

# Roadmap

C:

```
car *c = malloc(sizeof(car));
c->miles = 100;
c->gals = 17;
float mpg = get_mpg(c);
free(c);
```

Java:

```
Car c = new Car();
c.setMiles(100);
c.setGals(17);
float mpg =
    c.getMPG();
```

Assembly  
language:

```
get_mpg:
    pushq   %rbp
    movq   %rsp, %rbp
    ...
    popq   %rbp
    ret
```

Machine  
code:

```
0111010000011000
100011010000010000000010
1000100111000010
110000011111101000011111
```

Computer  
system:



Caches

OS:



Memory & data  
Integers & floats  
Machine code & C  
x86 assembly  
Procedures & stacks  
Arrays & structs  
**Memory & caches**  
Processes  
Virtual memory  
Memory allocation  
Java vs. C



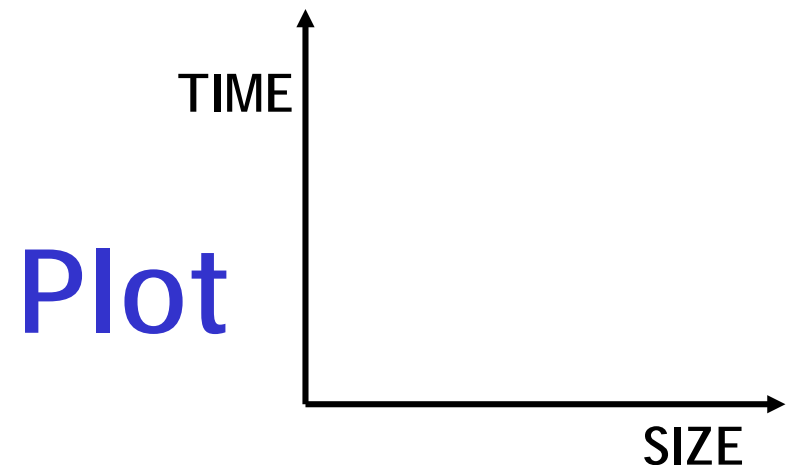
# Section 7: Memory and Caches

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches

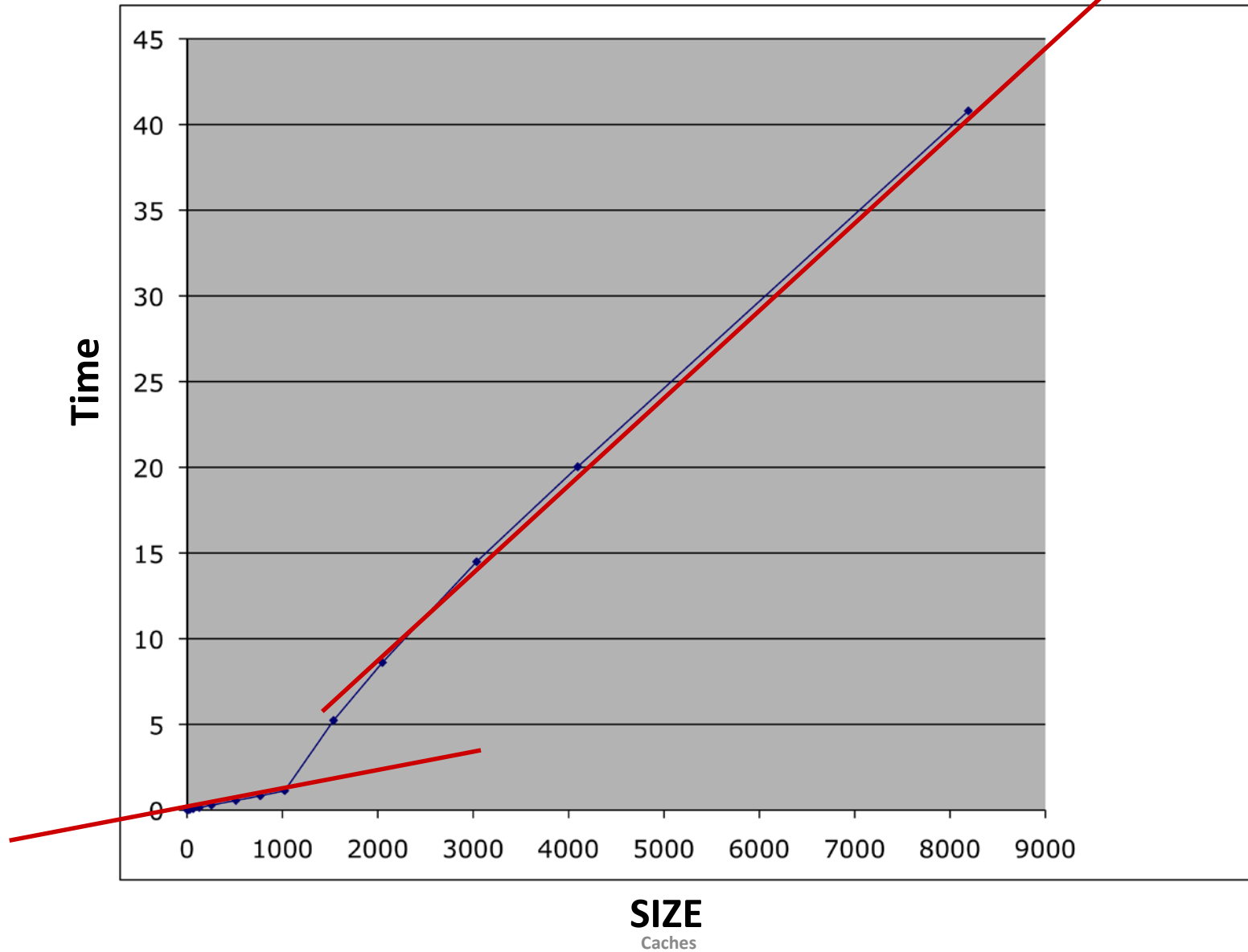
# How does execution time grow with SIZE?

```
int array[SIZE];
int A = 0;

for (int i = 0 ; i < 200000 ; ++ i) {
    for (int j = 0 ; j < SIZE ; ++ j) {
        A += array[j];
    }
}
```

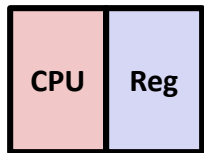


# Actual Data



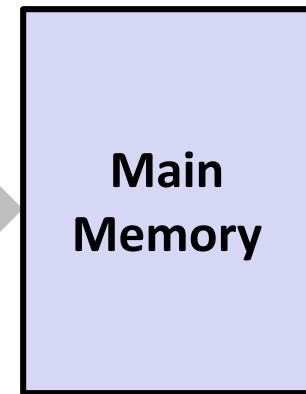
# Problem: Processor-Memory Bottleneck

Processor performance  
doubled about  
every 18 months



**Core 2 Duo:**  
Can process at least  
256 Bytes/cycle

Bus bandwidth  
evolved much slower

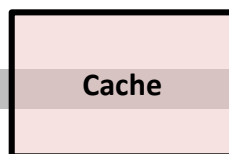
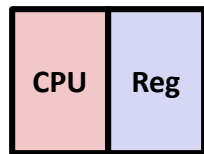


**Core 2 Duo:**  
Bandwidth  
2 Bytes/cycle  
Latency  
100 cycles

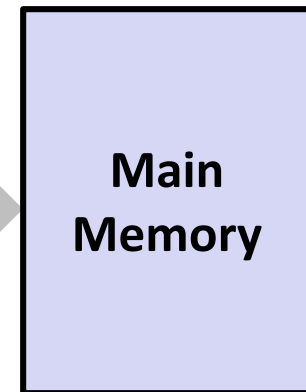
***Problem: lots of waiting on memory***

# Problem: Processor-Memory Bottleneck

Processor performance  
doubled about  
every 18 months



Bus bandwidth  
evolved much slower



**Core 2 Duo:**  
Can process at least  
256 Bytes/cycle

**Core 2 Duo:**  
Bandwidth  
2 Bytes/cycle  
Latency  
100 cycles

***Solution: caches***

# Cache

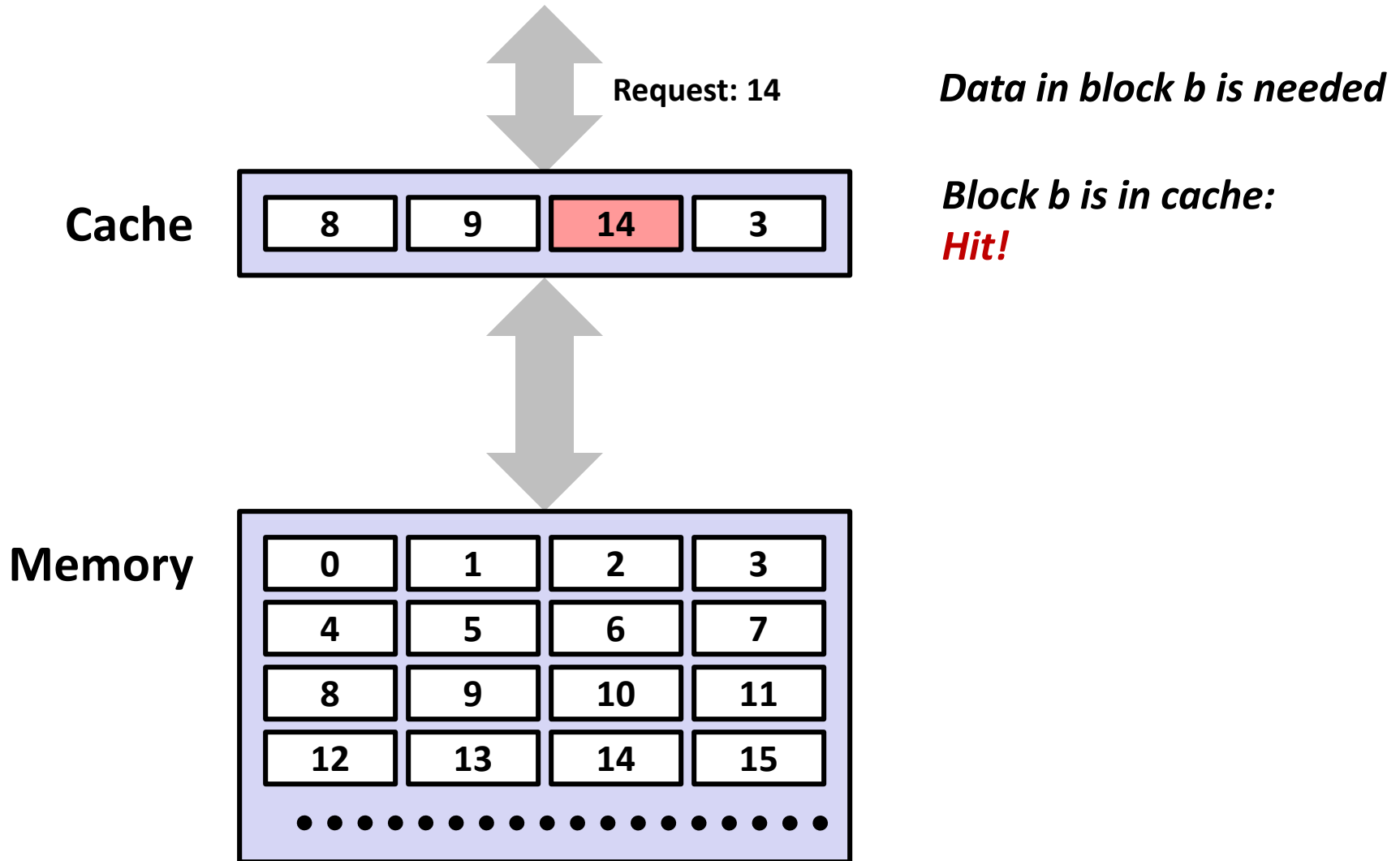
- **English definition:** a hidden storage space for provisions, weapons, and/or treasures
- **CSE definition:** computer memory with short access time used for the storage of frequently or recently used instructions or data (i-cache and d-cache)

more generally,

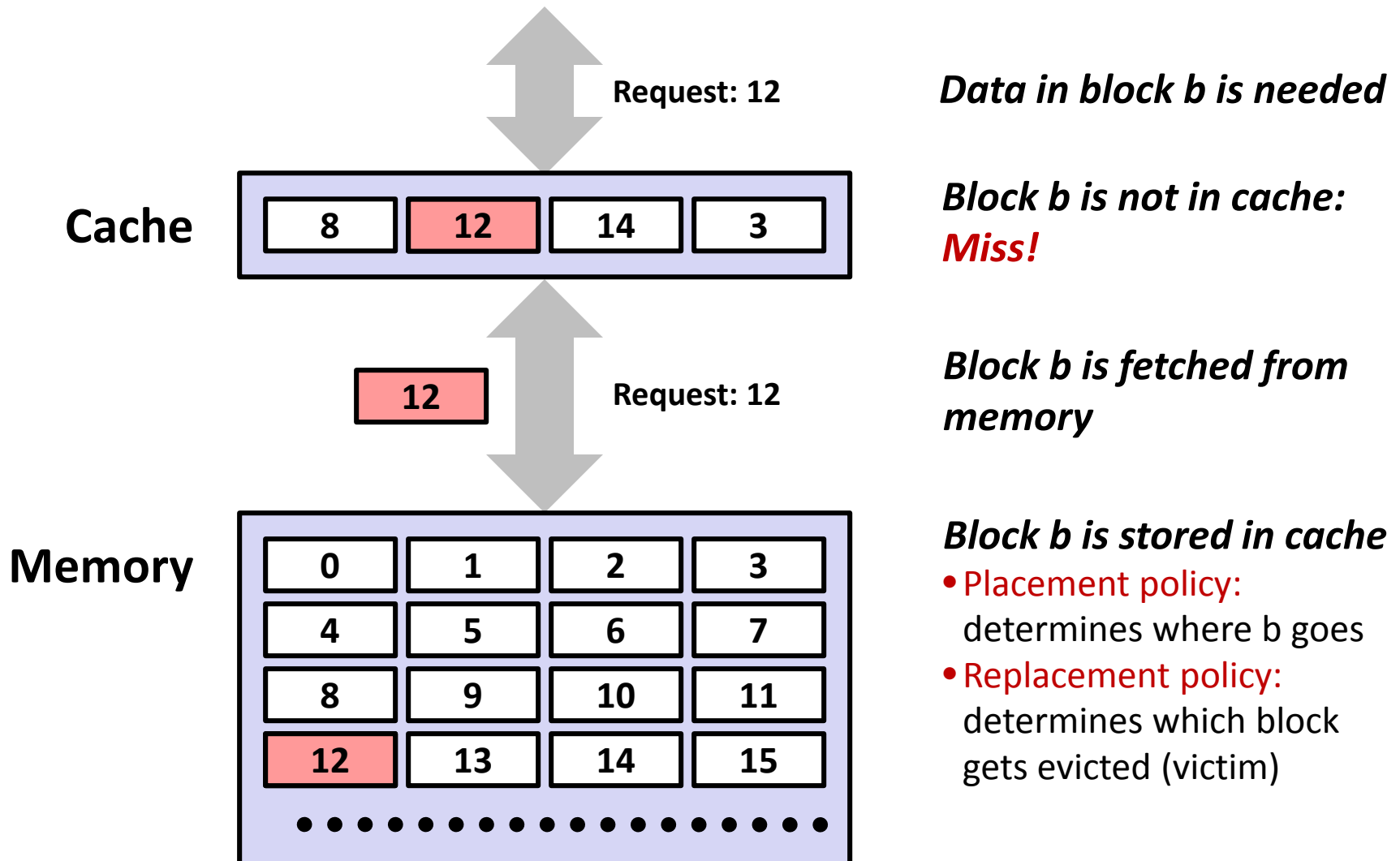
used to optimize data transfers between system elements with different characteristics (network interface cache, I/O cache, etc.)



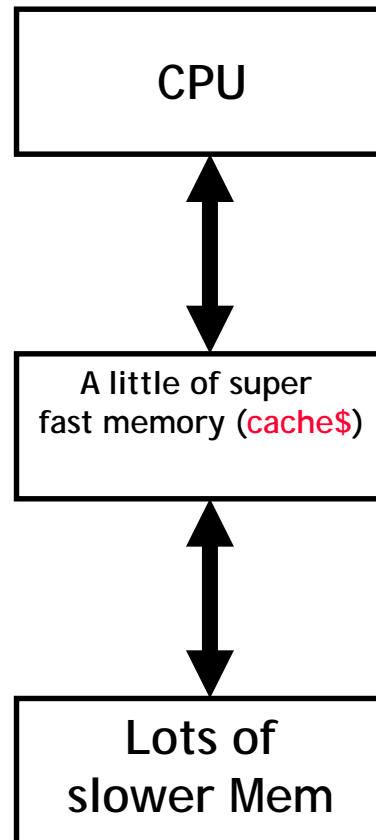
# General Cache Concepts: **Hit**



# General Cache Concepts: Miss



# Not to forget...



# Section 7: Memory and Caches

- ~~Cache basics~~
- Principle of locality
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches

# Why Caches Work

- **Locality:** Programs tend to use data and instructions with addresses near or equal to those they have used recently

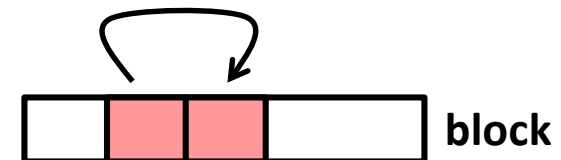
- **Temporal locality:**

- 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
- How do caches take advantage of this?



# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

## ■ Data:

- Temporal: **sum** referenced in each iteration
- Spatial: array **a[ ]** accessed in stride-1 pattern

# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

## ■ Data:

- Temporal: **sum** referenced in each iteration
- Spatial: array **a[ ]** accessed in stride-1 pattern

## ■ Instructions:

- Temporal: cycle through loop repeatedly
- Spatial: reference instructions in sequence

# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

## ■ Data:

- Temporal: `sum` referenced in each iteration
- Spatial: array `a[ ]` accessed in stride-1 pattern

## ■ Instructions:

- Temporal: cycle through loop repeatedly
- Spatial: reference instructions in sequence

- **Being able to assess the locality of code is a crucial skill for a programmer**

# Another Locality Example

```
int sum_array_3d(int a[M][N][N])
{
    int i, j, k, sum = 0;

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

- What is wrong with this code?
- How can it be fixed?

# Section 7: Memory and Caches

- ~~Cache basics~~
- ~~Principle of locality~~
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches

# Cost of Cache Misses

- **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:  $1 \text{ cycle} + 0.03 * 100 \text{ cycles} = 4 \text{ cycles}$
    - 99% hits:  $1 \text{ cycle} + 0.01 * 100 \text{ cycles} = 2 \text{ cycles}$
- **This is why “miss rate” is used instead of “hit rate”**

# 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

## ■ 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 hit times: 1 - 2 clock cycles for L1

## ■ Miss Penalty

- Additional time required because of a miss
- Typically 50 - 200 cycles

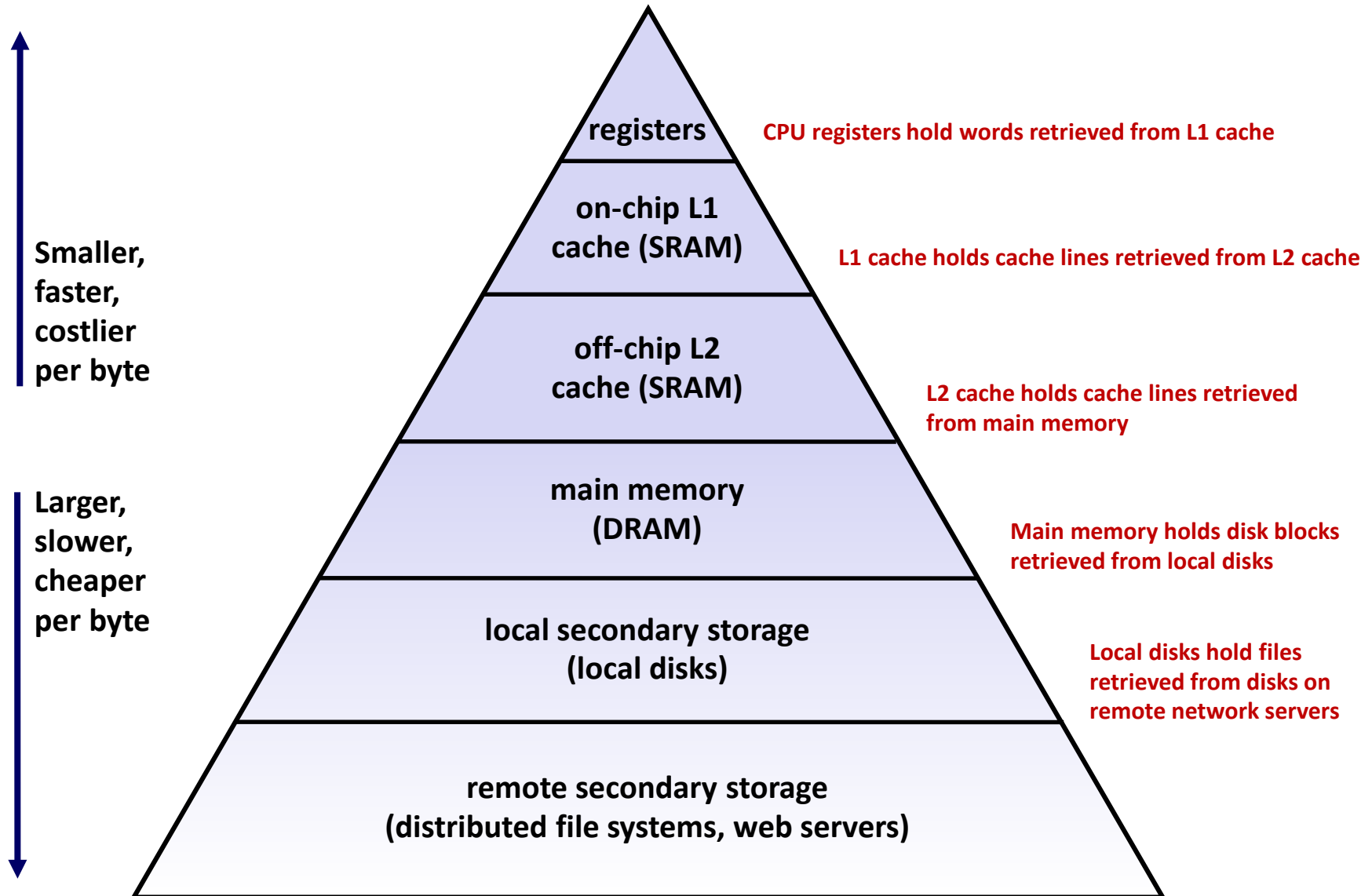
# Memory Hierarchies

- **Some fundamental and enduring properties of hardware and software systems:**
  - Faster storage technologies almost always cost more per byte and have lower capacity
  - The gaps between memory technology speeds are widening
    - True for: registers  $\leftrightarrow$  cache, cache  $\leftrightarrow$  DRAM, DRAM  $\leftrightarrow$  disk, etc.
  - Well-written programs tend to exhibit good locality
- **These properties complement each other beautifully**
- **They suggest an approach for organizing memory and storage systems known as a memory hierarchy**

# Memory Hierarchies

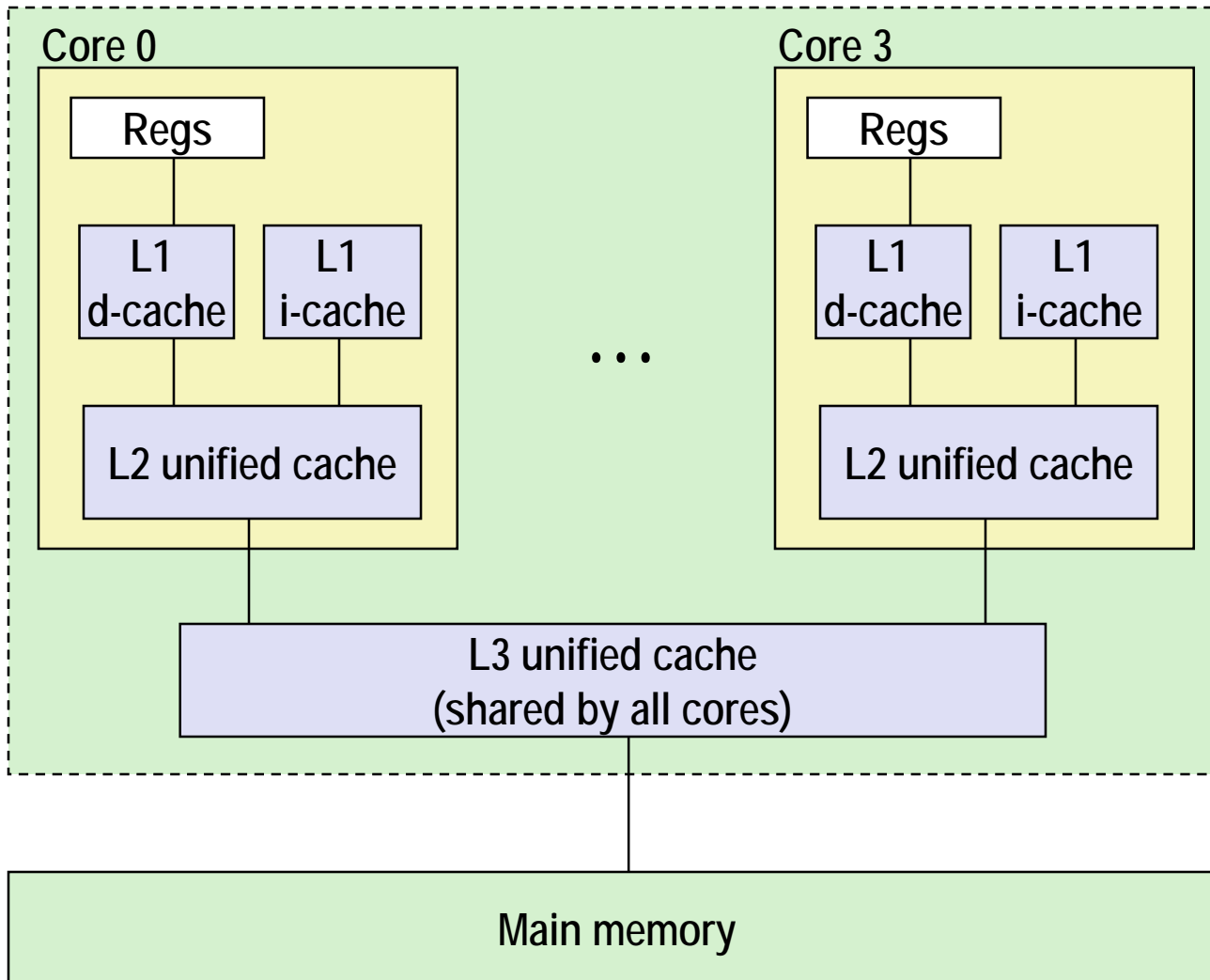
- **Fundamental idea of a memory hierarchy:**
  - Each level  $k$  serves as a cache for the larger, slower, level  $k+1$  below.
- **Why do memory hierarchies work?**
  - Because of locality, programs tend to access the data at level  $k$  more often than they access the data at level  $k+1$ .
  - Thus, the storage at level  $k+1$  can be slower, and thus larger and cheaper per bit.
- ***Big Idea:* The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.**

# An Example Memory Hierarchy



# Intel Core i7 Cache Hierarchy

Processor package



**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.

# Section 7: Memory and Caches

- ~~Cache basics~~
- ~~Principle of locality~~
- ~~Memory hierarchies~~
- ~~Cache organization~~
- Program optimizations that consider caches

# Optimizations for the Memory Hierarchy

- **Write code that has locality**
  - Spatial: access data contiguously
  - Temporal: make sure access to the same data is not too far apart in time
- **How to achieve?**
  - Proper choice of algorithm
  - Loop transformations

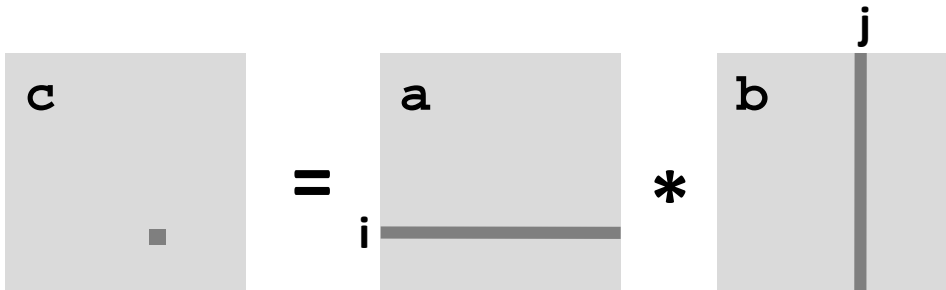
# Example: Matrix Multiplication

```

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];
}

```



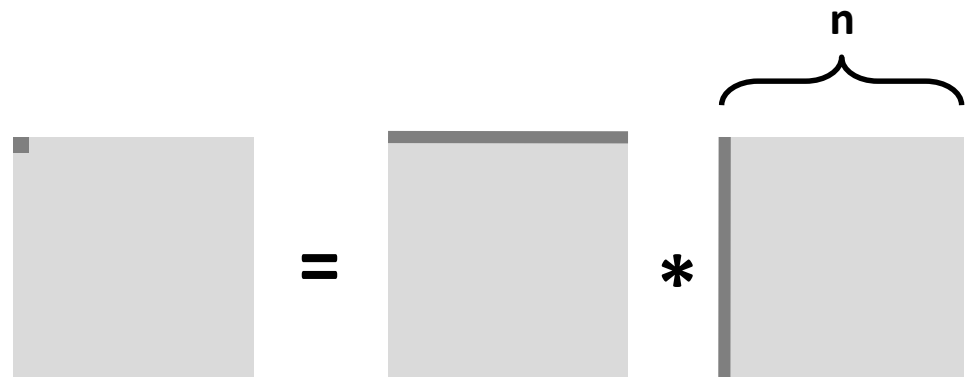
# Cache Miss Analysis

## ■ Assume:

- Matrix elements are doubles
- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )

## ■ First iteration:

- $n/8 + n = 9n/8$  misses  
(omitting matrix  $c$ )



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



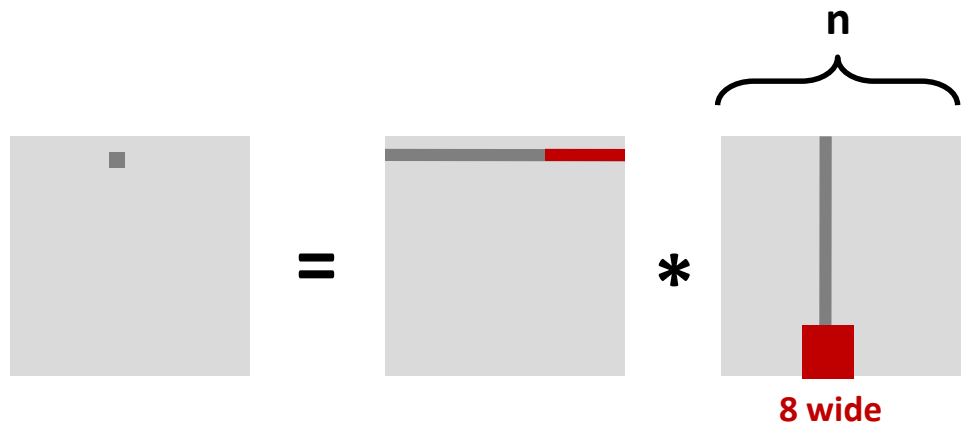
# Cache Miss Analysis

## ■ Assume:

- Matrix elements are doubles
- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )

## ■ Other iterations:

- Again:  
 $n/8 + n = 9n/8$  misses  
 (omitting matrix  $c$ )



## ■ Total misses:

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

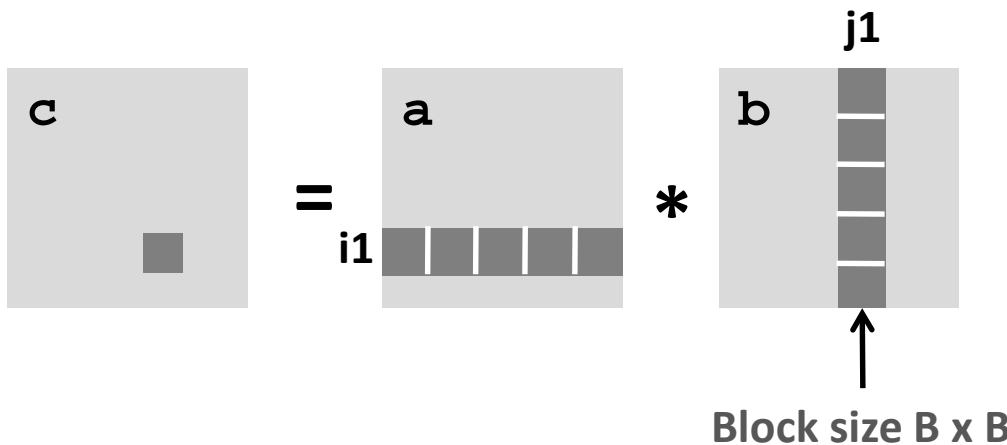
# Blocked Matrix Multiplication

```

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; i1++)
                    for (j1 = j; j1 < j+B; j1++)
                        for (k1 = k; k1 < k+B; k1++)
                            c[i1*n + j1] += a[i1*n + k1]*b[k1*n + j1];
}

```



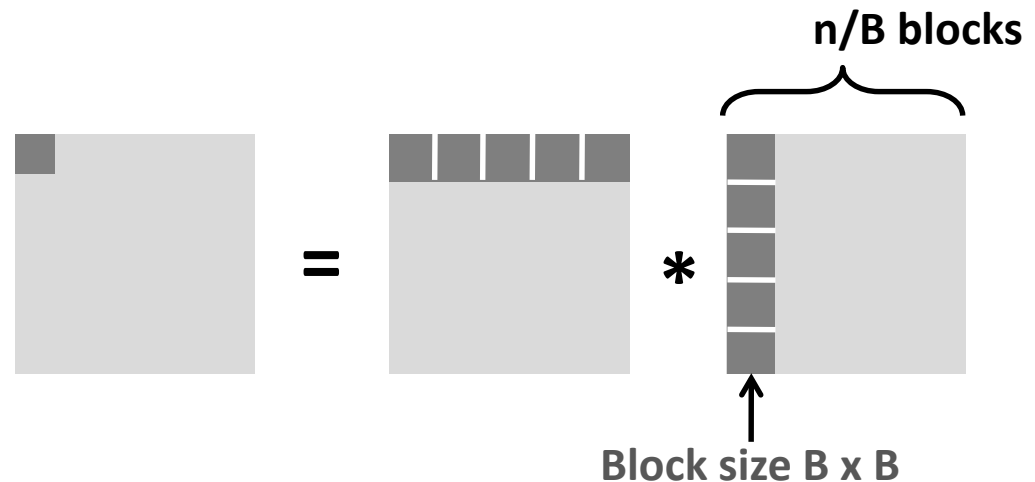
# Cache Miss Analysis

## ■ Assume:

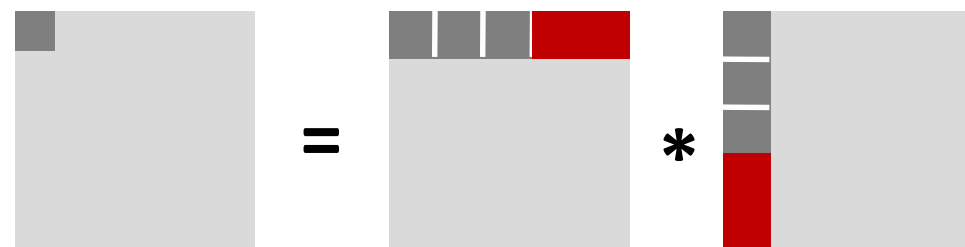
- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  fit into cache:  $3B^2 < C$

## ■ First (block) iteration:

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




- Afterwards in cache (schematic)



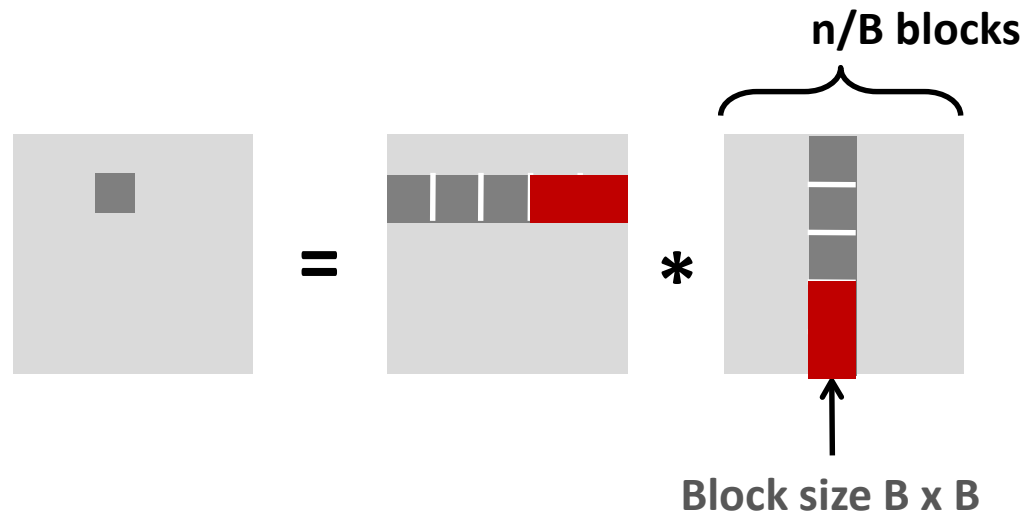
# Cache Miss Analysis

## ■ Assume:

- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  fit into cache:  $3B^2 < C$

## ■ Other (block) iterations:

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



## ■ Total misses:

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

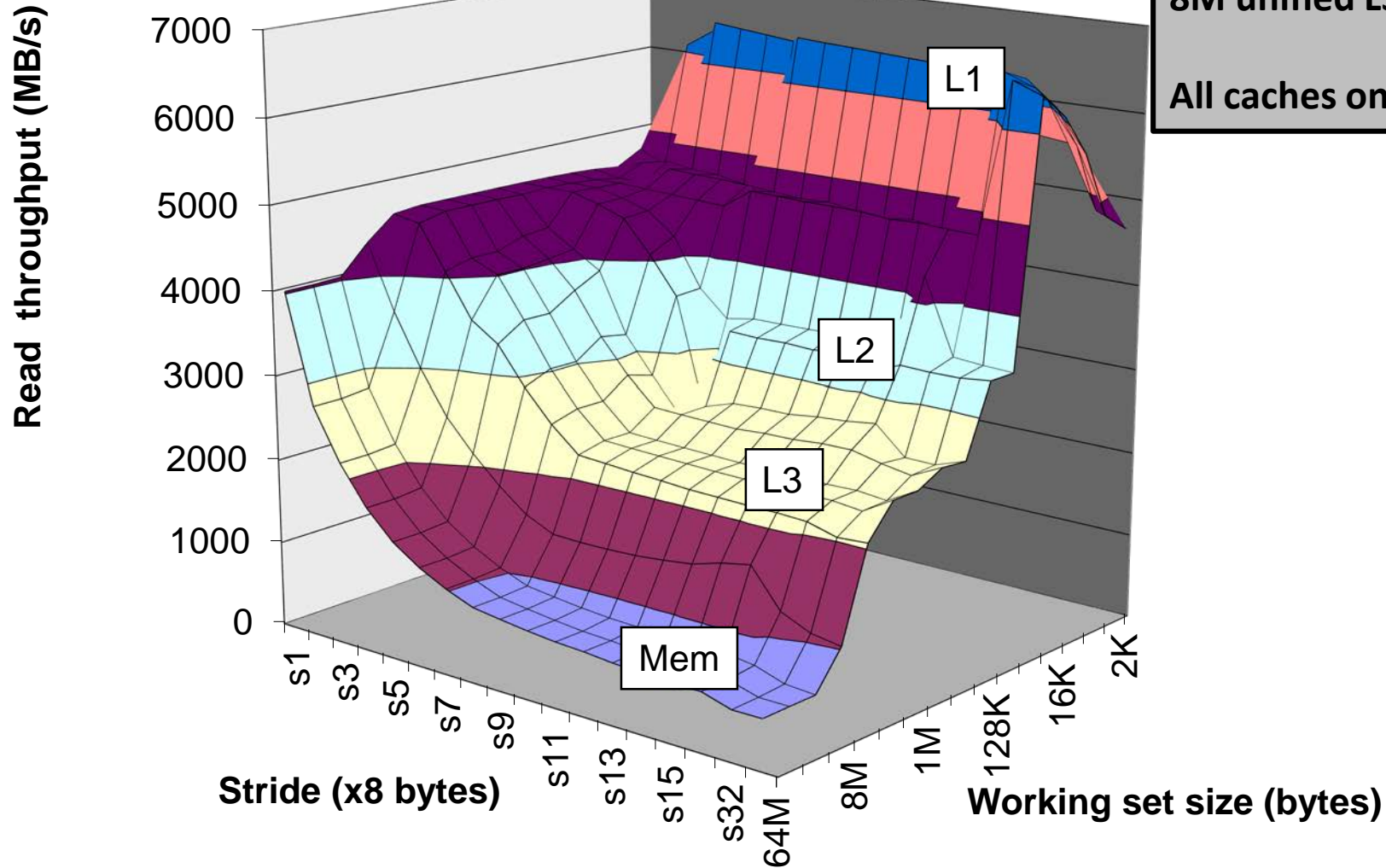
# Summary

- No blocking:  $(9/8) * n^3$
- Blocking:  $1/(4B) * n^3$
- If  $B = 8$  difference is  $4 * 8 * 9 / 8 = 36x$
- If  $B = 16$  difference is  $4 * 16 * 9 / 8 = 72x$
- Suggests 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

# Cache-Friendly Code

- **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)
    - Focus on inner loop code

# 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