication Example

► **Description:**

  ► Multiply N x N matrices

  ► Matrix elements are doubles (8 bytes)

  ► $O(N^3)$ total operations

  ► N reads per source element

  ► N values summed per destination

    ► but may be able to hold in register

*Variable* `sum` *held 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;
  }
}
```
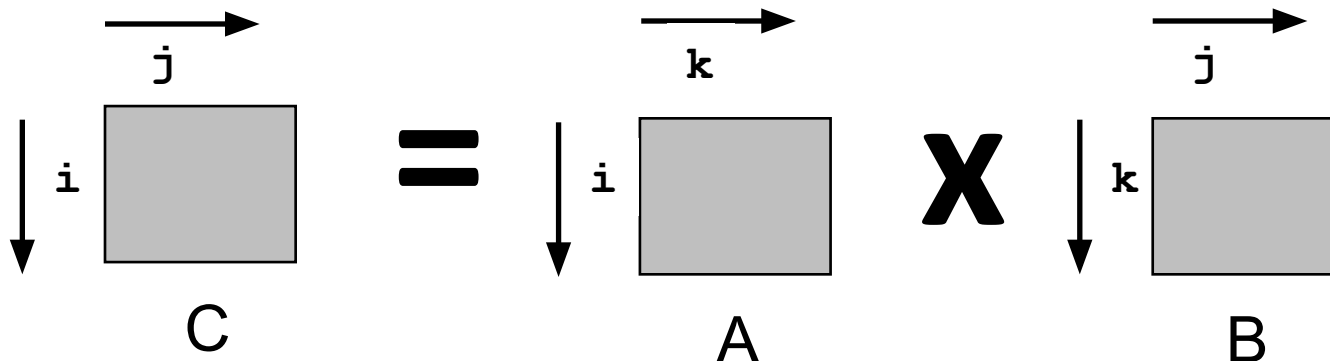*matmult/mm.c*

# Miss Rate Analysis for Matrix Multiply

► **Assume:**

  ► Block size = 32B (big enough for four doubles)

  ► 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

# 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) > sizeof($a_{ij}$) bytes, exploit spatial locality
    - miss rate = sizeof($a_{ij}$) / B
- **Stepping through rows in one column:**
  - ```
    for (i = 0; i < n; i++)
      sum += a[i][0];
    ```
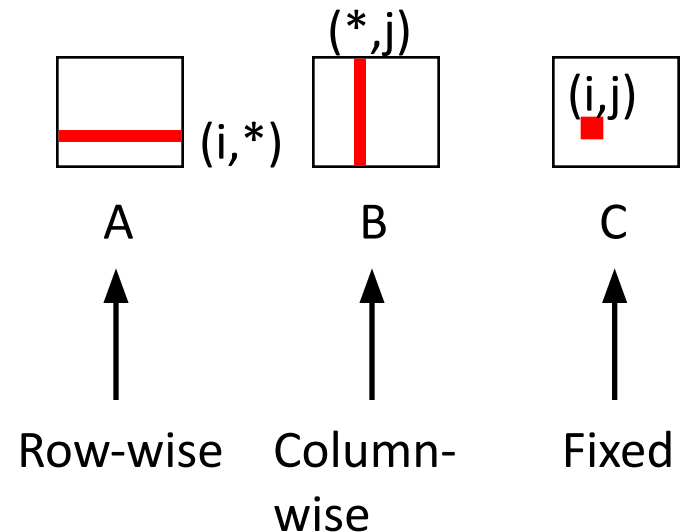  - accesses distant elements
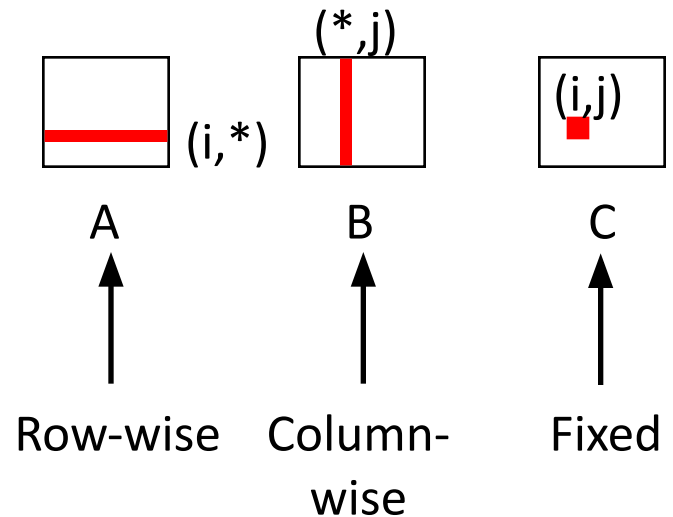  - no spatial locality!
    - miss rate = 1 (i.e. 100%)

# 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;
  }
}
                    matmult/mm.c
```

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
  }
}
                        matmult/mm.c
```

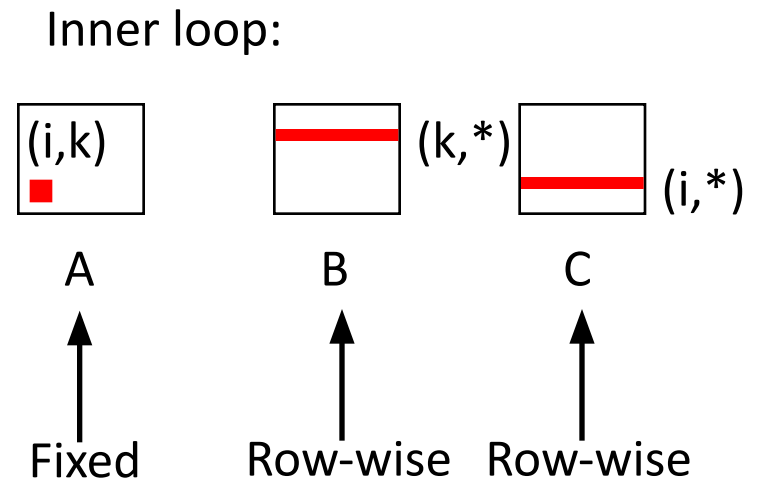Inner loop:



(*,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];
  }
}
                          matmult/mm.c
```
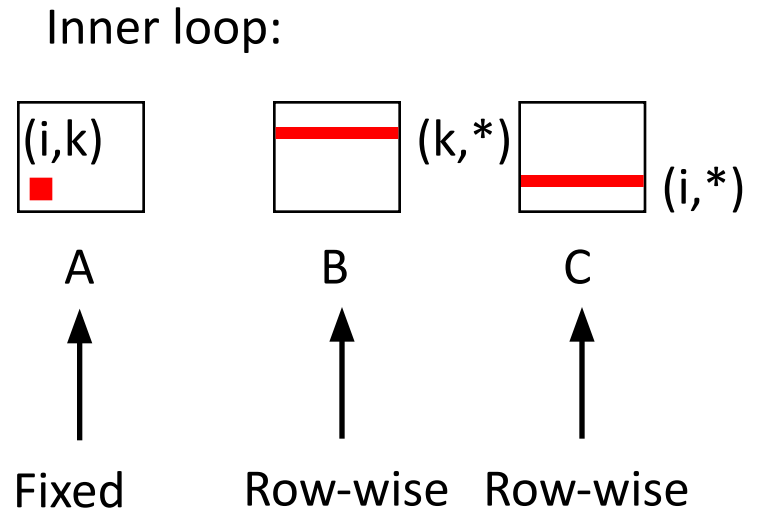
Inner loop:



| A | B | C |
|---|---|---|
| (i,k) | (k,*) | (i,*) |
| 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];
  }
}
                          matmult/mm.c
```
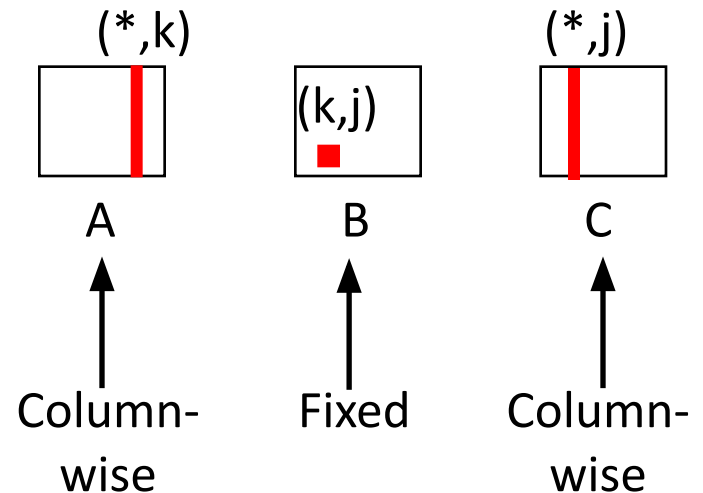
Inner loop:

| (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;
  }
}
                    matmult/mm.c
```
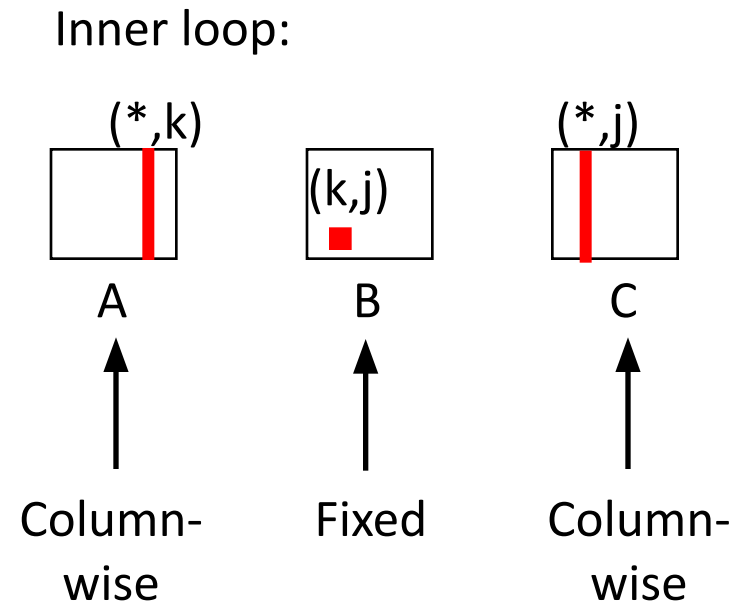
Inner loop:



| (*,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;
  }
}
                        matmult/mm.c
```

Inner loop:

```
        (*,k)                        (*,j)
      ┌──────┐      ┌──────┐      ┌──────┐
      │      │      │(k,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];
  }
}
```
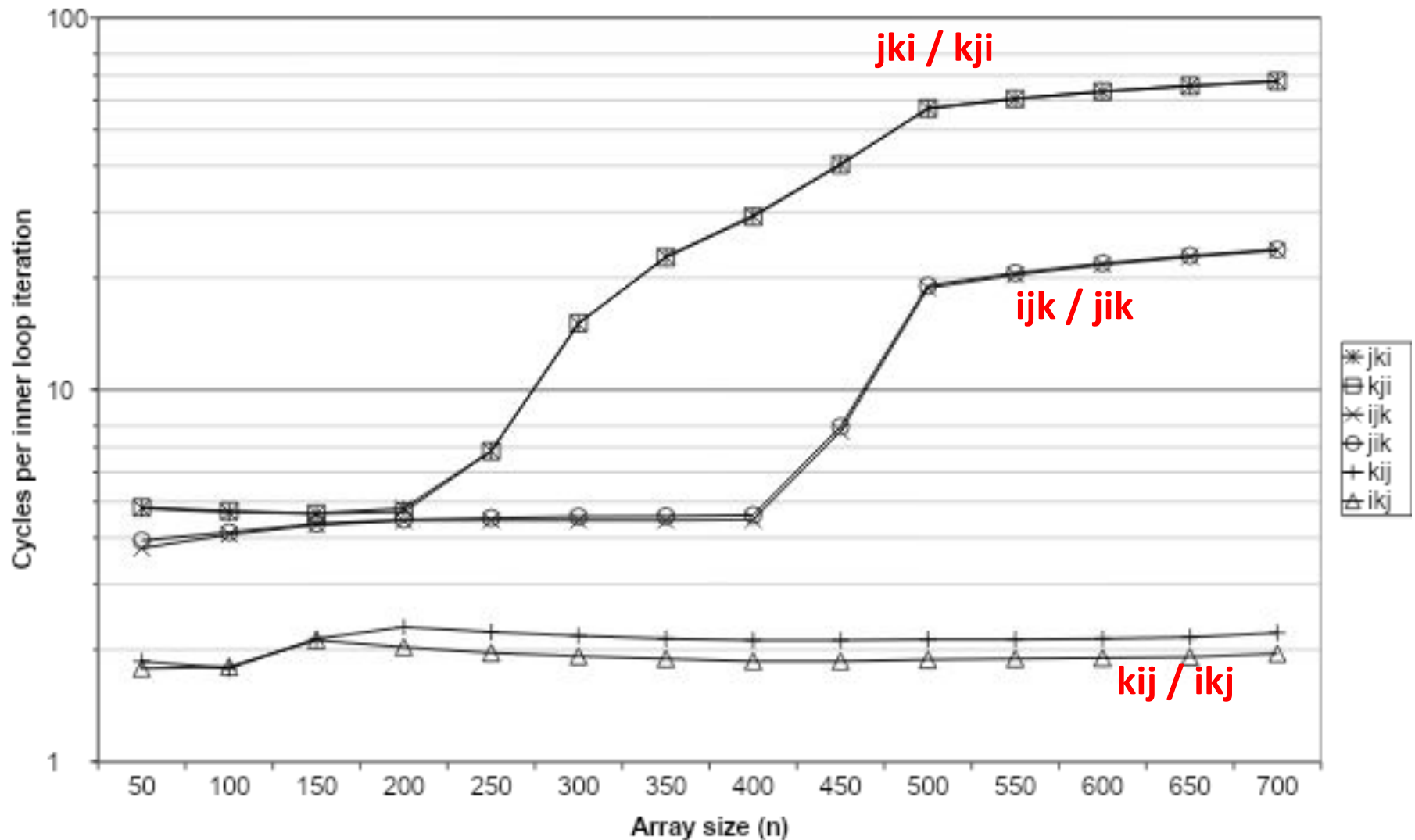
**kij (& ikj):**
- 2 loads, 1 store
- misses/iter = **0.5**

```
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):**
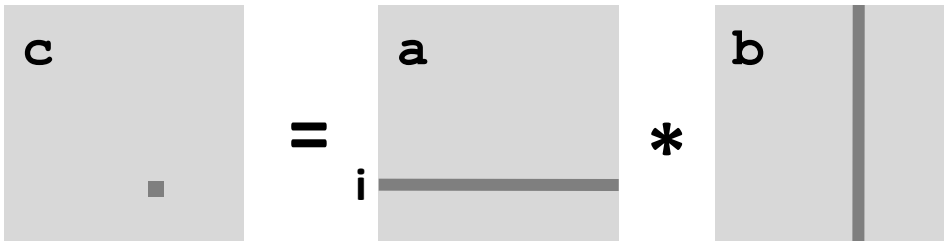- 2 loads, 1 store
- misses/iter = **2.0**

# Core i7 Matrix Multiply Performance

# Use blocking to improve temporal locality

# Example: Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);

/* 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];
}
```
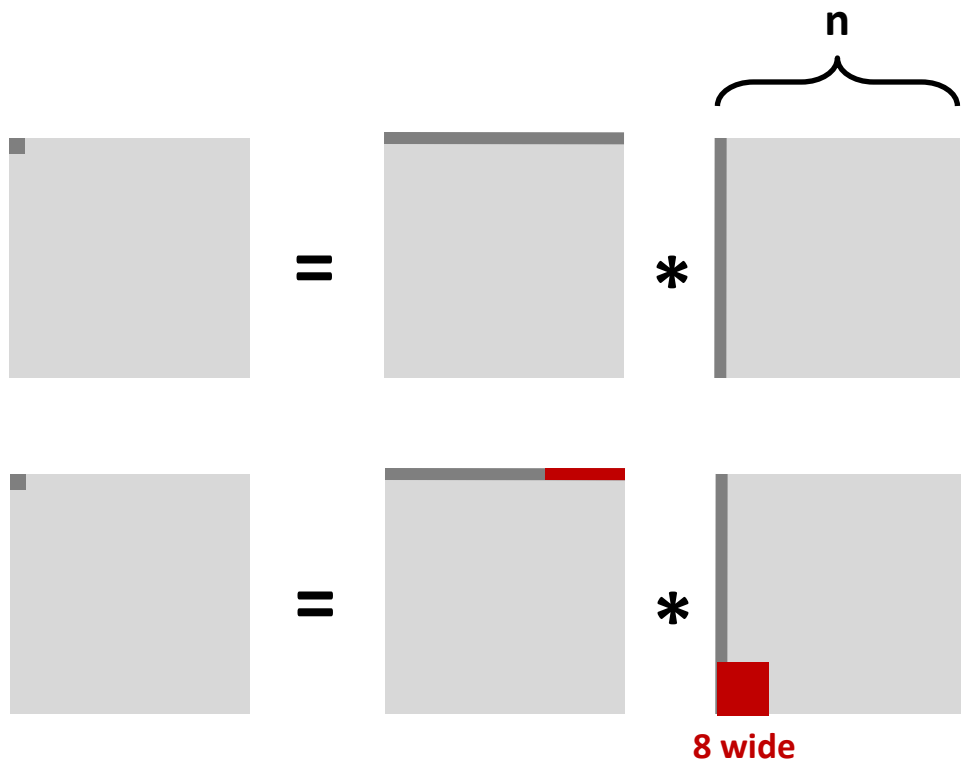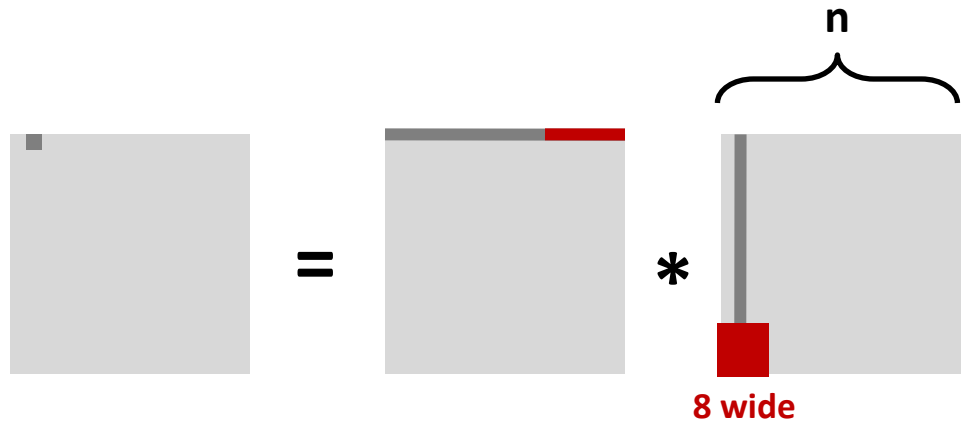
# Cache Miss Analysis

- **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)

# 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
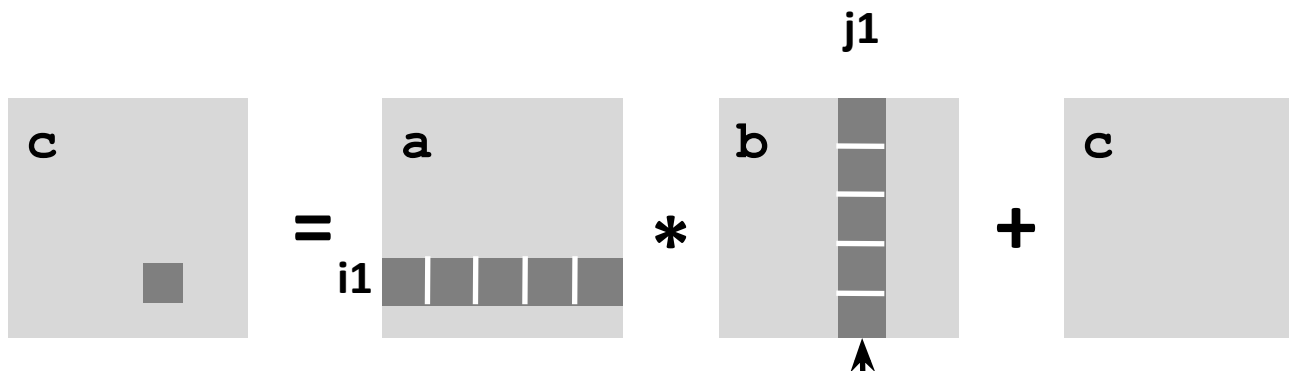
**n**

$=$    $*$

**8 wide**

▶ **Total misses:**

    ▶ $9n/8 * n^2 = (9/8) * n^3$

# Blocked Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);

/* 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];
}
                                            matmult/bmm.c
```
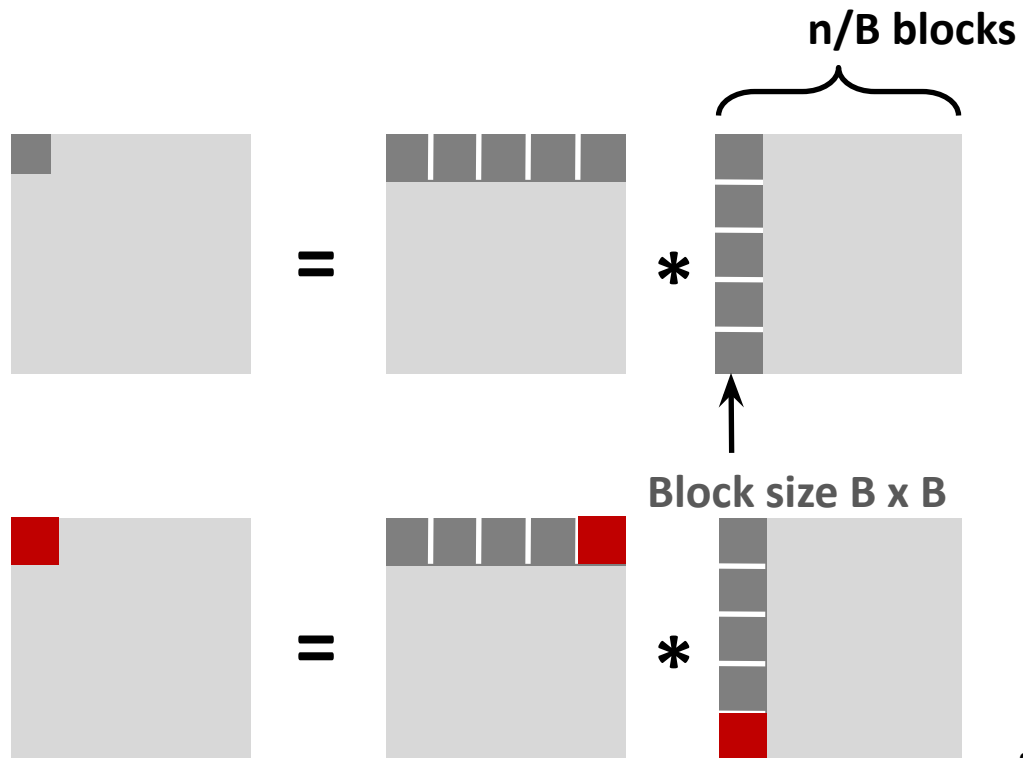


**Block size B x B**

# Cache Miss Analysis

▶ **Assume:**

    ▶ Cache block = 8 doubles

    ▶ Cache size C << 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)

    ▶ Afterwards in cache
    (schematic)

**n/B blocks**

$=$    $*$

**Block size B x B**

$=$    $*$

# Cache Miss Analysis

▶ **Assume:**

  ▶ Cache block = 8 doubles

  ▶ Cache size C << n (much smaller than n)
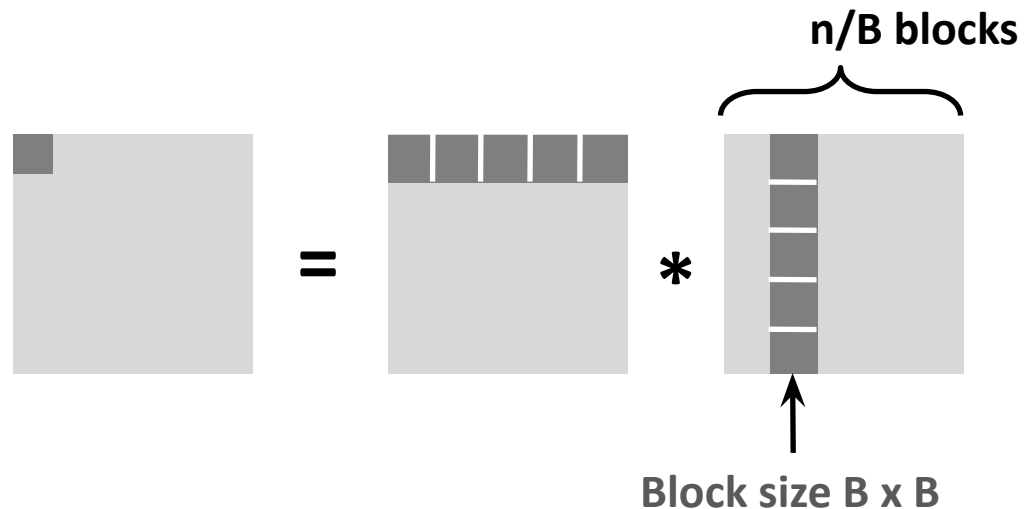
  ▶ Three blocks ▪ fit into cache: $3B^2 < C$

▶ **Second (block) iteration:**

  ▶ Same as first iteration

  ▶ $2n/B * B^2/8 = nB/4$

▶ **Total misses:**

  ▶ $nB/4 * (n/B)^2 = n^3/(4B)$

**n/B blocks**

= *

**Block size B x B**

# Blocking 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 elements used $O(n)$ times!

  - ► But program has to be written properly

# Cache Summary

▶ **Cache memories can have significant performance impact**

▶ **You can write your programs to exploit this!**

 ▶ Focus on the inner loops, where bulk of computations and memory accesses occur.

 ▶ Try to maximize spatial locality by reading data objects with sequentially with stride 1.

 ▶ Try to maximize temporal locality by using a data object as often as possible once it's read from memory.