

Thread-Level Parallelism

15-213: Introduction to Computer Systems
26th Lecture, Nov. 30, 2010

Instructors:

Randy Bryant and Dave O'Hallaron

Today

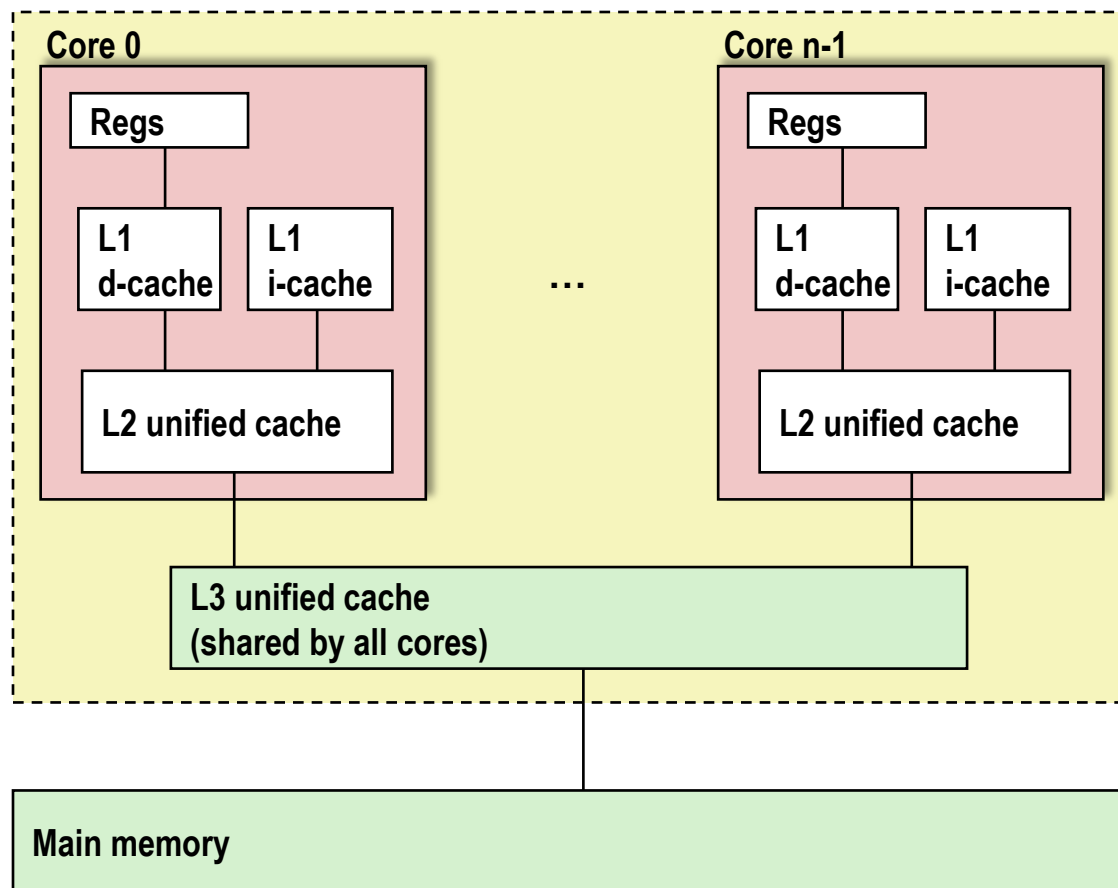
■ Parallel Computing Hardware

- Multicore
 - Multiple separate processors on single chip
- Hyperthreading
 - Replicated instruction execution hardware in each processor
- Maintaining cache consistency

■ Thread Level Parallelism

- Splitting program into independent tasks
 - Example: Parallel summation
 - Some performance artifacts

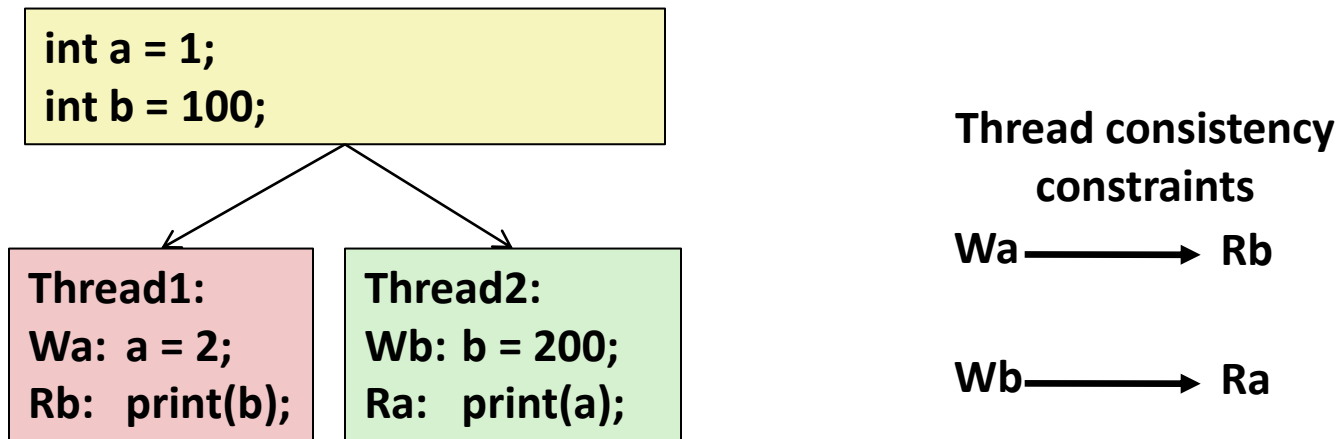
Multicore Processor



■ Intel Nehalem Processor

- E.g., Shark machines
- Multiple processors operating with coherent view of memory

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;
```

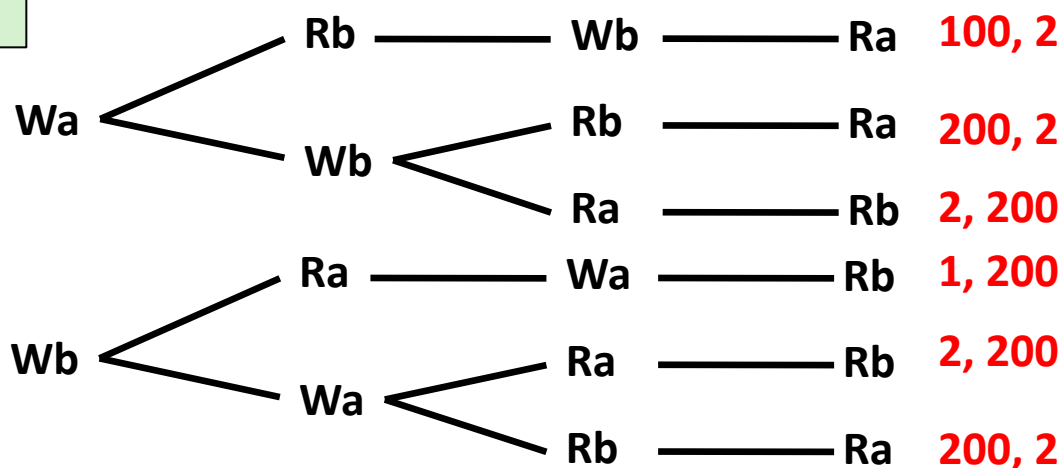
Thread1:
 Wa: a = 2;
 Rb: print(b);

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

Thread consistency
constraints

Wa ————— Rb

Wb ————— Ra

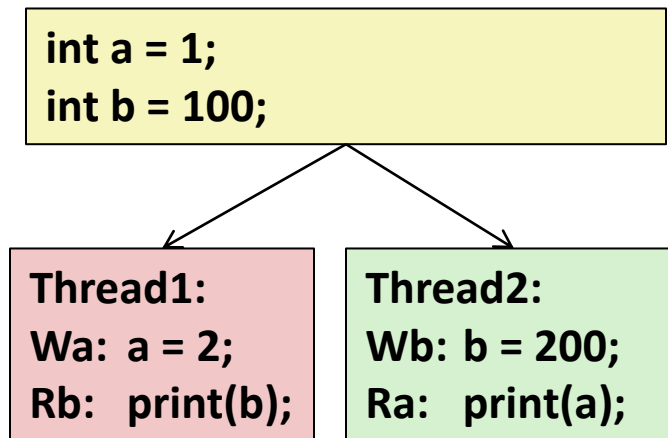
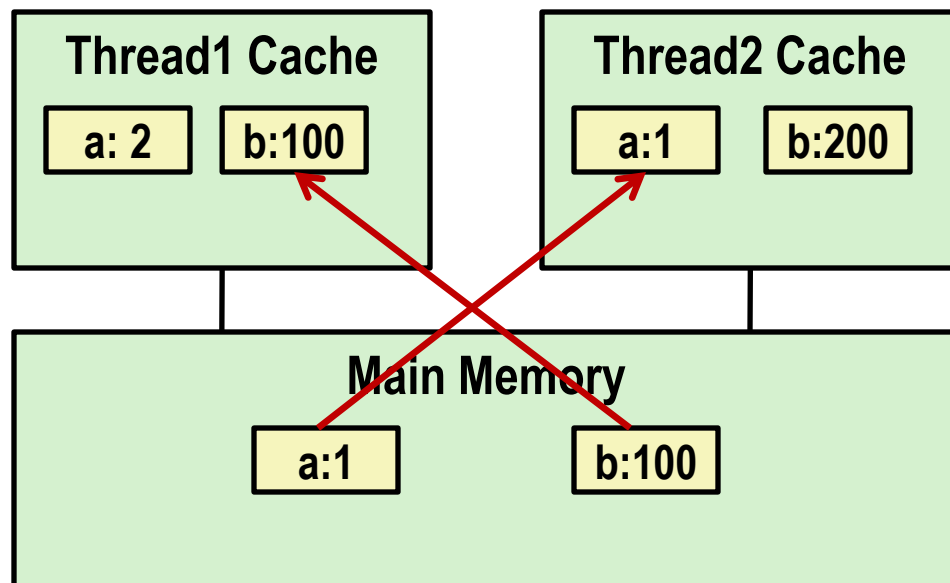


■ Impossible outputs

- **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



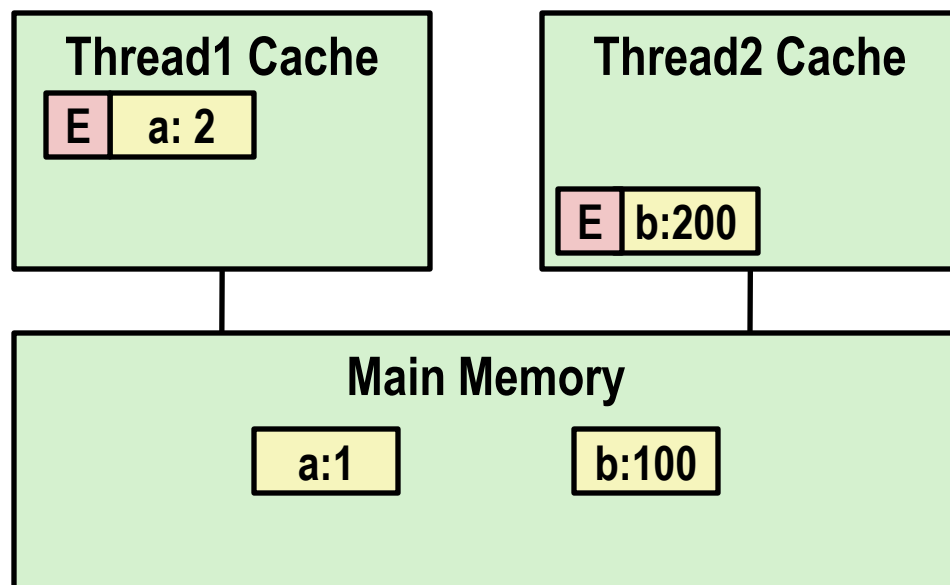
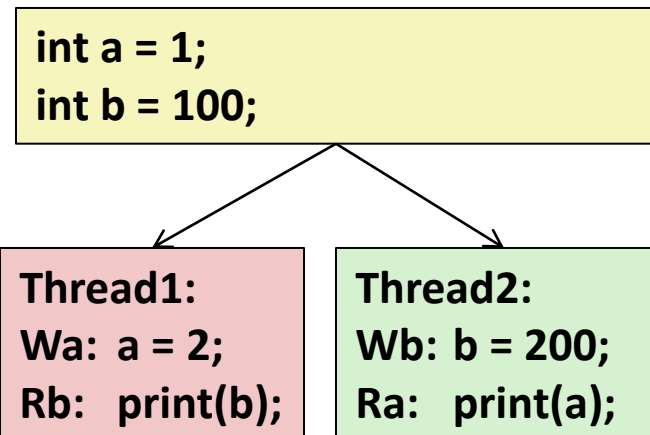
print 1

print 100

Snoopy Caches

■ Tag each cache block with state

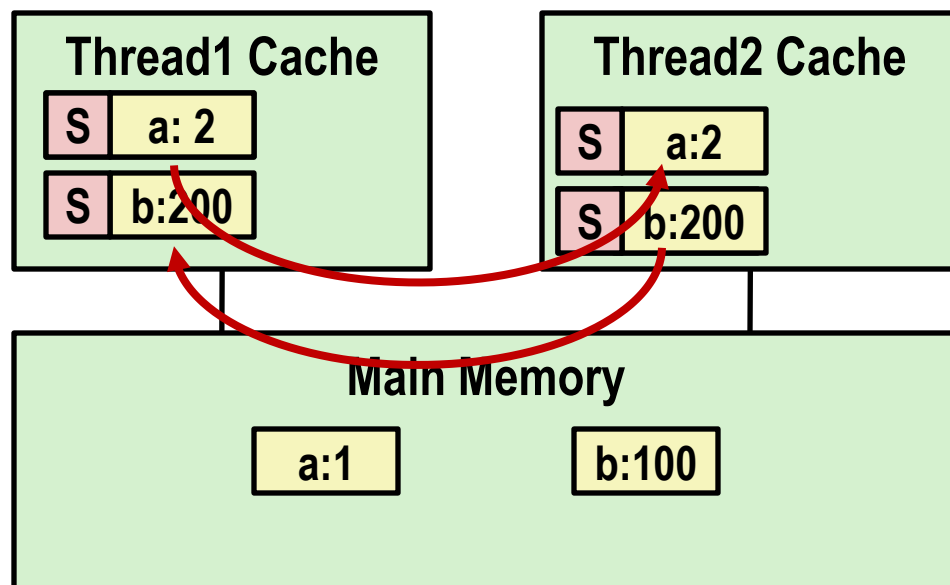
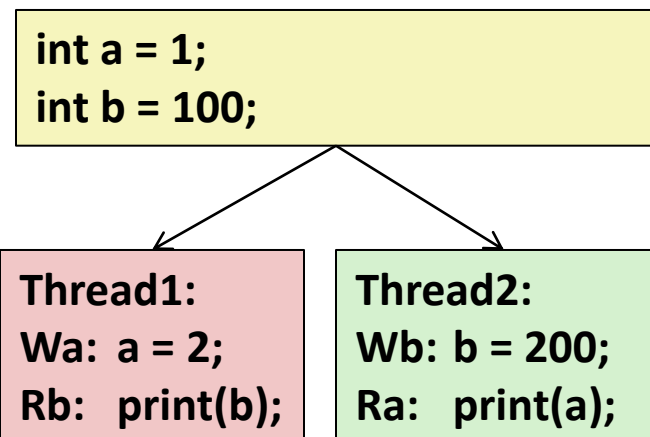
Invalid	Cannot use value
Shared	Readable copy
Exclusive	Writeable copy



Snoopy Caches

■ Tag each cache block with state

Invalid	Cannot use value
Shared	Readable copy
Exclusive	Writeable copy

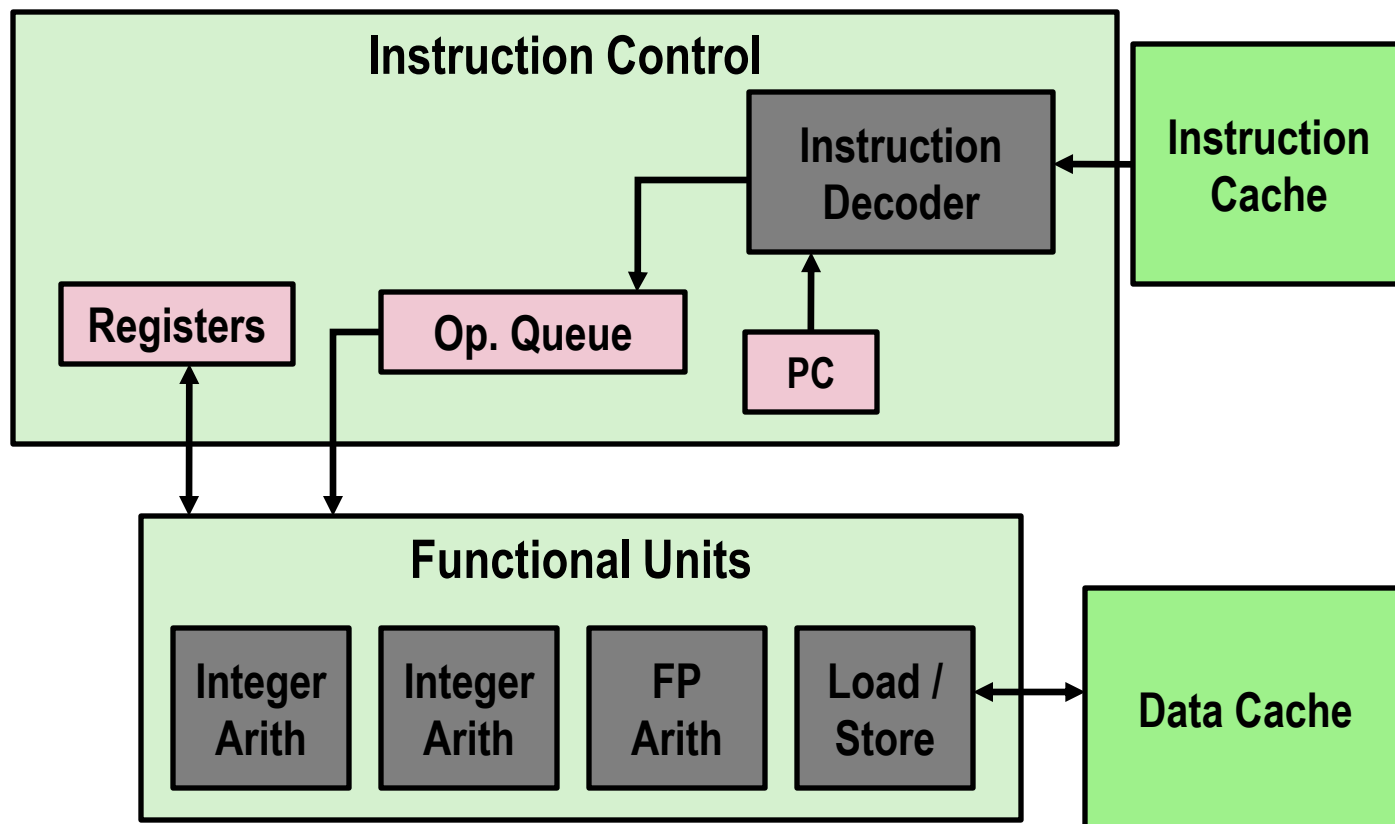


print 2

print 200

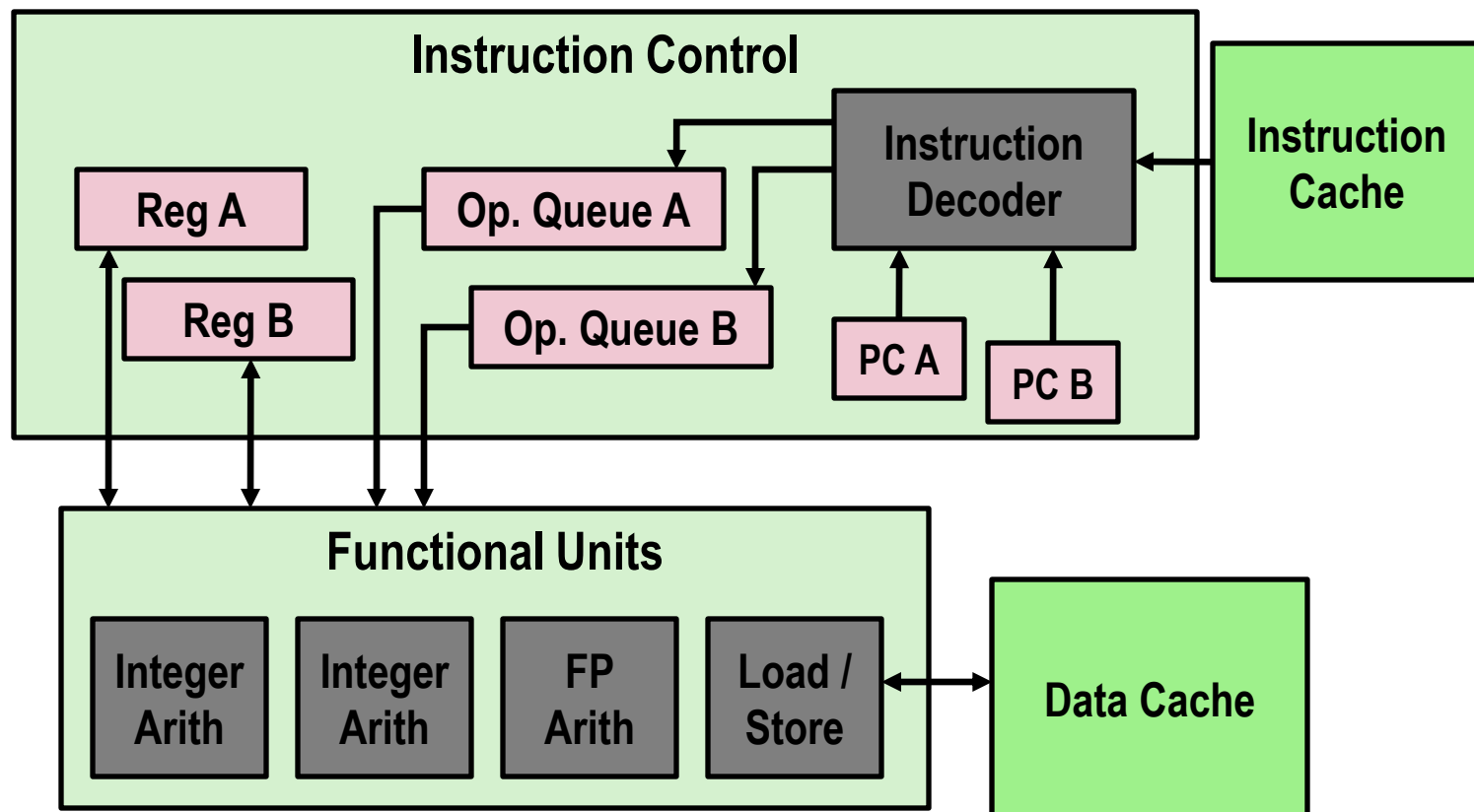
- When cache sees request for one of its E-tagged blocks
 - Supply value from cache
 - Set tag to S

Out-of-Order Processor Structure



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

Hyperthreading



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

Summary: Creating Parallel Machines

■ Multicore

- Separate instruction logic and functional units
- Some shared, some private caches
- Must implement cache coherency

■ Hyperthreading

- Also called “simultaneous multithreading”
- Separate program state
- Shared functional units & caches
- No special control needed for coherency

■ Combining

- Shark machines: 8 cores, each with 2-way hyperthreading
- Theoretical speedup of 16X
 - Never achieved in our benchmarks

Summation Example

- **Sum numbers 0, ..., N-1**
 - Should add up to $(N-1)*N/2$
- **Partition into K ranges**
 - $\lfloor N/K \rfloor$ values each
 - Accumulate leftover values serially
- **Method #1: All threads update single global variable**
 - 1A: No synchronization
 - 1B: Synchronize with pthread semaphore
 - 1C: Synchronize with pthread mutex
 - “Binary” semaphore. Only values 0 & 1

Accumulating in Single Global Variable: Declarations

```
typedef unsigned long data_t;
/* Single accumulator */
volatile data_t global_sum;

/* Mutex & semaphore for global sum */
sem_t semaphore;
pthread_mutex_t mutex;

/* Number of elements summed by each thread */
size_t nelems_per_thread;

/* Keep track of thread IDs */
pthread_t tid[MAXTHREADS];
/* Identify each thread */
int myid[MAXTHREADS];
```

Accumulating in Single Global Variable: Operation

```
nelems_per_thread = nelems / nthreads;
/* Set global value */
global_sum = 0;

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

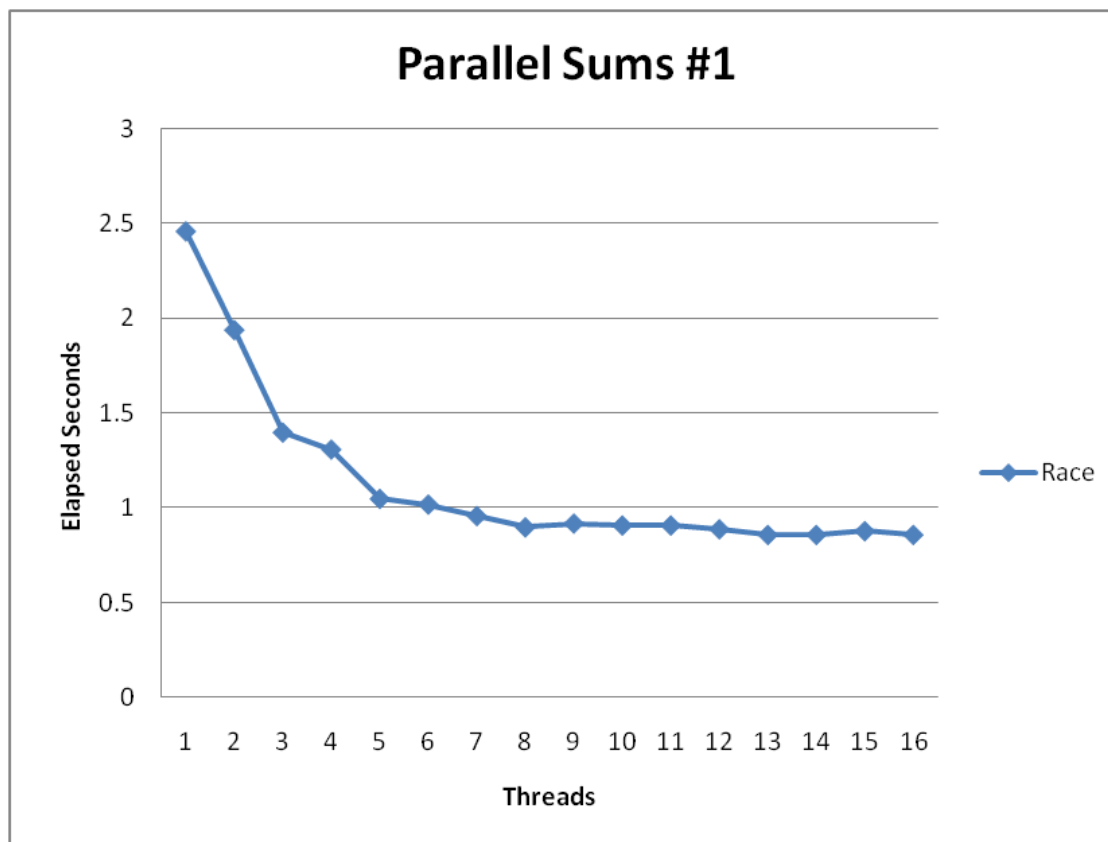
result = global_sum;
/* Add leftover elements */
for (e = nthreads * nelems_per_thread; e < nelems; e++)
    result += e;
```

Thread Function: No Synchronization

```
void *sum_race(void *vargp)
{
    int myid = *((int *)vargp);
    size_t start = myid * nelems_per_thread;
    size_t end = start + nelems_per_thread;
    size_t i;

    for (i = start; i < end; i++) {
        global_sum += i;
    }
    return NULL;
}
```

Unsynchronized Performance



- $N = 2^{30}$
- Best speedup = 2.86X
- Gets wrong answer when > 1 thread!

Thread Function: Semaphore / Mutex

Semaphore

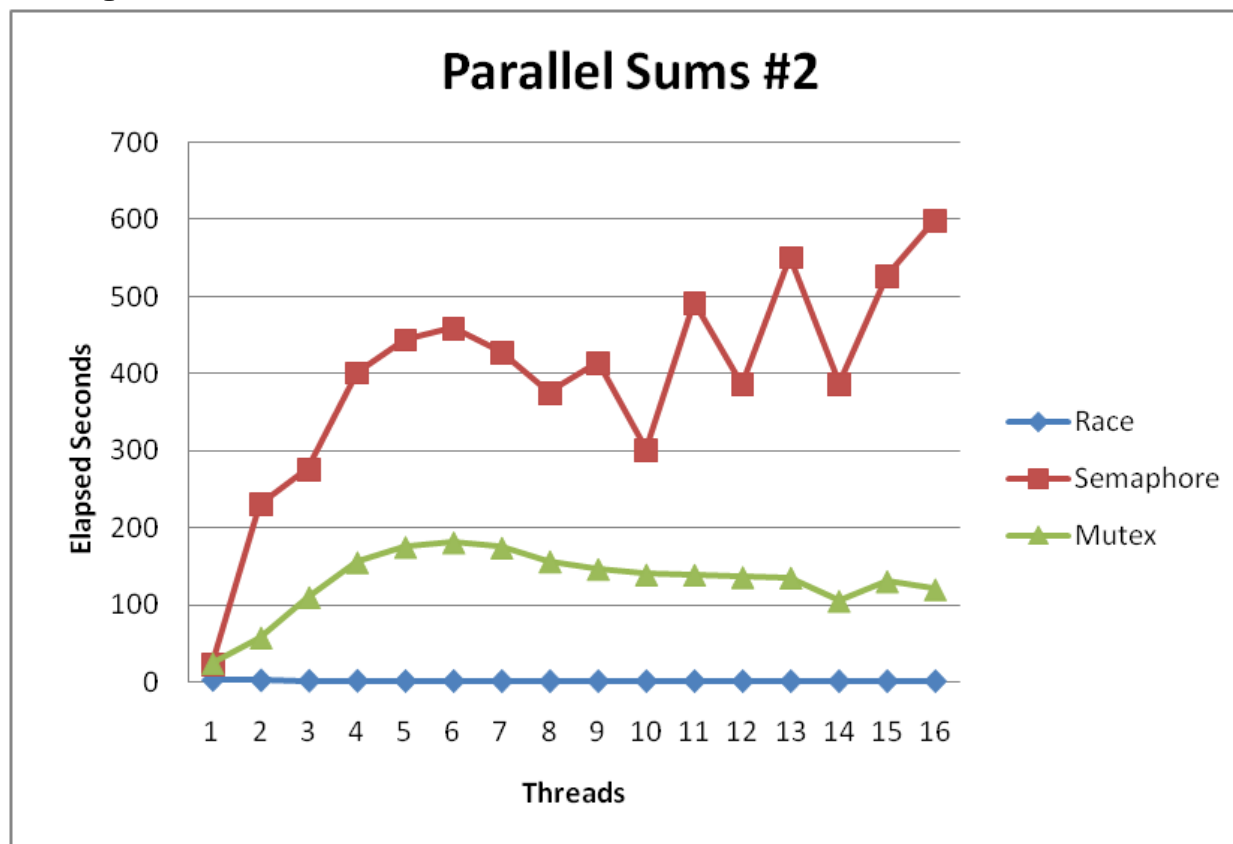
```
void *sum_sem(void *vargp)
{
    int myid = *((int *)vargp);
    size_t start = myid * nelems_per_thread;
    size_t end = start + nelems_per_thread;
    size_t i;

    for (i = start; i < end; i++) {
        sem_wait(&semaphore);
        global_sum += i;
        sem_post(&semaphore);
    }
    return NULL;
}
```

Mutex

```
pthread_mutex_lock(&mutex);
global_sum += i;
pthread_mutex_unlock(&mutex);
```

Semaphore / Mutex Performance



- **Terrible Performance**
 - 2.5 seconds → ~10 minutes
- **Mutex 3X faster than semaphore**
- **Clearly, neither is successful**

Separate Accumulation

- **Method #2: Each thread accumulates into separate variable**
 - 2A: Accumulate in contiguous array elements
 - 2B: Accumulate in spaced-apart array elements
 - 2C: Accumulate in registers

```
/* Partial sum computed by each thread */  
data_t psum[MAXTHREADS*MAXSPACING];  
/* Spacing between accumulators */  
size_t spacing = 1;
```

Separate Accumulation: Operation

```
nelems_per_thread = nelems / nthreads;

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

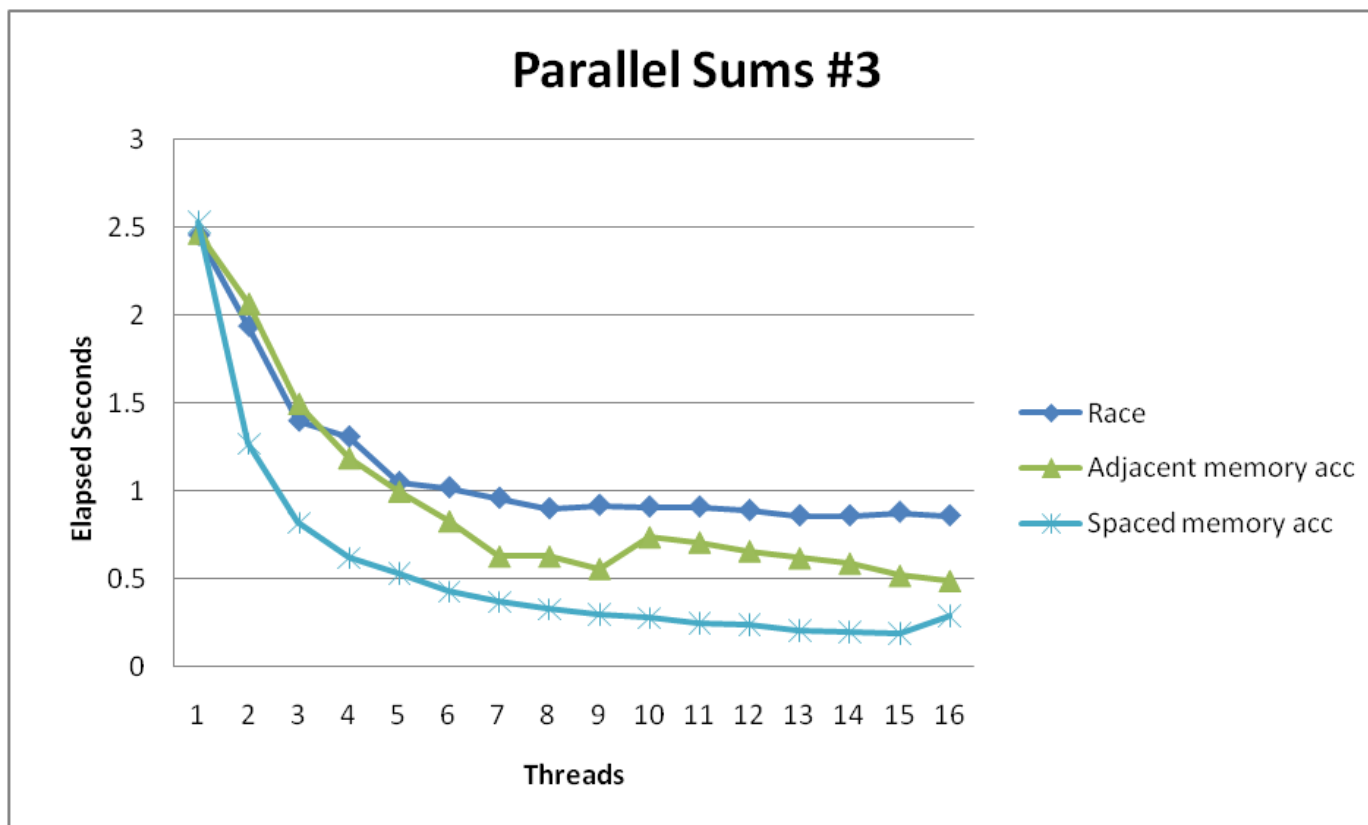
result = 0;
/* Add up the partial sums computed by each thread */
for (i = 0; i < nthreads; i++)
    result += psum[i*spacing];
/* Add leftover elements */
for (e = nthreads * nelems_per_thread; e < nelems; e++)
    result += e;
```

Thread Function: Memory Accumulation

```
void *sum_global(void *vargp)
{
    int myid = *((int *)vargp);
    size_t start = myid * nelems_per_thread;
    size_t end = start + nelems_per_thread;
    size_t i;

    size_t index = myid*spacing;
    psum[index] = 0;
    for (i = start; i < end; i++) {
        psum[index] += i;
    }
    return NULL;
}
```

Memory Accumulation Performance

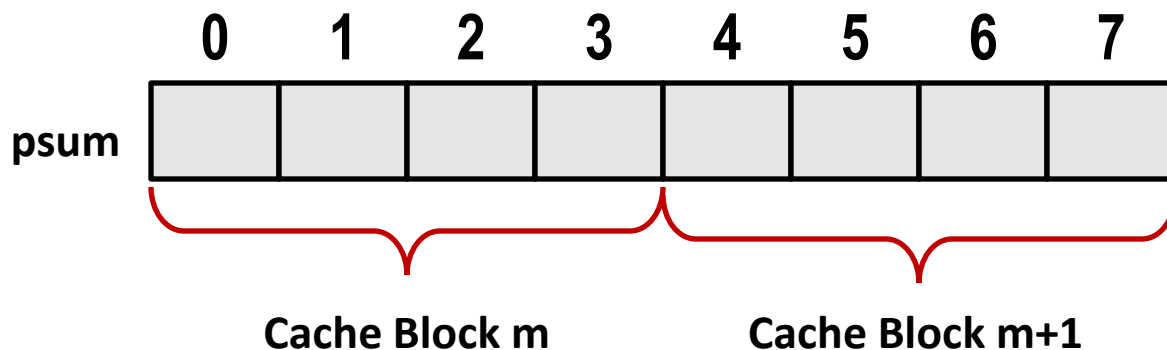


■ Clear threading advantage

- Adjacent speedup: 5 X
- Spaced-apart speedup: 13.3 X (Only observed speedup > 8)

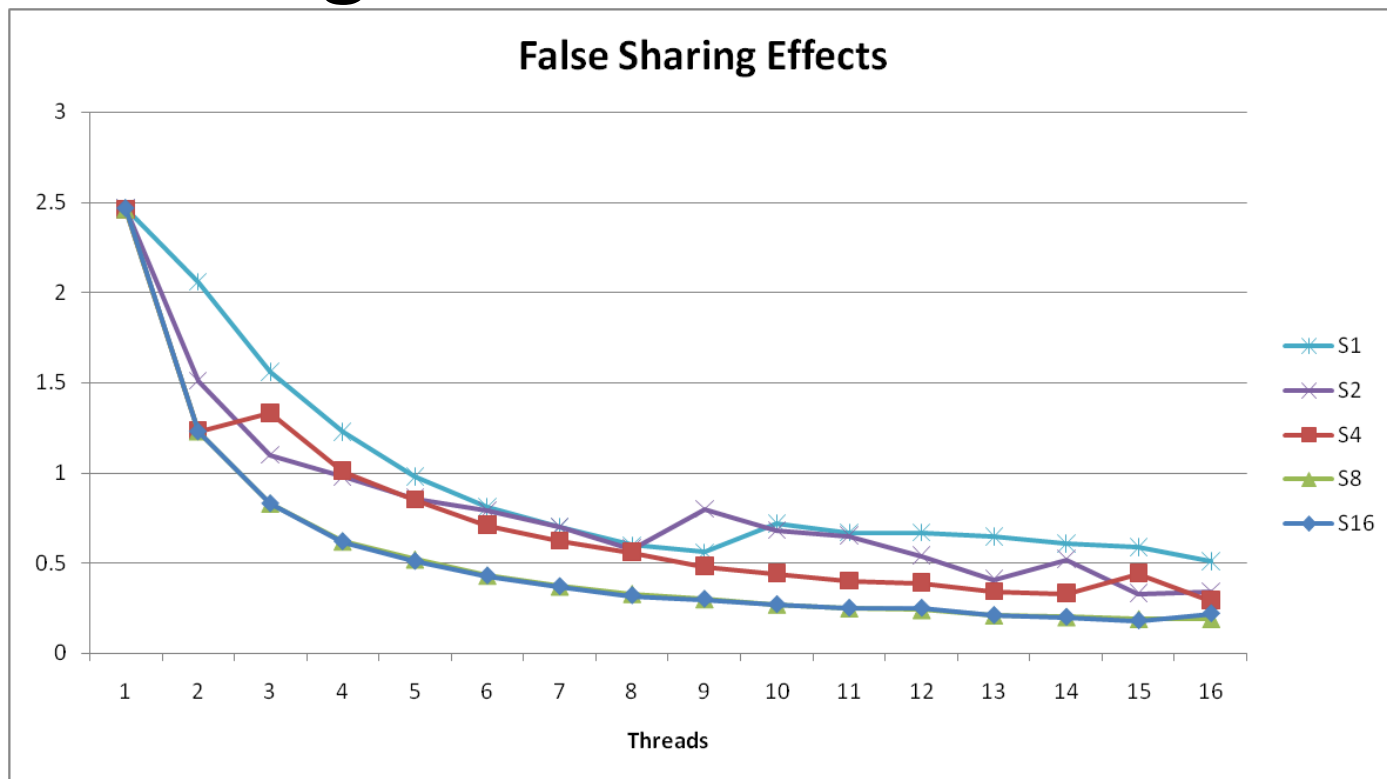
■ Why does spacing the accumulators apart matter?

False Sharing



- **Coherency maintained on cache blocks**
- **To update `psum[i]`, thread `i` must have exclusive access**
 - Threads sharing common cache block will keep fighting each other for access to block

False Sharing Performance

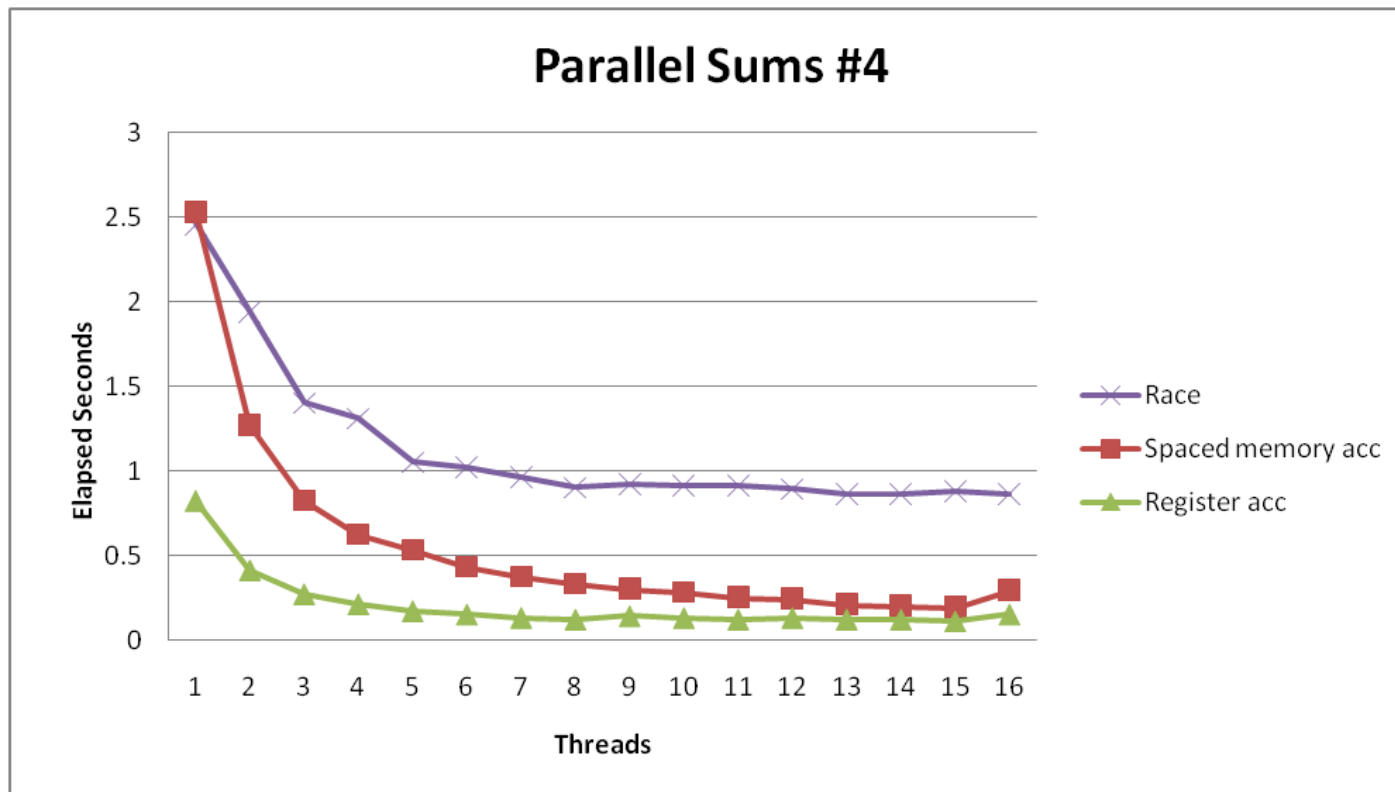


- Best spaced-apart performance 2.8 X better than best adjacent
- **Demonstrates cache block size = 64**
 - 8-byte values
 - No benefit increasing spacing beyond 8

Thread Function: Register Accumulation

```
void *sum_local(void *vargp)
{
    int myid = *((int *)vargp);
    size_t start = myid * nelems_per_thread;
    size_t end = start + nelems_per_thread;
    size_t i;
    size_t index = myid*spacing;
    data_t sum = 0;
    for (i = start; i < end; i++) {
        sum += i;
    }
    psum[index] = sum;    return NULL;
}
```

Register Accumulation Performance



- **Clear threading advantage**
 - Speedup = 7.5 X
- **2X better than fastest memory accumulation**