

# Introduction to Computer Systems

15-213, fall 2009

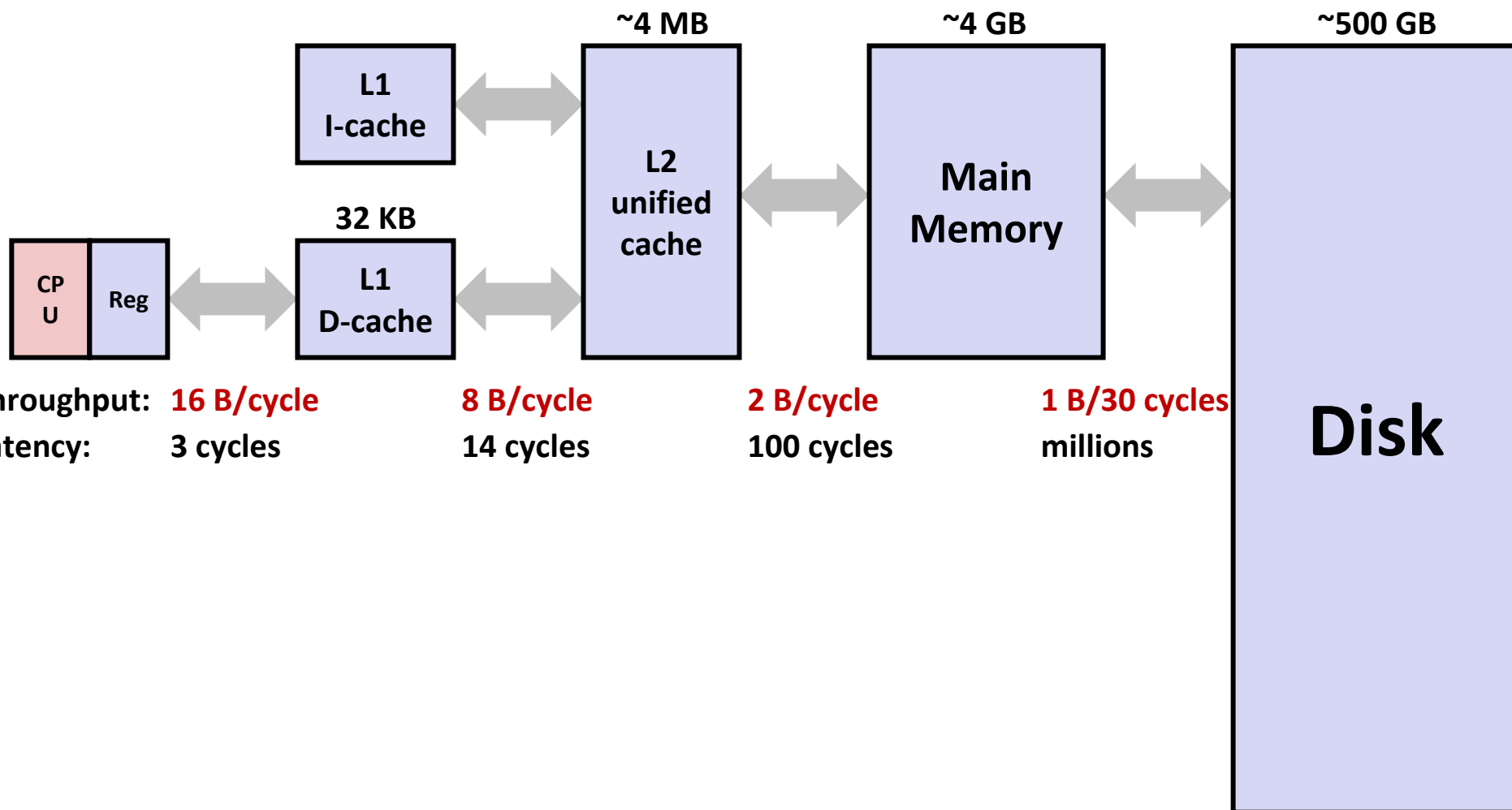
11<sup>th</sup> Lecture, Oct. 5<sup>th</sup>

## **Instructors:**

Majd Sakr and Khaled Harras

# Last Time

## ■ Memory hierarchy (Here: Core 2 Duo)

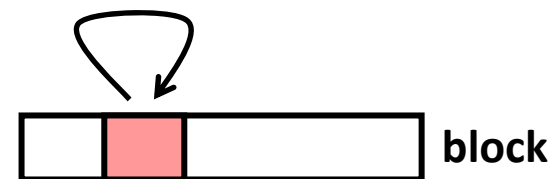


# Last Time

## ■ Locality

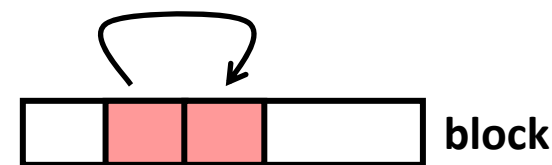
### ■ Temporal locality:

- Recently referenced items are likely to be referenced again in the near future



### ■ Spatial locality:

- Items with nearby addresses tend to be referenced close together in time

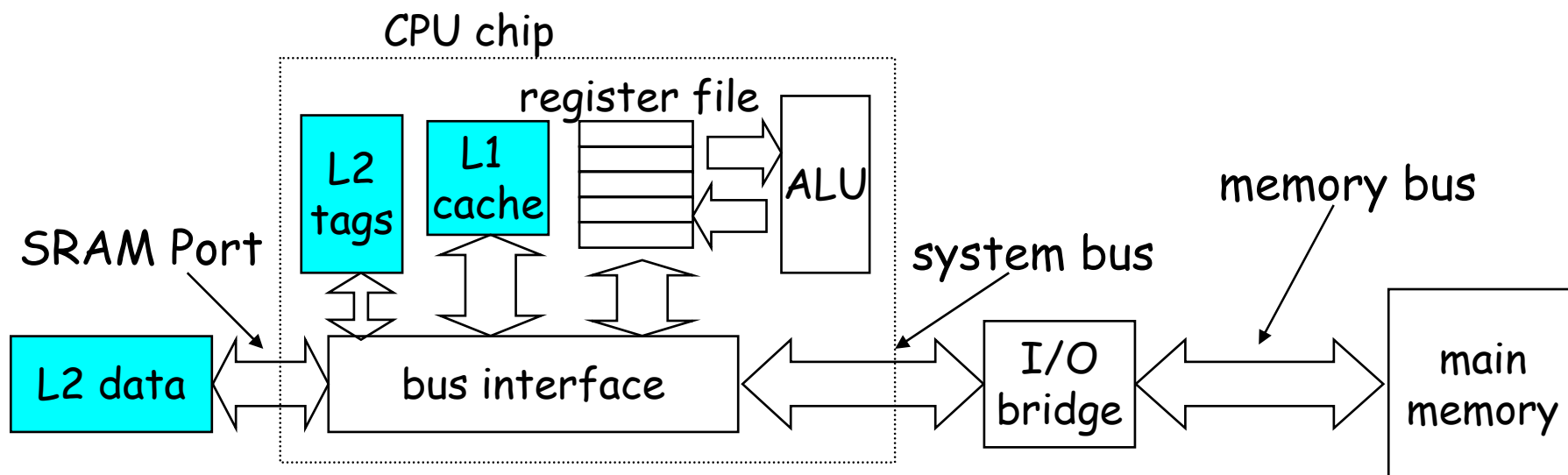


# Today

- **Cache organization**
- **Memory Wall**
- **Program optimization (Matrix Multiplication):**
  - Cache optimizations
  - Cache Miss Analysis

# Cache Memories

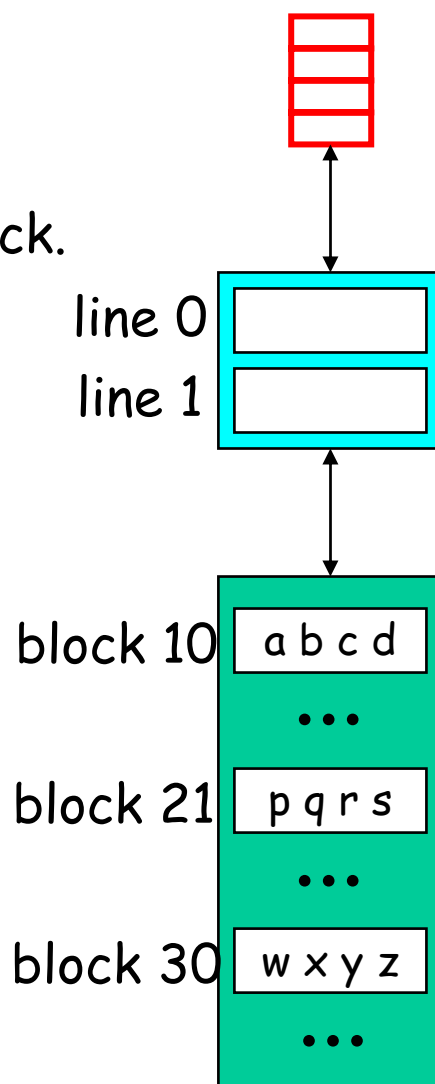
- Cache memories are small, fast SRAM-based memories managed automatically in hardware.
  - Hold frequently accessed blocks of main memory
- CPU looks first for data in L1, then in L2, then in main memory.
- Typical system structure:



# Inserting an L1 Cache Between the CPU and Main Memory

The transfer unit between the CPU **register file** and the **cache** is a 4-byte block.

The transfer unit between the cache and **main memory** is a 4-word block (16 bytes).

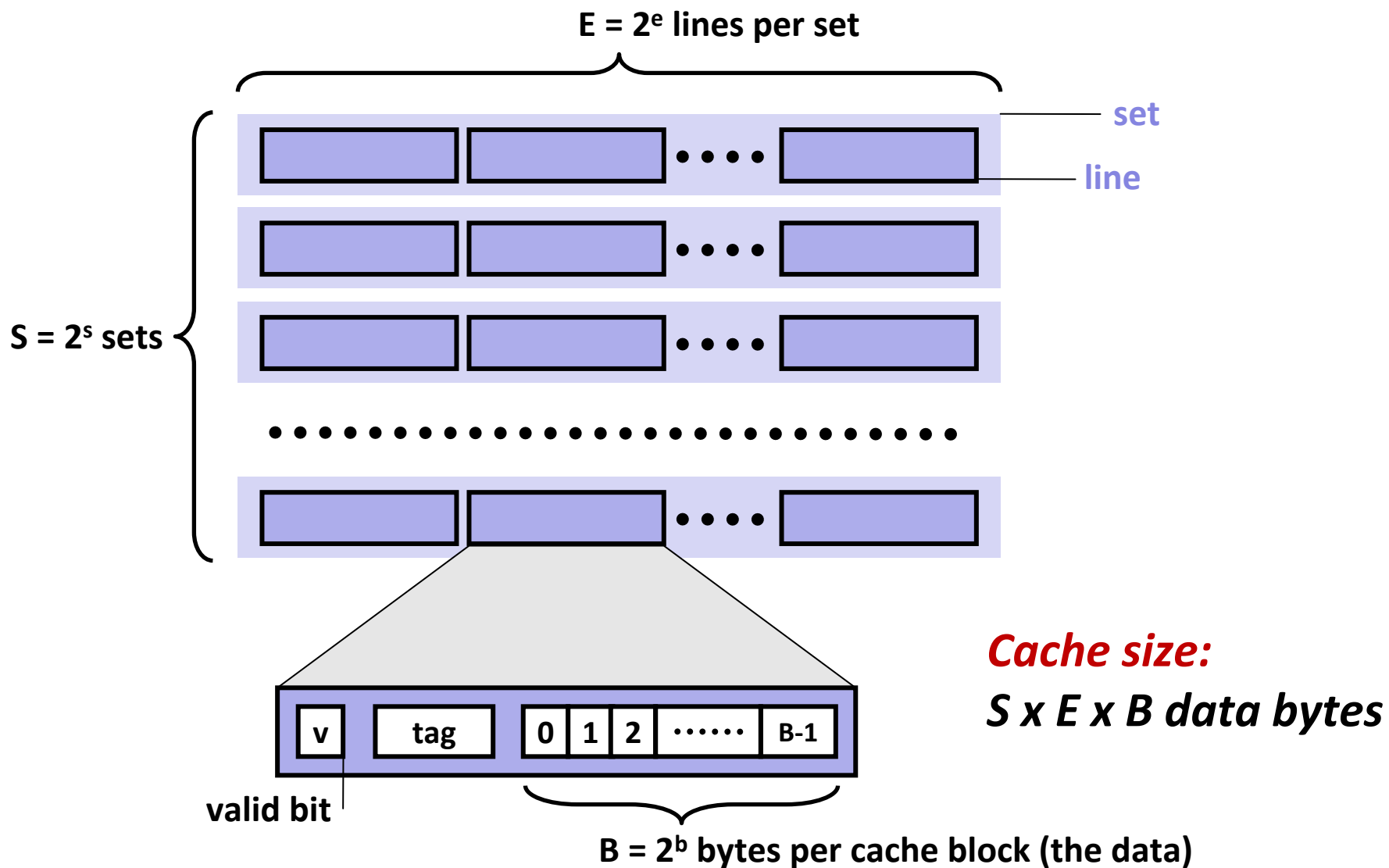


The tiny, very fast CPU **register file** has room for four 4-byte words.

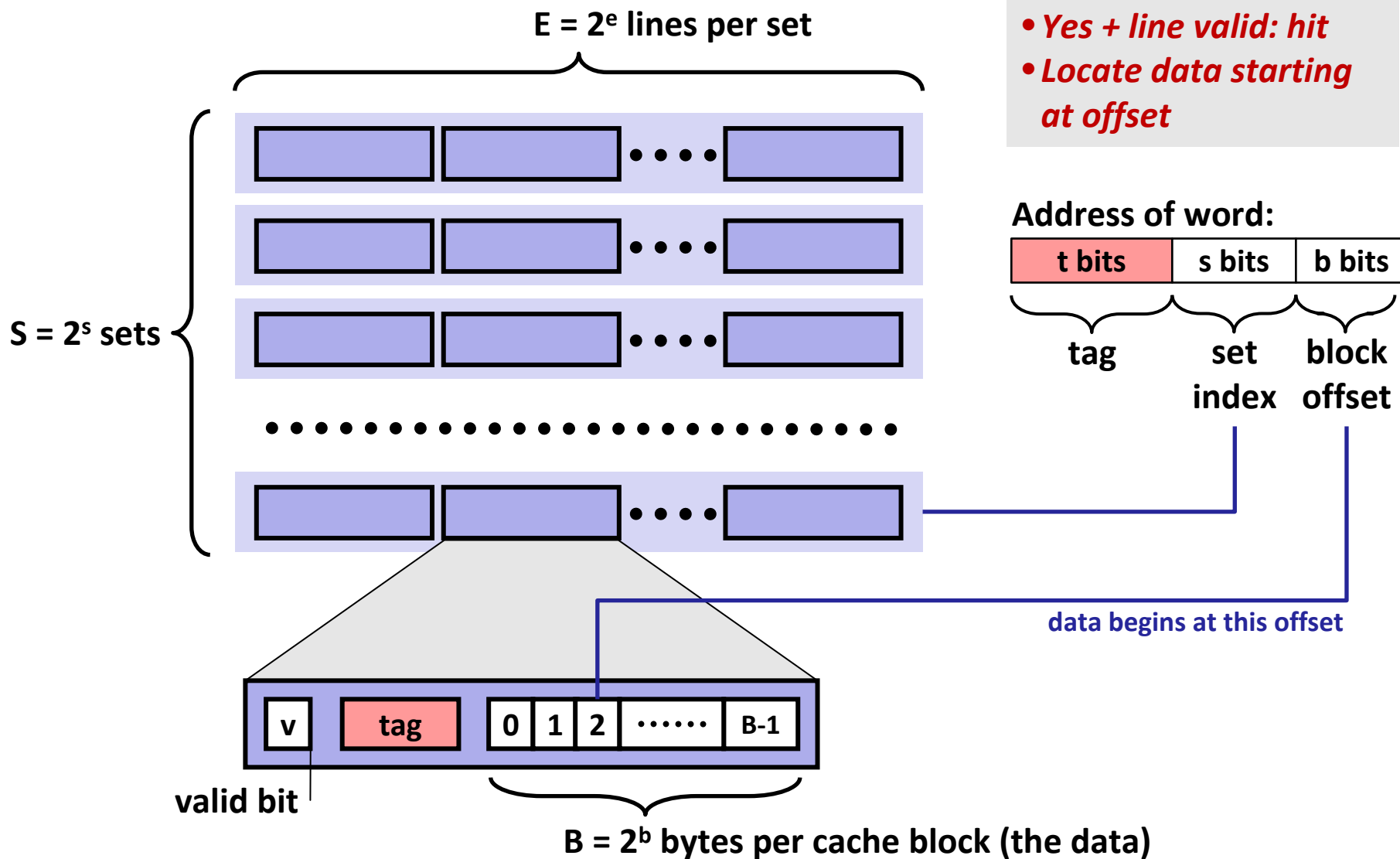
The small fast **L1 cache** has room for two 4-word blocks.

The big slow **main memory** has room for many 4-word blocks.

# General Cache Organization (S, E, B)



# Cache Read

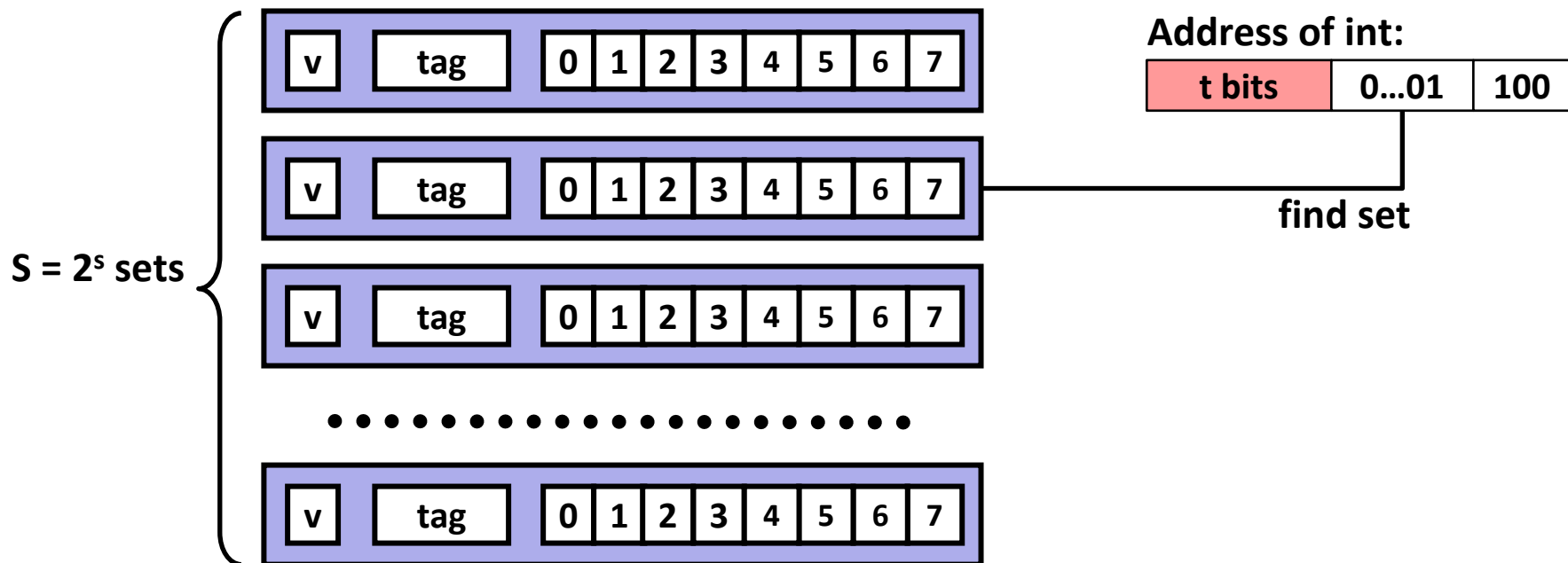




# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set

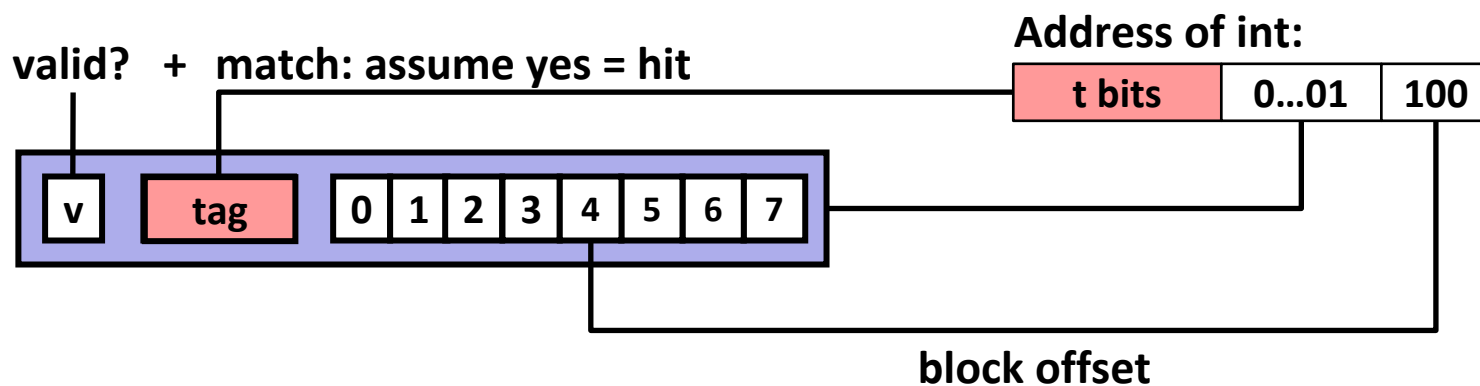
Assume: cache block size 8 bytes



# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set

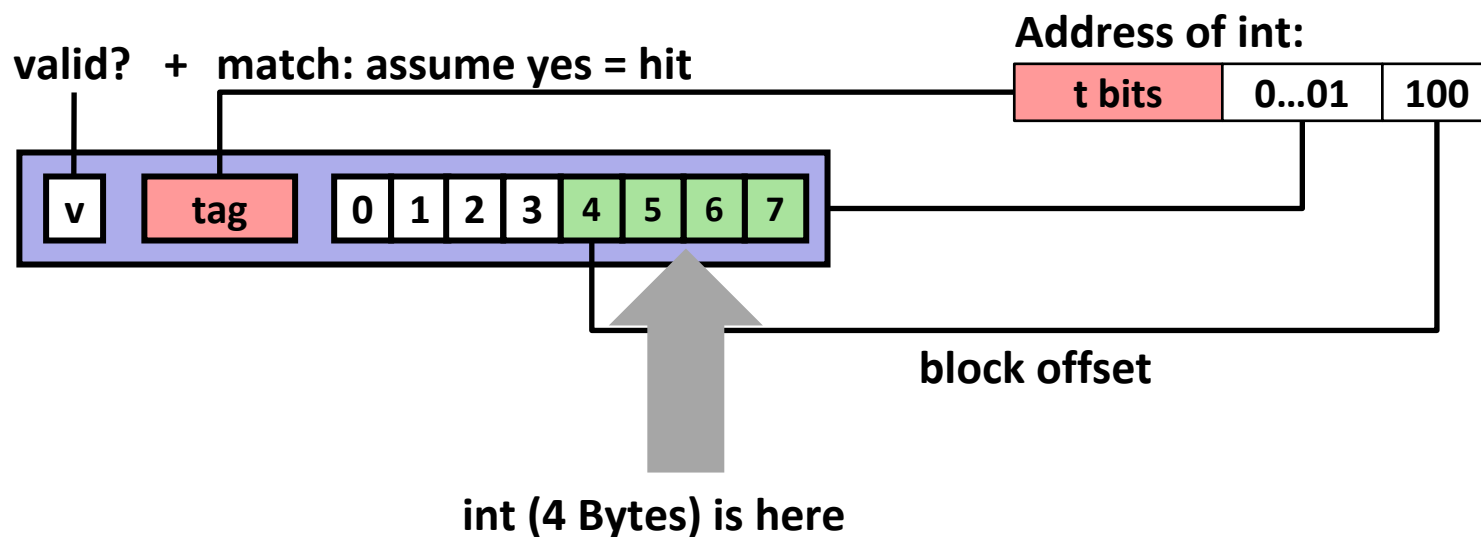
Assume: cache block size 8 bytes



# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set

Assume: cache block size 8 bytes



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

# Direct-Mapped Cache Simulation

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

$t=1$   $s=2$   $b=1$

X	XX	X
---	----	---

Address trace (reads):

0	[0000 <sub>2</sub> ],	miss
1	[0001 <sub>2</sub> ],	hit
7	[0111 <sub>2</sub> ],	miss
8	[1000 <sub>2</sub> ],	miss
0	[0000 <sub>2</sub> ]	miss

v tag data

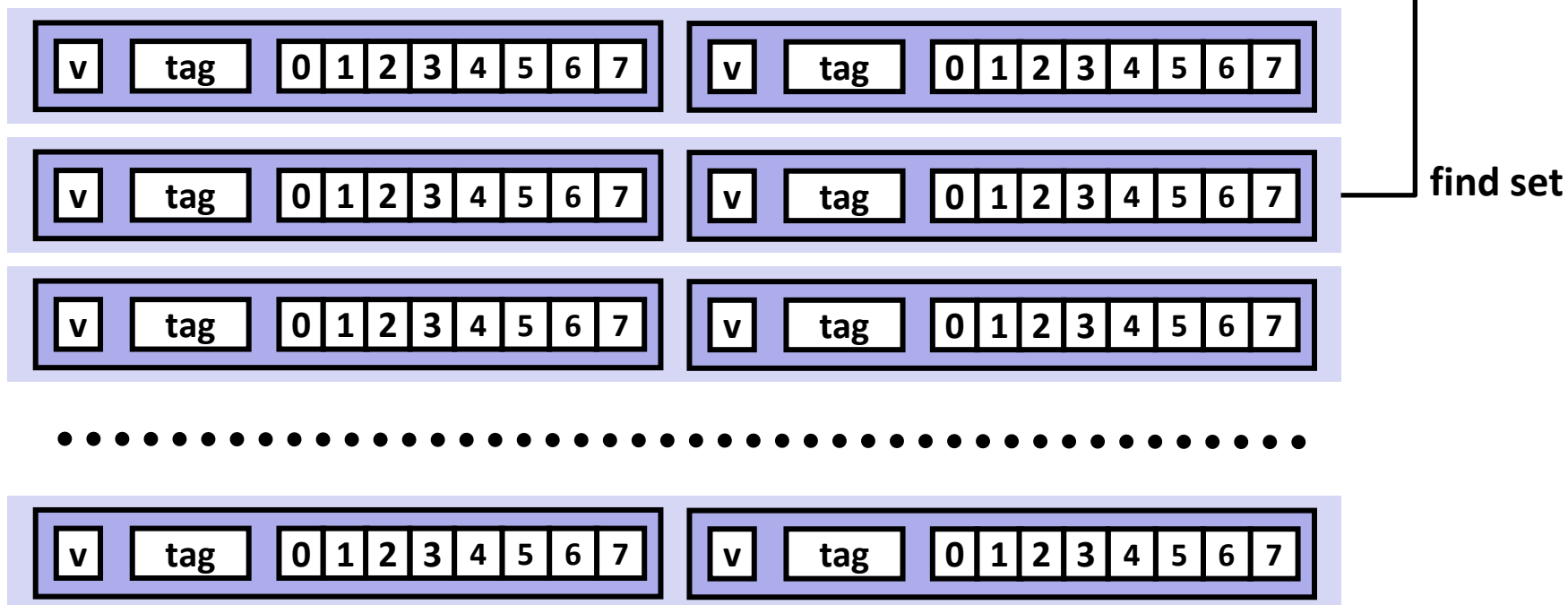
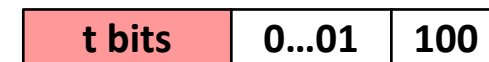
1	0	M[0-1]
1	0	M[6-7]

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

E = 2: Two lines per set

Assume: cache block size 8 bytes

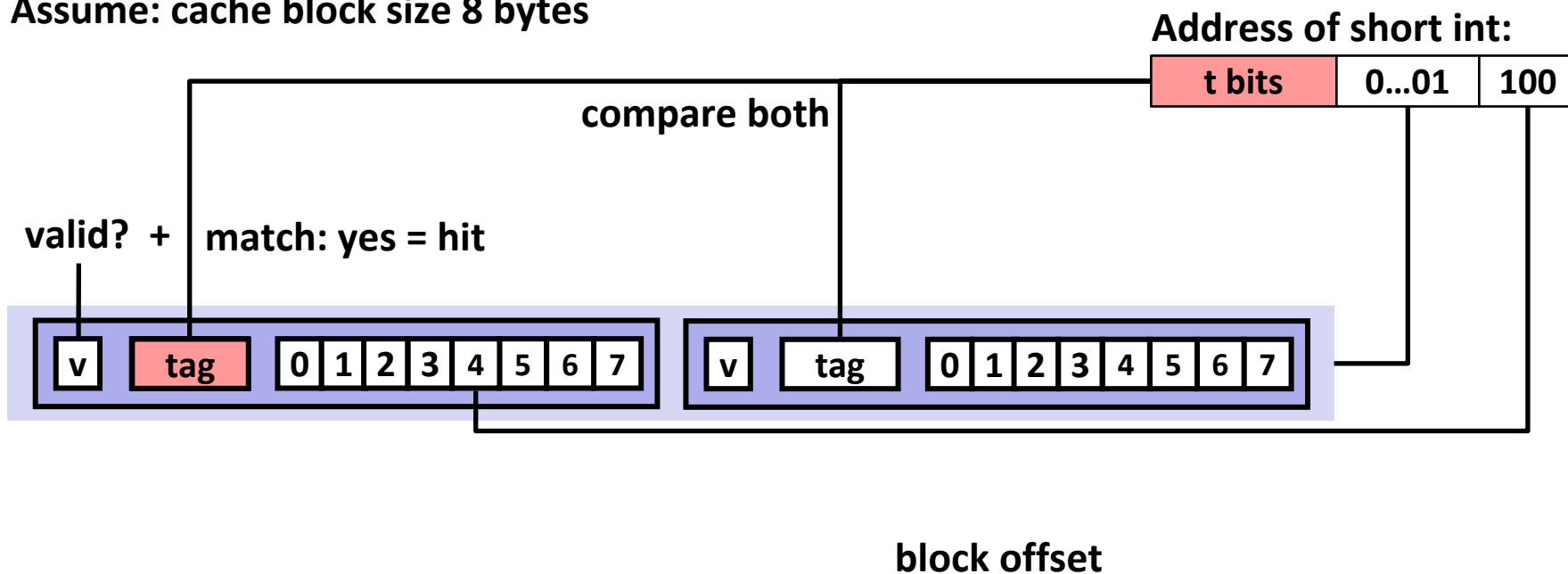
Address of short int:



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

E = 2: Two lines per set

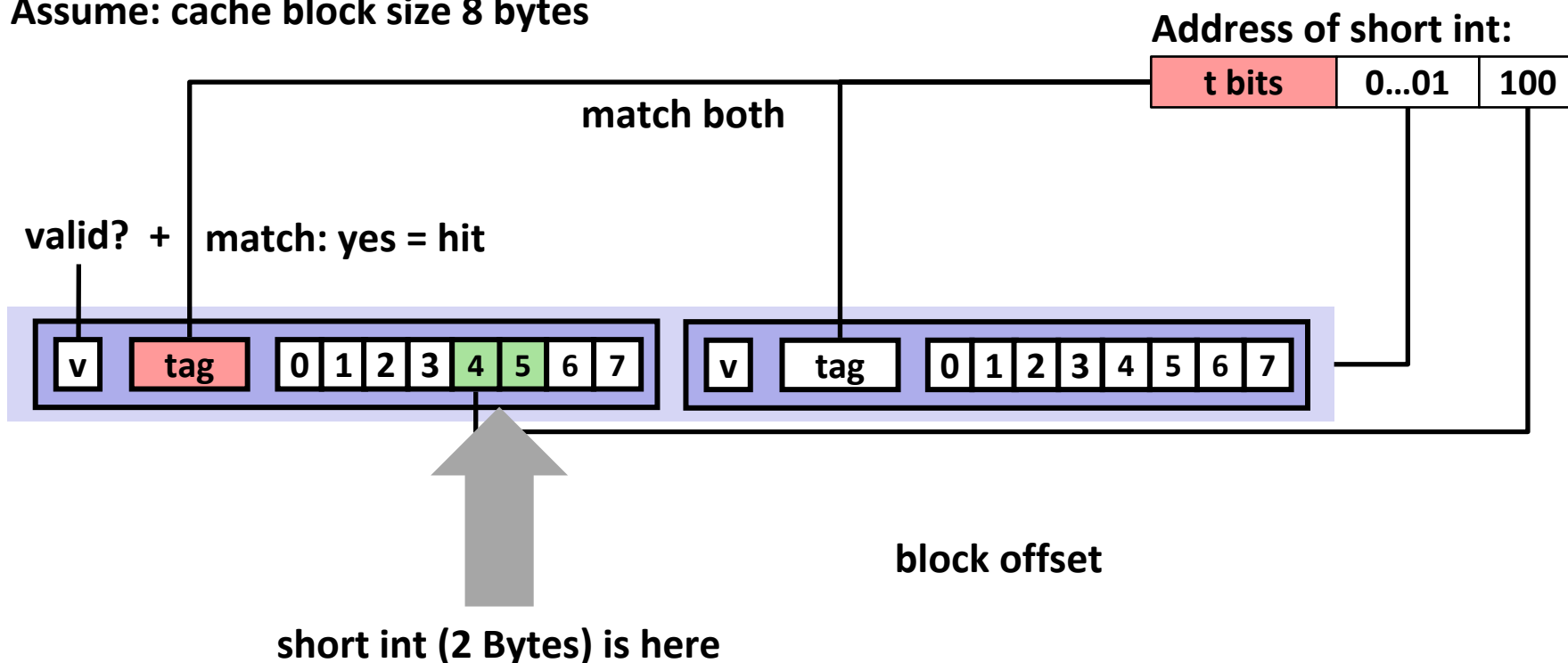
Assume: cache block size 8 bytes



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

E = 2: Two lines per set

Assume: cache block size 8 bytes



## No match:

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

# 2-Way Associative Cache Simulation

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

$t=2$   $s=1$   $b=1$

xx	x	x
----	---	---

Address trace (reads):

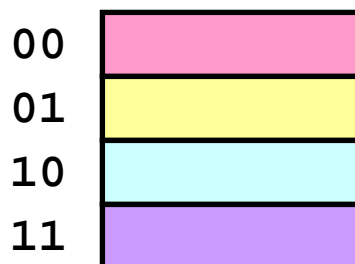
0	[0000 <sub>2</sub> ],	miss
1	[0001 <sub>2</sub> ],	hit
7	[0111 <sub>2</sub> ],	miss
8	[1000 <sub>2</sub> ],	miss
0	[0000 <sub>2</sub> ]	hit

v	tag	data
1	00	$M[0-1]$
1	10	$M[8-9]$
1	01	$M[6-7]$
0		



# Why Use Middle Bits as Index?

## 4-line Cache



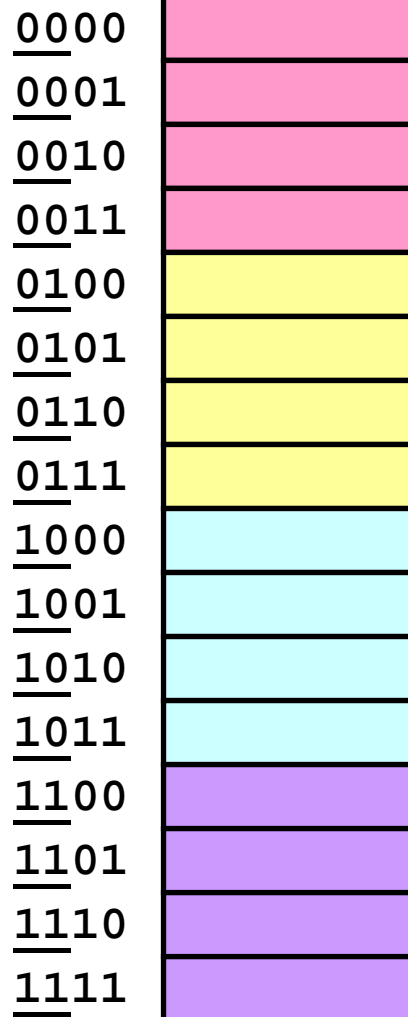
## ■ High-Order Bit Indexing

- Adjacent memory lines would map to same cache entry
- Poor use of spatial locality

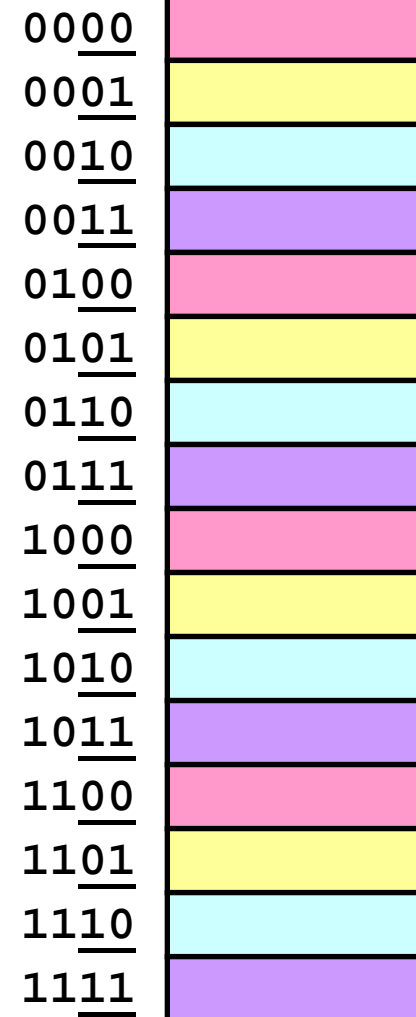
## ■ Middle-Order Bit Indexing

- Consecutive memory lines map to different cache lines
- Can hold  $S*B*E$ -byte region of address space in cache at one time

## High-Order Bit Indexing



## Middle-Order Bit Indexing



# 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)  
=  $1 - \text{hit rate}$
- 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

Aside for architects:

-Increasing cache size?

-Increasing block size?

-Increasing associativity?

## ■ Miss Penalty

- Additional time required because of a miss
  - typically 50-200 cycles for main memory (Trend: increasing!)

# Software Caches are More Flexible

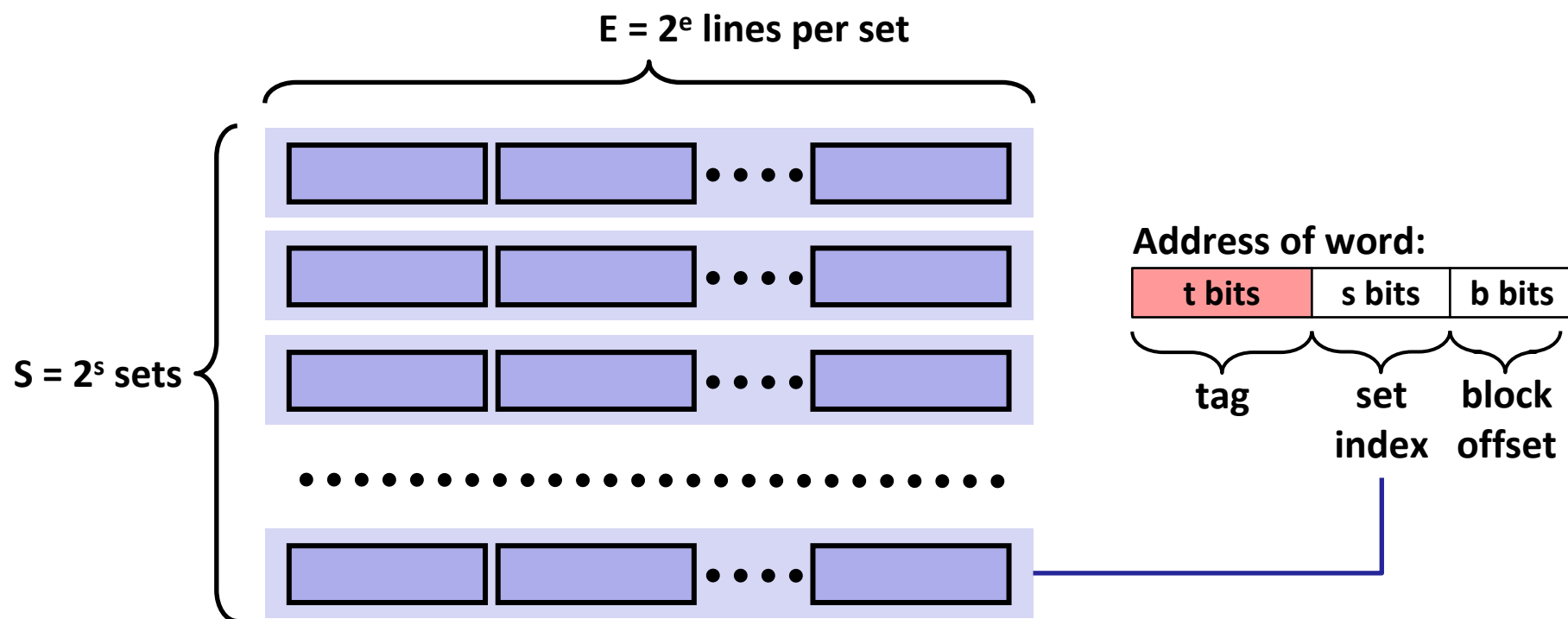
## ■ Examples

- File system buffer caches, web browser caches, etc.

## ■ Some design differences

- Almost always fully associative
  - so, no placement restrictions
  - index structures like hash tables are common
- Often use complex replacement policies
  - misses are very expensive when disk or network involved
  - worth thousands of cycles to avoid them
- Not necessarily constrained to single “block” transfers
  - may fetch or write-back in larger units, opportunistically

# Strided Access Question



- What happens if arrays are accessed in two-power strides?
- Example on the next slide

# The Strided Access Problem (Blackboard?)

- **Example: L1 cache, Core 2 Duo**
  - 32 KB, 8-way associative, 64 byte cache block size
  - What is S, E, B?
    - *Answer:*  $B = 2^6$ ,  $E = 2^3$ ,  $S = 2^6$ .
- **Consider an array of ints accessed at stride  $2^i$ ,  $i \geq 0$** 
  - What is the smallest  $i$  such that only one set is used?
    - *Answer:*  $i = 10$
  - What happens if the stride is  $2^9$ ?
    - *Answer:* two sets are used
- **Source of two-power strides?**
  - Example: Column access of 2-D arrays (images!)

# Writing Cache Friendly Code

- Repeated references to variables are good (**temporal locality**)
- Stride-1 reference patterns are good (**spatial locality**)
- Examples:
  - cold cache, 4-byte words, 4-word cache blocks

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

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

Miss rate = **1/4 = 25%**

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

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

Miss rate = **100%**

# Today

- Cache organization
- **Memory Wall**
- **Program optimization (Matrix Multiplication):**
  - Cache optimizations
  - Cache Miss Analysis



# 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 */
}
```

# Memory Mountain Main Routine

```
/* mountain.c - Generate the memory mountain. */
#define MINBYTES (1 << 10) /* Working set size ranges from 1 KB */
#define MAXBYTES (1 << 23) /* ... up to 8 MB */
#define MAXSTRIDE 16 /* Strides range from 1 to 16 */
#define MAXELEMS MAXBYTES/sizeof(int)

int data[MAXELEMS]; /* The array we'll be traversing */

int main()
{
    int size; /* Working set size (in bytes) */
    int stride; /* Stride (in array elements) */
    double Mhz; /* Clock frequency */

    init_data(data, MAXELEMS); /* Initialize each element in data to 1 */
    Mhz = mhz(0); /* Estimate the clock frequency */
    for (size = MAXBYTES; size >= MINBYTES; size >>= 1) {
        for (stride = 1; stride <= MAXSTRIDE; stride++)
            printf("%.1f\t", run(size, stride, Mhz));
        printf("\n");
    }
    exit(0);
}
```

# The Memory Mountain

Pentium III

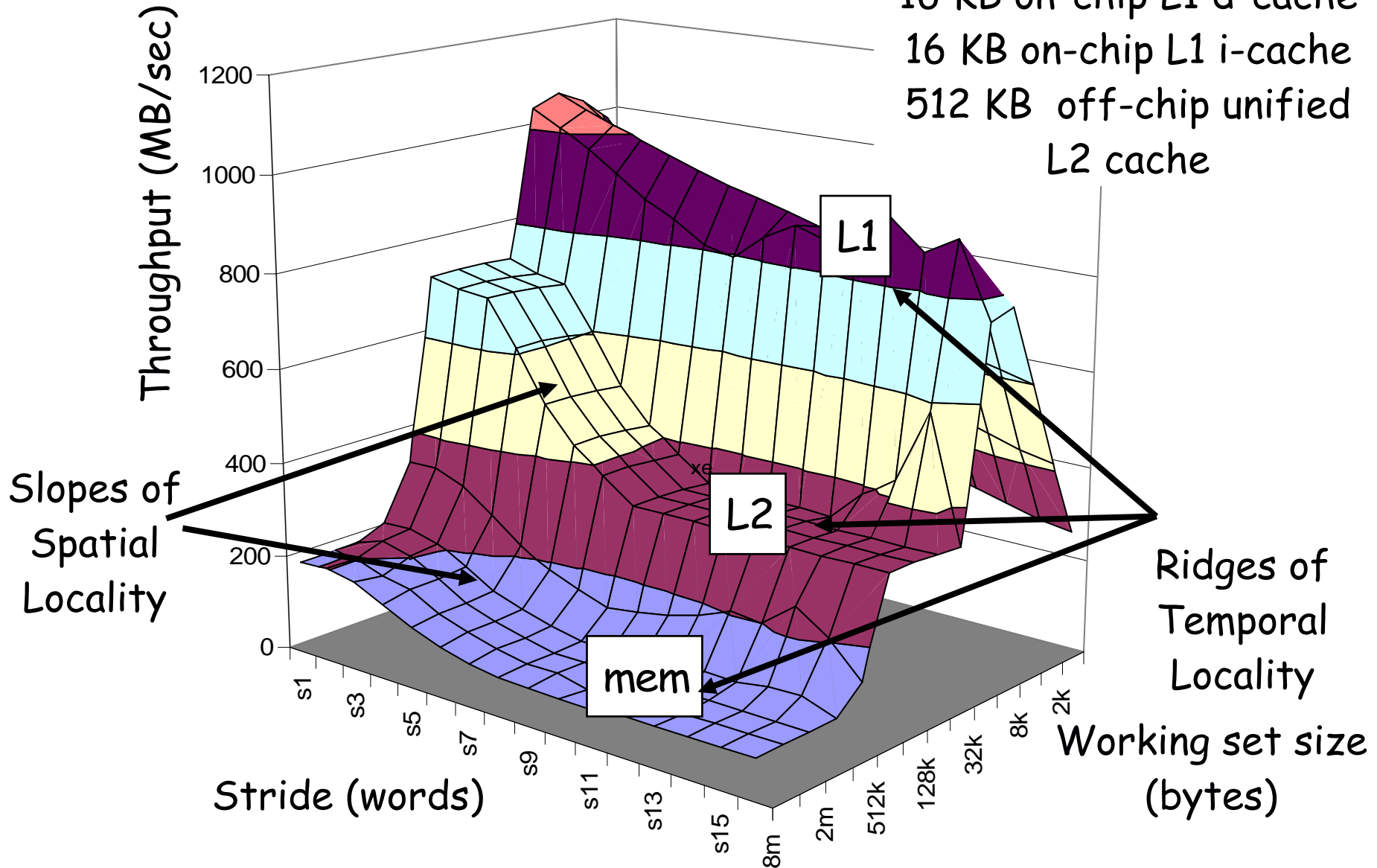
550 MHz

16 KB on-chip L1 d-cache

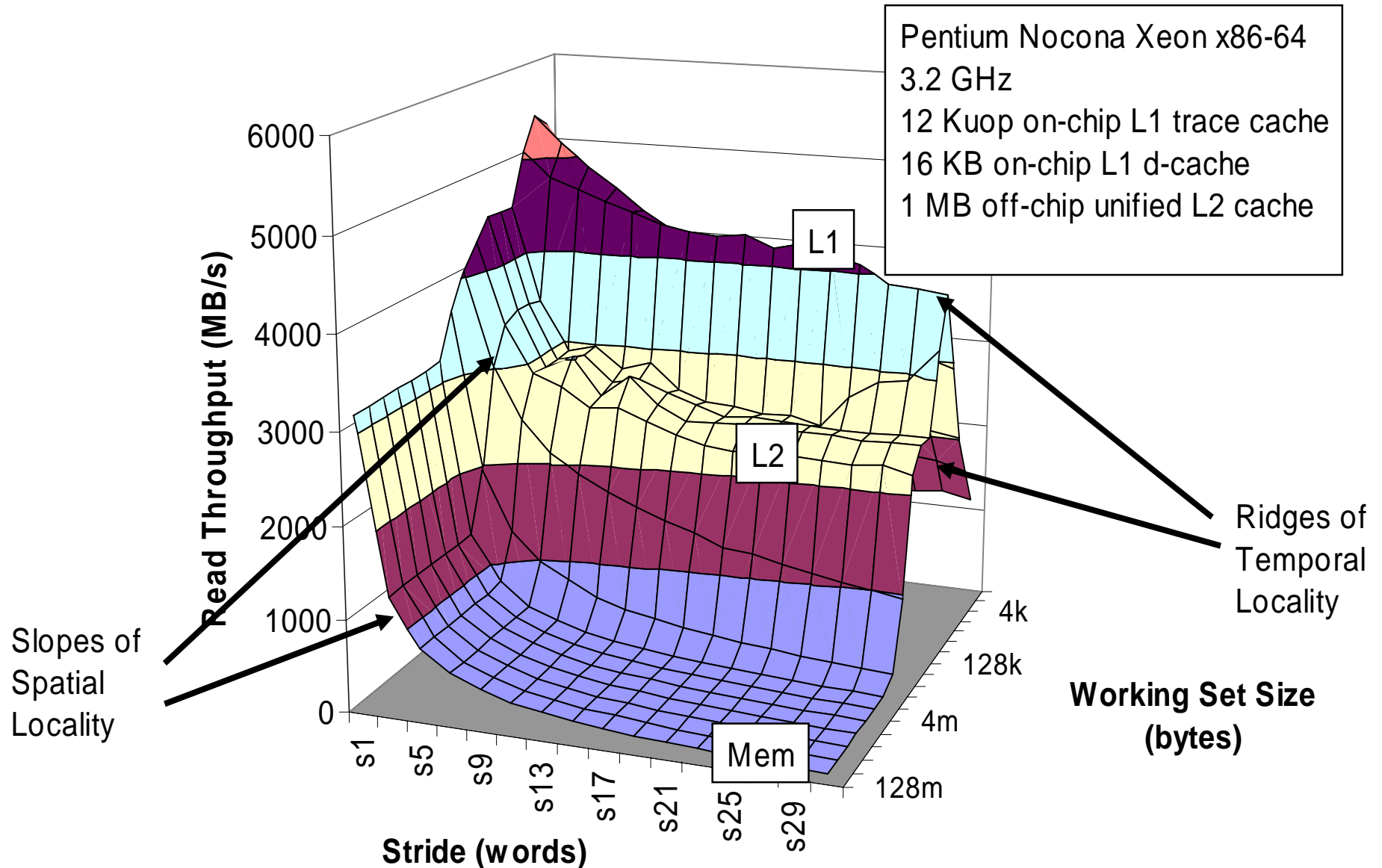
16 KB on-chip L1 i-cache

512 KB off-chip unified

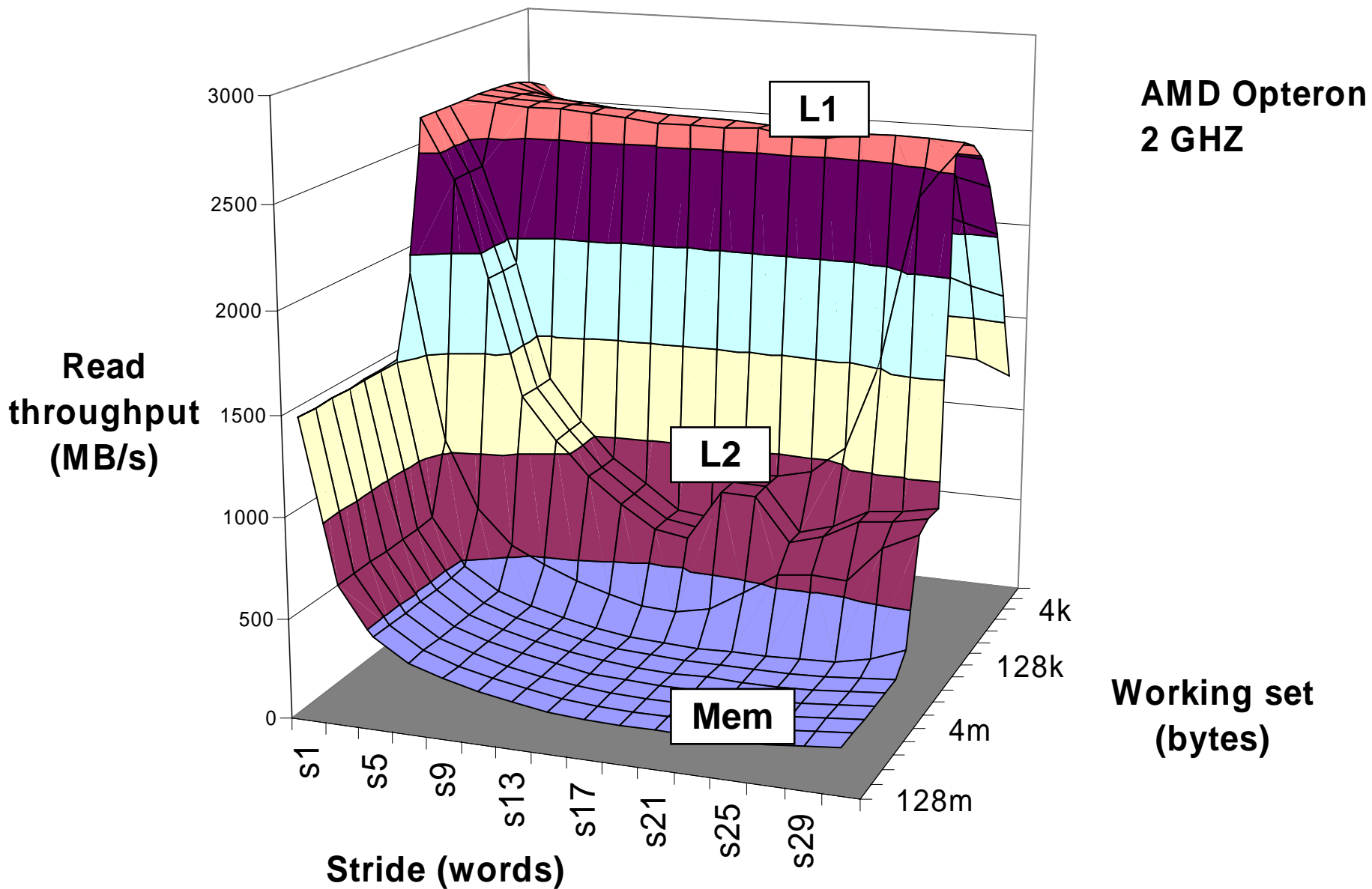
L2 cache



# X86-64 Memory Mountain

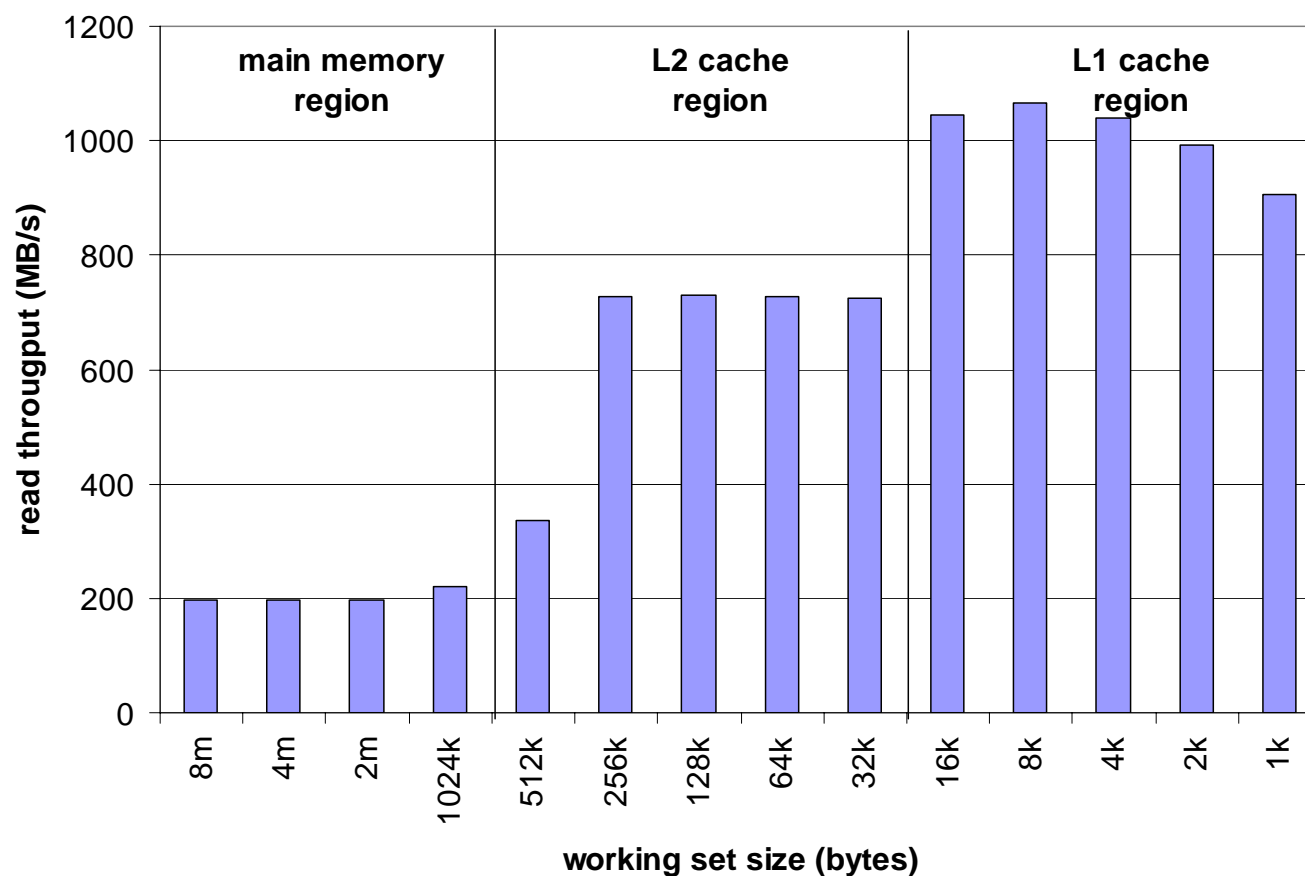


# Opteron Memory Mountain



# Ridges of Temporal Locality

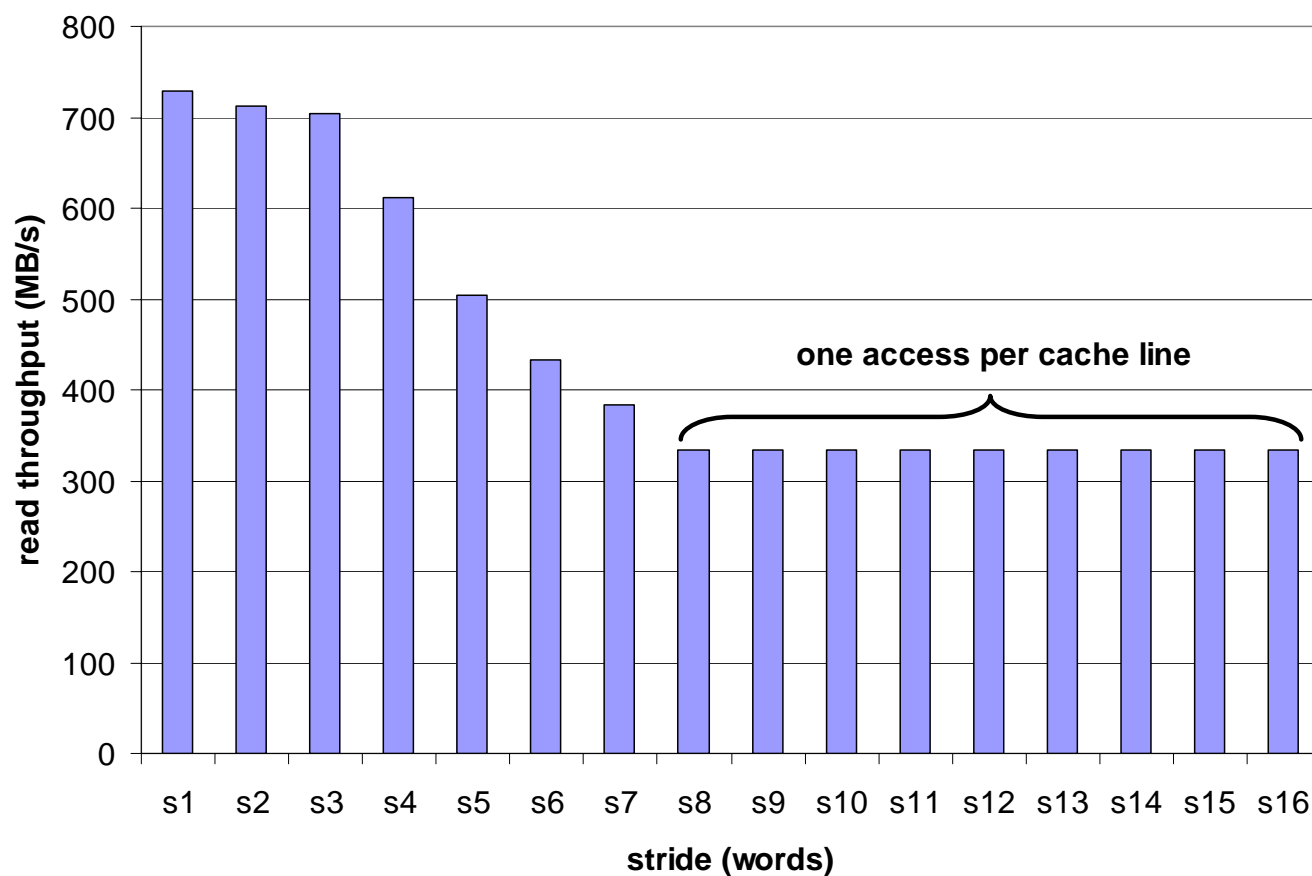
- Slice through the memory mountain with stride=1
  - illuminates read throughputs of different caches and memory



# A Slope of Spatial Locality

## ■ Slice through memory mountain with size=256KB

- shows cache block size.





# Today

- Cache organization
- Memory Wall
- **Program optimization (Matrix Multiplication):**
  - Cache optimizations
  - Cache Miss Analysis

# Matrix Multiplication Example

## ■ Major Cache Effects to Consider

- Total cache size
  - Exploit temporal locality and keep the working set small (e.g., use blocking)
- Block size
  - Exploit spatial locality

## ■ Description:

- Multiply  $N \times N$  matrices
- $O(N^3)$  total operations
- Accesses
  - $N$  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*

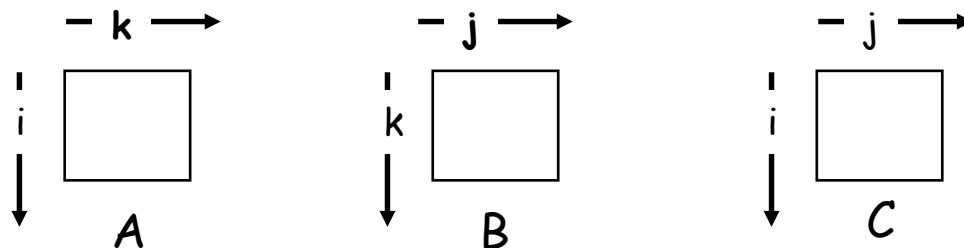
# Miss Rate Analysis for Matrix Multiply

## ■ 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



# Layout of C Arrays in Memory (review)

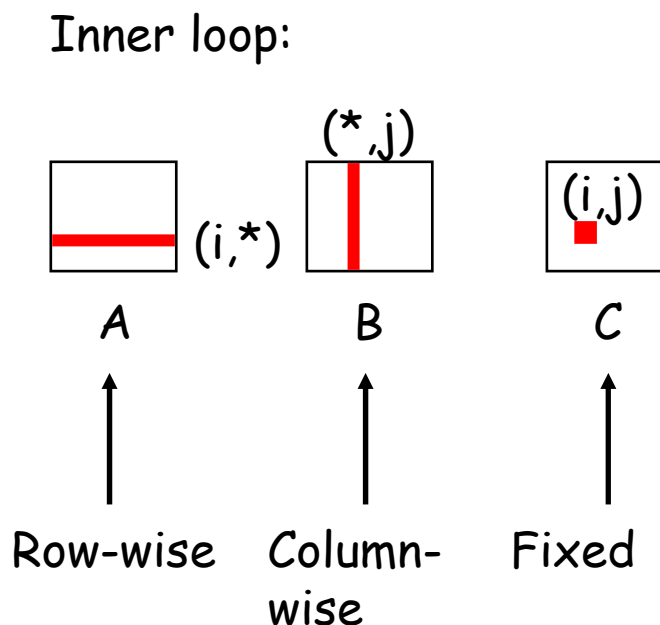
- **C arrays allocated in row-major order**
  - each row in contiguous memory locations
- **Stepping through columns in one row:**
  - `for (i = 0; i < N; i++)`  
`sum += a[0][i];`
  - accesses successive elements
  - if block size (B) > 4 bytes, exploit spatial locality
    - compulsory miss rate = 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;
  }
}

```



Misses per Inner Loop Iteration:

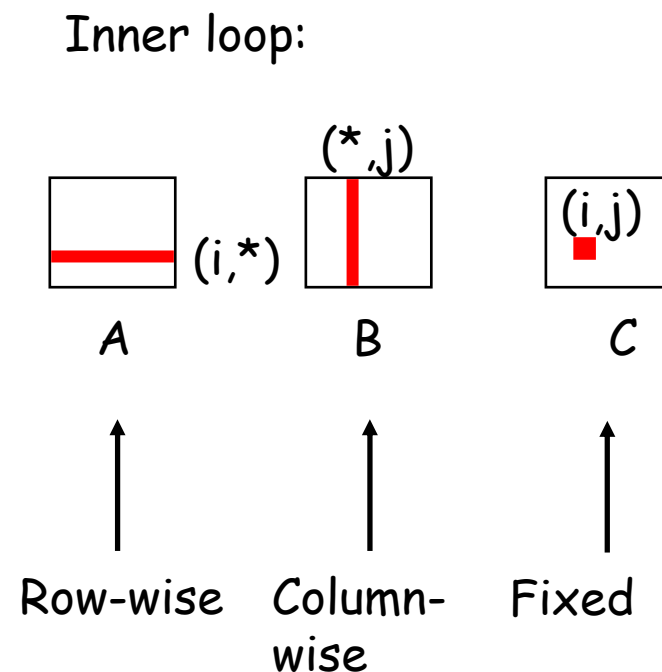
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

# 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
  }
}

```



Misses per Inner Loop Iteration:

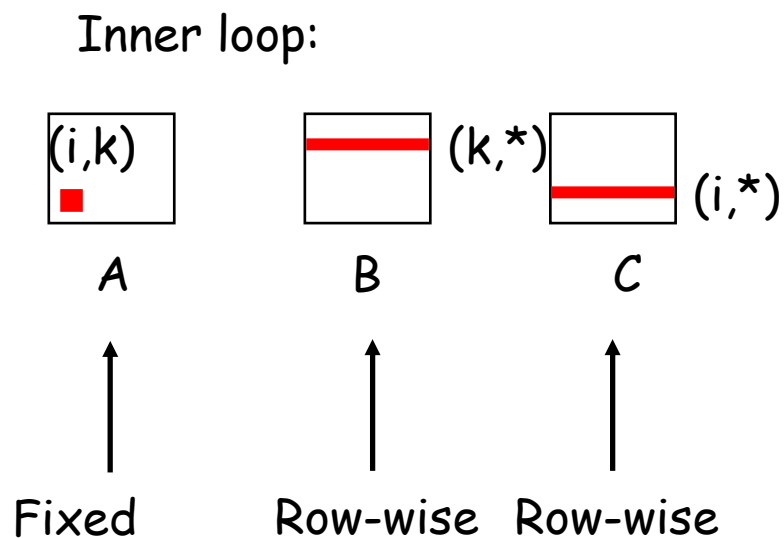
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

# 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];
  }
}

```



Misses per Inner Loop Iteration:

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

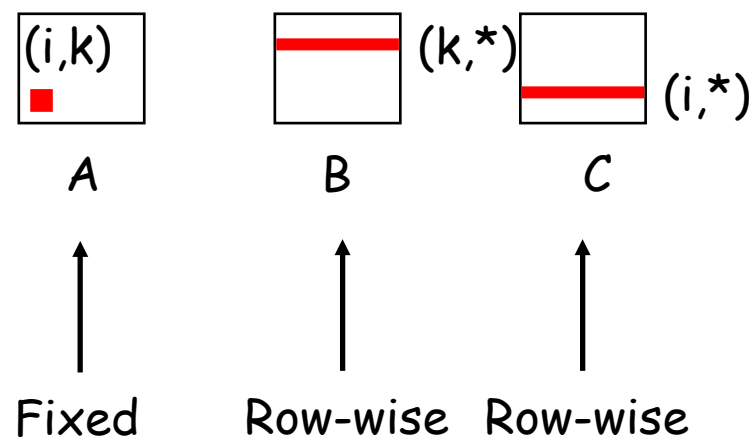
# Matrix Multiplication (ikj)

```

/* 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];
  }
}

```

Inner loop:



Misses per Inner Loop Iteration:

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



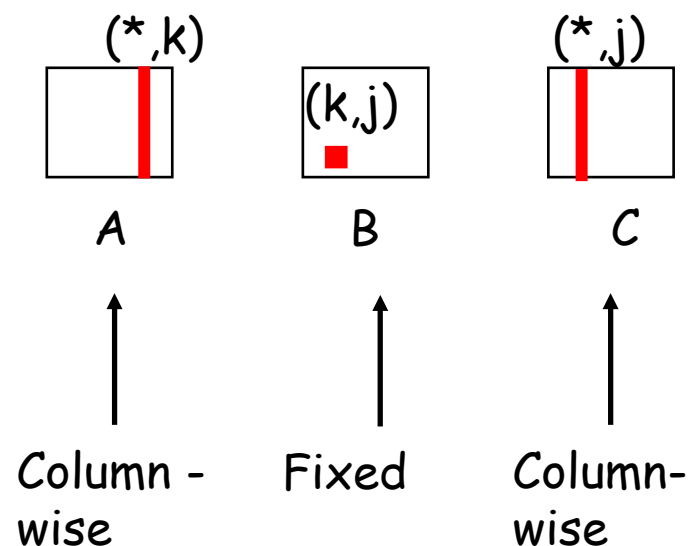
# Matrix Multiplication (jki)

```

/* 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:

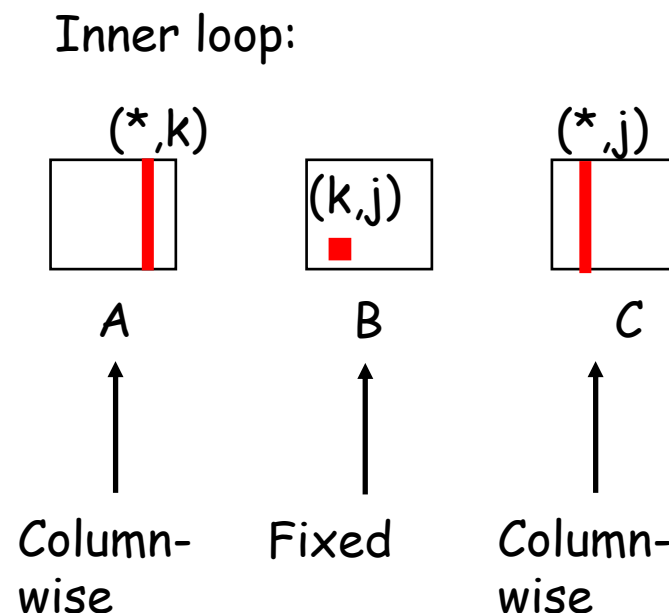
<u>A</u>	<u>B</u>	<u>C</u>
1.0	0.0	1.0

# 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:

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

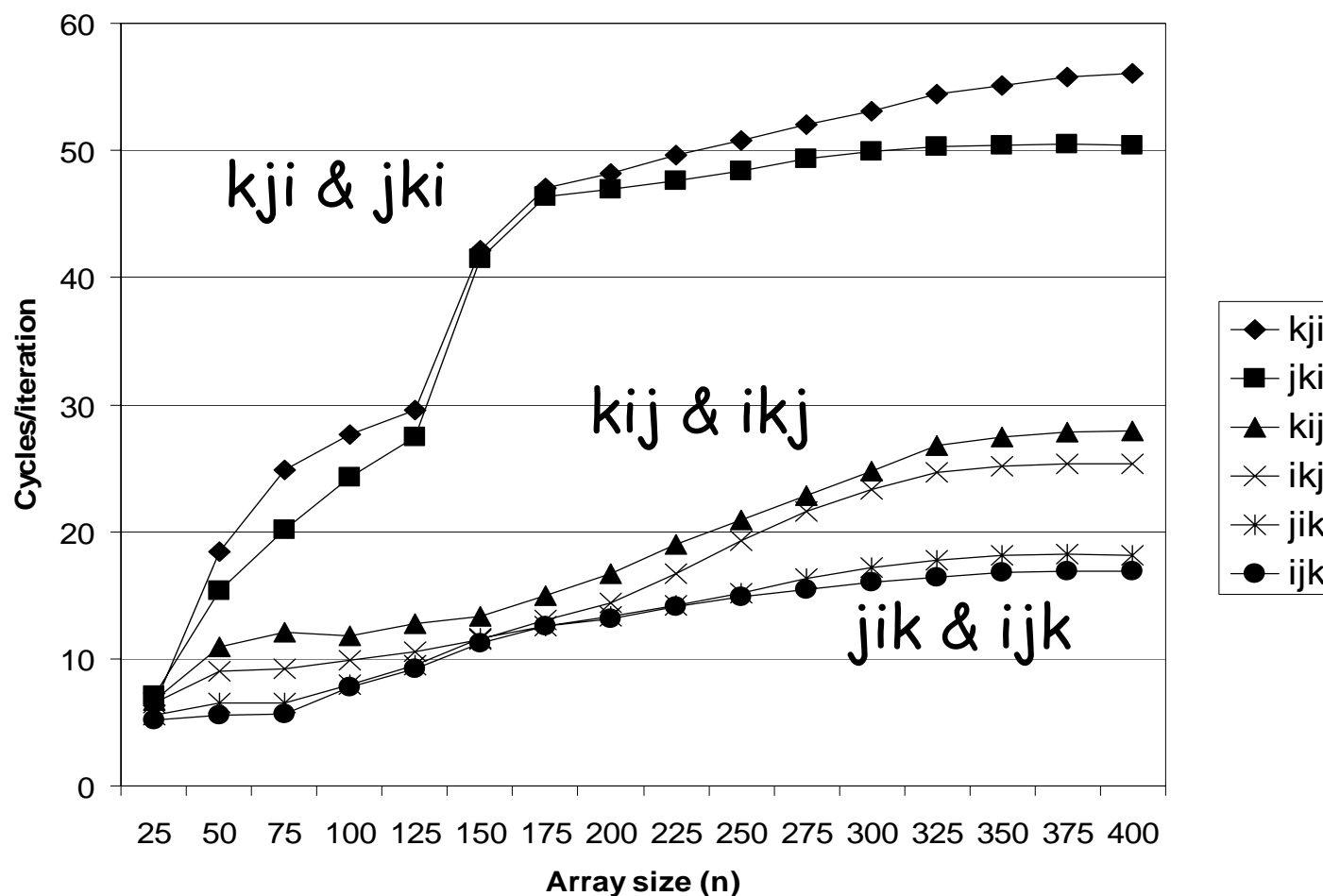
```

**jki (& kji):**

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

# Pentium Matrix Multiply Performance

- Miss rates are helpful but not perfect predictors.
  - Code scheduling matters, too.



# Improving Temporal Locality by Blocking

## ■ Example: Blocked matrix multiplication

- “block” (in this context) does not mean “cache block”.
- Instead, it means a sub-block within the matrix.
- Example:  $N = 8$ ; sub-block size = 4

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \times \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

Key idea: Sub-blocks (i.e.,  $A_{xy}$ ) can be treated just like scalars.

$$\begin{aligned} C_{11} &= A_{11}B_{11} + A_{12}B_{21} & C_{12} &= A_{11}B_{12} + A_{12}B_{22} \\ C_{21} &= A_{21}B_{11} + A_{22}B_{21} & C_{22} &= A_{21}B_{12} + A_{22}B_{22} \end{aligned}$$

# Blocked Matrix Multiply (bijk)

```
for (jj=0; jj<n; jj+=bsize) {  
  
    for (i=0; i<n; i++)  
        for (j=jj; j < min(jj+bsize,n); j++)  
            c[i][j] = 0.0;  
  
    for (kk=0; kk<n; kk+=bsize) {  
        for (i=0; i<n; i++) {  
            for (j=jj; j < min(jj+bsize,n); j++) {  
                sum = 0.0  
                for (k=kk; k < min(kk+bsize,n); k++) {  
                    sum += a[i][k] * b[k][j];  
                }  
                c[i][j] += sum;  
            }  
        }  
    }  
}
```

# Blocked Matrix Multiply Analysis

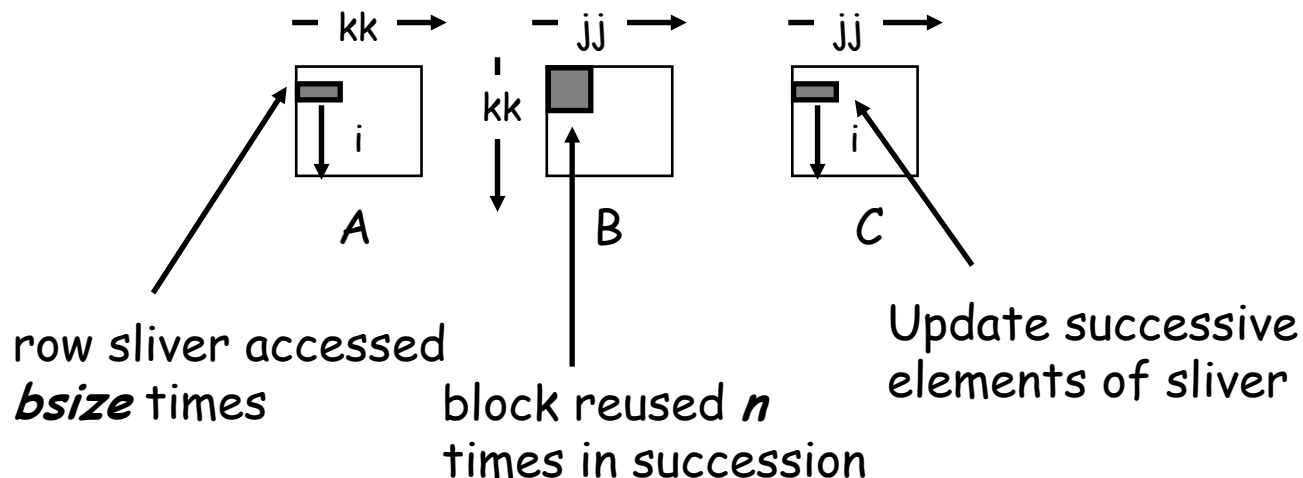
- Innermost loop pair multiplies a  $1 \times bsize$  sliver of  $A$  by a  $bsize \times bsize$  block of  $B$  and accumulates into  $1 \times bsize$  sliver of  $C$
- Loop over  $i$  steps through  $n$  row slivers of  $A$  &  $C$ , using same  $B$

```

for (i=0; i<n; i++) {
  for (j=jj; j < min(jj+bsize,n); j++) {
    sum = 0.0
    for (k=kk; k < min(kk+bsize,n); k++) {
      sum += a[i][k] * b[k][j];
    }
    c[i][j] += sum;
  }
}

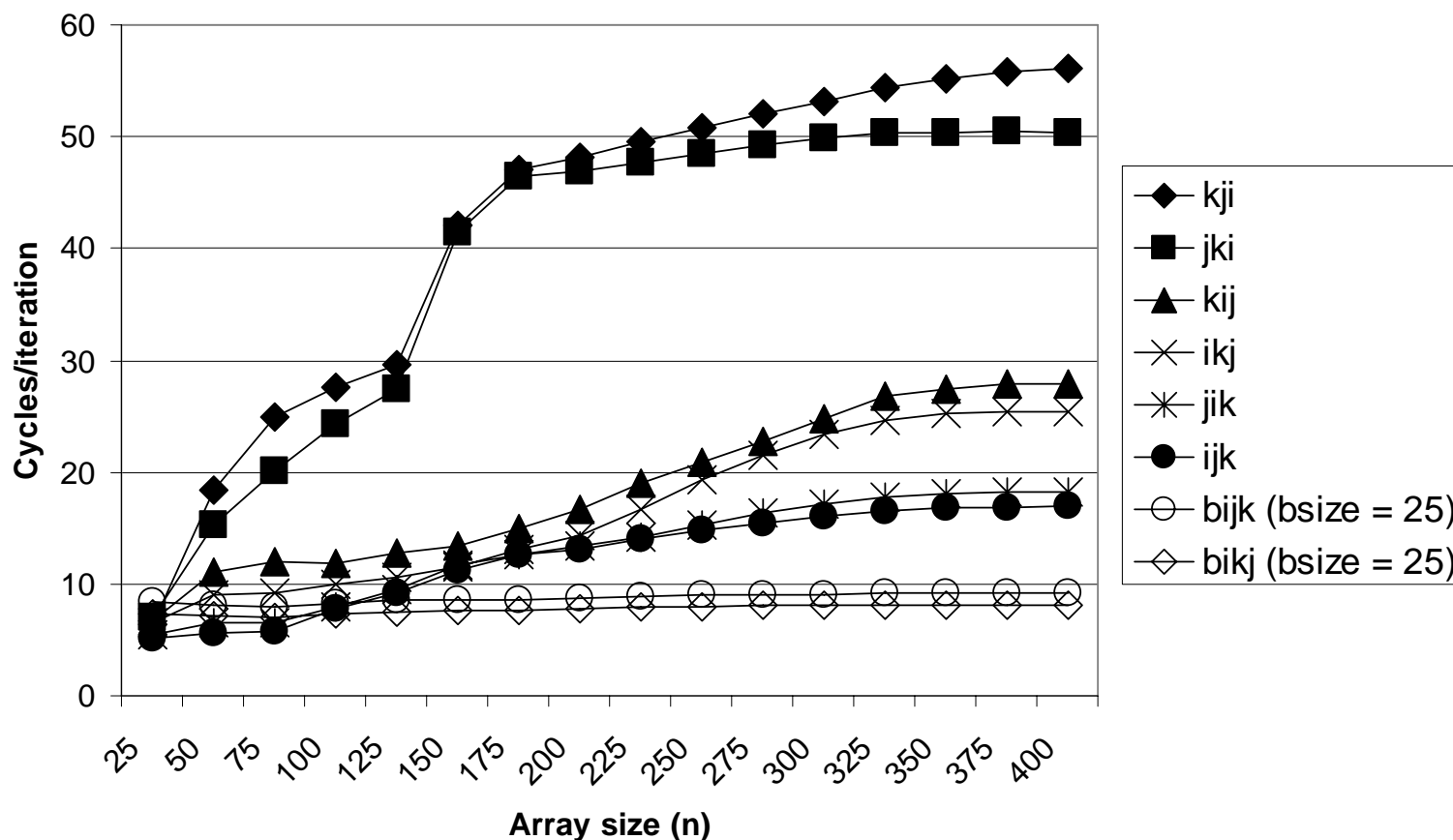
```

Innermost  
Loop Pair



# Pentium Blocked Matrix Multiply Performance

- **Blocking (bijk and bikj) improves performance by a factor of two over unblocked versions (ijk and jik)**
  - relatively insensitive to array size.





# 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)

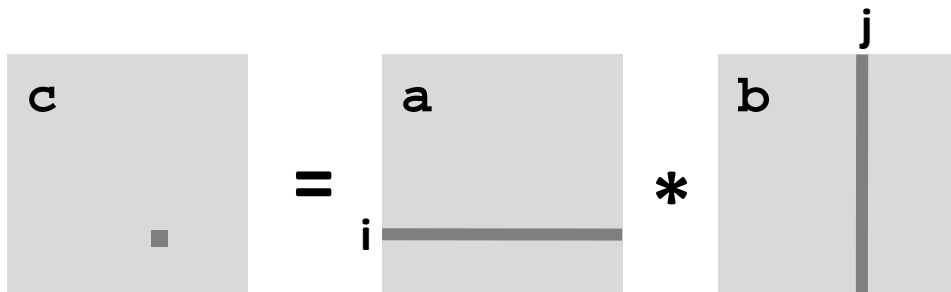
# Today

- Cache organization
- Memory Wall
- **Program optimization (Matrix Multiplication):**
  - Cache optimizations
  - Cache Miss Analysis

# 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];
}
```



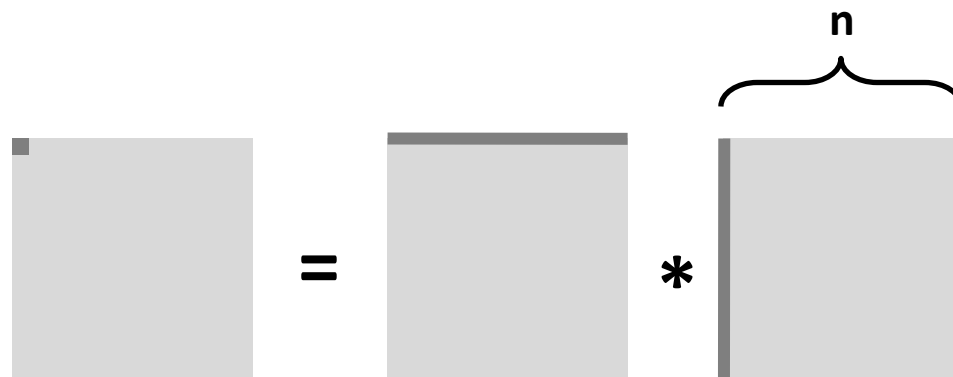
# Cache Miss Analysis

## ■ Assume:

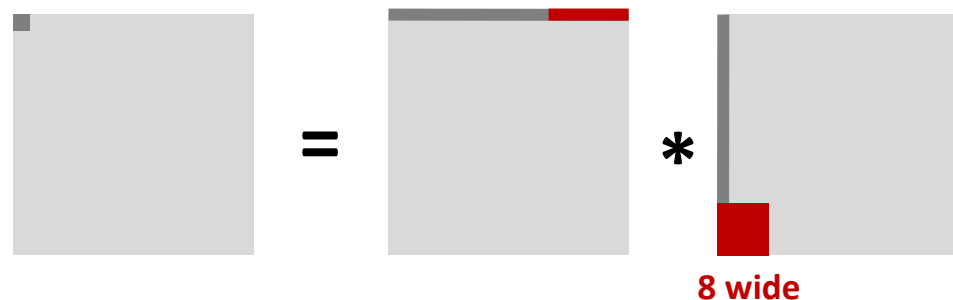
- Matrix elements are doubles
- Cache block = 8 doubles (64 B as in Core 2 Duo)
- Cache size  $C \ll n$  (much smaller than  $n$ )

## ■ First iteration:

- $n/8 + n = 9n/8$  misses



- Afterwards **in cache**:  
(schematic)



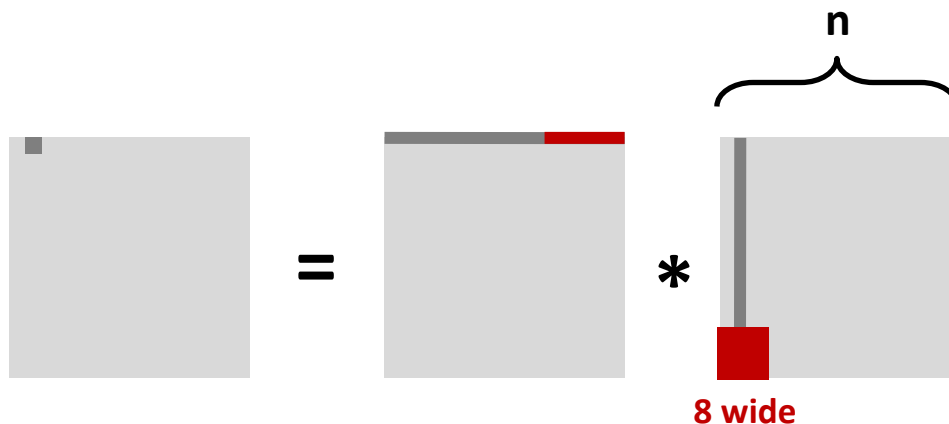
# Cache Miss Analysis

## ■ Assume:

- Matrix elements are doubles
- Cache block = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )

## ■ Second iteration:

- Again:  
 $n/8 + n = 9n/8$  misses



## ■ Total misses:

- $9n/8 * n^2 = (9/8) * n^3$

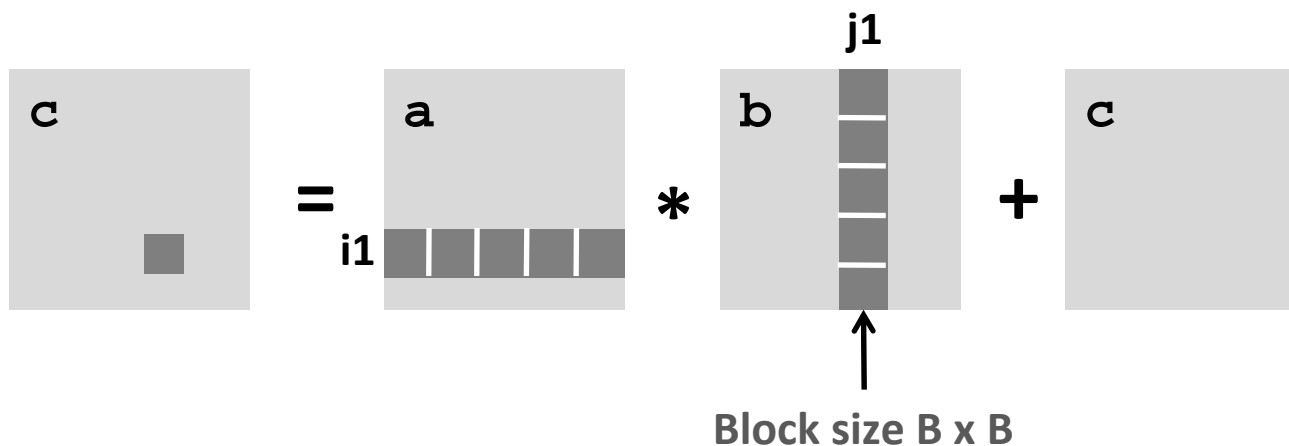
# Blocked Matrix Multiplication

```

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];
}

```



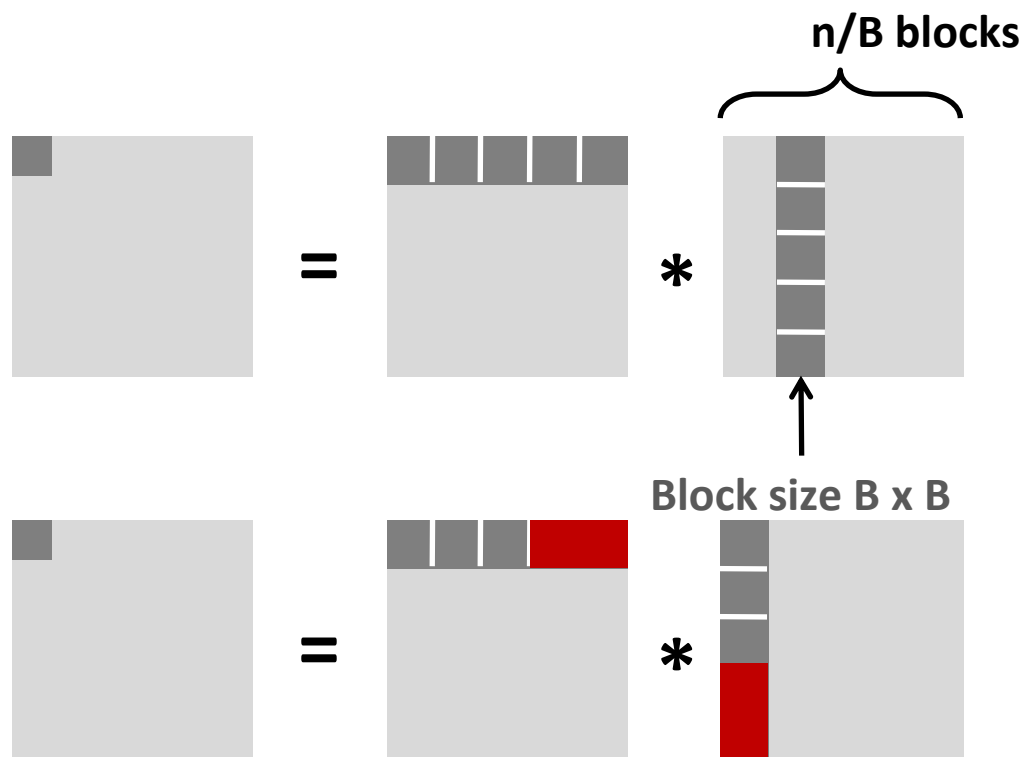
# Cache Miss Analysis

## ■ 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  $c$ )



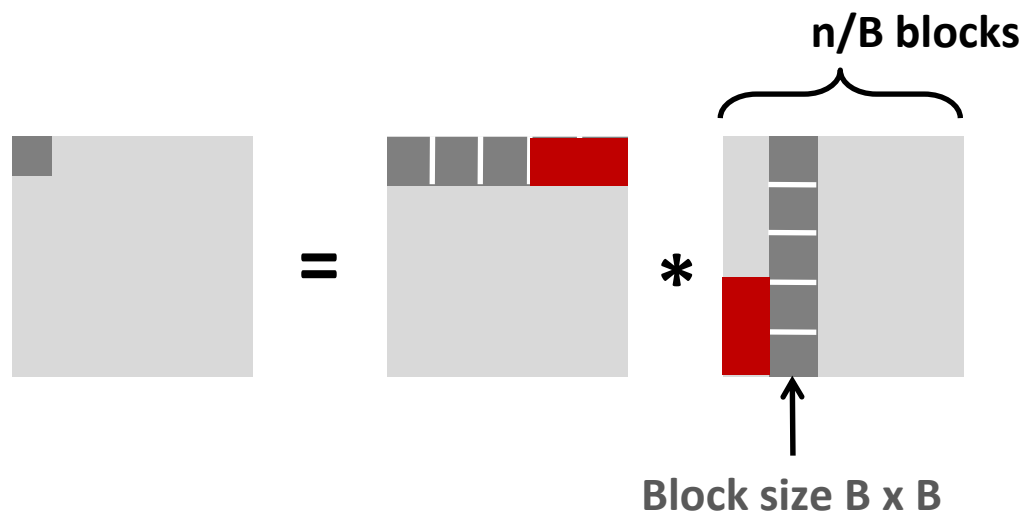
# Cache Miss Analysis

## ■ Assume:

- Cache block = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  $\blacksquare$  fit into cache:  $3B^2 < C$

## ■ Second (block) iteration:

- Same as first iteration
- $2n/B * B^2/8 = nB/4$



## ■ Total misses:

- $nB/4 * (n/B)^2 = n^3/(4B)$



# Summary

- **No blocking:**  $(9/8) * n^3$
- **Blocking:**  $1/(4B) * n^3$
  
- **Suggest largest possible block size B, but limit  $3B^2 < C!$**   
(can possibly be relaxed a bit, but there is a limit for B)
  
- **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