# Component-Based Lock Allocation

Richard Halpert     Chris Pickett     **Clark Verbrugge**

School of Computer Science, McGill University
{rhalpe,cpicke,clump}@sable.mcgill.ca

6th Workshop on Compiler-Driven Performance
CASCON
October 22nd, 2007

# Lock Allocation

- *Critical section*: piece of code that accesses shared state exclusively
- *Lock*: object that guards access to a critical section

- *Lock allocation*: mapping locks to critical sections

Sounds straightforward, but manual approaches are tricky!

# Race Conditions

```
class T1 extends Thread
{
  public static Object a;

  run ()
  {
    synchronized (T1.a)
      {
        Main.i++;
      }
  }
}
```

```
class T2 extends Thread
{
  public static Object b;

  run ()
  {
    synchronized (T2.b)
      {
        Main.i++;
      }
  }
}
```

race condition! ←——————→

```
class T1 extends Thread                          class T2 extends Thread
{                                                {
  public static Object a;                          public static Object b;

  run ()                                           run ()
  {                                                {
    synchronized (T1.a)                              synchronized (T2.b)
    {                                                {
      synchronized (T2.b)  ←—— deadlock! ——→           synchronized (T1.a)
      {                                                {
        Main.i++;                                          Main.i++;
      }                                                }
    }                                                }
  }                                                }
}                                                }
```

# Performance Degradation

```
class T1 extends Thread
{
  public static Object a;

  run ()
  {
    synchronized (T1.a)
    {
      synchronized (T2.b)
      {
        t1Work();
      }
    }
  }
}
```

```
class T2 extends Thread
{
  public static Object b;

  run ()
  {
    synchronized (T1.a)
    {
      synchronized (T2.b)
      {
        t2Work();
      }
    }
  }
}
```

←  performance
   degradation!  →

Our approach: *automatic* lock allocation

Goal: simplify concurrent programming

- Remove burden of manual allocation from programmer
- Aim to be *strictly* simpler: no extra language constructs
- Ideal result: automatic allocation performance matches or exceeds manual allocation performance

# Contributions

Our contributions:

- We investigate *component-based* lock allocation:
  - Coarse locking granularity
  - Construct a critical section interference graph
  - One lock per graph component
- Experiment with many static compiler analyses
- Show results for small and large Java benchmarks

The technique often performs well:

- Matches manual allocation performance on 2, 4, 8-way hardware for mtrt (SPEC JVM98), lusearch and xalan (DaCapo), and SPEC JBB2005.

# Outline

# Analysis Pipeline

```
class G {
    public static int X, Y;
}

class T1 extends Runnable {
    run() {
        synchronized(...) { // CS1
            G.Y = G.X;
        }
        synchronized(...) { // CS2
            G.X = G.X + 1;
        }
    }
}

class T2 extends Runnable {
    run() {
        synchronized(...) { // CS3
            int a = G.Y;
        }
    }
}
```

Initial Approximation

Interference Identification
    Thread-Local Objects Analysis
    Thread-Based Side Effect Analysis

Interference Pruning
    May Happen in Parallel Analysis

Component-Based Lock Allocation
    Static Locking
    Dynamic Locking

```
class G {
    public static int X, Y;
}

class T1 extends Runnable {
    run() {
        synchronized(...singletonObject...) { // CS1
            G.Y = G.X;
        }
        synchronized(...singletonObject...) { // CS2
            G.X = G.X + 1;
        }
    }
}

class T2 extends Runnable {
    run() {
        synchronized(...singletonObject...) { // CS3
            int a = G.Y;
        }
    }
}
```
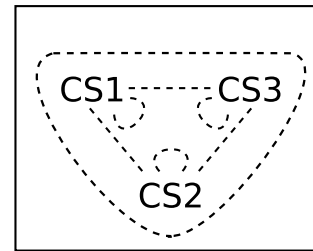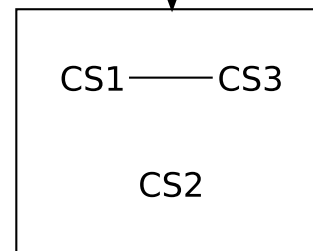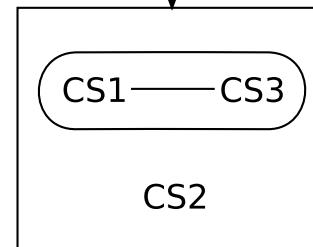
Initial Approximation

CS1 -------- CS3

CS2

Interference Identification
    Thread-Local Objects Analysis
    Thread-Based Side Effect Analysis

CS1 ——— CS3

CS2

Interference Pruning
    May Happen in Parallel Analysis

CS1 ——— CS3

CS2

Component-Based Lock Allocation
    Static Locking
    Dynamic Locking

CS1 ——— CS3

CS2

# Thread-Based Side Effect Analysis

```
class G {
    public static int X, Y;
}

class T1 extends Runnable {
    run() {
        synchronized(...) { // CS1
            G.Y = G.X;          Read from X, thread-shared
        }                       Write to Y, thread-shared
        synchronized(...) { // CS2
            G.X = G.X + 1;  Read from X, thread-shared
        }                   Write to X, thread-shared
    }
}

class T2 extends Runnable {
    run() {
        synchronized(...) { // CS3
            int a = G.Y;    Read from Y, thread-shared
        }
    }
}
```

Initial Approximation

CS1 ----- CS3
   CS2

Interference Identification
    Thread-Local Objects Analysis
    Thread-Based Side Effect Analysis

CS1 ----- CS3
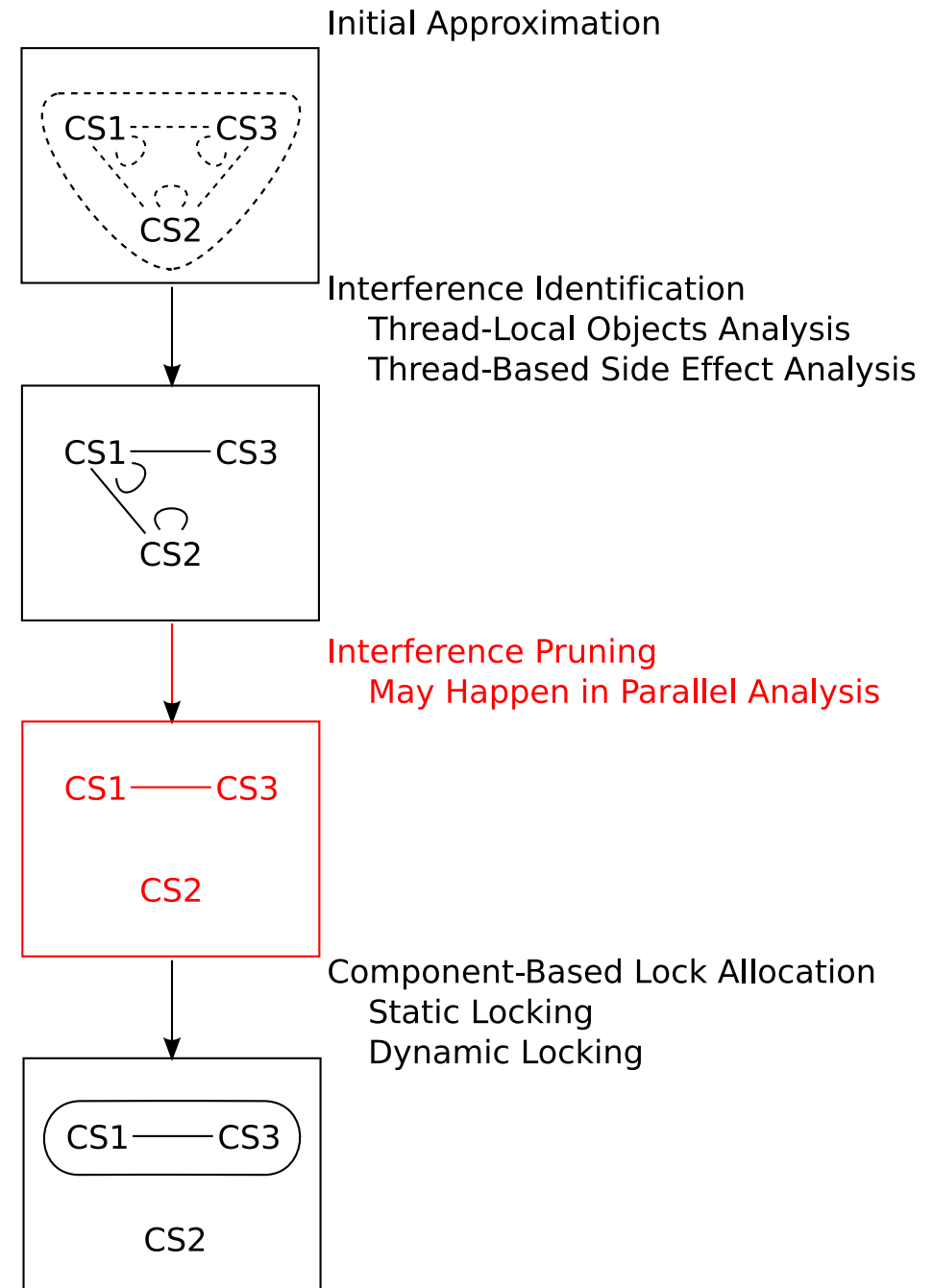      CS2

Interference Pruning
    May Happen in Parallel Analysis

CS1 ----- CS3
      CS2

Component-Based Lock Allocation
    Static Locking
    Dynamic Locking

CS1 ----- CS3
      CS2

# May Happen in Parallel Analysis

Find and apply MHP information
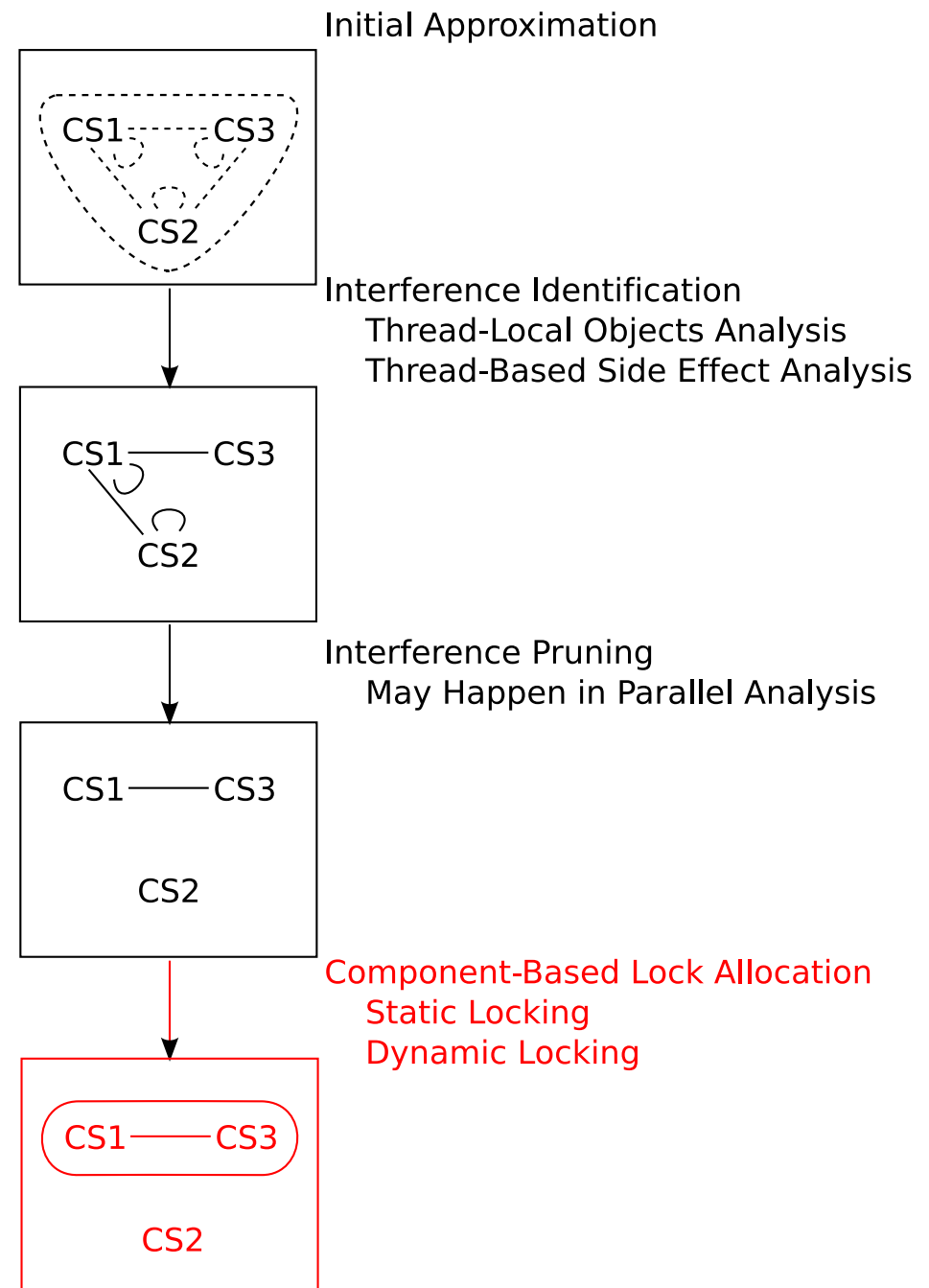
```
class G {
    public static int X, Y;
}

class T1 extends Runnable {
    run() {
        synchronized(...) { // CS1
            G.Y = G.X;
        }
        synchronized(...) { // CS2
            G.X = G.X + 1;
        }
    }
}

class T2 extends Runnable {
    run() {
        synchronized(...) { // CS3
            int a = G.Y;
        }
    }
}
```

Initial Approximation

Interference Identification
    Thread-Local Objects Analysis
    Thread-Based Side Effect Analysis

Interference Pruning
    May Happen in Parallel Analysis

Component-Based Lock Allocation
    Static Locking
    Dynamic Locking

# Component-Based Lock Allocation

Static Lock Allocation:
(Dynamic is the same in this case)

```
class G {
   public static int X, Y;
}

class T1 extends Runnable {
   run() {
      synchronized(LockObject1) { // CS1
         G.Y = G.X;
      }
      synchronized(...) { // CS2
         G.X = G.X + 1;
      }
   }
}

class T2 extends Runnable {
   run() {
      synchronized(LockObject1) { // CS3
         int a = G.Y;
      }
   }
   public static Object LockObject1 =
      new Object();
}
```

Initial Approximation



Interference Identification
    Thread-Local Objects Analysis
    Thread-Based Side Effect Analysis

Interference Pruning
    May Happen in Parallel Analysis

Component-Based Lock Allocation
    Static Locking
    Dynamic Locking

Build on an existing side-effect analysis

- Identify fields that are read & written
- Each field has a points-to set of possible base objects

Extend it to be thread-sensitive

- Approximate the thread-visible effects of library calls
- Exclude thread-local side effects

Use it to construct a critical section *interference graph*
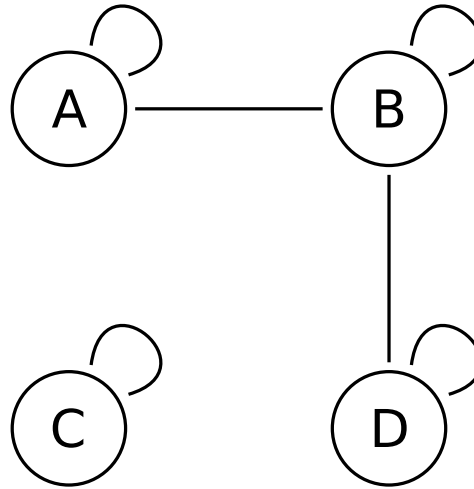
# Constructing an Interference Graph

```
class A {                        class B {
  public static int f;             public static int f;

  synchronized void a() {          synchronized void b() {
    A.f = B.f + 1;                   B.f = B.f + D.f;
  }                                }
}                                }



class C {                        class D {
  public static int f;             public static int f;

  synchronized void c() {          synchronized void d() {
    C.f = C.f + 1;                  D.f = D.f + 1;
  }                                }
}                                }
```

# Constructing an Interference Graph



Interference Graph

|   | A | B | C | D |
|---|---|---|---|---|
| A | 1 | 1 | 0 | 0 |
| B | 1 | 1 | 0 | 1 |
| C | 0 | 0 | 1 | 0 |
| D | 0 | 1 | 0 | 1 |

# Finding Thread-Local Objects

*Thread-local object*: object only read & written by a single thread
Similar to escape analysis

- Partition the heap into thread-shared and thread-local data
- Use information flow analysis to propagate thread-shared status
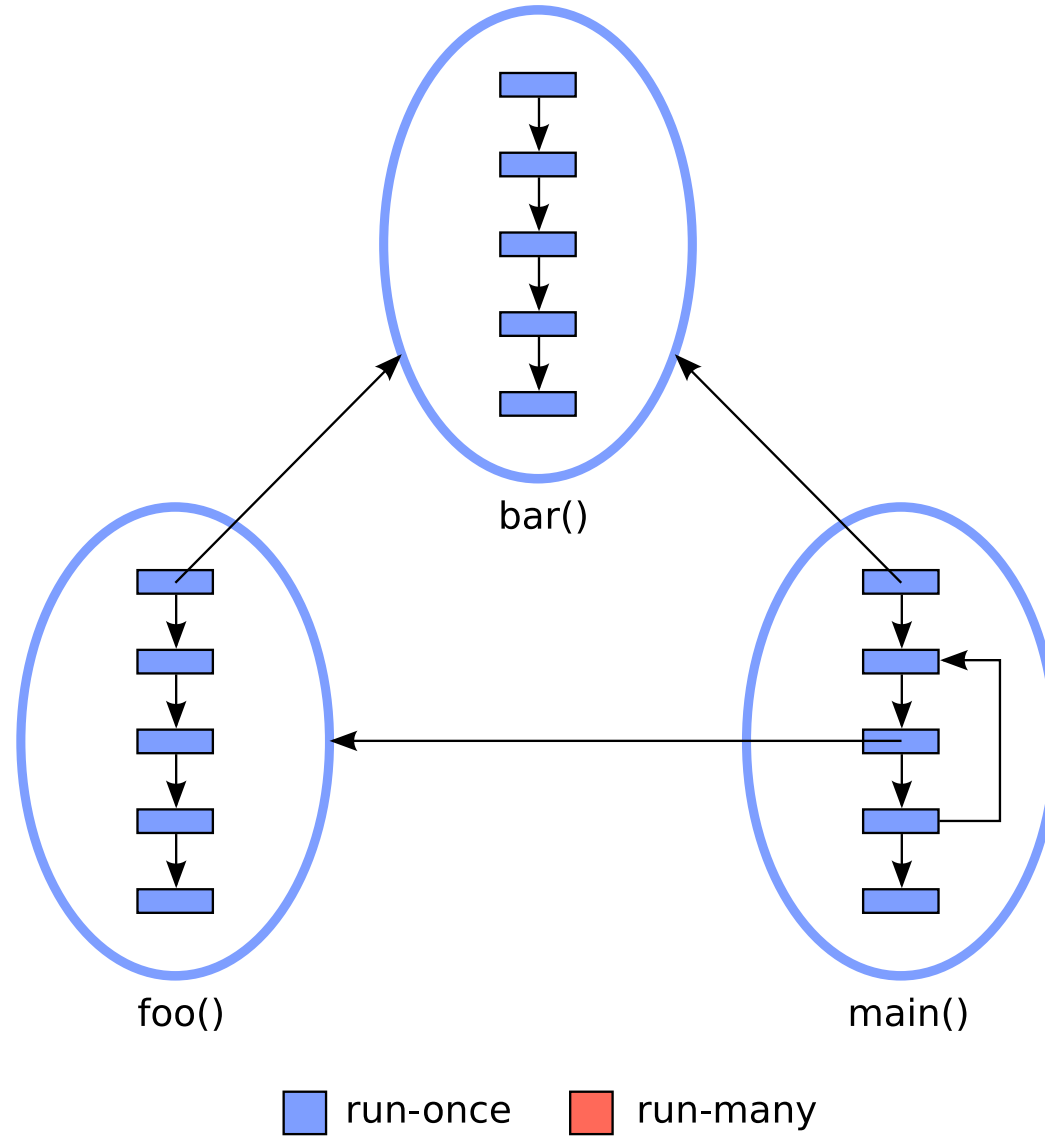
Values identified as thread-local do not require synchronized access

- MHP analysis finds methods that execute concurrently
- Several distinct steps:
    1. Identify run-once and run-many statements
    2. Identify run-once and run-many threads
    3. Categorize run-many threads as run-one-at-a-time or run-many-at-time
    4. Find methods that may happen in parallel based on thread reachability
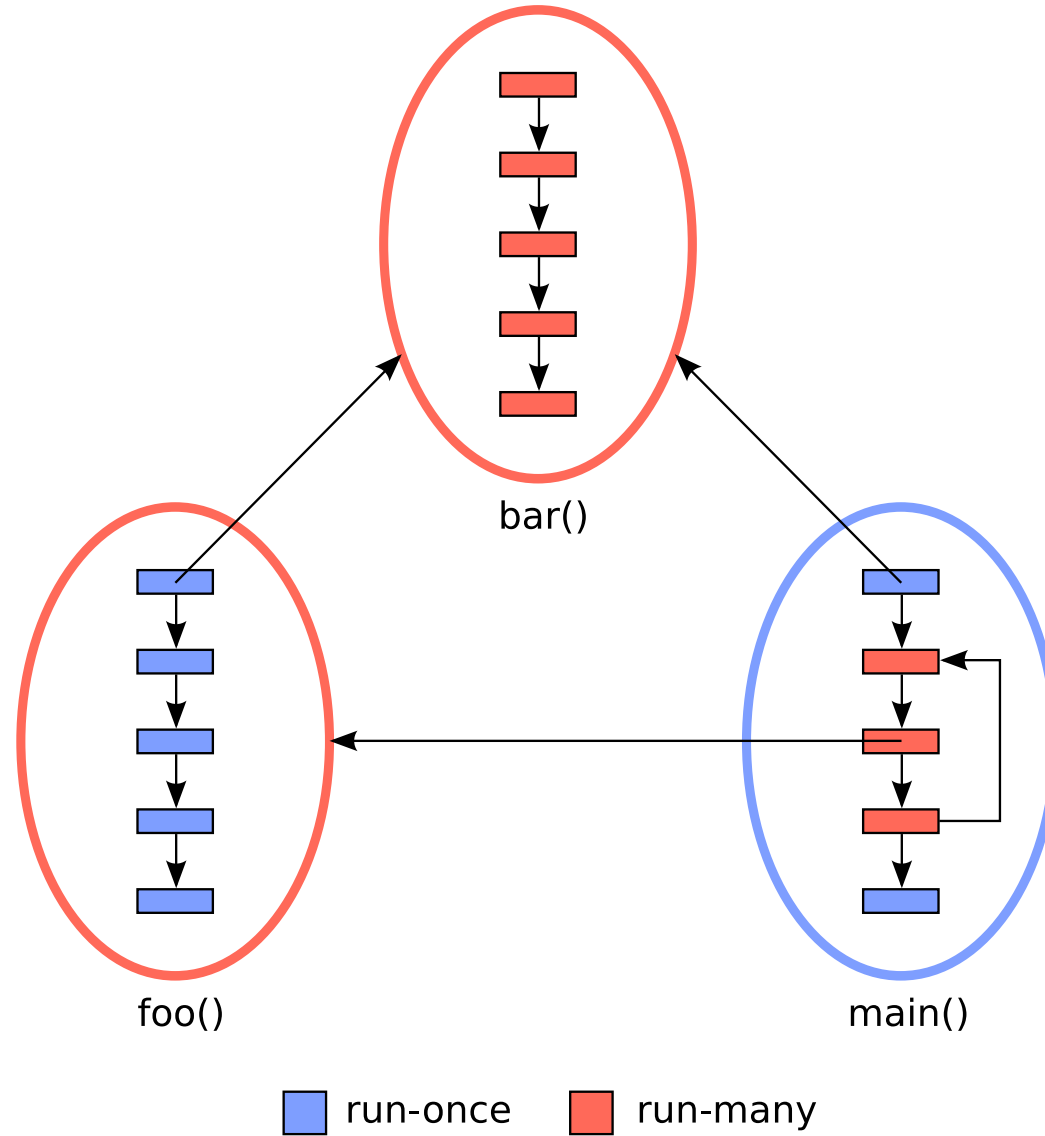- Critical sections that may not happen in parallel cannot interfere!
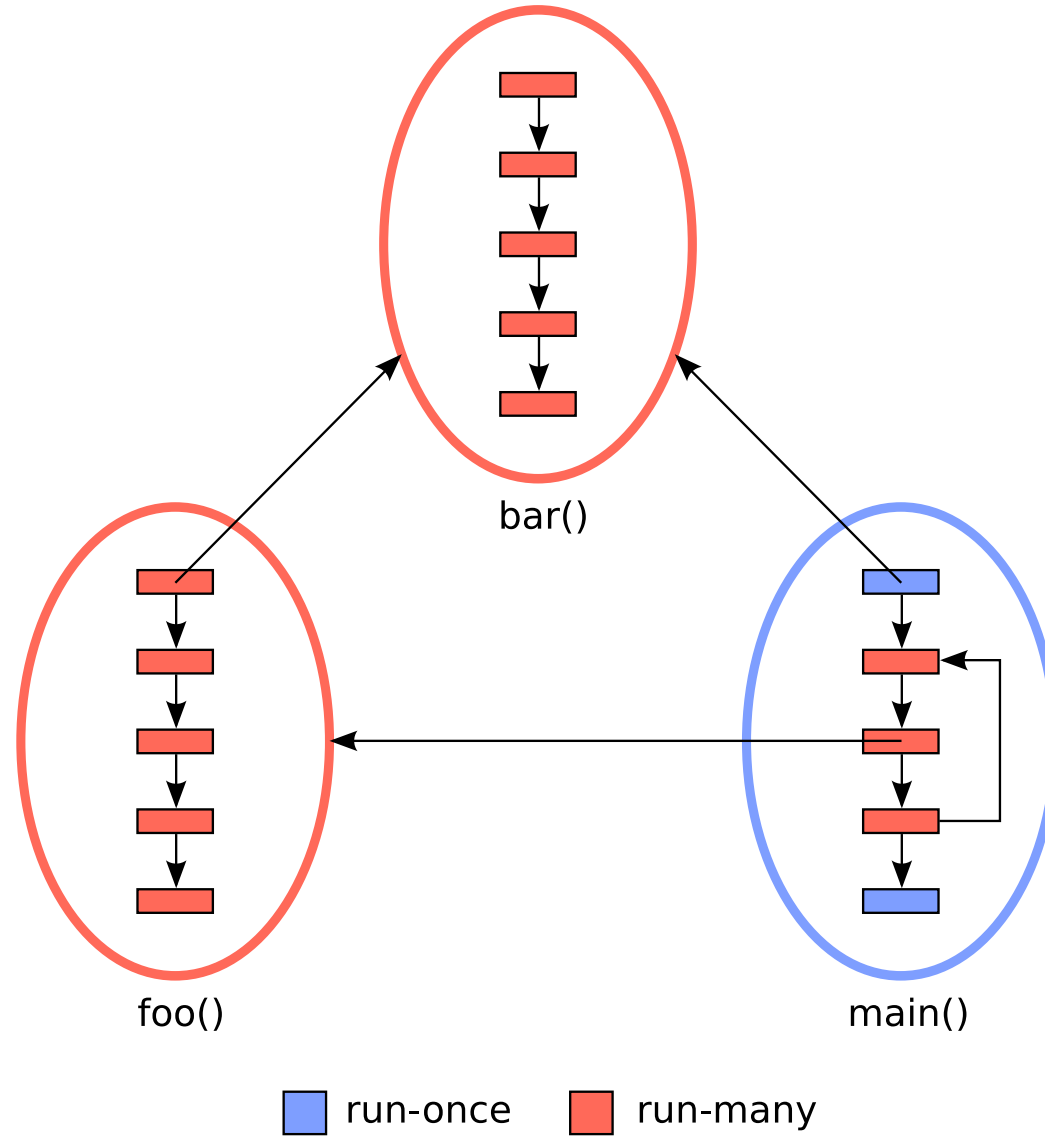
foo()

bar()

main()

■ run-once    ■ run-many

bar()

foo()

main()

run-once    run-many

bar()

foo()                                    main()

run-once          run-many

# Run-Once Run-Many Analysis



bar()

foo()

main()

run-once   run-many

# Run-Once Run-Many Analysis

bar()

foo()

main()

run-once   run-many

# Thread Categorization

Thread t1, t2, t3

int i

☐ run-once

☐ run-many

T1: ?

T2: ?

T3: ?

t1 = new T1()

t1.start()

t2 = new T2()

i = 0

if (i > 0)

F        T

t3 = new T3()

t2.start()

t3.start()

i = i + 1

if (i < 10)    T

F

Thread t1, t2, t3

int i

run-once

run-many

T1: ?

T2: ?

T3: ?

t1 = new T1()

t1.start()

t2 = new T2()

i = 0

if (i > 0)

F        T

t3 = new T3()

t2.start()

t3.start()

i = i + 1

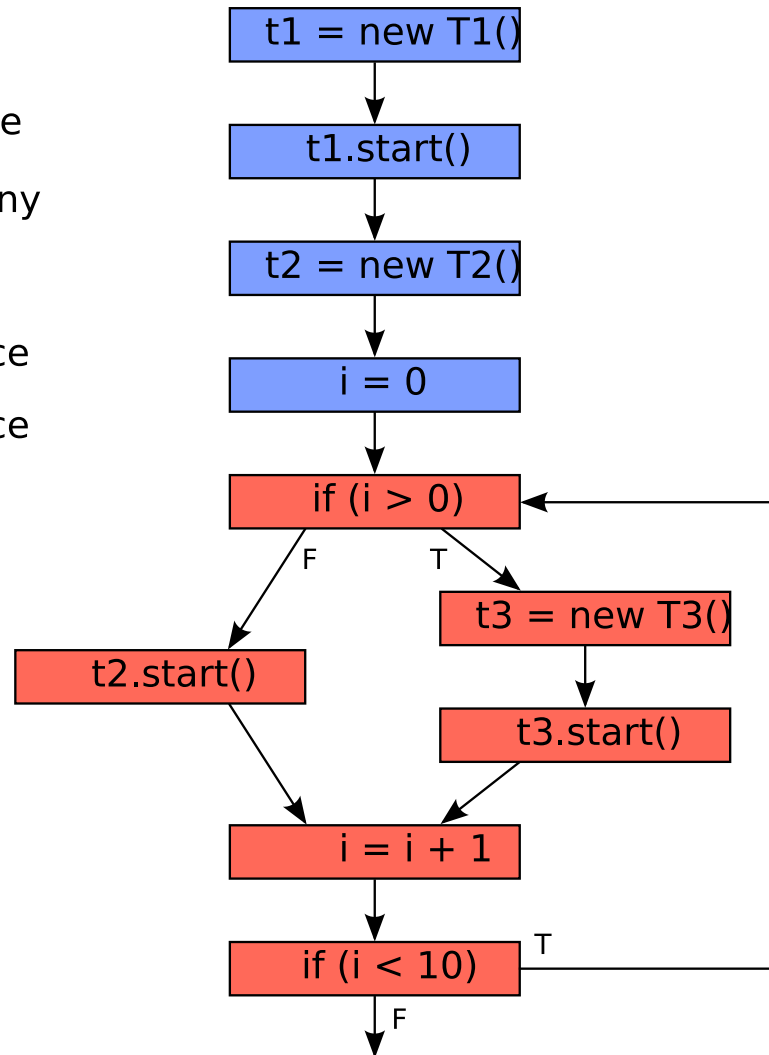if (i < 10)        T

F

Thread t1, t2, t3

int i

☐ run-once

☐ run-many

T1: run-once

T2: ?

T3: ?

# Thread Categorization

Thread t1, t2, t3

int i

■ run-once
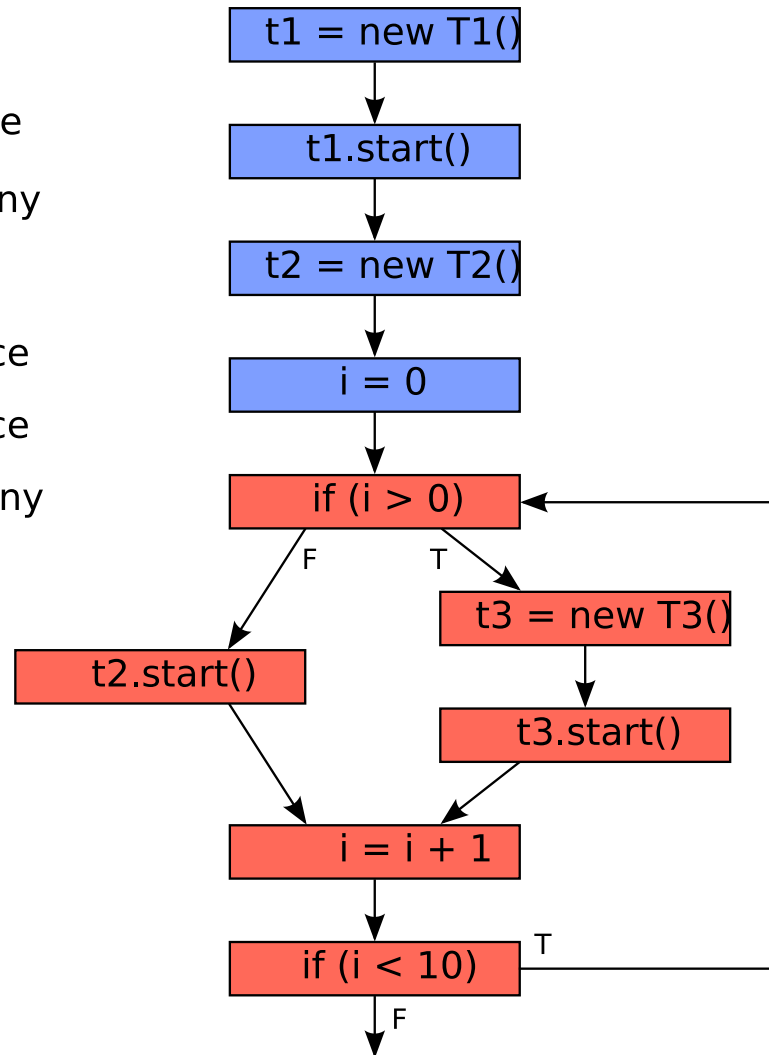
■ run-many

T1: run-once

T2: run-once

T3: ?

t1 = new T1()

t1.start()

t2 = new T2()

i = 0

if (i > 0)

F          T

t3 = new T3()

t2.start()

t3.start()

i = i + 1

if (i < 10)          T

F

# Thread Categorization

Thread t1, t2, t3

int i

☐ run-once

☐ run-many

T1: run-once

T2: run-once

T3: run-many

t1 = new T1()

t1.start()

t2 = new T2()

i = 0

if (i > 0)

F    T

t3 = new T3()

t2.start()

t3.start()

i = i + 1

if (i < 10)    T

F

# Finding run-one-at-a-time threads

Thread t1, t2, t3

int i

■ run-once

■ run-many

T1: run-once

T2: run-once

T3: ~~run-many~~
    one-at-a-time
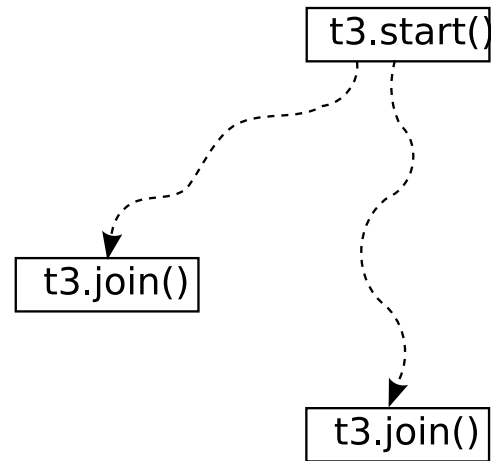
```
t1 = new T1()
      ↓
t1.start()
      ↓
t2 = new T2()
      ↓
i = 0
      ↓
if (i > 0)  ←──────────────┐
   F  ↓   T ↓              │
      │      t3 = new T3() │
      │          ↓         │
t2.start()   t3.start()    │
      │          ↓         │
      │      t3.join()     │
      ↓          ↓         │
   i = i + 1               │
      ↓                    │
   if (i < 10)  ── T ──────┘
      ↓ F
```

t3.start()

tA.join()

tC.join()

tB.join()

- For each start, consider all joins:

t3.start()

t3.join()

t3.join()
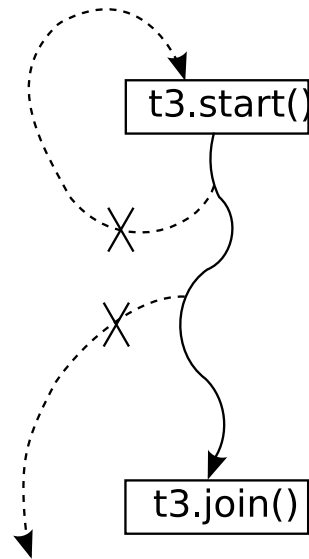
- For each start, consider all joins:
  - Any valid join receiver must alias start receiver

t3.start()

t3.join()

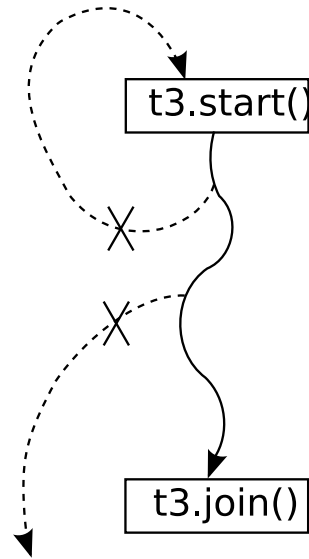- For each start, consider all joins:
  - Any valid join receiver must alias start receiver
  - Any valid join must post-dominate start
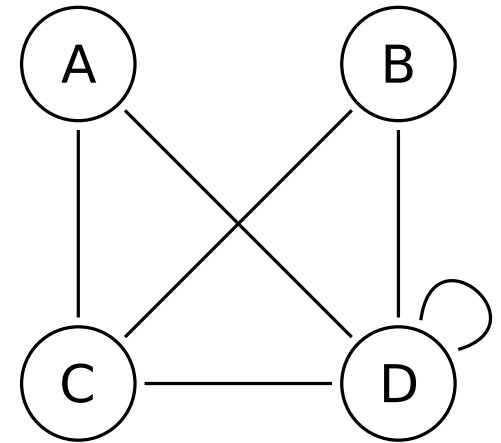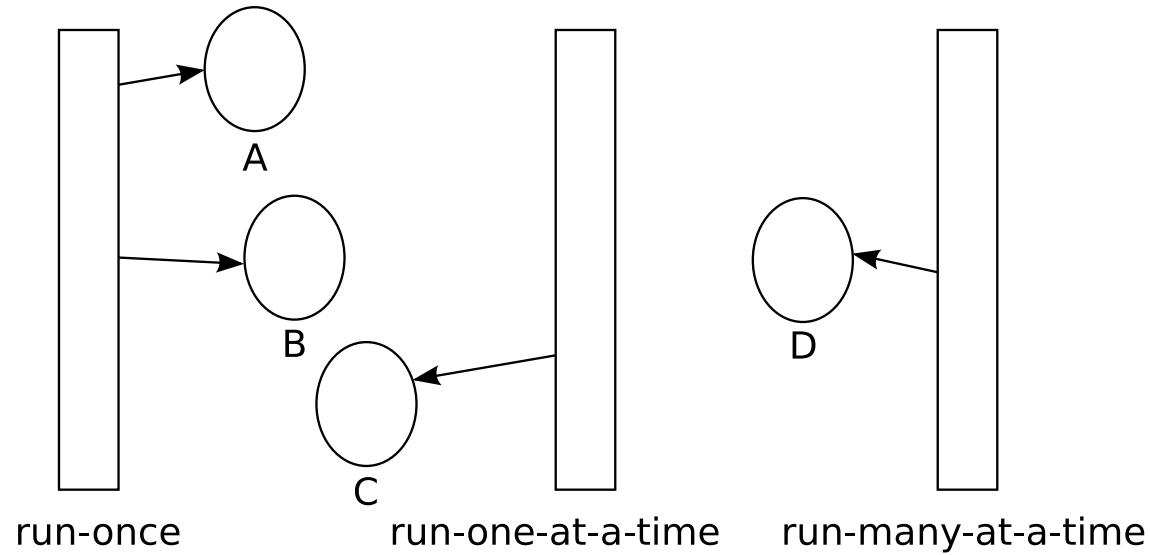
t3.start()

t3.join()

- For each start, consider all joins:
  - Any valid join receiver must alias start receiver
  - Any valid join must post-dominate start
  - And not have loops to start between the start and join...

```
        ┌──────────┐
        │ t3.start()│
        └──────────┘
              ✗
              ✗
        ┌──────────┐
        │ t3.join() │
        └──────────┘
```

- For each start, consider all joins:
  - Any valid join receiver must alias start receiver
  - Any valid join must post-dominate start
  - And not have loops to start between the start and join...
- If join is valid, check method validity:
  - Method must not be called recursively
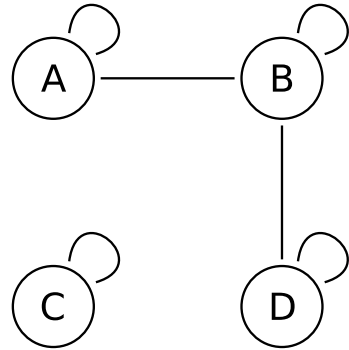  - Method must not happen in parallel with itself

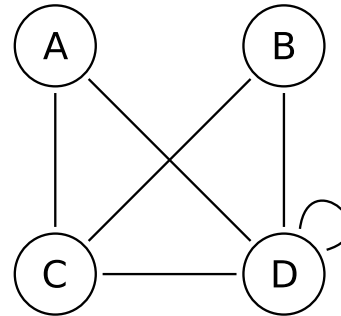run-once      run-one-at-a-time      run-many-at-a-time      MHP Information

|   | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 0 | 1 | 1 |
| B | 0 | 0 | 1 | 1 |
| C | 1 | 1 | 0 | 1 |
| D | 1 | 1 | 1 | 1 |

Interference Graph

MHP Information

Pruned Interference Graph

|   | A | B | C | D |
|---|---|---|---|---|
| A | 1 | 1 | 0 | 0 |
| B | 1 | 1 | 0 | 1 |
| C | 0 | 0 | 1 | 0 |
| D | 0 | 1 | 0 | 1 |

●

|   | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 0 | 1 | 1 |
| B | 0 | 0 | 1 | 1 |
| C | 1 | 1 | 0 | 1 |
| D | 1 | 1 | 1 | 1 |

=

|   | A | B | C | D |
|---|---|---|---|---|
| A |   |   |   |   |
| B |   |   |   |   |
| C |   |   |   |   |
| D |   |   |   |   |

A simple Hadamard product

Interference Graph

MHP Information

Pruned Interference Graph

|   | A | B | C | D |
|---|---|---|---|---|
| A | 1 | 1 | 0 | 0 |
| B | 1 | 1 | 0 | 1 |
| C | 0 | 0 | 1 | 0 |
| D | 0 | 1 | 0 | 1 |

$\bullet$

|   | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 0 | 1 | 1 |
| B | 0 | 0 | 1 | 1 |
| C | 1 | 1 | 0 | 1 |
| D | 1 | 1 | 1 | 1 |

$=$

|   | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 0 | 0 | 0 |
| B | 0 | 0 | 0 | 1 |
| C | 0 | 0 | 0 | 0 |
| D | 0 | 1 | 0 | 1 |

A simple Hadamard product

Three kinds of component-based lock allocation:

1. Singleton: a single static lock protects all components
2. Static: one static lock per component
3. Dynamic: attempt to use per-data structure locks for each component, otherwise static

Finally, isolated vertices with no self loops are *unlocked*

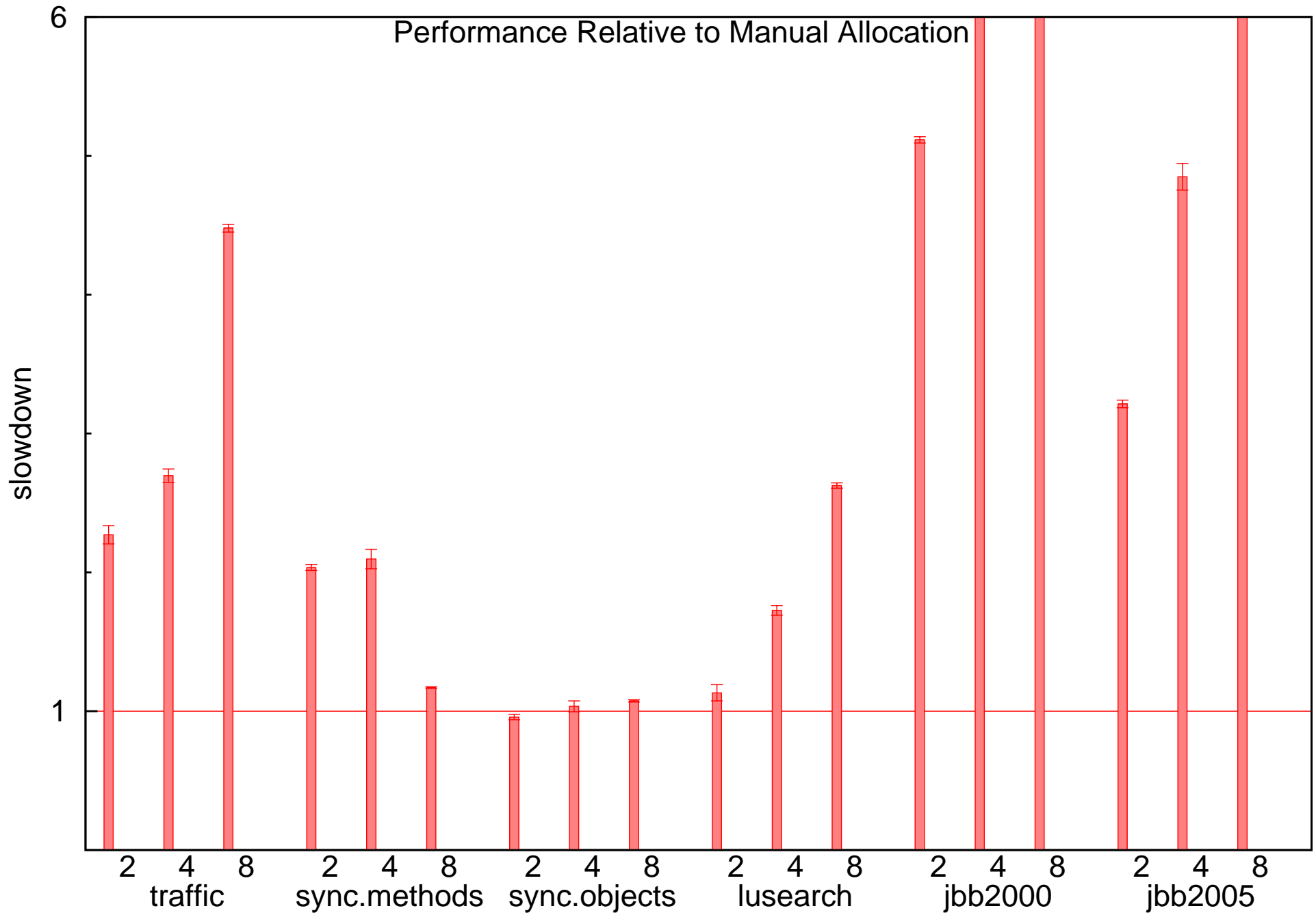# Outline

# Experimental Setup

For each benchmark, we do 13 experiments:

- control: original benchmark program
- singleton: single static lock for all critical sections
- 5 static locking allocations:
    1. CHA: class hierarchy analysis points-to and side effects
    2. Spark: context-insensitive points-to and side effects
    3. Spark-MHP: Spark with may happen in parallel [MHP] analysis
    4. Spark-TLO-MHP: Spark with both TLO and MHP
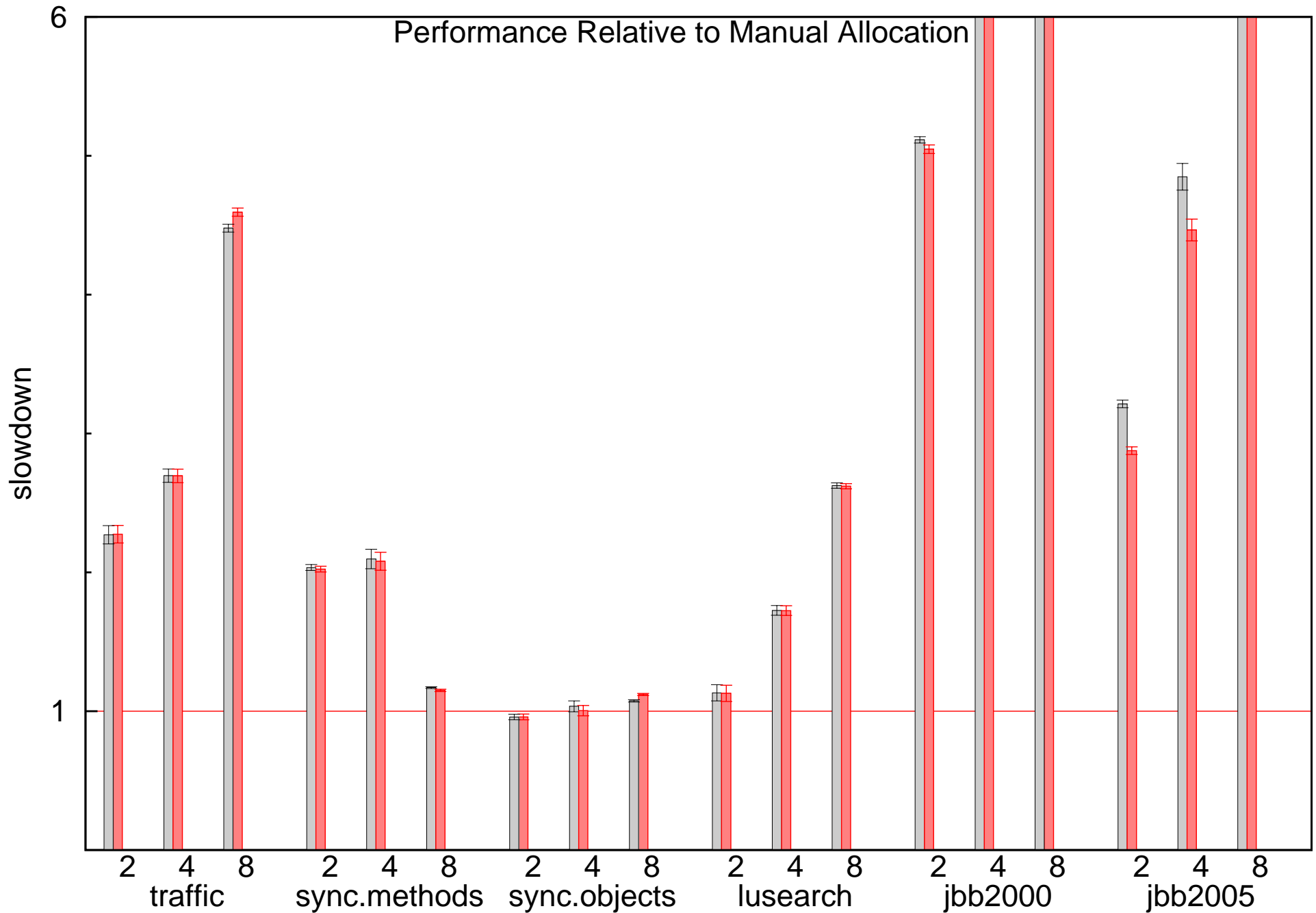- 5 analogous dynamic locking allocations

11 benchmarks: 5 micro, 6 standard
64-bit AMD Machines (dual, 4-way, 4-way dual), Sun JDK1.5

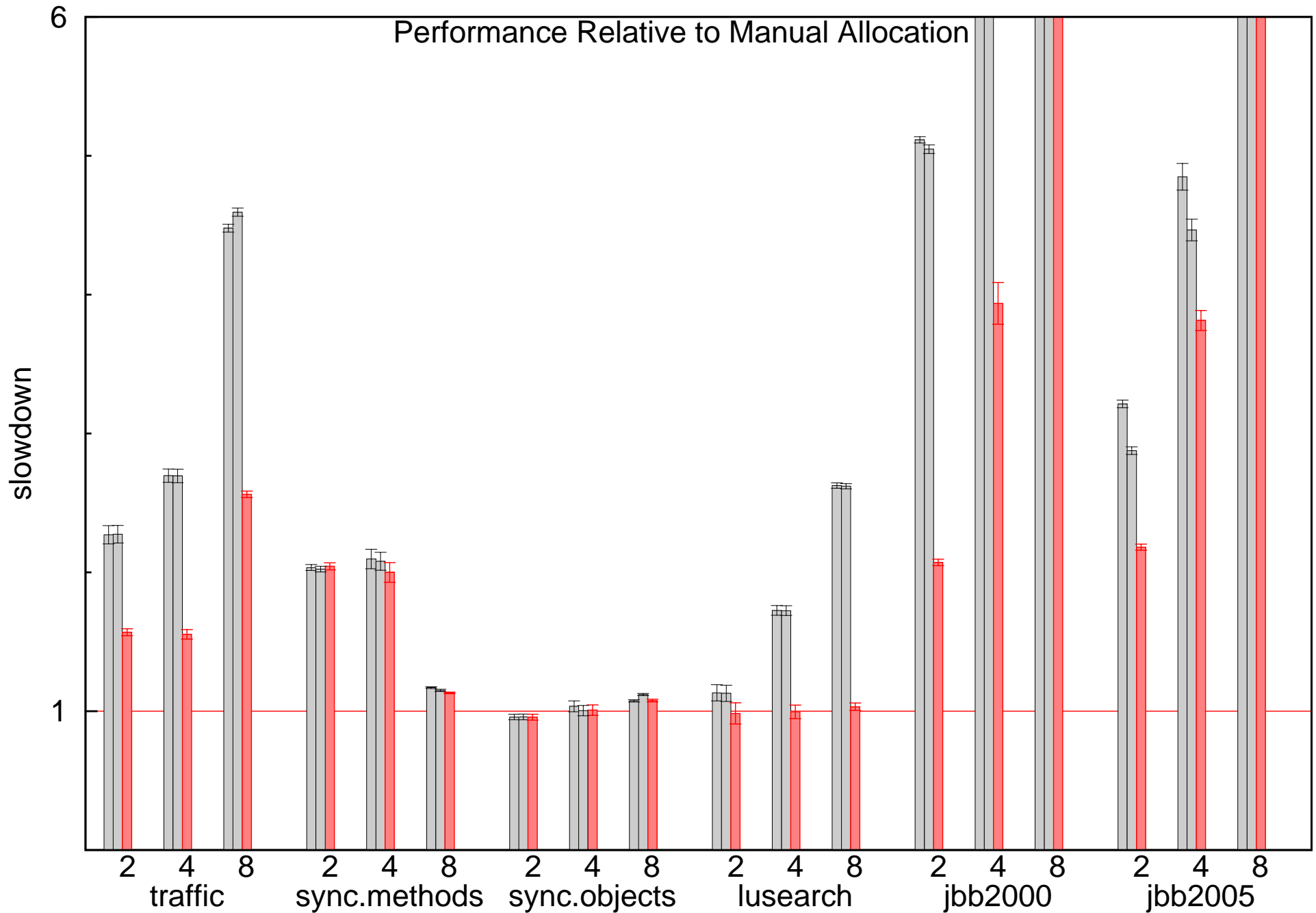# Relative Speedup of Using CHA



Performance Relative to Manual Allocation

slowdown

6

1

2  4  8        2  4  8        2  4  8        2  4  8        2  4  8        2  4  8

traffic      sync.methods   sync.objects    lusearch       jbb2000        jbb2005

Performance Relative to Manual Allocation

slowdown

6

1

2 4 8 | 2 4 8 | 2 4 8 | 2 4 8 | 2 4 8 | 2 4 8

traffic    sync.methods    sync.objects    lusearch    jbb2000    jbb2005

# Relative Speedup of Adding MHP Analysis



Performance Relative to Manual Allocation

slowdown

6

1

traffic    sync.methods    sync.objects    lusearch    jbb2000    jbb2005

# Relative Speedup of Adding TLO Analysis



Performance Relative to Manual Allocation

slowdown

6

1

| 2 | 4 | 8 | | 2 | 4 | 8 | | 2 | 4 | 8 | | 2 | 4 | 8 | | 2 | 4 | 8 | | 2 | 4 | 8 |

traffic    sync.methods    sync.objects    lusearch    jbb2000    jbb2005

# Relative Speedup of Using Dynamic Locking



Performance Relative to Manual Allocation

slowdown

6

1

2   4   8     2   4   8     2   4   8     2   4   8     2   4   8     2   4   8

traffic    sync.methods   sync.objects   lusearch    jbb2000    jbb2005

Performance Relative to Manual Allocation

# Outline

# Conclusions

- Singleton allocation is not generally viable
- Points-to analysis precision is important
- MHP analysis helps if it can split a larger component
- TLO analysis usually has a negligible effect
- Dynamic locking has a small impact; may degrade or improve performance
- Component-based allocation works surprisingly well for many benchmarks

# Future Work

- More precise compiler analyses
- Finer locking granularities
- Method synchronization
- Critical section inference
- Speculative locking and transactional memory

# Questions?

Thank you for your attention.

# Related Work

- May Happen in Parallel analysis for Java (Naumovich *et al.* '99, Li '04).
- Thread-sensitive points-to and escape analysis (Chang and Choi '04, Sǎlcianu and Rinard '01).
- Thread-local objects analysis for synchronization elimination (Ruf '00).
- Pessimistic atomic sections/transactions (McCloskey *et al.* '06, Hicks *et al.* '06).
- Lock allocation
    - Concurrency graph (Sreedhar, Zhang, *et al.* '05).
    - ILP-based optimal allocations (Sreedhar, Zhang, *et al.* '05, Emmi *et al.* '07).
- Static race detection (Naik *et al.*'06, and *many others*).
- Optimistic concurrency, transactional memory (*see Larus & Rajwar '06*).