

# Writing Cache Friendly Code

90/10

- **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**)

No HANDOUT

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

# Today

- **Performance impact of caches**
  - The memory mountain
  - Rearranging loops to improve spatial locality
  - Using blocking to improve temporal 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

```
/* 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 */
}
```

*Handwritten annotations:*

- elems* is boxed in red.
- stride* is underlined in red.
- elems* is boxed in red in the for loop.
- stride* is underlined in red in the for loop.
- data[i]* is underlined in red.
- spatial* is written in red with an arrow pointing to the for loop.
- Stride* is written in red above a diagram.
- The diagram shows a horizontal array of 10 cells. The first 4 cells are shaded with diagonal lines. A red double-headed arrow spans the first 4 cells.
- Don* is written in red below the arrow.

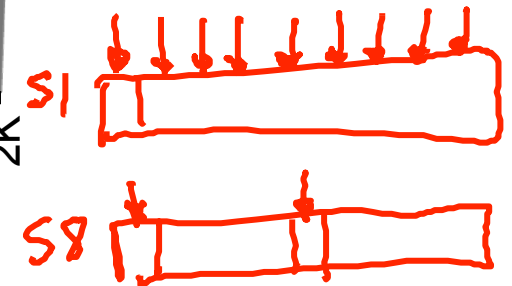
# The Memory Mountain

Read throughput (MB/s)



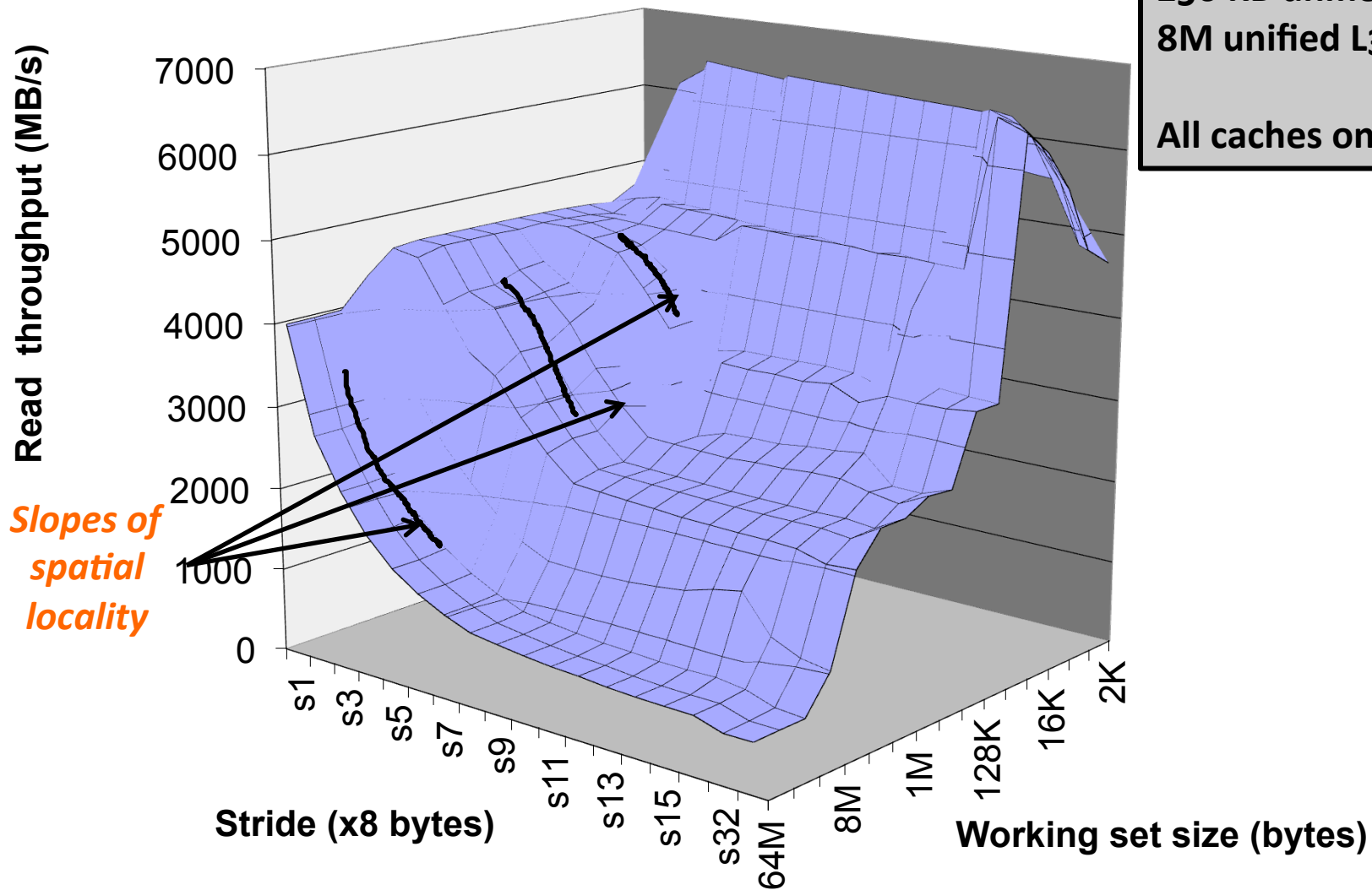
Intel Core i7  
32 KB L1 i-cache  
32 KB L1 d-cache  
256 KB unified L2 cache  
8M unified L3 cache

All caches on-chip

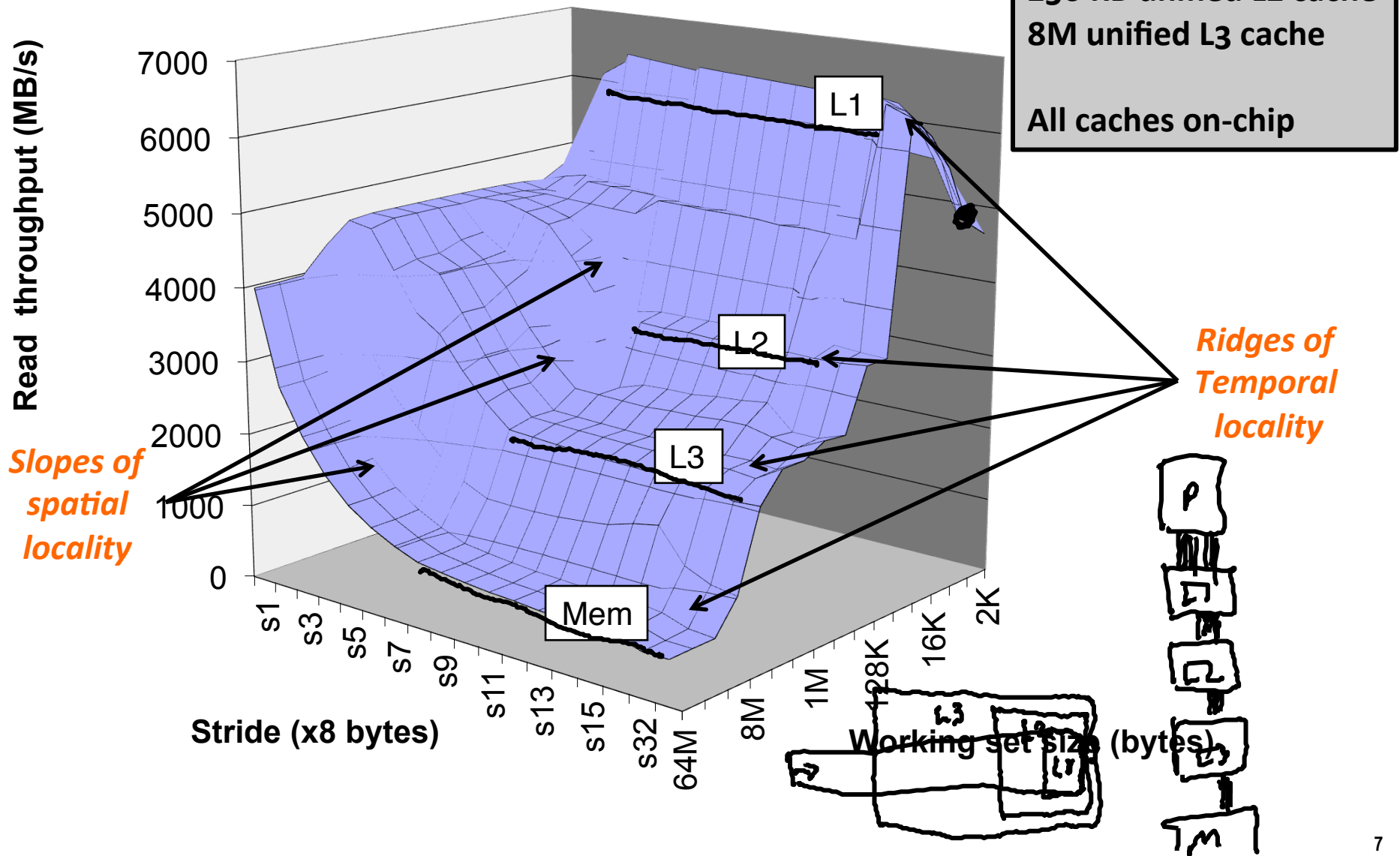


Working set size (bytes)

# The Memory Mountain



# The Memory Mountain



# Today

- **Performance impact of caches**
  - The memory mountain
  - Rearranging loops to improve spatial locality
  - Using blocking to improve temporal locality



# Miss Rate Analysis for Matrix Multiply

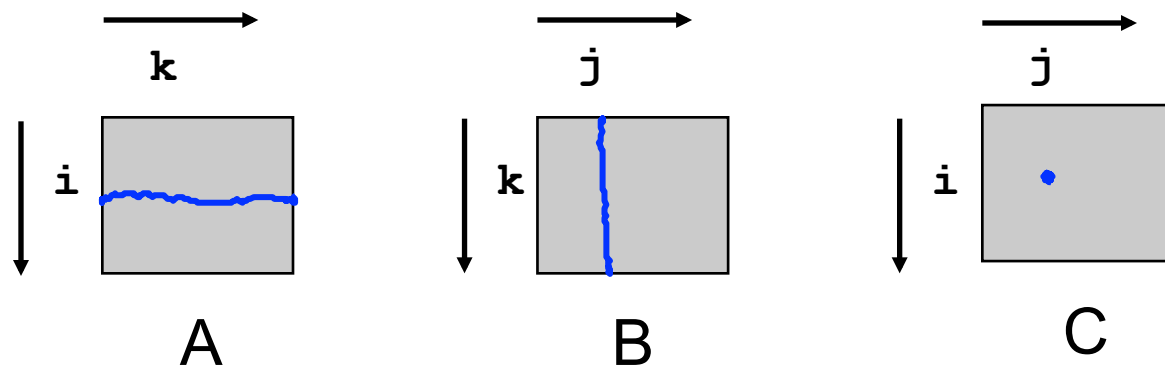
## ■ Assume:

- 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

$$C[i][j] = \sum_{k=1}^N A[i][k] \cdot B[k][j]$$



# Matrix Multiplication Example

## ■ Description:

- Multiply  $N \times N$  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*

# Layout of C Arrays in Memory

## ■ C arrays allocated in row-major order

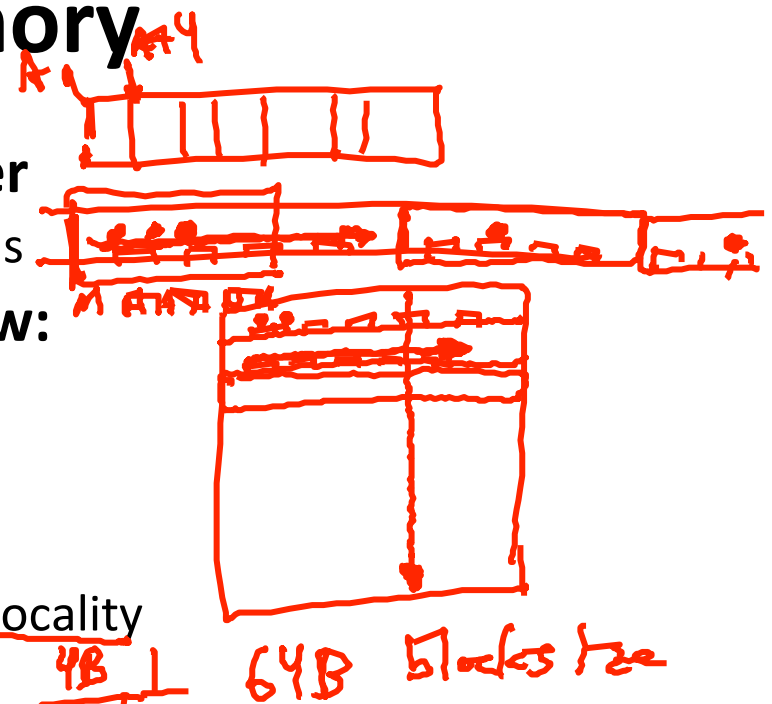
- each row in contiguous memory locations

## ■ 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%)



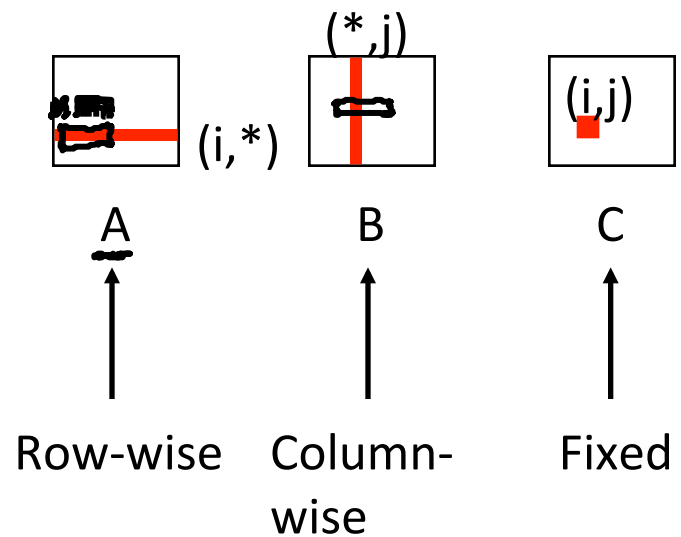
# Matrix Multiplication (ijk) 9 words/block

```

/* 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:



Misses per inner loop iteration:

$$\underline{0.25} + \underline{1.0} + 0.0 = 1.25 \text{ misses/iteration}$$

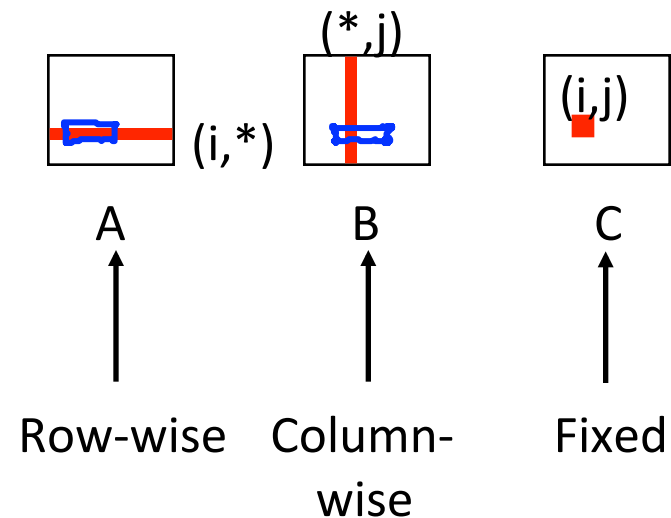
# Matrix Multiplication (jik)

9 words/block

```

/* 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

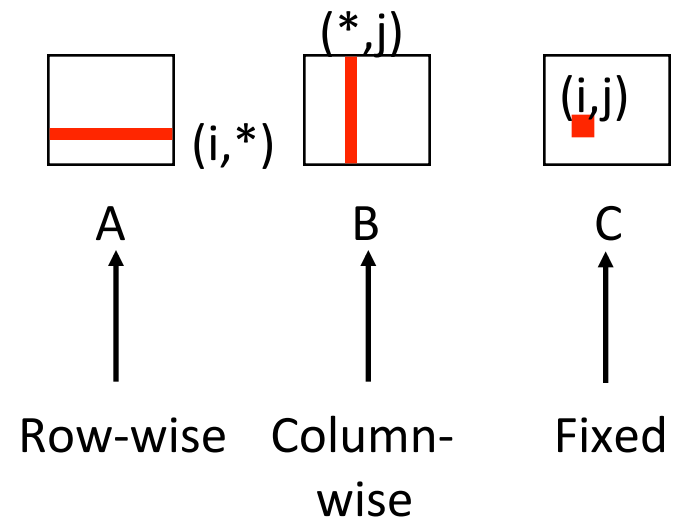
$$\underline{.25} + \underline{1.0} + \underline{0} = 1.25 \text{ misses/iter}$$

- a) 0
- b) .25
- c) .75
- d) 1.0
- e) 2.0

# 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:



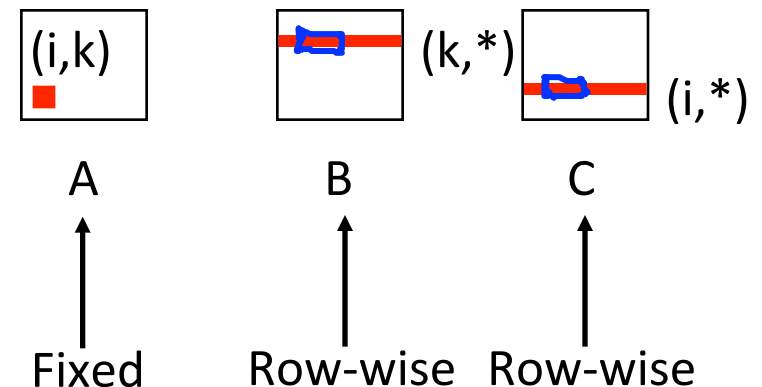
Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
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:



Misses per inner loop iteration:

A      B      C

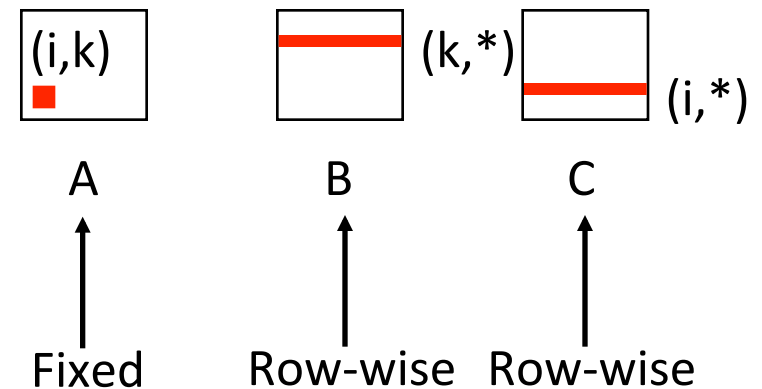
0      .25      .25 = .5 <sup>miss</sup> / iter

- a) 0
- b) .25
- c) .75
- d) 1.0
- e) 2.0

# 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:



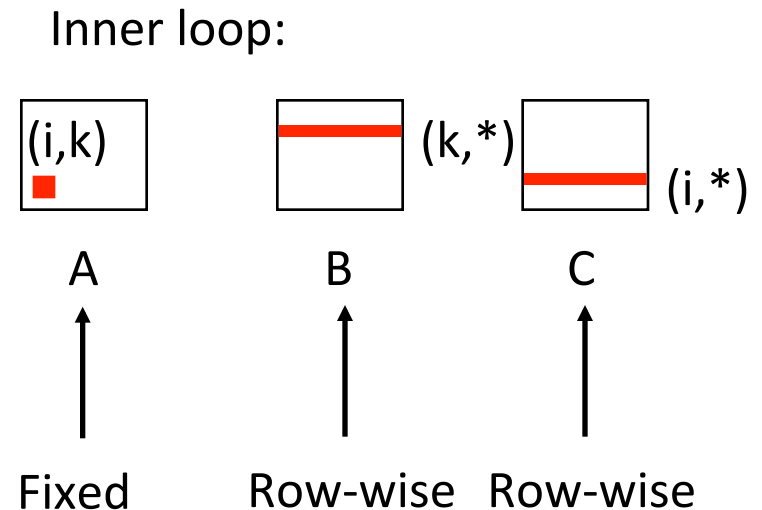
Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
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];
    }
}
```



Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
0.0	0.25	0.25

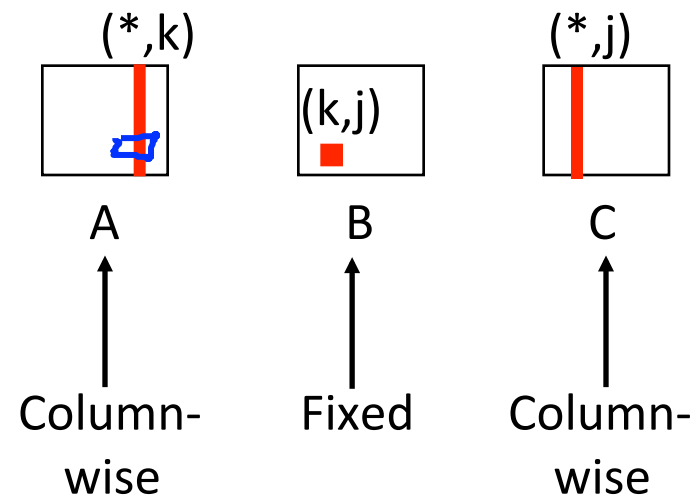
# Matrix Multiplication (jki)

```

/* 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:



Misses per inner loop iteration:

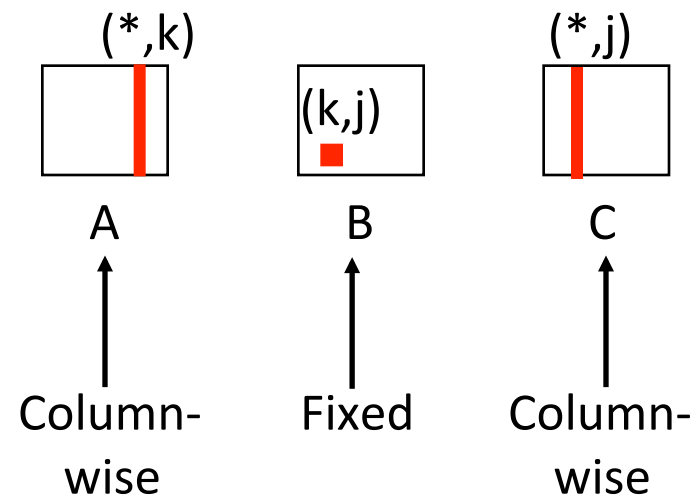
A                  B                  C  
1.0 + 0.0 + 1.0 = 2 misses / iteration

- a) 0
- b) .25
- c) .75
- d) 1.0
- e) 2.0

# 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:

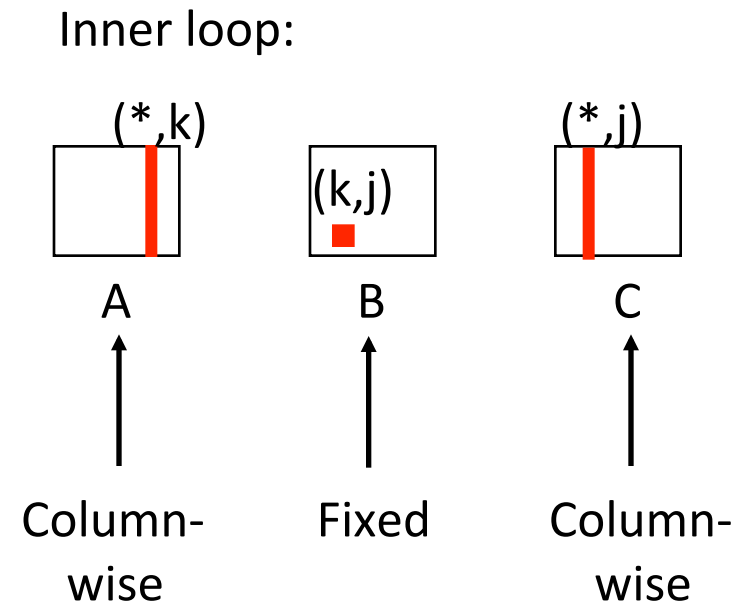


Misses per inner loop iteration:

<u>A</u>	<u>B</u>	<u>C</u>
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;
    }
}
```



Misses per inner loop iteration:

A  
1.0

B  
0.0

C  
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

