

Optimization III: Cache Memories

15-213/18-243: Introduction to Computer Systems

12th Lecture, 22 February 2010

Instructors:

Bill Nace and Gregory Kesden

The Exam

- **Tuesday, 2 March**
- **Closed book, closed notes, closed friend, open mind**
 - We will provide reference material
- **Quite unlike past exams**
- **Material from Lectures 1 - 9, Labs 1 & 2**
 - Representation of Integers, Floats
 - Machine code for control structures, procedures
 - Stack discipline
 - Layout of Arrays, Structs, Unions in memory
 - Floating point operations

Last Time

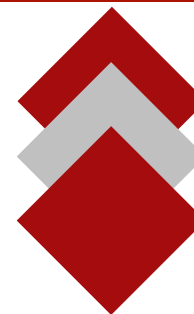
■ Program optimization

- Optimization blocker: Memory aliasing
- One solution: Scalar replacement of array accesses that are reused

```
for (i = 0; i < n; i++) {  
    b[i] = 0;  
    for (j = 0; j < n; j++)  
        b[i] += a[i*n + j];  
}
```

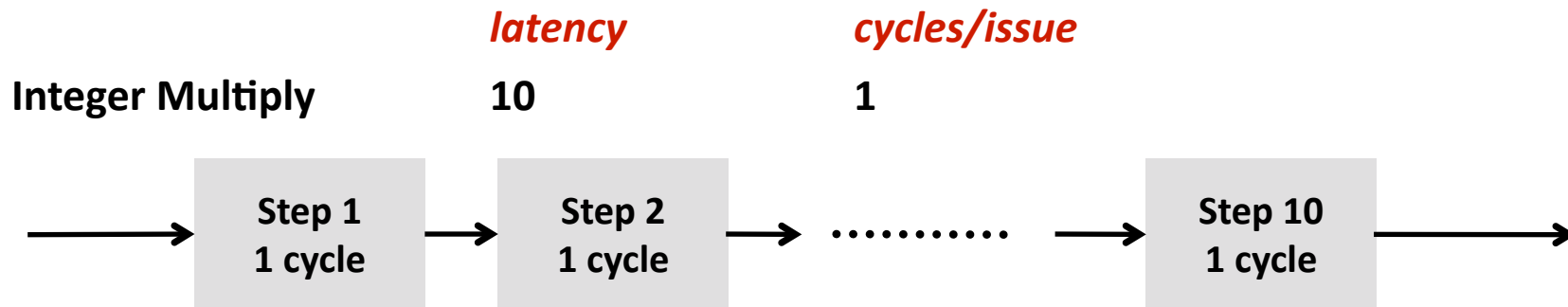
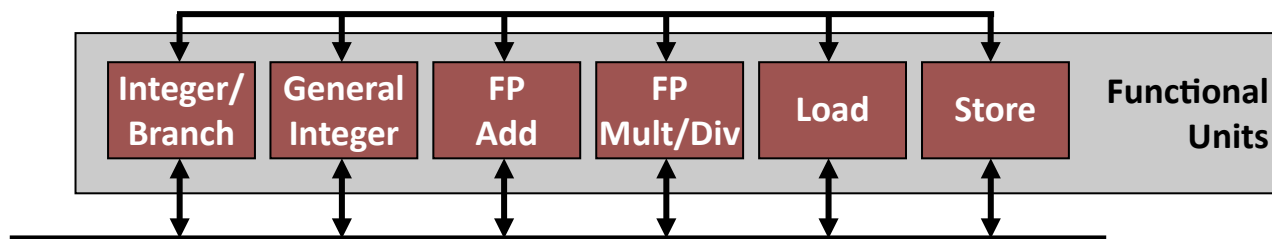


```
for (i = 0; i < n; i++) {  
    double val = 0;  
    for (j = 0; j < n; j++)  
        val += a[i*n + j];  
    b[i] = val;  
}
```



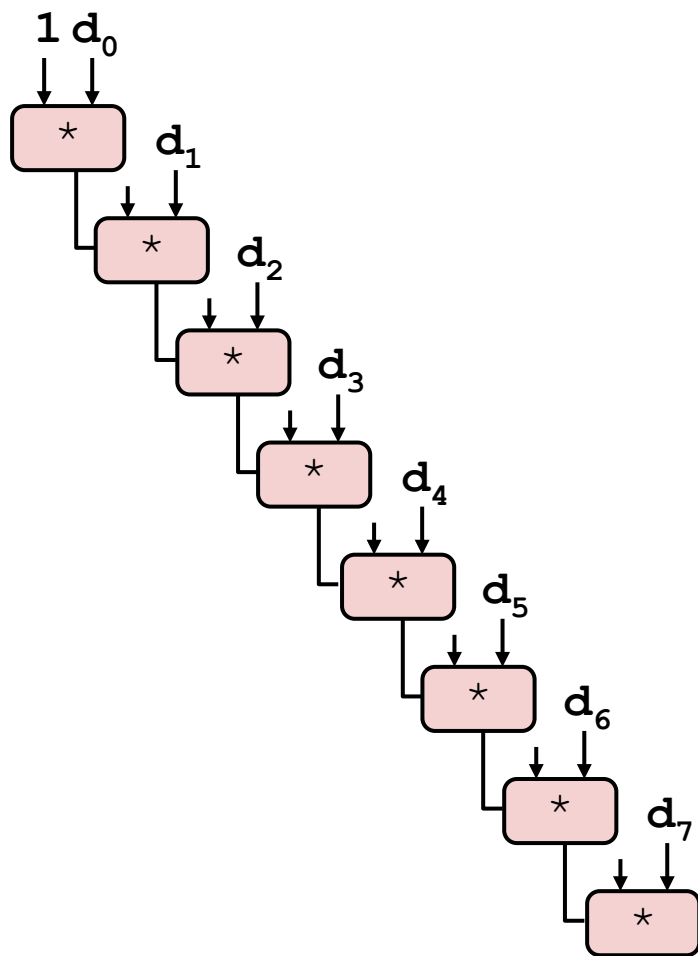
Last Time

- Instruction level parallelism
- Latency versus throughput

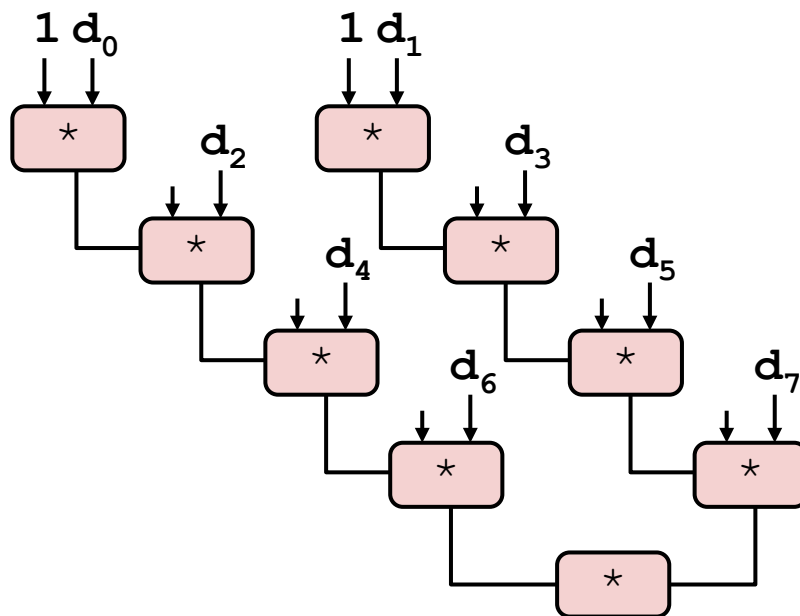


Last Time

■ Consequence

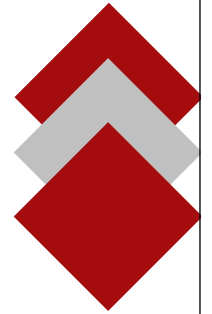


Twice as fast



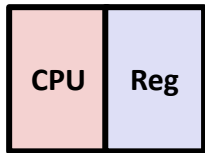
Today

- **Memory hierarchy, caches, locality**
- **Cache organization**
- **Program optimization:**
 - Cache optimizations

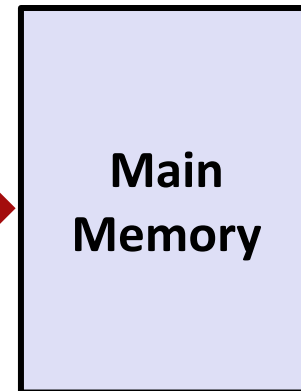


Problem: Processor-Memory Bottleneck

Processor performance
doubles about
every 18 months



Bus bandwidth
evolves much slower



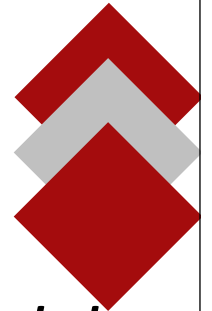
Core 2 Duo:
Can process at least
256 Bytes/cycle
(1 SSE two operand add and mult)

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

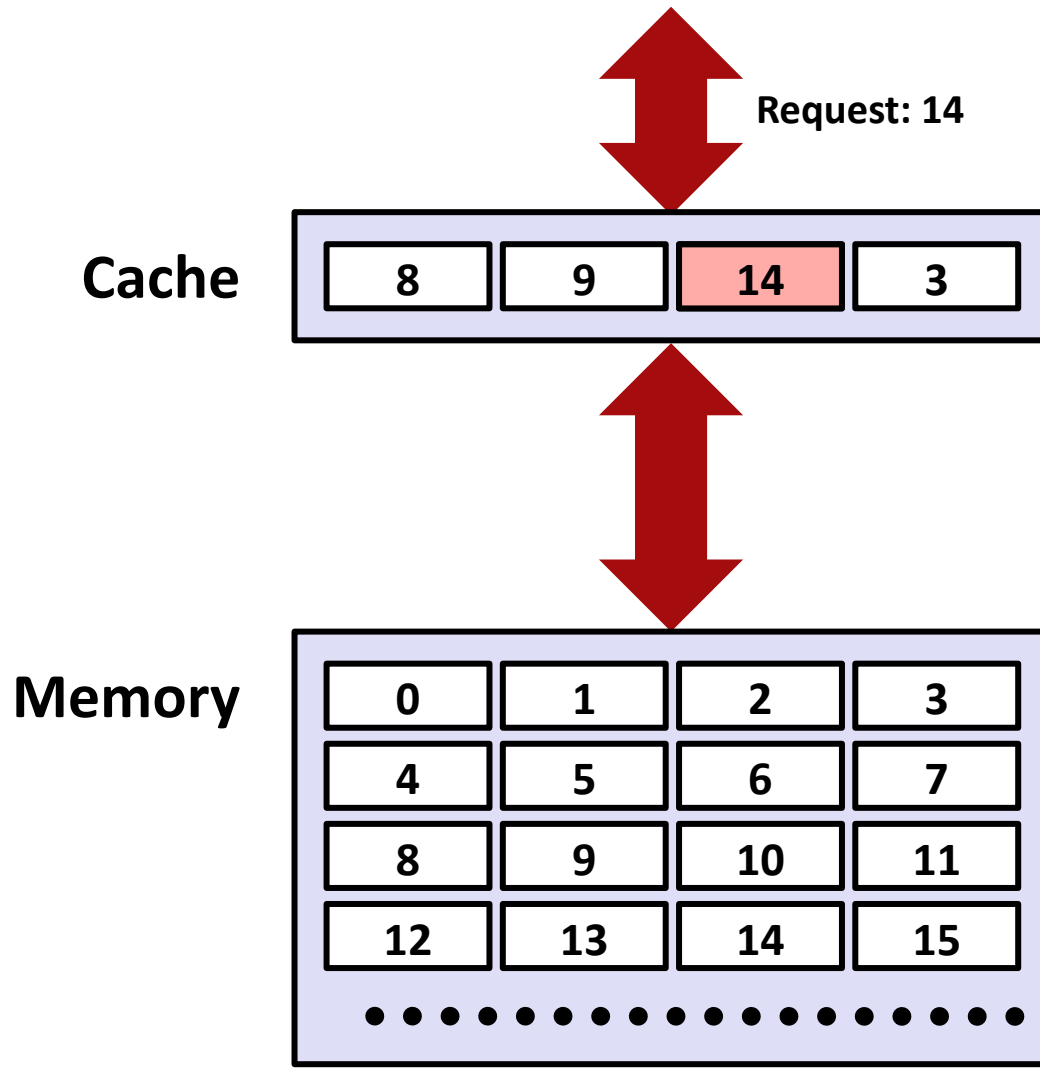
Solution: Caches

Cache

- **Definition:** Computer memory with short access time used for the storage of frequently or recently used instructions or data
 - From the French *caler* 'to hide'
 - Memory contents *hidden* close to their point of use, in hopes of using them again in the future
 - As opposed to buffers/variables/etc, which are not hidden and require programmer management

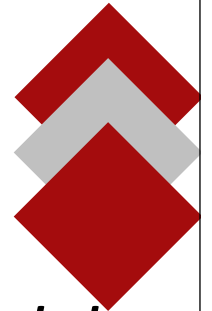


General Cache Concepts: Hit

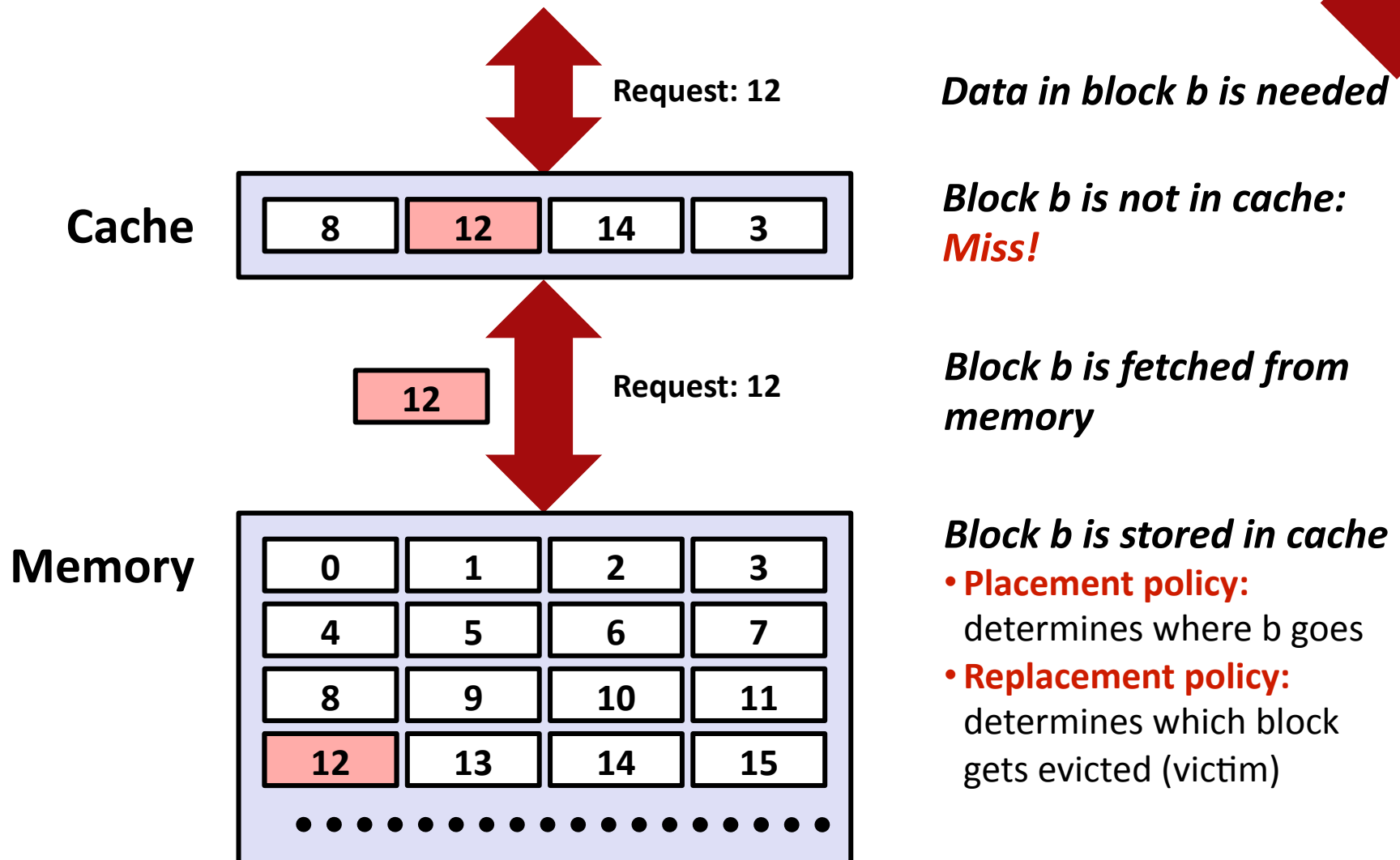


Data in block b is needed

Block b is in cache:
Hit!



General Cache Concepts: Miss



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

■ Miss Penalty

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

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

Types of Cache Misses

■ Cold (compulsory) miss

- Occurs on first access to a block

■ Capacity miss

- Occurs when the set of active cache blocks (working set) is larger than the cache size
- Must evict a victim to make space for replacement block

■ Conflict miss

- Most hardware caches limit blocks to a small subset (sometimes a singleton) of the available cache slots
 - e.g., block i must be placed in slot $(i \bmod 4)$
- Conflict misses occur when the cache is large enough, but multiple data objects all map to the same slot
 - e.g., referencing blocks 0, 8, 0, 8, ... would miss every time

Why Caches Work

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

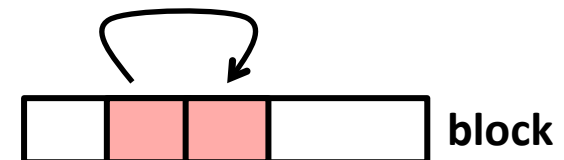
- **Temporal locality:**

- 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



- **Fundamental Law (almost): Programs with good locality run faster than programs with poor locality**

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**

Locality Example #1

```
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;
}
```

Locality Example #2

```
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;
}
```

Locality Example #3

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

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

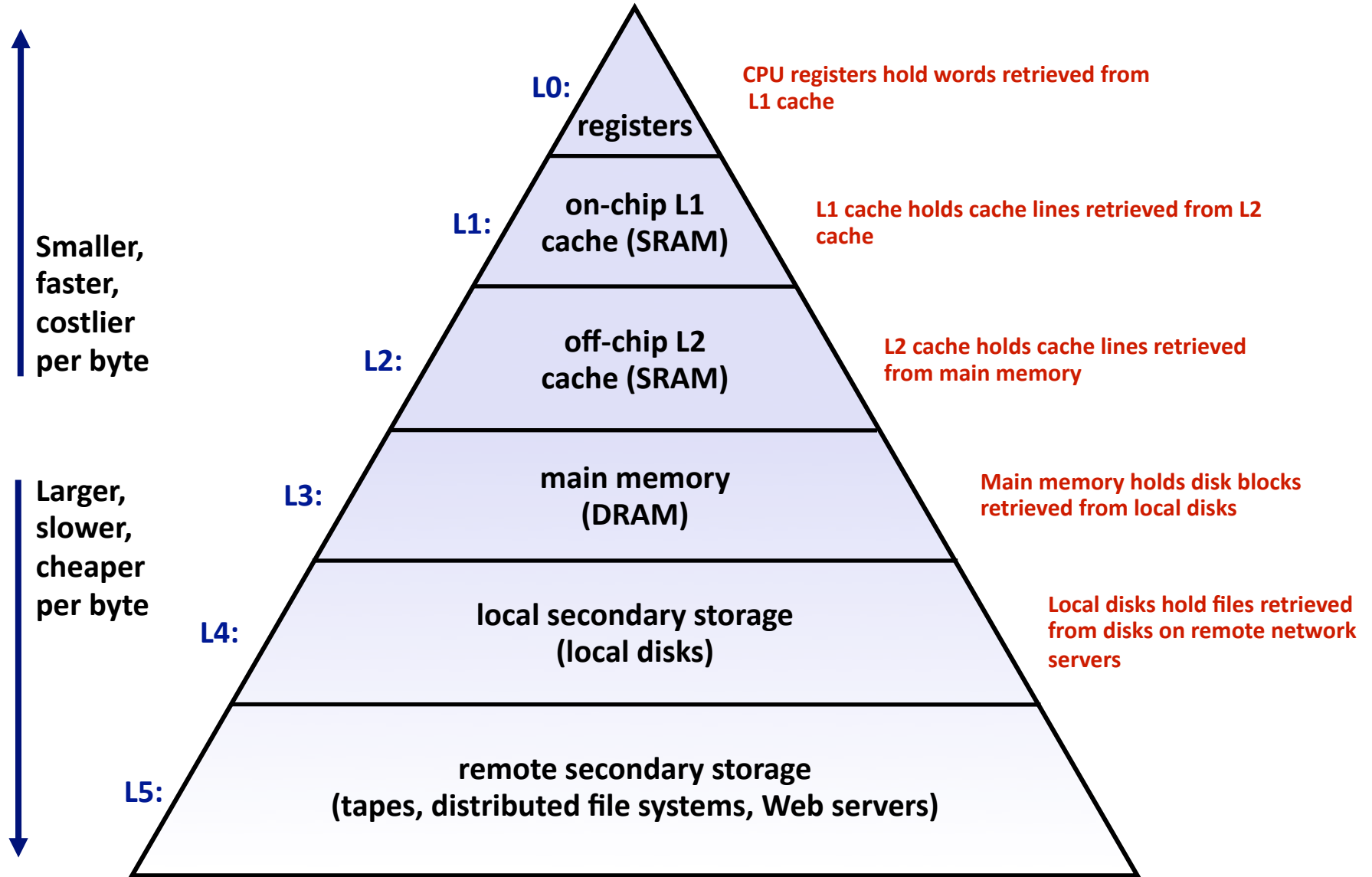
    return sum;
}
```

- How can it be fixed?

Memory Hierarchies

- **Fundamental and enduring properties of hardware:**
 - Faster storage technologies almost always cost more per byte and have lower capacity
 - The gaps between memory technology speeds are widening
 - True of registers \leftrightarrow DRAM, DRAM \leftrightarrow disk, etc.
- **Fundamental and enduring property of software:**
 - 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****

An Example Memory Hierarchy



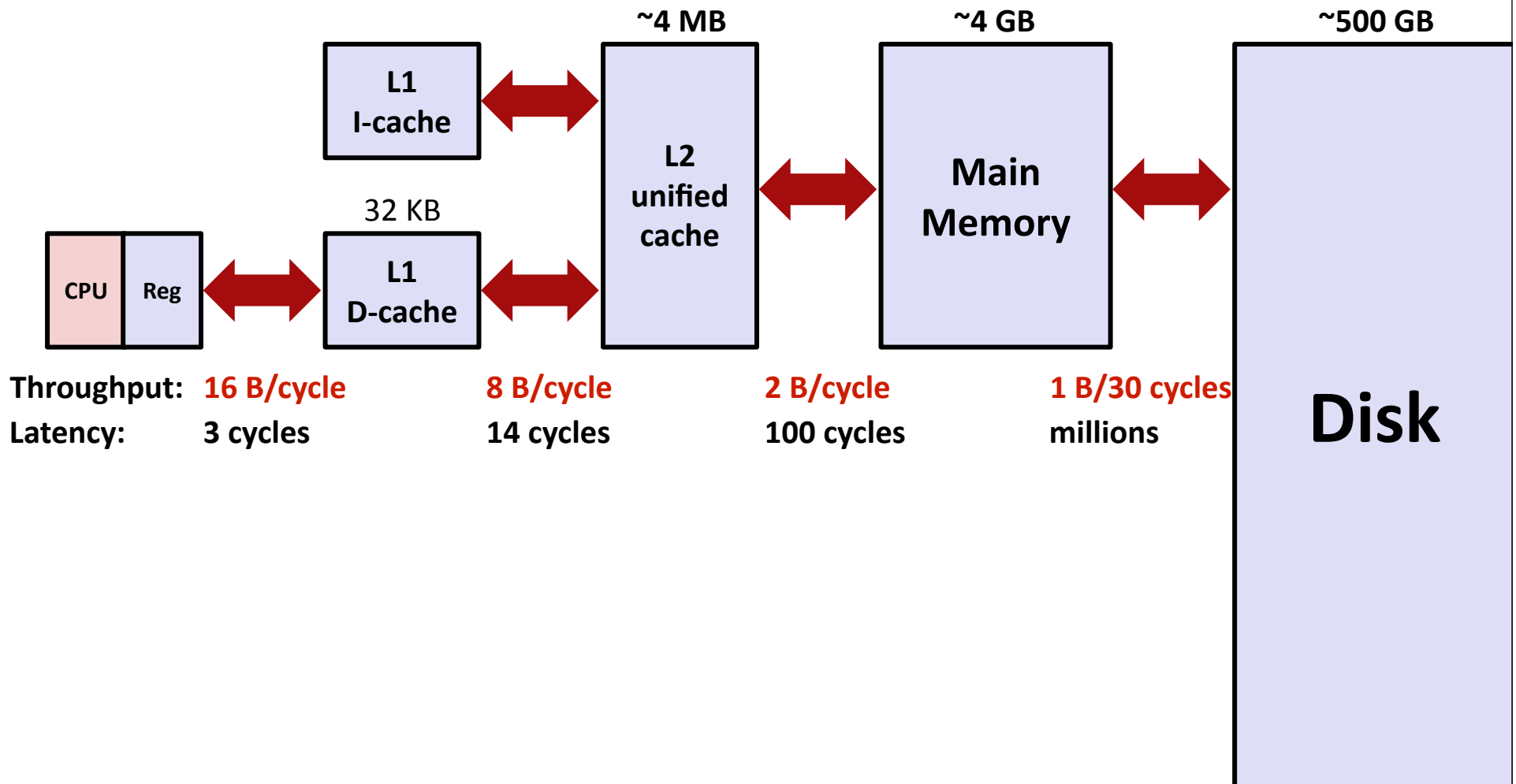
Examples of Caching in the Hierarchy

Cache Type	What is Cached?	Where is it Cached?	Latency (cycles)	Managed By
Registers	4-byte words	CPU core	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware
L1 cache	64-bytes block	On-Chip L1	1	Hardware
L2 cache	64-bytes block	Off-Chip L2	10	Hardware
Virtual Memory	4-KB page	Main memory	100	Hardware+OS
Buffer cache	Parts of files	Main memory	100	OS
Network buffer cache	Parts of files	Local disk	10,000,000	AFS/NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

Memory Hierarchy: Core 2 Duo

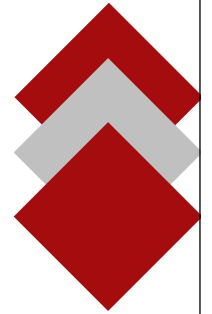
Not drawn to scale

L1/L2 cache: 64 B blocks

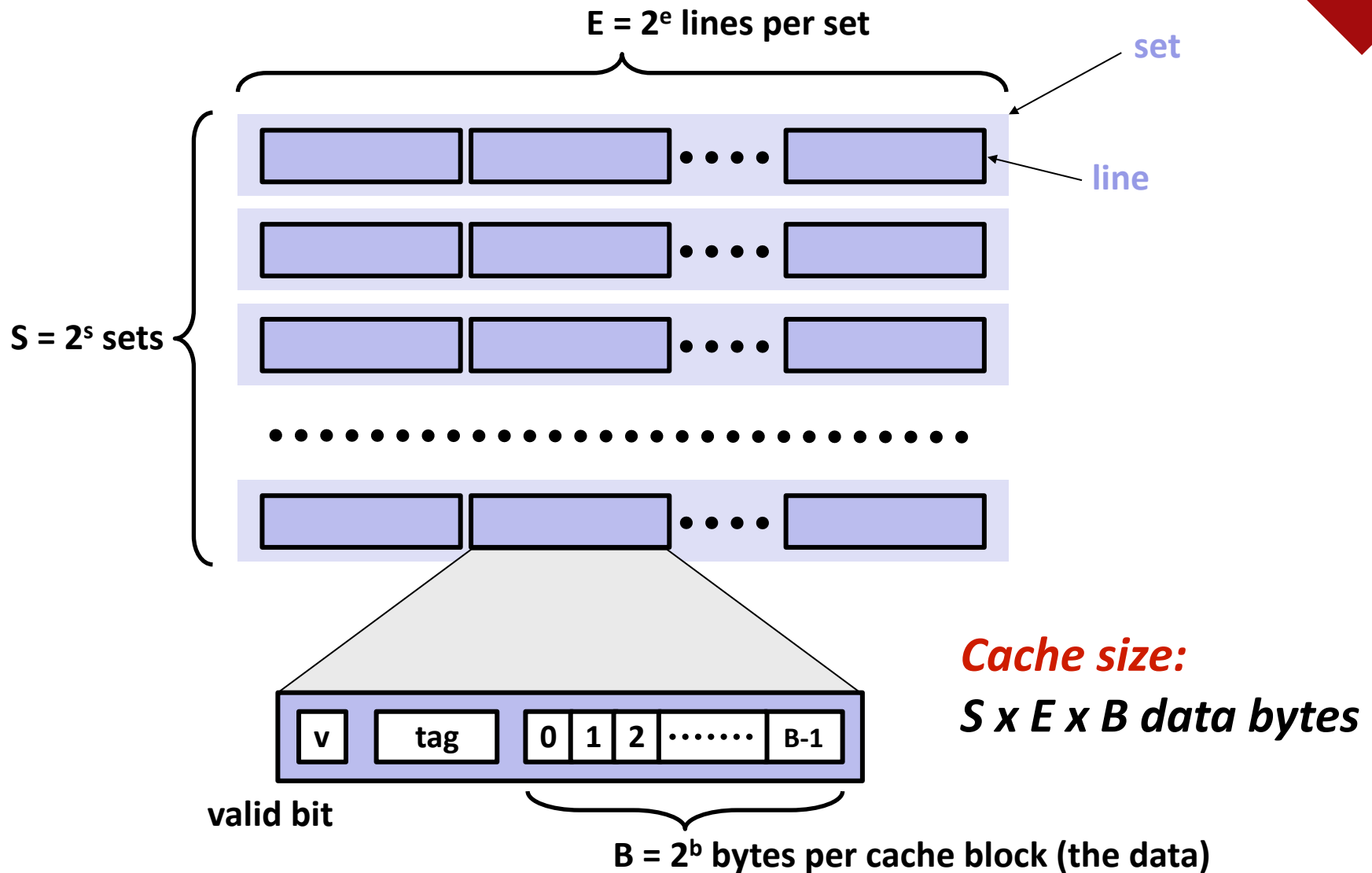


Today

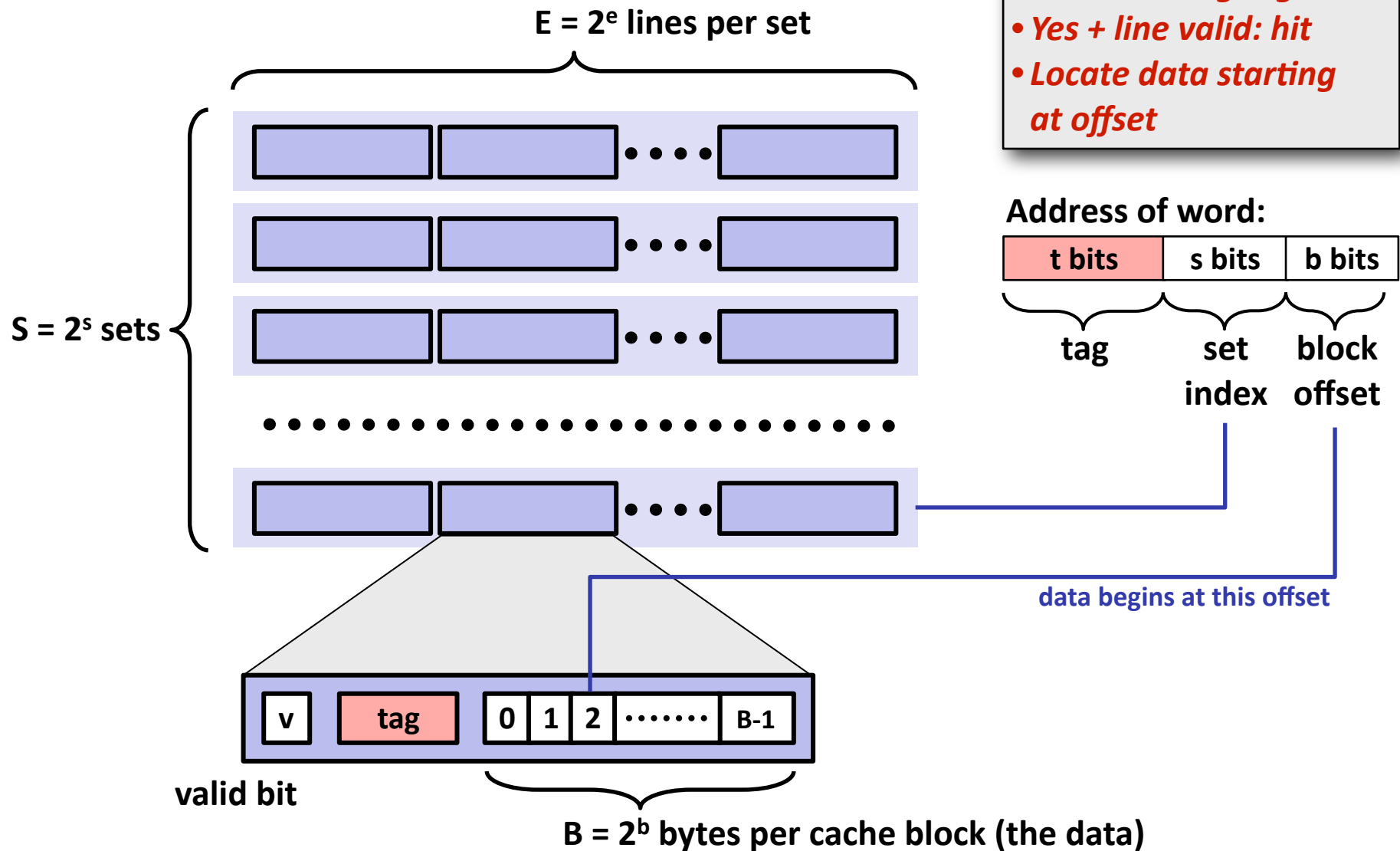
- Memory hierarchy, caches, locality
- **Cache organization**
- Program optimization:
 - Cache optimizations



General Cache Organization (S, E, B)



Cache Read

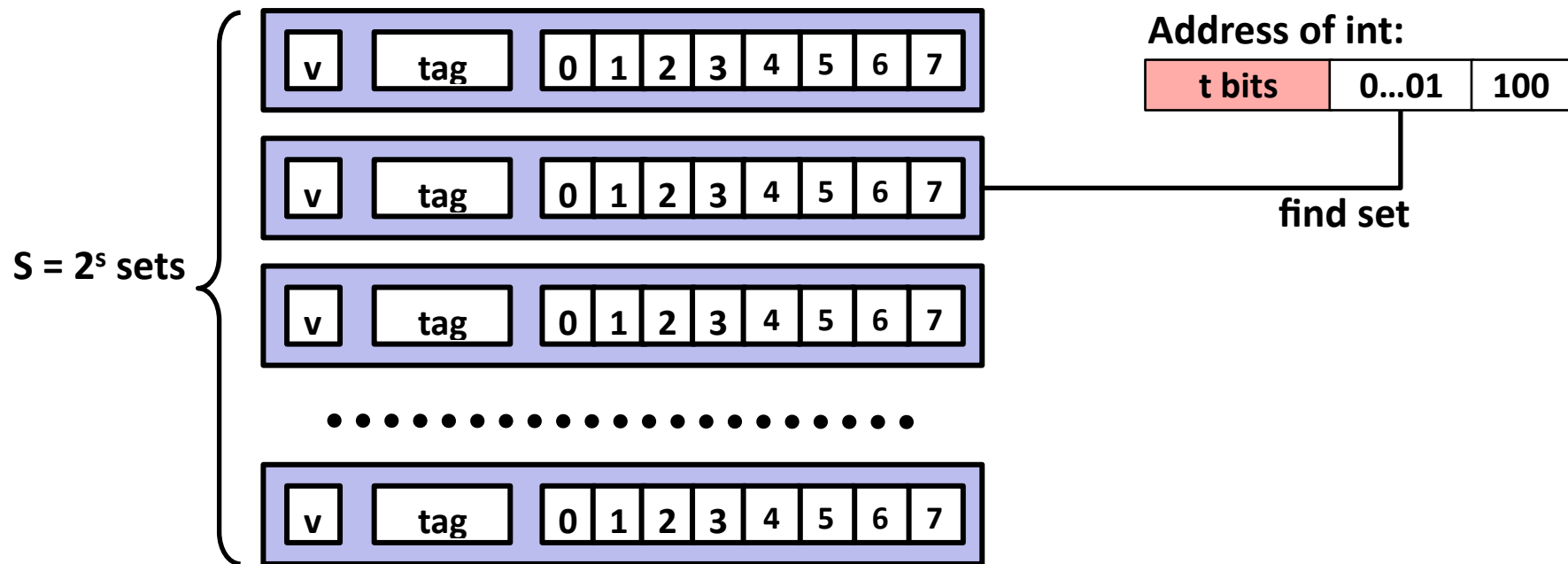


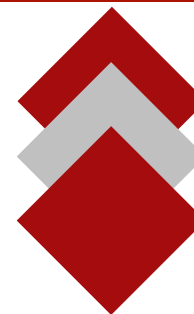
- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*
- *Locate data starting at offset*

Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set

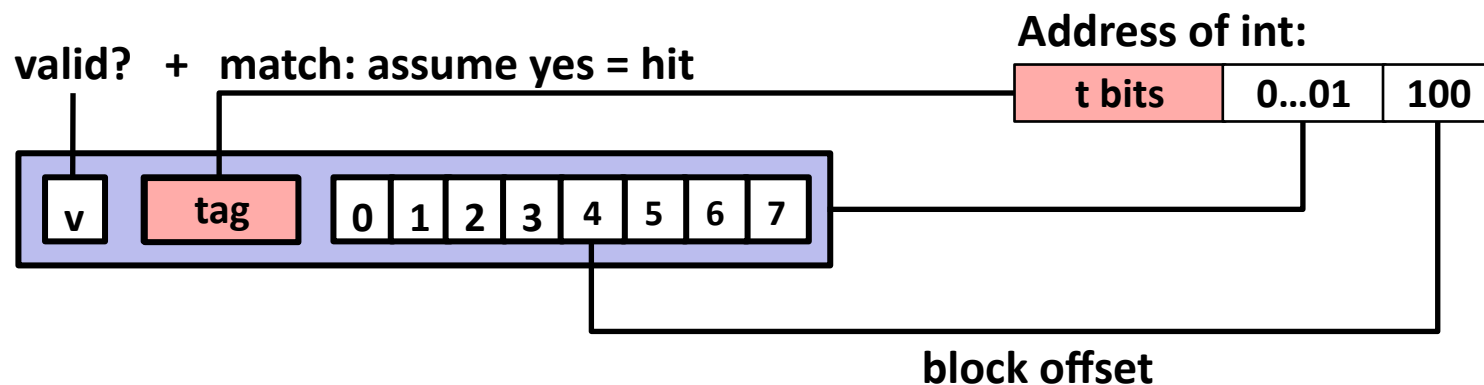
Assume: cache block size 8 bytes





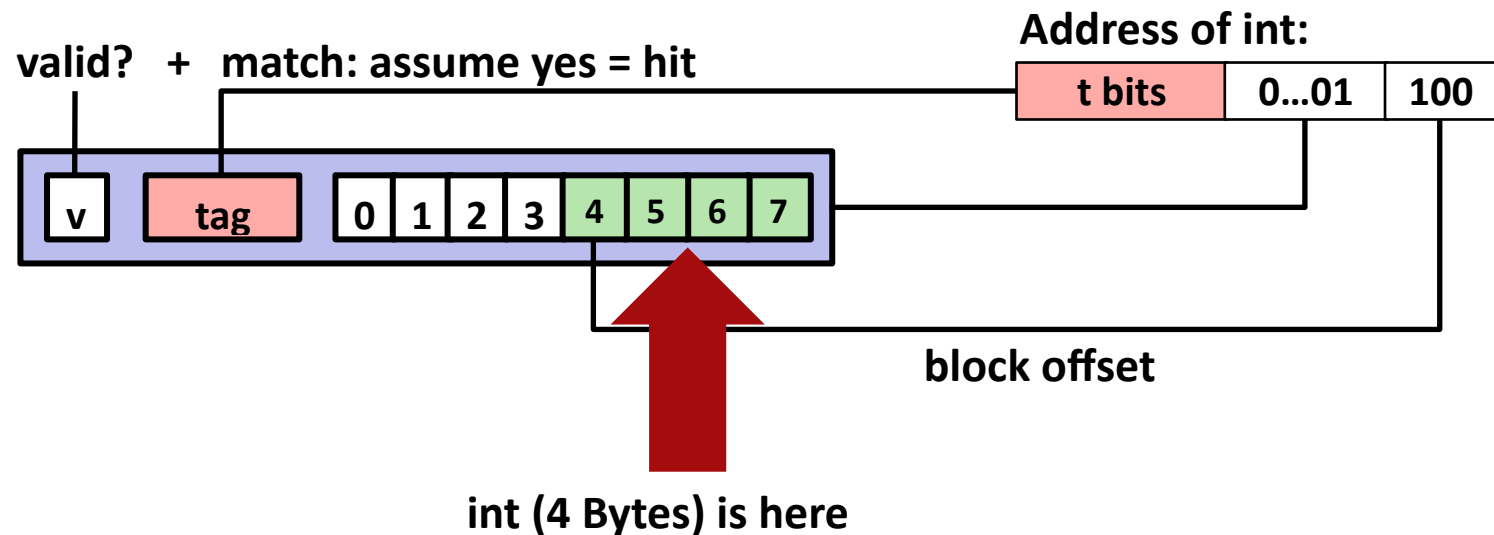
Example: Direct Mapped Cache (E = 1)

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

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

Ignore the variables sum, i, j

assume: cold (empty) cache,
a[0][0] goes here

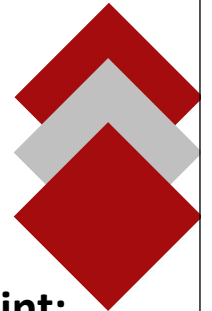


32 B = 4 doubles

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

$E = 2$: Two lines per set

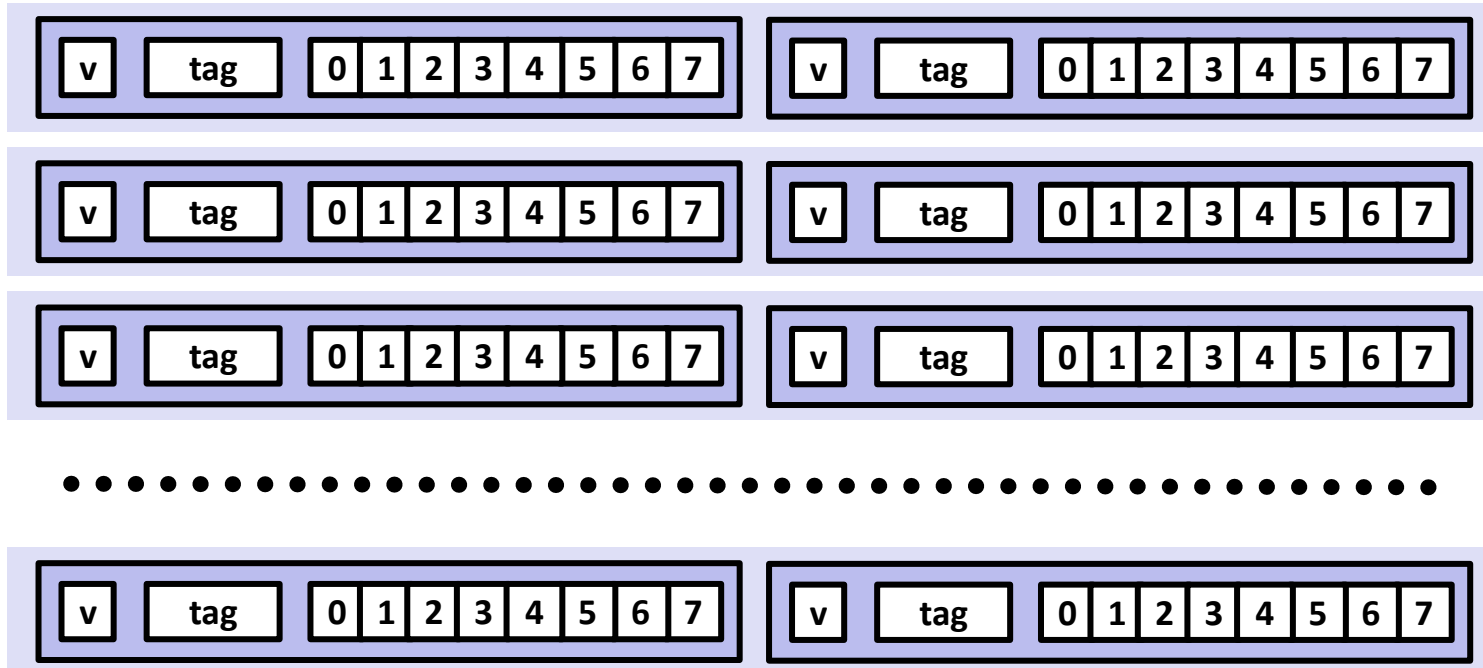
Assume: cache block size 8 bytes



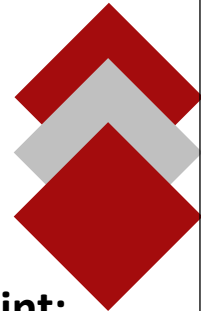
Address of short int:

t bits	0...01	100
--------	--------	-----

find set

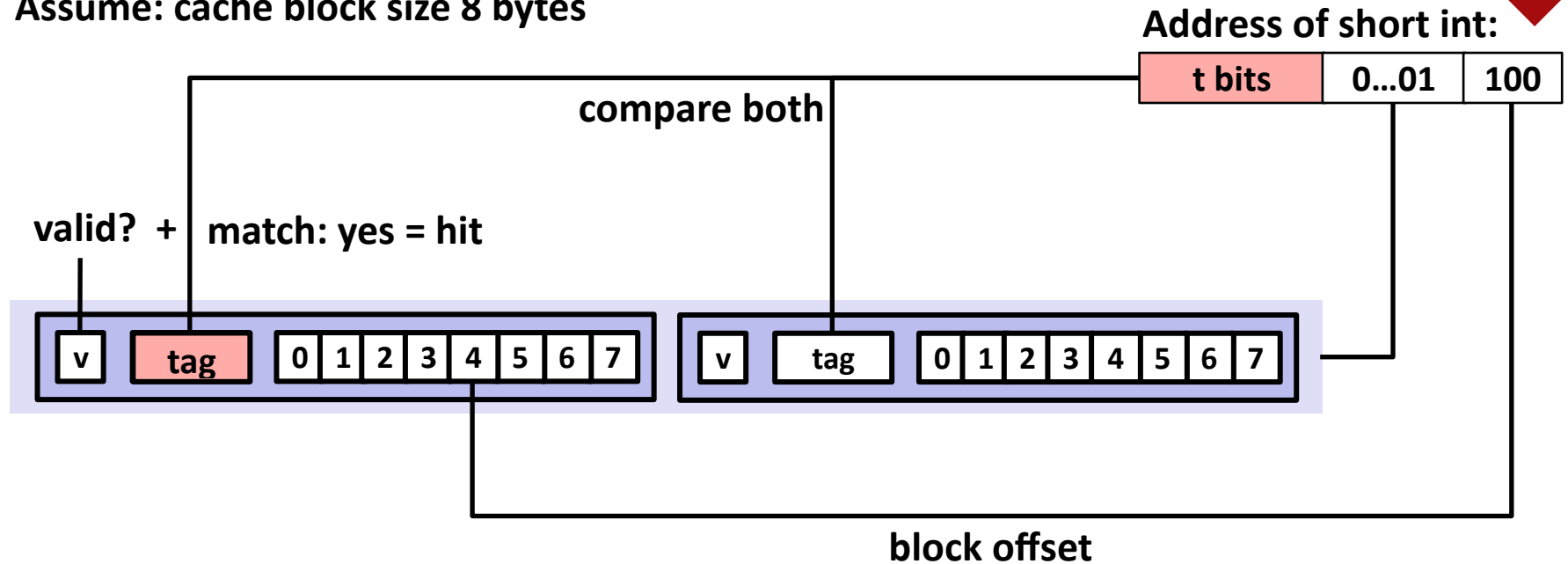


E-way Set Associative Cache (Here: E = 2)



E = 2: Two lines per set

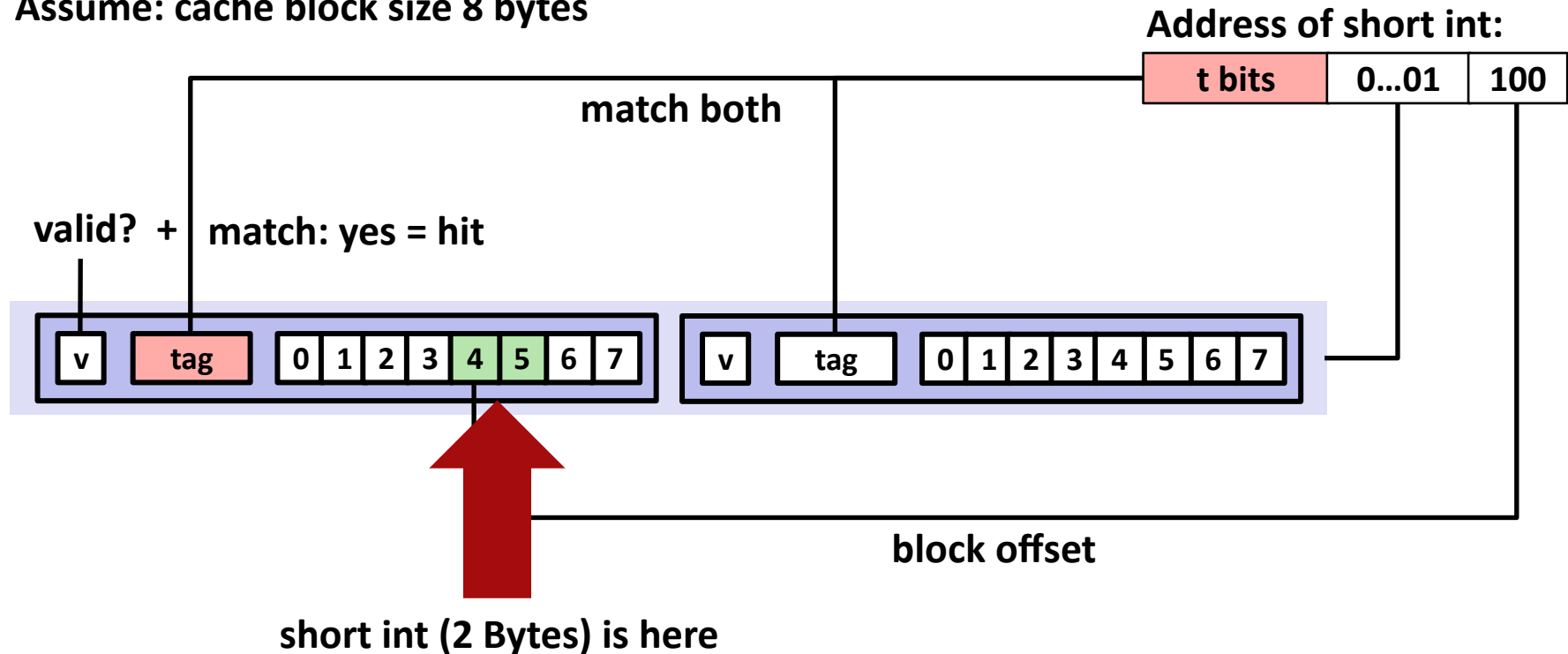
Assume: cache block size 8 bytes



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), ...

Example

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

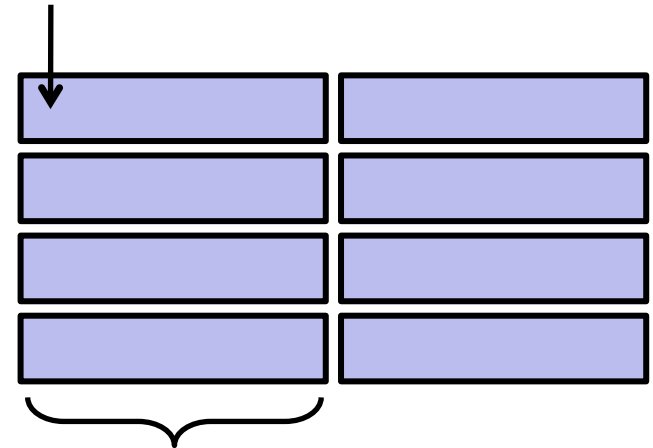
    for (i = 0; i < 16; i++)
        for (j = 0; j < 16; j++)
            sum += a[i][j];
    return sum;
}
```

```
int sum_array_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;
}
```

Ignore the variables sum, i, j

assume: cold (empty) cache,
a[0][0] goes here



32 B = 4 doubles

What about writes?

■ Multiple copies of data exist:

- L1, L2, Main Memory, Disk

■ What to do when a write-hit occurs?

- 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 when a write-miss occurs?

- 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**

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

Today

- Memory hierarchy, caches, locality
- Cache organization
- **Program optimization:**
 - Cache optimizations

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?

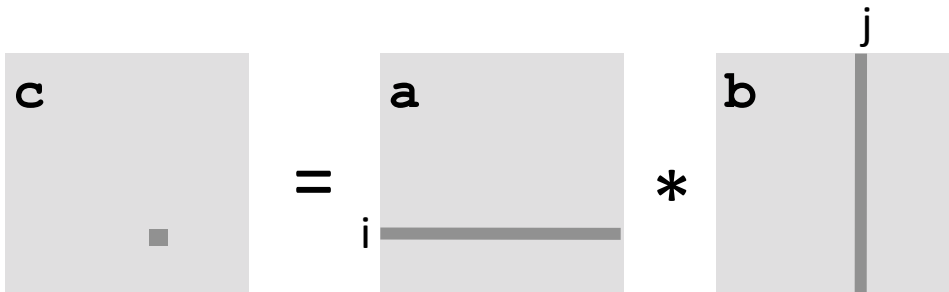
- Proper choice of algorithm
- Loop transformations

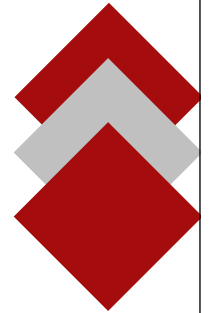
■ Cache versus register level optimization:

- In both cases locality desirable
- Register space much smaller + requires scalar replacement to exploit temporal locality
- Register level optimizations include exhibiting instruction level parallelism (conflicts with 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];  
}
```





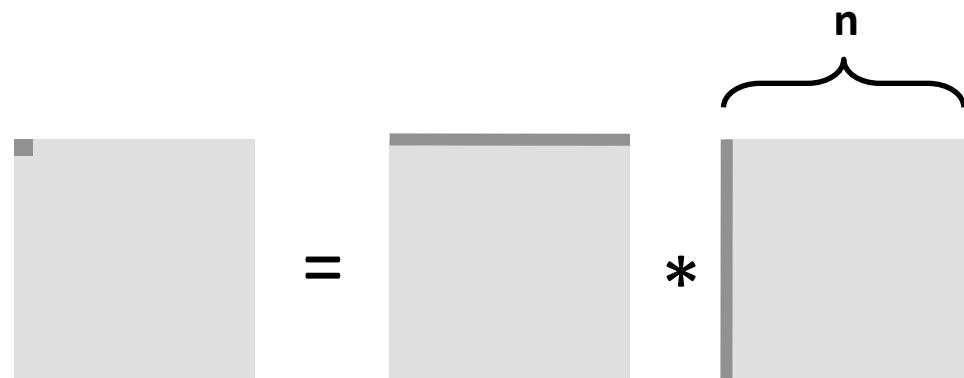
Cache Miss Analysis

■ Assume:

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

■ First iteration:

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



- Afterwards in cache:
(schematic)



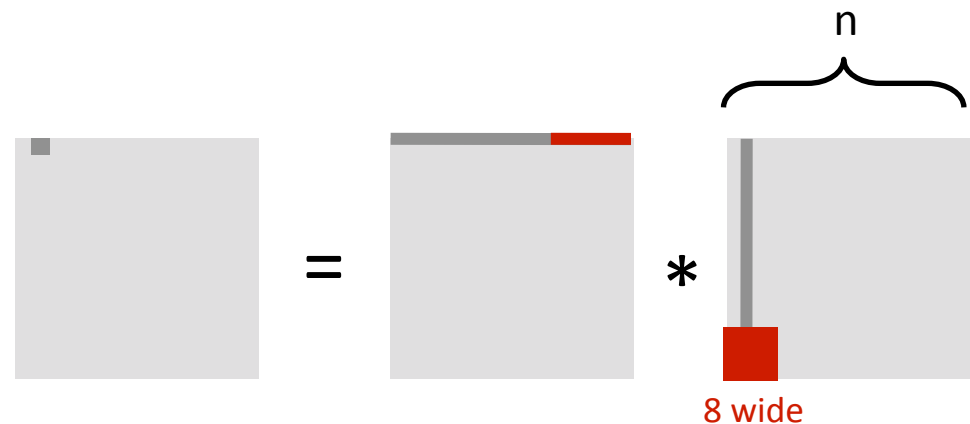
Cache Miss Analysis

■ 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:

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

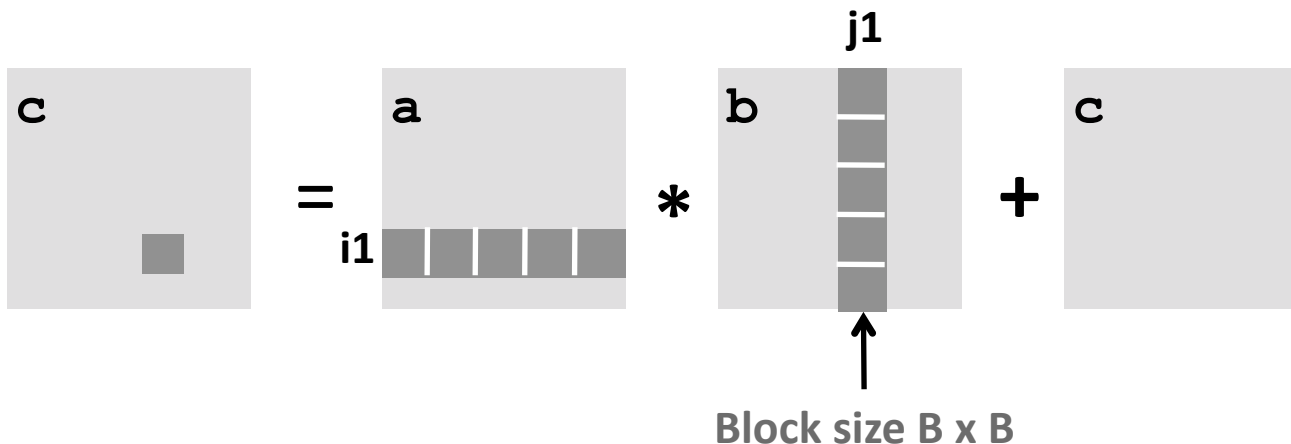
Blocked Matrix Multiplication

```

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];
}

```



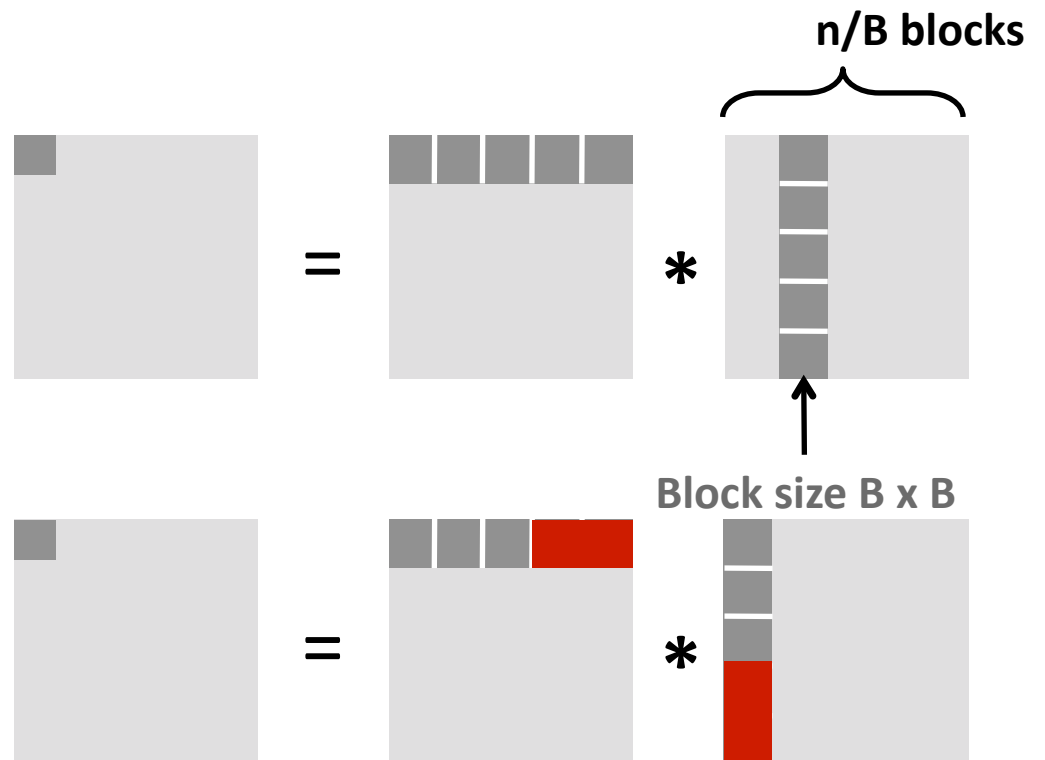
Cache Miss Analysis

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)



- Afterwards in cache
(schematic)

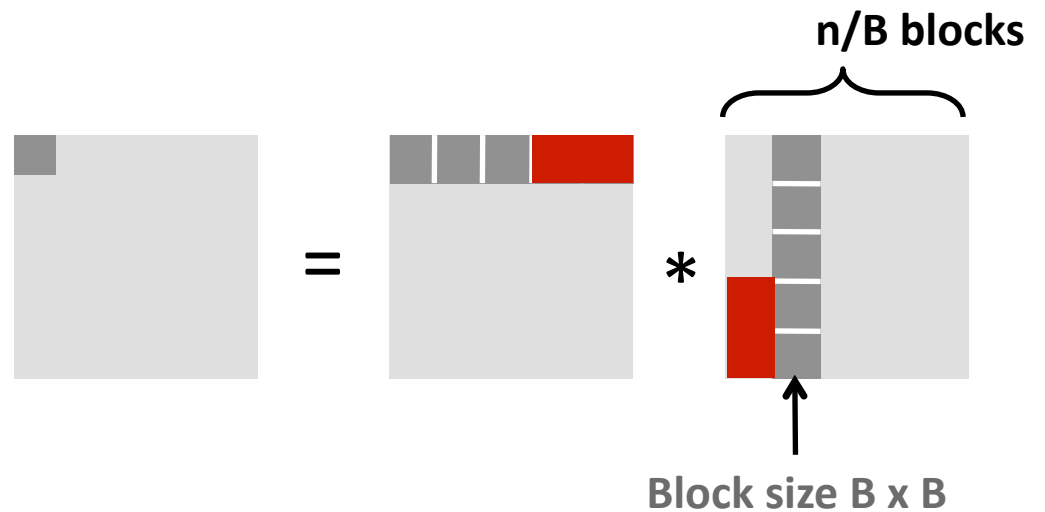
Cache Miss Analysis

■ Assume:

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

■ Second (block) iteration:

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



■ Total misses:

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

Matrix Multiplication: The Bottom Line

- 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

Summary

- **Memory hierarchy, caches, locality**
- **Cache organization**
- **Program optimization:**
 - Cache optimizations

- **Next Time: Exceptions**
 - Traps, Faults, Aborts, Interrupts
 - Processes
 - Context Switch
 - `fork`, `exit`