

# Memory Systems

## Lec09 – Cache Memories and Cache Coherence

Chin-Fu Nien (粘徹夫)

# Module 2: System & Software (con't)

# Cache Memories

The content of this part is mainly from:

Randal E. Bryant and David R. O'Hallaron, "Computer Systems: A Programmer's Perspective," 3/e.

(本節內容改自Prof. Randal E. Bryant and David R. O'Hallaron 8<sup>th</sup> Lectures課程講義)

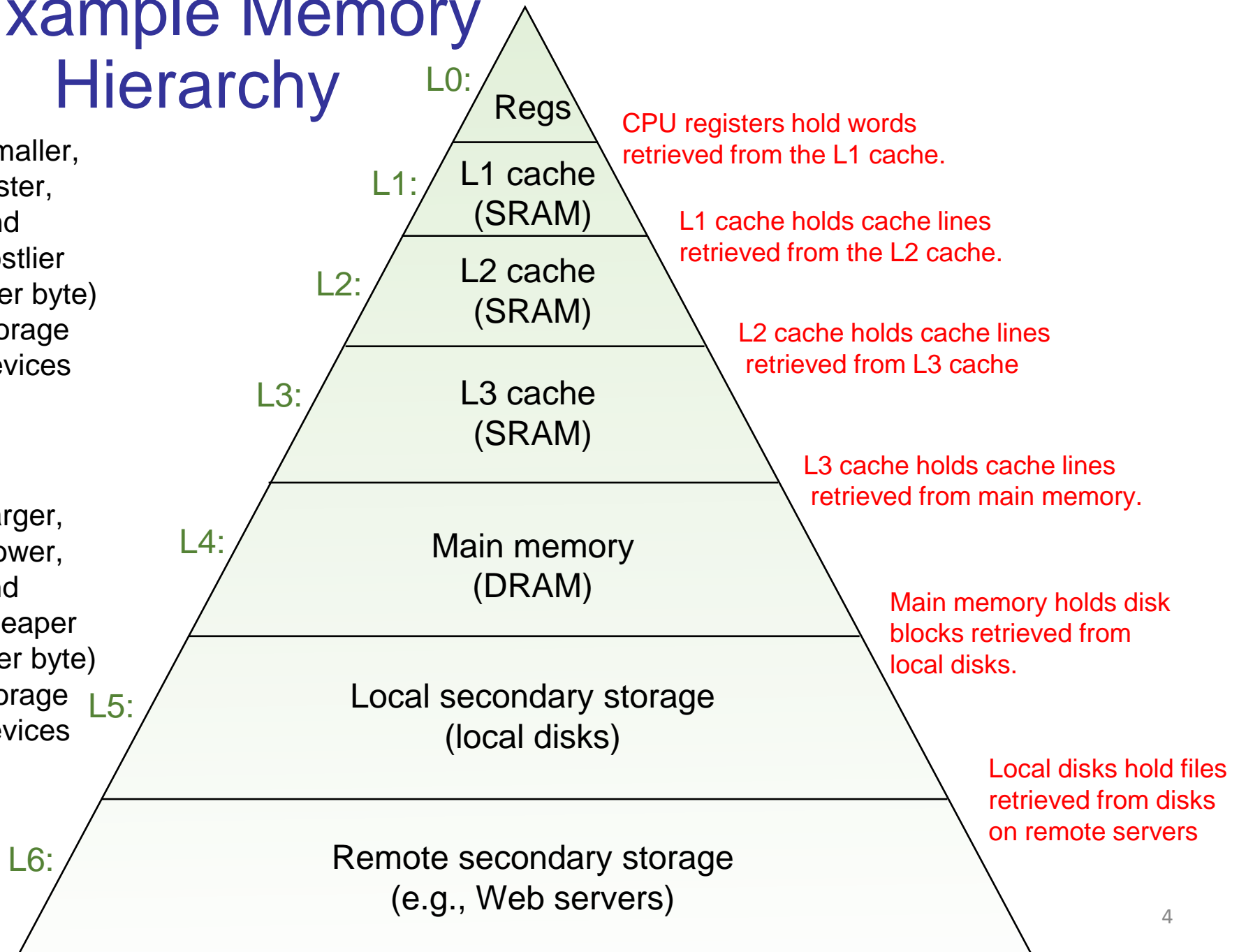
# Today

- Cache memory organization and operation
- Performance impact of caches
  - The memory mountain
  - Rearranging loops to improve spatial locality
  - Using blocking to improve temporal locality

# Example Memory Hierarchy

↑  
Smaller,  
faster,  
and  
costlier  
(per byte)  
storage  
devices

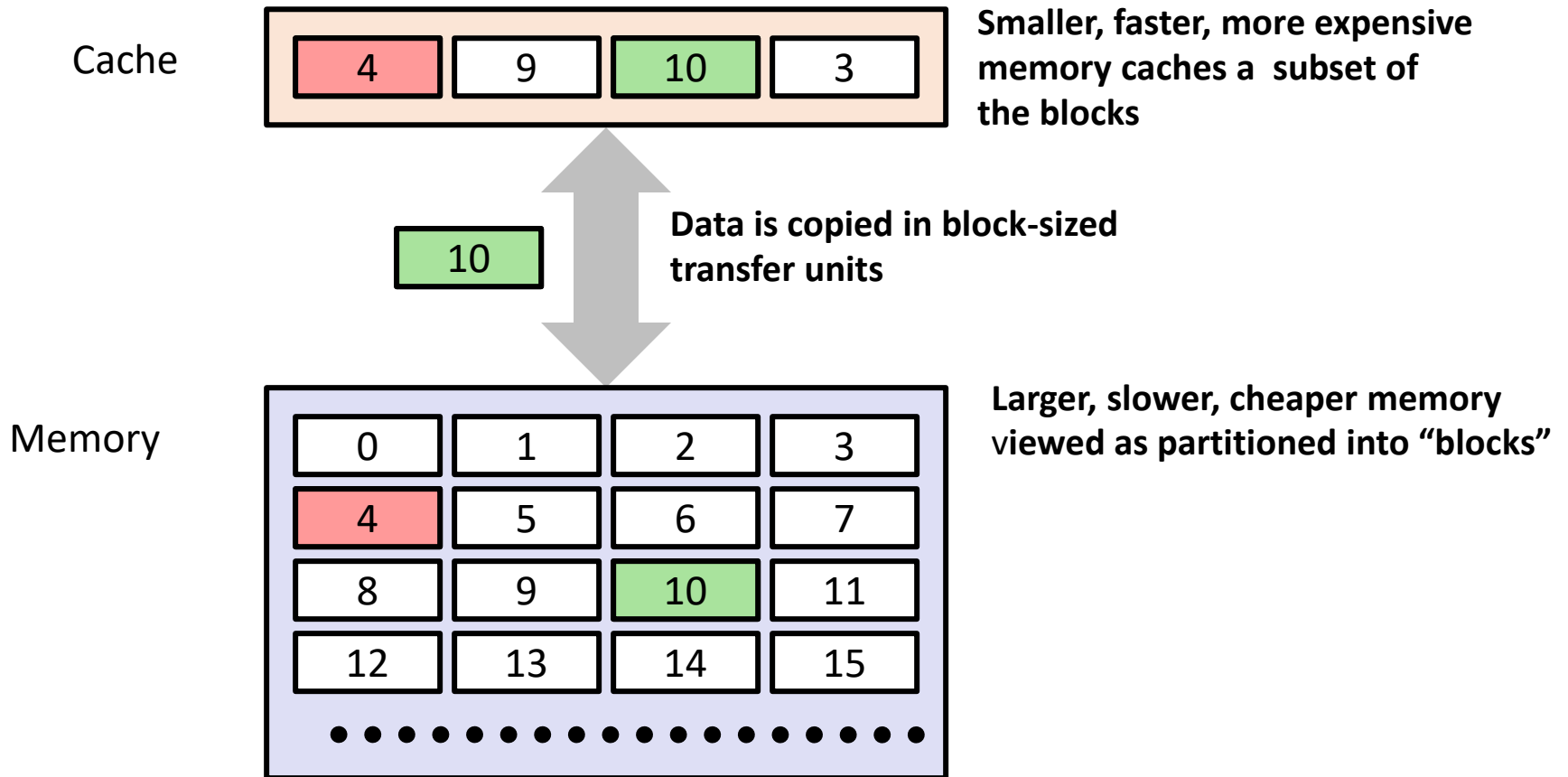
↓  
Larger,  
slower,  
and  
cheaper  
(per byte)  
storage  
devices



# Cache Line / Cache Block

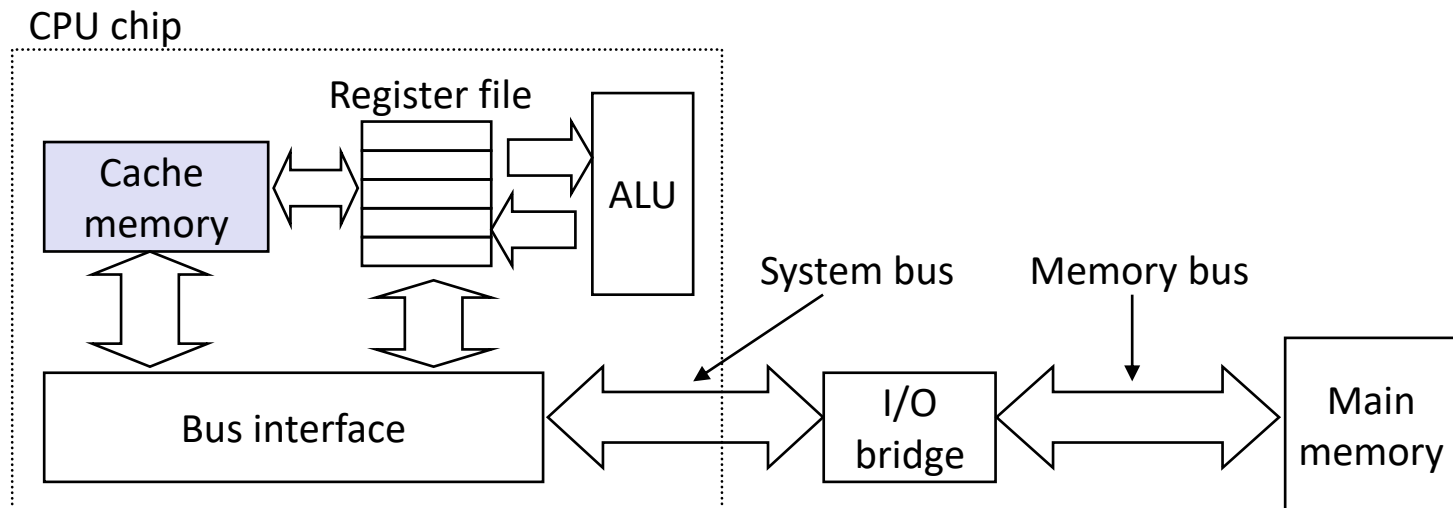
- A cache is typically read/write by a unit called “cache line” or “cache block”
- The size for a cache line (cache block) is CPU-dependent
  - Typical sizes are 16, 32, 64, and 128 Bytes.
  - Most common case: 64B.

# General Cache Concept



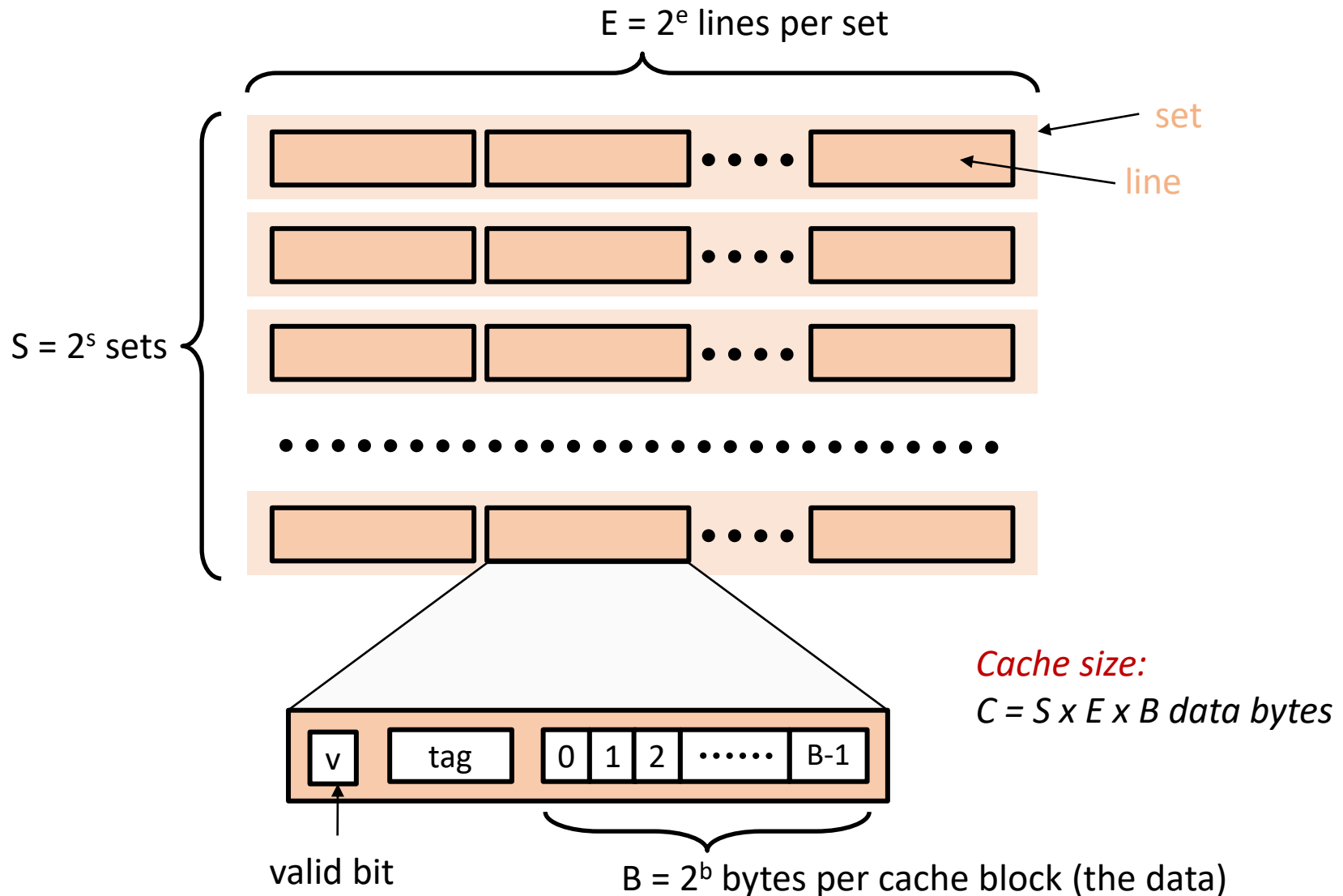
# Cache Memories

- **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 cache
- Typical system structure:

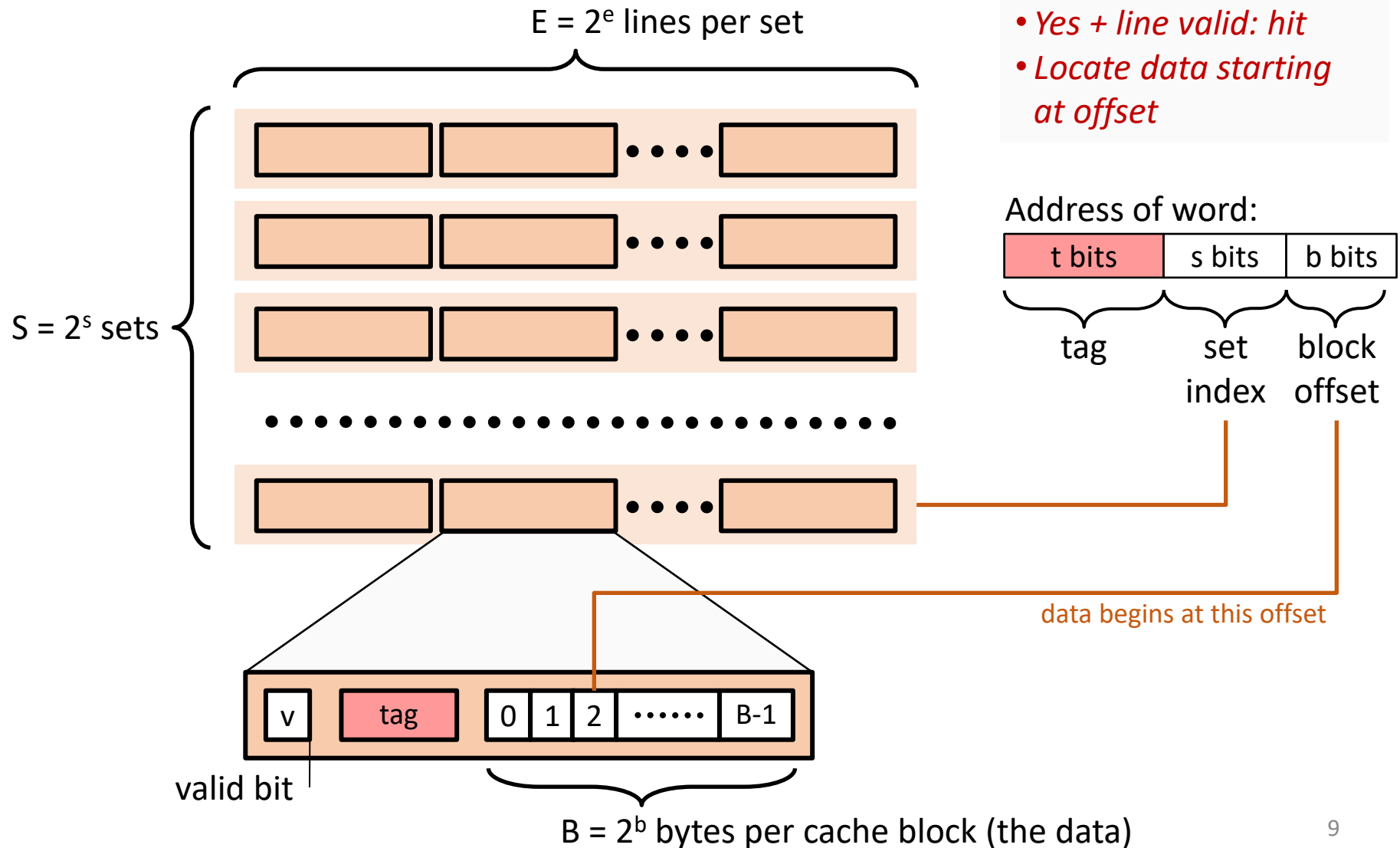




# General Cache Organization (S, E, B)

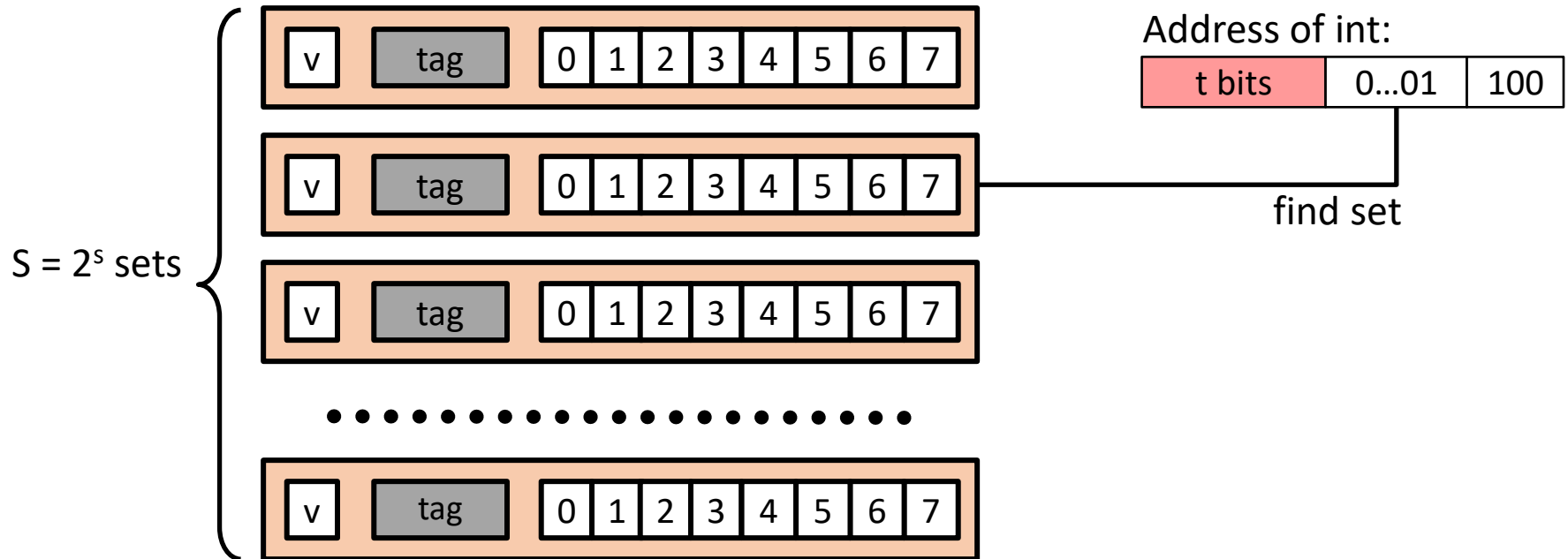


# Cache Read



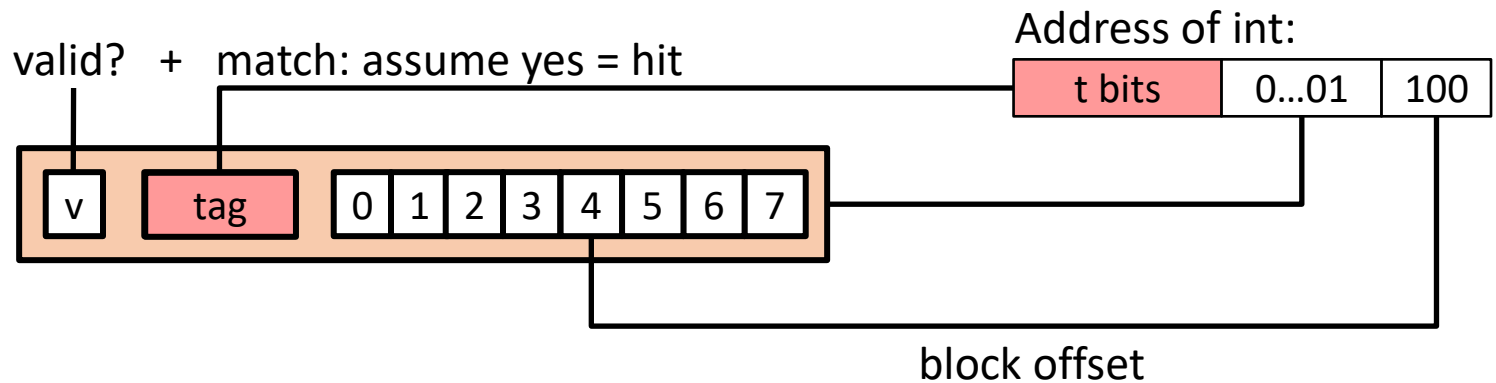
# Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set  
Assume: cache block size 8 bytes



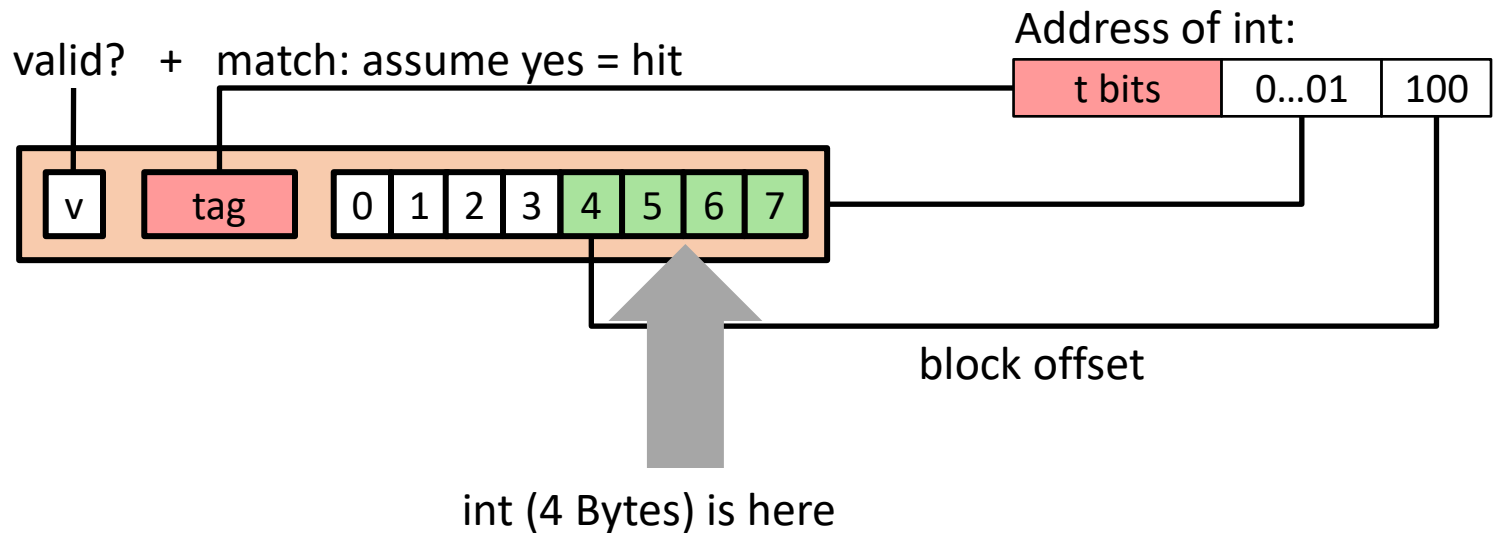
# Example: Direct Mapped Cache (E = 1)

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



If tag doesn't match: old line is evicted and replaced

# Direct-Mapped Cache Simulation

t=1	s=2	b=1
x	xx	x

M=16 bytes (4-bit addresses), B=2 bytes/block,  
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

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	Block
Set 0	1	0	M[0-1]
Set 1			
Set 2			
Set 3	1	0	M[6-7]

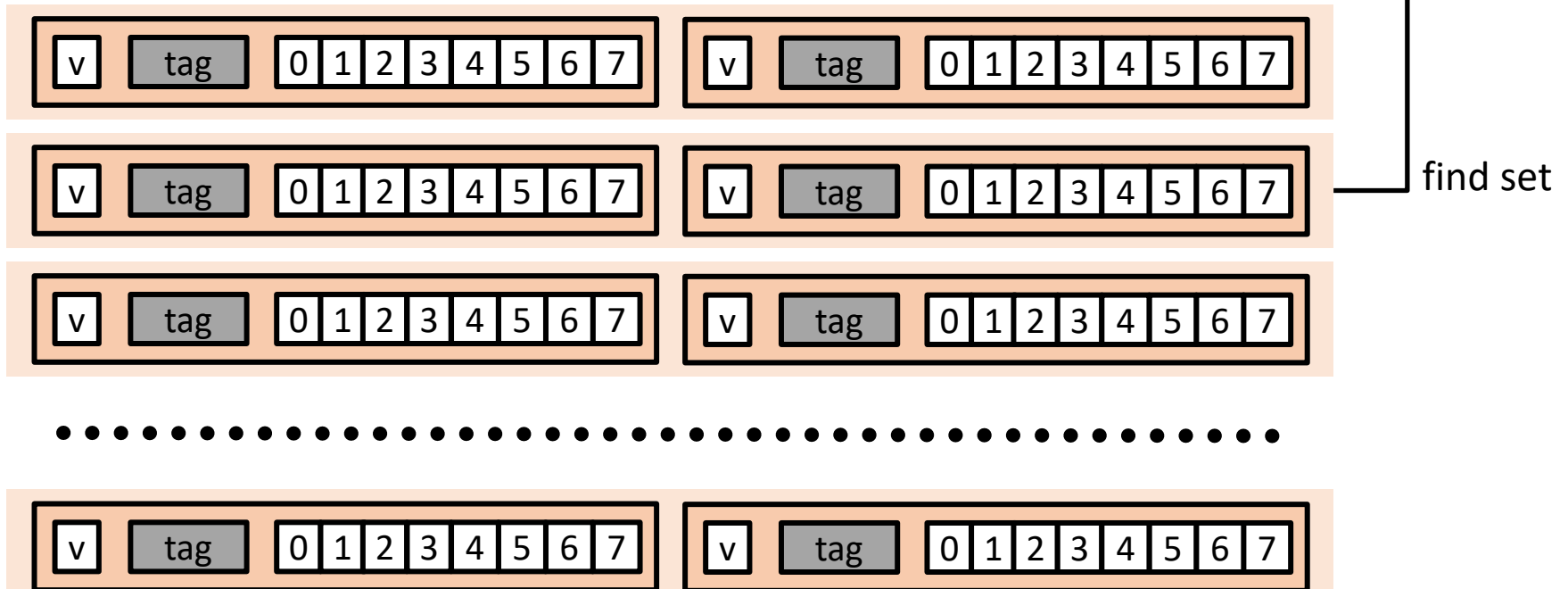
# E-way Set Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size 8 bytes

Address of short int:

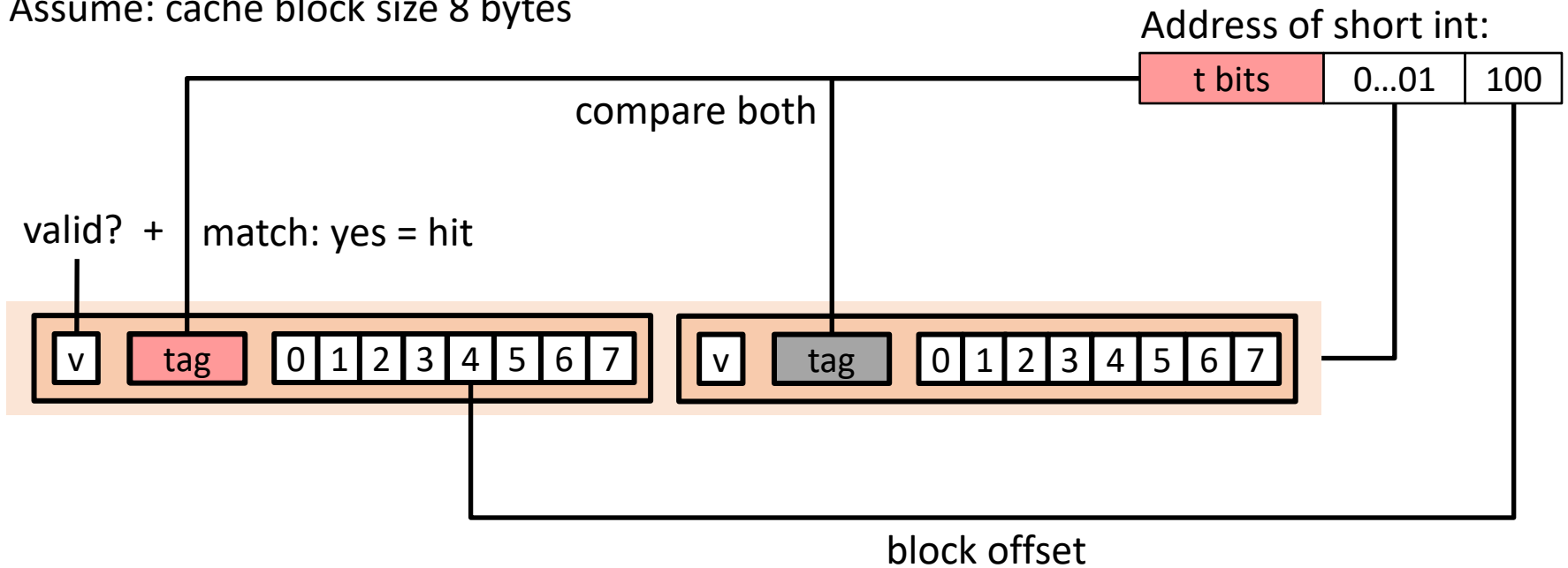
t bits	0...01	100
--------	--------	-----



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

E = 2: Two lines per set

Assume: cache block size 8 bytes

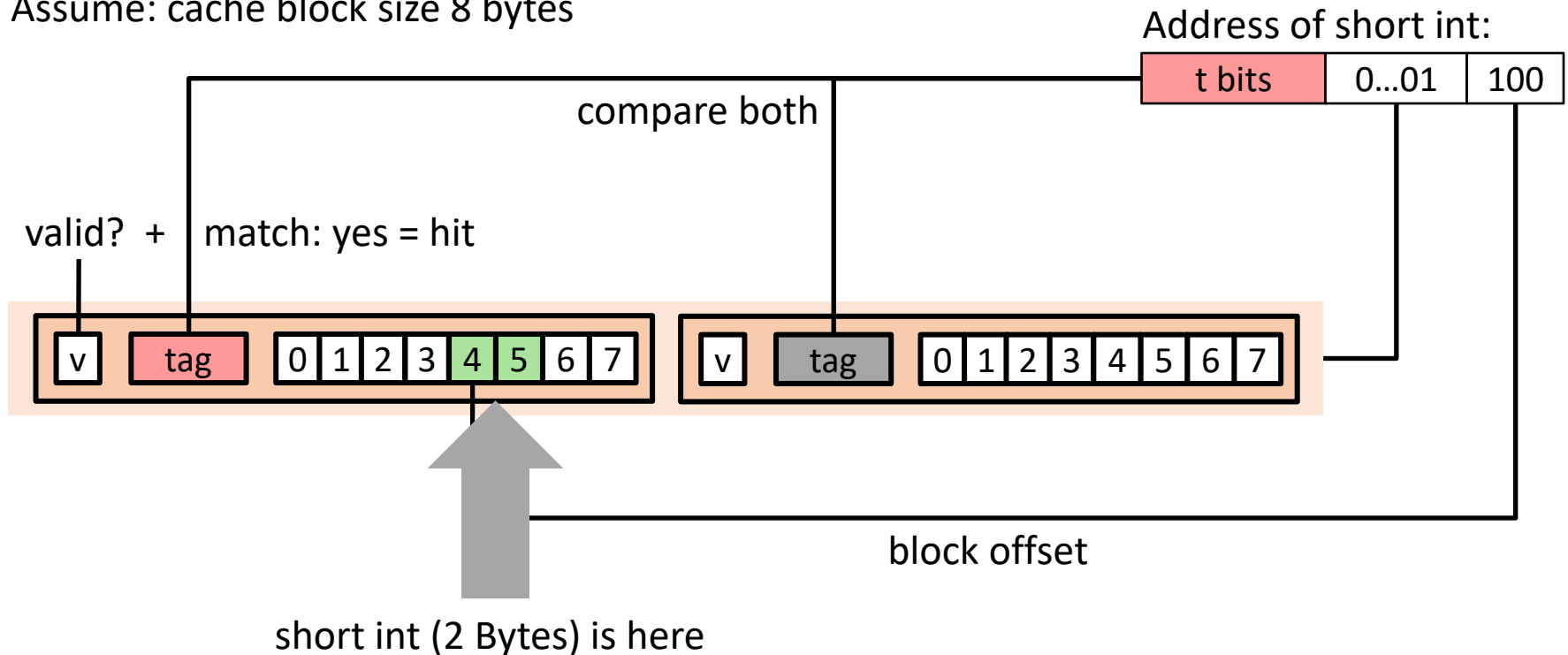




# E-way Set Associative Cache (Here: E = 2)

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 Set Associative Cache Simulation

t=2	s=1	b=1
xx	x	x

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

Address trace (reads, one byte per read):

0	[00 <u>00</u> <sub>2</sub> ],	miss
1	[00 <u>01</u> <sub>2</sub> ],	hit
7	[01 <u>11</u> <sub>2</sub> ],	miss
8	[10 <u>00</u> <sub>2</sub> ],	miss
0	[00 <u>00</u> <sub>2</sub> ]	hit

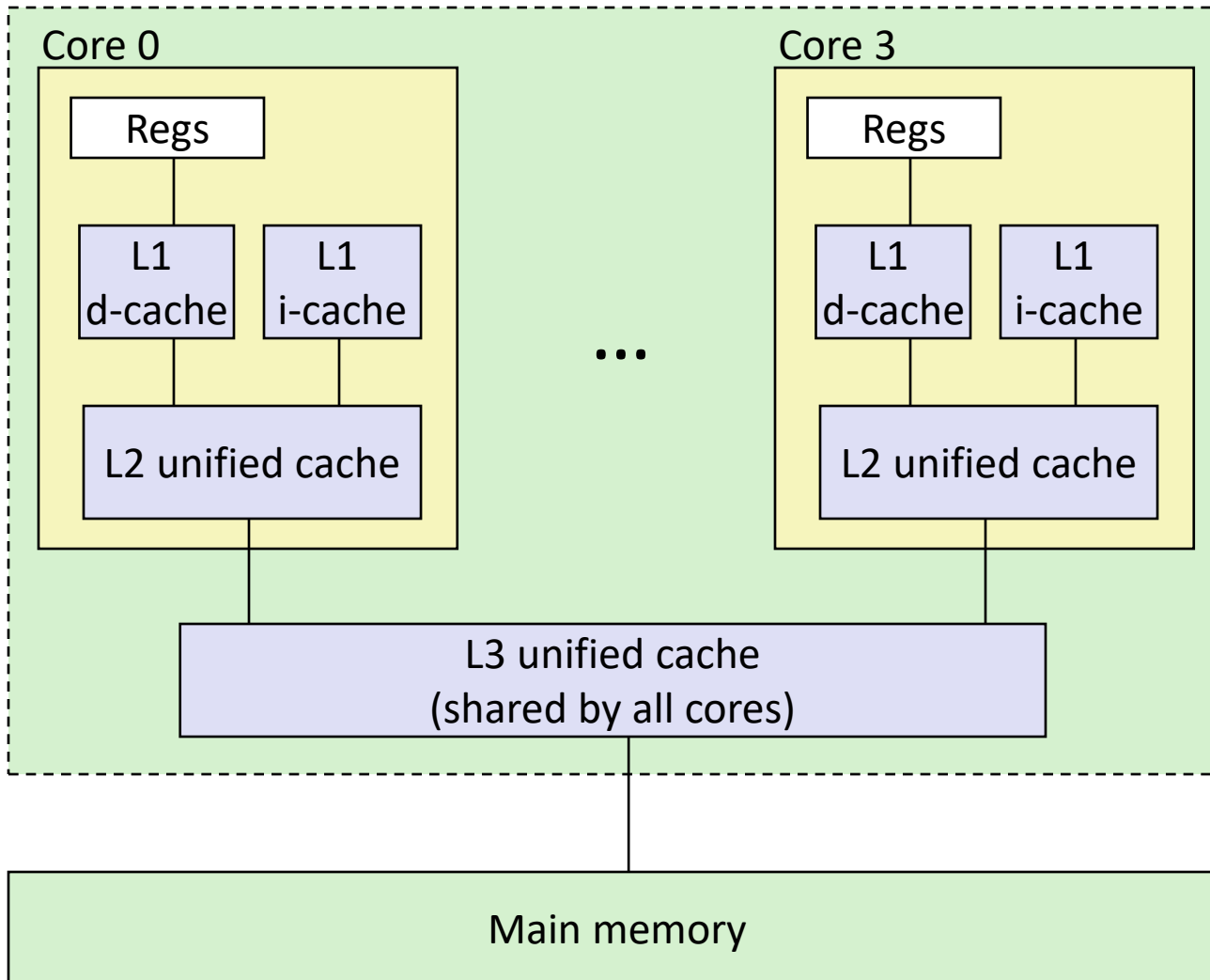
	v	Tag	Block
Set 0	1	00	M[0-1]
	1	10	M[8-9]
Set 1	1	01	M[6-7]
	0		

# 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**

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

# 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 (**temporal locality**)
  - Stride-1 reference patterns are good (**spatial locality**)

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

# Today

- Cache organization and operation
- **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

```
long data[MAXELEMS]; /* Global array to traverse */
```

```
/* test - Iterate over first "elems" elements of  
 *   array "data" with stride of "stride", using  
 *   using 4x4 loop unrolling.  
 */
```

```
int test(int elems, int stride) {  
    long i, sx2=stride*2, sx3=stride*3, sx4=stride*4;  
    long acc0 = 0, acc1 = 0, acc2 = 0, acc3 = 0;  
    long length = elems, limit = length - sx4;
```

```
/* Combine 4 elements at a time */
```

```
for (i = 0; i < limit; i += sx4) {  
    acc0 = acc0 + data[i];  
    acc1 = acc1 + data[i+stride];  
    acc2 = acc2 + data[i+sx2];  
    acc3 = acc3 + data[i+sx3];  
}
```

```
/* Finish any remaining elements */
```

```
for (; i < length; i++) {  
    acc0 = acc0 + data[i];  
}  
return ((acc0 + acc1) + (acc2 + acc3));
```

```
}
```

*mountain/mountain.c*

Call `test()` with many combinations of `elems` and `stride`.

For each `elems` and `stride`:

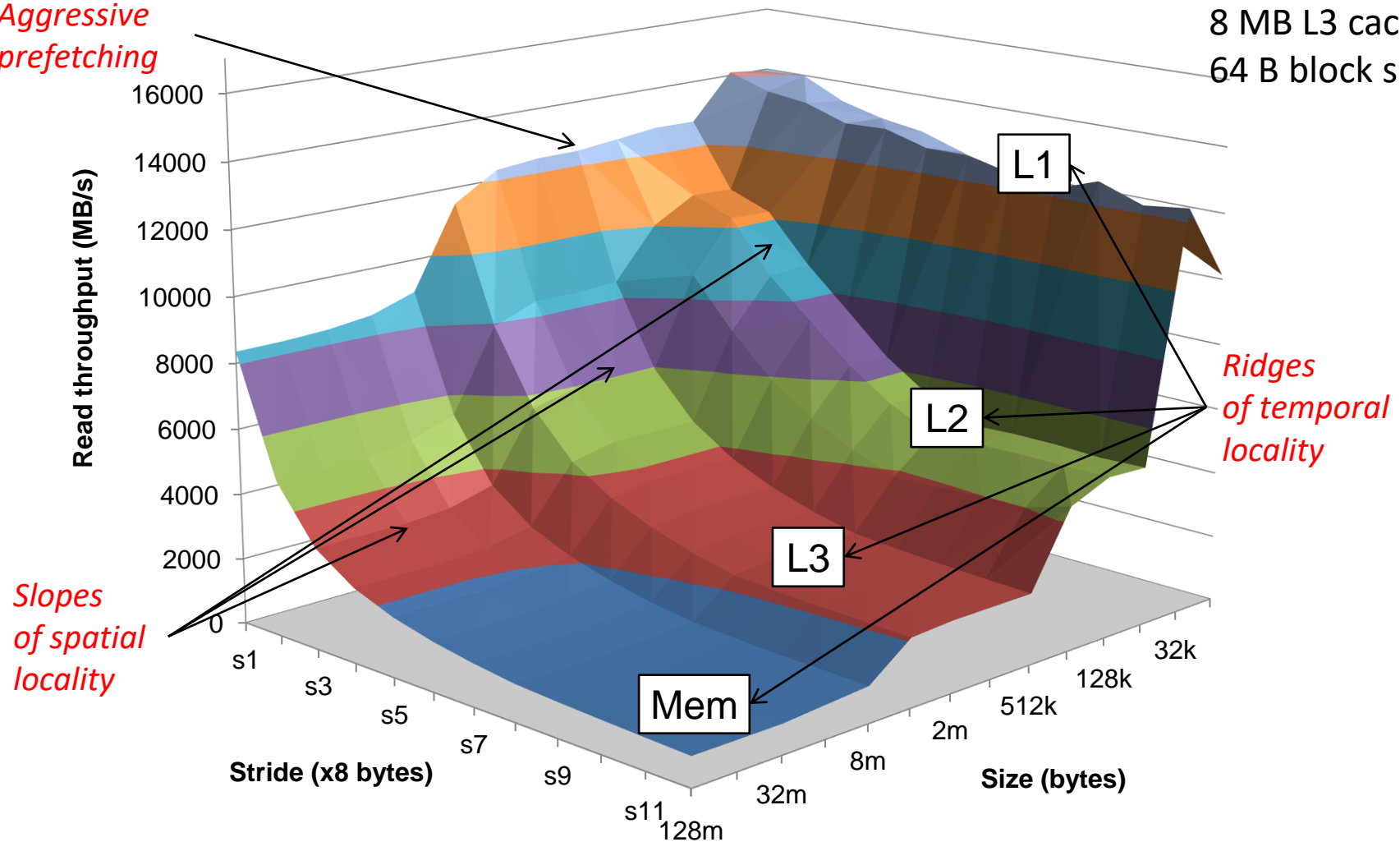
1. Call `test()` once to warm up the caches.

2. Call `test()` again and measure the read throughput (MB/s)

# The Memory Mountain

Core i7 Haswell  
2.1 GHz  
32 KB L1 d-cache  
256 KB L2 cache  
8 MB L3 cache  
64 B block size

*Aggressive  
prefetching*



# Today

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

# Matrix Multiplication Example

- Description:

- Multiply  $N \times N$  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

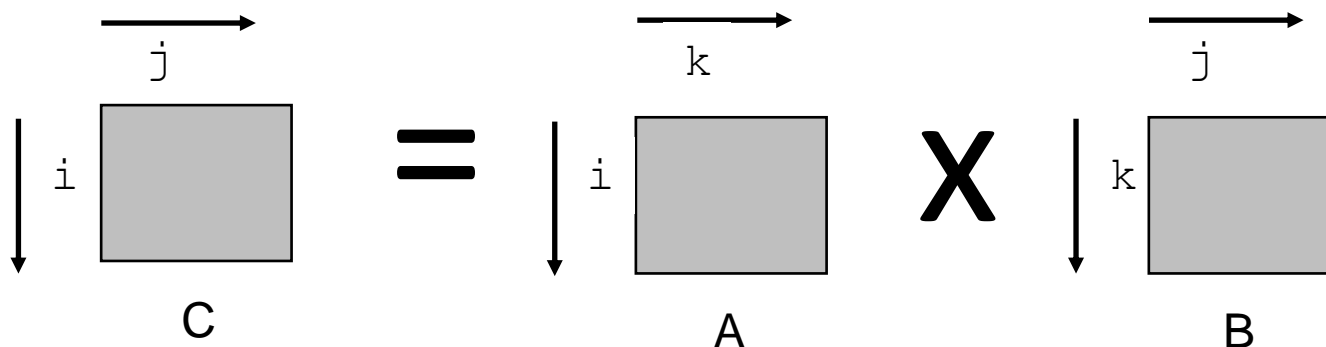
```
/* ijk */
for (i=0; i<n; i++) {
    for (j=0; j<n; j++) {
        sum = 0.0;
        for (k=0; k<n; k++)
            sum += a[i][k] * b[k][j];
        c[i][j] = sum;
    }
}
```

*Variable `sum`  
held in register*

*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)

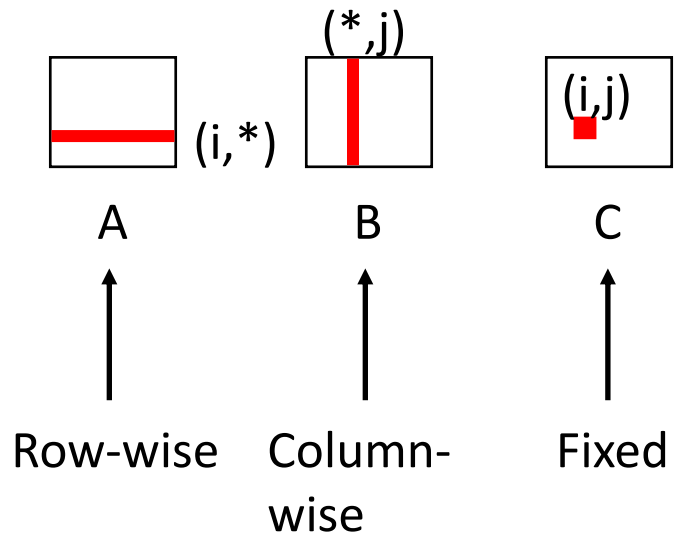
- C arrays allocated in row-major order
  - each row in contiguous memory locations
- Stepping through columns in one row:
  - `for (i = 0; i < N; i++)`  
    `sum += a[0][i];`
  - accesses successive elements
  - if block size ( $B$ ) > `sizeof(aij)` bytes, exploit spatial locality
    - miss rate =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



Misses per inner loop iteration:

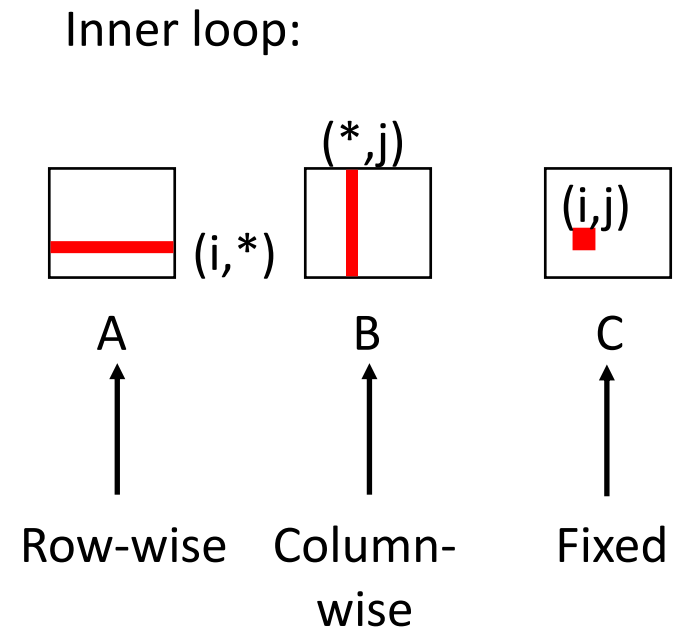
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0



# 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
  }
}
```

*matmult/mm.c*



Misses per inner loop iteration:

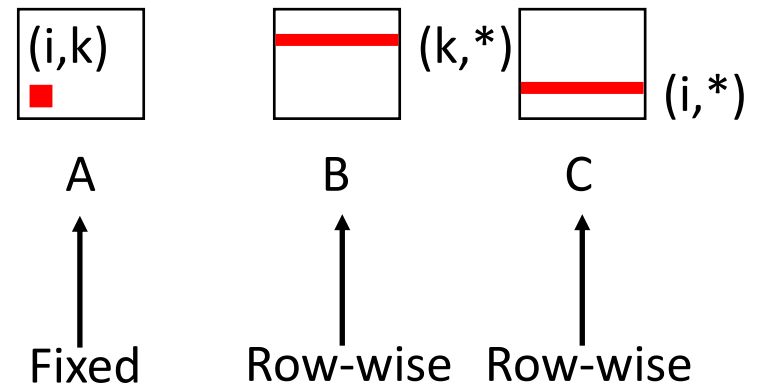
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

# 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];
    }
}
```

*matmult/mm.c*

Inner loop:



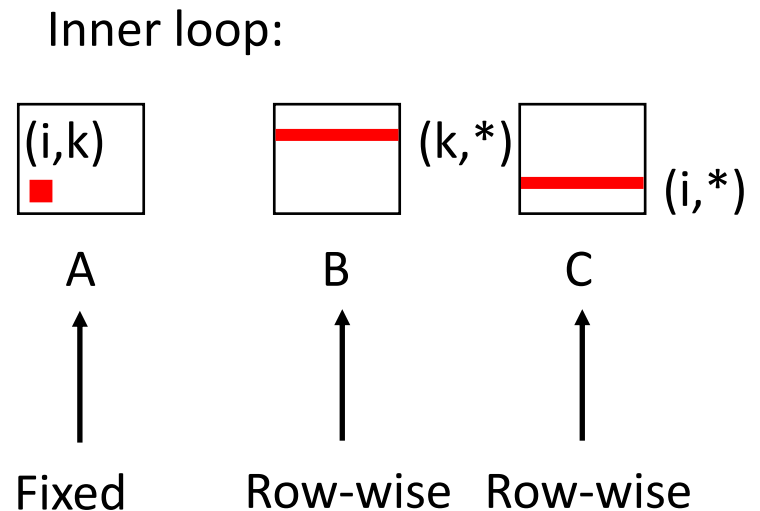
Misses per inner loop iteration:

<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];  
    }  
}
```

*matmult/mm.c*



Misses per inner loop iteration:

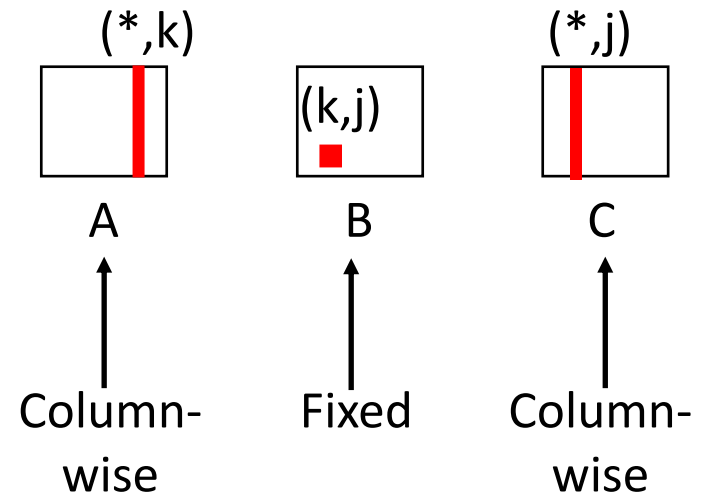
<u>A</u>	<u>B</u>	<u>C</u>
0.0	0.25	0.25

# 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;
    }
}
```

*matmult/mm.c*

Inner loop:



Misses per inner loop iteration:

A  
1.0

B  
0.0

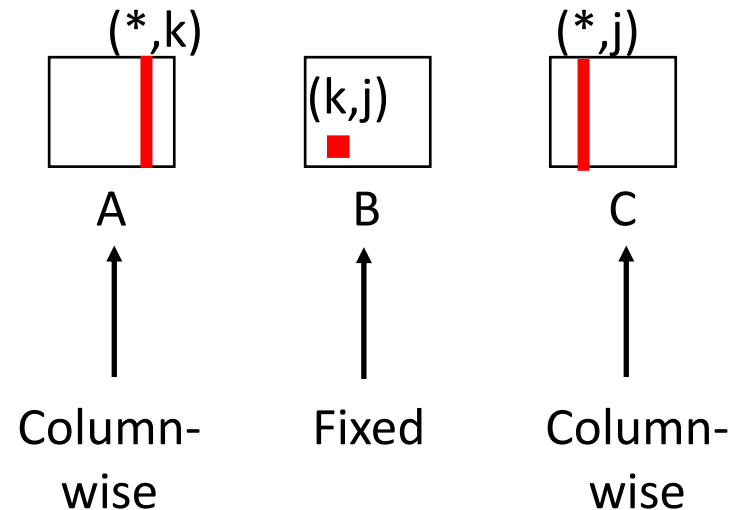
C  
1.0

# 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:



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):

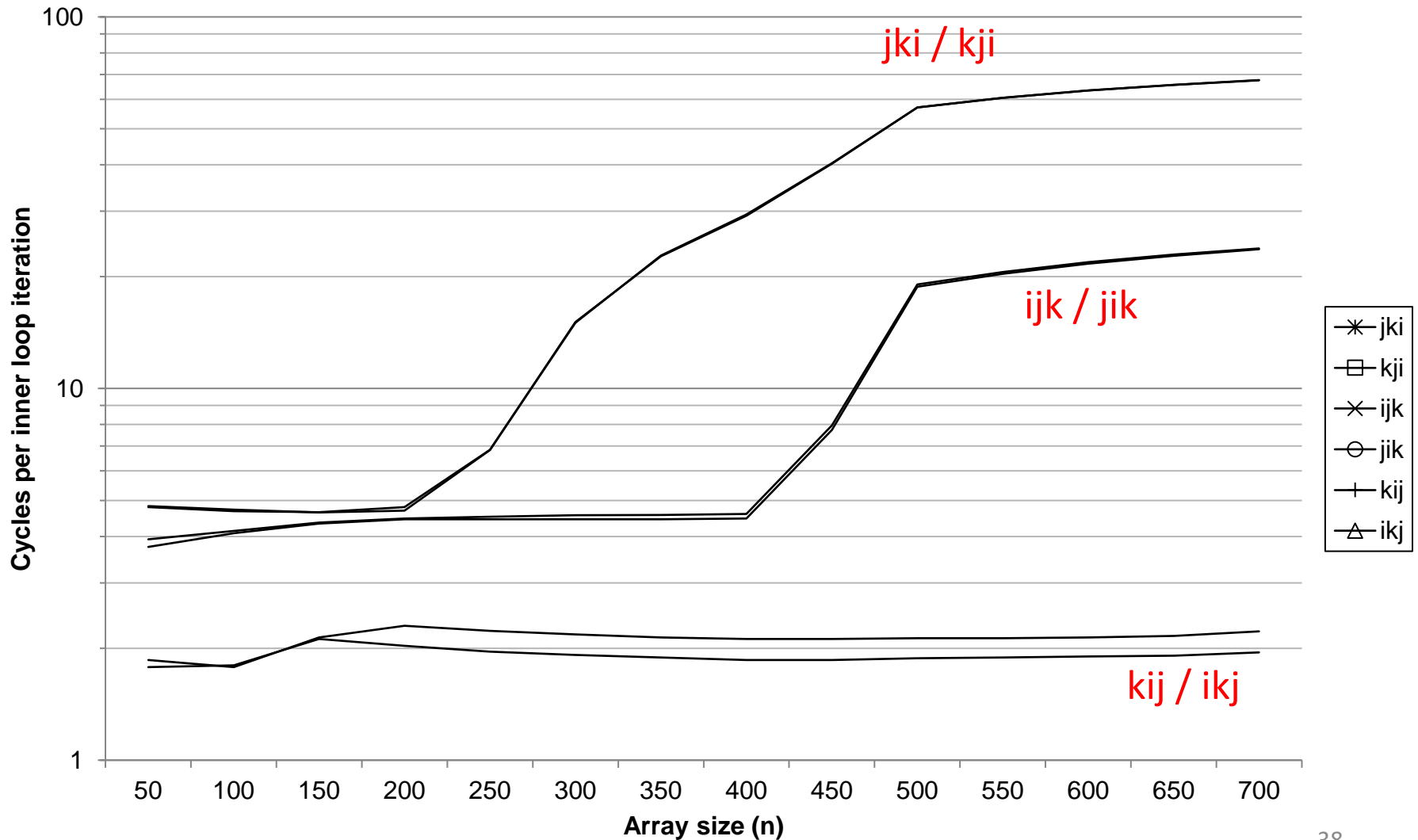
- 2 loads, 1 store
- misses/iter = 0.5

```
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



# Today

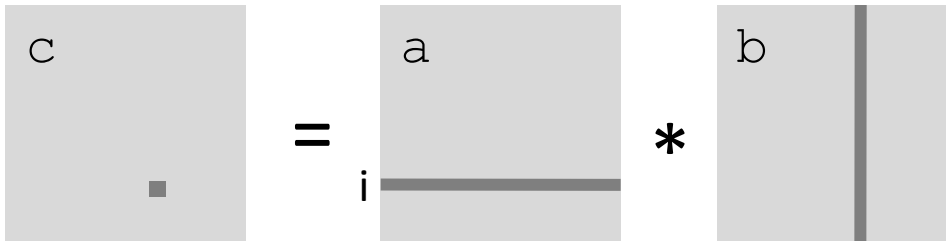
- Cache organization and operation
- Performance impact of caches
  - The memory mountain
  - Rearranging loops to improve spatial locality
  - Using blocking to improve temporal locality



# 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];
}
```



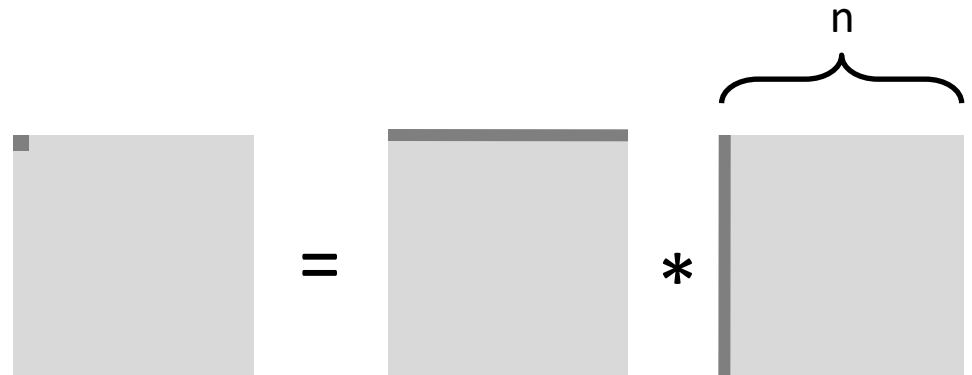
# Cache Miss Analysis

- Assume:

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

- First iteration:

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



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

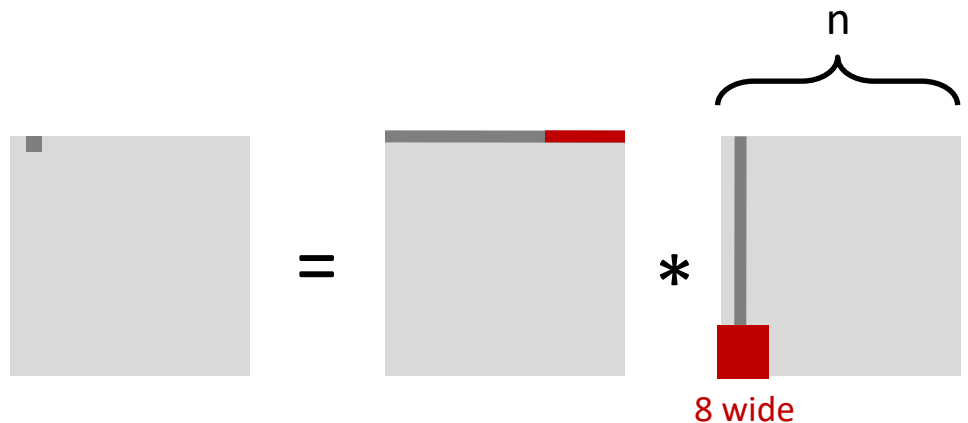


# Cache Miss Analysis

- 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:

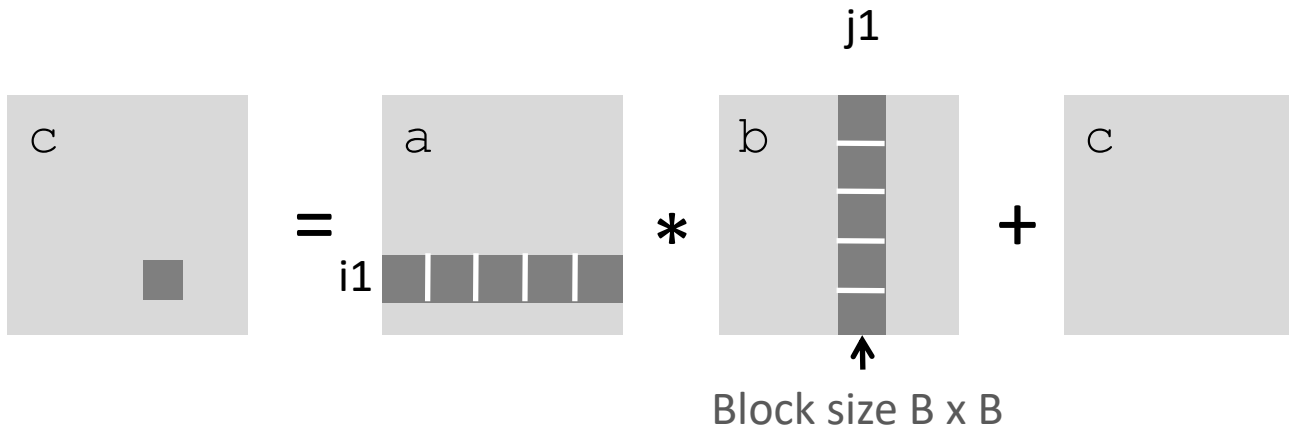
- $9n/8 * n^2 = (9/8) * n^3$

# 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];
}
```

*matmult/bmm.c*



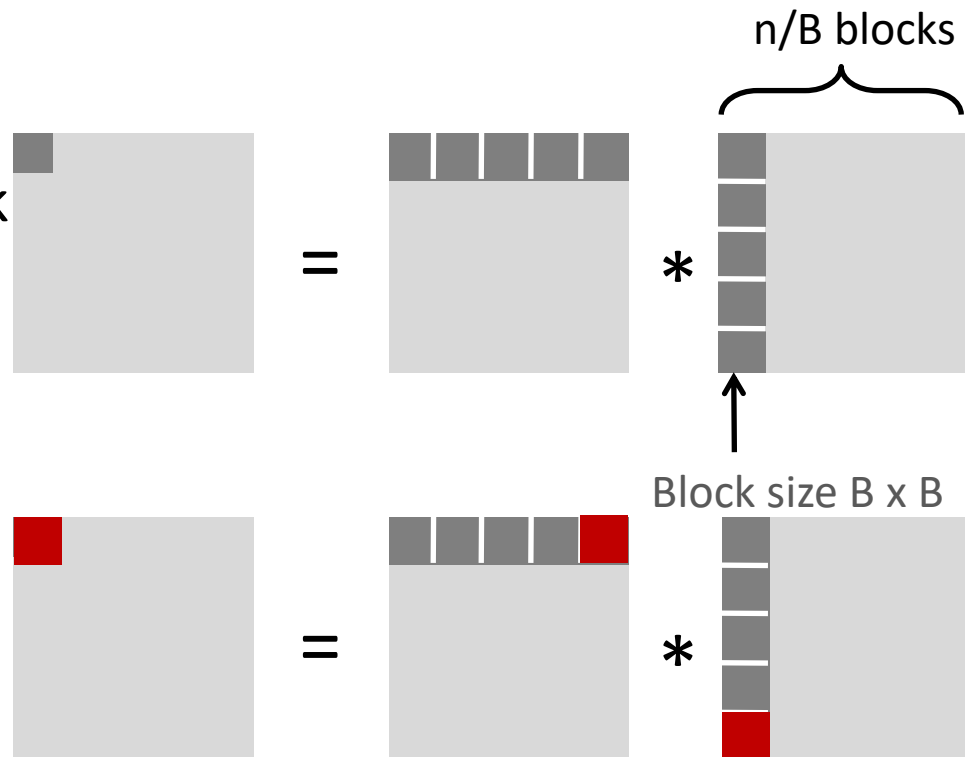
# Cache Miss Analysis

- 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)

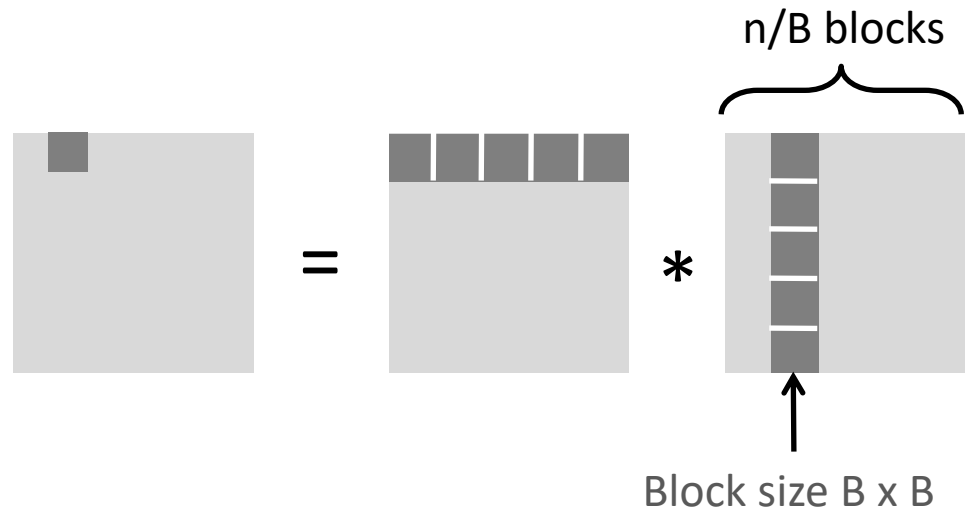
# Cache Miss Analysis

- 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:

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

# 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.



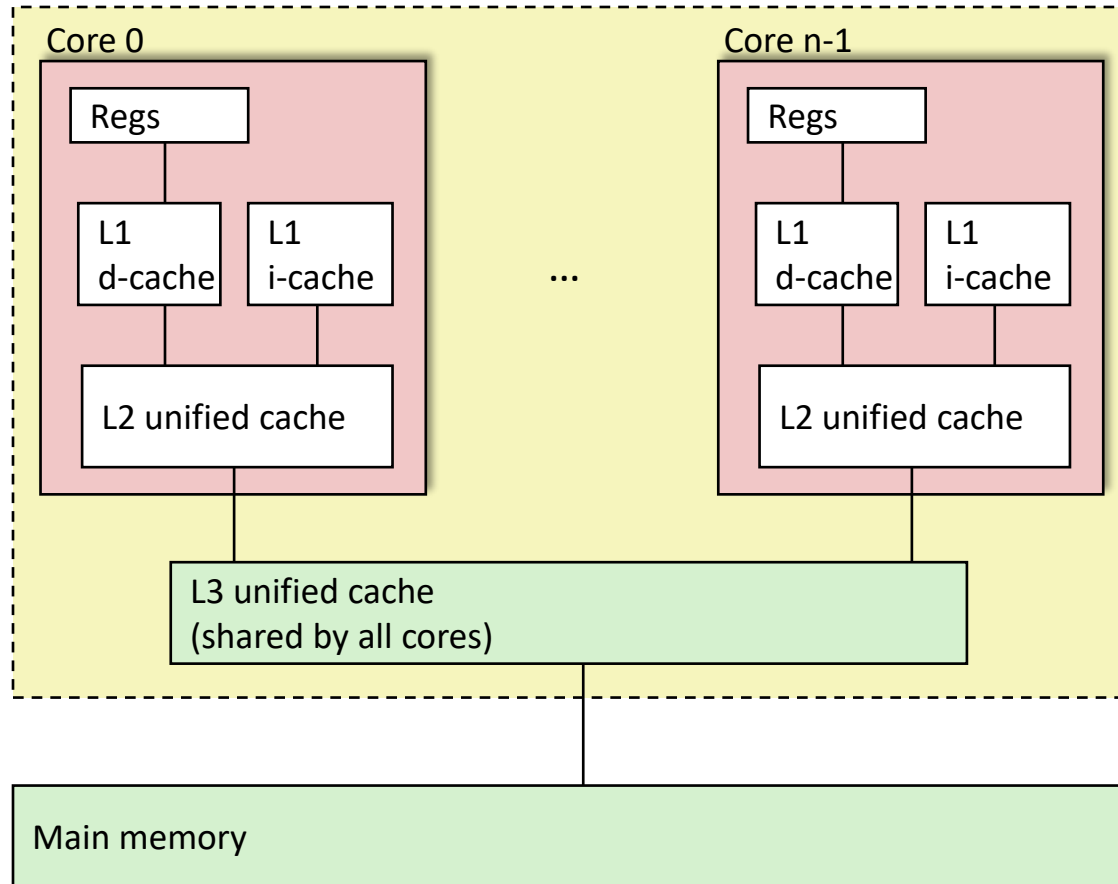
# Parallel Programming and Cache Coherence

The content of this part is mainly from:

Randal E. Bryant and David R. O'Hallaron, "Computer Systems: A Programmer's Perspective," 3/e.

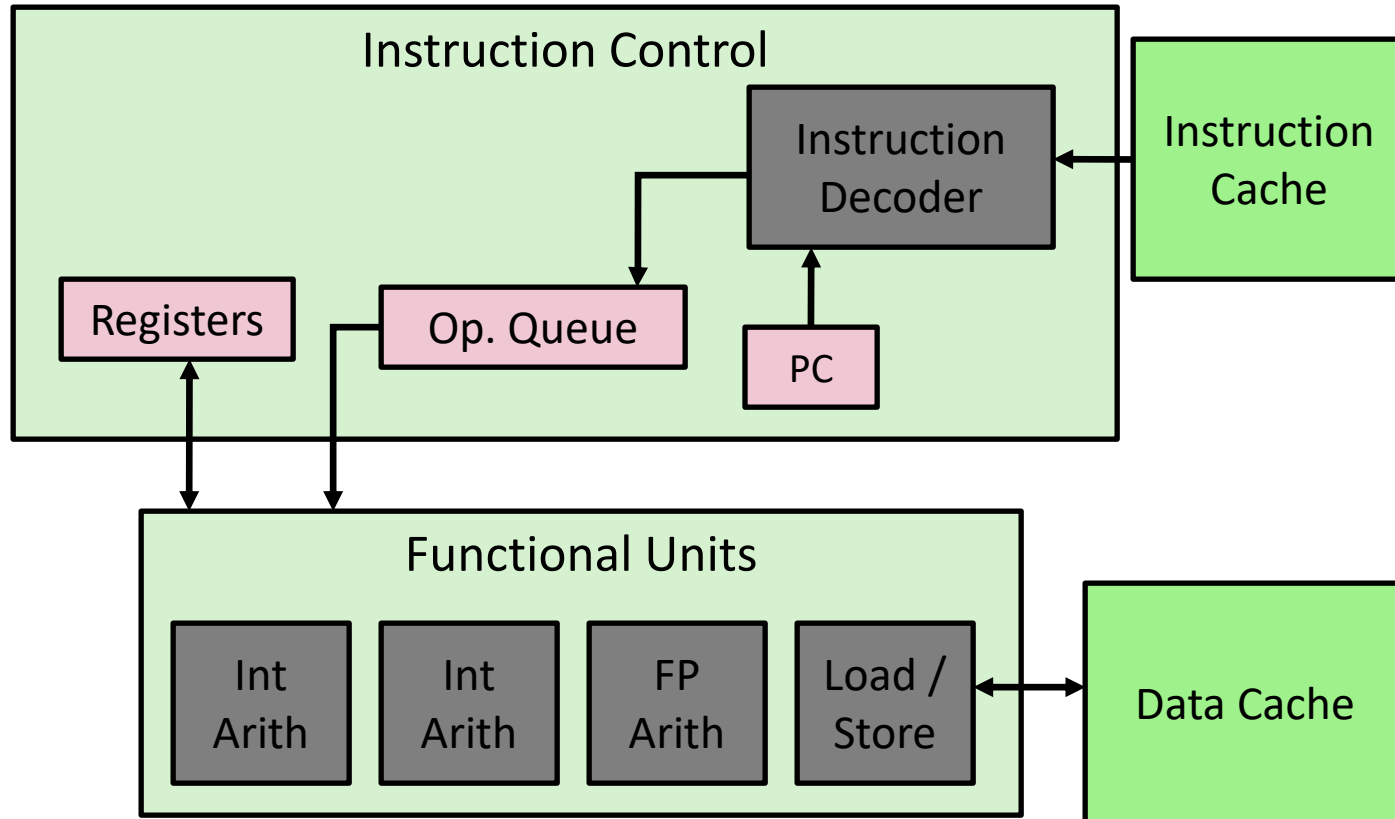
(本節內容改自Prof. Randal E. Bryant and David R. O'Hallaron 26<sup>th</sup> Lectures課程講義)  
(原課程名稱為Thread-Level Parallelism)

# Typical Multicore Processor



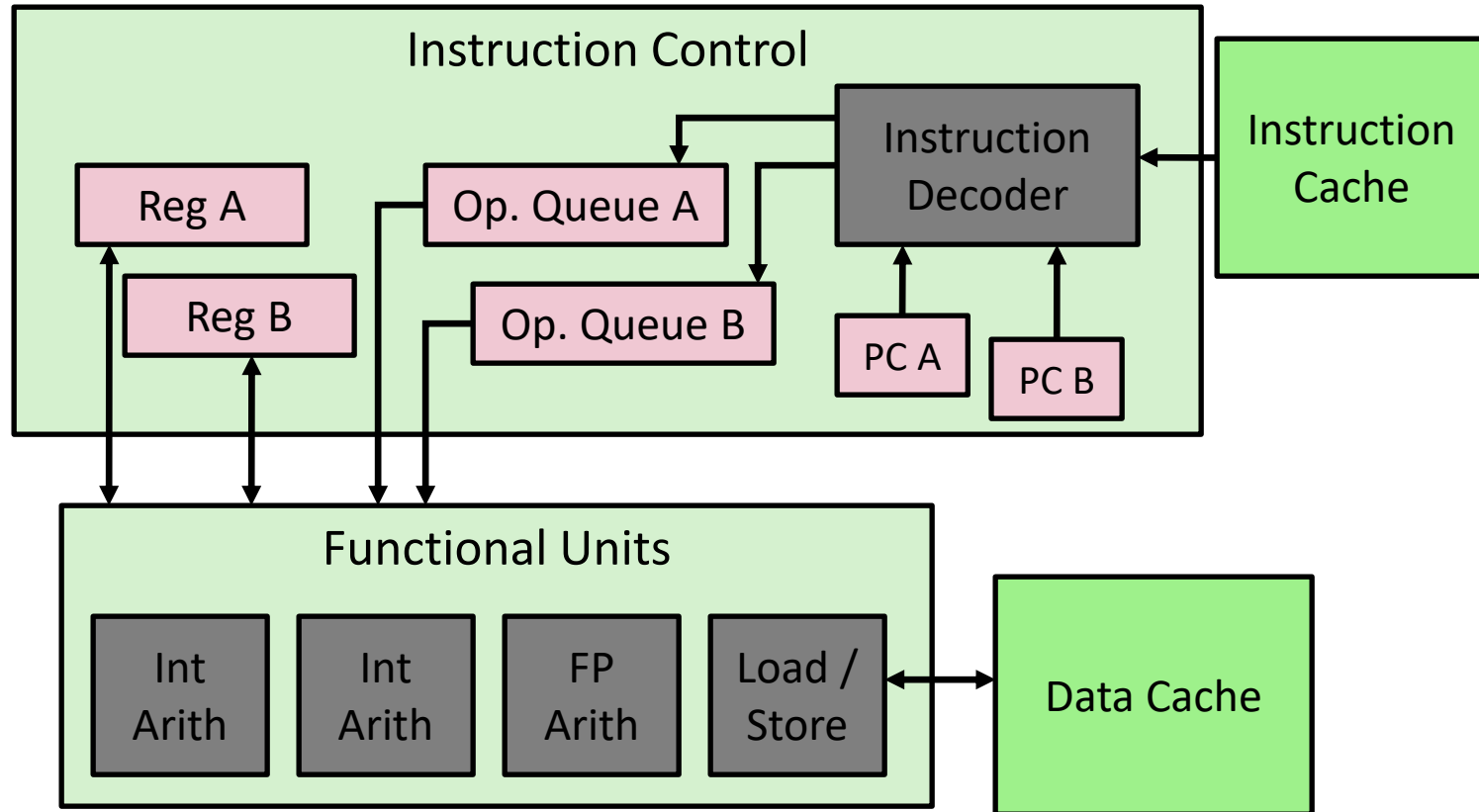
- Multiple processors operating with coherent view of memory

# Out-of-Order Processor Structure



- Instruction control dynamically converts program into stream of operations
- Operations mapped onto functional units to execute in parallel

# Hyperthreading Implementation



- Replicate enough instruction control to process K instruction streams
- K copies of all registers
- Share functional units

# Benchmark Machine

- Get data about machine from `/proc/cpuinfo`
- Shark Machines
  - Intel Xeon E5520 @ 2.27 GHz
  - Nehalem, ca. 2010
  - 8 Cores
  - Each can do 2x hyperthreading

# Example 1: Parallel Summation

- Sum numbers  $0, \dots, n-1$ 
  - Should add up to  $((n-1)*n)/2$
- Partition values  $1, \dots, n-1$  into  $t$  ranges
  - $\lfloor n/t \rfloor$  values in each range
  - Each of  $t$  threads processes 1 range
  - For simplicity, assume  $n$  is a multiple of  $t$
- Let's consider different ways that multiple threads might work on their assigned ranges in parallel

# First attempt: psum-mutex

- Simplest approach: Threads sum into a global variable protected by a semaphore mutex.

```
void *sum_mutex(void *vargp); /* Thread routine */

/* Global shared variables */
long gsum = 0;                /* Global sum */
long nelems_per_thread;      /* Number of elements to sum */
sem_t mutex;                  /* Mutex to protect global sum */

int main(int argc, char **argv)
{
    long i, nelems, log_nelems, nthreads, myid[MAXTHREADS];
    pthread_t tid[MAXTHREADS];

    /* Get input arguments */
    nthreads = atoi(argv[1]);
    log_nelems = atoi(argv[2]);
    nelems = (1L << log_nelems);
    nelems_per_thread = nelems / nthreads;
    sem_init(&mutex, 0, 1);
```

psum-mutex.c

## psum-mutex (cont)

- Simplest approach: Threads sum into a global variable protected by a semaphore mutex.

```
/* Create peer threads and wait for them to finish */
for (i = 0; i < nthreads; i++) {
    myid[i] = i;
    Pthread_create(&tid[i], NULL, sum_mutex, &myid[i]);
}
for (i = 0; i < nthreads; i++)
Pthread_join(tid[i], NULL);

/* Check final answer */
if (gsum != (nelems * (nelems-1))/2)
    printf("Error: result=%ld\n", gsum);

return 0;
}
```

psum-mutex.c



# psum-mutex Thread Routine

- Simplest approach: Threads sum into a global variable protected by a semaphore mutex.

```
/* Thread routine for psum-mutex.c */
void *sum_mutex(void *vargp)
{
    long myid = *((long *)vargp);          /* Extract thread ID */
    long start = myid * nelems_per_thread; /* Start element index */
    long end = start + nelems_per_thread;  /* End element index */
    long i;

    for (i = start; i < end; i++) {
        P(&mutex);
        gsum += i;
        V(&mutex);
    }
    return NULL;
}
```

psum-mutex.c

# psum-mutex Performance

- Shark machine with 8 cores,  $n=2^{31}$

Threads (Cores)	1 (1)	2 (2)	4 (4)	8 (8)	16 (8)
psum-mutex (secs)	51	456	790	536	681

## ■ Nasty surprise:

- Single thread is very slow
- Gets slower as we use more cores

## Next Attempt: psum-array

- Peer thread `i` sums into global array element `psum[i]`
- Main waits for theads to finish, then sums elements of `psum`
- Eliminates need for mutex synchronization

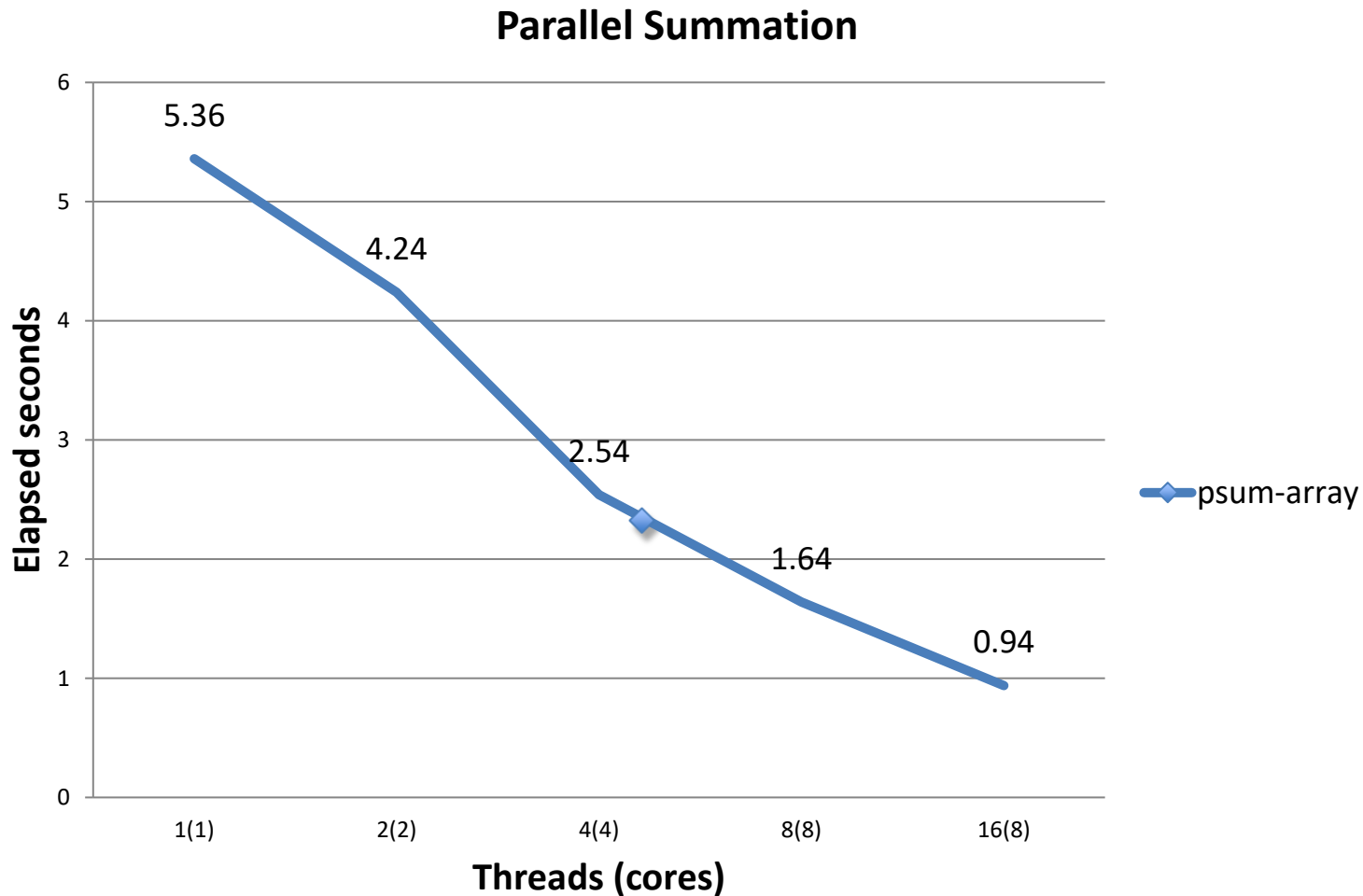
```
/* Thread routine for psum-array.c */
void *sum_array(void *vargp)
{
    long myid = *((long *)vargp);          /* Extract thread ID */
    long start = myid * nelems_per_thread; /* Start element index */
    long end = start + nelems_per_thread;  /* End element index */
    long i;

    for (i = start; i < end; i++) {
        psum[myid] += i;
    }
    return NULL;
}
```

psum-array.c

# psum-array Performance

- Orders of magnitude faster than psum-mutex



# Next Attempt: psum-local

- Reduce memory references by having peer thread i sum into a local variable (register)

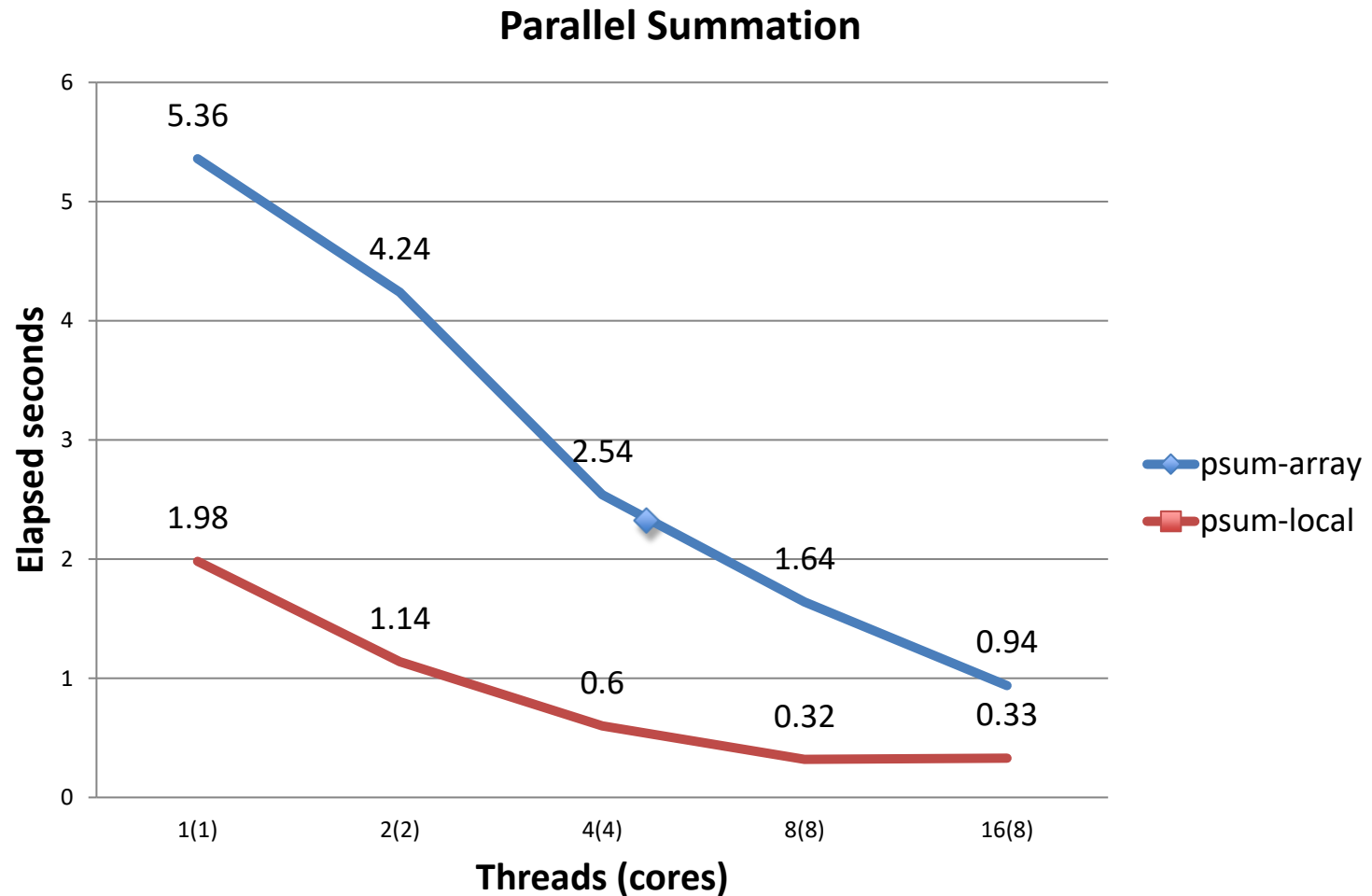
```
/* Thread routine for psum-local.c */
void *sum_local(void *vargp)
{
    long myid = *((long *)vargp);          /* Extract thread ID */
    long start = myid * nelems_per_thread; /* Start element index */
    long end = start + nelems_per_thread;  /* End element index */
    long i, sum = 0;

    for (i = start; i < end; i++) {
        sum += i;
    }
    psum[myid] = sum;
    return NULL;
}
```

psum-local.c

# psum-local Performance

- Significantly faster than psum-array



# Characterizing Parallel Program Performance

- $p$  processor cores,  $T_k$  is the running time using  $k$  cores
- **Def. Speedup:**  $S_p = T_1 / T_p$ 
  - $S_p$  is *relative speedup* if  $T_1$  is running time of parallel version of the code running on 1 core.
  - $S_p$  is *absolute speedup* if  $T_1$  is running time of sequential version of code running on 1 core.
  - Absolute speedup is a much truer measure of the benefits of parallelism.
- **Def. Efficiency:**  $E_p = S_p / p = T_1 / (pT_p)$ 
  - Reported as a percentage in the range (0, 100].
  - Measures the overhead due to parallelization

# Performance of `psum-local`

Threads (t)	1	2	4	8	16
Cores (p)	1	2	4	8	8
Running time ( $T_p$ )	1.98	1.14	0.60	0.32	0.33
Speedup ( $S_p$ )	1	1.74	3.30	6.19	6.00
Efficiency ( $E_p$ )	100%	87%	82%	77%	75%

- Efficiencies OK, not great
- Our example is easily parallelizable
- Real codes are often much harder to parallelize
  - e.g., parallel quicksort later in this lecture



# Amdahl's Law

- Gene Amdahl (Nov. 16, 1922 – Nov. 10, 2015)
- Captures the difficulty of using parallelism to speed things up.
- Overall problem
  - $T$  Total sequential time required
  - $p$  Fraction of total that can be sped up ( $0 \leq p \leq 1$ )
  - $k$  Speedup factor
- Resulting Performance
  - $T_k = pT/k + (1-p)T$ 
    - Portion which can be sped up runs  $k$  times faster
    - Portion which cannot be sped up stays the same
  - Least possible running time:
    - $k = \infty$
    - $T_\infty = (1-p)T$

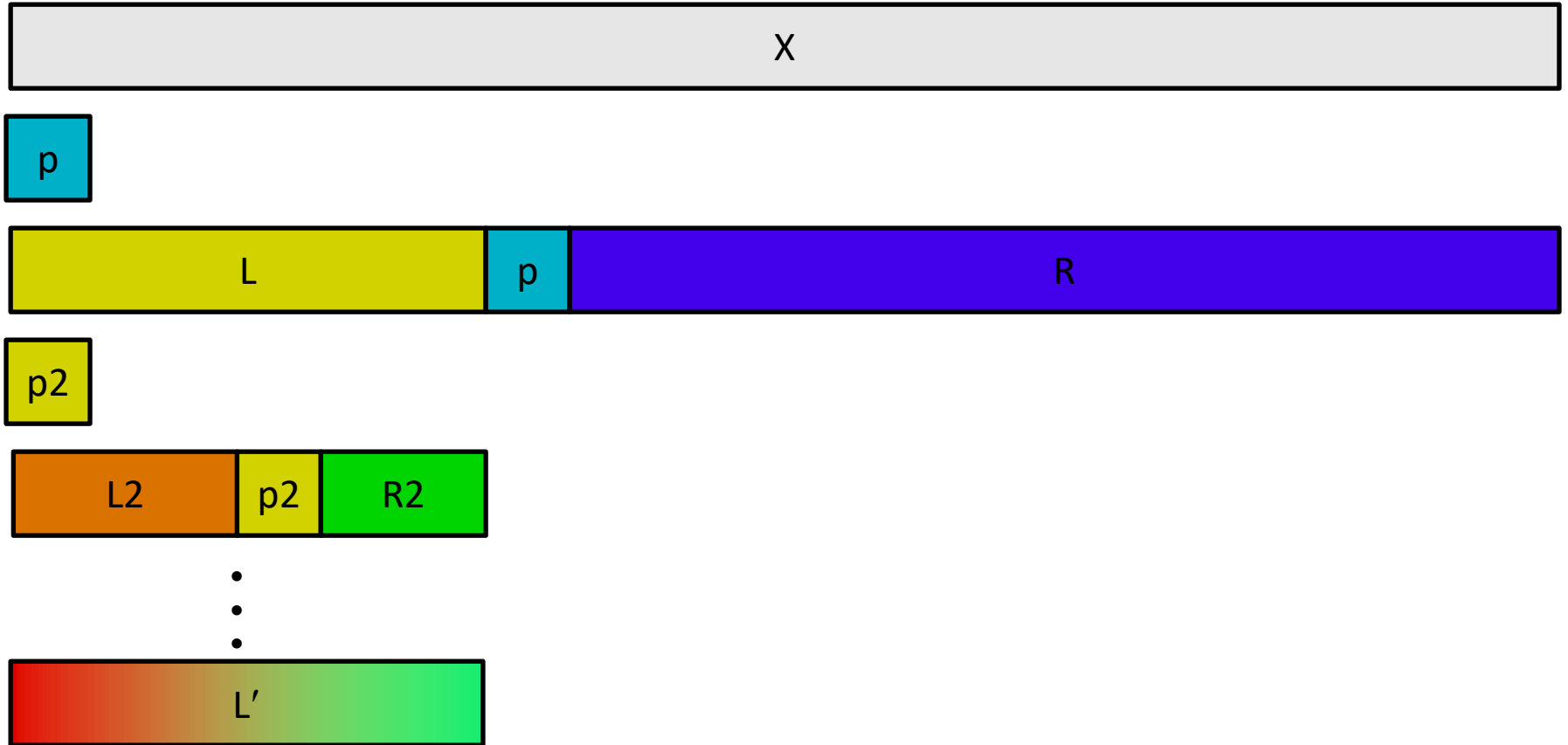
# Amdahl's Law Example

- Overall problem
  - $T = 10$  Total time required
  - $p = 0.9$  Fraction of total which can be sped up
  - $k = 9$  Speedup factor
- Resulting Performance
  - $T_9 = 0.9 * 10/9 + 0.1 * 10 = 1.0 + 1.0 = 2.0$
  - Least possible running time:
    - $T_{\infty} = 0.1 * 10.0 = 1.0$

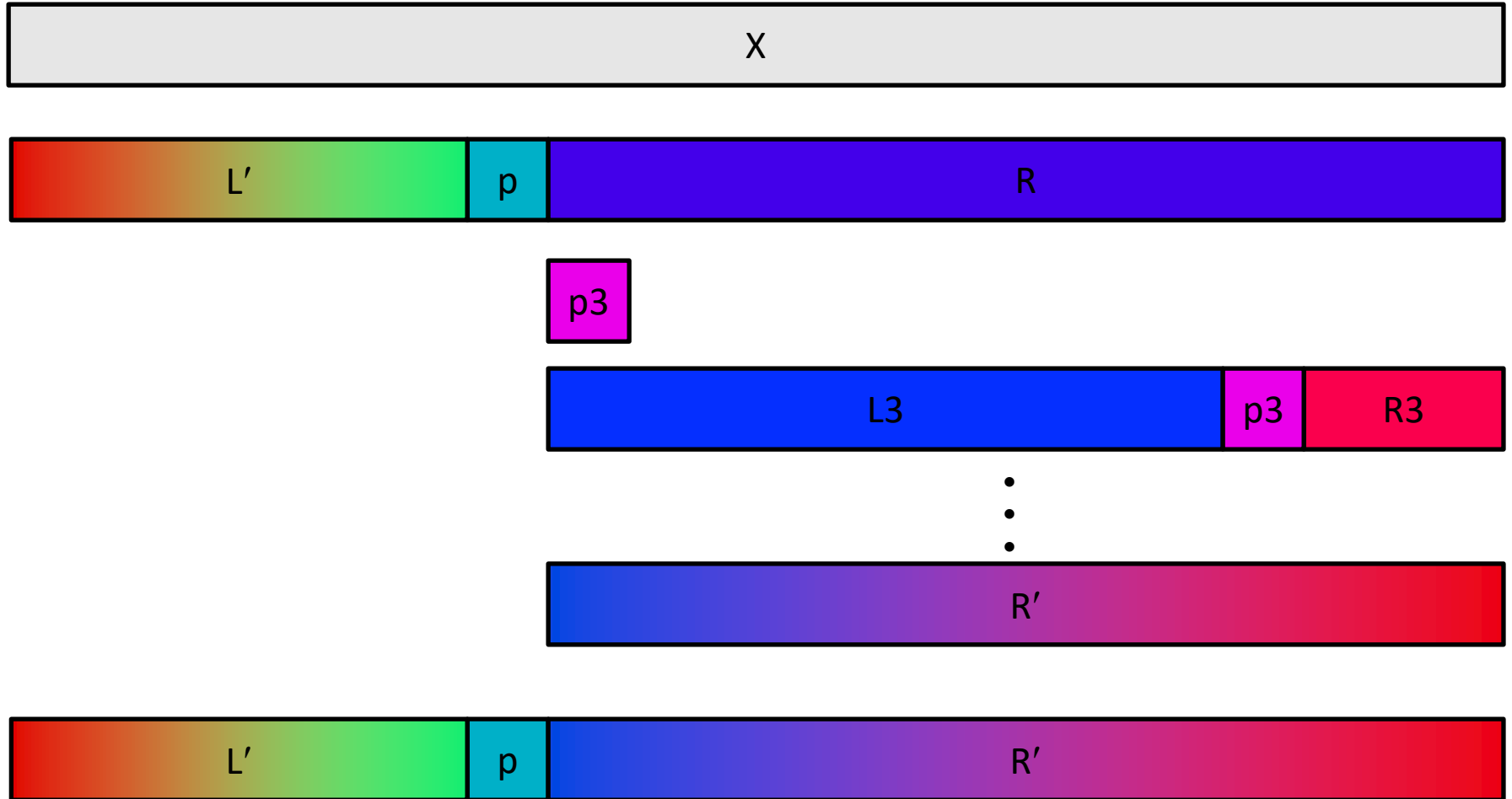
# A More Substantial Example: Sort

- Sort set of  $N$  random numbers
- Multiple possible algorithms
  - Use parallel version of quicksort
- Sequential quicksort of set of values  $X$ 
  - Choose “pivot”  $p$  from  $X$
  - Rearrange  $X$  into
    - $L$ : Values  $\leq p$
    - $R$ : Values  $\geq p$
  - Recursively sort  $L$  to get  $L'$
  - Recursively sort  $R$  to get  $R'$
  - Return  $L' : p : R'$

# Sequential Quicksort Visualized



# Sequential Quicksort Visualized



# Sequential Quicksort Code

```
void qsort_serial(data_t *base, size_t nele) {
    if (nele <= 1)
        return;
    if (nele == 2) {
        if (base[0] > base[1])
            swap(base, base+1);
        return;
    }

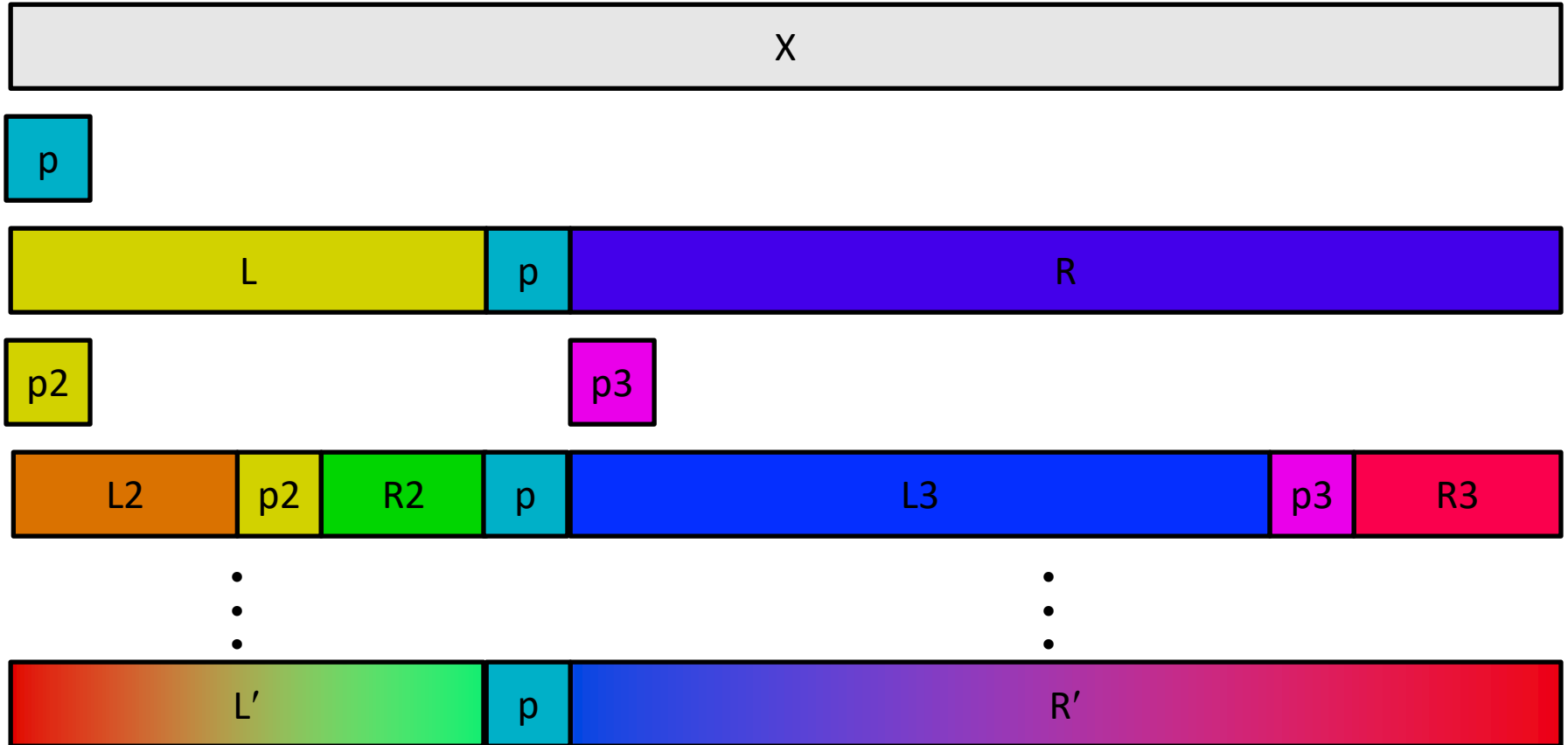
    /* Partition returns index of pivot */
    size_t m = partition(base, nele);
    if (m > 1)
        qsort_serial(base, m);
    if (nele-1 > m+1)
        qsort_serial(base+m+1, nele-m-1);
}
```

- Sort nele elements starting at base
  - Recursively sort L or R if has more than one element

# Parallel Quicksort

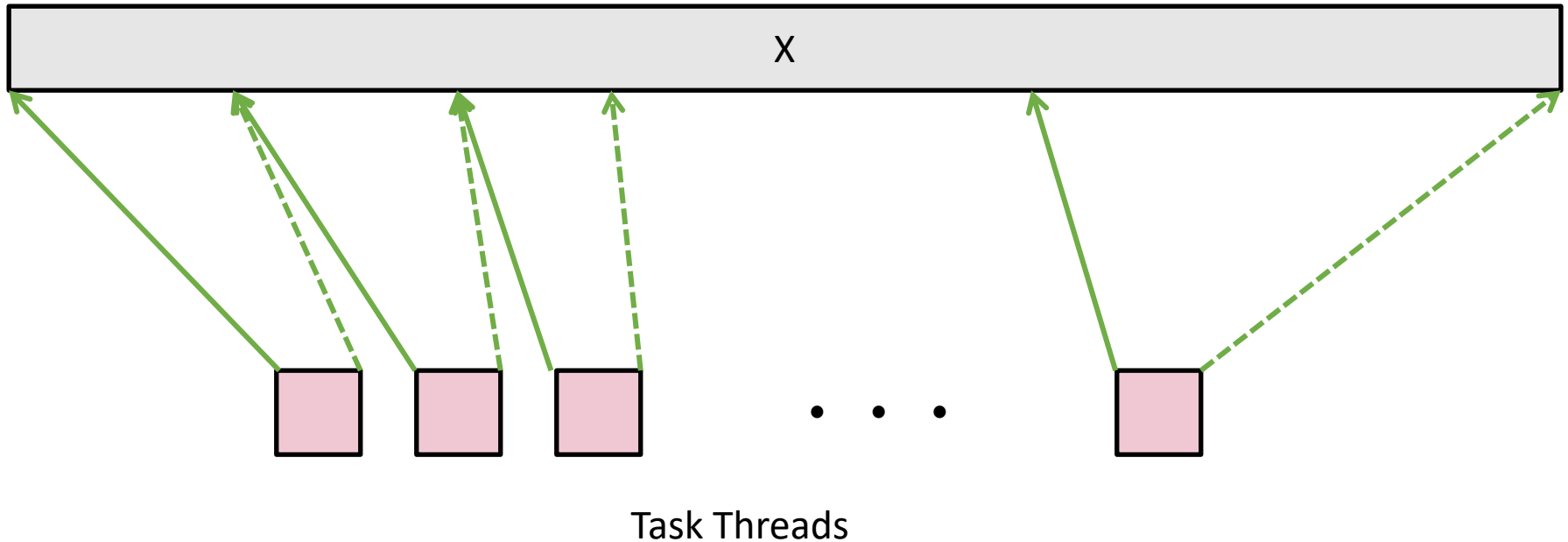
- Parallel quicksort of set of values  $X$ 
  - If  $N \leq N_{\text{thresh}}$ , do sequential quicksort
  - Else
    - Choose “pivot”  $p$  from  $X$
    - Rearrange  $X$  into
      - $L$ : Values  $\leq p$
      - $R$ : Values  $\geq p$
    - Recursively spawn separate threads
      - Sort  $L$  to get  $L'$
      - Sort  $R$  to get  $R'$
    - Return  $L' : p : R'$

# Parallel Quicksort Visualized



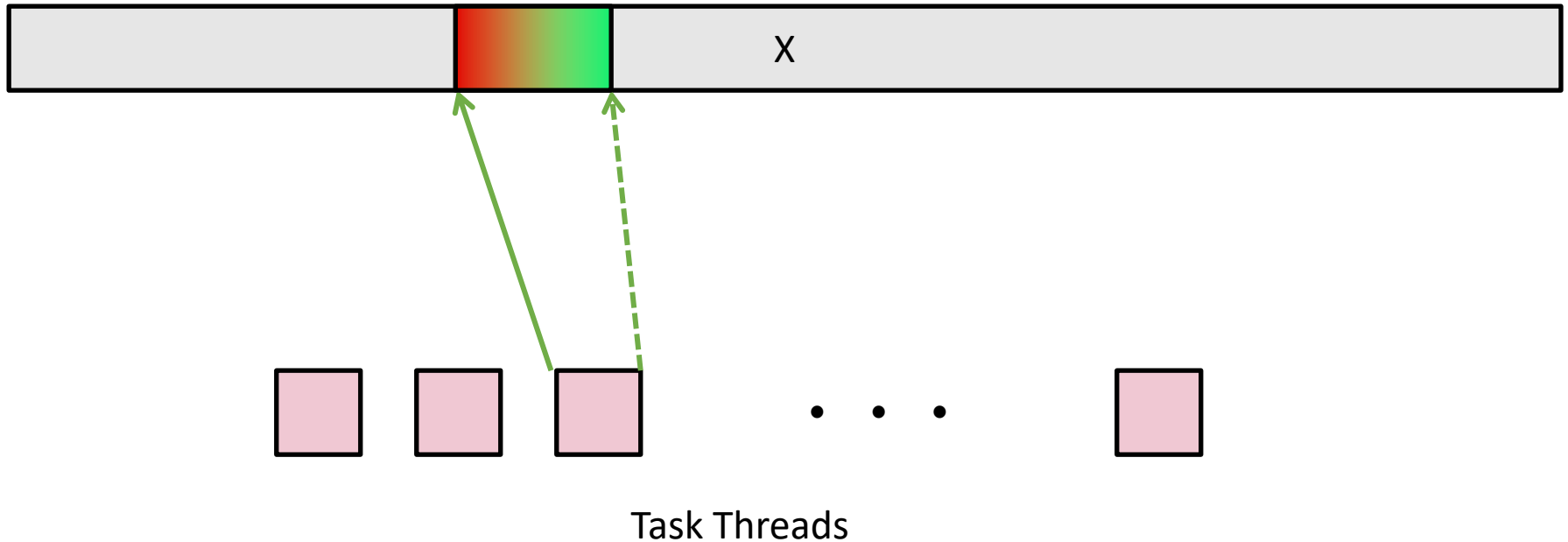


# Thread Structure: Sorting Tasks



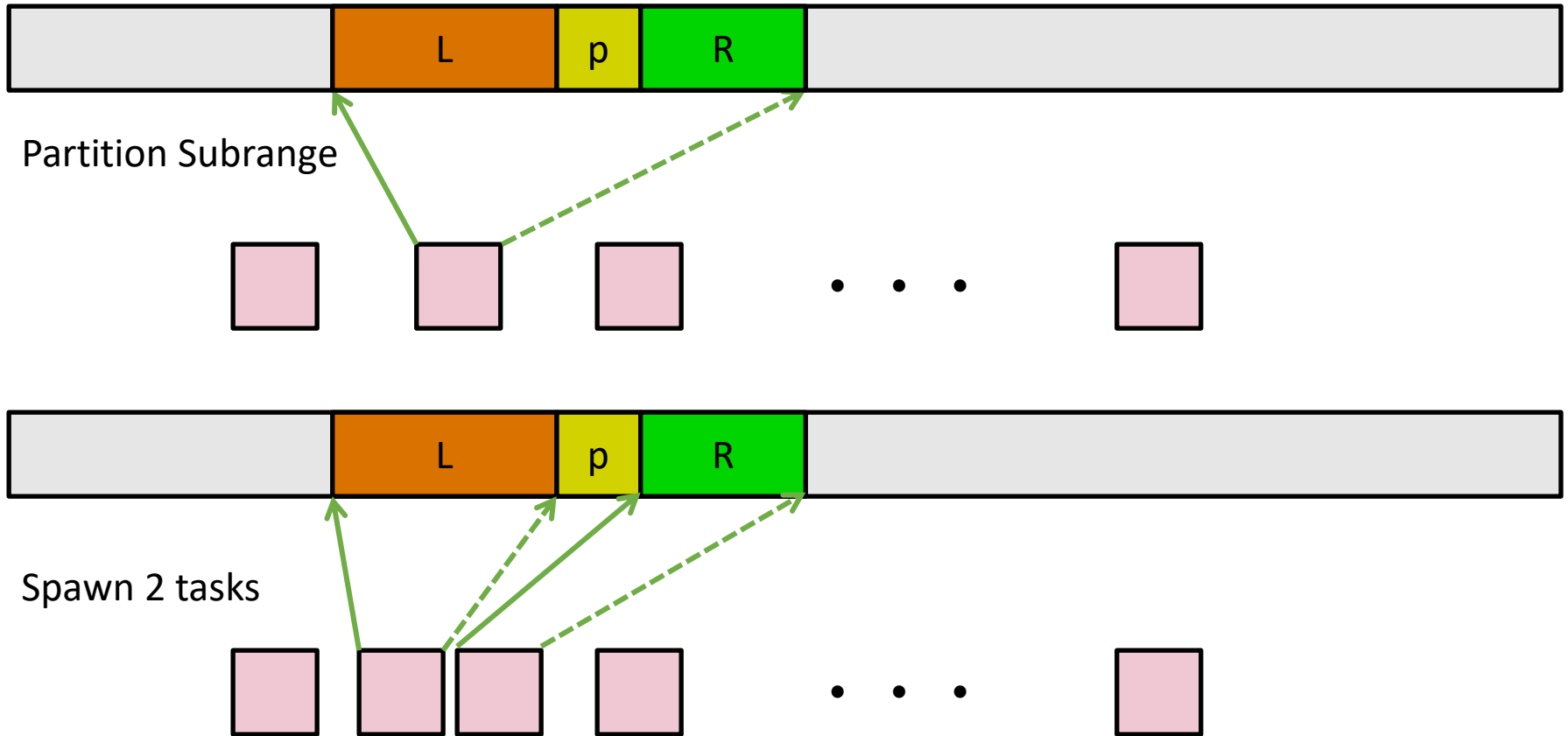
- Task: Sort subrange of data
  - Specify as:
    - **base**: Starting address
    - **nele**: Number of elements in subrange
- Run as separate thread

# Small Sort Task Operation



- Sort subrange using serial quicksort

# Large Sort Task Operation



# Top-Level Function (Simplified)

```
void tqsort(data_t *base, size_t nele) {  
    init_task(nele);  
    global_base = base;  
    global_end = global_base + nele - 1;  
    task_queue_ptr tq = new_task_queue();  
    tqsort_helper(base, nele, tq);  
    join_tasks(tq);  
    free_task_queue(tq);  
}
```

- Sets up data structures
- Calls recursive sort routine
- Keeps joining threads until none left
- Frees data structures

# Recursive sort routine (Simplified)

```
/* Multi-threaded quicksort */
static void tqsort_helper(data_t *base, size_t nele,
                          task_queue_ptr tq) {
    if (nele <= nele_max_sort_serial) {
        /* Use sequential sort */
        qsort_serial(base, nele);
        return;
    }
    sort_task_t *t = new_task(base, nele, tq);
    spawn_task(tq, sort_thread, (void *) t);
}
```

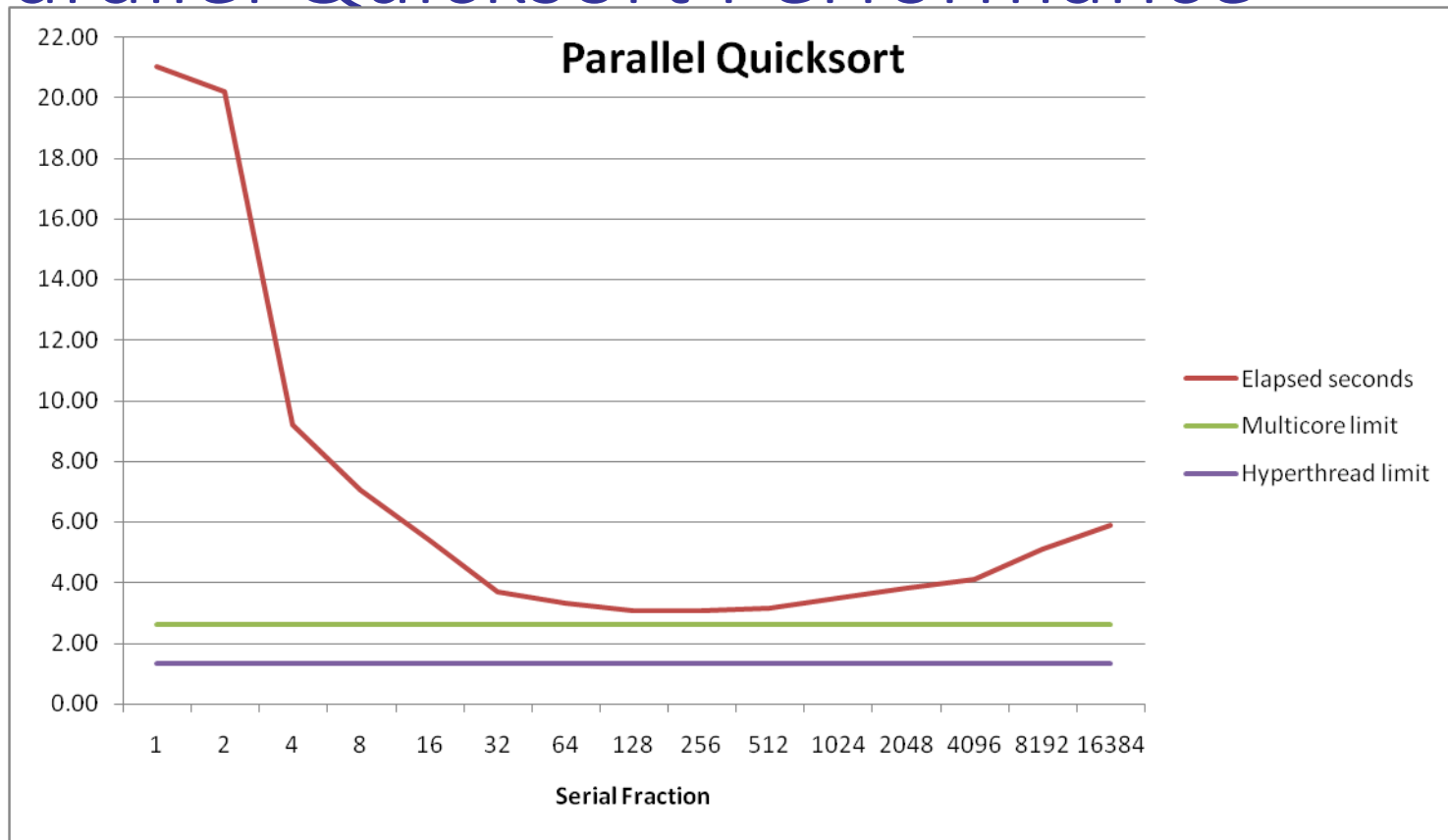
- Small partition: Sort serially
- Large partition: Spawn new sort task

# Sort task thread (Simplified)

```
/* Thread routine for many-threaded quicksort */
static void *sort_thread(void *vargp) {
    sort_task_t *t = (sort_task_t *) vargp;
    data_t *base = t->base;
    size_t nele = t->nele;
    task_queue_ptr tq = t->tq;
    free(vargp);
    size_t m = partition(base, nele);
    if (m > 1)
        tqsort_helper(base, m, tq);
    if (nele-1 > m+1)
        tqsort_helper(base+m+1, nele-m-1, tq);
    return NULL;
}
```

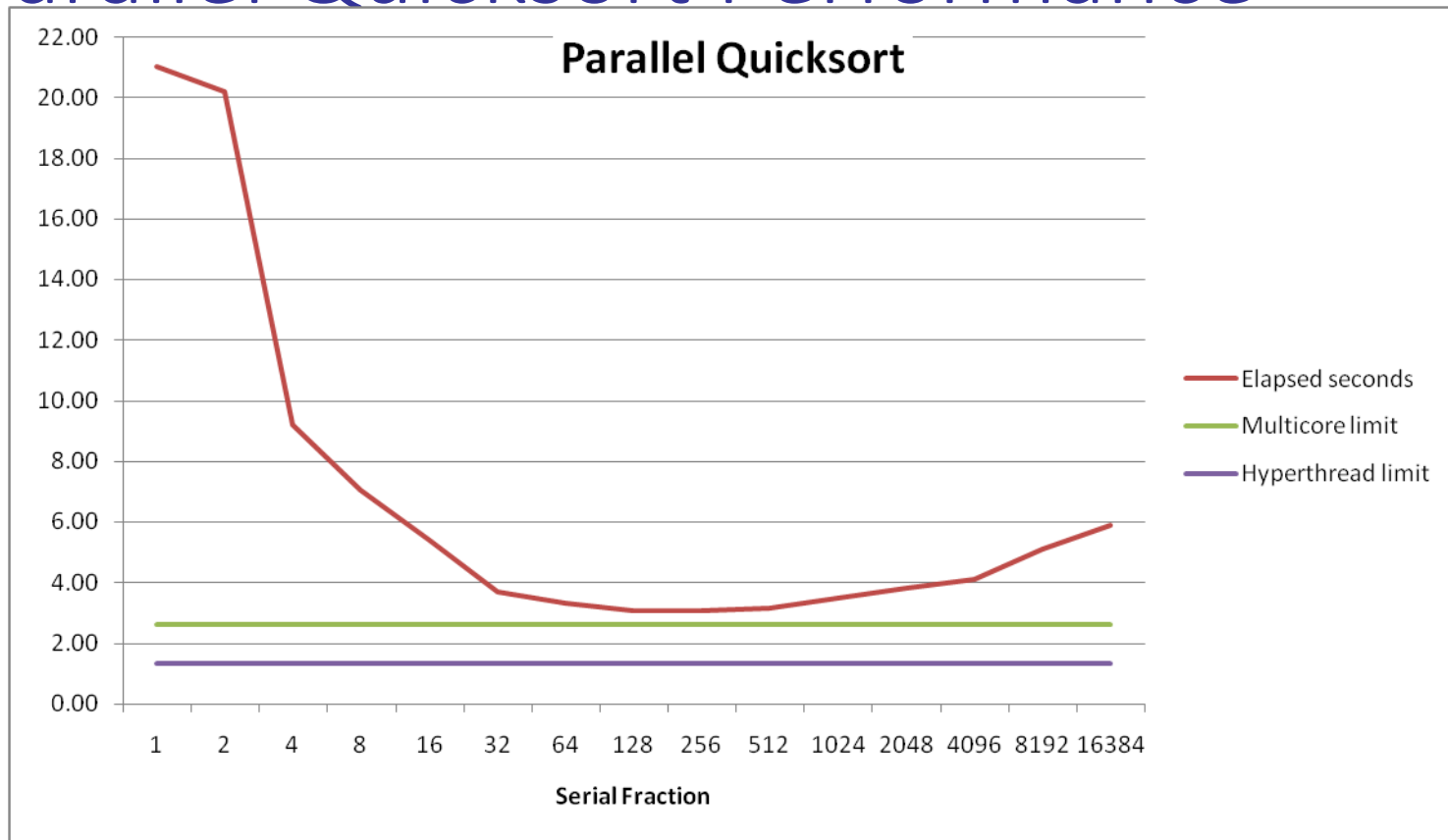
- Get task parameters
- Perform partitioning step
- Call recursive sort routine on each partition

# Parallel Quicksort Performance



- Serial fraction: Fraction of input at which do serial sort
- Sort  $2^{27}$  (134,217,728) random values
- Best speedup = 6.84X

# Parallel Quicksort Performance



- Good performance over wide range of fraction values
  - F too small: Not enough parallelism
  - F too large: Thread overhead + run out of thread memory



# Amdahl's Law & Parallel Quicksort

- Sequential bottleneck
  - Top-level partition: No speedup
  - Second level:  $\leq 2X$  speedup
  - $k^{\text{th}}$  level:  $\leq 2^{k-1}X$  speedup
- Implications
  - Good performance for small-scale parallelism
  - Would need to parallelize partitioning step to get large-scale parallelism
    - Parallel Sorting by Regular Sampling
      - H. Shi & J. Schaeffer, J. Parallel & Distributed Computing, 1992

# Parallelizing Partitioning Step

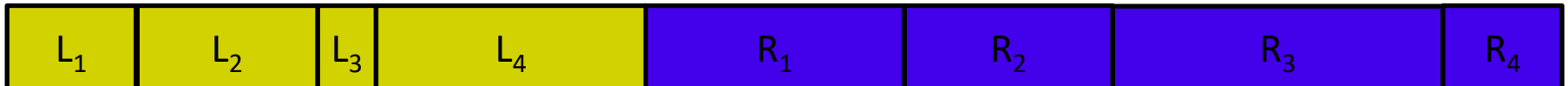


$p$

Parallel partitioning based on global  $p$



Reassemble into partitions



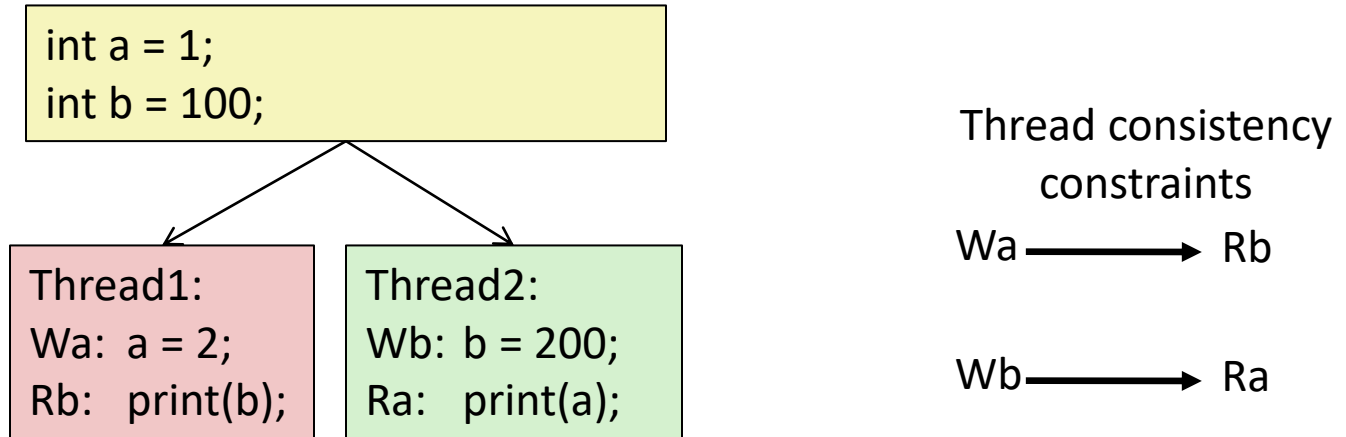
# Experience with Parallel Partitioning

- Could not obtain speedup
- Speculate: Too much data copying
  - Could not do everything within source array
  - Set up temporary space for reassembling partition

# Lessons Learned

- Must have parallelization strategy
  - Partition into  $K$  independent parts
  - Divide-and-conquer
- Inner loops must be synchronization free
  - Synchronization operations very expensive
- Beware of Amdahl's Law
  - Serial code can become bottleneck
- You can do it!
  - Achieving modest levels of parallelism is not difficult
  - Set up experimental framework and test multiple strategies

# Memory Consistency



- What are the possible values printed?
  - Depends on memory consistency model
  - Abstract model of how hardware handles concurrent accesses
- Sequential consistency
  - Overall effect consistent with each individual thread
  - Otherwise, arbitrary interleaving

# Sequential Consistency Example

```
int a = 1;
int b = 100;
```

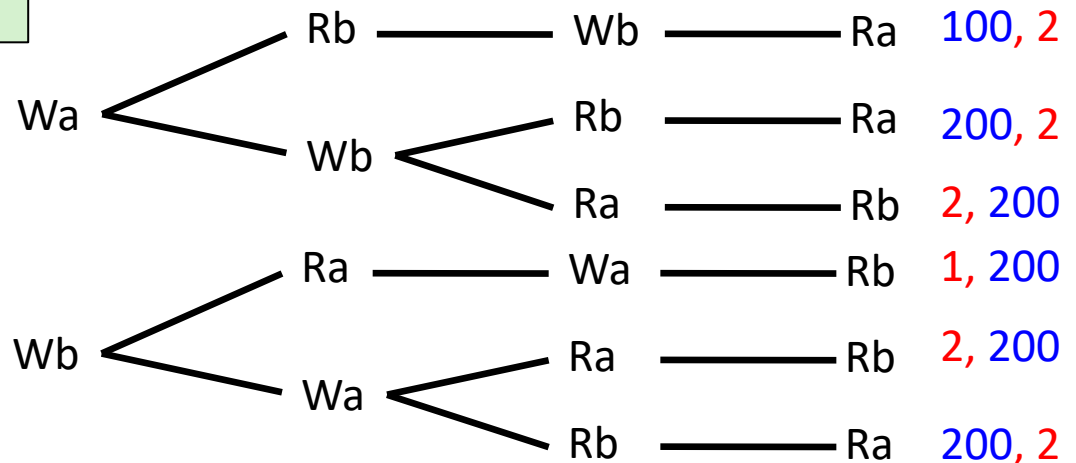
Thread1:  
Wa: a = 2;  
Rb: **print(b);**

Thread2:  
Wb: b = 200;  
Ra: **print(a);**

Thread consistency  
constraints

Wa ————— Rb

Wb ————— Ra

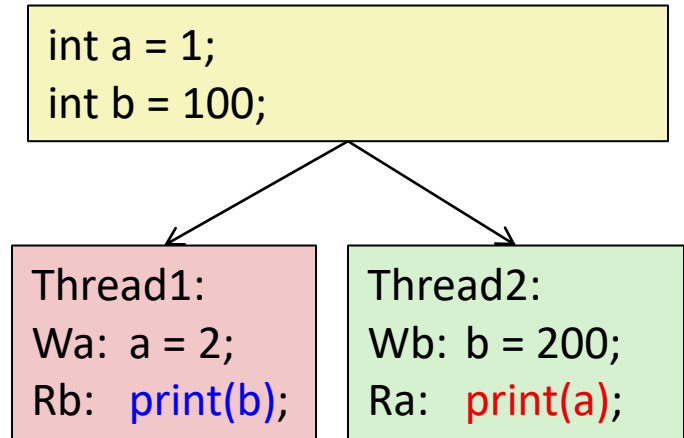
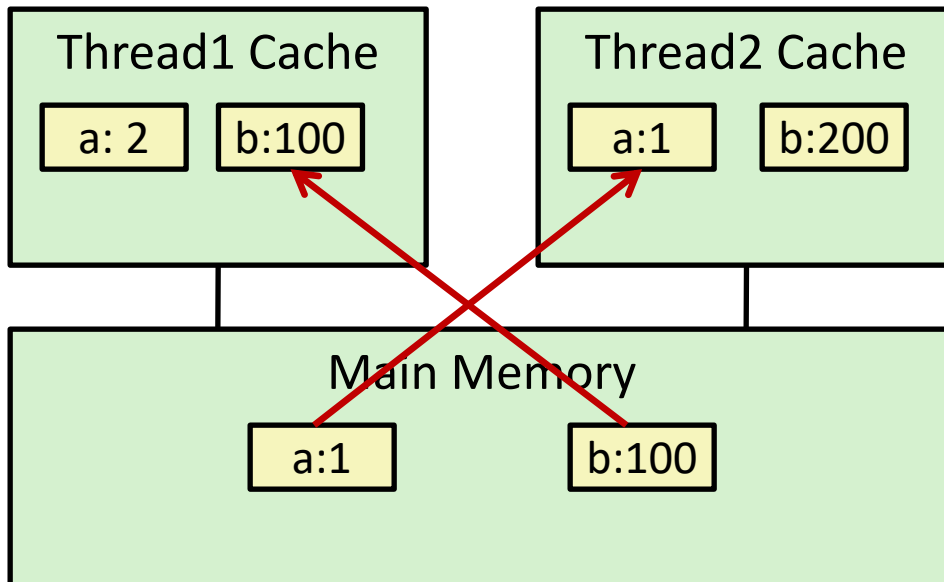


- Impossible outputs

- **100, 1** and **1, 100**
- Would require reaching both Ra and Rb before Wa and Wb

# Non-Coherent Cache Scenario

- Write-back caches, without coordination between them



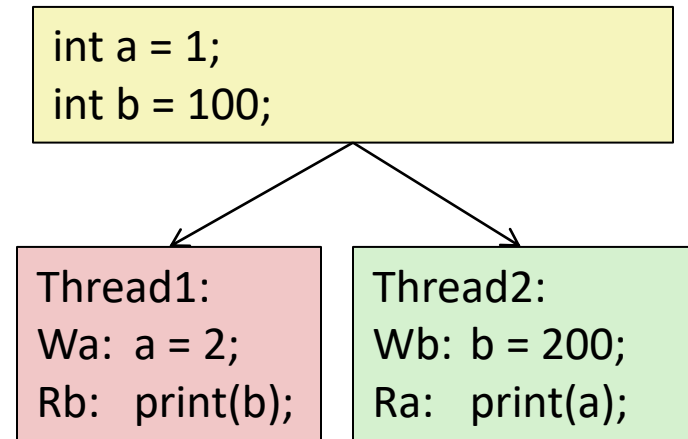
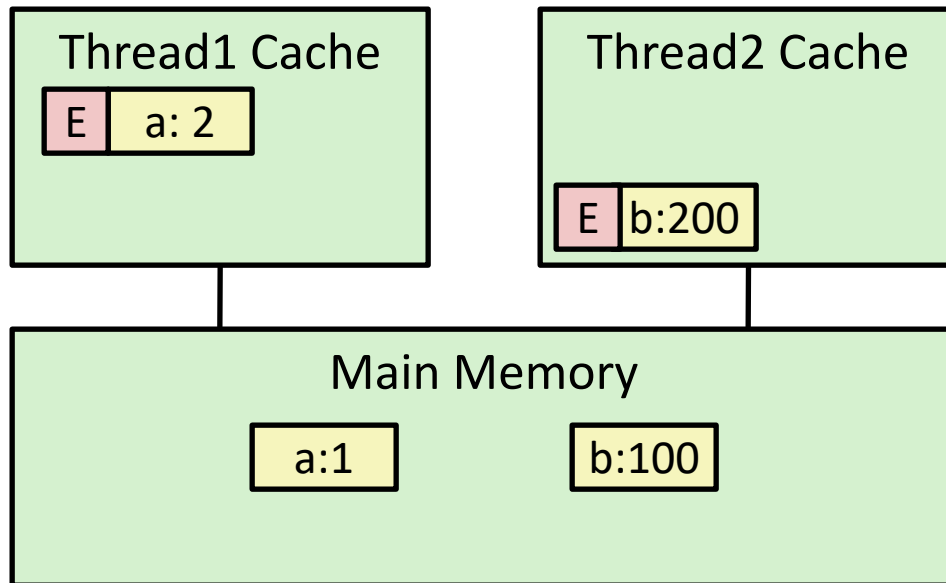
print 1

print 100

# Snoopy Caches

- Tag each cache block with state

Invalid	Cannot use value
Shared	Readable copy
Exclusive	Writeable copy

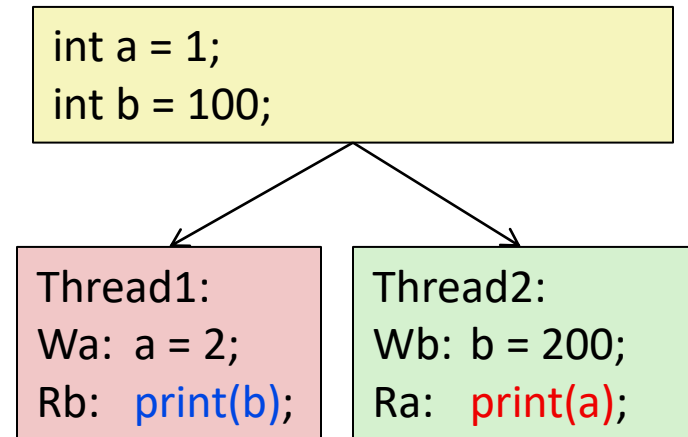
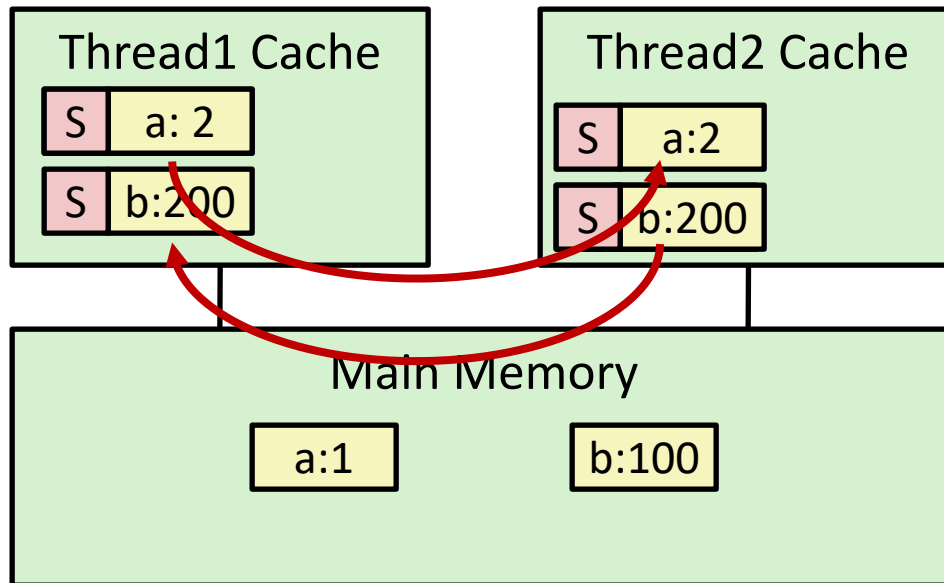




# Snoopy Caches

- Tag each cache block with state

Invalid	Cannot use value
Shared	Readable copy
Exclusive	Writeable copy



print 2

print 200

- When cache sees request for one of its E-tagged blocks
  - Supply value from cache
  - Set tag to S

# MESI Protocol for Snoopy Caches

**Modified:** The local processor has modified the cache line. This also implies it is the only copy in any cache.

**Exclusive:** The cache line is not modified but known to not be loaded into any other processor's cache.

**Shared:** The cache line is not modified and might exist in another processor's cache.

**Invalid:** The cache line is invalid, i.e., unused.

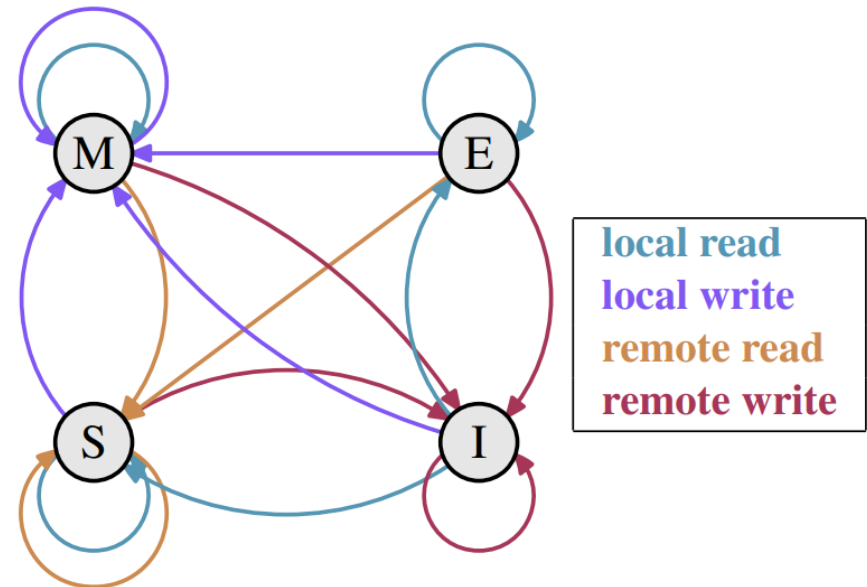


Figure 3.18: MESI Protocol Transitions