# Thread-Level Parallelism

15-213 / 18-213: Introduction to Computer Systems
26th Lecture, Apr. 26, 2012

**Instructors:**

Todd Mowry and Anthony Rowe

1

---

# Today

- **Parallel Computing Hardware**
  - Multicore
    - Multiple separate processors on single chip
  - Hyperthreading
    - Multiple threads executed on a given processor at once
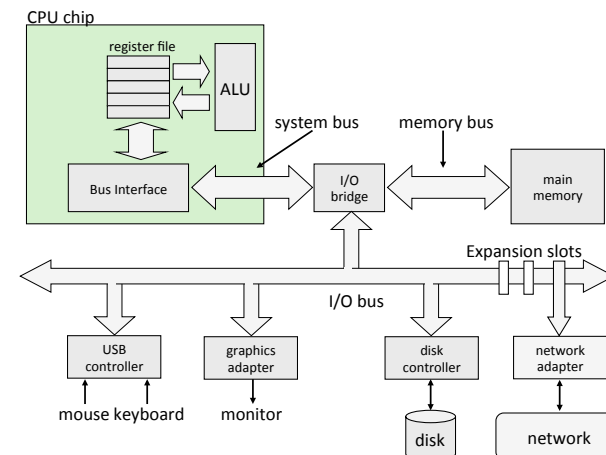- **Thread-Level Parallelism**
  - Splitting program into independent tasks
    - Example: Parallel summation
    - Some performance artifacts
  - Divide-and conquer parallelism
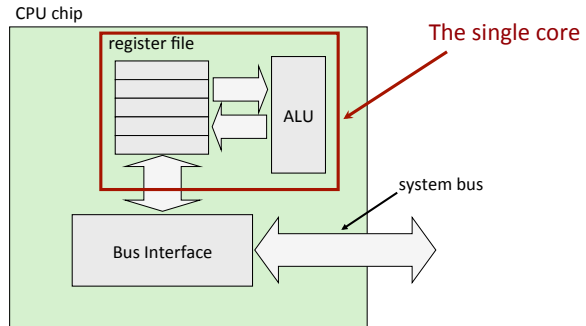    - Example: Parallel quicksort

2

---

# Why Multi-Core?

- Traditionally, single core performance is improved by increasing the clock frequency...
- ...and making deeply pipelined circuits...
- Which leads to...
  - Heat problems
  - Speed of light problems
  - Difficult design and verification
  - Large design teams
  - Big fans, heat sinks
  - Expensive air-conditioning on server farms
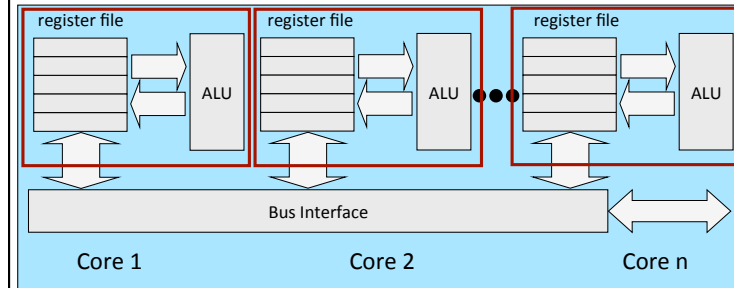- Increasing clock frequency no longer the way to go forward
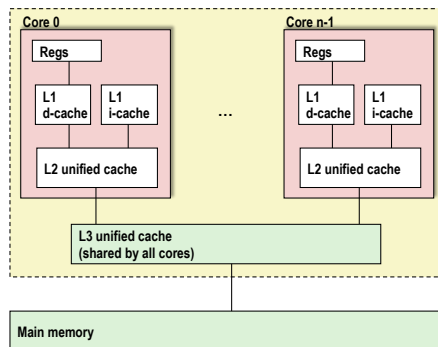
---

# Single Core Computer



---

## Single Core Processor (chip)

CPU chip

The single core

register file

ALU

system bus

Bus Interface

## Multi-Core Architecture

- Somewhat recent trend in computer architecture
- Replicate many cores on a single die

register file

ALU

register file

ALU

register file

ALU

Bus Interface

Core 1

Core 2

Core n

Multi-core Chip

## Multi-core Processor

Core 0

Regs

L1 d-cache

L1 i-cache

L2 unified cache

...

Core n-1

Regs

L1 d-cache

L1 i-cache

L2 unified cache

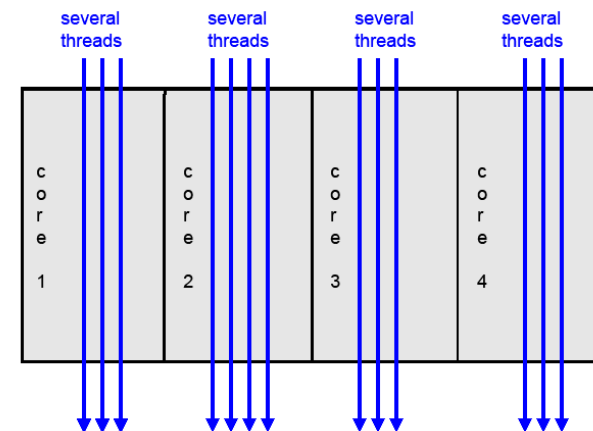L3 unified cache (shared by all cores)

Main memory

- **Intel Nehalem Processor**
  - E.g., Shark machines
  - Multiple processors operating with coherent view of memory

7

## Within each core, threads are time-sliced (just like on a uniprocessor)

several threads

several threads

several threads

several threads

core 1

core 2

core 3

core 4

**2**

## Interaction With the Operating System

- OS perceives each core as a separate processor
- OS scheduler maps threads/processes to different cores
- Most major OS support multi-core today:
  - Mac OS X, Linux, Windows, …

## Flavors of Parallelism

- Instruction Level Parallelism (ILP)
- Thread Level Parallelism (TLP)
- Simultaneous Multi-Threading (SMT)
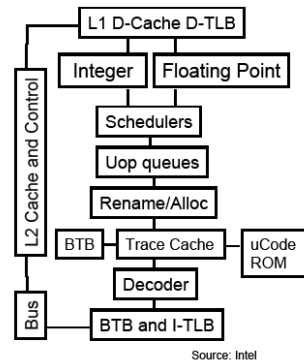
## Instruction-Level Parallelism

- Parallelism at the machine-instruction level
- Achieved in the processor with
  - Pipeline
  - Re-ordered instructions
  - Split into micro-instructions
  - Aggressive branch prediction
  - Speculative execution
- ILP enabled rapid increases in processor performance
  - Has since plateaued

## Thread-level Parallelism

- Parallelism on a coarser scale
- Server can serve each client in a separate thread
  - Web server, database server
- Computer game can do AI, graphics, physics, UI in four different threads
- Single-core superscalar processors cannot fully exploit TLP
  - Thread instructions are interleaved on a coarse level with other threads
- Multi-core architectures are the next step in processor evolution: explicitly exploiting TLP
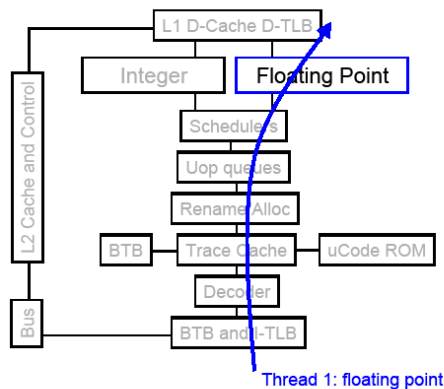
## Simultaneous Multithreading (SMT)

- Complimentary technique to multi-core
- Addresses the stalled pipeline problem
  - Pipeline is stalled waiting for the result of a long operation (float?)
  - ... or waiting for data to arrive from memory (long latency)
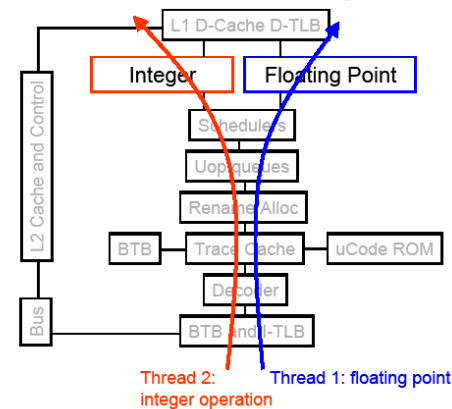- Other execution units are idle



Source: Intel

## SMT

- Permits multiple independent threads to execute SIMULTANEOUSLY on the SAME core
- Weaving together multiple "threads"
- Example: if one thread is waiting for a floating point operation to complete, another thread can use the integer units

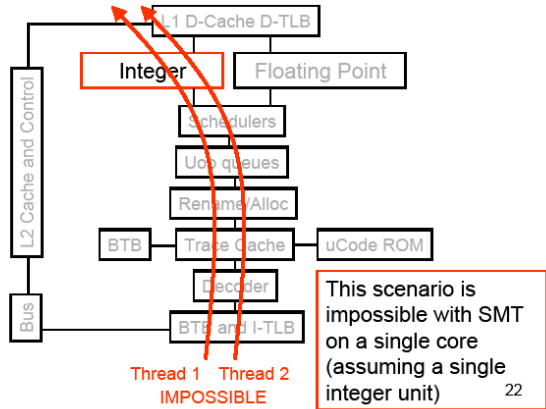# Without SMT, only a single thread can run at any given time



Thread 1: floating point

19

# SMT processor: both threads can run concurrently



Thread 2: integer operation      Thread 1: floating point

21

4

# But: Can't simultaneously use the same functional unit



L1 D-Cache D-TLB

Integer    Floating Point

L2 Cache and Control

Schedulers

Uop queues

Rename/Alloc

BTB    Trace Cache    uCode ROM

Decoder

Bus    BTB and I-TLB

Thread 1    Thread 2
IMPOSSIBLE

This scenario is impossible with SMT on a single core (assuming a single integer unit)    22
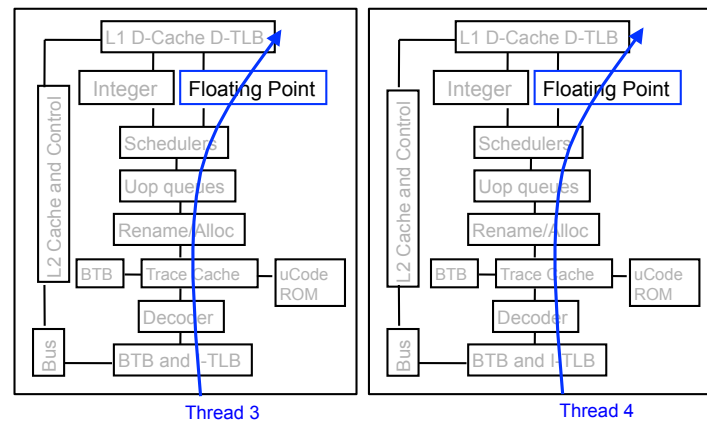
---

## SMT is not a "true" parallel processor

- Enables better threading (e.g. up to 30%)
- OS and applications perceive each simultaneous thread as a separate "virtual processor"
- The chip has only a single copy of each resource
- Compare to multi-core:
  - Each core has its own copy of resources
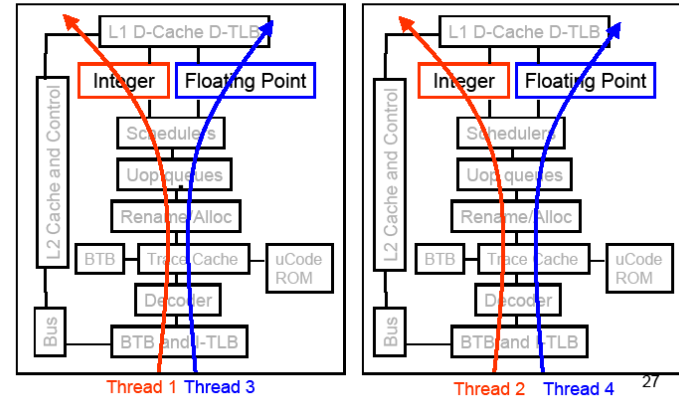
---

## Multi-core: Threads run on separate cores



L1 D-Cache D-TLB

Integer    Floating Point

L2 Cache and Control

Schedulers

Uop queues

Rename/Alloc

BTB    Trace Cache    uCode ROM

Decoder

Bus    BTB and I-TLB

Thread 1

L1 D-Cache D-TLB

Integer    Floating Point

L2 Cache and Control

Schedulers

Uop queues

Rename/Alloc

BTB    Trace Cache    uCode ROM

Decoder

Bus    BTB and I-TLB

Thread 2

---

## Multi-core: Threads run on separate cores



L1 D-Cache D-TLB

Integer    Floating Point

L2 Cache and Control

Schedulers

Uop queues

Rename/Alloc

BTB    Trace Cache    uCode ROM

Decoder

Bus    BTB and I-TLB

Thread 3

L1 D-Cache D-TLB

Integer    Floating Point

L2 Cache and Control

Schedulers

Uop queues

Rename/Alloc

BTB    Trace Cache    uCode ROM

Decoder

Bus    BTB and I-TLB

Thread 4

**5**

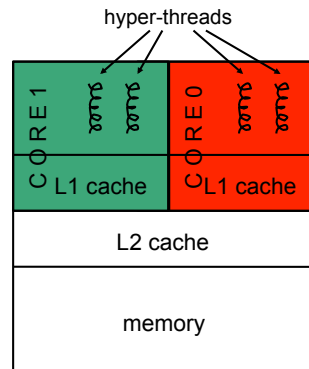## Combining Multi-core and SMT

- Cores can be SMT-enabled (or not)
- The different combinations:
  - Single-core, non-SMT: standard uniprocessor
  - Single-core, with SMT
  - Multi-core, non-SMT
  - Multi-core, with SMT: our fish machines
- The number of SMT threads is determined by hardware design
  - 2, 4 or sometimes 8 simultaneous threads
- Intel calls them "Hyper-threads"

---

## SMT Dual-core: all four threads can run concurrently



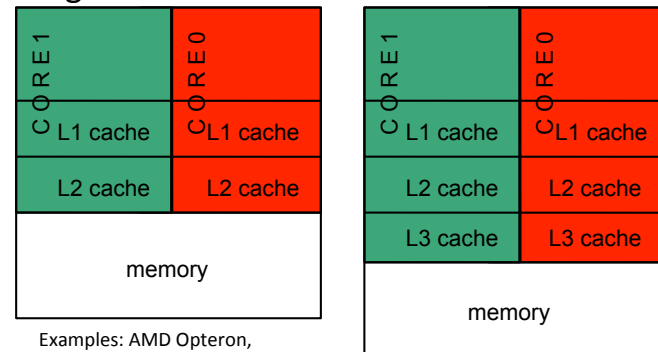Thread 1   Thread 3                    Thread 2   Thread 4    27

---

## SMT/Multi-Core and the Memory Hierarchy

- SMT is a sharing of pipeline resources
  - Thus all caches are shared
- Multi-core chips:
  - L1 caches are private (i.e. each core has its own L1)
  - L2 cache private in some architectures, shared in others
  - Main memory is always shared
- Example: Fish machines
  - Dual-core Intel Xeon processors
  - Each core is hyper-threaded
  - Private L1, shared L2 caches



hyper-threads

---

## Designs with Private L2 Caches



Examples: AMD Opteron, AMD Athlon, Intel Pentium D

Example: Intel Itanium 2

Quad Core 2 Duo shares L2 in pairs of cores

**6**

## Private vs Shared Cache

- Advantages of Private Cache
  - Closer to the core, so faster access
  - No contention for core access -- no waiting while another core accesses
- Advantages of Shared Cache
  - Threads on different cores can share same cache data
  - More cache space is available if a single (or a few) high-performance threads run
- Cache Coherence Problem
  - The same memory value can be stored in multiple private caches
  - Need to keep the data consistent across the caches
  - Many solutions exist
    - Invalidation protocol with bus snooping, ...

## Exploiting parallel execution

- **So far, we've used threads to deal with I/O delays**
  - e.g., one thread per client to prevent one from delaying another
- **Multi-core CPUs offer another opportunity**
  - Spread work over threads executing in parallel on N cores
  - Happens automatically, if many independent tasks
    - e.g., running many applications or serving many clients
  - Can also write code to make one big task go faster
    - by organizing it as multiple parallel sub-tasks
- **Shark machines can execute 16 threads at once**
  - 8 cores, each with 2-way hyperthreading
  - Theoretical speedup of 16X
    - never achieved in our benchmarks

26

## Summation Example

- **Sum numbers 0, ..., N-1**
  - Should add up to (N-1)*N/2
- **Partition into K ranges**
  - ⌊N/K⌋ 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

27

## 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];
```

28

7

## 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;
```

29

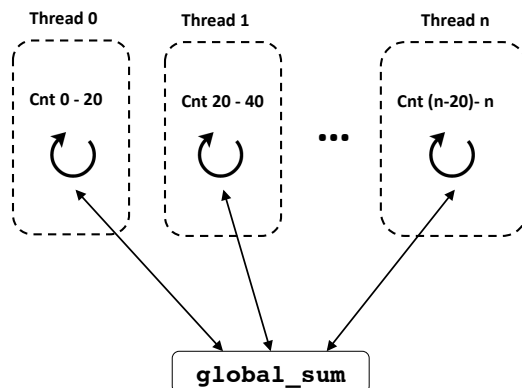## Thread Function: No Synchronization

```
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;
}
```
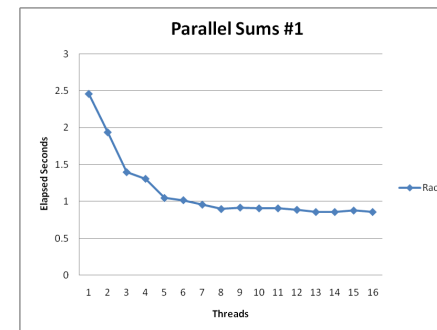
30

## Accumulating in Single Global Variable: Illustration



31

## Unsynchronized Performance



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

32

**8**

## 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;
}
```
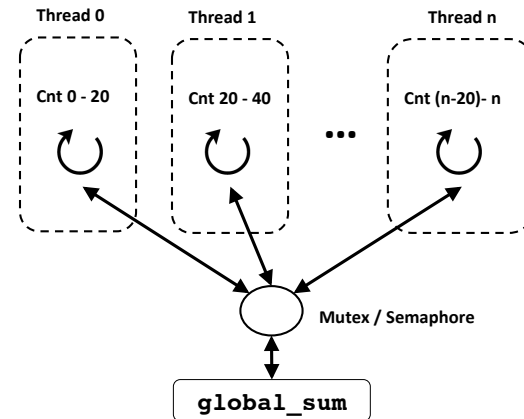
**Mutex**

```
pthread_mutex_lock(&mutex);
global_sum += i;
pthread_mutex_unlock(&mutex);
```
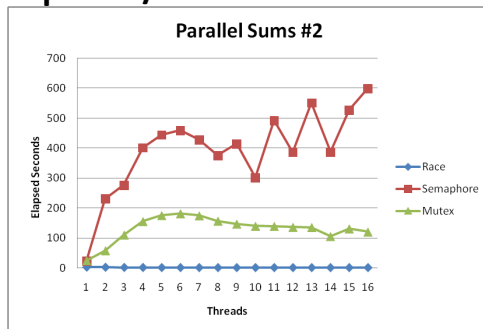
33

## Accumulating with Mutex / Semaphore



34

## Semaphore / Mutex Performance



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

35

## 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;
```

36

9

## 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;
```

37

## Thread Function: Memory Accumulation
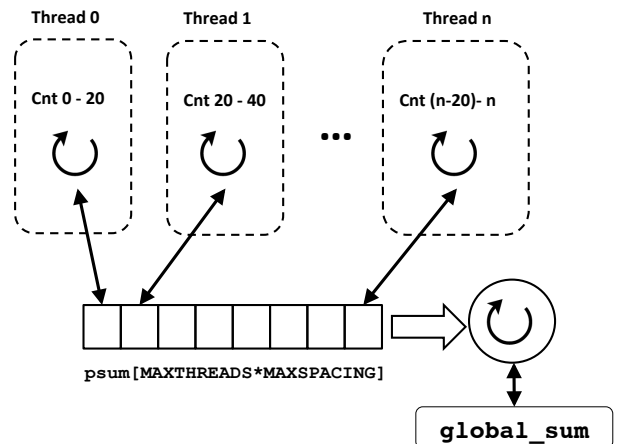
```
void *sum_global(void *vargp)
{
    int myid = *((int *)vargp);
    size_t start = myid * nelems_per_thread;
    size_t end = start + nelems_per_thread;
    size_t i;

    size_t index = myid*spacing;
    psum[index] = 0;
    for (i = start; i < end; i++) {
        psum[index] += i;
    }
    return NULL;
}
```
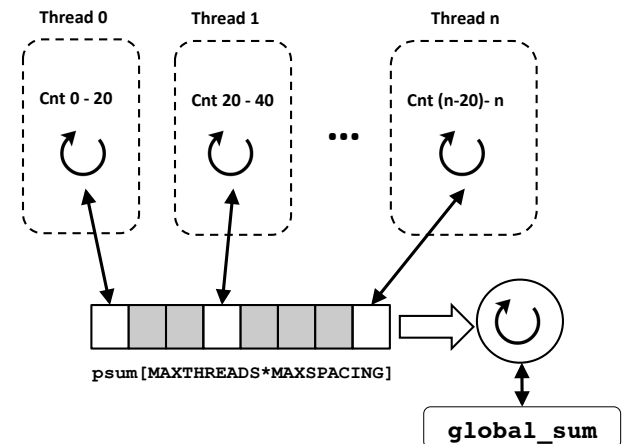
38

## Accumulating into memory (no spacing)



**Thread 0** Cnt 0 - 20    **Thread 1** Cnt 20 - 40    ...    **Thread n** Cnt (n-20)- n

`psum[MAXTHREADS*MAXSPACING]`

`global_sum`

39

## Accumulating into memory (spacing)



**Thread 0** Cnt 0 - 20    **Thread 1** Cnt 20 - 40    ...    **Thread n** Cnt (n-20)- n

`psum[MAXTHREADS*MAXSPACING]`

`global_sum`

40

10

## Memory Accumulation Performance

**Parallel Sums #3**



- **Clear threading advantage**
  - Adjacent speedup: 5 X
  - Spaced-apart speedup: 13.3 X (Only observed speedup > 8)
- **Why does spacing the accumulators apart matter?**

41

## 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

42

## False Sharing Performance

**False Sharing Effects**



  - 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

43

## 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;
}
```
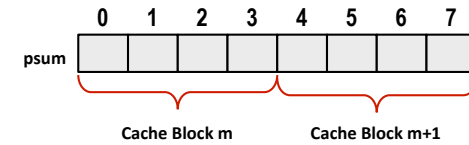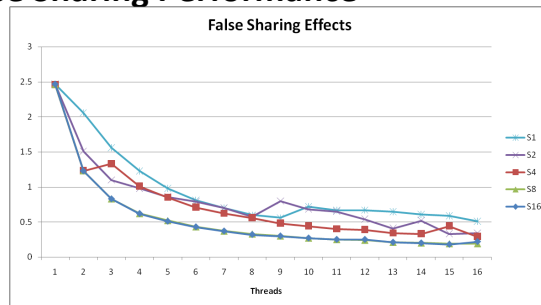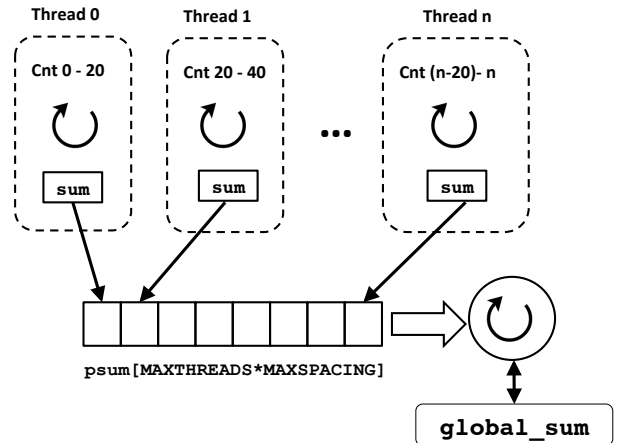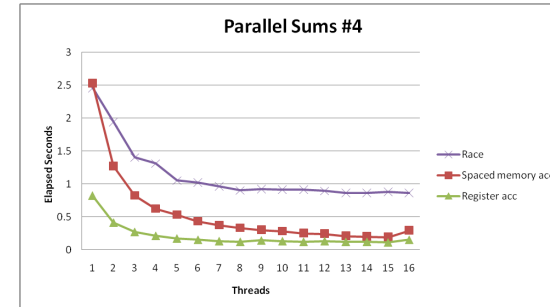
44

11

## Accumulating into register

**Thread 0**  **Thread 1**  **Thread n**

Cnt 0 - 20    Cnt 20 - 40    Cnt (n-20)- n

↻    ↻    ...    ↻

`sum`    `sum`    `sum`

`psum[MAXTHREADS*MAXSPACING]`

↻

`global_sum`

45

## Register Accumulation Performance

**Parallel Sums #4**

Elapsed Seconds — Threads

— Race
— Spaced memory acc
— Register acc

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

46

## Amdahl's Law

- **Overall problem**
  - T    Total 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
  - Maximum possible speedup
    - $k = \infty$
    - $T_\infty = (1-p)T$

47

## 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**
  - $T_9 = 0.9 * 10/9 + 0.1 * 10 = 1.0 + 1.0 = 2.0$
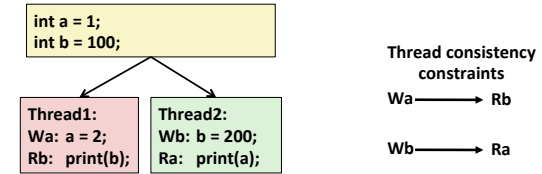  - Maximum possible speedup
    - $T_\infty = 0.1 * 10.0 = 1.0$

48

## Memory Consistency

- **There are different memory consistency models**
  - Abstract model of how hardware handles concurrent accesses
- **Most systems provide "sequential consistency"**
  - Overall effect consistent with each individual thread
  - But, the threads can be interleaved in any way
    - like when one-thread-at-a-time, but with constant interleaving
- **So, no correctness effects**
  - But, there can be performance effects
    - related to keeping cached values consistent
    - copying data from one cache to another is sorta like a cache miss
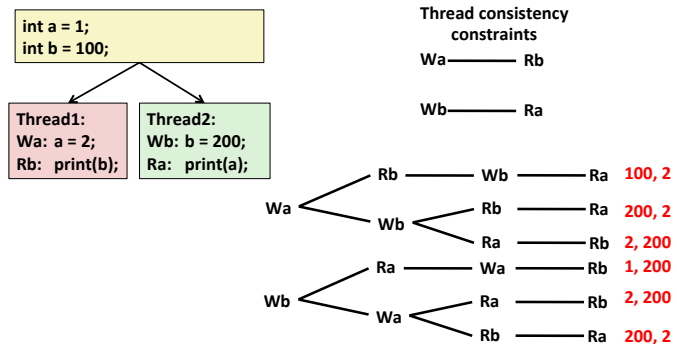
## Memory Consistency

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

| Thread1: | Thread2: |
| Wa:  a = 2; | Wb: b = 200; |
| Rb:  print(b); | Ra:  print(a); |

**Thread consistency constraints**

Wa ———→ Rb

Wb ———→ Ra

- **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: | Thread2: |
| Wa: a = 2; | Wb: b = 200; |
| Rb:  print(b); | Ra:  print(a); |

**Thread consistency constraints**

Wa ——— Rb

Wb ——— Ra

```
         ┌─ Rb ——— Wb ——— Ra   100, 2
      Wa ┤        ┌─ Rb ——— Ra  200, 2
         └─ Wb ───┤
                  └─ Ra ——— Rb  2, 200
         ┌─ Ra ——— Wa ——— Rb    1, 200
      Wb ┤        ┌─ Ra ——— Rb  2, 200
         └─ Wa ───┤
                  └─ Rb ——— Ra  200, 2
```
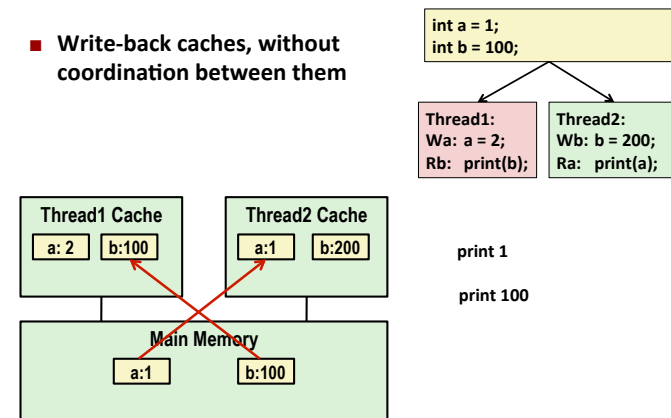
- **Impossible outputs**
  - 100, 1 and 1, 100
  - Would require reaching both Ra and Rb before Wa and Wb

## Non-Coherent Cache Scenario

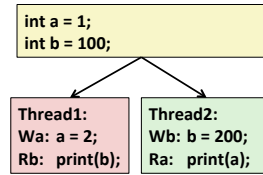- **Write-back caches, without coordination between them**
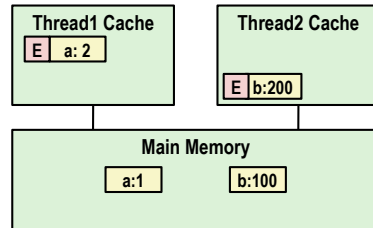
```
int a = 1;
int b = 100;
```

| Thread1: | Thread2: |
| Wa: a = 2; | Wb: b = 200; |
| Rb:  print(b); | Ra:  print(a); |

**Thread1 Cache**  a: 2  b:100

**Thread2 Cache**  a:1  b:200

print 1

print 100

**Main Memory**  a:1  b:100

## Snoopy Caches

- **Tag each cache block with state**

  Invalid     Cannot use value

  Shared     Readable copy

  Exclusive     Writeable copy

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

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

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

**Thread1 Cache**

| E | a: 2 |

**Thread2 Cache**

| E | b:200 |

**Main Memory**

| a:1 | | b:100 |

53

---

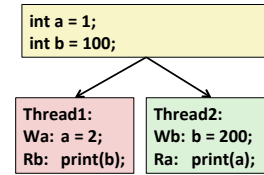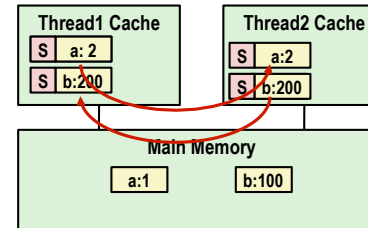## Snoopy Caches

- **Tag each cache block with state**

  Invalid     Cannot use value

  Shared     Readable copy

  Exclusive     Writeable copy

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

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

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

**Thread1 Cache**

| S | a: 2 |
| S | b:200 |

**Thread2 Cache**

| S | a:2 |
| S | b:200 |

print 2

print 200

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

**Main Memory**

| a:1 | | b:100 |

54

---

## 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

55

14