

# Lightweight Data Race Detection for Production Runs

Ohio State CSE Technical Report #OSU-CISRC-1/15-TR01; last updated August 2016

Swarnendu Biswas

Ohio State University  
biswass@cse.ohio-state.edu

Man Cao

Ohio State University  
caoma@cse.ohio-state.edu

Minjia Zhang

Ohio State University  
zhanminj@cse.ohio-state.edu

Michael D. Bond

Ohio State University  
mikebond@cse.ohio-state.edu

Benjamin P. Wood

Wellesley College  
bpw@cs.wellesley.edu

## Abstract

To detect real data races, program analysis must target production runs. However, sound and precise data race detection adds too much run-time overhead for use in production systems. Even existing approaches that provide soundness *or* precision incur significant limitations.

This work addresses the need for soundness (no missed races) and precision (no false races) by introducing new production-time analyses that address each need separately. (1) *Precise* data race detection is useful for developers, who fix bugs but loathe false positives. We introduce a novel, precise analysis called *LiteCollider* that provides low, bounded run-time overhead. (2) *Sound* race detection benefits analyses and tools whose correctness relies on knowledge of *all* potential data races. We present a novel, sound but efficient approach called *Caper* that combines static and dynamic analysis to detect every true race (plus false races) in observed runs.

Our evaluation, which compares implementations of *LiteCollider* and *Caper* with our implementations of closely related prior work, shows that *LiteCollider* and *Caper* are efficient and effective, and compare favorably with state-of-the-art approaches that target the same problem. These results suggest that *LiteCollider* and *Caper* enable practical race detection that is precise and sound, respectively, ultimately leading to more reliable software systems.

## 1. Introduction

In a multithreaded, shared-memory program execution, a *data race* occurs when two accesses are *conflicting* (two threads access the same variable with at least one write) and *concurrent* (not ordered by synchronization operations) [3]. Data races often directly or indirectly lead to concurrency bugs [45, 58]. The presence of data races—whether accidental or intentional—can affect an execution by crashing, hanging, or silently corrupting data [34, 39, 54]. The Therac-

25 disaster [43], the Northeastern electricity blackout of 2003 [69], and the mismatched NASDAQ Facebook share prices [59] were all due to race conditions. Data races will only become more problematic as software systems become increasingly parallel to scale with parallel hardware.

Data races not only tend to be associated with bugs, but they have weak or undefined semantics in modern shared-memory languages and systems [2, 14, 15, 48]. The C++ memory model provides no guarantees for racy programs [15]. A C++ program with apparently “benign” data races can behave erroneously, as a result of compiler transformations or being ported to a different architecture [13]. The Java memory model attempts to preserve memory and type safety even in the presence of data races [48], but researchers have shown that common compiler optimizations permit violations of Java’s memory model [16, 65].

**Detecting data races.** Program analyses can detect data races, but there exists a fundamental tradeoff between *soundness* (no missed races) and *precision* (no false race reports). Most static and some dynamic analyses provide soundness<sup>1</sup> but report many false races [22, 25, 28, 52, 53, 56, 57, 60, 62, 71, 72, 75], which developers find unacceptable [8, 49]. On the other hand, precise dynamic analyses generally cannot detect races in executions *other than* the current execution [10, 17, 29, 33, 40, 49]. This tradeoff is compounded by a second challenge: the *occurrence* of a data race is sensitive to thread interleavings, program inputs, and execution environments—so data races can remain undetected during in-house testing, even for extensively tested programs [69], and occur unexpectedly in production runs [43, 59, 69].

To find data races that manifest only in production, race detection must target production runs. However, in practice, race detectors see little use in production runs, due to high run-time overheads [49]. This paper addresses an open chal-

<sup>1</sup>A dynamic analysis is *sound* if it misses no races in observed executions.

lenge: devising *sound and precise* data race detection analyses that are *efficient* enough for production systems.

**Our approach.** A key, high-level insight of this work is that although sound and precise data race detection are too inefficient for production, *separate analyses that each provide soundness or precision are still beneficial*. Detecting precise (real) data races enables developers to find and fix software bugs and thus improves software reliability in the long run. On the other hand, detecting a sound overapproximation of data races can help simplify and optimize other dynamic analyses such as record & replay [42], atomicity checking [32], and software transactional memory systems [66]. Section 2 further motivates these benefits.

Mirroring this idea of separating precision and soundness, our approach decouples data race detection as two complementary, lightweight analyses that maintain and refine a precise underapproximation and a sound overapproximation of all data races over the course of many production runs.

To get a precise underapproximation, we introduce a novel dynamic analysis called *LiteCollider* (Section 4). Each run of *LiteCollider* takes one potential data race as input, and tries to detect whether the potential race occurs in the run. *LiteCollider* provides low, *bounded* run-time overhead, and it does not rely on any hardware support. In contrast, the closest related work adds unbounded, unscalable run-time overhead [40] (Sections 2.1 and 7.2.1).

For a sound overapproximation of data races from observed runs, we introduce a novel approach called *Caper* (Section 5). *Caper* combines static and dynamic analysis, starting with a set of potential races obtained from sound (no missed races) static analysis. During each analyzed run, *Caper* detects new potential races that prior runs have not already detected, by using a variant of dynamic escape analysis. *Caper* compares favorably with prior work that provides soundness, which provides worse performance or precision than *Caper* (Sections 2, 7.3, and 8).

We have implemented *LiteCollider* and *Caper* in a high-performance Java virtual machine [5] (Section 6). Our evaluation compares *LiteCollider* and *Caper* empirically with closely related work (Section 7). Overall, *LiteCollider* provides better performance (lower, scalable, and bounded overhead) and the same race coverage as prior work. *Caper* provides substantially better precision than static analysis, while providing low enough overhead for production runs, unlike other known approaches that are dynamically sound (no missed races in observed executions). These results suggest that these analyses can be integrated into and employed continuously in production settings, providing the benefits of precise underapproximation and sound overapproximation of observed data races.

## 2. Background and Motivation

Data race detection that is either precise or sound has distinct benefits. Prior approaches that provide precision or soundness have serious drawbacks.

### 2.1 Detecting Real Data Races Only

To target production environments, where the key constraint is run-time overhead, prior work has employed *sampling* of data race detection analysis, trading coverage for lower overhead [17, 29, 40, 49]. However, sampling-based race detection approaches suffer from run-time overhead that is high, unscalable, and unbounded. *LiteRace* and *Pacer* sample race detection analysis but instrument all program accesses, adding high baseline overhead even when the sampling rate is miniscule [17, 49]. *RaceMob* applies sampling by limiting its analysis to a single pair of static accesses per execution, limiting its instrumentation overhead [40], but even its optimized analysis incurs high, unscalable, unbounded overhead, as we show empirically in Section 7.2.1.

While most sampling-based approaches track the happens-before relation [41], *DataCollider* exposes and detects simultaneous, conflicting accesses—a sufficient condition for a data race [29]. Our paper refers to this kind of analysis as *collision analysis*. (Prior work has employed collision analysis in a non-sampling context [63].) To expose and detect simultaneous, conflicting accesses, *DataCollider* periodically pauses a thread’s execution at a potentially racy access; in the meantime, other threads detect conflicting accesses to the same variable. *DataCollider* avoids heavyweight instrumentation by using hardware debug registers to monitor memory locations [29]. However, debug registers are hardware specific and thus unportable; architectures often have only a few debug registers, limiting the number of memory locations that can be monitored simultaneously. The overhead of setting debug watchpoints with inter-processor interrupts increases with the number of cores, so *DataCollider* may not scale well to many cores [70]. Furthermore, although developers or users can adjust *DataCollider*’s run-time overhead by adjusting the sampling rate, it does not guarantee bounded run-time overhead.

In summary, existing precise data race detection analyses have serious limitations that make them unsuitable for production settings. This work asks the question: what is the best use of the production setting for detecting data races precisely? In particular, is it possible to design a precise, portable analysis that adds low, bounded overhead?

### 2.2 Detecting All Data Races in Observed Executions

While sampling-based approaches are precise, they are inherently unsound, missing data races that occur in production runs. This situation is unacceptable for analyses and systems whose correctness relies on soundly knowing *all* data races, such as record & replay systems [42], atomicity violation detectors [32], and software transactional memory [66]. For example, *Chimera* provides sound multithreaded record

& replay by conservatively tracking ordering between all potential data races, as identified by sound static analysis [42]. Although whole-program static analysis can provide a sound set of potential data races, its precision does not scale with program size and complexity, and it suffers from many false positives [28, 52, 53, 60, 72]. Chimera (and other approaches) could provide significantly better performance by using a more precise set of potential data races.

This work asks the question: Is it possible to exploit the production setting in order to detect potential data races soundly, i.e., without missing any true races that occur in production runs? Such an approach must add low overhead, and it must have “reasonable” precision (i.e., significantly better than static analysis) in order to be useful. To our knowledge, prior work does not provide a solution that both improves over the precision of static analysis and performs well enough for all-the-time use in production environments.

### 3. Overview

This paper addresses the two challenges motivated by the previous section with two complementary, lightweight production analyses. Section 4 presents *LiteCollider*, a precise analysis that detects only true data races. Section 5 presents *Caper*, a sound approach that detects all data races in observed runs. These two approaches—one precise but unsound, the other sound but imprecise—are inherently complementary. They maintain and refine an underapproximation and overapproximation, respectively, of the set of real data races in observed executions.

### 4. LiteCollider: Precise Data Race Detection

This section describes a new analysis called *LiteCollider* that targets minimal production-time race detection. *LiteCollider* limits its analysis to a single potential data race per program execution, and it bounds the run-time overhead it adds. In practice, *LiteCollider* could get each potential race from a set of races identified by developers or by another analysis that identifies potential races [40, 63]. While it may seem unsatisfying to try to detect only one potential race per execution, we note that current practice is to perform *no* race detection in production!

*LiteCollider* provides several properties aside from precision: bounded time and space overhead, scalability (i.e., time and space overheads remain stable with more threads), and portability. Like *DataCollider*, *LiteCollider* employs collision analysis, but *DataCollider* cannot bound its overhead, nor is it portable [29] (Section 2.1). On the other hand, sampling-based analyses that track happens-before provide precision and portability, but they cannot provide low, scalable, or bounded overhead [17, 40, 49] (Section 2.1).

#### 4.1 How the Analysis Works

A production run executing *LiteCollider* takes as input a single *potential* data race, which is an ordered pair of program sites,  $\langle s_1, s_2 \rangle$ . A *site* is a unique static program location (e.g.,

a method and bytecode index in Java). *LiteCollider* limits its analysis to these two sites (or one site if  $s_1 = s_2$ ). When  $s_1 \neq s_2$ , *LiteCollider* considers  $\langle s_1, s_2 \rangle$  and  $\langle s_2, s_1 \rangle$  to be distinct, and it uses separate runs to try to detect them.

Before a thread  $T$  executes an access to memory location  $m$  at  $s_1$ , the analysis potentially *samples* the current access by *waiting*, i.e., pausing the current thread for some time. The analysis updates global analysis state to indicate that  $T$  is waiting at  $s_1$  to access  $m$ .

When a thread  $T'$  executes an access to memory location  $m$  at  $s_2$ , the analysis checks whether some thread  $T$  is already waiting at  $s_1$  and is accessing the same memory location  $m$ . If so, the analysis has definitely detected a data race, and it reports the stack traces of  $T$  and  $T'$ .

Aside from the cost of waiting at some instances of  $s_1$ , the analysis can achieve very low overhead because it instruments only two (or one, if  $s_1 = s_2$ ) static sites and—like other collision analyses [29, 63]—avoids instrumenting synchronization operations. Although updating global analysis state may seem to be a potential scalability bottleneck, these updates occur when at least one thread is waiting—and the fraction of time spent waiting is bounded, as discussed next.

**Instrumentation overhead.** Although *LiteCollider* bounds how long it spends waiting, it always executes instrumentation at  $s_1$  and  $s_2$ , which is lightweight in the common case (particularly because it uses optimizations described shortly) but could incur nontrivial overhead for very frequent accesses. *LiteCollider*’s design includes the ability to control this cost at run time by detecting very frequent accesses (e.g., if  $s_1$  or  $s_2$  is in a hot, tight loop) and triggering dynamic re-compilation or code patching to remove the instrumentation.

#### 4.2 Sampling Policy

*LiteCollider*’s sampling policy bounds the total amount of time spent waiting, and it prioritizes instances of  $s_1$  more likely to be involved in data races.

**Budgeting the overhead of waiting.** *LiteCollider* is a “best-effort” approach that seeks to detect data races without affecting production performance too much. As in prior work called *QVM* that targets a specified maximum overhead by throttling the analysis [6], *LiteCollider* targets a maximum overhead  $r_{max}$  specified by developers or users. To enforce this maximum, *LiteCollider* keeps track throughout execution of (1) wall-clock time of the *ongoing* execution,  $t_{total}$ , and (2) wall-clock time *so far* during which one or more threads are waiting, i.e., taking into account overlapped waits. *LiteCollider* ensures the following throughout execution:

$$\frac{t_{waited}}{t_{total}} \leq r_{max}$$

This condition is conservative because a thread waiting for time  $t$  does not necessarily extend total execution time by  $t$ .

LiteCollider computes the following probability for whether to take a sample at  $s_1$ :<sup>2</sup>

$$P_{budget} = 1 - \frac{t_{waited} + t_{delay}}{(t_{total} + t_{delay}) \times r_{max}}$$

This computation is based on the fraction of time that will have been spent waiting *after* waiting for a planned amount of time  $t_{delay}$ . Note that  $P_{budget}$  will be close to 0 if LiteCollider is near its maximum overhead.  $P_{budget}$  will be close to 1 if LiteCollider is substantially under-budget.

**Prioritizing accesses.** Some static sites execute only once per run, while others execute millions of times. To account for this uncertainty and variability, it is important to wait at the first instance of a site (it might be the only one!) but less important to wait at later instances. These ideas echo the intuition behind the *cold-region hypothesis*, which postulates that the likelihood of detecting bugs in some part of a program is inversely proportional to the execution frequency of that part of the program [21, 49]. As in prior work on sampling-based race detection [49], we find that it is important to prioritize sampling of each *thread*’s initial access(es) (instead of the global execution’s initial accesses) at  $s_1$ .

LiteCollider thus samples an instance of  $s_1$  at a rate inversely proportional to  $freq(s_1, T)$ , the execution frequency of  $s_1$  by the current thread  $T$  so far:

$$P_{freq} = \frac{1}{freq(s_1, T)}$$

**Overall sampling probability.** LiteCollider’s sampling policy combines  $P_{budget}$  and  $P_{freq}$ :

$$P_{LiteCollider} = P_{budget} \times P_{freq}$$

A remaining issue is that soon after an execution starts, LiteCollider will have virtually no budget ( $P_{budget} \approx 0$ ), so very early accesses will go unsampled, potentially missing some races consistently. To detect such races, it seems unavoidable that an execution must risk exceeding its budget. LiteCollider assumes a minimum total running time  $t_{min}$  and uses  $max(t_{min}, t_{total})$  in place of  $t_{total}$  when computing  $P_{budget}$ . Developers can estimate a reasonable, small value for  $t_{min}$  using their knowledge of the program and likely execution scenarios. In our experiments,  $t_{min} = 1$  second.

### 4.3 Optimizations

LiteCollider’s low instrumentation overhead and good scalability (Section 7.2) rely on the following optimizations.

<sup>2</sup> Although LiteCollider could choose which dynamic accesses to sample using a deterministic function, it instead uses *randomness*, waiting at each dynamic access with some probability  $P$  as described below. Using randomness increases the chances of finding a race in the long run, over multiple executions, although in practice LiteCollider detects most races consistently from run to run.

**Sampling check.** LiteCollider’s instrumentation at  $s_1$  uses the following optimized two-part check, which is equivalent to sampling with probability  $P_{LiteCollider}$ . The instrumentation first computes  $P_{freq}$ , which is thread local and cheap. Then, with probability  $P_{freq}$ —that is, infrequently in the common case—the instrumentation computes  $P_{budget}$  and samples the access with probability  $P_{budget}$ .

**Avoiding scalability bottlenecks.** Whenever instrumentation at  $s_1$  waits, it acquires a lock on global analysis state, records that the current thread  $T$  is waiting at site  $s_1$  on memory location  $m$ , and it increments a global count of waiting threads. The instrumentation then performs a non-busy timed wait on the global lock. Instrumentation at  $s_2$  checks if another thread is waiting on the same memory location  $m$  at  $s_1$ . To do this efficiently,  $s_2$  first checks the global counter to see if any threads are waiting, *avoiding* remote cache misses in the common case when no threads are waiting. It only acquires the global lock and accesses global analysis state if at least one thread is waiting.

## 5. Caper: Detecting All Potential Data Races

This section focuses on the problem of detecting a sound overapproximation of *all* data races that occur in production runs. Although whole-program static analysis can provide a sound set of potential races, it reports *many* false races (Sections 7.3). Can a sound, low-overhead approach provide significantly better precision than static analysis alone?

### 5.1 Caper Overview

We introduce a novel approach called *Caper* that detects a set of potential data races that includes all true data races from observed production runs. Caper combines static and dynamic analyses. Caper initially runs a static analysis (e.g., [28, 52, 53, 60, 72]) to produce a set of *statically possible race pairs*,  $spPairs$ . This set consists of unordered pairs  $\langle s_1, s_2 \rangle$  (i.e.,  $\langle s_1, s_2 \rangle = \langle s_2, s_1 \rangle$ ) where  $s_1$  and  $s_2$  are each a static program location. The set  $spPairs$  need not be particularly precise—and in fact we find that a state-of-the-art static analysis reports many false races [53] (Section 7.3).

During each program execution, Caper’s *dynamic* analysis identifies dynamically possible race pairs,  $dpPairs$ , from  $spPairs$ . Each invocation of Caper copies newly identified possible race pairs from  $spPairs$  to  $dpPairs$  ( $dpPairs \subseteq spPairs$ ). At any given time,  $dpPairs$  is a sound overapproximation (no missed races) of every data race observed so far across all analyzed program executions.

Caper incurs low run-time overhead by using two insights. First, Caper refines  $spPairs$  and  $dpPairs$  on each successive program execution, so that in steady state, production runs are unlikely to detect any new pairs in  $spPairs$  to move to  $dpPairs$ . This feature allows Caper to optimize for *not* detecting new pairs in  $spPairs$  that should move to  $dpPairs$ . Second, Caper employs lightweight, sound dynamic escape analysis, described shortly. Although Caper’s

run-time overhead in steady state is low, it is *not* bounded (unlike LiteCollider). Having unbounded overhead seems inherent to ensuring soundness (short of “giving up” and declaring all *spPairs* to be in *dpPairs*).

## 5.2 Caper’s Dynamic Analysis

Caper’s dynamic analysis builds on *dynamic escape analysis* (DEA), which identifies objects that have potentially become shared (accessed by two or more threads). Caper’s DEA instruments all accesses to soundly track the escape property, while Caper’s sharing analysis instruments sites in *spPairs*, i.e., every site  $s$  such that  $\exists s' \mid \langle s, s' \rangle \in spPairs \setminus dpPairs$ . The instrumentation checks if the object accessed at  $s$  has escaped; if so, it marks  $s$  as “escaped.” Caper maintains a set of dynamically escaped sites called *deSites* such that

$$deSites = \{s \mid (\exists s' \mid \langle s, s' \rangle \in spPairs) \wedge s \text{ escaped in an observed run}\}$$

Caper removes  $\langle s_1, s_2 \rangle$  from *spPairs* and adds it to *dpPairs* for future invocations of Caper if both  $s_1$  and  $s_2$  are in *deSites*, i.e.,  $dpPairs = \{\langle s_1, s_2 \rangle \mid s_1 \in deSites \wedge s_2 \in deSites\}$ . Future runs of Caper instrument only the sites that are in  $spPairs \setminus dpPairs$ , but DEA continues to monitor all memory accesses.

**Caper’s reachability-based DEA.** Caper’s form of DEA is based on the idea that an object is escaped if it is transitively reachable from another escaped object. Reachability-based analysis conservatively identifies the first time an object becomes reachable by some thread other than its allocating thread, using these rules:

- Each object is initialized to NOT\_ESCAPED state on allocation. Thread objects (e.g., java.lang.Thread objects in Java) are initialized to ESCAPED state.
- At each reference-type store to a global variable (C.sf = p), the object referenced by p becomes ESCAPED.
- At each reference-type store to an instance field (q.f = p), the object referenced by p becomes ESCAPED if the object referenced by q is ESCAPED.
- Whenever an object becomes ESCAPED, mark all objects transitively reachable from the object as ESCAPED.

**Soundness.** Importantly, this reachability-based DEA is *sound* for detecting *shortest* data races, which are data races that are *not* dependent on some other race. In contrast, a non-shortest data race  $r$  “depends” on some other data race  $r'$ , meaning that eliminating  $r'$  necessarily eliminates  $r$ . Consider the following example, where thread T1’s local variable  $x$  initially refers to an unescaped object  $o$ :

```
// T1:           // T2:
x.f = ...; // s1
C.sf = x; // s2

y = C.sf; // s3
... = y.f; // s4
```

T1’s access to  $o.f$  through  $x$  occurs while the object  $o$  is still unescaped. Object  $o$  becomes reachable by other threads when T1 publishes a reference to it in the escaped field, C.sf.<sup>3</sup> T2 acquires this reference via C.sf and uses it to access  $o$ ’s field  $f$ .

T1 and T2 race on C.sf at  $s_2$  and  $s_3$ . This is a shortest race: it does not depend on any other races. T1 and T2 also perform unsynchronized accesses to  $o.f$  at  $s_1$  and  $s_4$ , but this race depends on the other, “shorter” race.  $s_1$  and  $s_4$  can only race on  $o.f$  if  $s_2$  and  $s_3$  also race on C.sf. Fixing the  $s_2$ – $s_3$  race (e.g., by making C.sf a volatile field in Java) necessarily fixes the  $s_1$ – $s_4$  race.

In contrast to Caper, prior DEA-based data race detection is generally unsound [23, 35, 44, 56, 57, 62] (Section 8). To our knowledge, our work is the first to apply a form of DEA as a fully sound filter for data race detection.<sup>4</sup> More significantly, to our knowledge, Caper is the first to use DEA to prune the results of static data race detection.

## 6. Implementation

We have implemented LiteCollider and Caper’s dynamic analysis in Jikes RVM 3.1.3 [4, 5], a Java virtual machine that has performance competitive with commercial JVMs [10]. We have also implemented two analyses from prior work that are closely related to LiteCollider and Caper. We will make these implementations publicly available.

### 6.1 Detecting Data Races Precisely

**LiteCollider.** LiteCollider’s implementation closely follows the design, described in detail in Section 4. LiteCollider’s design includes the ability to dynamically remove the lightweight instrumentation from  $s_1$  and  $s_2$  in the event that these sites are very frequent, e.g., in a hot, tight loop (Section 4.1). We have not implemented this feature since it is not needed for our experiments.

Instrumenting different accesses in each production run, as in LiteCollider, naturally lends itself to just-in-time-compiled languages such as Java: the dynamic compiler instruments only those accesses targeted by the current run. An implementation in an ahead-of-time compiled language such as C++ could use dynamic code patching, distribute a different binary to each production site, or (at higher run-time cost) instrument all potentially racy sites [40].

**LiteHB.** We have implemented an analysis that we call *LiteHB* that is based on RaceMob’s happens-before analysis [40].<sup>5</sup> Like LiteCollider, LiteHB takes as input a single potential data race  $\langle s_1, s_2 \rangle$ , and tries to detect whether it manifests in the current execution.

<sup>3</sup>We chose a static field for simplicity. Any escaped field suffices.

<sup>4</sup>The TRaDe data race detector applies a sound DEA, but extends it with a reprivatization optimization that is unsound for data race detection, allowing TRaDe to miss some true data races [23].

<sup>5</sup>Although RaceMob’s happens-before analysis also includes waiting at some accesses [40], we omit this feature from LiteHB in order to measure the cost of (optimized) happens-before analysis alone.

LiteHB implements RaceMob’s optimized tracking of the happens-before relationship [40]. In particular, LiteHB (1) only starts tracking happens-before after some thread executes  $s_1$  and (2) stops tracking happens-before when all threads can no longer race with a prior instance of  $s_1$ , i.e., when all threads’ vector clocks *happen after* the clocks for all executed instances of  $s_1$ .

Our LiteHB implementation adds the following metadata to each object: a header word to track the vector clock if the object’s lock is acquired; a header word that points to an array of metadata, which is used if the object is an array that is accessed by LiteHB’s instrumentation; and a word of metadata for each field, which is used if the field is accessed by LiteHB’s instrumentation. With additional engineering effort, LiteHB could avoid adding some of this metadata; in particular, it could limit per-field metadata to those fields that  $s_1$  and  $s_2$  might access. Instead, our evaluation of LiteHB accounts for the costs of extraneous metadata.

## 6.2 Detecting Potential Data Races Soundly

**Caper.** As described in Section 5, Caper starts with statically overapproximating the set of potential data races. For Caper’s static analysis, we use the publicly available static data race detector *Chord*, revision 1825 [53].<sup>6</sup> Chord’s full analysis is *unsound* because it uses a may-alias analysis for locks.<sup>7</sup> We thus *disable* Chord’s static lockset analysis entirely, along with a few other unsound options such as ignoring accesses in constructors, and we enable a Chord feature that resolves reflective calls by executing the program. The result is a sound analysis that identifies potential races based solely on (static) thread escape analysis and fork-join analysis. Although Chord analyzes each program together with the Java libraries, it analyzes a different library implementation than what Jikes RVM uses, so our experiments detect potential races in application code only. Caper uses the set reported by Chord as its *spPairs* (Section 5). Caper’s dynamic analysis works as described in Section 5.

**Dynamic alias analysis.** For comparison with Caper, we have implemented a *dynamic alias analysis*, which detects all pairs of aliasing sites. It reports  $\langle s_1, s_2 \rangle$  if and only if accesses at  $s_1$  and  $s_2$  by threads  $T_1$  and  $T_2$  at memory locations  $m_1$  and  $m_2$  such that  $T_1 \neq T_2 \wedge m_1 = m_2$ .

Prior work uses dynamic alias analysis as a filter for race detection, by checking one potentially aliasing pair per run [40]. Dynamic alias analysis is less imprecise than Caper but adds high run-time overhead (Section 7.3). It offers a different point of comparison than static analysis, which is highly imprecise but adds no overhead. To our knowledge, there exists no sound analysis that has (1) better precision than Caper and (2) overhead low enough for production.

<sup>6</sup> <http://www.cc.gatech.edu/~naik/chord.html>

<sup>7</sup> We have confirmed with Naik that there is no available sound implementation of Chord that uses conditional must-not alias analysis [52].

## 7. Evaluation

This section evaluates LiteCollider and Caper’s efficiency and effectiveness, and compares with competing approaches.

### 7.1 Methodology

**Benchmarks.** Our evaluation analyzes and executes benchmarked versions of large, real applications. Our experiments execute the DaCapo Benchmarks [11], versions 2006-10-MR2 and 9.12-bach, which we distinguish using names suffixed by 6 and 9; and fixed-workload versions of SPECjbb2000 and SPECjbb2005.<sup>8</sup> We omit programs with only one thread or that Jikes RVM cannot execute. We also omit eclipse6 since Chord’s static analysis fails when analyzing it; and jython9 and pmd9 because Chord reports no potential data races for them.<sup>9</sup> We run the *large* workload size of the DaCapo benchmarks. Table 1(a) (page 8) reports total and maximum active threads for each program.

**Platform.** For each of our implementations, we build a high-performance configuration of Jikes RVM (FastAdaptive) that adaptively optimizes the application at run time and uses the default high-performance garbage collector, which adjusts the heap size automatically.

The experiments execute on an AMD Opteron 6272 system with eight 8-core 2.0 GHz processors (64 cores total), running 64-bit RedHat Enterprise Linux 6.7, kernel 2.6.32.

### 7.2 Detecting Data Races Precisely

This section evaluates the run-time performance and race detection coverage of LiteCollider, compared with LiteHB.

**Methodology.** Each run of LiteCollider or LiteHB takes as an input a potential data race  $\langle s_1, s_2 \rangle$ . In a real production setting, such potential races could be identified by developers or by another analysis (Section 4). Potential races could be generated by a *sound* analysis such as static analysis, Caper, or dynamic alias analysis; then, across many production runs, LiteCollider and LiteHB have at least the potential to detect all production-time data races. *RaceMob* uses potential races identified by sound static analysis as input to its LiteHB-like analysis (after first filtering them using dynamic alias analysis that checks one potential race per run) [40].

However, to keep our experiments manageable (especially since we run multiple configurations and trials), we evaluate LiteCollider and LiteHB on a relatively small set of potential races: the set of access pairs identified by dynamic alias analysis *that also violate the dynamic lockset property*, i.e., no common lock is held for both accesses; prior work introduces heavyweight dynamic analyses that check the lockset property [22, 25, 56, 57, 62, 71]. This methodology mirrors a likely real-world scenario: *prioritization* of certain potential races as input to LiteCollider or LiteHB, e.g., through

<sup>8</sup> <http://users.cecs.anu.edu.au/~steveb/research/research-infrastructure/pjbb2005>

<sup>9</sup> pmd9 has data races [10], but Chord does not recognize its use of threads.

additional heavyweight but less-imprecise analysis than Capser run at *testing* time. Note that generating the set of potential data races only at testing time would miss some data races that manifest in production only. Our experiments use the same input for identifying potential races, so the potential races include all known data races (all races that manifest in executed runs).

### 7.2.1 Performance

**Run-time overhead.** Figure 1 shows the overhead added over uninstrumented execution (unmodified Jikes RVM) for configurations of LiteCollider and LiteHB. Each bar is the mean execution time over 30 randomly selected potential races (from those identified by dynamic alias analysis that violate the lockset property). If a program has fewer than 30 potential races, we use multiple trials of each potential race.

The LiteCollider configurations are for three target maximum overheads,  $r_{max} = 0\%$ ,  $5\%$ , and  $10\%$ . All configurations use wait times of  $t_{delay} = 10$  ms. The  $r_{max} = 0\%$  configuration shows overhead without any waiting, measuring the cost of instrumentation at accesses alone. This overhead, which LiteCollider cannot measure or bound with its sampling model, is consistently low (always  $<5\%$ ). For  $r_{max} = 5\%$  and  $10\%$ , LiteCollider stays under the target overhead across all programs, as expected. Since LiteCollider instruments only one pair of accesses, the instrumentation overhead is generally low. The average overheads for  $r_{max} = 5\%$  and  $10\%$  are  $1.0\%$  and  $1.6\%$ , respectively. It is unsurprising that actual overheads are generally less than  $r_{max}$ , since LiteCollider’s model conservatively assumes that pausing one thread by  $t_{delay}$  slows the entire program by  $t_{delay}$ .

The *LiteHB* and *FullHB* configurations are variants of LiteHB. *LiteHB* is the default LiteHB algorithm that uses RaceMob’s optimized happens-before tracking; it adds  $20\%$  run-time overhead on average. *LiteHB (only HB tracking)* shows overhead incurred by *optimized* happens-before tracking only: it adds no instrumentation to accesses, except for instrumentation at  $s_1$  that enables happens-before tracking if it is currently disabled. This configuration isolates the cost of optimized tracking of happens-before, which incurs  $13\%$  average overhead.

We note that both LiteHB configurations include space overhead, as well as cache pressure and GC costs, from per-object and per-field metadata—but most of the per-field metadata could be avoided with additional engineering effort (Section 6.1). We find that per-field metadata *alone* adds run-time overhead of  $7\%$  (results not shown), so the true cost of optimized happens-before analysis is as low as  **$13\%$** .

*FullHB (only HB tracking)* performs *unoptimized* happens-before tracking by performing vector clock computations at every synchronization operation; it adds no instrumentation at accesses. This configuration, which shows the benefit of RaceMob’s happens-before optimization [40], adds  $18\%$  average overhead. The last configuration, *FullHB*, shows the overhead of instrumenting  $s_1$  and  $s_2$  and performing unopti-

mized happens-before tracking, and adds  $37\%$  overhead on average. These results show that RaceMob’s optimization is indeed beneficial in our experiments.

For sunflow9, both LiteHB and FullHB add *especially high* overheads ( $162$  and  $526\%$ , respectively). As the results indicate, most of this overhead actually comes from instrumentation at accesses, not from tracking the happens-before relationship. In sunflow9, instrumentation at mostly read-shared accesses must perform *updates* to per-variable clocks, leading to many remote cache misses. Such high, *unpredictable* overheads prohibit use of happens-before-analysis-based race detection on production systems.

We have also measured the *space* overhead incurred by LiteCollider and LiteHB (full results omitted). LiteCollider, which adds only a small amount of global metadata, adds  $1\%$  overhead on average and at most  $5\%$  overhead for any program; much of that overhead is actually a side effect of slowing execution, which leads to more dynamic compiler work. For LiteHB, if we subtract out all space overhead due to per-object and per-field metadata (optimistically assuming all such metadata could be avoided; Section 6.1), space overhead is  $5\%$  on average and at most  $29\%$  for any program.

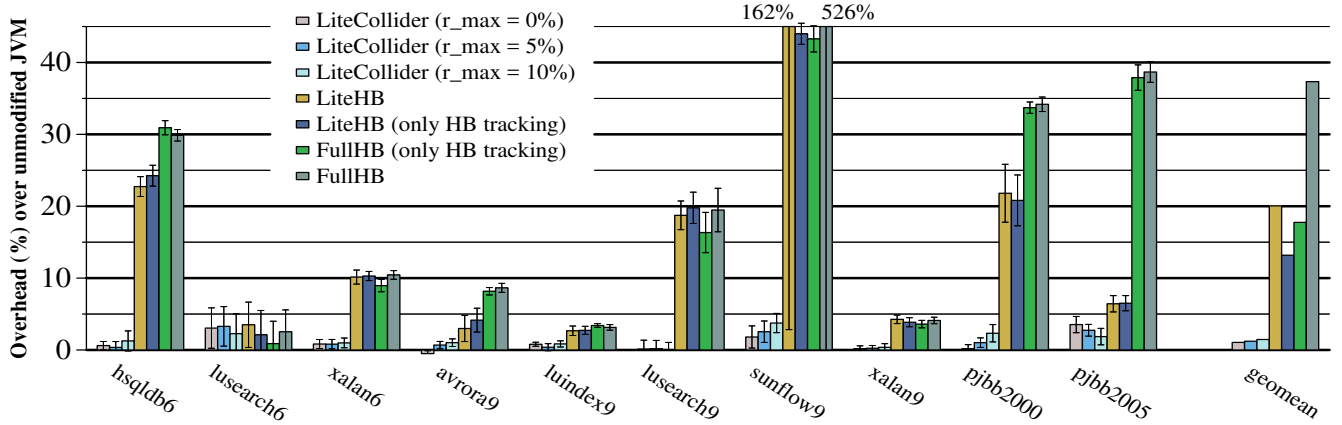
**Scalability.** Production software often has many threads and will likely have more as future systems provide more cores [46]. To compare how LiteCollider and LiteHB scale with more threads, we evaluate the three programs that support spawning a configurable number of application threads: lusearch9, sunflow9, and xalan9 (Table 1(a)). Figure 2 plots execution time for 1–64 application threads, using configurations from Figure 1. The *Unmodified JVM* configuration shows that lusearch9 and sunflow9 naturally scale with more threads, while xalan9 anti-scales starting at 16 threads.

*LiteCollider* ( $r_{max} = 5\%$ ) not only adds low overhead, but its overhead is largely unaffected by the number of threads for all programs. Note that in the plots, the *Unmodified JVM* and *LiteCollider* lines are difficult to distinguish. In contrast, for lusearch9 and sunflow9, *LiteHB* not only adds high overhead, but its overhead grows with more threads.<sup>10</sup> As discussed earlier, most of sunflow9’s LiteHB overhead is due to conflicts from updating per-variable clock metadata—a cost that increases as the number of threads increases. For lusearch9, since most of its overhead comes from tracking the happens-before relation, which requires vector clock computations that take linear time in the number of threads. For xalan9, LiteHB always incurs very low overhead and has no noticeable scalability issues.

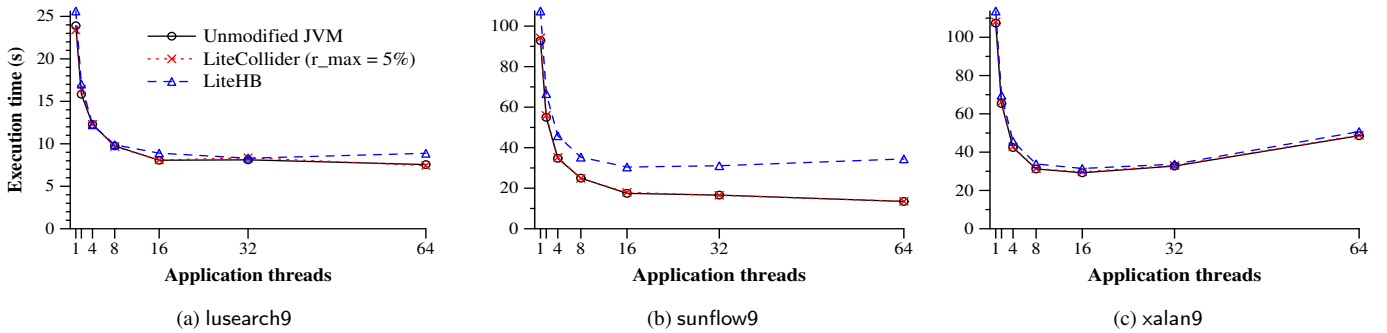
These scalability results imply that for even larger thread counts, LiteHB will add additional overhead, while LiteCollider will provide consistently low overhead.

**Comparison with RaceMob’s reported results.** Although RaceMob reports very low run-time overhead for its happens-

<sup>10</sup>For 32 threads, the *Unmodified JVM* and *LiteCollider* configurations experience anomalous performance for lusearch9. We have attributed this anomaly to Linux thread scheduling decisions [7].



**Figure 1.** Run-time overhead of configurations of LiteCollider and LiteHB. The error bars represent 95% confidence intervals.



**Figure 2.** Execution times of LiteCollider ( $r_{max} = 5\%$ ) and LiteHB for 1–64 application threads. The execution times for 64 threads correspond to the results in Figure 1. The legend applies to all graphs.

Java program	Total	Max live	Avg. accesses per SFR
hsqldb6	402	102	26
lusearch6	65	65	156
xalan6	9	9	21
avrora9	27	27	553
luindex9	2	2	718
lusearch9	$c$	$c$	201
sunflow9	$2 \times c$	$c$	1,030,000
xalan9	$c$	$c$	53
pjbb2000	37	9	7
pjbb2005	9	9	15

(a) Java programs

C/C++ program	Threads	Avg. accesses per SFR		
		$n=8$	$n=16$	$n=32$
blackscholes	$1 + n$	9,150,000	4,750,000	2,290,000
bodytrack	$2 + n$	63,600	57,400	47,800
canneal	$1 + n$	5,470,000	2,746,000	1,370,000
dedup	$3 + 3n$	36,300	36,300	35,900
ferret	$3 + 4n$	630,000	514,000	388,000
fluidanimate	$1 + n$	131	99	68
raytrace	$1 + n$	5,820,000	3,030,000	1,550,000
streamcluster	$1 + 2n$	4,320	2,260	1,250
swaptions	$1 + n$	83,000,000	41,600,000	20,800,000
vips	$3 + n$	105,000	81,100	55,800
x264	$1 + 2 \times frames$	208,000	202,000	202,000

(b) C/C++ programs

**Table 1.** Spawned threads and average executed memory accesses per synchronization-free region (SFR), rounded to three significant figures and the nearest integer, for Java and C/C++ programs. (a) The table shows *Total* threads created and *Max live* at any time, which is dependent on the core count  $c$  (64 in our experiments) for three programs. (b)  $n$  is PARSEC’s minimum threads parameter. *frames* is the input-size-dependent number of frames processed by x264.

before analysis evaluated on C/C++ programs (2.3% on average) [40], our implementation of optimized happens-before analysis, evaluated on Java programs, has significantly higher overhead: 13% on average (after subtracting out overhead from per-field metadata). Although there

are many implementation and experimental differences that might lead to differences, we suspect three primary reasons. First, our programs generally have more threads (Table 1). Second, Java programs tend to have much more frequent synchronization operations. Table 1 compares the sizes,



measured in executed memory accesses, of synchronization-free regions for C/C++ and the Java programs that we evaluate, i.e., the ratio of total memory accesses to total synchronization operations. The C/C++ programs are the PARSEC benchmark suite [9], version 3.0-beta-20150206 (excluding *freqmine* since it uses OpenMP for its parallelization and *facesim* which does to run with our tool), with the *simmedium* input size; we count their synchronization operations and memory accesses by modifying a Pintool from prior work [24, 47]. With the exception of one Java program (*sunflow9*), synchronization operations are several orders of magnitude more frequent in Java programs than most C/C++ programs. A related issue is that the RaceMob authors report that their evaluated C/C++ programs execute synchronization barriers, which optimized happens-before analysis can exploit by tracking happens-before immediately after a barrier [40]. In contrast, the Java programs we evaluate execute few or no synchronization barriers.

As evidence of these claims, we note that the RaceMob authors implement and evaluate a sampling-based happens-before analysis from prior work called *Pacer* [17]. Their *Pacer* implementation adds low overhead compared to the original Java-based implementation and evaluation of *Pacer*, suggesting that happens-before analysis is relatively inexpensive for RaceMob’s evaluated programs. For synchronization-intensive programs, their *Pacer* implementation adds high overhead, but RaceMob adds low overhead, suggesting that RaceMob’s happens-before optimization is particularly effective for the evaluated programs.

**Comparison with other happens-before tracking.** The state-of-the-art sound and precise data race detection analysis *FastTrack* incurs 340% average overhead for a comparable implementation in Jikes RVM [10], or 750% for the *FastTrack* authors’ implementation and evaluation [33]. Industry-standard tools such as Intel’s Inspector XE,<sup>11</sup> Google’s ThreadSanitizer v2 [64], and Helgrind [55] are largely based on happens-before analysis. They add high run-time overheads and are suitable only for testing runs.

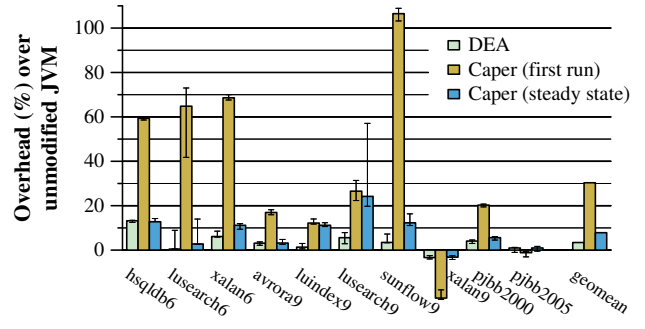
The sampling-based happens-before race detector called *Pacer* avoids analysis at most operations, but still requires instrumenting all operations [17]. For sampling rates of 0%, 1%, and 3%, *Pacer*’s implementation in Jikes RVM incurs 33%, 52%, and 86% run-time overhead on average [17]—significantly higher than LiteCollider’s overhead.

### 7.2.2 Detecting Real Data Races

**Empirical comparison.** Table 2 reports how many (true) data races LiteCollider and LiteHB detect. For each potential race pair from the sound set of lockset-violating access pairs, we execute 10 trials each of LiteCollider and LiteHB. The first column for each analysis reports distinct races reported at least once across 10 trials. To account for uncertainty about the repeatability of detecting a race reported in just

	LiteCollider		LiteHB		Overlap	
hsqldb6	10	(10)	10	(10)	10	(10)
lusearch6	0	(0)	0	(0)	0	(0)
xalan6	17	(17)	17	(17)	17	(17)
avrorra9	7	(6)	7	(7)	7	(6)
luindex9	2	(2)	2	(2)	2	(2)
lusearch9	0	(0)	0	(0)	0	(0)
sunflow9	2	(2)	2	(2)	2	(2)
xalan9	7	(7)	7	(7)	7	(7)
pjbb2000	7	(7)	7	(7)	7	(7)
pjbb2005	1	(1)	1	(1)	1	(1)
<b>Total</b>	53	(52)	53	(53)	53	(52)

**Table 2.** Data races reported by LiteCollider and LiteHB. The two numbers are distinct races reported at least once and (in parentheses) at least twice, out of 10 trials. *Overlap* is distinct races reported by both LiteCollider and LiteHB.



**Figure 3.** Run-time performance of Caper.

one trial, we also report, in parentheses, races reported in at least two trials. We run LiteCollider configured with a target maximum run-time overhead of  $r_{max} = 5\%$  and a fixed wait time of  $t_{delay} = 10$  ms. The table omits 30 additional access pairs (all in *sunflow9* or *xalan9*) reported by LiteHB, which are non-shortest, or dependent, races (Section 5.2).

LiteCollider detects all of the races detected by LiteHB, with the caveat that LiteCollider detects one race in *avrorra9* in only one of the trials. LiteCollider’s per-trial coverage is dependent on its (randomized) sampling policy, which could benefit from further work. We have verified that the data races reported by LiteCollider and LiteHB match those reported by a publicly available implementation of *FastTrack* in Jikes RVM from prior work [10].

Thus, despite using sampling-based collision analysis, LiteCollider provides essentially the same coverage as happens-before analysis for the evaluated programs. Furthermore, LiteCollider provides significantly lower, bounded overhead and better scalability than LiteHB.

### 7.3 Detecting Potential Data Races Soundly

This section evaluates the performance and precision of Caper, compared with static analysis and dynamic alias analysis.

<sup>11</sup> <https://software.intel.com/en-us/intel-inspector-xe>

### 7.3.1 Performance

**Static analysis.** Caper initially employs static analysis to generate the set  $spPairs$ . Static analysis needs to execute only once (or whenever the code changes) and does not affect run time, so its performance is not crucial. Chord takes at most 30 minutes to analyze any program.

**Caper’s dynamic analysis.** Figure 3 shows the overhead added by Caper’s dynamic analysis over unmodified Jikes RVM. Each bar is the mean of 25 trials, with 95% confidence intervals centered at the mean. The first bar, *DEA*, shows the overhead of reachability-based dynamic escape analysis alone, which incurs 3% overhead on average. *Caper (first run)* Caper’s dynamic analysis when it runs for the first time for a program, i.e., when  $dpPairs = \emptyset$  and  $deSites = \emptyset$ . Under these conditions, the analysis finds many newly escaped sites to add to  $deSites$ , incurring 27% average overhead.

In contrast, *Caper (steady state)*, represents performance during production, when prior testing (and production) runs have added nearly all escaped sites to  $deSites$ . Under these conditions, Caper’s dynamic analysis elides instrumentation at sites in  $deSites$ , and it incurs little or no overhead from instrumented sites being added to  $deSites$ . On average, *Caper (steady state)* incurs 9% run-time overhead. These results suggest that Caper in steady state is efficient enough for many production environments.

The results for *xalan9* are unintuitive. We find that adding any instrumentation to *xalan9* improves performance relative to the baseline. We are investigating this issue, which we believe is due to Linux thread scheduling decisions.

**Dynamic alias analysis.** We find that our implementation of dynamic alias analysis slows program execution by 13X on average (results not shown). Admittedly, this implementation is not heavily optimized; we have implemented it mainly for measuring its precision. But dynamic alias analysis is inherently a heavyweight analysis because it must track every accessed combination of thread, site, and variable, i.e.,  $\langle t, s, x \rangle$ , in order to detect future aliasing in the run.

### 7.3.2 Precision and Effectiveness

**Precision.** Table 3 shows potential races (unique, unordered access pairs) reported by the static data race detection analysis Chord, Caper, and dynamic alias analysis. The table counts a potential race if it is reported at least once across 10 trials. Although it is undecidable whether a potential race is real or not, we assume that the vast majority of potential races detected are false races; as evidence, we note that predictive analyses have been able to expose at most dozens, not thousands, of additional races beyond those found by precise dynamic analysis [30, 38, 68].

The *Static analysis* columns show potential race pairs identified by *unsound* and *sound* versions of Chord. For most programs, the *sound* analysis reports tens of thousands of potential races. In comparison, Chord’s *unsound* analysis, which is less imprecise since it uses Chord’s may-alias lock-

	Known races	Static analysis		Caper	Dyn. alias analysis
		unsound	sound		
hsqldb6	10	13,749 (5)	212,205	1,612	757
lusearch6	0	395 (0)	4,692	302	292
xalan6	17	70,263 (7)	83,488	1,241	581
avro9	7	1,301 (5)	61,193	19,941	570
luindex9	2	6,015 (0)	10,257	192	193
lusearch9	0	441 (0)	7,303	34	39
sunflow9	2	24,616 (0)	28,587	200	1,086
xalan9	7	13,335 (0)	20,036	1,861	600
pjbb2000	7	12,708 (0)	29,604	11,243	1,679
pjbb2005	1	682 (0)	2,552	984	447

**Table 3.** Potential data races reported by two variants of static analysis, Caper, and dynamic alias analysis. *Known races* are data races reported by a precise detector on the large workload (see Table 2). For *unsound* static analysis, the number in parentheses is known data races *missing* from the set of reported potential races.

	Static analysis	Caper	Dyn. alias analysis
hsqldb6	97%	44%	71%
lusearch6	73%	57%	52%
xalan6	74%	20%	8%
avro9	99%	91%	54%
luindex9	33%	4%	<1%
lusearch9	84%	<1%	<1%
sunflow9	65%	14%	53%
xalan9	62%	28%	<1%
pjbb2000	72%	39%	35%
pjbb2005	97%	26%	8%

**Table 4.** Percentage of *dynamic memory accesses* that three sound analyses identify as part of a potential data race.

set analysis [53], reports three times fewer potential races on average. However, the *unsound* analysis is demonstrably *unsound*: we find that it misses 17 of the 53 real data races identified by LiteCollider and LiteHB. We have verified that the potential races reported by *sound* static analysis, Caper, and dynamic escape analysis include all known true races.

The last two columns of Table 3 show the precision of Caper and dynamic alias analysis. Notably, Caper usually provides substantially better precision (often 1–2 orders of magnitude fewer potential races) than static analysis for all but two programs (*avro9* and *pjbb2000*). For the most part, dynamic alias analysis provides better precision than Caper, which is unsurprising since dynamic alias analysis detects aliasing and Caper detects reachability-based escaping. However, in some cases, Caper actually reports fewer potential races than dynamic alias analysis because Caper does not count a site as escaped if it accesses a non-escaped object, even if the object later becomes escaped. This effect is particularly pronounced for *sunflow9*, whose main thread initializes data before the data escapes and is accessed by worker threads. Importantly, Caper’s approach is sufficient both for detecting all true data races (Section 5.2) and for clients that need sound knowledge of data races or sharing [32, 42, 66] (Section 2.2).

**Effectiveness.** To estimate the effectiveness of Caper in optimizing other dynamic analyses that must account for potential data races [32, 42, 66], we compute how many *dynamic (executed) accesses* are identified as being part of a potential data race. We expect this value to be proportional to a client dynamic analysis’s run-time overhead. Table 4 shows the percentage of dynamic accesses that static analysis, Caper, and dynamic alias analysis identify as being part of a potential data race. Each percentage is computed as:

$$\frac{\sum_{s|(\exists s'|(s,s')\in\text{potentialRaces})} \text{freq}(s)}{\sum_s \text{freq}(s)}$$

where *potentialRaces* is the set of potential races identified by the analysis (e.g., *dpPairs* for Caper), and *freq(s)* is the dynamic execution frequency of site *s*.

The table shows that static analysis’s high imprecision leads to most executed accesses being potentially racy. Caper improves precision substantially over static analysis for all programs but *avro9*. Although dynamic alias analysis usually identifies the most executed accesses as data race free, Caper provides better precision for *sunflow9* and *hsqldb6*, due to analysis differences discussed above.

In summary, these results suggest that Caper is an efficient and effective approach that, compared with alternate approaches, provides a reasonable tradeoff in *balancing performance and precision* in detecting potential data races.

## 8. Related Work

Previous sections compared our LiteCollider analysis with closely related work on sampling-based data race detectors including DataCollider and RaceMob’s happens-before analysis [29, 40] (Sections 2 and 7.2.1), and Caper with static analysis and dynamic alias analysis. This section compares LiteCollider and Caper with other existing work.

**Dynamic escape analysis.** Dynamic escape analysis (DEA) identifies when thread-private data becomes shared (Section 5.2). DEA can be implemented based on either first shared access [35, 57, 62], reachability from shared context [23], or variants on these approaches [44, 56].

First-shared-access-based analysis marks each field with its original accessing thread [35, 56, 57, 62].<sup>12</sup> These analyses check at each access whether the current accessing thread is the same as the original accessing thread, thereby detecting the first shared access and marking the field shared. Access-based analysis is *unsound* as a filter for more expensive data race checks. Since no race detection metadata is recorded for accesses to a field before its first shared access, the data race detector cannot determine whether *pre-sharing* accesses race with *post-sharing* accesses on this field.

Reachability-based DEA has the potential to be a sound filter for data race detection, as we show in Section 5.2.

<sup>12</sup>Nishiyama’s DEA marks an object escaped when a second thread *obtains a direct reference to the object*, which is still unsound for race detection [56].)

However, existing data race detection techniques that make use of reachability-based DEA may miss data races. The *TRaDe* race detection analysis uses sound, reachability-based DEA, but pairs it with an unsound optimization for privatization of escaped objects [23]. The *SOS* race detection analysis uses a less precise, but sound, reachability-based DEA as part of a larger *unsound* optimization to detect stationary objects [44]. Reachability-based DEA has been deployed soundly in the implementation of thread-local heaps [26] and software transactional memory [66].

Recent work introduces a field-precise dynamic analysis to precisely track sharing across threads [37], but its overhead is too high for production.

**Optimizing data race detection.** Wester et al. parallelize happens-before and lockset analyses by using a speculation-based technique called *uniparallelism* [73]. Low overhead relies on having extra available cores for speculative execution. Veeraraghavan et al. employ uniparallelism to infer data races, based on different outcomes under deterministic synchronization schedules [70].

Several analyses detect conflicts between regions of code, in order to detect the subset of true data races that may violate *region serializability* in the current execution [10, 27, 46, 50, 67], trading coverage for performance. These techniques either require custom hardware [46, 50, 67] or incur substantial slowdowns [10, 27].

Custom hardware can accelerate data race detection by adding on-chip memory for tracking vector clocks or locksets, and extending cache coherence to identify shared accesses [24, 74, 76]. However, manufacturers have been reluctant to change already-complex memory and cache subsystems substantially to support race detection.

**Increasing coverage.** Predictive analysis and model checking enhance race detection coverage without sacrificing precision [30, 38, 51, 68]. These analyses are inherently heavyweight and unsuited for production, but they cannot find all real races due to the following limitations. Predictive analysis finds data races in an execution *other than* the observed execution [30, 38, 68], but coverage is still limited largely by the observed execution. (*Racageddon* uses a combination of concolic testing and predictive analysis [30].) Model checkers such as *CHESS* [51] can explore many thread interleavings and/or inputs, but they suffer state-space explosion and do not scale well to large programs.

**Estimating harmfulness.** Prior work tries to infer which data races are most likely to be harmful (e.g., crash the program, hurt performance, or corrupt data) [19, 20, 29, 34, 39, 54, 63]. Several approaches expose errors by exposing rare but allowable behaviors, such as weak memory model behaviors, at racy accesses [19, 20, 34, 54, 63]. Prioritizing races is complementary to our work, which seeks to expose and detect data races. Researchers have argued persuasively that essentially all data races are problematic because lan-

languages like Java and C/C++ provide virtually no guarantees for them [2, 14, 15, 48] (Section 1).

**Languages and types.** Researchers have designed new languages or extended existing languages to ensure programs are data race free by construction. Race-free languages require a new, constrained programming model [12, 36, 61]. Type-based approaches require programmer annotations to assist type inference [1, 31]; writing annotations is tedious and constraining. Boyapati et al. add ownership types to an existing programming language, so that well-typed programs are guaranteed to be race free [18].

## 9. Conclusion

Data races are elusive, manifesting unexpectedly in production runs. Sound and precise race detection is too expensive for production runs, so we attack soundness and precision separately, introducing complementary approaches that maintain a precise underapproximation and sound overapproximation of true data races—useful inputs both for developers and for other analyses and tools. LiteCollider is a novel, precise race detector that adds bounded run-time overhead and outperforms its closest competitor without sacrificing coverage. Caper combines static and dynamic analyses in a novel way to provide significantly better precision than static analysis alone, without the high overhead of more-precise dynamic analyses, without missing true data races in observed executions. Together, these contributions and results demonstrate how to leverage production runs for data race detection effectively, ultimately improving the reliability and performance of production software systems.

## Acknowledgments

We thank Baris Kasikci for detailed feedback on the text; and Dan Grossman, Mayur Naik, and Aritra Sengupta for discussions and advice.

## References

- [1] M. Abadi, C. Flanagan, and S. N. Freund. Types for Safe Locking: Static Race Detection for Java. *TOPLAS*, 28(2):207–255, 2006.
- [2] S. V. Adve and H.-J. Boehm. Memory Models: A Case for Rethinking Parallel Languages and Hardware. *CACM*, 53:90–101, 2010.
- [3] S. V. Adve, M. D. Hill, B. P. Miller, and R. H. B. Netzer. Detecting Data Races on Weak Memory Systems. In *ISCA*, pages 234–243, 1991.
- [4] B. Alpern, C. R. Attanasio, J. J. Barton, M. G. Burke, P. Cheng, J.-D. Choi, A. Cocchi, S. J. Fink, D. Grove, M. Hind, S. F. Hummel, D. Lieber, V. Litvinov, M. Mergen, T. Ngo, J. R. Russell, V. Sarkar, M. J. Serrano, J. Shepherd, S. Smith, V. C. Sreedhar, H. Srinivasan, and J. Whaley. The Jalapeño Virtual Machine. *IBM Systems Journal*, 39(1):211–238, 2000.
- [5] B. Alpern, S. Augart, S. M. Blackburn, M. Butrico, A. Cocchi, P. Cheng, J. Dolby, S. Fink, D. Grove, M. Hind, K. S. McKinley, M. Mergen, J. E. B. Moss, T. Ngo, and V. Sarkar. The Jikes Research Virtual Machine Project: Building an Open-Source Research Community. *IBM Systems Journal*, 44:399–417, 2005.
- [6] M. Arnold, M. Vechev, and E. Yahav. QVM: An Efficient Runtime for Detecting Defects in Deployed Systems. In *OOPSLA*, pages 143–162, 2008.
- [7] E. D. Berger, T. Yang, T. Liu, and G. Novark. Grace: Safe Multithreaded Programming for C/C++. In *OOPSLA*, pages 81–96, 2009.
- [8] A. Bessey, K. Block, B. Chelf, A. Chou, B. Fulton, S. Hallem, C. Henri-Gros, A. Kamsky, S. McPeak, and D. Engler. A Few Billion Lines of Code Later: Using Static Analysis to Find Bugs in the Real World. *CACM*, 53(2):66–75, 2010.
- [9] C. Bienia, S. Kumar, J. P. Singh, and K. Li. The PARSEC Benchmark Suite: Characterization and Architectural Implications. In *PACT*, pages 72–81, 2008.
- [10] S. Biswas, M. Zhang, M. D. Bond, and B. Lucia. Valor: Efficient, Software-Only Region Conflict Exceptions. In *OOPSLA*, pages 241–259, 2015.
- [11] S. M. Blackburn, R. Garner, C. Hoffman, A. M. Khan, K. S. McKinley, R. Bentzur, A. Diwan, D. Feinberg, D. Frampton, S. Z. Guyer, M. Hirzel, A. Hosking, M. Jump, H. Lee, J. E. B. Moss, A. Phansalkar, D. Stefanović, T. VanDrunen, D. von Dincklage, and B. Wiedermann. The DaCapo Benchmarks: Java Benchmarking Development and Analysis. In *OOPSLA*, pages 169–190, 2006.
- [12] R. L. Bocchino, Jr., V. S. Adve, S. V. Adve, and M. Snir. Parallel Programming Must Be Deterministic by Default. In *HotPar*, pages 4–9, 2009.
- [13] H.-J. Boehm. How to miscompile programs with “benign” data races. In *HotPar*, 2011.
- [14] H.-J. Boehm. Position paper: Nondeterminism is Unavoidable, but Data Races are Pure Evil. In *RACES*, pages 9–14, 2012.
- [15] H.-J. Boehm and S. V. Adve. Foundations of the C++ Concurrency Memory Model. In *PLDI*, pages 68–78, 2008.
- [16] H.-J. Boehm and B. Demsky. Outlawing Ghosts: Avoiding Out-of-Thin-Air Results. In *MSPC*, pages 7:1–7:6, 2014.
- [17] M. D. Bond, K. E. Coons, and K. S. McKinley. Pacer: Proportional Detection of Data Races. In *PLDI*, pages 255–268, 2010.
- [18] C. Boyapati, R. Lee, and M. Rinard. Ownership Types for Safe Programming: Preventing Data Races and Deadlocks. In *OOPSLA*, pages 211–230, 2002.
- [19] J. Burnim, K. Sen, and C. Stergiou. Testing Concurrent Programs on Relaxed Memory Models. In *ISSTA*, pages 122–132, 2011.
- [20] M. Cao, J. Roemer, A. Sengupta, and M. D. Bond. Prescient Memory: Exposing Weak Memory Model Behavior by Looking into the Future. In *ISMM*, pages 99–110, 2016.
- [21] T. M. Chilimbi and M. Hauswirth. Low-Overhead Memory Leak Detection Using Adaptive Statistical Profiling. In *ASPLOS*, pages 156–164, 2004.
- [22] J.-D. Choi, K. Lee, A. Loginov, R. O’Callahan, V. Sarkar, and M. Sridharan. Efficient and Precise Datarace Detection for Multithreaded Object-Oriented Programs. In *PLDI*, pages 258–269, 2002.

- [23] M. Christiaens and K. De Bosschere. TRaDe, A Topological Approach to On-the-fly Race Detection in Java Programs. In *Symposium on Java Virtual Machine Research and Technology Symposium*, pages 15–15, 2001.
- [24] J. Devietti, B. P. Wood, K. Strauss, L. Ceze, D. Grossman, and S. Qadeer. RADISH: Always-On Sound and Complete Race Detection in Software and Hardware. In *ISCA*, pages 201–212, 2012.
- [25] A. Dinning and E. Schonberg. Detecting Access Anomalies in Programs with Critical Sections. In *PADD*, pages 85–96, 1991.
- [26] D. Doligez and X. Leroy. A Concurrent, Generational Garbage Collector for a Multithreaded Implementation of ML. In *POPL*, pages 113–123, 1993.
- [27] L. Effinger-Dean, B. Lucia, L. Ceze, D. Grossman, and H.-J. Boehm. IFRit: Interference-Free Regions for Dynamic Data-Race Detection. In *OOPSLA*, pages 467–484, 2012.
- [28] D. Engler and K. Ashcraft. RacerX: Effective, Static Detection of Race Conditions and Deadlocks. In *SOSP*, pages 237–252, 2003.
- [29] J. Erickson, M. Musuvathi, S. Burckhardt, and K. Olynyk. Effective Data-Race Detection for the Kernel. In *OSDI*, pages 1–16, 2010.
- [30] M. Eslamimehr and J. Palsberg. Race Directed Scheduling of Concurrent Programs. In *PPoPP*, pages 301–314, 2014.
- [31] C. Flanagan and S. N. Freund. Type Inference Against Races. *SCP*, 64(1):140–165, 2007.
- [32] C. Flanagan and S. N. Freund. Atomizer: A Dynamic Atomicity Checker for Multithreaded Programs. *SCP*, 71(2):89–109, 2008.
- [33] C. Flanagan and S. N. Freund. FastTrack: Efficient and Precise Dynamic Race Detection. In *PLDI*, pages 121–133, 2009.
- [34] C. Flanagan and S. N. Freund. Adversarial Memory for Detecting Destructive Races. In *PLDI*, pages 244–254, 2010.
- [35] C. Flanagan and S. N. Freund. The RoadRunner Dynamic Analysis Framework for Concurrent Programs. In *PASTE*, pages 1–8, 2010.
- [36] D. Grossman. Type-Safe Multithreading in Cyclone. In *TLDI*, pages 13–25, 2003.
- [37] J. Huang. Scalable Thread Sharing Analysis. In *ICSE*, pages 1097–1108, 2016.
- [38] J. Huang, P. O. Meredith, and G. Rosu. Maximal Sound Predictive Race Detection with Control Flow Abstraction. In *PLDI*, pages 337–348, 2014.
- [39] B. Kasicki, C. Zamfir, and G. Candea. Data Races vs. Data Race Bugs: Telling the Difference with Portend. In *ASPLOS*, pages 185–198, 2012.
- [40] B. Kasicki, C. Zamfir, and G. Candea. RaceMob: Crowdsourced Data Race Detection. In *SOSP*, pages 406–422, 2013.
- [41] L. Lamport. Time, Clocks, and the Ordering of Events in a Distributed System. *CACM*, 21(7):558–565, 1978.
- [42] D. Lee, P. M. Chen, J. Flinn, and S. Narayanasamy. Chimera: Hybrid Program Analysis for Determinism. In *PLDI*, pages 463–474, 2012.
- [43] N. G. Leveson and C. S. Turner. An Investigation of the Therac-25 Accidents. *IEEE Computer*, 26(7):18–41, 1993.
- [44] D. Li, W. Srisa-an, and M. B. Dwyer. SOS: Saving Time in Dynamic Race Detection with Stationary Analysis. In *OOPSLA*, pages 35–50, 2011.
- [45] S. Lu, S. Park, E. Seo, and Y. Zhou. Learning from Mistakes: A Comprehensive Study on Real World Concurrency Bug Characteristics. In *ASPLOS*, pages 329–339, 2008.
- [46] B. Lucia, L. Ceze, K. Strauss, S. Qadeer, and H.-J. Boehm. Conflict Exceptions: Simplifying Concurrent Language Semantics with Precise Hardware Exceptions for Data-Races. In *ISCA*, pages 210–221, 2010.
- [47] C.-K. Luk, R. Cohn, R. Muth, H. Patil, A. Klauser, G. Lowney, S. Wallace, V. J. Reddi, and K. Hazelwood. Pin: Building Customized Program Analysis Tools with Dynamic Instrumentation. In *PLDI*, pages 190–200, 2005.
- [48] J. Manson, W. Pugh, and S. V. Adve. The Java Memory Model. In *POPL*, pages 378–391, 2005.
- [49] D. Marino, M. Musuvathi, and S. Narayanasamy. LiteRace: Effective Sampling for Lightweight Data-Race Detection. In *PLDI*, pages 134–143, 2009.
- [50] D. Marino, A. Singh, T. Millstein, M. Musuvathi, and S. Narayanasamy. DRFx: A Simple and Efficient Memory Model for Concurrent Programming Languages. In *PLDI*, pages 351–362, 2010.
- [51] M. Musuvathi and S. Qadeer. Iterative Context Bounding for Systematic Testing of Multithreaded Programs. In *PLDI*, pages 446–455, 2007.
- [52] M. Naik and A. Aiken. Conditional Must Not Aliasing for Static Race Detection. In *POPL*, pages 327–338, 2007.
- [53] M. Naik, A. Aiken, and J. Whaley. Effective Static Race Detection for Java. In *PLDI*, pages 308–319, 2006.
- [54] S. Narayanasamy, Z. Wang, J. Tigani, A. Edwards, and B. Calder. Automatically Classifying Benign and Harmful Data Races Using Replay Analysis. In *PLDI*, pages 22–31, 2007.
- [55] N. Nethercote and J. Seward. Valgrind: A Framework for Heavyweight Dynamic Binary Instrumentation. In *PLDI*, pages 89–100, 2007.
- [56] H. Nishiyama. Detecting Data Races using Dynamic Escape Analysis based on Read Barrier. In *VMRT*, pages 127–138, 2004.
- [57] R. O’Callahan and J.-D. Choi. Hybrid Dynamic Data Race Detection. In *PPoPP*, pages 167–178, 2003.
- [58] J. Ouyang, P. M. Chen, J. Flinn, and S. Narayanasamy. ...and region serializability for all. In *HotPar*, 2013.
- [59] PCWorld. Nasdaq’s facebook glitch came from race conditions, 2012. [http://www.pcworld.com/article/255911/nasdaqs\\_facebook\\_glitch\\_came\\_from\\_race\\_conditions.html](http://www.pcworld.com/article/255911/nasdaqs_facebook_glitch_came_from_race_conditions.html).
- [60] P. Pratikakis, J. S. Foster, and M. Hicks. LOCKSMITH: Context-Sensitive Correlation Analysis for Race Detection. In *PLDI*, pages 320–331, 2006.
- [61] M. C. Rinard and M. S. Lam. The Design, Implementation, and Evaluation of Jade. *TOPLAS*, 20:483–545, 1998.
- [62] S. Savage, M. Burrows, G. Nelson, P. Sobalvarro, and T. Anderson. Eraser: A Dynamic Data Race Detector for Multithreaded Programs. In *SOSP*, pages 27–37, 1997.
- [63] K. Sen. Race Directed Random Testing of Concurrent Programs. In *PLDI*, pages 11–21, 2008.

- [64] K. Serebryany, A. Potapenko, T. Iskhodzhanov, and D. Vyukov. Dynamic Race Detection with LLVM Compiler. pages 110–114, 2012.
- [65] J. Ševčík and D. Aspinall. On Validity of Program Transformations in the Java Memory Model. In *ECOOP*, pages 27–51, 2008.
- [66] T. Shpeisman, V. Menon, A.-R. Adl-Tabatabai, S. Balensiefer, D. Grossman, R. L. Hudson, K. F. Moore, and B. Saha. Enforcing Isolation and Ordering in STM. In *PLDI*, pages 78–88, 2007.
- [67] A. Singh, D. Marino, S. Narayanasamy, T. Millstein, and M. Musuvathi. Efficient Processor Support for DRFx, a Memory Model with Exceptions. In *ASPLOS*, pages 53–66, 2011.
- [68] Y. Smaragdakis, J. Evans, C. Sadowski, J. Yi, and C. Flanagan. Sound Predictive Race Detection in Polynomial Time. In *POPL*, pages 387–400, 2012.
- [69] U.S.–Canada Power System Outage Task Force. Final Report on the August 14th Blackout in the United States and Canada. Technical report, Department of Energy, 2004.
- [70] K. Veeraraghavan, P. M. Chen, J. Flinn, and S. Narayanasamy. Detecting and Surviving Data Races using Complementary Schedules. In *SOSP*, pages 369–384, 2011.
- [71] C. von Praun and T. R. Gross. Object Race Detection. In *OOPSLA*, pages 70–82, 2001.
- [72] J. W. Voun, R. Jhala, and S. Lerner. RELAY: Static Race Detection on Millions of Lines of Code. In *ESEC/FSE*, pages 205–214, 2007.
- [73] B. Wester, D. Devesery, P. M. Chen, J. Flinn, and S. Narayanasamy. Parallelizing Data Race Detection. In *ASPLOS*, pages 27–38, 2013.
- [74] B. P. Wood, L. Ceze, and D. Grossman. Low-Level Detection of Language-Level Data Races with LARD. In *ASPLOS*, pages 671–686, 2014.
- [75] X. Xie and J. Xue. Acculock: Accurate and Efficient Detection of Data Races. In *CGO*, pages 201–212, 2011.
- [76] P. Zhou, R. Teodorescu, and Y. Zhou. HARD: Hardware-Assisted Lockset-based Race Detection. In *HPCA*, pages 121–132, 2007.