



Multicore Programming

From Threads to Transactional Memory



Yuan Lin

Sun Microsystems

Ali-Reza Adl-Tabatabai

Intel

Christos Kozyrakis

Stanford University

Bratin Saha

Intel



HotChips-18 Tutorial, August 20th, 2006

Tutorial Motivation & Goals



- Motivation

- All future processors will be multi-core chips
- How do we write efficient code for >1 core



- Goals

1. Introduction to multithreaded programming
 - The common practice in industry
 - Models, tools, and challenges
2. Introduction to transactional memory
 - A research technology for easier parallel programming
 - Overview, use examples, and implementation



Presenters



- **Yuan Lin**, Sun Microsystems (Ph.D. UIUC)
 - Senior Staff Engineer, Scalable Systems Group
 - Compilers and tools for multithreaded applications



- **Ali-Reza Adl-Tabatabai**, Intel (Ph.D. CMU)
 - Principal Engineer, Programming Systems Lab
 - Compilers & runtimes for future Intel architectures



- **Christos Kozyrakis**, Stanford University (Ph.D. UC Berkeley)
 - Assistant Professor, Electrical Eng. & Computer Science
 - Architectures & programming models for transactional memory
- **Bratin Saha**, Intel (Ph.D. Yale)
 - Senior Staff Researcher, Programming Systems Lab
 - Design of highly scalable runtimes for multi-core processors

Agenda



- Multithreaded Programming

BREAK

- Transactional Memory (TM)
 - TM Introduction
 - TM Implementation Overview
 - Hardware TM Techniques
 - Software TM Techniques



Q&A

- Save your questions for the end of each topic
- Additional Q&A at the end of the tutorial

Extended bibliography attached at the end of the handout

Multithreaded Programming

Challenges, current practice, and languages/tools support

Yuan Lin

Systems Group

Sun Microsystems

Multi-core Architectures

- What?
 - > A chip contains some number of "cores", each of which can run some number of strands ("hardware threads").
 - > Different cores or chips may share caches at different levels. All caches are coherent.
 - > All strands have access to the same shared memory.
- Why?
 - > More and more difficult to get good ROI by improving single-core performance only.
 - > Design complexity, power consumption, heat production, ...
 - > Throughput is as important as speed.

Multi-core Architectures

- 2 strand/core, 1 core/chip
 - > Intel Pentium 4 with Hyperthreading
- 1 strand/core, 2 cores/chip
 - > Intel Dual Core Xeon, AMD Dual Core Opteron, Sun UltraSPARC IV
- 2 strand/core, 2 cores/chip
 - > IBM Power 5
- 4 strand/core, 8 cores/chip
 - > Sun UltraSPARC T1



Example: UltraSPARC T1

- Eight cores (individual execution pipelines) per chip.
- Four strands share a pipeline in each core.
- Different strands are scheduled on the pipeline in round-robin order.
- The slot of a stalled strand is given to the next strand.
- Four strands on a core share L1 cache, all strands share L2 cache.
- A single FPU is shared by all cores.
- OS sees the each strand as a processor. OS schedules LWPs on strands. HW schedules strands in the core.

Class of Applications

- Multi-process applications
 - > e.g. Oracle database, SAP, PeopleSoft, ...
- Multi-threaded applications
 - > e.g. Siebel CRM, Sybase engine, ...
- Single-threaded applications
 - > Cannot be directly benefitted from multi-core architectures
 - > Solution: make it multi-threaded or multi-process

Multithreaded Programming

- Are developers ready now?
 - > Language developers
 - > Compiler / tool developers
 - > Library / runtime developers
 - > OS / VM developers
 - > Application developers
- What are the challenges?

Topics

- Multi-core Architectures
- Brief Overview of Pthreads, OpenMP and Java
- The Challenges

Topics

- Multi-core Architectures
- Brief Overview of Pthreads, OpenMP and Java
- The Challenges

Current Popular Parallel Programming Languages and APIs

- Shared Memory
 - > Memory is “shared” unless declared “private”.
 - > Accessing shared memory by direct read or write.
 - > Examples: Pthreads, OpenMP, Java, C#
 - > Closer to multi-core architectures than the rest two.
- Global Address Space
 - > Memory is “private” unless declared “shared”.
 - > Accessing shared memory by direct read or write.
 - > Examples: UPC, Co-array Fortran
- Message Passing
 - > No shared memory. Exchange data via messages.
 - > Example: MPI

Current Popular Parallel Programming Languages and APIs

- Shared Memory
 - > Memory is “shared” unless declared “private”.
 - > Accessing shared memory by direct read or write.
 - > Examples: Pthreads, OpenMP, Java, C#
 - > Closer to multi-core architectures than the rest two.
- Global Address Space
 - > Memory is “private” unless declared “shared”.
 - > Accessing shared memory by direct read or write.
 - > Examples: UPC, Co-array Fortran
- Message Passing
 - > No shared memory. Exchange data via messages.
 - > Example: MPI

POSIX Threads (Pthreads)

- IEEE POSIX 1003.1c-1995 standard:
A standardized C language thread API
 - > thread creation and termination
 - > synchronization
 - > scheduling
 - > data management
 - > process interaction
- Now incorporated into Single UNIX Specification, Version 3.
 - > barriers, read/write locks, spin locks, etc.

Pthreads – a simple example

```
g = g + compute(1) + compute(2)
```

Pthreads – a simple example

```
g = g + compute(1) + compute(2)
```

```
main()
{
    pthread_t tid;
    void *status;
    int part;

    part = 1;
    pthread_create(&tid, NULL, work, (void *)part);

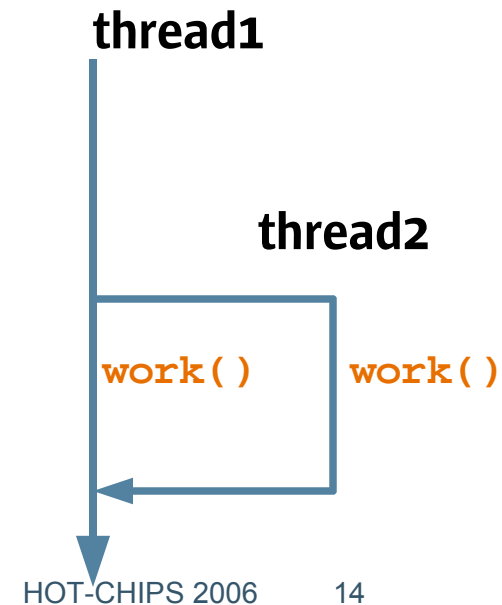
    part = 2;
    work((void *)part);

    pthread_join(tid, &status);
}
```

Pthreads – a simple example

```
g = g + compute(1) + compute(2)
```

```
main()  
{  
    pthread_t tid;  
    void *status;  
    int part;  
  
    part = 1;  
    pthread_create(&tid, NULL, work, (void *)part);  
  
    part = 2;  
    work((void *)part);  
  
    pthread_join(tid, &status);  
}
```



Pthreads – a simple example

```
g = g + compute(1) + compute(2)
```

```
void *work(void *arg)
{
    int result = compute((int)arg);
    pthread_mutex_lock(&lock);
    g += result;
    pthread_mutex_unlock(&lock);
    return NULL;
}

main()
{
    pthread_t tid;
    void *status;
    int part;

    part = 1;
    pthread_create(&tid, NULL, work, (void *)part);

    part = 2;
    work((void *)part);

    pthread_join(tid, &status);
}
```


Pthreads – a simple example

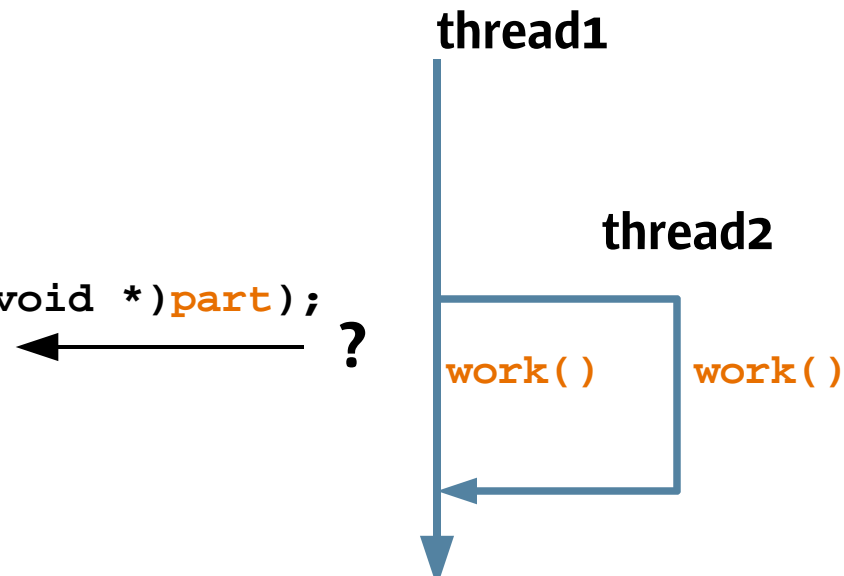
- Understand Concurrency: what is the status of thread 2 when pthread_create() returns?
 - > a) thread 2 has not started executing work().
 - > b) thread 2 is executing work().
 - > c) thread 2 has done work() and returned from work().

```
main()
{
    pthread_t tid;
    void *status;
    int part;

    part = 1;
    pthread_create(&tid, NULL, work, (void *)part);

    part = 2;
    work((void *)part);

    pthread_join(tid, &status);
}
```



Pthreads – a simple example

- Understand Concurrency: what is the status of thread 2 when pthread_create() returns?
 - > a) thread 2 has not started executing work().
 - > b) thread 2 is executing work().
 - > c) thread 2 has done work() and returned from work().

Could be any of the three!

```
main()
{
    pthread_t tid;
    void *status;
    int part;

    part = 1;
    pthread_create(&tid, NULL, work, (void *)part);

    part = 2;
    work((void *)part);

    pthread_join(tid, &status);
}
```

thread1

thread2

work()

work()

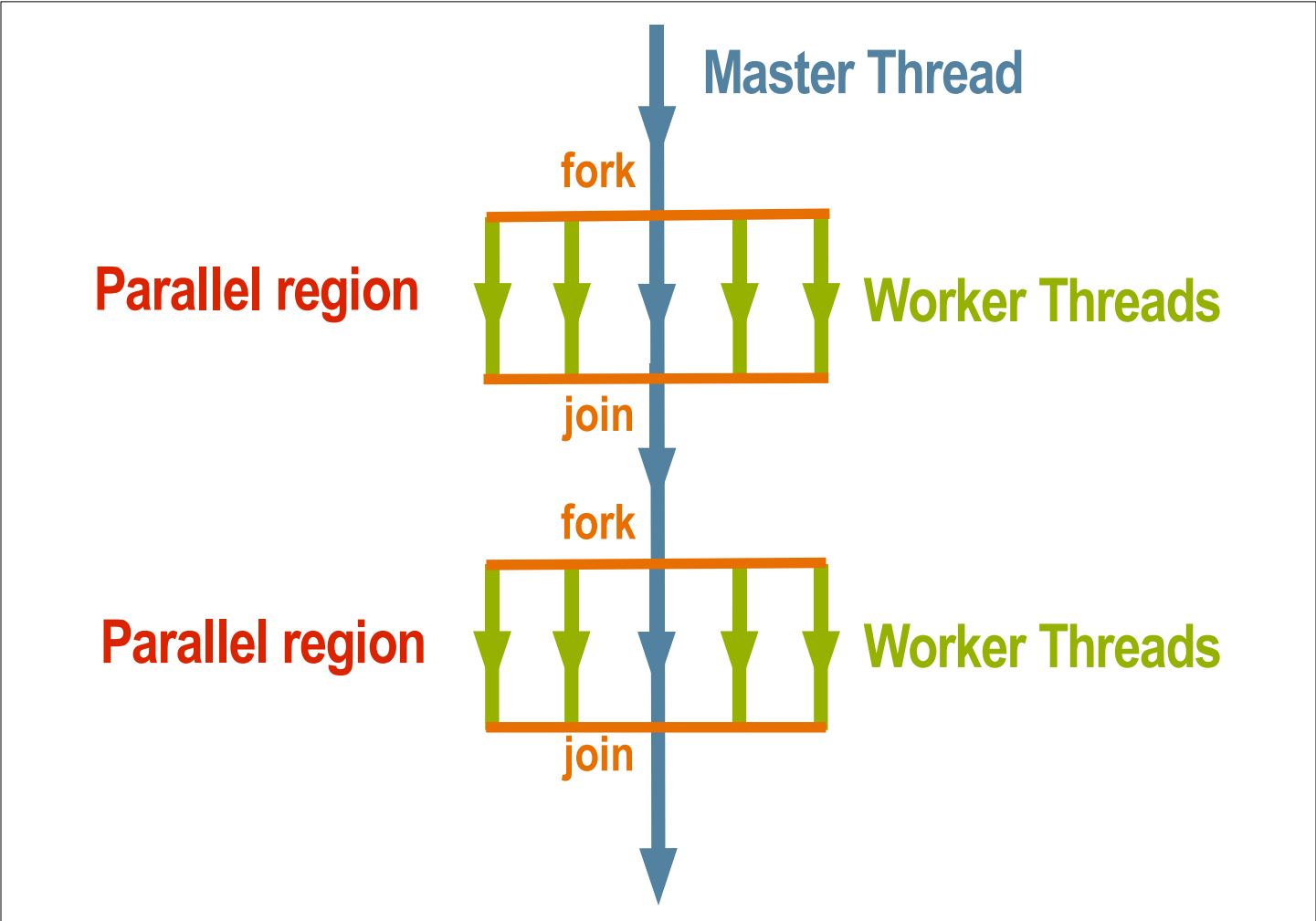
The Beauty of Pthreads

- Comprehensive Set of Functions and Specifications
 - > Thread management
 - > Synchronizations
 - > Scheduling
 - > Data management
 - > Interaction with process, signals, ...
- ✓ Flexible
- ✓ Building block of higher level runtime

OpenMP

- An API specification for writing shared memory parallel programs in C, C++ and Fortran
- Consists of:
 - > Compiler directives
 - > Runtime routines
 - > Environment variables
- Specification maintained by the OpenMP Architecture Review Board <http://www.openmp.org>
- Latest Specification: Version 2.5
- Language committee meetings for V3.0: since 9/2005

Thread Execution Model: fork and join



OpenMP – a simple example

```
g = g + compute(1) + compute(2)
```

```
for (int i=1; i<=2; i++) {  
    int result;  
  
    result = compute(i);  
    g += result;  
}
```

OpenMP – a simple example

```
g = g + compute(1) + compute(2)
```

```
#pragma omp parallel for reduction(+:g)
for (int i=1; i<=2; i++) {
    int result;

    result = compute(i);
    g += result;
}
```

The Beauty of OpenMP

- Possible to write a sequential program and its “functionally-equivalent” parallel version in the same code.

```
#pragma omp parallel for reduction(+:g)
for (int i=1; i<=2; i++) {
    int result;

    result = compute(i);
    g += result;
}
```

- ✓ makes debugging easier
- ✓ allows incremental parallelization

Java (J2SE)

- Basic threading support in the original spec
 - > new Thread(), Object.wait, notify, notifyAll, synchronized
- Threading related updates for J2SE5
 - > JSR-133: java memory model
 - > JSR-166: java.util.concurrent package
 - > Task scheduling framework
 - > Concurrent collections
 - > Atomic variables
 - > Synchronizers (e.g. barriers, semaphores, latches, ...)
 - > Locks
 - > Timing

Java – a simple example

```
g = g + compute(1) + compute(2)
```

```
import java.util.concurrent.atomic.*;
public class A implement Runnable {
    private int g;
    private AtomicInteger id = new AtomicInteger(0);
    ...
    public void run() {
        int myid, myg;
        myid = id.incrementAndGet();
        myg = compute(myid);
        synchronized(this) {
            g += myg;
        }
    }
}
```

```
public void go() {
    Thread t[];
    t = new Thread[2];
    id = 1;
    for (i=0; i<2; i++) {
        t[i] = new Thread(this);
        t[i].start();
    }
    for (i=0; i<2; i++) {
        try {
            t[i].join();
        }
        catch (InterruptedException iex){}
    }
}
}
```

Java – a simple example

```
g = g + compute(1) + compute(2)
```

```
import java.util.concurrent.atomic.*;
public class A implement Runnable {
    private int g;
    private AtomicInteger id = new AtomicInteger(0);
    ...
    public void run() {
        int myid, myg;
        myid = id.incrementAndGet();
        myg = compute(myid);
        synchronized(this) {
            g += myg;
        }
    }
}
```

**thread creation
start, and join**

```
public void go() {
    Thread t[];
    t = new Thread[2];
    id = 1;
    for (i=0; i<2; i++) {
        t[i] = new Thread(this);
        t[i].start();
    }
    for (i=0; i<2; i++) {
        try {
            t[i].join();
        }
        catch (InterruptedException iex){}
    }
}
```

Java – a simple example

```
g = g + compute(1) + compute(2)
```

```
import java.util.concurrent.atomic.*;
public class A implements Runnable {
    private int g;
    private AtomicInteger id = new AtomicInteger(0);
    ...
    public void run() {
        int myid, myg;
        myid = id.incrementAndGet();
        myg = compute(myid);
        synchronized(this) {
            g += myg;
        }
    }
}
```

atomic variables

```
public void go() {
    Thread t[];
    t = new Thread[2];
    id = 1;
    for (i=0; i<2; i++) {
        t[i] = new Thread(this);
        t[i].start();
    }
    for (i=0; i<2; i++) {
        try {
            t[i].join();
        }
        catch (InterruptedException iex){}
    }
}
```

Java – a simple example

```
g = g + compute(1) + compute(2)
```

```
import java.util.concurrent.atomic.*;
public class A implement Runnable {
    private int g;
    private AtomicInteger id = new AtomicInteger(0);
    ...
    public void run() {
        int myid, myg;
        myid = id.incrementAndGet();
        myg = compute(myid);
        synchronized(this) {
            g += myg;
        }
    }
}
```

synchronized block

```
public void go() {
    Thread t[];
    t = new Thread[2];
    id = 1;
    for (i=0; i<2; i++) {
        t[i] = new Thread(this);
        t[i].start();
    }
    for (i=0; i<2; i++) {
        try {
            t[i].join();
        }
        catch (InterruptedException iex){}
    }
}
}
```

The Beauty of Java

- A well studied and specified memory model.
- A set of easy to use, reliable, and performant concurrency utilities with many commonly used data structures.

Topics

- Multi-core Architectures
- Brief Overview of Pthreads, OpenMP and Java
- The Challenges

Challenges in Multithreaded Programming

1. Finding and creating concurrent tasks
2. Mapping tasks to threads
3. Defining and implementing synchronization protocols
4. Dealing with race conditions
5. Dealing with deadlocks
6. Dealing with memory model
7. Composing parallel tasks
8. Achieving scalability
9. Achieving portable & predictable performance
10. Recovering from errors
11. Dealing with all single thread issues

For Each Challenge

- What is it? 
- How do programmers deal with it now?
- Any help from Pthreads, OpenMP or Java? **P/O/J**
- Any tools that can help now? 

Find and Create Concurrent Tasks

- Concurrency at Data Level
 - > e.g. 6 crying babies waiting for diaper-change
- Concurrency at Function Level
 - > e.g. changing diaper for a baby, complaining to the spouse, while playing sudoku in mind
- Granularity

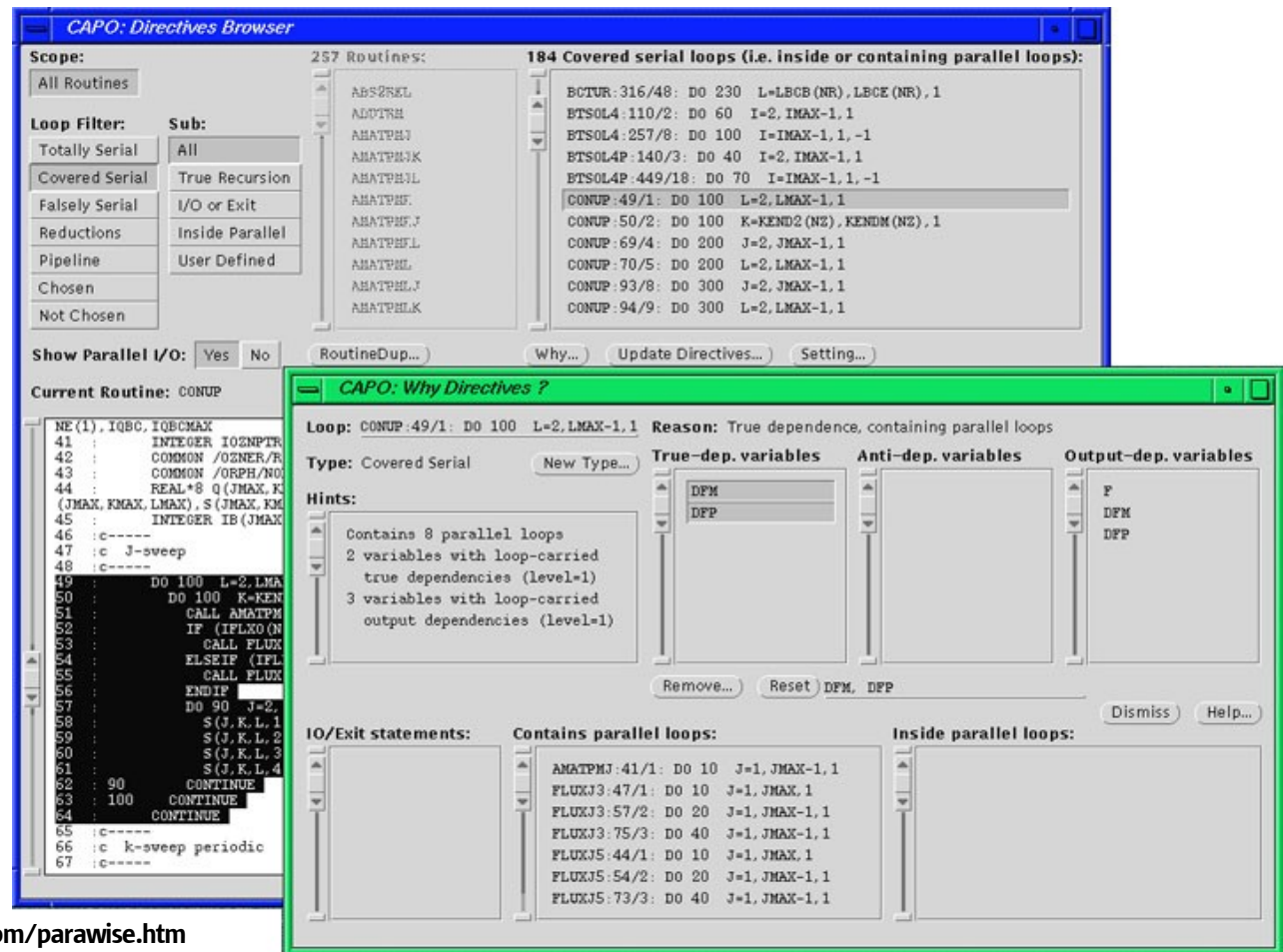


Automatic Parallelization

- Available in many compilers
 - > SGI, Sun, Intel, IBM, ...
 - > `f90 -xautopar old-serial.f`
- Long history
- Mostly loop based parallelization
- Limited adoption

Parallelization Assistance Tool

- ParaWise from Parallel Software Products, Inc.
 - > a semi-automatic parallelization tool for Fortran programs
 - > Long history
 - > Limited adoption



The image shows two screenshots of the CAPO (Cray Automated Parallelization Optimizer) interface.

The top screenshot is the **CAPO: Directives Browser**. It displays a list of 257 routines and 184 covered serial loops. The "Current Routine" is `CONUP`. The interface includes a "Loop Filter" section with buttons for "All Routines", "Totally Serial", "Covered Serial", "Falsely Serial", "Reductions", "Pipeline", "Chosen", and "Not Chosen". The "Sub:" section has buttons for "All", "True Recursion", "I/O or Exit", "Inside Parallel", and "User Defined". The "Show Parallel I/O:" section has "Yes" and "No" buttons. The "Current Routine: CONUP" section shows the following Fortran code:

```

NE(1), IQBC, IQBCMAX
41 : INTEGER IOENPTR
42 : COMMON /QENR/R
43 : COMMON /ORPH/NO
44 : REAL*8 Q(JMAX, K
(JMAX, KMAX, LMAX), S(JMAX, KM
45 : INTEGER IB(JMAX
46 : c-----
47 : c J-sweep
48 : c-----
49 : DO 100 L=2, LMA
DO 100 K=KEM
51 : CALL AMATPM
52 : IF (IPLXO(N
53 : CALL FLUX
54 : ELSEIF (IPL
55 : CALL FLUX
56 : ENDIF
57 : DO 90 J=2,
58 : S(J, K, L, 1
59 : S(J, K, L, 2
60 : S(J, K, L, 3
61 : S(J, K, L, 4
62 : 90 CONTINUE
63 : 100 CONTINUE
64 : CONTINUE
65 : c-----
66 : c k-sweep periodic
67 : c-----
  
```

The bottom screenshot is the **CAPO: Why Directives?** window. It provides details for a specific loop: `Loop: CONUP:49/1: DO 100 L=2, LMAX-1,1`. The "Reason" is "True dependence, containing parallel loops". The "Type" is "Covered Serial". The "Hints" section lists: "Contains 8 parallel loops", "2 variables with loop-carried true dependencies (level=1)", and "3 variables with loop-carried output dependencies (level=1)". The "True-dep. variables" section lists `DPM` and `DPP`. The "Anti-dep. variables" and "Output-dep. variables" sections are empty. The "IO/Exit statements:" section is empty. The "Contains parallel loops:" section lists several nested loops with their respective indices and limits. The "Inside parallel loops:" section is empty. The window includes "Remove...", "Reset", "Dismiss", and "Help..." buttons.



Parallelization Assistance Tool

- A valuable tool may not need to be fancy.
- Dieter an Mey (RWTH, Aachen University) :
 - > “I consider the following a very valuable assistance tool: a tool that can find all the places in the source code where global variables are declared and all the places where the global variables are referenced, and show these places with a different color.”
 - > “I got the functionality from the Foresys tool from Simulog for Fortran and I found it quite useful. It would be nice to have such features in an IDE. ”

Find and Create Concurrent Tasks **P/O/J**

- Pthread
 - > No direct support.
- Java
 - > No direct support.
 - > The use of Executor/Thread Pools makes the concept of task explicit. (See details later.)

Find and Create Concurrent Tasks **P/O/J**

- OpenMP

- > Use *worksharing* loops for basic data level parallelism

```
#pragma omp for
for (i=1; i<n; i++)
    work(data[i]);
```

- > Use *worksharing* sections for basic function level parallelism

```
#pragma omp sections
{
    #pragma omp section
    func1();
    #pragma omp section
    func2();
}
```

Tasks and Threads Mapping

- Threads vs Tasks
 - > Thread: parallel execution engines provided by the OS
 - > Tasks: concurrency unit in program logic
- A program can have tens of thousands of concurrent tasks.
 - > e.g. players in a game server, components in a simulation system, requests in a web server
- The system can only supply hundreds of threads.
 - > e.g. because of resource limitation
- There is overhead associated with creating and destroying threads.

Tasks and Threads Mapping

- Use well-designed data structures to hold task state.
- Deploy a task-to-thread model to deliver scalable and predictable performance.
 - > master/slave
 - > Master thread dispatches slave threads from a thread-pool to work on tasks. The threads in the pool are pre-allocated.
 - > pipeline
 - > Each thread performs a specific operation on a data item, and passes the data to the next thread in the pipeline.
 - > task pool
 - > Each thread execute a task from a pool and put new task generated to the pool.
 - > ...

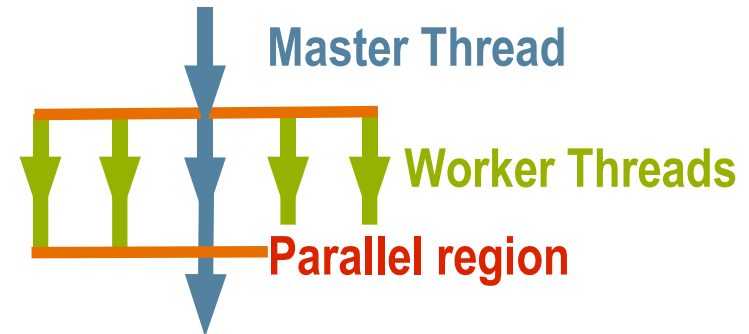
Tasks and Threads Mapping

P/O/J

- Pthreads : No direct support

- OpenMP :

- > Basic master/slave model



- > Three loop scheduling methods

- > Static



- > Dynamic



- > Guided



← 16 iterations, 4 threads →

Tasks and Threads Mapping

P/O/J

- Java
 - > Thread pool utilities in J2SE5 provide pre-built thread mechanism to run tasks to a bounded number of threads.
 - > To use thread pool,
 - > instantiate an implementation of the ExecutorService interface and hand it a set of tasks, or
 - > use configurable implementations: ThreadPoolExecutor and ScheduledThreadPoolExecutor
 - > Use Callable and Future interfaces to get result from tasks.

```
class NetworkService {
    private final ServerSocket serverSocket;
    private final ExecutorService pool;

    public NetworkService(int port, int poolSize) throws IOException {
        serverSocket = new ServerSocket(port);
        pool = Executors.newFixedThreadPool(poolSize);
    }
    public void serve() {
        try {
            for (;;) {
                pool.execute(new Handler(serverSocket.accept()));
            }
        } catch (IOException ex) {
            pool.shutdown();
        }
    }
}

class Handler implements Runnable {
    private final Socket socket;
    Handler(Socket socket) { this.socket = socket; }
    public void run() {
        // read and service request
    }
}
```

Thread pool

```
class NetworkService {
    private final ServerSocket serverSocket;
    private final ExecutorService pool;

    public NetworkService(int port, int poolSize) throws IOException {
        serverSocket = new ServerSocket(port);
        pool = Executors.newFixedThreadPool(poolSize);
    }
    public void serve() {
        try {
            for (;;) {
                pool.execute(new Handler(serverSocket.accept()));
            }
        } catch (IOException ex) {
            pool.shutdown();
        }
    }
}
```

```
class Handler implements Runnable {
    private final Socket socket;
    Handler(Socket socket) { this.socket = socket; }
    public void run() {
        // read and service request
    }
}
```

Task class

```
class NetworkService {
    private final ServerSocket serverSocket;
    private final ExecutorService pool;

    public NetworkService(int port, int poolSize) throws IOException {
        serverSocket = new ServerSocket(port);
        pool = Executors.newFixedThreadPool(poolSize);
    }
    public void serve() {
        try {
            for (;;) {
                pool.execute(new Handler(serverSocket.accept()));
            }
        } catch (IOException ex) {
            pool.shutdown();
        }
    }
}

class Handler implements Runnable {
    private final Socket socket;
    Handler(Socket socket) { this.socket = socket; }
    public void run() {
        // read and service request
    }
}
```

Create thread pool

```

class NetworkService {
    private final ServerSocket serverSocket;
    private final ExecutorService pool;

    public NetworkService(int port, int poolSize) throws IOException {
        serverSocket = new ServerSocket(port);
        pool = Executors.newFixedThreadPool(poolSize);
    }
    public void serve() {
        try {
            for (;;) {
                pool.execute(new Handler(serverSocket.accept()));
            }
        } catch (IOException ex) {
            pool.shutdown();
        }
    }
}

class Handler implements Runnable {
    private final Socket socket;
    Handler(Socket socket) { this.socket = socket; }
    public void run() {
        // read and service request
    }
}

```

Generate new tasks

```
class NetworkService {
    private final ServerSocket serverSocket;
    private final ExecutorService pool;

    public NetworkService(int port, int poolSize) throws IOException {
        serverSocket = new ServerSocket(port);
        pool = Executors.newFixedThreadPool(poolSize);
    }
    public void serve() {
        try {
            for (;;) {
                pool.execute(new Handler(serverSocket.accept()));
            }
        } catch (IOException ex) {
            pool.shutdown();
        }
    }
}

class Handler implements Runnable {
    private final Socket socket;
    Handler(Socket socket) { this.socket = socket; }
    public void run() {
        // read and service request
    }
}
```

Assign task to thread pool

```

class NetworkService {
    private final ServerSocket serverSocket;
    private final ExecutorService pool;

    public NetworkService(int port, int poolSize) throws IOException {
        serverSocket = new ServerSocket(port);
        pool = Executors.newFixedThreadPool(poolSize);
    }
    public void serve() {
        try {
            for (;;) {
                pool.execute(new Handler(serverSocket.accept()));
            }
        } catch (IOException ex) {
            pool.shutdown();
        }
    }
}

class Handler implements Runnable {
    private final Socket socket;
    Handler(Socket socket) { this.socket = socket; }
    public void run() {
        // read and service request
    }
}

```

Shutdown thread pool

Defining and Implementing Synchronizations

- Understand the application logic and design the high-level synchronization scheme
 - > e.g. traffic light
- Use the right synchronization primitives at the right place.
 - > e.g. need to coordinate threads so they are executed in phases
 - > use condition variables?
 - > use barriers?
 - > join threads and create new threads?

Defining and Implementing Synchronizations

P/O/J

- Pthreads:
 - > Mutex locks, condition variables, semaphores, barriers, reader-writer locks
- OpenMP:
 - > Locks, barriers, critical sections, ordered regions
- Java
 - > Synchronized method/block, wait/notify
 - > Condition variables, semaphores, barriers, countdown latch, exchanger, reader-writer locks

Race Conditions

- Race conditions occur when different threads access shared data without explicit synchronization.
- Race conditions can cause programs to behave in ways unexpected by the programmer.
- Bugs caused by race conditions are notoriously difficult to debug.

Race Conditions – An Example

Assume $a = b = 0$;

Thread 1:

```
a = a + 5;  
  
lock (a_lock);  
    b = b + 10;  
unlock (a_lock);  
  
barrier();
```

Thread 2:

```
lock (a_lock);  
    b = b - 10;  
unlock (a_lock);  
  
a = a - 5;  
  
barrier();
```

What's the value of 'a' and 'b' after the barrier?

Race Conditions – An Example

Thread 1:

```
a = a + 5;  
  
lock (a_lock);  
    b = b + 10;  
unlock (a_lock);  
  
barrier();
```

Assume a = b = 0;

Thread 2:

```
lock (a_lock);  
    b = b - 10;  
unlock (a_lock);  
  
a = a - 5;  
  
barrier();
```

We have a = 0 and b = 0.

Race Conditions – An Example

Assume $a = b = 0$;

Thread 1:

```
a = a + 5;  
  
lock (a_lock);  
    b = b + 10;  
unlock (a_lock);  
  
barrier();
```

Thread 2:

```
lock (a_lock);  
    b = b - 10;  
unlock (a_lock);  
  
a = a - 5;  
  
barrier();
```

'b' is 0.
'a' may be 0, 5, or -5.

Two Kinds of Race Conditions (1/2)

- Data Race
 - > Concurrent accesses (at least one is a write) to shared memory are not protected by critical sections. Atomicity of the accesses is lost.

Thread 1

```
lock();  
acnt1 = acnt1 + delta_x;  
acnt2 = acnt2 - delta_x;  
unlock();
```

Thread 2

```
lock();  
acnt1 = acnt1 + delta_y;  
acnt2 = acnt2 - delta_y;  
unlock();
```

Two Kinds of Race Conditions (2/2)

- General Race
 - > The order of accesses to shared memory is not enforced by synchronization.

Thread 1

```
put_job_in_queue(job)  
post();
```

Thread 2

```
wait();  
job = get_job_from_queue();
```

Data Race vs General Race

- A data race is also a general race.
- A data race can be fixed by using critical sections or adding synchronizations.
- A general race can be fixed by adding synchronizations to enforce ordering.
- Watch out for general races when the shared-memory accesses should occur in specific order.
- Watch out for data races in asynchronous programs.
- Tools are available to detect data races.



Data Race Detection Tools

- Static Detection
 - > + Can be fast and consume little memory.
 - > + Does not affect the behavior of program.
 - > + Is not affected by input data and scheduling
 - > + Can be used to check OS kernels and device drivers
 - > - False positives
- Example: LockLint from Sun Studio
 - > Reports data races and deadlocks due to inconsistent use of locking techniques.
 - > Originates from *WARLOCK*, which was designed to detect such errors in Solaris kernels and device drivers.



Data Race Detection Tools

- Runtime Detection
- Examples:
 - > Helgrind from Valgrind
 - > HP Visual Threads
 - > Intel Thread Checker
 - > Sun Data Race Detection Tool



Data Race Detection Tools

Sun Studio Analyzer [test.1.er]

File View Timeline Help

Races Race Source Experiments

Total Races: 4

Race #1, Vaddr :0xffbfeec4

- Access 1: Read, main -- MP doall from line 42 [_\$d1A42.main] + 0x00000060, line 45 in "omp_prime.c"
- Access 2: Write, main -- MP doall from line 42 [_\$d1A42.main] + 0x0000008C, line 46 in "omp_prime.c"

Total Traces: 2

Race #2, Vaddr :0xffbfeec4

- Access 1: Write, main -- MP doall from line 42 [_\$d1A42.main] + 0x0000008C, line 46 in "omp_prime.c"
- Access 2: Write, main -- MP doall from line 42 [_\$d1A42.main] + 0x0000008C, line 46 in "omp_prime.c"

Total Traces: 1

Race #3, Vaddr :(Multiple Addresses)

- Access 1: Write, main -- MP doall from line 42 [_\$d1A42.main] + 0x0000007C, line 45 in "omp_prime.c"
- Access 2: Write, main -- MP doall from line 42 [_\$d1A42.main] + 0x0000007C, line 45 in "omp_prime.c"

Total Traces: 1

Race #4, Vaddr :0x21418

- Access 1: Read, is_prime + 0x00000074, line 18 in "omp_prime.c"
- Access 2: Write, is_prime + 0x00000114, line 21 in "omp_prime.c"

Total Traces: 1

Summary Race Details

Data for Selected Race

Id: Race #1

Vaddr: 0xffbfeec4

Access 1

Type: Read

main -- MP doall from line 42 [_\$d1A42.main] +

Access 2

Type: Write

main -- MP doall from line 42 [_\$d1A42.main] +

What to Do After a Data-race is Found?

- 1) Check whether it is a false positive.
 - > A false positive is a reported data race that actually does not exist in the program.
 - > No tool is perfect.
 - > For example, the tool may not understand the synchronizations in your program.

What to Do After a Data-race is Found?

- 1) Check whether it is a false positive.
 - > A false positive is a reported data race that actually does not exist in the program.
 - > No tool is perfect.
 - > For example, the tool may not understand the synchronizations in your program.
- 2) Check whether it is a benign race.
 - > The programmer may put data race there in order to get better performance.
 - > Yes? Are you sure? Check again!
 - > Programs with intentional data races are very difficult to get right.
 - > Are you relying on sequential consistency?

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.

```
int pthread_create(pthread_t *thread,  
                  const pthread_attr_t *attr,  
                  void *(*start_routine)(void*),  
                  void *arg);
```

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.

```
int pthread_create(pthread_t *thread,  
                  const pthread_attr_t *attr,  
                  void *(*start_routine)(void*),  
                  void *arg);
```

```
void *work(void *arg)  
{  
    int myid = *(int *)arg;  
    data[myid] = update(data[myid]);  
    return NULL;  
}  
  
for (i=0; i<THREADS; i++)  
    pthread_create(tids[i], NULL, work, &i);
```

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.

```
int pthread_create(pthread_t *thread,  
                  const pthread_attr_t *attr,  
                  void *(*start_routine)(void*),  
                  void *arg);
```

```
void *work(void *arg)  
{  
    int myid = *(int *)arg;  
    data[myid] = update(data[myid]);  
    return NULL;  
}
```

```
for (i=0; i<THREADS; i++)  
    pthread_create(tids[i], NULL, work, &i);
```

Data race reported

Bug!

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.
 - > Most tools detect *certain* kind of data races only!

thread 1

```
acct1 += x;  
  
acct2 -= x;
```

thread 2

```
acct1 += y;  
  
acct2 -= y;
```

thread 3

```
print acct1, acct2
```

Data Race! Accounts are not balanced!
“acct1 + acct2” is not constant.

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.
 - > Most tools detect *certain* kind of data races only!

thread 1

```
lock();  
acct1 += x;  
unlock();  
lock();  
acct2 -= x;  
unlock();
```

thread 2

```
lock();  
acct1 += y;  
unlock();  
lock();  
acct2 -= y;  
unlock();
```

thread 3

```
lock();  
print acct1, acct2  
unlock();
```

**A wrong fix. No data race is reported.
But accounts are not balanced at certain stages!
“account1 + account2” is not constant.**

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.
 - > Most tools detect *certain* kind of data races only!

thread 1

```
lock();  
acnt1 += x;  
  
acnt2 -= x;  
unlock();
```

thread 2

```
lock();  
acnt1 += y;  
  
acnt2 -= y;  
unlock();
```

thread 3

```
lock();  
print acnt1, acnt2  
unlock();
```

A correct fix.

What to Do After a Data-race is Found?

- 3) Fix the bug, not the data race!
 - > A data race could be a bug or caused by a bug.
 - > Most tools detect *certain* kind of data races only!

- 4) Run the detection tool again after code change!
 - > A tool may not be able to detect all data races by design.
 - > The code change may introduce or reveal another data race.

Making Global Variables Private

- Pthreads: thread specific data
 - > `pthread_key_create()`, `pthread_key_delete()`
 - > `pthread_setspecific()`, `pthread_getspecific()`
- OpenMP: “threadprivate” directive


```
int total_for_this_thread;
#pragma omp threadprivate (total_for_this_thread)
```
- Java: `java.lang.ThreadLocal` class
 - > `private static ThreadLocal<int> totoal_for_this_thread;`
- Thread Local Storage
 - > Provided by compiler and runtime linker


```
__thread int total_for_this_thread;
```

Use MT-Safe Routines

- Four levels of MT safe attributes for library interfaces.
- 1) Unsafe
 - > Contains global and static data that are not protected. User should make sure only one thread at time to execute the call.

Unsafe Function

```

ctime
localtime
asctime
gmtime
ctermid
getlogin
rand
readdir
strtok
tmpnam
  
```

Reentrant counterpart

```

ctime_r
localtime_r
asctime_r
gmtime_r
ctermid_r
getlogin_r
rand_r
readdir_r
strtok_r
tmpnam_r
  
```

Use MT-Safe Routines

- 2) Safe
 - > Global and static data are protected. Might not provide any concurrency between calls made by different threads.
 - > Example: malloc in libc(3c)
- 3) MT-Safe
 - > Safe and can provide a reasonable amount of concurrency.
 - > Example: malloc in libmtmalloc(3LIB).

Use MT-Safe Routines

- 4) Async-signal-safe
 - > Can be safely called from a signal handler.
 - > Example:
 - > Not async-signal-safe: malloc(), pthread_getspecific()
 - > Async-signal-safe: open(), read()

Deadlocks

- Two or more competing actions are waiting for the other to finish. None can finish.
- A common misunderstanding: a deadlock must involve mutex locks.
 - > Deadlocks may be caused by mis-use of synchronizations or communications, such as barriers.
- Livelock
 - > Two or more competing actions continually change their state in response to changes in the other actions. None will finish.

Necessary Conditions for a Deadlock

1. Mutual exclusion: a resource is exclusively owned by the process to which it is assigned.
2. Hold and wait: processes request new resources while holding some resources.
3. No preemption: only the process holding a resource can release it.
4. Circular wait: there is a circular chain of processes in which each process waits for a resource that the next process in the chain holds.

Deadlock: Prevention and Avoidance

- Prevention
 - > Use protocols to ensure at least one of the four conditions does not hold.
 - > Example: enforce lock hierarchy to ensure circular-wait never happens.
- Avoidance
 - > Use global resource information to decide whether the current request can be satisfied or delayed.
 - > Example: Banker's algorithm

Deadlock: Prevention and Avoidance

- It may be easy to apply the text book rules to deadlocks that involve obvious resources.
- In many applications, the four conditions are subtle, especially when the deadlock is caused by communication.
- Always examine the program logic!

Deadlocks: Example 1

thread 1:

Lock(a);

Lock(b);

thread 2:

Lock(b);

Lock(c);

thread 3:

Lock(c);

Lock(a);

The program may or may not deadlock in a particular run.

Deadlocks: Example 2

thread tid1:

```
pthread_join(tid2, NULL);
```

thread tid2:

```
pthread_join(tid1, NULL);
```

Deadlocks: Example 3

```
#pragma omp parallel num_threads(4)
{
    if (omp_thread_num() < 3) {
        #pragma omp barrier
    }
    else {
        #pragma omp barrier
    }
}
```

Illegal OpenMP program!

Deadlocks: Example 4

```
int flag = 0;
```

```
thread t1: while (!flag); ← May loop forever.  
          pthread_exit(0);
```

```
thread t2: flag = 1;  
          pthread_join(t1, NULL);
```

Deadlocks: Example 4

thread t1: `while (!flag);` ← **May loop forever.**
`pthread_exit(0);`

↓ **Compiler optimization**

```
reg = flag;
if (!reg) {
    while (1);
}
pthread_exit(0);
```

Deadlocks

P/O/J

- Pthreads
 - > Some error checking on locks, but not deadlocks.
 - > If the mutex type is PTHREAD_MUTEX_ERRORCHECK and a thread attempts to relock a mutex that it has already locked, an error will be returned.
 - > If the mutex type is PTHREAD_MUTEX_ERRORCHECK and a thread attempts to unlock a mutex that it has not locked or a mutex which is unlocked, an error will be returned.
- OpenMP
 - > Makes one kind of deadlock illegal.
 - > Barriers in a parallel region must be encountered by all threads in the team in the same order.

Deadlocks

P/O/J

- Java
 - > Structured synchronized method/block makes the critical regions clearer than unstructured locks.
 - > Locks associated with *synchronized* keyword are always released when the thread leaves (due to exception or not) the synchronized scope.
 - > This automatic-releasing could be a bad thing for error recovery since it may leave the data in an inconsistent state that is not protected by locks.



Deadlock Detection Tools

- Static Checking
 - > LockLint from Sun
- Runtime Checking
 - > Thread Checker from Intel
 - > OpenMP Runtime Error Detection from Sun

```
> setenv SUNW_MP_WARN TRUE
> a.out
WARNING (libmtnsk): Threads at barrier from different directives.
  Thread at barrier from init.c:12.
  Thread at barrier from forbidden.c:17.
Possible Reasons:
Worksharing constructs not encountered by all threads in the
team in the same order.
Incorrect placement of barrier directives.
```

Memory Consistency Model

- The order in which memory operations will appear to be executed to a programmer.
- What affects the memory consistency?
 - > Language
 - > Compiler
 - > Hardware

Example: Lazy Initialization

- The Sequential Version

```
typedef struct
{
    int data1;
    int data2;
    ...
} A;
```

```
A *init_single_A()
{
    static A *single_A;
    if (single_A == NULL) {
        single_A = malloc(sizeof(A));
        single_A->data1 = ...;
        ...
    }
    return single_A;
}
```

It does not work in mt applications.

Example: Lazy Initialization

- Multi-threaded Versions

```
A *init_single_A()
{
    static A *single_A;
    lock();
    if (single_A == NULL) {
        single_A = malloc(sizeof(A));
        single_A->data1 = ...;
        ...
    }
    unlock();
    return single_A;
}
```

Example: Lazy Initialization

- Multi-threaded Versions

```
A *init_single_A()
{
    static A *single_A;
    lock();
    if (single_A == NULL) {
        single_A = malloc(sizeof(A));
        single_A->data1 = ...;
        ...
    }
    unlock();
    return single_A;
}
```

Not efficient

```
A *init_single_A()
{
    static A *single_A;
    A *temp = single_A;
    if (temp == NULL) {
        lock();
        if (single_A == NULL) {
            temp = malloc(sizeof(A));
            temp->data1 = ...;
            ...
            single_A = temp;
        }
        unlock();
    }
    return temp;
}
```

Example: Lazy Initialization

- Multi-threaded Versions

```
A *init_single_A()
{
    static A *single_A;
    lock();
    if (single_A == NULL) {
        single_A = malloc(sizeof(A));
        single_A->data1 = ...;
        ...
    }
    unlock();
    return single_A;
}
```

Not efficient

“Double-checked Locking”

```
A *init_single_A()
{
    static A *single_A;
    A *temp = single_A;
    if (temp == NULL) {
        lock();
        if (single_A == NULL) {
            temp = malloc(sizeof(A));
            temp->data1 = ...;
            ...
            single_A = temp;
        }
        unlock();
    }
    return temp;
}
```

May be broken

Double-checked Locking

```
A *init_single_A()
{
    static A *single_A;
    A *temp = single_A;
    if (temp == NULL) {
        lock();
        if (single_A == NULL) {
            temp = malloc(sizeof(A));
            temp->data1 = ...;
            ...
            single_A = temp;
        }
        unlock();
    }
    return temp;
}
```

- > The compiler may reorder these two writes.
- > Even if the compiler does not reorder them, a thread on another processor may perceive the two writes in a different order.
- > Therefore, a thread on another processor may read wrong value of `single_A->data1`.

```
A *p = init_single_A();
... = p->data1;
```

Double-checked Locking

```
A *init_single_A()
{
    static A *single_A;
    A *temp = single_A;
    if (temp == NULL) {
        lock();
        if (single_A == NULL) {
            temp = malloc(sizeof(A));
            temp->data1 = ...;
            ...
            memory_barrier();
            single_A = temp;
        }
        unlock();
    }
    return temp;
}
```

- > A possible fix.
- > Still broken on some architectures, e.g. Alpha.

Memory Consistency Model

P/O/J

- Pthreads
 - > No formal specification
 - > Shared accesses must be synchronized by calling pthread synchronization functions.
- C++/C
 - > Assumes single thread program execution.
 - > “volatile” restricts compiler optimization, but it does not address the memory consistency issue.
 - > Memory model for multithreaded C++ is being worked on.

Memory Consistency Model

P/O/J

- OpenMP
 - > Detailed clarification. No formal specification.
 - > Each thread has a temporary view of shared memory.
 - > A flush operation restricts the ordering of memory operations and synchronizes a thread's temporary view with shared memory.
 - > All threads must observe any two flush operations with overlapping variable lists in sequential order.

Memory Consistency Model

P/O/J

- Java: revised and clarified by JSR-133
 - > Volatile variables
 - > Final variables
 - > Immutable objects (objects whose fields are only set in their constructor)
 - > Thread- and memory-related JVM functionality and APIs such as class initialization, asynchronous exceptions, finalizers, thread interrupts, and the sleep, wait, and join methods of class Thread

Memory Consistency Model

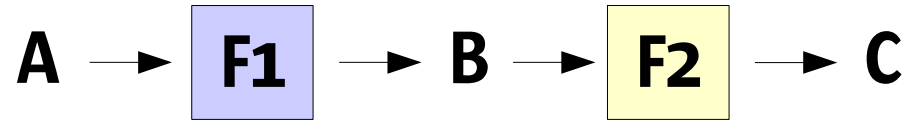
- Avoid writing codes that have deliberate data races. It is tricky and difficult to understand and debug.

Composing Parallel Tasks

- Modular design is common in software development.
- How to compose modules that are multi-threaded?
- Three models of composition
 - > Sequential composition
 - > Parallel composition
 - > Nested composition

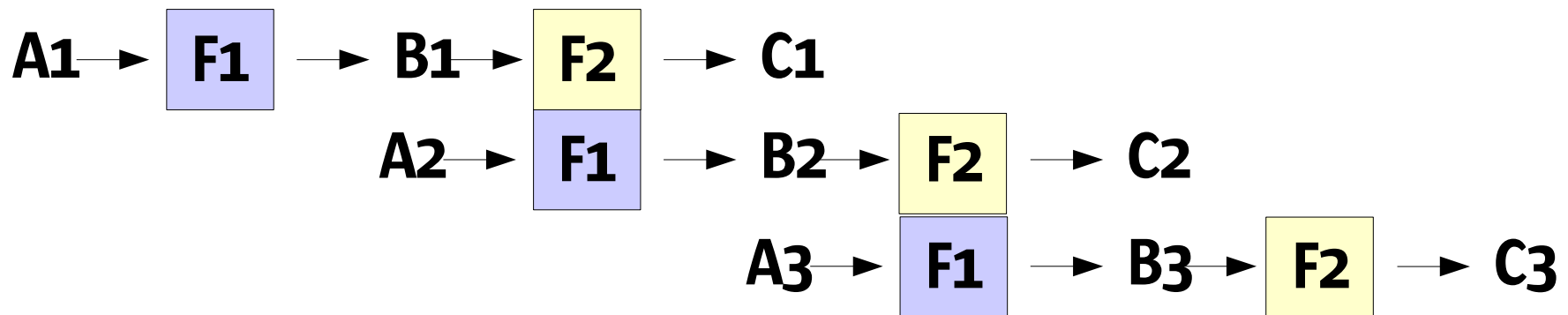
Composition: sequential vs parallel

- Sequential composition:



> n threads for F1 and n threads for F2

- Parallel composition:



> n1 threads for F1 and n2 threads for F2, n1+n2=n

Composition: sequential vs parallel

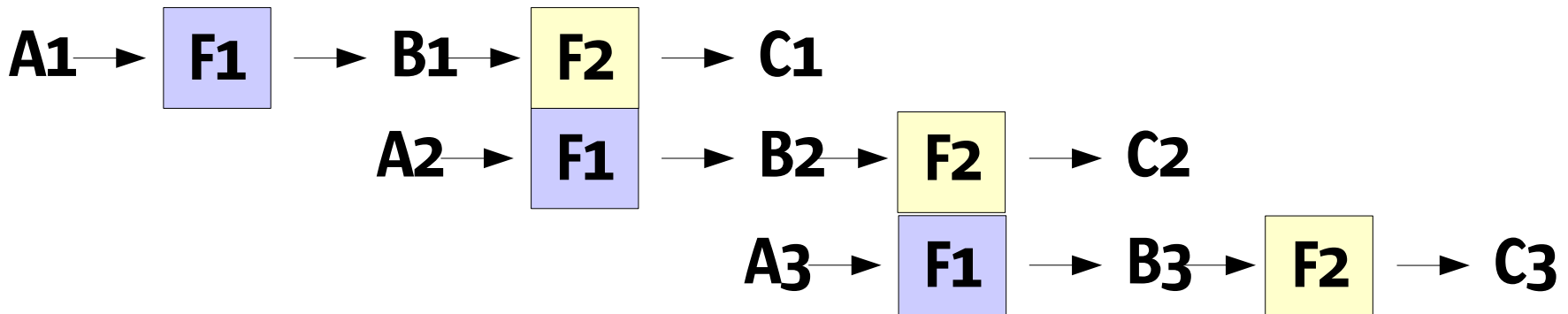
- Sequential composition



> n threads for F1 and

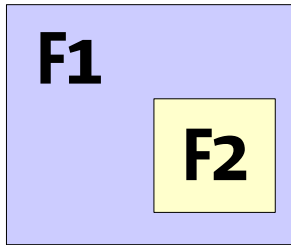
+ may improve resource utilization by overlapping F1 and F2
 - may introduce extra communication, data migration, or thread migration

- Parallel composition:



> n1 threads for F1 and n2 threads for F2, n1+n2=n

Nested Composition



n_1 thread for F1 and n_2 threads for F2, $n_1 * n_2 = n$

- Straight-forward approach: outer-module and inner module create their own threads.
 - > Threads explosion.
- More complicated approach: outer-module and inner module get threads from the same thread pool.
 - > Need a common threading frame work.
 - > Implies API interface change or external global state
- Locks are not composable in general!
 - > Dead locks

Composing Parallel Tasks

P/O/J

- Pthreads and Java: No direct support
- OpenMP:
 - > Support nested composition of parallel tasks through nested parallel regions.

```
#omp parallel for
for (i=0; i<n; i++) {
    #omp parallel for
    for (j=0; j<m; j++) {
        ...
    }
}
```

- > Inflexible and possible low utilization of threads

Scalability

- Performance does not scale with the number of threads used.
- Common reason 1: system oversubscribed
 - > e.g. 16 hardware threads, 32 software threads
- Common reason 2: not enough concurrency
 - > Amdahl's law
- Common reason 3: lock contention
 - > Frequent locking, unlocking activities
 - > Long lock holding time



plockstat(1M)

- A utility in Solaris 10 that gathers and displays user-level locking statistics.
- Uses plockstat DTrace provider.
- Three types of lock events can be traced.
 - > Contention events - probes for user level lock contention
 - > Hold events - probes for lock acquiring, releasing etc.
 - > Error events - error conditions.



plockstat(1M) - example

```
>plockstat ./a.out
```

Mutex block

Count	nsec	Lock	Caller
863	3822626	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x5c
900	3423154	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x20
340	2969890	a.out`mutex	a.out`_\$_p1A12.main+0x30
15835	2894962	a.out`mutex	a.out`_\$_p1A12.main+0x30
415	2456968	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x20
296	2208085	a.out`mutex	a.out`_\$_p1A12.main+0x30
441	2129956	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x5c
14032	2054952	a.out`mutex	a.out`_\$_p1A12.main+0x30
4	42900	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc
1	14800	libc.so.1`_uberdata	libmtdsk.so.1`threads_fini+0x174



plockstat(1M) - example

```
>plockstat ./a.out
```

Mutex block

Count	nsec	Lock	Caller
863	3822626	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x5c
900	3423154	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x20
340	2969890	a.out`mutex	a.out`_\$p1A12.main+0x30
15835	2894962	a.out`mutex	a.out`_\$p1A12.main+0x30
415	2456968	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x20
296	2208085	a.out`mutex	a.out`_\$p1A12.main+0x30
441	2129956	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x5c
14032	2054952	a.out`mutex	a.out`_\$p1A12.main+0x30
4	42900	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc
1	14800	libc.so.1`_uberdata	libmtdsk.so.1`threads_fini+0x174

libc_malloc_lock

a.out`mutex



plockstat(1M) - example

(plockstat output cont.)

Mutex spin

Count	nsec	Lock	Caller
4225	13553	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x20
72525	13518	a.out`mutex	a.out`_\$p1A12.main+0x30
1711	13399	a.out`mutex	a.out`_\$p1A12.main+0x30
7	11600	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc
8111	10778	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x5c

Mutex unsuccessful spin

Count	nsec	Lock	Caller
1315	31766	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x20
636	31008	a.out`mutex	a.out`_\$p1A12.main+0x30
29867	30875	a.out`mutex	a.out`_\$p1A12.main+0x30
31	30590	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc
1304	29273	libc.so.1`libc_malloc_lock	a.out`_\$p1A12.main+0x5c
7	26320	libc.so.1`_uberdata	libmtsk.so.1`threads_fini+0x174
19	24328	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc



plockstat(1M) - example

(plockstat output cont.)

Mutex spin

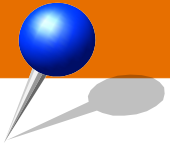
Count	nsec	Lock	Caller
4225	13553	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x20
72525	13518	a.out`mutex	a.out`_\$_p1A12.main+0x30
1711	13399	a.out`mutex	a.out`_\$_p1A12.main+0x30
7	11600	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc
8111	10778	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x5c

Mutex unsuccessful spin

Count	nsec	Lock	Caller
1315	31766	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x20
636	31008	a.out`mutex	a.out`_\$_p1A12.main+0x30
29867	30875	a.out`mutex	a.out`_\$_p1A12.main+0x30
31	30590	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc
1304	29273	libc.so.1`libc_malloc_lock	a.out`_\$_p1A12.main+0x5c
7	26320	libc.so.1`_uberdata	libmtsk.so.1`threads_fini+0x174
19	24328	libc.so.1`_uberdata	libc.so.1`_thr_exit_common+0xbc

libc_malloc_lock

a.out`mutex



plockstat(1M) - example

- libc.so.1`libc_malloc_lock
 - > Use libmtmalloc or libumem
- a.out`mutex
 - > Resize the critical sections
 - > scope
 - > lock holding time
 - > Replace critical sections with atomic operations

Atomic Operations

- atomic_ops
 - > atomic_ops package from HP
 - > atomic_ops(3c), available on Solaris 10
 - atomic_add_16(), atomic_or_32(), atomic_add_32_nv(), ...
- OpenMP: “atomic” directive

```
#pragma omp atomic
a++;
```
- Java: atomic variable classes
 - > Nine flavors of atomic variables.

Portable & Predictable Performance

- Use the same program on different platforms and achieve consistent performance.
- Ideal: performance is determined only by the algorithm used.
- What affects performance?
 - > Algorithm
 - > Runtime (VM)
 - > Compiler
 - > OS
 - > HW



Performance Analysis Tools

- Sun Studio Performance Analyzer
- Intel Vtune and Thread Profiler
- Quest Software JProbe
- Borland Optimizeit
- ...
- OS level tools: dtrace, mpstat, lockstat, ...



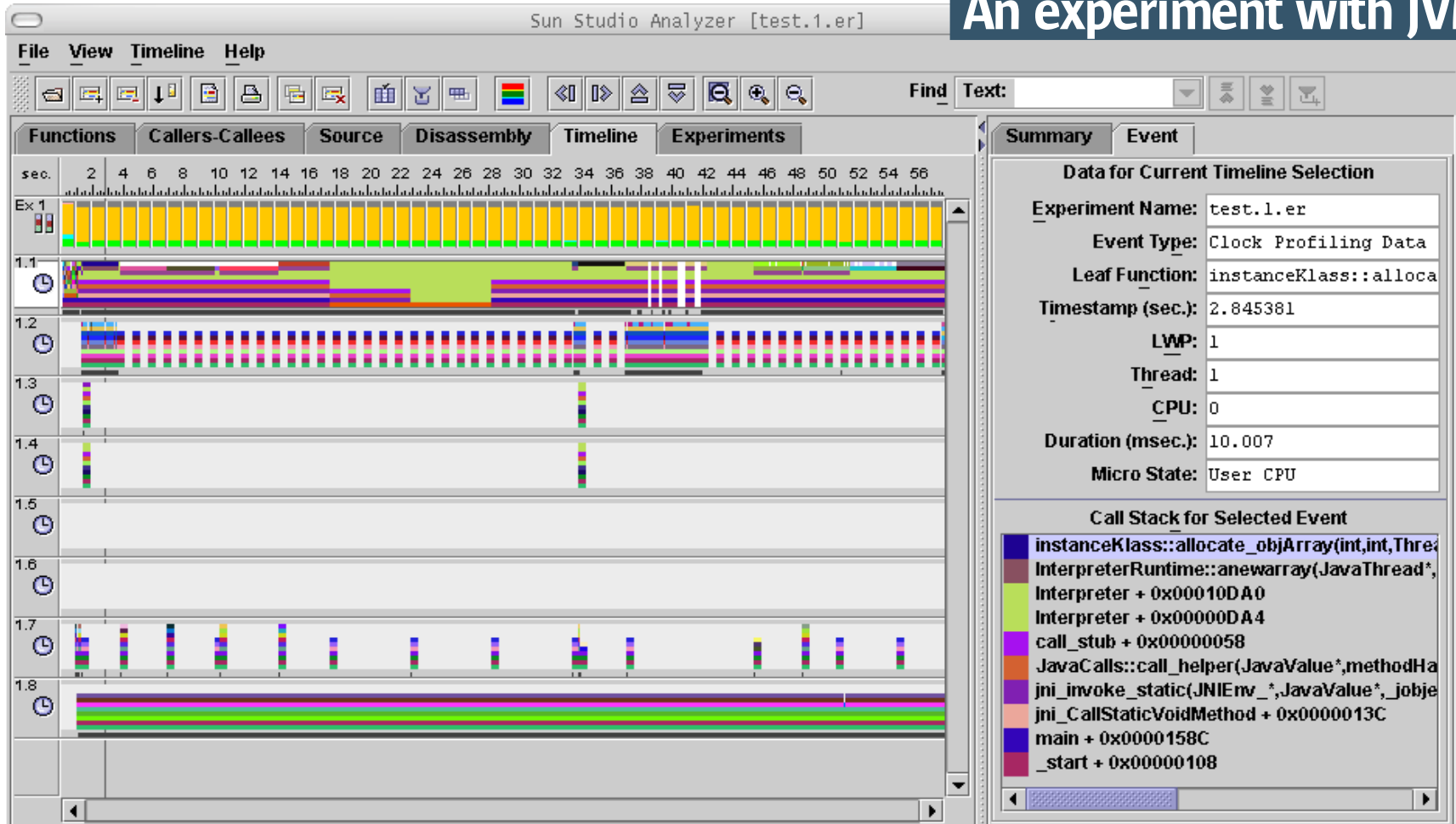
Sun Studio Performance Analyzer

- General purpose application level profiler
- Post-mortem analysis
- MT aware
- Supports dynamic libraries
- Clock profiling
- Hardware counter profiling
- Support for OpenMP



Sun Studio Performance Analyzer

An experiment with JVM



Single threaded Java application. JVM is multi-threaded.



Sun Studio Performance Analyzer

An experiment with JVM

The screenshot displays the Sun Studio Analyzer interface for a test named 'test.1.er'. The main window shows a 'Timeline' view with a horizontal axis representing time in seconds (0 to 56). A yellow box highlights a specific thread (1.1) which is executing user application code. To the right, a 'Summary' panel provides details for the current timeline selection:

Data for Current Timeline Selection	
Experiment Name:	test.1.er
Event Type:	Clock Profiling Data
Leaf Function:	instanceClass::alloca
Timestamp (sec.):	2.845381
LWP:	1
Thread:	1
CPU:	0
Duration (msec.):	10.007
Micro State:	User CPU

Below the summary, a 'Call Stack for Selected Event' is shown, listing the following frames from top to bottom:

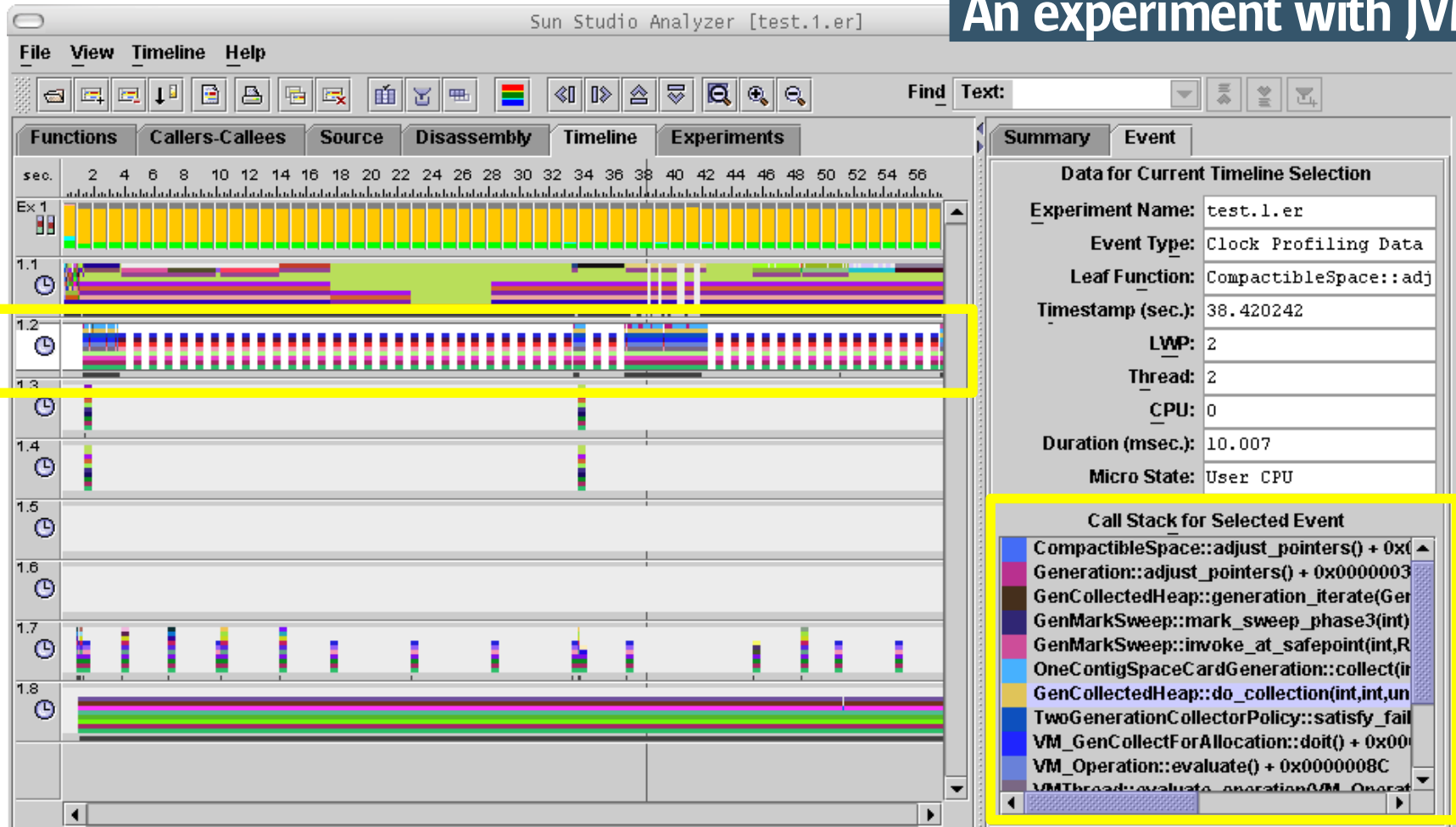
- instanceClass::allocate_objArray(int,int,Thread)
- InterpreterRuntime::anewarray(JavaThread*,int)
- Interpreter + 0x00010DA0
- Interpreter + 0x00000DA4
- call_stub + 0x00000058
- JavaCalls::call_helper(JavaValue*,methodHandle,Thread)
- jni_invoke_static(JNIEnv*,JavaValue*,_jobject)
- jni_CallStaticVoidMethod + 0x0000013C
- main + 0x0000158C
- _start + 0x00000108

This thread executes user application.



Sun Studio Performance Analyzer

An experiment with JVM

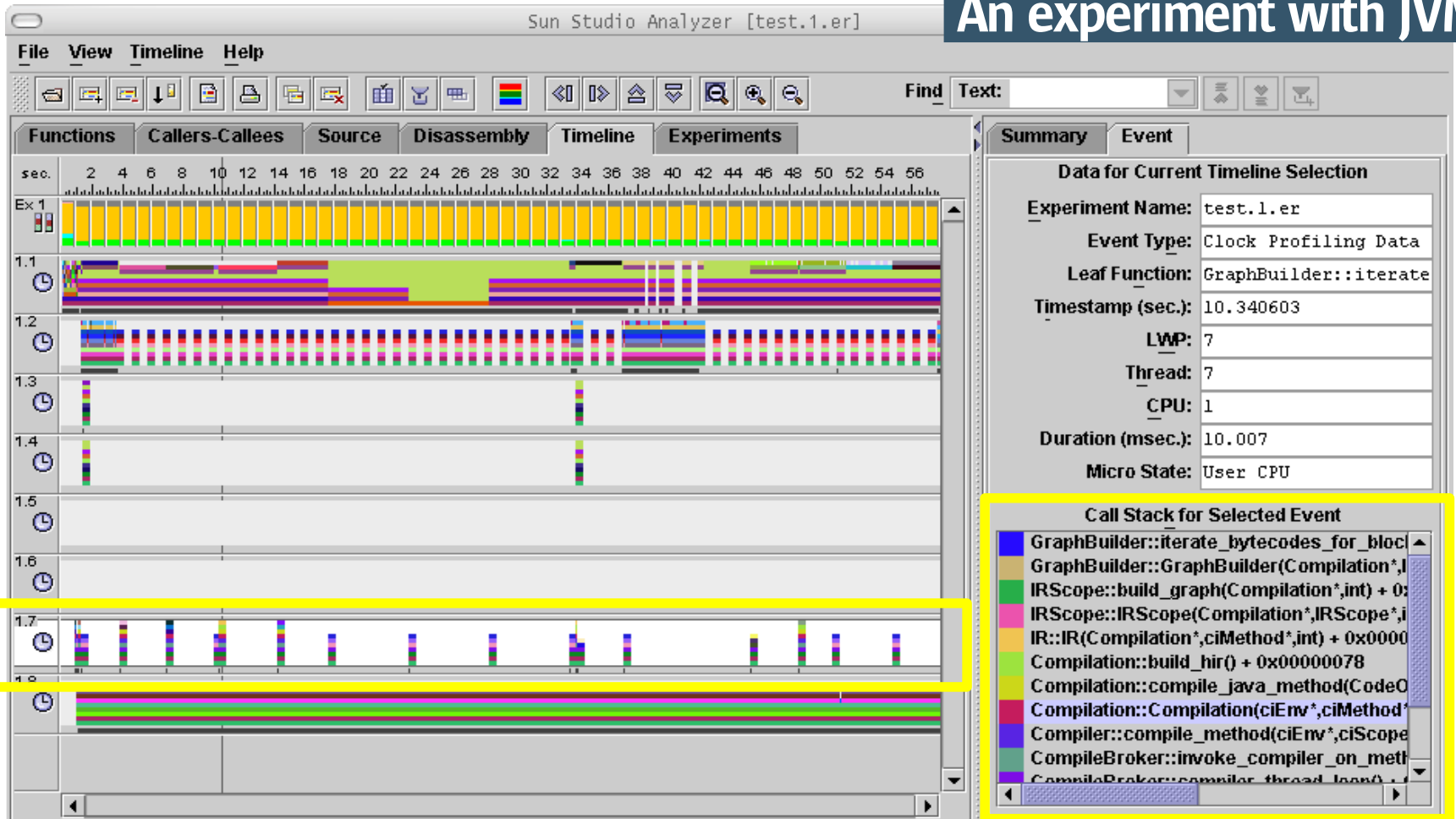


This thread does garbage collection.



Sun Studio Performance Analyzer

An experiment with JVM

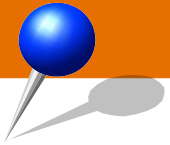


This thread does JIT compilation.



Performance Tools: A Challenge

- Be aware of the threading model and provide relevant information.
 - > Identify tasks, in addition to threads
 - > Provide state information at different abstraction levels
 - > Goals
 - > Identify serial bottlenecks
 - > Help resolve load imbalances



Sun Studio Performance Analyzer

Sun Studio Analyzer [tfs.16x16_e6900_hwprof_ds.er]

File View Timeline Help

Find Text:

Functions Callers-Callees Source Disassembly Experiments Threads Seconds

Display Mode: Text Graphical

User CPU (sec.)	OMP Work (sec.)	OMP Wait (sec.)	Name
1 255.810	489.470	815.540	<Total>
71.810	80.120	1.940	Thread 1
79.860	36.570	44.960	Thread 2
79.760	24.820	56.710	Thread 3
80.090	12.100	69.430	Thread 4
79.870	12.940	68.590	Thread 5
79.680	18.500	63.030	Thread 6
78.240	46.110	35.420	Thread 7
79.640	12.650	68.880	Thread 8
79.260	16.370	65.160	Thread 9
77.550	47.550	33.980	Thread 10
77.420	56.010	25.520	Thread 11
78.200	35.420	46.110	Thread 12
79.260	14.280	67.250	Thread 13
79.310	12.770	68.760	Thread 14
76.860	47.710	33.820	Thread 15
79.000	15.550	65.980	Thread 16

An OpenMP Experiment



Sun Studio Performance Analyzer

User CPU (sec.)	OMP Work (sec.)	OMP Wait (sec.)	Name
1 255.810	489.470	815.540	<Total>
71.810	80.120	1.940	Thread 1
79.860	36.570	44.960	Thread 2
79.760	24.820	56.710	Thread 3
80.090	12.100	69.430	Thread 4
79.870	12.940	68.590	Thread 5
79.680	18.500	63.030	Thread 6
78.240	46.110	35.420	Thread 7
79.640	12.650	68.880	Thread 8
79.260	16.370	65.160	Thread 9
77.550	47.550	33.980	Thread 10
77.420	56.010	25.520	Thread 11
78.200	35.420	46.110	Thread 12
79.260	14.280	67.250	Thread 13
79.310	12.770	68.760	Thread 14
76.860	47.710	33.820	Thread 15
79.000	15.550	65.980	Thread 16

User-CPU time for all threads is the same.
But that does not mean the OpenMP workload for the threads is balanced.

An OpenMP Experiment



Sun Studio Performance Analyzer

Sun Studio Analyzer [tfs.16x16_e6900_hwprof_ds.er]

File View Timeline Help

Functions Callers-Callees Source Disassembly Experiments Threads Seconds

Display Mode: Text Graphical

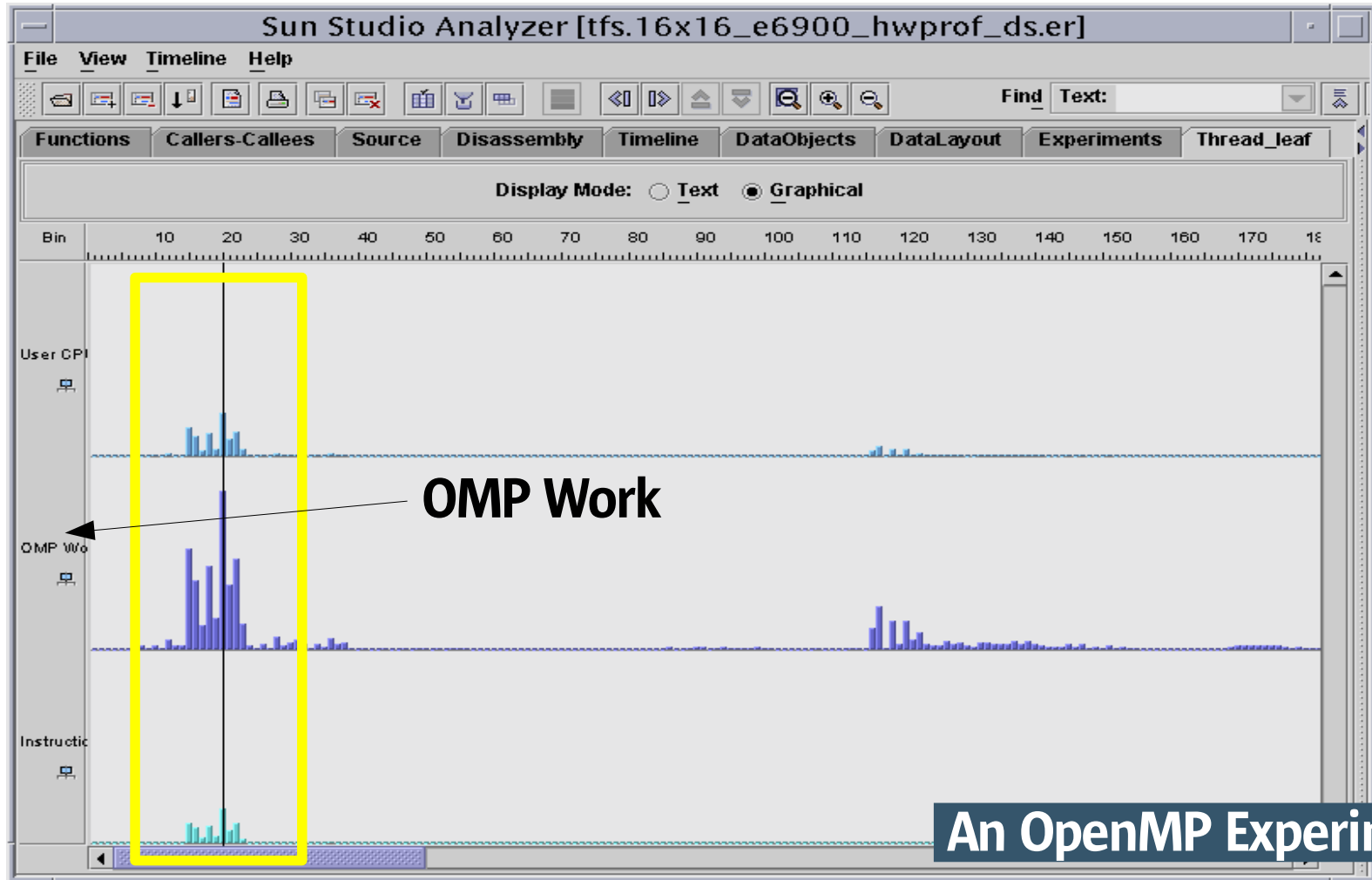
User CPU (sec.)	OMP Work (sec.)	OMP Wait (sec.)	Name
1 255.810	489.470	815.540	<Total>
71.810	80.120	1.940	Thread 1
79.860	36.570	44.960	Thread 2
79.760	24.820	56.710	Thread 3
80.090	12.100	69.430	Thread 4
79.870	12.940	68.590	Thread 5
79.680	18.500	63.030	Thread 6
78.240	46.110	35.420	Thread 7
79.640	12.650	68.880	Thread 8
79.260	16.370	65.160	Thread 9
77.550	47.550	33.980	Thread 10
77.420	56.010	25.520	Thread 11
78.200	35.420	46.110	Thread 12
79.260	14.280	67.250	Thread 13
79.310	12.770	68.760	Thread 14
76.860	47.710	33.820	Thread 15
79.000	15.550	65.980	Thread 16

The OMP Work metrics shows the load is not balanced.

An OpenMP Experiment



Sun Studio Performance Analyzer



Error Recovery

- Catch the errors and restore the application to a stable state.
- Error recovery becomes more difficult in mt applications because it may be an asynchronous problem!
 - > Notify other threads
 - > Who? How? When?
 - > Update global structures
 - > Deadlocks or infinite loops in recovering process!
- A common approach: use check points



MT- capable debuggers

- Basic Features
 - > Examining the thread states
 - > Examining thread specific data
 - > Setting breakpoints for a specific thread
 - > Control the execution of the threads
- Advanced Features
 - > Integrated with deadlock/data race detection tools
 - > Examining synchronization status
 - > Be aware of threading model
- Tools
 - > Dbx from Sun, TotalView from Etnus, GDB, ...



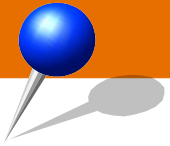
dbx – Examine Status

- Print a list of all known threads
 - > (dbx) threads
 - o t@1 a l@1 ?() breakpoint in work()
 - *>t@2 a l@2 work() breakpoint in work()
 - t@3 a l@3 work() running in _thr_setup()
- Switch the viewing context to another thread
 - > (dbx) thread t@3



dbx – Set Breakpoints

- A break point can be set specifically to a particular thread.
 - > `(dbx) stop at 430 -thread t@2`
- Breaking is synchronous - “stop the world”
 - > When any thread stops, all other threads *sympathetically* stop.



dbx - Control Execution

- Keep the given thread from ever running
> (dbx) thread -suspend t@2
- Resume a thread
> (dbx) thread -resume t@2
- Single step a thread
> (dbx) step t@2
> (dbx) next t@2

Summary

- Multi-core architecture provides a new thread-rich environment.
- Challenges in multi-threaded programming need to be addressed at all levels.
 - > languages
 - > compilers and tools
 - > libraries and utilities
 - > VM and OS
 - > programming practice

To Learn More ...

- Pthread
 - > “Programming with POSIX Threads”, David Butenhof
- OpenMP
 - > www.openmp.org
- Java
 - > “Java Concurrency in Practice”, Brian Goetz et al
 - > JSR-133
 - > <http://www.cs.umd.edu/~pugh/java/memoryModel/>
 - > JSR-166
 - > <http://gee.cs.oswego.edu/dl/concurrency-interest/index.html>

To Learn More ...

- “Developing and Tuning Applications on UltraSPARC T1 Chip Multithreading Systems”, Denis Sheahan
 - > <http://www.sun.com/blueprints/1205/819-5144.html>
- Dtrace
 - > <http://www.sun.com/bigadmin/content/dtrace/>
- Threading Methodolgy: Principles and Practices
 - > <http://www.intel.com/cd/software/products/asmo-na/eng/threading/219349.htm>
- Using the Sun Studio Data-Race Detection Tool
 - > <http://developers.sun.com/prodtech/cc/downloads/drdt/using.html>

Thank you!

yuan.lin@sun.com

Agenda

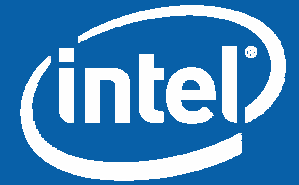
Multithreaded Programming

Transactional Memory (TM)

- TM Introduction
- TM Implementation Overview
- Hardware TM Techniques
- Software TM Techniques



Q&A



Transactional Memory Introduction

Ali-Reza Adl-Tabatabai
Programming Systems Lab
Intel Corporation

Transactional memory definition

Memory transaction: A sequence of memory operations that either execute completely (commit) or have no effect (abort)

An “all or nothing” sequence of operations

- On commit, all memory operations appear to take effect as a unit (all at once)
- On abort, none of the stores appear to take effect

Transactions run in isolation

- Effects of stores are not visible until transaction commits
- No concurrent conflicting accesses by other transactions

Similar to database ACID properties

Transactional memory language construct

The basic **atomic** construct:

```
lock(L); x++; unlock(L);    →    atomic {x++;}
```

Declarative – user simply specifies, system implements “under the hood”

Basic atomic construct universally proposed

- HPCS languages (Fortress, X10, Chapel) provide atomic in lieu of locks
- Research extensions to languages – Java, C#, Atomos, CaML, Haskell, ...

Lots of recent research activity

- Transactional memory language constructs
- Compiling & optimizing atomic
- Hardware & software implementations of transactional memory



Example: Java 1.4 HashMap

Fundamental data structure

- Map: Key \rightarrow Value

```
public Object get(Object key) {  
    int idx = hash(key);           // Compute hash  
    HashEntry e = buckets[idx];   // to find bucket  
    while (e != null) {           // Find element in bucket  
        if (equals(key, e.key))  
            return e.value;  
        e = e.next;  
    }  
    return null;  
}
```

Not thread safe: don't pay lock overhead if you don't need it

Synchronized HashMap

Java 1.4 solution: Synchronized layer

- Convert any map to thread-safe variant
- Explicit locking – user specifies concurrency

```
public Object get(Object key)
{
    synchronized (mutex) // mutex guards all accesses to map m
    {
        return m.get(key);
    }
}
```

Coarse-grain synchronized HashMap:

- Thread-safe, easy to program
- Limits concurrency → poor scalability
 - E.g., 2 threads can't access disjoint hashtable elements

Transactional HashMap

Transactional layer via an 'atomic' construct

- Ensure all operations are atomic
- Implicit atomic directive – system discovers concurrency

```
public Object get(Object key)
{
    atomic                // System guarantees atomicity
    {
        return m.get(key);
    }
}
```

Transactional HashMap:

- Thread-safe, easy to program
- Good scalability

Transactions: Scalability

Concurrent read operations

- Basic locks do not permit multiple readers
 - Reader-writer locks
- Transactions automatically allow multiple concurrent readers

Concurrent access to disjoint data

- Programmers have to manually perform fine-grain locking
 - Difficult and error prone
 - Not modular
- Transactions automatically provide fine-grain locking

ConcurrentHashMap

Java 5 solution: Complete redesign

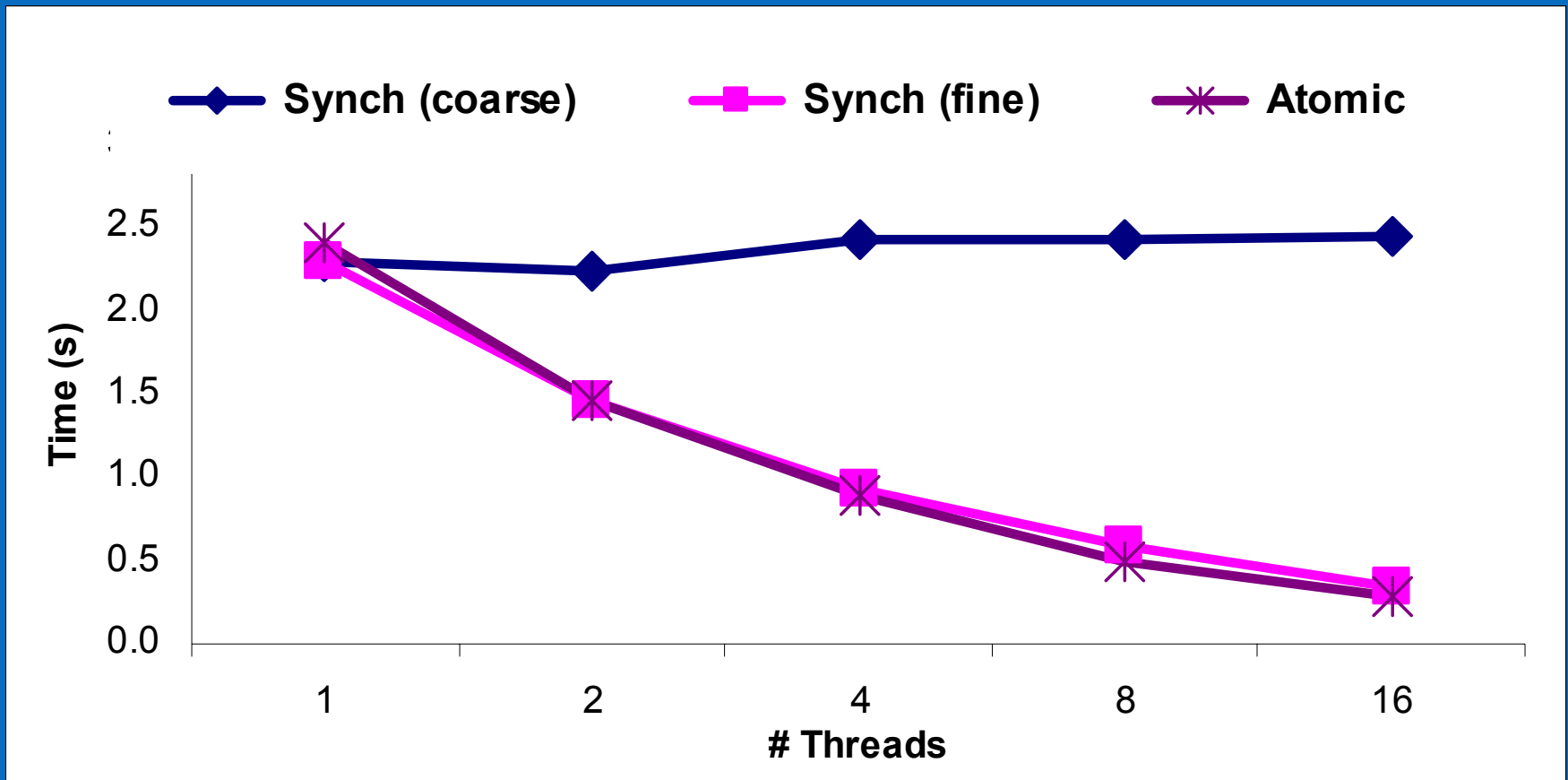
```
public Object get(Object key) {
    int hash = hash(key);
    // Try first without locking...
    Entry[] tab = table;
    int index = hash & (tab.length - 1);
    Entry first = tab[index];
    Entry e;

    for (e = first; e != null; e = e.next) {
        if (e.hash == hash && eq(key, e.key)) {
            Object value = e.value;
            if (value != null)
                return value;
            else
                break;
        }
    }
    ...

    // Recheck under synch if key not there or interference
    Segment seg = segments[hash & SEGMENT_MASK];
    synchronized(seg) {
        tab = table;
        index = hash & (tab.length - 1);
        Entry newFirst = tab[index];
        if (e != null || first != newFirst) {
            for (e = newFirst; e != null; e = e.next) {
                if (e.hash == hash && eq(key, e.key))
                    return e.value;
            }
        }
        return null;
    }
}
```

Fine-grain locking & concurrent reads: complicated & error prone

HashMap performance



Transactions scales as well as fine-grained locks

Transactional memory benefits

As easy to use as coarse-grain locks

Scale as well as fine-grain locks

Composition:

- Safe & scalable composition of software modules

Example: A bank application

Bank accounts with names and balances

- HashMap is natural fit as building block

```
class Bank {
    ConcurrentHashMap accounts;
    ...
    void deposit(String name, int amount) {
        int balance = accounts.get(name);           // Get the current balance
        balance = balance + amount;                // Increment it
        accounts.put(name, balance);               // Set the new balance
    }
    ...
}
```

Not thread-safe – Even with ConcurrentHashMap

Thread safety

Suppose Fred has \$100

T0: deposit("Fred", 10)

- `bal = acc.get("Fred") <- 100`
- `bal = bal + 10`
- `acc.put("Fred", bal) -> 110`

T1: deposit("Fred", 20)

- `bal = acc.get("Fred") <- 100`
- `bal = bal + 20`
- `acc.put("Fred", bal) -> 120`

Fred has \$120. \$10 lost.

Traditional solution: Locks

```
class Bank {
    ConcurrentHashMap accounts;
    ...
    void deposit(String name, int amount) {
        synchronized(accounts) {
            int balance = accounts.get(name);           // Get the current balance
            balance = balance + amount;                 // Increment it
            accounts.put(name, balance);                // Set the new balance
        }
    }
    ...
}
```

Thread-safe – but no scaling

- ConcurrentHashMap does not help
- Performance requires redesign from scratch & fine-grain locking

Transactional solution

```
class Bank {
    HashMap accounts;
    ...
    void deposit(String name, int amount) {
        atomic {
            int balance = accounts.get(name);           // Get the current balance
            balance = balance + amount;                 // Increment it
            accounts.put(name, balance);                 // Set the new balance
        }
    }
    ...
}
```

Thread-safe – and it scales!

Safe composition + performance



Transactional memory benefits

As easy to use as coarse-grain locks

Scale as well as fine-grain locks

Safe and scalable composition

Failure atomicity:

- Automatic recovery on errors



Traditional exception handling

```
class Bank {
  Accounts accounts;
  ...
  void transfer(String name1, String name2, int amount) {
    synchronized(accounts) {
      try {
        accounts.put(name1, accounts.get(name1)-amount);
        accounts.put(name2, accounts.get(name2)+amount);
      }
      catch (Exception1) {...}
      catch (Exception2) {...}
    }
    ...
  }
}
```

Manually catch all exceptions and determine what needs to be undone

Side effects may be visible to other threads before they are undone

Failure recovery using transactions

```
class Bank {  
    Accounts accounts;  
    ...  
    void transfer(String name1, String name2, int amount) {  
        atomic {  
            accounts.put(name1, accounts.get(name1)-amount);  
            accounts.put(name2, accounts.get(name2)+amount);  
        }  
    }  
    ...  
}
```

System rolls back updates on an exception
Partial updates not visible to other threads

Challenges in parallel programming

Finding independent tasks

Mapping tasks to threads

Defining & implementing synchronization protocol

Race conditions

Deadlock avoidance

Memory model

Composing parallel tasks

Scalability

Portable & predictable performance

Recovering from errors

... Single thread issues

→ Transactions address a lot of parallel programming problems



Challenges in parallel programming

Finding independent tasks

Mapping tasks to threads

Defining & implementing
synchronization protocol

Race conditions

Deadlock avoidance

Memory model

Composing parallel tasks

Scalability

Portable & predictable
performance

Recovering from errors

... Single thread issues

→ But not a silver bullet

Summary

Transactions provide many benefits over locks

- Automatic fine-grain concurrency
- Automatic read concurrency
- Deadlock avoidance
- Eliminates locking protocols
- Automatic failure recovery

→ Safe & scalable composition of thread-safe software modules

Challenge: How to implement transactions efficiently?



Agenda

- ❑ Multithreaded Programming

- ❑ Transactional Memory (TM)
 - TM Introduction
 - **TM Implementation Overview** ←
 - **Hardware TM Techniques**
 - Software TM Techniques

- ❑ Q&A



Transactional Memory Implementation Overview

Christos Kozyrakis

Computer Systems Laboratory
Stanford University

<http://csl.stanford.edu/~christos>



TM Implementation Requirements

- ❑ TM implementation must provide atomicity and isolation
 - Without sacrificing concurrency

- ❑ Basic implementation requirements
 - Data versioning
 - Conflict detection & resolution

- ❑ Implementation options
 - Hardware transactional memory (HTM)
 - Software transactional memory (STM)
 - Hybrid transactional memory



Data Versioning

- ❑ Manage uncommitted (new) and committed (old) versions of data for concurrent transactions

1. Eager or undo-log based

- Update memory location directly; maintain undo info in a log
- + Faster commit, direct reads (SW)
- Slower aborts, no fault tolerance, weak atomicity (SW)

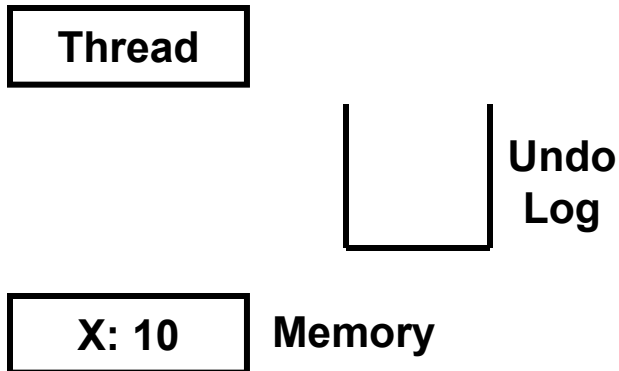
2. Lazy or write-buffer based

- Buffer writes until commit; update memory location on commit
- + Faster abort, fault tolerance, strong atomicity (SW)
- Slower commits, indirect reads (SW)

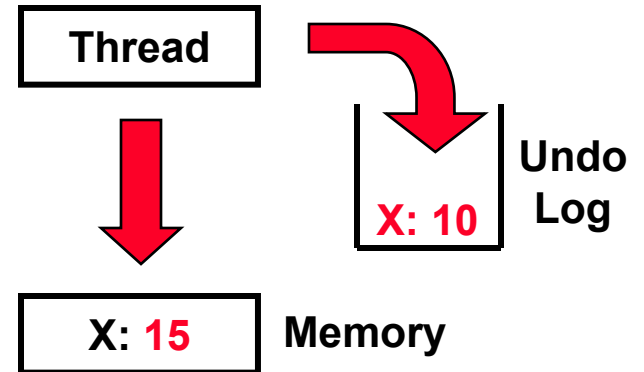


Eager Versioning Illustration

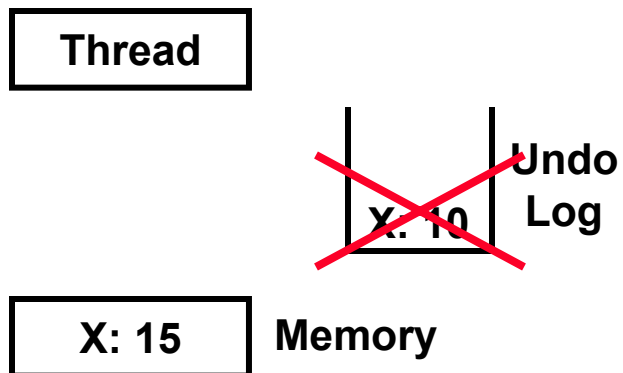
Begin Xaction



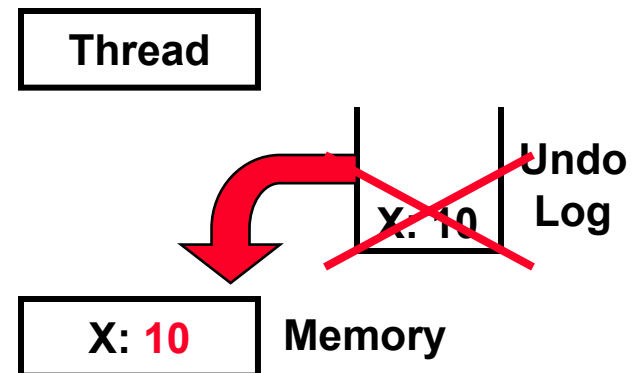
Write X ← 15



Commit Xaction



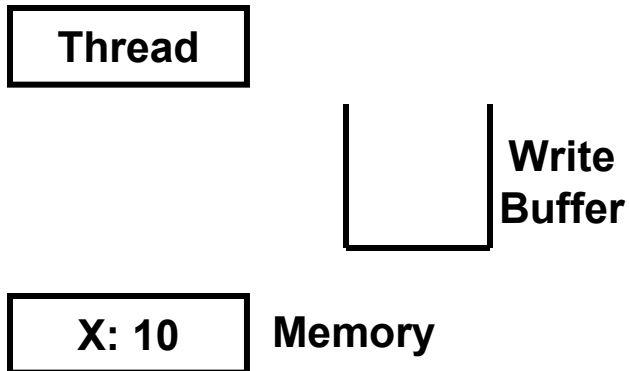
Abort Xaction



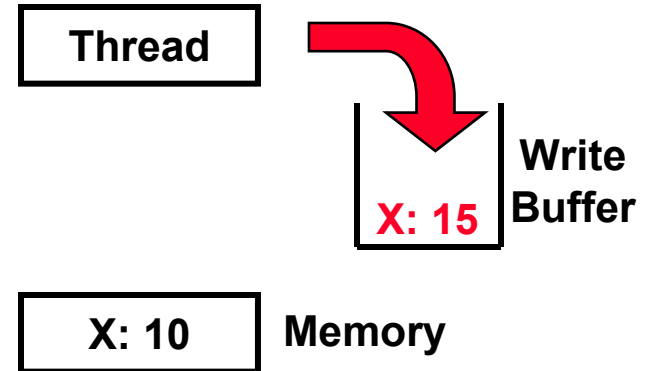


Lazy Versioning Illustration

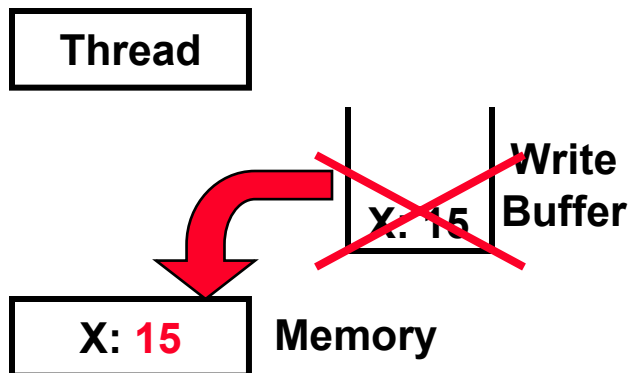
Begin Xaction



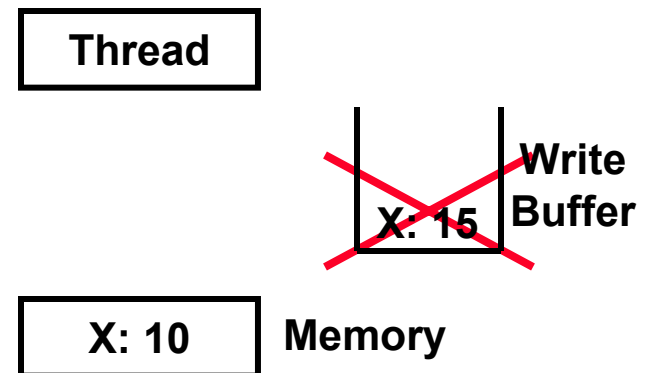
Write X ← 15



Commit Xaction



Abort Xaction





Conflict Detection

- ❑ Detect and handle conflicts between transaction
 - Read-Write and (often) Write-Write conflicts
 - For detection, a transactions tracks its read-set and write-set

- 1. Eager or encounter or pessimistic
 - Check for conflicts during loads or stores
 - HW: check through coherence lookups
 - SW: checks through locks and/or version numbers
 - Use contention manager to decide to stall or abort

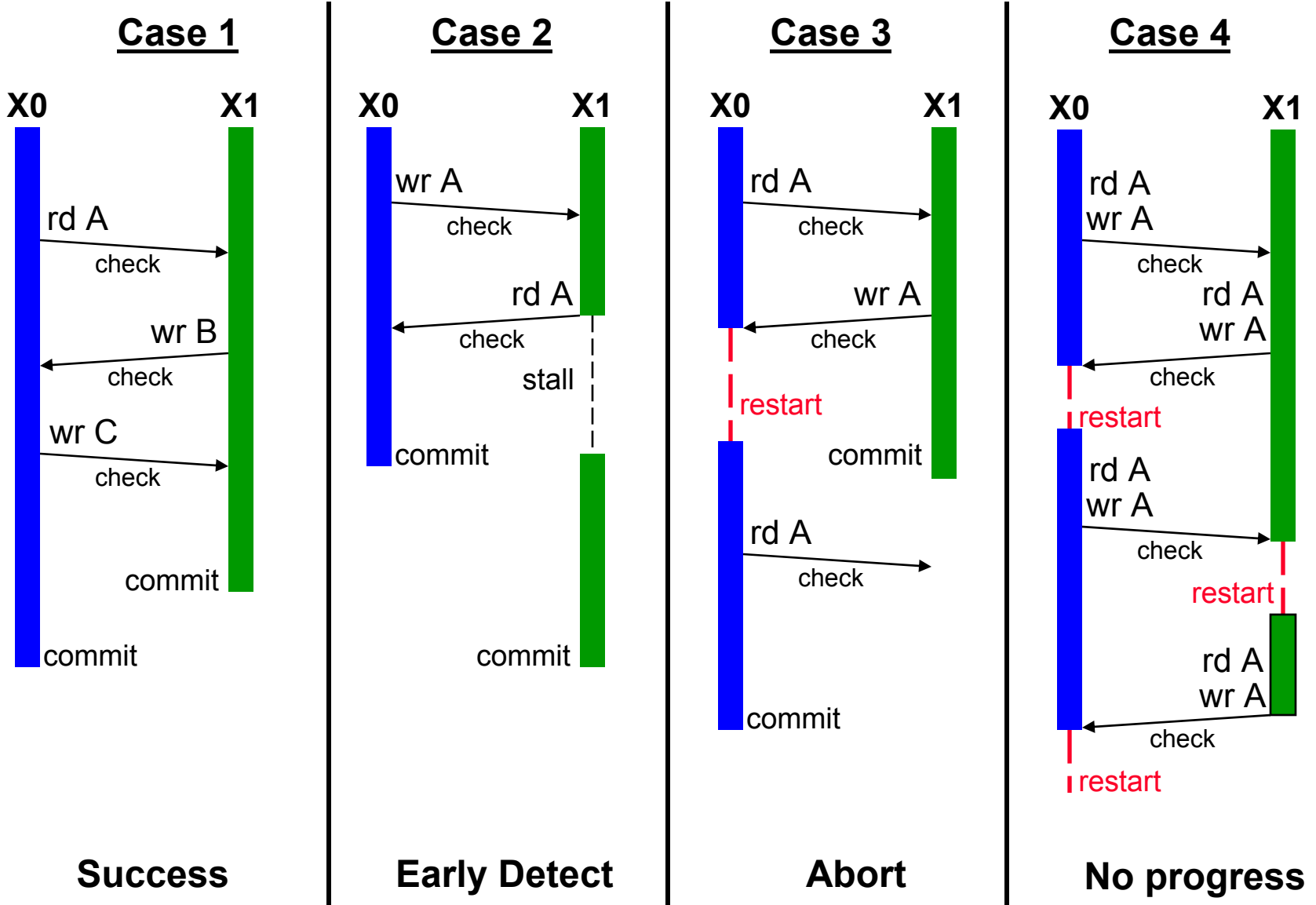
- 2. Lazy or commit or optimistic
 - Detect conflicts when a transaction attempts to commit
 - HW: write-set of committing transaction compared to read-set of others
 - Committing transaction succeeds; others may abort
 - SW: validate write-set and read-set using locks and/or version numbers

- ❑ Can use separate mechanism for loads & stores (SW)



Pessimistic Detection Illustration

TIME ↓

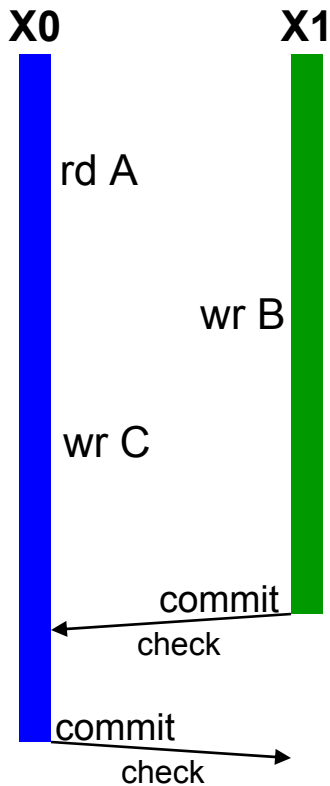




Optimistic Detection Illustration

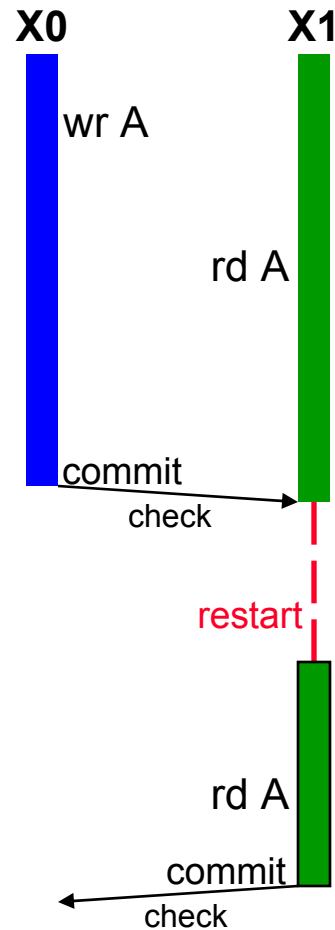
TIME
↓

Case 1



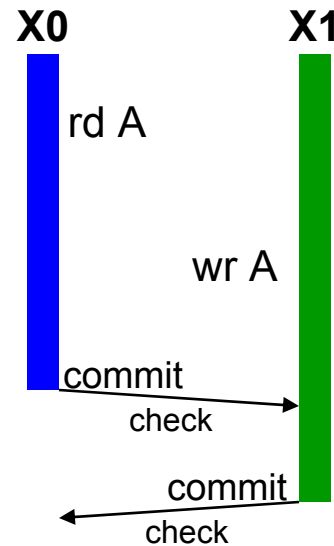
Success

Case 2



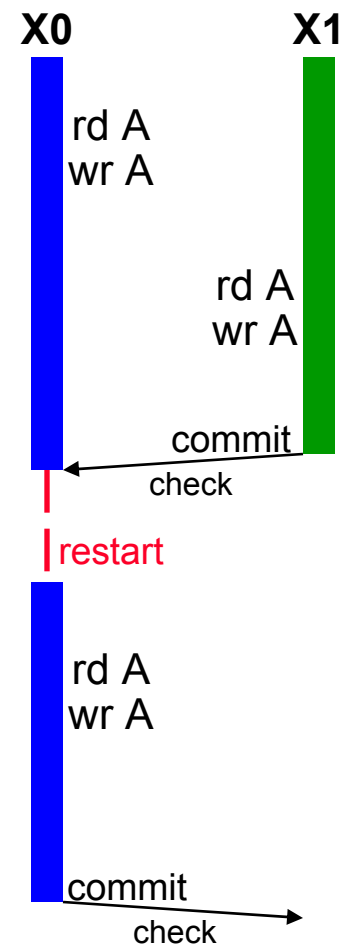
Abort

Case 3



Success

Case 4



Forward progress



Conflict Detection Tradeoffs

1. Eager or encounter or pessimistic

- + Detect conflicts early
 - Lower abort penalty, turn some aborts to stalls
- No forward progress guarantees, more aborts in some cases
- Locking issues (SW), fine-grain communication (HW)

2. Lazy or commit or optimistic

- + Forward progress guarantees
- + Potentially less conflicts, no locking (SW), bulk communication (HW)
- Detects conflicts late



Implementation Space

		Version Management	
		Eager	Lazy
Conflict Detection	Pessimistic	HW: UW LogTM SW: Intel McRT, MS-STM	HW: MIT LTM, Intel VTM SW: MS-OSTM
	Optimistic	HW: -- SW: --	HW: Stanford TCC SW: Sun TL/2

[This is just a subset of proposed implementations]

- No convergence yet
- Decision will depend on
 - Application characteristics
 - Importance of fault tolerance & strong atomicity
 - Success of contention managers, implementation complexity
- May have different approaches for HW, SW, and hybrid



Conflict Detection Granularity

- ❑ Object granularity (SW/hybrid)
 - + Reduced overhead (time/space)
 - + Close to programmer's reasoning
 - False sharing on large objects (e.g. arrays)
 - Unnecessary aborts
- ❑ Word granularity
 - + Minimize false sharing
 - Increased overhead (time/space)
- ❑ Cache line granularity
 - + Compromise between object & word
 - + Works for both HW/SW
- ❑ Mix & match → best of both worlds
 - Word-level for arrays, object-level for other objects, ...



Advanced Implementation Issues

- ❑ Atomicity with respect to non-transactional code
 - Weak atomicity: non-committed transaction state is visible
 - Strong atomicity: non-committed transaction state not visible

- ❑ Nested transactions
 - Common approach: subsume within outermost transaction
 - Recent: nested version management & conflict detection

- ❑ Support for PL & OS design
 - Conditional synchronization, exception handling, ...
 - Key mechanisms: 2-phase commit, commit/abort handlers, open nesting

See paper by McDonald et.al at ISCA'06



HTM: Hardware Transactional Memory Implementations

Christos Kozyrakis

Computer Systems Laboratory
Stanford University

<http://csl.stanford.edu/~christos>



Why Hardware Support for TM

❑ Performance

- Software TM starts with a 40% to 2x overhead handicap

❑ Features

- Works for all binaries and libraries wo/ need to recompile
- Forward progress guarantees
- Strong atomicity
- Word-level conflict detection

❑ How much HW support is needed?

- This is the topic of ongoing research
- All proposed HTMs are essentially hybrid
 - Add flexibility by switching to software on occasion



HTM Implementation Mechanisms

□ Data versioning in caches

- Cache the write-buffer or the undo-log
- Zero overhead for both loads and stores
- Works with private, shared, and multi-level caches

□ Conflict detection through cache coherence protocol

- Coherence lookups detect conflicts between transactions
- Works with snooping & directory coherence

□ Notes

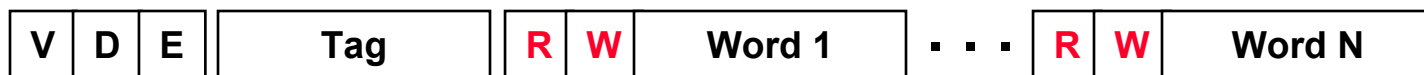
- HTM support similar to that for thread-level speculation (TLS)
 - Some HTMs support both TM and TLS
- Virtualization of hardware resources discussed later



HTM Design

Cache lines annotated to track read-set & write set

- R bit: indicates data read by transaction; set on loads
- W bit: indicates data written by transaction; set on stores
 - R/W bits can be at word or cache-line granularity
- R/W bits gang-cleared on transaction commit or abort
- For eager versioning, need a 2nd cache write for undo log

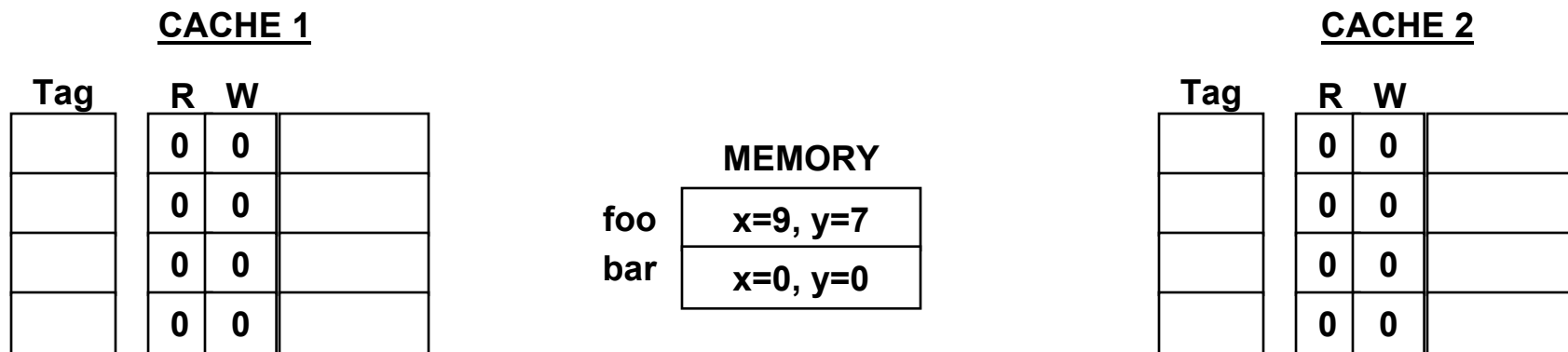


Coherence requests check R/W bits to detect conflicts

- E.g. shared request to W-word is a read-write conflict
- E.g. exclusive request to W-word is a write-write conflict
- E.g. exclusive request to R-word is a write-read conflict



HTM Example



T1 atomic {
 bar.x = foo.x;
 bar.y = foo.y;
}

T2 atomic {
 t1 = bar.x;
 t2 = bar.y;
}

- ❑ T1 copies **foo** into **bar**
- ❑ T2 should read [0, 0] or should read [9,7]
- ❑ Assume HTM system with lazy versioning & optimistic detection



HTM Example (1)

CACHE 1

Tag	R	W	
foo.x	1	0	9
bar.x	0	1	9
	0	0	
	0	0	

MEMORY

foo	x=9, y=7
bar	x=0, y=0

CACHE 2

Tag	R	W	
	0	0	
	0	0	
	0	0	
	0	0	

T1 atomic {
 bar.x = foo.x; ←
 bar.y = foo.y;
 }

T2 atomic { ←
 t1 = bar.x;
 t2 = bar.y;
 }

□ Both transactions make progress independently



HTM Example (2)

CACHE 1

Tag	R	W	
foo.x	1	0	9
bar.x	0	1	9
	0	0	
	0	0	

MEMORY

foo	x=9, y=7
bar	x=0, y=0

CACHE 2

Tag	R	W	
bar.x	1	0	0
t1	0	1	0
	0	0	
	0	0	

T1 atomic {
 bar.x = foo.x; ←
 bar.y = foo.y;
 }

T2 atomic {
 t1 = bar.x; ←
 t2 = bar.y;
 }

□ Both transactions make progress independently



HTM Example (3)

CACHE 1

Tag	R	W	
foo.x	1	0	9
bar.x	0	1	9
foo.y	1	0	7
bar.y	0	1	7

MEMORY

foo	x=9, y=7
bar	x=0, y=0

CACHE 1

Tag	R	W	
bar.x	1	0	0
t1	0	1	0
	0	0	
	0	0	

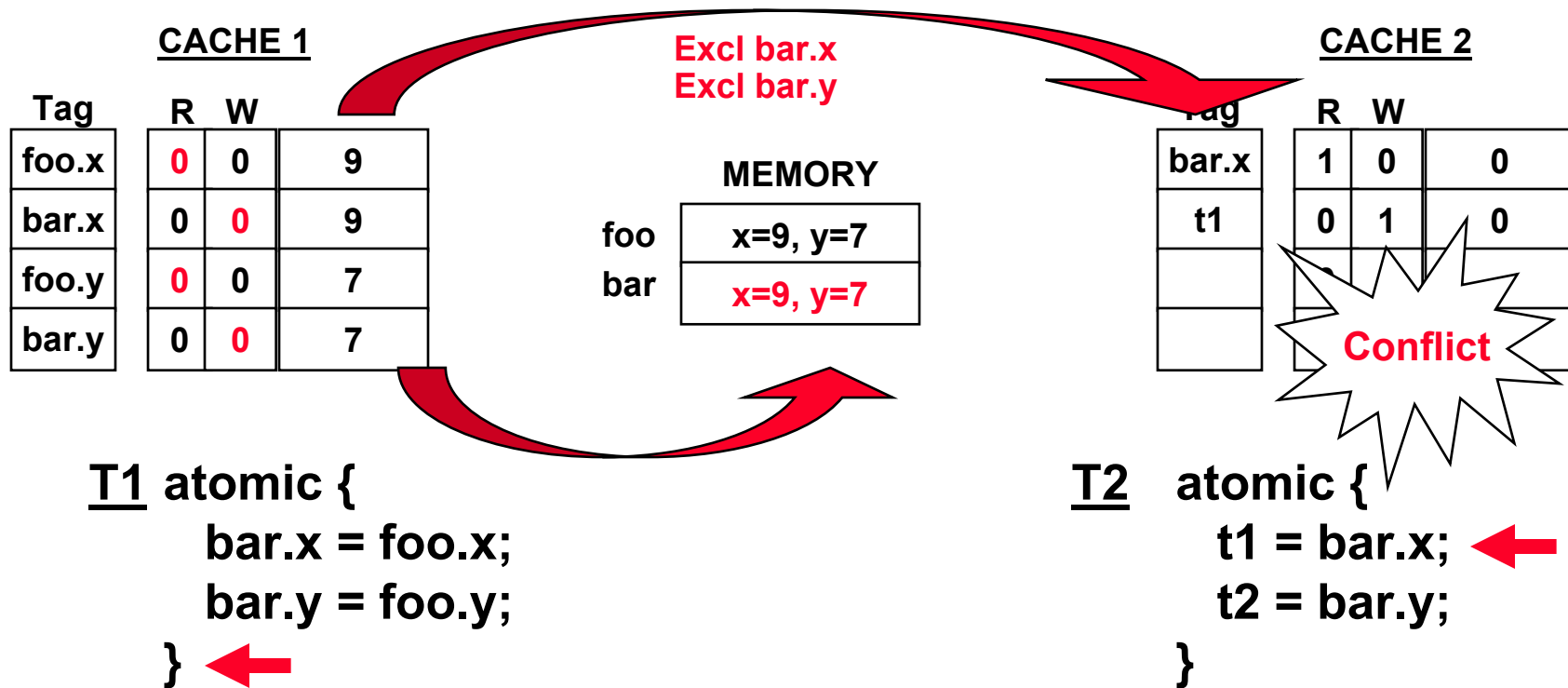
T1 atomic {
 bar.x = foo.x;
 bar.y = foo.y; ←
 }

T2 atomic {
 t1 = bar.x; ←
 t2 = bar.y;
 }

□ Transaction T1 is now ready to commit



HTM Example (3)



T1 updates shared memory

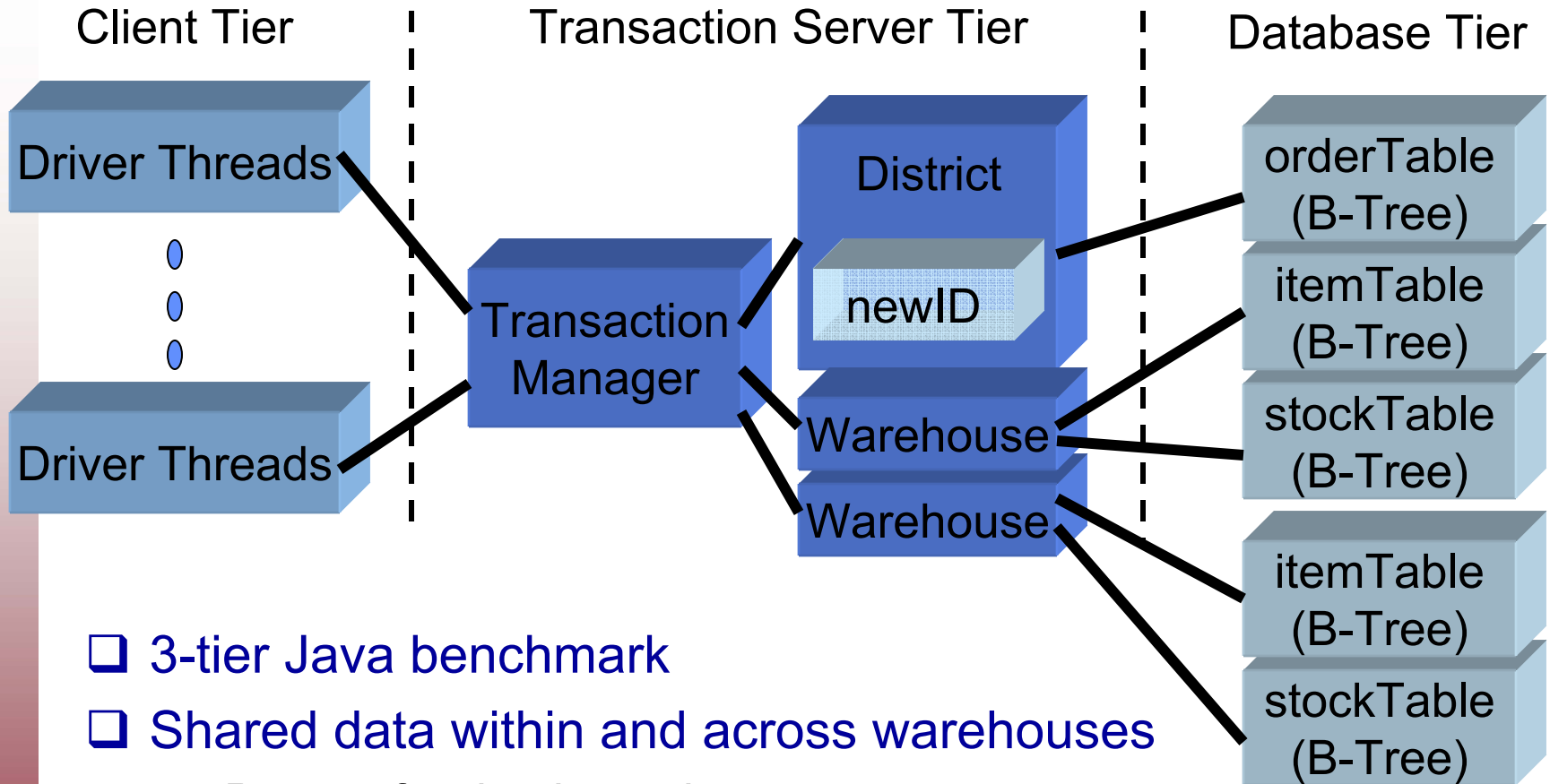
- R/W bits are cleared
- This is a logical update, data may stay in caches as dirty

Exclusive request for bar.x reveals conflict with T2

- T2 is aborted & restarted; all R/W cache lines are invalidated
- When it reexecutes, it will read [9,7] without a conflict



Performance Example: SpecJBB2000



- ❑ 3-tier Java benchmark
- ❑ Shared data within and across warehouses
 - B-trees for database tier
- ❑ Can we parallelize the actions within a warehouse?
 - Orders, payments, delivery updates, etc



Sequential Code for NewOrder

```
TransactionManager::go() {  
    // 1. initialize a new order transaction  
    newOrderTx.init();  
    // 2. create unique order ID  
    orderId = district.nextOrderId(); // newID++  
    order = createOrder(orderId);  
    // 3. retrieve items and stocks from warehouse  
    warehouse = order.getSupplyWarehouse();  
    item = warehouse.retrieveItem(); // B-tree search  
    stock = warehouse.retrieveStock(); // B-tree search  
    // 4. calculate cost and update node in stockTable  
    process(item, stock);  
    // 5. record the order for delivery  
    district.addOrder(order); // B-tree update  
    // 6. print the result of the process  
    newOrderTx.display();  
}
```

❑ Non-trivial code with complex data-structures

- Fine-grain locking → difficult to get right
- Coarse-grain locking → no concurrency



Transactional Code for NewOrder

```
TransactionManager::go() {  
    atomic { // begin transaction  
        // 1. initialize a new order transaction  
        // 2. create a new order with unique order ID  
        // 3. retrieve items and stocks from warehouse  
        // 4. calculate cost and update warehouse  
        // 5. record the order for delivery  
        // 6. print the result of the process  
    } // commit transaction  
}
```

- ❑ Whole NewOrder as one atomic transaction
 - 2 lines of code changed
- ❑ Also tried nested transactional versions
 - To reduce frequency & cost of violations



HTM Performance

❑ Simulated 8-way CMP with TM support

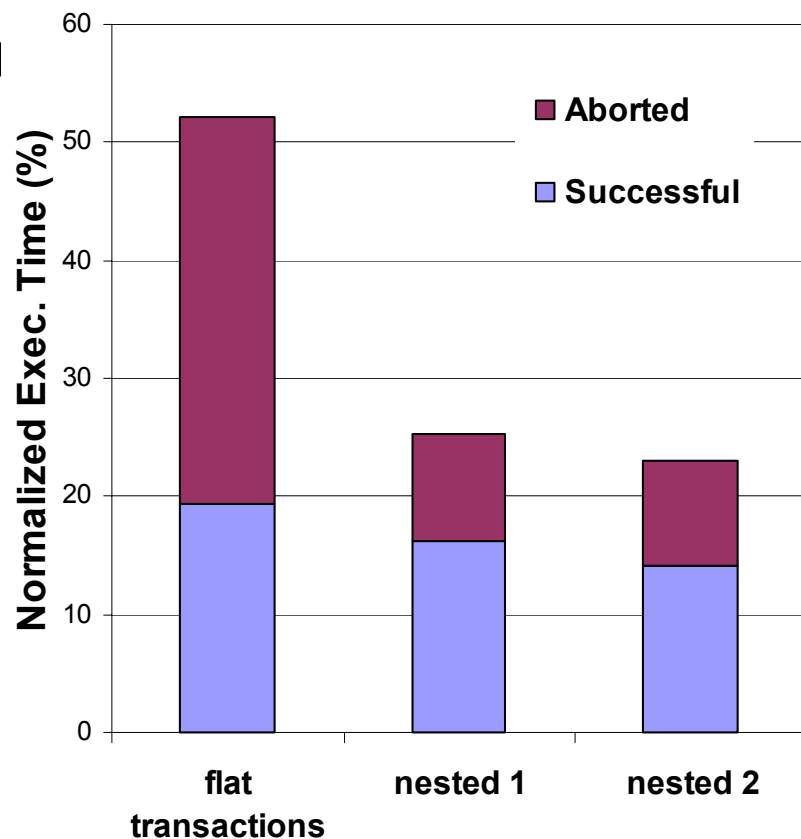
- Stanford's TCC architecture
- Lazy versioning and optimistic conflict detection

❑ Speedup over sequential

- Flat transactions: 1.9x
 - Code similar to coarse-grain locks
 - Frequent aborted transactions due to dependencies
- Nested transactions: 3.9x to 4.2x
 - Reduced abort cost OR
 - Reduced abort frequency

❑ See paper in [WTW'06] for details

- <http://tcc.stanford.edu>





HTM Virtualization (1)

❑ Hardware TM resources are limited

- What if cache overflows? → Space virtualization
- What if time quanta expires? → Time virtualization
- HTM + interrupts, paging, thread migrations, ...

❑ HTM virtualization approaches

1. Dual TM implementation [Intel@PPoPP'06]

- Start transaction in HTM; switch to STM on overflow
- Carefully handle interactions between HTM & STM transactions
- Typically requires 2 versions of the code

2. Hybrid TM [Sun@ASPLOS'06]

- HTM design is optional
- Hash-based techniques to detect interaction between HTM & STM



HTM Virtualization Approaches (cont)

3. Virtualized TM [Intel@ISCA'05]

- Map write-buffer/undo-log and read-/write-set to virtual memory
 - They become unbounded; they can be at any physical location
- Caches capture working set of write-buffer/undo-log
 - Hardware and firmware handle misses, relocation, etc

4. eXtended TM [Stanford@ASPLOS'06]

- Use OS virtualization capabilities (virtual memory)
 - On overflow, use page-based TM → no HW/firmware needed
 - Overflow either all transaction state or just a part of it
- Works well when most transactions are small
 - See common case study at HPCA'06
- Smart interrupt handling
 - Wait for commit Vs. abort transaction Vs. virtualize transaction



Coarse-grain or Bulk HTM Support

❑ Concept

- Track read and write addresses using signatures
 - Bloom filters implemented in hardware
- Process sets of addresses at once using signature operations
 - To manage versioning and to detect conflicts
- Adds 2Kbits per signature, 300 bits compressed

❑ Tradeoffs

- + Conceptually simpler design
 - Decoupled from cache design and coherence protocol
- Inexact operations can lead to false conflicts
 - May lead to degradation
 - Depends on application behavior and signature mechanism

❑ See paper by Ceze et.al at ISCA'06



Transactional Coherence

□ Key observation

- Coherence & consistency only needed at transaction boundaries

□ Transactional Coherence

- Eliminate MESI coherence protocol
- Coherence based on R/W bits
- All coherence communication at commit points

□ Bulk coherence creates hybrid between shared-memory and message passing

□ See TCC papers at [ISCA'04], [ASPLOS'04], & [PACT'05]

```
foo() {  
    work1();  
    atomic {  
        a.x = b.x;  
        a.y = b.y;  
    }  
    work2();  
}
```



Hardware TM Summary

- ❑ High performance + compatibility with binary code,...

- ❑ Common characteristics
 - Data versioning in caches
 - Conflict detection through the coherence protocol

- ❑ Active research area; current research topics
 - Support for PL and OS development (see paper [ISCA'06])
 - Two-phase commit, transactional handlers, nested transactions
 - Development and comparison of various implementations
 - Hybrid TM systems
 - Scalability issues

Agenda

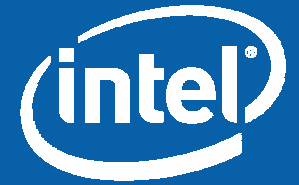
Multithreaded Programming

Transactional Memory (TM)

- TM Introduction
- TM Implementation Overview
- Hardware TM Techniques
- **Software TM Techniques**



Q&A



Software Transactional Memory

Bratin Saha
Programming Systems Lab
Intel Corporation

Outline

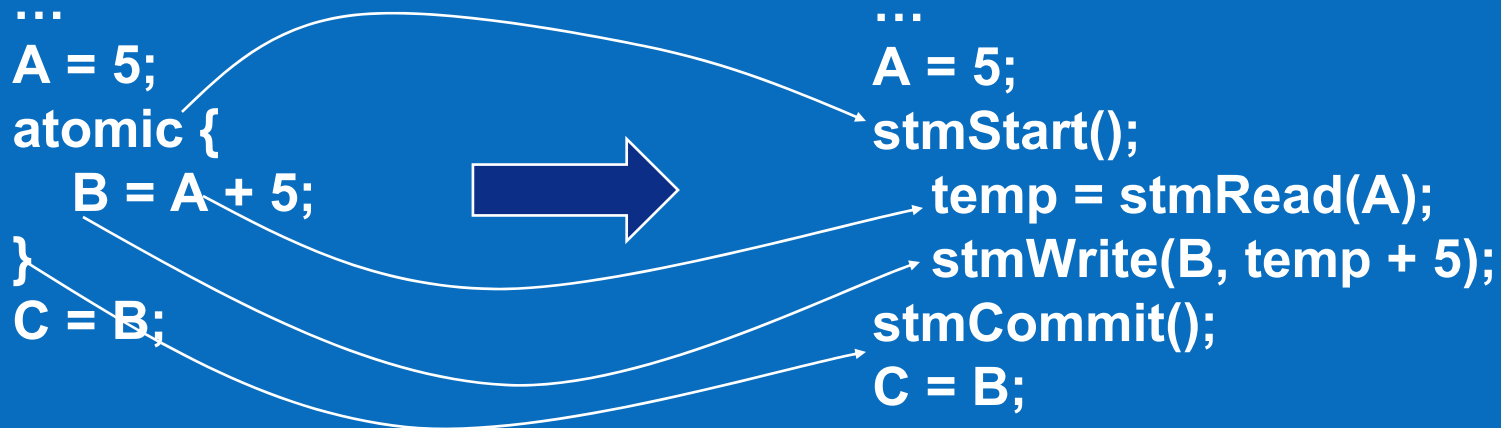
→ Software Transactional Memory

- Translating a language construct
- Runtime support
- Compiler support

→ Hybrid Transactional Memory

→ Open issues & conclusions

Compiling Atomic



Transactional memory accessed via STM read & write functions

- Compiler inserts appropriate calls
- Code generation, control flow, optimizations in later slides

STM tracks accesses & detects data conflicts

Runtime Data Structures

Per-thread

- Transaction Descriptor
 - Read set, write set, & log
 - For validation, commit, & rollback
- Transaction Memento
 - Checkpoint of transaction descriptor
 - For nesting & partial rollback

Per-data

- Transaction Record (TxR)
 - Pointer-sized field guarding shared data
 - Track transactional state of data
 - **Shared**: Read-only access by multiple readers
 - **Exclusive**: write-only access by single owner

Mapping Data to Transaction Records

Every data item has an associated transaction record

Object
granularity
(Java/C#)

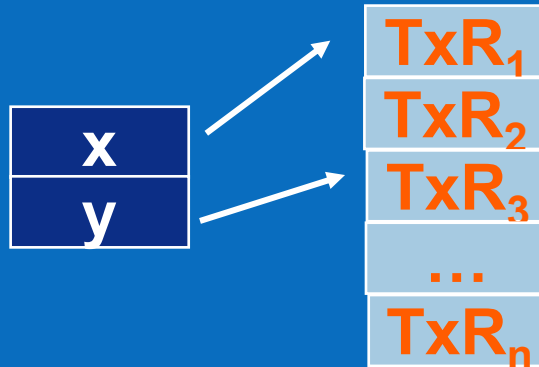
```
class Foo {  
  int x;  
  int y;  
}
```



TxR embedded in object

Cache line
or word
granularity
(C/C++)

```
struct Foo {  
  int x;  
  int y;  
}
```



Address-based hash
into global TxR table

Implementing Atomicity: Example

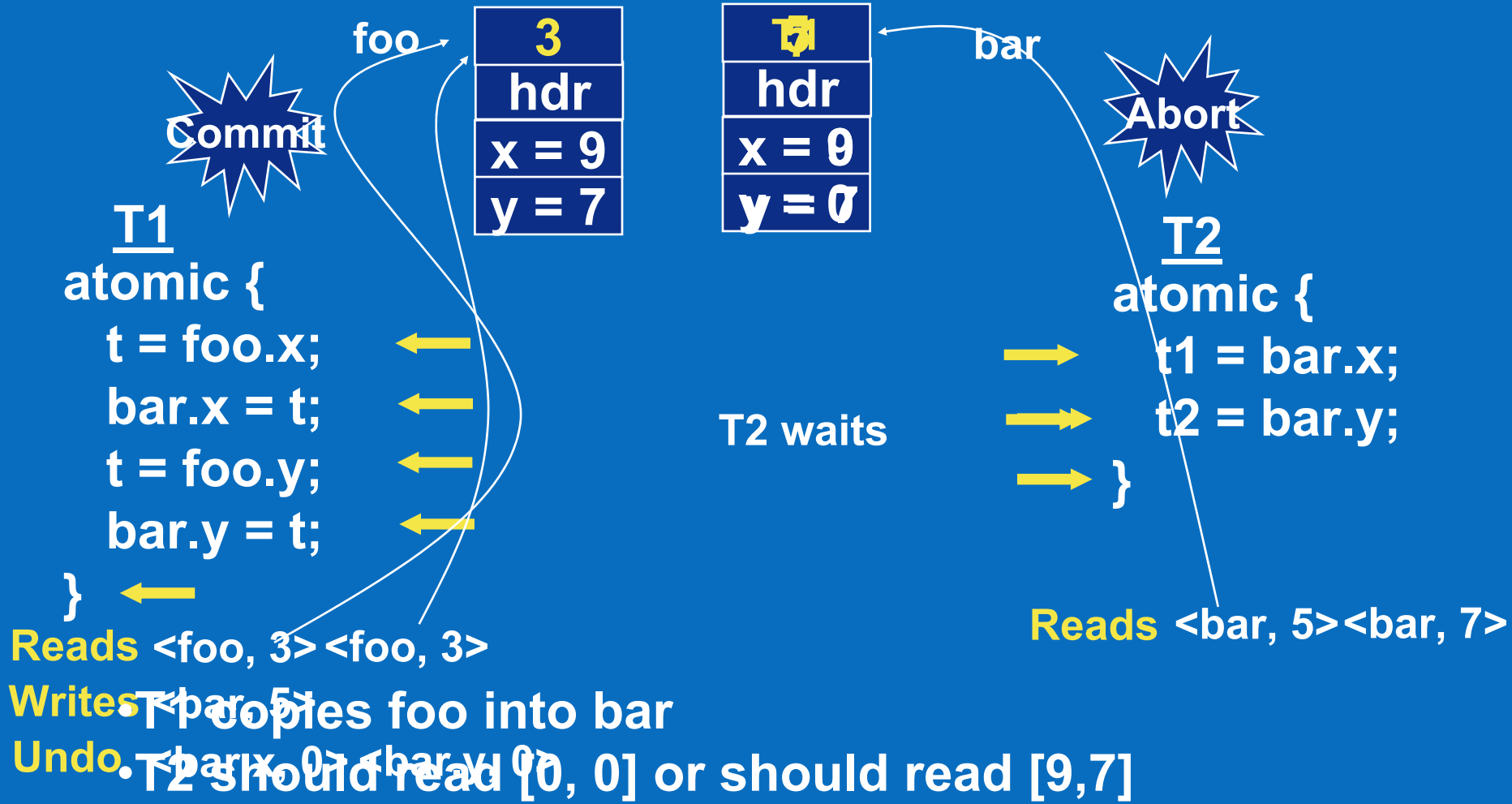
We will show one way to implement atomicity in a STM

Uses two phase locking for writes

Uses optimistic concurrency for reads

Illustrates how transaction records are used

Example



Ensuring Atomicity: Options

Memory Ops → Mode ↓	Reads	Writes
Pessimistic Concurrency	Read lock on TxR (reader-writer lock or reader list)	
Optimistic Concurrency	Use versioning on TxR	



Ensuring Atomicity: Options

Memory Ops → Mode ↓	Reads	Writes
Pessimistic Concurrency	- Caching effects - Lock operations	
Optimistic Concurrency	+ Caching effects + Avoids lock operations	

See Saha et al. PPOPP '06 paper for quantitative results



Ensuring Atomicity: Options

Memory Ops → Mode ↓	Reads	Writes
Pessimistic Concurrency		Write lock on TxR
Optimistic Concurrency		Buffer writes & acquire locks at commit



Ensuring Atomicity: Options

Memory Ops → Mode ↓	Reads	Writes
Pessimistic Concurrency		+ In place updates + Fast commits + Fast reads
Optimistic Concurrency		- Slow commits - Reads have to search for latest value

See Saha et al. PPOPP '06 paper for quantitative results



Java Virtual Machine Support

On-demand cloning of methods called inside transactions

JIT compiler automatically inserts read/write barriers

- Maps barriers to first class opcodes in intermediate representation
- Good compiler representation → greater optimization opportunities
- Determine conflict detection granularity on per-type basis

Garbage collection support

- Enumeration of references in STM data structures
- Filtering to remove redundant log entries
- Mappings are valid across moving GC



Representing Read/Write Barriers

Coarse-grain barriers hide redundant locking/logging

atomic {

 a.x = t1

 a.y = t2

 if(a.z == 0) {

 a.x = 0

 a.z = t3

 }

}

...

stmWr(&a.x, t1)

stmWr(&a.y, t2)

if(**stmRd**(&a.z) != 0) {

stmWr(&a.x, 0);

stmWr(&a.z, t3)

}

An STM IR for Optimization

Redundancies exposed:

```
atomic {  
    a.x = t1  
    a.y = t2  
    if(a.z == 0) {  
        a.x = 0  
        a.z = t3  
    }  
}
```

```
txnOpenForWrite(a)  
txnLogObjectInt(&a.x, a)  
a.x = t1  
txnOpenForWrite(a)  
txnLogObjectInt(&a.y, a)  
a.y = t2  
txnOpenForRead(a)  
if(a.z != 0) {  
    txnOpenForWrite(a)  
    txnLogObjectInt(&a.x, a)  
    a.x = 0  
    txnOpenForWrite(a)  
    txnLogObjectInt(&a.z, a)  
    a.z = t3  
}
```

Optimized Code

Fewer & cheaper STM operations

```
atomic {  
    a.x = t1  
    a.y = t2  
    if(a.z == 0) {  
        a.x = 0  
        a.z = t3  
    }  
}
```

```
txnOpenForWrite(a)  
txnLogObjectInt(&a.x, a)  
a.x = t1  
txnLogObjectInt(&a.y, a)  
a.y = t2  
if(a.z != 0) {  
    a.x = 0  
    txnLogObjectInt(&a.z, a)  
    a.y = t3  
}
```

Compiler Optimizations for Transactions

Standard optimizations

- CSE, Dead-code-elimination, ...
- Careful IR representation exposes opportunities and enables optimizations with almost no modifications
- Subtle in presence of nesting

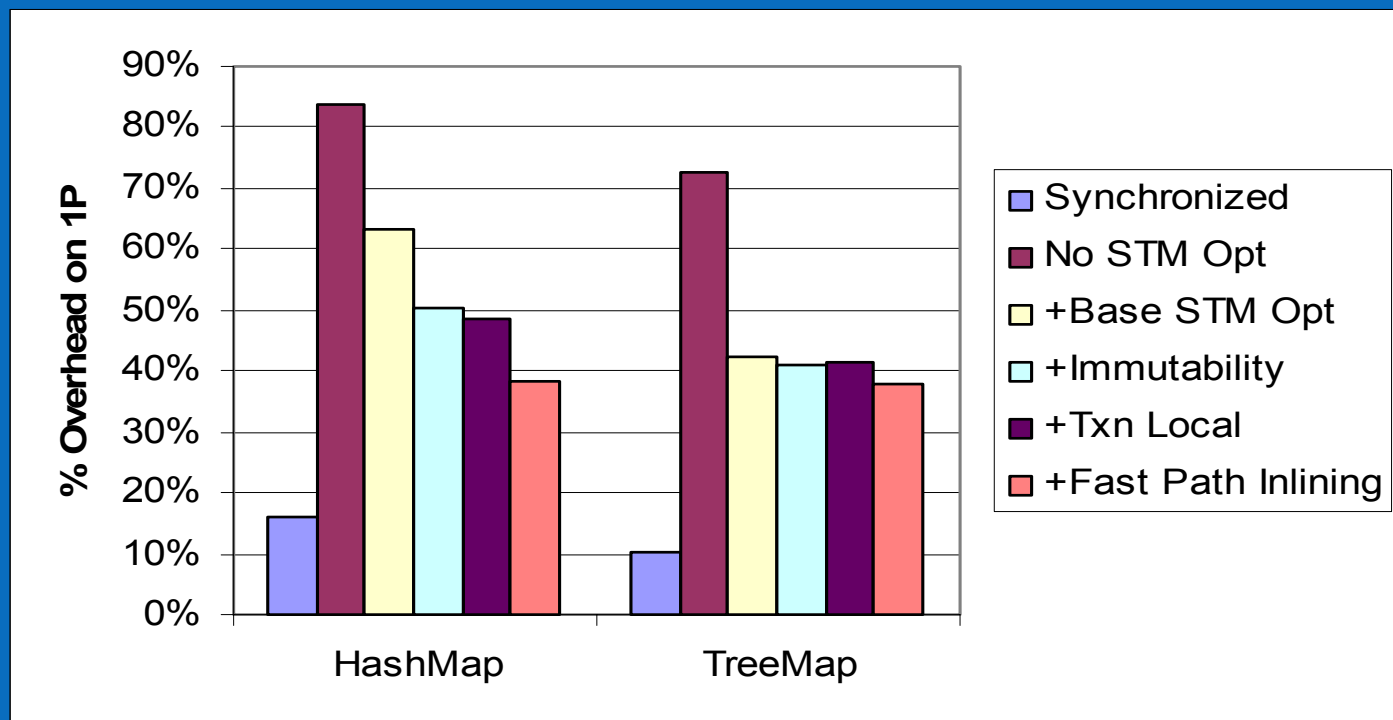
STM-specific optimizations

- Immutable field / class detection & barrier removal (vtable/String)
- Transaction-local object detection & barrier removal
- Partial inlining of STM fast paths to eliminate call overhead



Effect of Compiler Optimizations

1P overheads over thread-unsafe baseline



Prior STMs typically incur $\sim 2x$ on 1P

With compiler optimizations:

- < 40% over no concurrency control
- < 30% over synchronization

Hybrid TM: Combining HTM with STM

General approach:

- Try transaction using HTM first
- Fall back on STM if HTM aborts
- Atomic blocks multiversed for HTM & STM execution

Accelerates simple transaction

- Small
- Flat transactions

STM-STM conflicts detected by the STM machinery

HTM-HTM conflicts detected by the HTM machinery

HTM-STM conflicts requires additional code in the HTM code path



Hybrid TM: Basic Mechanism

HTMReadBarrier(*addr*)

check transaction record for *addr* is not locked by a SW transaction

if (transaction record free)

read the address

else

abort

HTMWriteBarrier(*addr*)

check transaction record for *addr* is not locked by a SW transaction

if (transaction record is free)

perform the write

increment version number to indicate HTM modification

else

abort

HTM check ensures no concurrent SW TM modification

Research challenges

Performance

- Right mix of HW & SW components
- Good diagnostics & contention management

Semantics

- I/O & communication
- Nested parallelism

Debugging & performance analysis tools

System integration



Conclusions

Multi-core architectures: an inflection point in mainstream SW development

Navigating inflection requires new parallel programming abstractions

Transactions are a better synchronization abstraction than locks

- Software engineering and performance benefits

Lots of research on implementation and semantics issues

- Great progress, but there are still open problems

