

BARIS KASIKCI

TECHNIQUES FOR DETECTION, ROOT CAUSE
DIAGNOSIS AND CLASSIFICATION OF
IN-PRODUCTION CONCURRENCY BUGS

TECHNIQUES FOR DETECTION, ROOT CAUSE
DIAGNOSIS AND CLASSIFICATION OF
IN-PRODUCTION CONCURRENCY BUGS

BARIS KASIKCI

Baris Kasikci: *Techniques for Detection, Root Cause Diagnosis and Classification of In-Production Concurrency Bugs*

Everybody who learns concurrency thinks they understand it, ends up finding mysterious races they thought weren't possible, and discovers that they didn't actually understand it yet after all.

— Herb Sutter

Any fool can know. The point is to understand.

— Albert Einstein

RÉSUMÉ

Les bogues de concurrences sont au cœur des pires problèmes que rencontrent les programmes en production. Les bogues de concurrence narguent les développeurs de logiciels en restant élusifs pendant des semaines voir des mois avant que les développeurs ne les trouvent et puissent les fixer. Les bogues de concurrence empêchent ainsi les développeurs de créer de nouveaux systèmes et programmes à un rythme plus élevé que ce qui est possible actuellement.

La détection, le diagnostic de la cause d'un bogue de concurrence et sa classification sont un défi spécialement compliqué pour les logiciels en production. Ces bogues sont difficiles à détecter à faible coût, les mécanismes usuels de détection requérant de lourdes analyses dynamiques, ce qui les rend inutilisables en production. Même une fois détectés, il est compliqué de déterminer quel problème de concurrence a la plus forte probabilité de causer une défaillance du programme, car de telles défaillances sont le résultats d'interactions complexes entre les entrées du programme, l'ordonnancement des processus, et le modèle de mémoire sous-jacent. Explorer de telles interactions demande des analyses demandant beaucoup de temps de calculs, ce qui n'est pas non plus utilisable pour des programmes en production. Même lorsque les bogues de concurrences qui peuvent probablement causer une défaillance du programme sont détectés, les développeurs ont encore besoin de trouver une explication de comment une telle défaillance a été causée. Cela demande de récolter un nombre important d'informations sur l'exécution du programme en production et d'exécuter de plus amples analyses, ce qui est difficile à faire sans impacter les performances à l'exécution.

Cette thèse développe des techniques pour une détection efficace et économique des bogues de concurrence pour des logiciels en production ainsi qu'une analyse de la cause de ceux-ci et leur classification. L'idée générale de cette thèse est de procurer au développeur des techniques qui lui permettront de mieux raisonner et comprendre le comportement des programmes concurrents. Cette thèse se base sur les principes suivant :

- Réduire l'impact des analyses en production au minimum : en décomposant sa solution en deux parties distincte, les analyses les plus lourdes en interne, et un minimum d'analyses en production, cette thèse permet d'obtenir une balance appropriée entre les analyses statiques et dynamiques des programmes.
- Éviter d'utiliser des techniques qui ne sont pas largement déployées, en limitant la portée des technologies qu'elle utilise à celles qui sont disponibles aujourd'hui. Cette approche frugale

- pousse cette thèse à développer de nouvelles techniques capable de résoudre les problèmes présentés sans utiliser de solutions sur mesure qui pourrait rendre impossible leur implémentation.
- Utiliser des techniques qui vont parfois à l’encontre de la sagesse populaire. Cette approche hérétique permet de surpasser certains problèmes non résolus autrement.

Utilisant les principes ci-dessus, cette dissertation propose:

1. Des techniques pour détecter des bogues de concurrence en production, en combinant des analyses statiques en interne et une analyse dynamique en production.
2. Une technique pour identifier la cause d’un bogue en production.
3. Une technique qui, pour un accès concurrent, la classifie selon ses conséquences potentielles, permettant aux développeurs de répondre à certaines questions comme “Cet accès concurrent cause-t-il un plantage, ou bloque-t-il le programme ?”, “Est-ce que cet accès concurrent peut-il changer la sortie du programme ?”, “Est ce que cet accès concurrent n’a aucun effet observable ?”.

Nous avons construit une série d’outils qui implémentent toutes les techniques mentionnées ci-dessus. Nous montrons que ces outils que nous avons développés dans cette thèse sont efficace, ont un impact faible sur les performances des logiciels et ont une haute précision.

Mots clés: bogues de concurrences, accès concurrent, violation d’atomicité, analyse statique, analyse dynamique

ABSTRACT

Concurrency bugs are at the heart of some of the worst bugs that plague in-production software. Concurrency bugs stymie software development by remaining elusive for weeks or even months before developers can identify and fix them. Consequently, concurrency bugs hinder developers' efforts for building new systems and programs with a faster pace than what is possible today.

Detection, root cause diagnosis and classification of concurrency bugs is especially challenging for in-production software. Concurrency bugs are hard to detect with low overhead, because common detection mechanisms require heavyweight analyses, and therefore they are not suited for in-production usage. Even when detected, it is hard to determine which concurrency bugs are more likely to cause program failures, because such consequences arise as a result of complicated interactions of program inputs, thread schedules, and the underlying memory model. Exploring such interactions require compute-intensive analyses, which are also not suited for in-production software. Even when concurrency bugs that are likely to lead to failures are detected, developers still need to come up with an explanation of how they caused a given failure. This requires gathering a significant amount of runtime information from in-production program execution and performing further analyses, which is hard to do without undue runtime performance overhead.

This dissertation develops techniques for effective and efficient detection, root cause diagnosis and classification of concurrency bugs for in-production software. The high level goal of this dissertation is to provide developers with techniques that will allow them to better reason about and understand the behavior of concurrent programs. It builds upon the following fundamental tenets:

- It decouples its solution techniques in order to perform most of the heavyweight analyses in-house and resorts to minimal in-production analysis. This separation of concerns allows this dissertation to strike an appropriate balance between static program analysis and dynamic program analysis.
- It eschews relying on technologies that are not widely deployed, by limiting the scope of the technologies it leverages to what is available today. This frugal approach pushes the limits of the dissertation to come up with novel techniques that are capable of solving the outlined problems without resorting to custom solutions that may be infeasible to implement in the real world.
- It employs techniques that sometimes oppose conventional wisdom. This heretical approach allows the dissertation to overcome key open problems.

Using the aforementioned key tenets, this dissertation proposes:

1. Techniques to detect concurrency bugs in-production by combining in-house static analysis and in-production dynamic analysis.
2. A technique to identify the root cause of concurrency bugs.
3. A technique that given a data race, classifies it based on its potential consequence, allowing developers to answer questions regarding the data race such as “can the data race cause a crash or a hang?”, “can the data race cause a change in program output?”, “does the data race have no observable effect?”.

We build a toolchain that implements all the aforementioned techniques. We show that the tools we develop in this dissertation are effective, incur low runtime performance overhead, have high accuracy and precision.

Keywords: Concurrency bugs, data race, atomicity violation, static analysis, dynamic analysis

PUBLICATIONS

This dissertation primarily builds upon the ideas presented in the following publications:

- Baris Kasikci et al. “Failure Sketching: A Technique for Automated Root Cause Diagnosis of In-Production Failures.” In: *Symp. on Operating Systems Principles*. 2015
- Baris Kasikci, Cristian Zamfir, and George Candea. “Automated Classification of Data Races Under Both Strong and Weak Memory Models.” In: *ACM Transactions on Programming Languages and Systems* 37.3 (2015)
- Baris Kasikci et al. “Efficient Tracing of Cold Code Via Bias-Free Sampling.” In: *USENIX Annual Technical Conf.* 2014
- Baris Kasikci, Cristian Zamfir, and George Candea. “RaceMob: Crowdsourced Data Race Detection.” In: *Symp. on Operating Systems Principles*. 2013
- Baris Kasikci, Cristian Zamfir, and George Candea. “Data Races vs. Data Race Bugs: Telling the Difference with Portend.” In: *Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*. 2012

ACKNOWLEDGMENTS

CONTENTS

i	SETTING THE STAGE	1
1	INTRODUCTION	3
1.1	Problem Definition	3
1.2	Challenges	5
1.2.1	The Runtime Performance Overhead Challenge	5
1.2.2	The Accuracy Challenge	7
1.2.3	The In-Production Challenge	8
1.3	Overview of Prior Solution Attempts	9
1.4	Attempts at Addressing the Overhead Challenge	9
1.5	Attempts at Addressing the Accuracy Challenge	10
1.6	Attempts at Addressing the In-Production Challenge	10
1.7	Solution Overview	11
1.8	Summary of Contributions	12
1.9	Summary of Results	13
2	BACKGROUND AND RELATED WORK	15
2.1	Concurrency Bug Surveys	15
2.2	Data Race Detection	15
2.2.1	Data Race Definition	15
2.2.2	Data Race Detection Literature	16
2.2.3	Static Data Race Detection	16
2.2.4	Dynamic Data Race Detection	16
2.2.5	Mixed Static-Dynamic Data Race Detection	18
2.2.6	Detecting Data Races In Production	18
2.3	Atomicity and Order Violation Detection	18
2.4	Data Race Classification	18
2.5	Root Cause Diagnosis of Concurrency Bugs	19
2.5.1	Delta Debugging	19
2.5.2	Statistical Techniques	19
ii	ELIMINATING CONCURRENCY BUGS FROM IN-PRODUCTION SYSTEMS	21
3	RACEMOB: DETECTING DATA RACES IN PRODUCTION	23
3.1	RaceMob's Design Overview	23
3.2	Static Data Race Detection	24
3.3	Dynamic Data Race Validation	25
3.3.1	Dynamic Context Inference	25
3.3.2	On-Demand Data Race Detection	26
3.3.3	Schedule Steering	28
3.4	Crowdsourcing the Validation	28
3.5	Reaching a Verdict	30
3.6	RaceMob's Implementation Details	31

4	ROOT CAUSE DIAGNOSIS OF IN-PRODUCTION FAILURES USING FAILURE SKETCHES	33
4.1	Gist's Design Overview	34
4.2	Static Slice Computation	36
4.3	Slice Refinement	38
4.3.1	Adaptive Slice Tracking	39
4.3.2	Tracking Control Flow	40
4.3.3	Tracking Data Flow	41
4.4	Identifying the Root Cause	43
4.5	Gist's Implementation Details	45
5	PORTEND: CLASSIFYING DATA RACES DURING DEVELOPMENT AND TESTING	47
5.1	A Fine-Grained Way to Classify Data Races	49
5.2	Portend's Design Overview	51
5.3	Single-Path Analysis	57
5.4	Multi-Path Analysis	58
5.5	Symbolic Output Comparison	60
5.6	Multi-Schedule Analysis	62
5.7	Portend's Implementation Details	63
6	EVALUATION	67
6.1	RaceMob's Evaluation	67
6.1.1	Experimental Setup	67
6.1.2	Effectiveness	68
6.1.3	Efficiency	69
6.1.4	Comparison to Existing Data Race Detectors	71
6.2	Root Cause Diagnosis Results	72
6.2.1	Experimental Setup	72
6.2.2	Automated Generation of Sketches	73
6.2.3	Accuracy of Failure Sketches	74
6.2.4	Efficiency	76
6.3	Data Race Classification Results	79
6.3.1	Experimental Setup	79
6.3.2	Effectiveness	81
6.3.3	Accuracy and Precision	82
6.3.4	Efficiency	84
6.3.5	Comparison to Existing Data Race Detectors	86
6.3.6	Efficiency and Effectiveness of Symbolic Memory Consistency Modeling	88
6.3.7	Memory Consumption of Symbolic Memory Consistency Modeling	90
iii	WRAPPING UP	93
7	FUTURE WORK	95
8	CONCLUSIONS	97
	BIBLIOGRAPHY	99

LIST OF FIGURES

Figure 1	Example of a switch statement adapted from [12]	4
Figure 2	False negatives in happened-before (HB) dynamic race detectors: the race on x is not detected in Execution 1, but it is detected in Execution 2.	8
Figure 3	RaceMob’s crowdsourced architecture: A static detection phase, run on the hive, is followed by a dynamic validation phase on users’ machines.	24
Figure 4	Minimal monitoring in DCI: For this example, DCI stops tracking synchronization operations as soon as each thread goes once through the barrier.	27
Figure 5	The state machine used by the hive to reach verdicts based on reports from program instances. Transition edges are labeled with validation results that arrive from instrumented program instances; states are labeled with RaceMob’s verdict.	31
Figure 6	The failure sketch of pbzip2 bug.	34
Figure 7	The architecture of Gist	35
Figure 8	Adaptive slice tracking in Gist	39
Figure 9	Example of control (a) and data (b) flow tracking in Gist. Solid horizontal lines are program statements, circles are basic blocks.	42
Figure 10	Four common atomicity violation patterns (RWR, WWR, RWW, WRW). Adapted from [4].	43
Figure 11	A sample execution failing at the second read in T_1 (a), and three potential concurrency errors: a RWR atomicity violation (b), 2 WR data races (c-d).	44
Figure 12	Portend taxonomy of data races.	49
Figure 13	High-level architecture of Portend. The six shaded boxes indicate new code written for Portend, whereas clear boxes represent reused code from KLEE [17] and Cloud9 [16].	52
Figure 14	Increasing levels of completeness in terms of paths and schedules: [a. single-pre/single-post] \ll [b. single-pre/multi-post] \ll [c. multi-pre/multi-post].	53

Figure 15	Simplified example of a harmful race from Ctrace [71] that would be classified as harmless by classic race classifiers. 55
Figure 16	Portend prunes paths during symbolic execution. 59
Figure 17	A program to illustrate the benefits of symbolic output comparison 62
Figure 18	Breakdown of average overhead into instrumentation-induced overhead and detection-induced overhead. 71
Figure 19	RaceMob scalability: Induced overhead as a function of the number of application threads. 72
Figure 20	The failure sketch of Curl bug #965. 73
Figure 21	The failure sketch of Apache bug #21287. The grayed-out components are not part of the ideal failure sketch, but they appear in the sketch that Gist automatically computes. 74
Figure 22	Accuracy of Gist, broken down into relevance accuracy and ordering accuracy. 76
Figure 23	Contribution of various techniques to Gist’s accuracy. 76
Figure 24	Gist’s average runtime performance overhead across all runs as a function of tracked slice size. 77
Figure 25	Tradeoff between slice size and the resulting accuracy and latency. Accuracy is in percentage, latency is in the number of failure recurrences. 78
Figure 26	Comparison of the full tracing overheads of Mozilla rr and Intel PT. 79
Figure 27	Breakdown of the contribution of each technique toward Portend’s accuracy. We start from single-path analysis and enable one by one the other techniques: ad-hoc synchronization detection, multi-path analysis, and finally multi-schedule analysis. 83
Figure 28	Simplified examples for each race class from real systems. (a) and (b) are from ctrace, (c) is from memcached and (d) is from pbzip2. The arrows indicate the pair of racing accesses. 84
Figure 29	Change in classification time with respect to number of preemptions and number of dependent branches for some of the races in Table 6. Each sample point is labeled with race id. 86
Figure 30	Portend’s accuracy with increasing values of k. 86

Figure 31	Micro-benchmarks used for evaluating SMC.	89
Figure 32	Running time of Portend-weak and Portend-seq	91
Figure 33	Memory usage of Portend-weak and Portend-seq	91

LIST OF TABLES

Table 1	Data race detection with RaceMob. The static phase reports <i>Race candidates</i> (row 2). The dynamic phase reports verdicts (rows 3-10). <i>Causes hang</i> and <i>Causes crash</i> are races that caused the program to hang or crash. <i>Single order</i> are true races for which either the primary or the alternate executed (but not both) with no intervening synchronization; <i>Both orders</i> are races for which both executed without intervening synchronization.	67
Table 2	Runtime overhead of race detection as a percentage of uninstrumented execution. Average overhead is 2.32%, and maximum overhead is 4.54%.	70
Table 3	Race detection results with RaceMob, ThreadSanitizer (TSAN), and RELAY. Each cell shows the number of reported races. The data races reported by RaceMob and TSAN are all true data races. The only true data races among the ones detected by RELAY are the ones in the row "RaceMob". To the best of our knowledge, two of the data races that cause a hang in SQLite were not previously reported.	70
Table 4	Programs analyzed with Portend. Source lines of code are measured with the <code>cloc</code> utility.	80
Table 5	"Spec violated" races and their consequences.	81

Table 6	Summary of Portend’s classification results. We consider two races to be distinct if they involve different accesses to shared variables; the same race may be encountered multiple times during an execution—these two different aspects are captured by the <i>Distinct races</i> and <i>Race instances</i> columns, respectively. Portend uses the stack traces and the program counters of the threads making the racing accesses to identify distinct races. The last 5 columns classify the distinct races. The <i>states same/differ</i> columns show for how many races the primary and alternate states were different after the race, as computed by the Record/Replay Analyzer [77]. 82
Table 7	Portend’s classification time for the 93 races in Table 6. 85
Table 8	Accuracy for each approach and each classification category, applied to the 93 races in Table 6. “Not-classified” means that an approach cannot perform classification for a particular class. 87
Table 9	Portend’s effectiveness in bug finding and state coverage for two memory model configurations: sequential memory consistency mode and Portend’s weak memory consistency mode. 90

Part I

SETTING THE STAGE

In this part, we define the problem tackled in this dissertation along with the associated challenges for solving it, and prior solution attempts. We give a brief overview of the solution we propose, followed by a thorough treatment of related work on detection, root cause diagnosis, and classification of concurrency bugs.

INTRODUCTION

In this chapter, we elaborate on the definition of the problem addressed in this dissertation (§1.1); we describe the challenges of detection, root cause diagnosis and classification of concurrency bugs for in-production software (§1.2); we summarize prior attempts at solving the problem (§1.3). Finally, we give an overview of the solution we propose in this dissertation (§1.7).

1.1 PROBLEM DEFINITION

Concurrency bugs such as data races, atomicity violations, and deadlocks are at the root of many software problems [66]. These problems have led to losses of human lives [61], caused massive material losses [98], and triggered various security vulnerabilities [23, 37]. Perhaps more subtly, concurrency bugs hinder reasoning about the behavior of concurrent programs because of their sporadic occurrence and unpredictable effects.

Concurrency bugs proliferated in modern software after the advent of multicore processors. As hardware became increasingly parallel, developers wrote more programs that tried to leverage such parallelism by relying on concurrency. Since then, multithreading and parallel programming became widespread. Concurrency is desirable for getting more performance out of parallel hardware, but it comes with a cost: concurrent programs are harder to write correctly, and therefore it is easier to make mistakes (e.g., data races, deadlocks, etc) when writing such programs.

During the transition to the multicore era (early 2000s), mainstream programming languages were not designed to support concurrent programming natively, which caused a rise in concurrency bugs. For example, C and C++, which were among the most popular programming languages when this transition happened [101], were specified as single-threaded languages [14], without reference to the semantics of threads.

Concurrency was added to these mainstream languages through libraries (e.g. Pthreads [39] and Windows threads [106]), which added informal constructs that developers could use to restrict access to shared variables (e.g., `pthread_mutex_lock`). These constructs were informal, because they did not change the nature of the C/C++ compilers that were inherently oblivious to concurrency.

Despite the presence of libraries attempting to add concurrency support to C/C++, associated compilers would generate code as if

```

1 unsigned x;
2 ...
3 if (x < 5) {
4     ... code that doesn't change x ...
5     switch (x) {
6         case 0:
7             ...
8         case 1:
9             ...
10        case 2:
11            ...
12        case 3:
13    }
14 }

```

Figure 1 – Example of a switch statement adapted from [12]

the programs were single-threaded, thereby occasionally violating the intended semantics of concurrent programs. For instance, consider the program snippet in Fig. 1, where a compiler could compile the program to emit a branch table for the switch statement and omit bounds check for the branch table because it already knows that $x < 5$. If the resulting program loads x twice once on line 3 and once on line 5, and x is modified between the two loads by another thread (i.e., there is a data race on x), the program may take a wild branch and will most probably crash.

Atomicity violations, data races, and deadlocks can all cause similar subtle behavior and cause software to fail in hard-to-predict circumstances [13]. Moreover, the subtle behavior such bugs complicate reasoning about concurrent programs.

It is challenging to fix concurrency bugs as it is, however if such failures only occur in production, the problem is exacerbated. This is because developers traditionally rely on reproducing failures in order to understand the associated bugs and fix them. However, if such bugs only recur in production and cannot be reproduced in-house, diagnosing the root cause and fixing them is truly hard. In [86], developers noted: “We don’t have tools for the once every 24 hours bugs in a 100 machine cluster.” An informal poll on Quora [85] asked “What is a coder’s worst nightmare,” and the most popular answer was “The bug only occurs in production and can’t be replicated locally,”.

To address these problems, this dissertation introduces techniques for the detection, root cause diagnosis, and classification of concurrency bugs that occur in production. We introduce techniques that are applicable to concurrency bugs at large, however we focus on concurrency bugs that occur in production, because such bugs present additional challenges as we describe in the next section (§1.2).

Intuitively, a root cause is the real reason behind a failure, we talk about root causes in detail in (§2).

1.2 CHALLENGES

Researchers and practitioners have observed that concurrency bugs are hard to detect and fix [36, 47, 53, 54, 86]. In this section, we first explain the fundamental challenges of the detection root cause diagnosis and classification of concurrency bugs, namely the runtime performance overhead challenge (§1.2.1) and the accuracy challenge (§1.2.2). We then elaborate on why performing these tasks is harder in production (§1.2.3).

1.2.1 The Runtime Performance Overhead Challenge

The runtime tracing that is required for the detection root cause diagnosis, and classification of concurrency bugs incurs high runtime performance overhead. In this section, we discuss the challenges that arise from runtime overheads of techniques and tools that perform dynamic program analysis, because purely static analysis has no runtime overhead.

Dynamic concurrency bug detection, whether it is the detection of data races, atomicity violations, or deadlocks, is expensive. This is because concurrency bug detection requires monitoring memory accesses and/or synchronization operations and performing intensive computations at runtime (e.g., building a happens-before relationship [57] graph for data race detection).

For instance, dynamic data race detection needs to monitor many memory accesses and synchronization operations, therefore it incurs high runtime overhead (as high as $200\times$ in industrial-strength tools like Intel Parallel Studio [41]). The lion's share of instrumentation overhead is due to monitoring memory reads and writes, reported to account for as much as 96% of all monitored operations [29].

Similarly, atomicity violation detectors incur high overheads (up to $45\times$ in the case of state-of-the-art detector AVIO-S [65] and up to $65\times$ in the case of SVD [111]). The overhead of atomicity violation detection stems from tracking updates to each monitored memory access and performing the necessary checks for determining whether a given access constitutes an atomicity violation or not.

With regards to the classification of concurrency bugs, prior work mostly focused on data race classification [45, 47, 52, 53, 77, 100], because data race detectors tend to report many data races. The abundance of data races pushes developers to understand which data races have higher impact, in order to prioritize their fixing.

From a programming languages standpoint, attempting to classify data races is only meaningful for some languages. One such language is the Java programming language. The Java memory model [69] defines semantics for programs with data races, because Java must support the execution of untrusted sandboxed code, and such code could

Although static analysis tools do not impose runtime overhead, they suffer from false positives, which is another challenge we discuss in (§1.2.2)

By classification, we mean the classification of true positives (i.e., real bugs). Identification of false positives (i.e., reports that do not correspond to real bugs) is considered separately in this dissertation

contain data races. Therefore, attempting to classify data races in Java is a meaningful endeavor from a programming languages point of view.

On the other hand, recent C [44] and C++ [43] standards do not provide meaningful semantics for programs involving data races. As a consequence, C and C++ compilers are allowed to perform optimizations on code with data races that may transform seemingly benign data races into harmful ones [13].

From a practical standpoint, developers may choose to prioritize the fixing of data races (and they do so) regardless of the implications of language standards. This happens primarily because modern multithreaded software tends to have a large number of data races, and it may be impractical to try to fix all the data races in a given program at once. For example, Google’s Thread Sanitizer [94] reports over 1,000 unique data races in Firefox when the browser starts up and loads <http://bbc.co.uk>.

Another reason why developers sometimes choose to not fix all data races is because synchronizing all racing memory accesses would introduce performance overheads that may be considered unacceptable. For example, developers have not fixed a race that can lead to lost updates in memcached for a year—ultimately finding an alternate solution—because, it leads to a 7% drop in throughput [72]. Performance implications led to 23 data races in Internet Explorer and Windows Vista being purposely left unfixed [77]. Similarly, several data races have been left unfixed in the Windows 7 kernel, because fixing those races did not justify the associated costs [47].

Classifying data races according to their potential consequences requires more computationally-intensive analyses than mere data race detection, and therefore imposes significant runtime overhead. For example, a state of the art data race classification tool, Record/Replay analyzer [77], incurs $45\times$ runtime overhead when performing data race classification. In order to classify data races based on their severity, not only data races need to be detected, but further analyses need to be enabled to monitor data races’ effects on the program state and output. Moreover, accurate classification of data races requires exploring multiple program paths and schedules to gain sufficient confidence in the classification results, and this further increases the runtime performance overhead.

Finally, root cause diagnosis of concurrency bugs requires tracking memory accesses and certain relations among memory accesses (e.g., their execution order), and therefore incur large runtime overhead.

The overhead of root cause diagnosis of concurrency bugs is further exacerbated because of the need to execute a program multiple times in order to isolate the failing thread schedules and inputs [68]. For example the *Delta Debugging* algorithm [20, 116]—a state of the art technique for isolating bug inducing inputs and thread schedules—

requires gathering execution information from several dozens (50 to 100) of runs before honing in on bugs' root causes. Another state of the art concurrency bug isolation technique CBI [46], reports overheads as high as $460\times$.

1.2.2 The Accuracy Challenge

Static detection of concurrency bugs works without actually running programs, therefore it does not incur any runtime performance overhead, however this comes at the expense of false positives (i.e., bug reports that do not correspond to actual bugs). False positives arise, because static analysis does not reason about the program's full runtime execution context.

False positives in static analysis of concurrency bugs arise because of four main reasons: first, static detectors perform some approximations such as conflating program paths during analysis or constraining the analysis to be intraprocedural (as opposed to interprocedural) in order to scale to large code bases. Second, static analyzers cannot always accurately infer which program contexts are multithreaded. Third, static analyzers typically model the semantics of lock/unlock synchronization primitives but not other primitives, such as barriers, semaphores, or wait/notify constructs. Finally, static analyzers cannot accurately determine whether two memory accesses alias or not.

Static classification of concurrency bugs typically relies on heuristics, and therefore inherently has false positives as is the case with most heuristic approach. For instance, DataCollider [47] prunes data race reports that appear to correspond to updates of statistics counters and to read-write conflicts involving different bits of the same memory word, or that involve variables known to developers to have intentional races (e.g., a "current time" variable is read by many threads while being updated by the timer interrupt). Updates on a statistics counter might be considered harmless for the cases investigated by DataCollider, but if a counter gathers critical statistics related to resource consumption in a language runtime, classifying a race on such a counter as harmless may be incorrect. More importantly, even data races that developers consider harmless may become harmful (e.g., cause a crash or a hang) for a different compiler and hardware combination [13].

Dynamic detectors and classifiers [38, 41, 94] tend to report fewer false positives. Developers prefer tools that have fewer false positives, because they do not have the time to cherry-pick true positives (i.e., reports corresponding to real bugs) in the presence of false positives [9].

Dynamic root cause diagnosis techniques [3, 4, 46, 56, 63, 88, 102] typically rely on statistical analysis for isolating the root causes of bugs, and therefore are susceptible to false positives. These tech-

We are not aware of static root cause diagnosis schemes. However, because root cause diagnosis shares common elements with detection and classification, we cautiously speculate that it will suffer

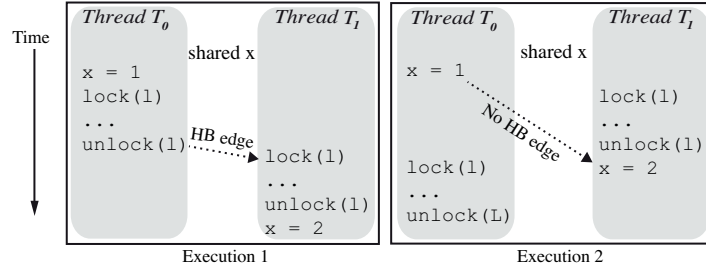


Figure 2 – False negatives in happened-before (HB) dynamic race detectors: the race on x is not detected in Execution 1, but it is detected in Execution 2.

niques gather execution information from multiple failing and successful executions to determine the key differences between those executions. The accuracy of statistical analysis hinges on the number of samples gleaned, therefore dynamic root cause diagnosis techniques can have false positives if they cannot monitor a sufficiently large sample of executions.

On the other hand of the spectrum are false negatives (i.e., real bug reports that are missed). False negatives can be an artifact of the approximations used in static analysis, or they may occur because a certain analysis (static or dynamic) cannot analyze a certain portion of the code.

False negatives are typical of dynamic detection, root cause diagnosis, and classification, because dynamic analysis can only operate on executions it witnesses, which is typically only a tiny subset of a program's possible executions.

False negatives also arise because of fortuitous events. For example, while monitoring a subset of executions, dynamic data race detectors may incorrectly infer happened-before relationships that are mere artifacts of the witnessed thread interleaving. To illustrate this point, consider Fig. 2. In execution 1, the accesses to the shared variable x are ordered by an accidental happened-before relationship (due to a fortuitous ordering of the acquire/release order of locks) that masks the true data race. Therefore, a precise dynamic detector would not flag this as a data race. However, this program does have a race, which becomes visible under a different schedule. This is shown in execution 2, where there is no happened-before relationship between accesses to x ; a precise dynamic detector would have reported a data race only if it witnessed this latter thread schedule.

An exception to this is data race detection using causal precedence [96], which can accurately predict data races that do not occur during actual program executions

1.2.3 The In-Production Challenge

Any in-production detection, classification, and root cause diagnosis technique needs to incur very low performance overhead and minimally perturb real-user executions. The overhead challenge (§1.2.1)

is exacerbated in production, because users will not tolerate performance degradation—even if it comes with increased reliability. Solutions that perturb the actual behavior of production runs nondeterministically may mask the bug frequently but not always, and thus make it harder to detect the bug and remove the potential for (even occasional) failure [74].

Moreover, a great challenge is posed by bugs that only recur in production and cannot be reproduced in-house. The ability to reproduce failures is essential for detecting, classifying and diagnosing the root causes of bugs. A recent study at Google [86] revealed that developers’ ability to reproduce bugs is crucial to fixing them. However, in practice, it is not always possible to reproduce bugs, and practitioners report that it takes weeks to fix hard-to-reproduce concurrency bugs [36].

Overall, detection, classification and root cause diagnosis of concurrency bugs pose significant challenges. In particular, it is hard to efficiently and accurately perform these tasks. These challenges are further exacerbated in production.

1.3 OVERVIEW OF PRIOR SOLUTION ATTEMPTS

In this section, we present an overview of prior attempts at addressing the challenges of detection, root cause diagnosis and classification of concurrency bugs. We briefly summarize how an existing technique attempts to address an aforementioned challenge. We also elaborate on the shortcomings of prior attempts.

1.4 ATTEMPTS AT ADDRESSING THE OVERHEAD CHALLENGE

In order to reduce the runtime performance overhead, certain dynamic concurrency bug detectors combine static analysis with dynamic detection. For example, Goldilocks [25] uses thread escape analysis [75] to reduce the set of memory accesses that needs to be monitored at runtime. A similar approach was proposed earlier by Choi et al. [21], using a variant of escape analysis. Certain approaches [1, 89] introduce a type system to reduce the overhead of data race and atomicity violation detection. Despite these assisting static analyses, existing concurrency bug detectors still incur overheads that make them impractical for in-production use.

Another way in which existing dynamic detectors and root cause diagnosis techniques address the overhead challenge is sampling. Sampling based concurrency bug detection and root cause diagnosis tracks synchronization operations whenever sampling is enabled. For instance, sampling-based data race detectors [15, 70] reduce runtime performance overhead. Although sampling reduces runtime overhead, this comes at the expense of false negatives: since these detec-

We haven't come across attempts to overcome the overhead challenge for classification of concurrency bugs.

tors do not monitor all runtime events, they inevitably miss certain bugs.

Another common way prior work copes with the overhead challenge of detection and root cause diagnosis of concurrency bugs is through the usage of customized hardware. HARD [118] uses special hardware support for data race detection. LBRA/LCRA [3] uses hardware extensions to diagnose root causes of bugs. These techniques indeed alleviate the overhead challenge. Alas, the hardware support they introduce is not implemented and deployed in the real world.

1.5 ATTEMPTS AT ADDRESSING THE ACCURACY CHALLENGE

In order to deal with false positives, dynamic tools employ filtering, which is typically unsound. Unsound filtering can filter out true positives along with false positives. Although this type of filtering reduces false positives, it cannot eliminate all of them. For example, even after attempting to filter out false data race reports, RacerX still has 37% - 46% false positives.

False negatives in concurrency bug detection can be trivially reduced (or even eliminated) by flagging more bug reports, but this will come at the expense of increased false positives. For instance, a data race detector could report a data race for every pair of memory accesses in a program. This strategy will eliminate all false negatives, but it will likely introduce a lot of false positives. Static concurrency bug detection tools (e.g., RELAY [103]) do not go to such extremes, nevertheless they rely on unsound techniques such as using inaccurate but fast alias analysis to flag as many bugs as possible (i.e., reduce false negatives), and consequently suffer from false positives (84%).

We explain hybrid data race detectors in detail in §2.2.4.3

We haven't come across attempts to overcome the accuracy challenge for classification of concurrency bugs.

Hybrid data race detectors [79] overcome the false negatives due to fortuitous happened-before relationships by combining two of the main dynamic data race detection algorithms, namely happened-before based data race detection and lockset-based data race detection. Although this combination reduces false negatives, it can introduce false positives due to the imprecise nature of lockset-based data race detection.

1.6 ATTEMPTS AT ADDRESSING THE IN-PRODUCTION CHALLENGE

Recall from §1.2.3 that the *in production challenge* exacerbates the *overhead challenge*, therefore prior work used similar methods to cope with the in-production challenge as it did for the overhead challenge. Below, we outline a few of the techniques that prior work used to deal with the in production challenge in addition to the techniques

used to cope with the overhead challenge (which we talked about in §1.4).

To alleviate the aggravated overhead challenge, prior work employs a variant of sampling. In particular, a common way prior work addresses the in-production challenge is through collaborative approaches like CCI [46] and CBI [63] that rely on monitoring executions at multiple user endpoints. There are two outstanding issues with collaborative approaches: although they reduce runtime overhead per user endpoint for which they perform detection or root cause diagnosis, the reduced overhead is still not suitable (up to $9\times$) for most in-production environments. Second, because these collaborative approaches sample a subset of the executions—in order to not impose overhead for every execution they monitor—they may miss rare failures that only recur in-production. This happens because sampling further reduces the probability of encountering failures that rarely recur in the first place.

1.7 SOLUTION OVERVIEW

In this section, we present an overview of the solution we propose to the problem we defined in §1.1.

In this dissertation, we address the challenge of in-production detection and root cause diagnosis of concurrency bugs by first employing deep static program analysis offline, and subsequently performing lightweight dynamic analysis online at user endpoints. Static analysis and dynamic analysis work synergistically in a feedback loop: static analysis reduces the overhead of ensuing dynamic analysis and dynamic analysis improves the accuracy of static analysis.

More specifically, with regards to data race detection, our key objective is to have a good data race detector that can be used (1) always-on and (2) in production. This is why we use static analysis to reduce the number of memory accesses that need to be monitored at runtime, thereby reducing overhead by up to two orders of magnitude compared to existing sampling-based techniques. Because we don't rely on sampling an execution during data race detection, our data race detection ends up being more accurate.

We attack the problem of detection of atomicity violations and root cause diagnosis of failures due to concurrency bugs using a technique we call failure sketching. Failure sketching is a technique that automatically produces a high level execution trace called the *failure sketch* that includes the statements that lead to a failure and the differences between the properties of failing and successful program executions. We show in this dissertation that these differences, which are commonly accepted as pointing to root causes [63, 88, 116], indeed point to the root causes of the failures we evaluated (§??). Identifying root

causes of failures also allows detecting the bugs that are associated with those failures.

Addressing the challenge of data race classification requires first addressing the challenge of in-production data race detection, which we do via hybrid static-dynamic analysis as we just mentioned.

We then do the classification entirely offline, because classification is a computationally expensive process: multiple program paths and schedules need to be explored in order to understand the consequences of a data race, and it is not possible to do such analyses in-production without incurring prohibitive runtime performance overheads or utilizing many more resources.

In this dissertation, we do not introduce new hardware mechanisms that would conveniently solve the aforementioned challenges. There are two key reasons for this: (1) not inventing our own custom hardware that solves a challenge we are facing allows us to come up with novel contributions (detailed below in §1.8); (2) the techniques we develop are broadly applicable, because they do not depend on a hardware feature that is not deployed in the real world.

1.8 SUMMARY OF CONTRIBUTIONS

This dissertation introduces **the first data race detector that can both be used always-on in production and provides good accuracy.**

Data race detection with low overhead has been a longstanding problem. Because data race detection is very costly, to our knowledge, prior work has not attempted to explore data race detection in-production. In this dissertation we tackle this problem via:

- A two-phase static-dynamic approach for detecting data races in real world software in a way that is more accurate than the state of the art.
- A new algorithm for dynamically detecting data races on-demand, which has lower overhead than state-of-the-art dynamic detectors, including those based on sampling.
- A crowdsourcing framework that, unlike traditional testing, taps directly into real user executions to detect data races.

The second contribution of this dissertation is **failure sketching, a low overhead, automated technique to automatically build failure sketches, which succinctly represent a failure’s root cause.**

Root cause diagnosis of in-production failures—especially failures due to concurrency bugs—has long been explored. To our knowledge there is no prior work that can perform root cause diagnosis of in-production failures without resorting to custom hardware or system state checkpointing infrastructure. In this dissertation, we achieve this via:

- A hybrid static-dynamic approach that combines in-house static program analysis with in-production collaborative and adaptive dynamic analysis.
- A first and practical demonstration of how Intel Processor Trace, a new technology that started shipping in early 2015 Broadwell processors [42], can be used to perform root cause diagnosis.

The third and final contribution of this dissertation is **a technique to automatically classify data races based on their potential consequences.**

Prior work on data race classification has not been accurate, either because it relied on heuristics or because the abstraction-level of the classification criteria was not correctly identified. In this dissertation, we solve the classification problem via:

- A four-category taxonomy of data races that is finer grain, more precise and, we believe, more useful than what has been employed by the state of the art.
- A technique for predicting the consequences of data races that combines multi-path and multi-schedule analysis with symbolic program-output comparison to achieve high accuracy in consequence prediction, and thus classification of data races according to their severity.
- Symbolic memory consistency modeling, a technique that can be used to model various architectural memory models in a principled way in order to perform data race classification under those memory models.

1.9 SUMMARY OF RESULTS

We built prototypes of all the techniques we present in the thesis and evaluate them. In this section, we give an overview of our evaluation results. Later in §6, we detail these evaluation results.

We evaluated RaceMob, our in-production data race detector on ten different systems, including Apache, SQLite, and Memcached. It found 106 real data races while incurring an average runtime overhead of 2.32% and a maximum overhead of 4.54%. Three of the data races hang SQLite, four data races crash Pbzp2, and one data race in Aget causes data corruption. Of all the 841 data race candidates found during the static detection phase, RaceMob labeled 77% as likely false positives. Compared to three state-of-the-art data race detectors [15, 94, 103] and two concurrency testing tools [53, 92], RaceMob has lower overhead and better accuracy than all of them.

We evaluated, Gist, our root cause diagnosis prototype using 11 failures from 7 different programs including Apache, SQLite, and Memcached. The Gist prototype managed to automatically build failure sketches with an average accuracy of 96% for all the failures while in-

curing an average performance overhead of 3.74%. On average, Gist incurs $166\times$ less runtime performance overhead than a state-of-the-art record/replay system.

We evaluated our data race classification prototype Portend, we applied Portend to 93 data race reports from 7 real-world applications—it classified 99% of the detected data races accurately in less than 5 minutes per race on average. Compared to state-of-the-art race classifiers, Portend is up to 89% more accurate in predicting the consequences of data races (§6.3.7). This improvement comes from Portend’s ability to perform multi-path and multi-thread schedule analysis, as well as Portend’s fine grained classification scheme. We found not only that multi-path multi-schedule analysis is critical for high accuracy, but also that the “post-race state comparison” approach used in state-of-the-art classifiers does not work well on our real-world programs, despite being perfect on simple microbenchmarks (§6.3).

BACKGROUND AND RELATED WORK

In this chapter, we first briefly review surveys that examine concurrency bug characteristics, we then review the literature on the detection, classification and root cause diagnosis of concurrency bugs.

2.1 CONCURRENCY BUG SURVEYS

Failure recovery study collected only 12 bugs [18].

As concurrent programming gained more momentum, concurrency bugs proliferated and new surveys were made.

The first comprehensive study is due to Shan Lu [66, 64].

Another study conducted at Microsoft [36] revealed interesting results: 72% considers these bugs as hard to very hard. Respondents say that fixing such bugs can take days (63.4) to weeks (8.3) to months (0.9). 66% responders deal with these bugs. Total of 684 people answered the survey. Most bugs are of high severity (1 or 2) 84%.

A more recent study [87] in the context of root cause diagnosis of bugs determined that although most of the bugs can be reproduced in-production by running the program with the same set of inputs (82%), the remainder of the bugs had non-deterministic behavior. One of the conclusions of this study was that determining fault-triggering inputs for concurrency bugs and reproducing failures due to concurrency bugs is significantly harder than for other bugs.

2.2 DATA RACE DETECTION

2.2.1 Data Race Definition

Two memory accesses are conflicting if they access a shared memory location and at least one of the two accesses is a write. A data race occurs when two threads make a conflicting access, and these accesses are not ordered by a *happened-before* relationship [57]— if memory effects of an operation O_1 in a process P_1 becomes visible to a process P_2 before P_2 performs O_2 , we say O_1 happened before O_2 .

The terms *data race* and *race* are often incorrectly used interchangeably. There is a subtle yet important distinction between these terms that has garnered attention from both the academic [48] and practical [6] community.. A data race is a condition which can be formally specified as in the previous paragraph. This formal specification al-

lows the detection of a data race to be automated. A race condition on the other hand is a flaw

2.2.2 *Data Race Detection Literature*

Detection of data races garnered a lot of attention, because it is actually possible to detect data races with high accuracy. However, it is not possible to build a meaningful atomicity violation detector or an order violation detector that does not have false positives. The reason behind this is that data races are precisely defined programming errors, whereas atomicity and order violations are high-level semantic bugs that do not have a universally accepted definition that would allow their automated detection with high accuracy.

Data race detection can be broadly classified into three classes: static data race detection, dynamic data race detection, and mixed static-dynamic data race detection.

2.2.3 *Static Data Race Detection*

Static data race detectors [26, 103, 50, 83, 76, 51, 30] analyze the program source code without executing it. They reason about multiple program paths at once, and thus typically do not miss data races (i.e., have low rate of false negatives) [79]. Static detectors also run fast and scale to large code bases. The problem is that static data race detectors tend to have many false positives, i.e., produce reports that do not correspond to real data races (e.g., 84% of data races reported by RELAY are not true data races [103]). This can send developers on a wild goose chase, making the use of static detectors potentially frustrating and expensive.

Static data race detectors employ various strategies:

rccjava relies on the type system .

2.2.4 *Dynamic Data Race Detection*

Dynamic data race detectors [41, 94] typically monitor memory accesses and synchronization at runtime, and determine if the observed accesses race with each other. Such detectors can achieve low rates of false positives. Alas, dynamic detectors miss all the data races that are not seen in the directly observed execution (i.e., they have false negatives), and these can be numerous. The instrumentation required to observe all memory accesses makes dynamic detectors incur high runtime overheads (200× for Intel Thread Checker [41], 30× for Google ThreadSanitizer [94]). As a result, dynamic detectors are not practical for in-production use, rather only during testing—this deprives them of the opportunity to observe real user executions, thus missing data races that only occur in real user environments. Some dynamic data

race detectors employ sampling [15, 70] to decrease runtime overhead, but this comes at the cost of further false negatives.

Below, we present a policy-centric classification of data race detection algorithms. Whenever applicable, we also talk about various mechanisms with which these data race detection policies are implemented.

2.2.4.1 *Happened Before Relationship-Based Data Race Detection*

Happened-before relationship [57] based-detectors [91, 94, 38, 53, 52] track the happened-before relationships among memory accesses for a program execution and based on those relationships they detect data races.

More specifically, if two memory accesses access the same memory location, at least one of the accesses is a write and there is no happened-before relationship between the two accesses, these detectors flag a data race.

Dynamic detectors that use happens-before relationships do not have false positives as long as they are aware of the synchronization mechanisms employed by the developer.

2.2.4.2 *Lockset-Based Data Race Detection*

Locksets describe the locks held by a program at any given point in the execution. Consider the example program on Figure where steps in the execution of two threads are enumerated along the flow of time and time flows downward.

Lockset-based data race detection [90, 118] checks whether all shared memory accesses follow a consistent *locking discipline*. A locking discipline is a policy that ensures the absence of data races. A trivial locking discipline would require all memory accesses in a program to be protected by the same lock.

The simple locking discipline that Eraser [90], namely the first lockset-based data race detector, uses states that every shared variable access should be protected by some lock. In other words, Eraser requires any thread accessing a given shared variable to hold a common lock while it is performing the access.

Lockset-based data race detection works as follows:

Eraser’s simple locking discipline is overly strict, and as a result, it would cause Eraser to report many false positives.

To lower the number of false positives, Eraser employs several refinements. First, Eraser will not report data races due to the initialization of data races, which is frequently done without holding a lock.

HARD [118] is a hardware implementation of the lockset-based data race detection. HARD uses bloom filters [bloom:filter] to store the lockset and uses bitwise logic operations to implement the aforementioned set operations. HARD was able to detect 54 data races out

of 60 randomly-injected data races in six SPLASH-2 applications (20% more than happens-before based data race detection) with overheads ranging between 0.1% to 2.6%.

2.2.4.3 Hybrid Data Race Detection

This is as opposed to calling the combination of static and dynamic data race detection as hybrid.

Perhaps unfortunately, in the literature of data race detection, hybrid data race detection implies data race detection that combines two major dynamic data race detection algorithms, namely dynamic data race detection and static data race detection. This is because the first piece of work that called a data race detection algorithm “hybrid” [79] combined these two major dynamic data race detection algorithms.

The concept of hybrid data race detection idea is due to , however the first implementation of a hybrid data race detector is Eraser.

Hybrid data race detectors [79, 109, 115] combine lockset-based data race detection with happened before relationship-based data race detection.

Hybrid data race detection improves the accuracy of data race detection by reporting fewer false positives than lockset-based data race detection and fewer false negatives than happened-before data race detection [79].

2.2.5 Mixed Static-Dynamic Data Race Detection

The first data race detector to combine

2.2.6 Detecting Data Races In Production

Litecollider [10] and RaceMob [54]. RaceMob is the first to do so to our knowledge.

2.3 ATOMICITY AND ORDER VIOLATION DETECTION

Velodrome [31]

Exposing concurrency bugs is also in this category. SKI

DoubleChecker, AVIO, AtomAid AtomTracker, AtomFuzzer

2.4 DATA RACE CLASSIFICATION

Data race detectors (both dynamic and static), report a lot of data races. There is a need to classify these data races

2.5 ROOT CAUSE DIAGNOSIS OF CONCURRENCY BUGS

2.5.1 *Delta Debugging*

2.5.2 *Statistical Techniques*

Part II

ELIMINATING CONCURRENCY BUGS FROM IN-PRODUCTION SYSTEMS

In this part, we describe the design, implementation, and evaluation of the techniques and tools we developed in order to detect, perform root cause diagnosis of, and classify in-production concurrency bugs.

First, we present a technique to accurately detect data races in-production. We then present a general technique for diagnosing the root causes of in-production failures that may be caused by concurrency bugs, failing inputs, or rare program paths—although, we primarily focus on failures caused by concurrency bugs. Our technique for root cause diagnosis allows both detecting bugs and providing an explanation of how the failure happened. Then, we discuss how we can perform data race classification under arbitrary memory models. Finally, we present a comprehensive evaluation of all the prototypes we developed.

RACEMOB: DETECTING DATA RACES IN PRODUCTION

In this section, we present RaceMob, a way to combine static and dynamic data race detection to obtain both good accuracy and low runtime overhead. For a given program P , RaceMob first uses a static detection phase with few false negatives to find potential data races; in a subsequent dynamic phase, RaceMob crowdsources the validation of these alleged data races to user machines that are anyway running P . RaceMob provides developers with a dynamically updated list of data races, split into “confirmed true races”, “likely false positives”, and “unknown”—developers can use this list to prioritize their debugging attention. To minimize runtime overhead experienced by users of P , RaceMob adjusts the complexity of data race validation on-demand to balance accuracy and cost. By crowdsourcing validation, RaceMob amortizes the cost of validation and (unlike traditional testing) gains access to real user executions. RaceMob also helps discovering user-visible failures like crashes or hangs, and therefore helps developers to reason about the consequences of data races. To the best of our knowledge, RaceMob is the first data race detector that combines sufficiently low overhead to be always-on with sufficiently good accuracy to improve developer productivity.

3.1 RACEMOB’S DESIGN OVERVIEW

RaceMob is a crowdsourced, two-phase static–dynamic data race detector. It first statically detects potential data races in a program, then crowdsources the dynamic task of validating these potential data races to users’ sites. This validation is done using an on-demand data race detection algorithm. The benefits of crowdsourcing are twofold: first, data race validation occurs in the context of real user executions; second, crowdsourcing amortizes the per-user validation cost. Data race validation confirms true data races, thereby increasing the data race detection coverage.

The usage model is presented in Fig. 3. First, developers set up a “hive” service for their program P ; this hive can run centralized or distributed. The hive performs static data race detection on P and finds potential data races (§3.2); these go onto P ’s list of data races maintained by the hive, and initially each one is marked as “Unknown”. Then the hive generates an instrumented binary P' , which users download ① and use instead of the original P . The instrumentation in P' is commanded by the hive, to activate the validation of

We define data race detection coverage as the ratio of true data races found in a program by a detector to the total number of true data races in that program.

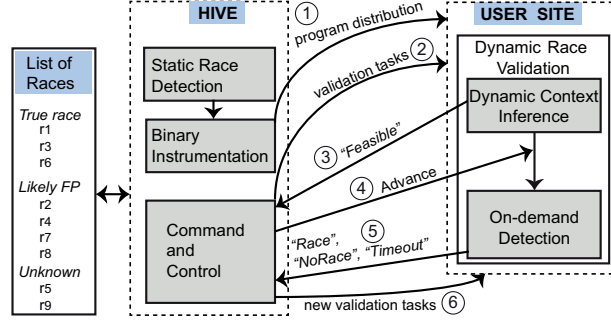


Figure 3 – RaceMob’s crowdsourced architecture: A static detection phase, run on the hive, is followed by a dynamic validation phase on users’ machines.

specific data races in P ②; different users will typically be validating different, albeit potentially overlapping, sets of data races from P (§3.3). The first phase of validation, called dynamic context inference (§3.3.1), may decide that a particular racing interleaving for data race r is feasible, at which point it informs the hive ③. At this point, the hive instructs all copies of P' that are validating r to advance r to the second validation phase ④. This second phase runs RaceMob’s on-demand detection algorithm (§3.3.2), whose result can be one of Race, NoRace, or Timeout ⑤. As results come in to the hive, it updates the list of data races: if a “Race” result came in from the field for data race r , the hive promotes r from “Unknown” to “True Race”; the other answers are used to decide whether to promote r from “Unknown” to “Likely False Positive” or not (§3.5). For data races with status “Unknown” or “Likely False Positive,” the hive redistributes “validation tasks” ⑥ among the available users (§3.4). We now describe each step in further detail.

3.2 STATIC DATA RACE DETECTION

RaceMob can use any static data race detector, regardless of whether it is complete or not. We chose RELAY, a lockset-based data race detector [103]. Locksets describe the locks held by a program at any given point in the program §2.2.4.2. RELAY performs its analysis bottom-up through the program’s control flow graph while computing function summaries that summarize which variables are accessed and which locks are held in each function. RELAY then composes these summaries to perform data race detection: it flags a data race whenever it sees at least two accesses to memory locations that are the same or may alias, and at least one of the accesses is a write, and the accesses are not protected by at least one common lock.

RELAY is complete (i.e., does not miss data races) if the program does not have inline assembly and does not use pointer arithmetic. RELAY may become incomplete if configured to perform file-based

alias analysis or aggressive filtering, but we disable these options in RaceMob. As suggested in [60], it might be possible to make RELAY complete by integrating program analysis techniques for assembly code [5] and by handling pointer arithmetic [105].

Based on the data race reports from RELAY, RaceMob instruments the suspected-racing memory accesses as well as all synchronization operations in the program. This instrumentation will later be commanded (in production) by RaceMob to perform on-demand data race detection.

The hive activates parts of the instrumentation on-demand when the program runs, in different ways for different users. The activation mechanism aims to validate as many data races as possible by uniformly distributing the validation tasks across the user population.

3.3 DYNAMIC DATA RACE VALIDATION

The hive instructs the instrumented programs for which memory accesses to perform data race validation. The validation task sent by the hive to the instrumented program consists of a data race candidate to validate and one desired order (of two possible) of the racing accesses. We call these possible orders the *primary* and the *alternate*, borrowing terminology from our earlier work [77].

The dynamic data race validation phase has three stages: dynamic context inference (§3.3.1), on-demand data race detection (§3.3.2), and schedule steering (§3.3.3). Instrumentation for each stage is present in all the programs, however stages 2 and 3 are toggled on/off separately from stage 1, which is always on. Next, we explain each stage in detail.

3.3.1 *Dynamic Context Inference*

Dynamic context inference (DCI) is a lightweight analysis that partly compensates for the inaccuracies of the static data race detection phase. RaceMob performs DCI to figure out whether the statically detected data races can occur in a dynamic program execution context.

DCI validates two assumptions made by the static data race detector about a race candidate. First, the static detector’s abstract analysis hypothesizes aliasing as the basis for some of its race candidates, and DCI looks for concrete instances that can validate this hypothesis. Second, the static detector hypothesizes that racing accesses are made by different threads, and DCI aims to validate this as well. Once these two hypotheses are confirmed, the user site communicates this to the hive, and the hive promotes the race candidate to the next phase.

Without a confirmation from DCI, the race remains in the “Unknown” state.

The motivation for DCI comes from our observation that the majority of the potential data races detected by static data race detection (53% in our evaluation) are false positives due to only alias analysis inaccuracies and the inability of static race detection to infer multi-threaded program contexts.

For every potential data race r with racing instructions r_1 and r_2 , made by threads T_1 and T_2 , respectively, DCI determines whether the potentially racing memory accesses to addresses a_1 and a_2 made by r_1 and r_2 , respectively, may alias with each other (i.e., $a_1 = a_2$), and whether these accesses are indeed made by different threads (i.e., $T_1 \neq T_2$). To do this, DCI keeps track of the address that each potentially racing instruction accesses, along with the accessing thread’s ID at runtime. Then, the instrumentation checks to see if *at least one* pair of accesses is executed. If yes, the instrumented program notifies the hive, which promotes r to the next stages of validation (on-demand data race detection and schedule steering) on all user machines where r is being watched. If no access is executed by any instrumented instance of the program, DCI continues watching r ’s potentially racing memory accesses until further notice.

DCI has negligible runtime overhead (0.01%) on top of the binary instrumentation overhead (0.77%); therefore, it is feasible to have DCI always-on. DCI’s memory footprint is small: it requires maintaining 12 bytes of information per potential racing instruction (8 bytes for the address, 4 bytes for the thread ID). DCI is sound because, for every access pair that it reports as being made from different threads and to the same address, DCI has clear concrete evidence from an actual execution. DCI is of course not guaranteed to be complete.

3.3.2 On-Demand Data Race Detection

In this section, we explain how on-demand data race detection works; for clarity, we restrict the discussion to a single potential data race.

On-demand race detection starts tracking happens-before relationships once the first potentially racing access is made, and it stops tracking once a happens-before relationship is established between the first accessing thread and all the other threads in the program (in which case a “NoRace” result is sent to the hive). Tracking also stops if the second access is made before such a happens-before relationship is found (in which case a “Race” result is sent to the hive).

Intuitively, RaceMob tracks a minimal number of accesses and synchronization operations. RaceMob needs to track both racing accesses to validate a potential data race. However, RaceMob does not need

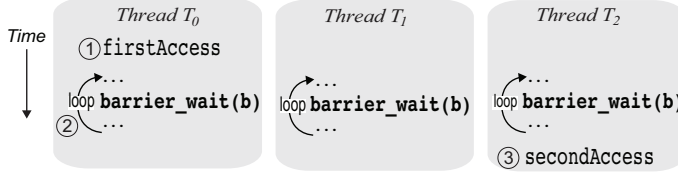


Figure 4 – Minimal monitoring in DCI: For this example, DCI stops tracking synchronization operations as soon as each thread goes once through the barrier.

to track any memory accesses other than the target racing accesses, because any other access is irrelevant to this data race.

Sampling-based data race detection (e.g., PACER [15]) adopts a similar approach to on-demand race detection by tracking synchronization operations whenever sampling is enabled. The drawback of PACER’s approach is that it may start to track synchronization operations too soon, even if the program is not about to execute a racing access. RaceMob avoids this by turning on tracking synchronization operations on-demand, when an access reported by the static race detection phase is executed.

RaceMob tracks synchronization operations—and thus, happens-before relationships—using an efficient, dynamic, vector-clock algorithm similar to DJIT⁺ [82]. We maintain vector clocks for each thread and synchronization operation, and the clocks are partially ordered with respect to each other.

We illustrate in Fig. 4 how the on-demand race detection stops tracking synchronization operations, using a simple example derived from the fmm program [95]: *firstAccess* executes in the beginning of the program in T_0 ①, and the program goes through a few thousand iterations of synchronization-intensive code ②. Finally, T_2 executes *secondAccess* ③. It is sufficient to keep track of the vector clocks of all threads only up until the first time they go through the *barrier_wait* statement, as this establishes a happens-before relationship between *firstAccess* in T_0 and any subsequent access in T_1 and T_2 . Therefore, on-demand data race detection stops keeping track of the vector clocks of threads T_0 , T_1 , and T_2 after they each go through *barrier_wait* once.

RaceMob distinguishes between a static racing instruction in the program and its dynamic instance at runtime, and it can enable on-demand data race detection for any dynamic instance of a potential racing access. However, practice shows that going beyond the first dynamic instances adds little value (§6.1).

Our experimental evaluation shows that on-demand data race detection reduces the overhead of dynamic race detection (§??).

3.3.3 Schedule Steering

The schedule steering phase further improves RaceMob’s data race detection coverage by exploring both the primary and the alternate executions of potentially racing accesses. This has the benefit of detecting races that may be hidden by accidental happens-before relationships (as discussed in §?? and Fig. 2).

Schedule steering tries to enforce the order of the racing memory accesses provided by the hive, i.e., either the primary or the alternate. Whenever the intended order is about to be violated (i.e., the undesired order is about to occur), RaceMob pauses the thread that is about to make the access, by using a wait operation with a timeout τ , to enforce the desired order. Every time the hive receives a “Timeout” from a user, it increments τ for that user (up to a maximum value), to more aggressively steer it toward the desired order, as described in the next section.

Prior work [77, 92, 53] used techniques similar to schedule steering to detect whether a *known* data race can cause a failure or not. RaceMob, however, uses schedule steering to increase the likelihood of encountering a suspected race and to improve data race detection coverage.

Our evaluation shows that schedule steering helps RaceMob to find two data races (one in Memcached and the other one in Pfscan) that would otherwise be missed. It also helps RaceMob uncover failures (e.g., data corruption, hangs, crashes) that occur as a result of data races, thereby enabling developers to reason about the consequences of data races and fix the important ones early, before they affect more users. However, users who do not wish to help in determining the consequences of data races can easily turn off schedule steering. We discuss the trade-offs involved in schedule steering in §??.

3.4 CROWDSOURCING THE VALIDATION

Crowdsourcing is defined as the practice of obtaining needed services, ideas, or content by soliciting contributions from a large group of people and especially from the online community. RaceMob gathers validation results from a large base of users and merges them to come up with verdicts for potential data races.

RaceMob’s main motivation for crowdsourcing data race detection is to access real user executions. This enables RaceMob, for instance, to detect the important but often overlooked class of input-dependent data races [53], i.e., races that occur only when a program is run with a particular input. RaceMob found two such races in Aget, and we detail them in §??. Crowdsourcing also enables RaceMob to leverage many executions to establish statistical confidence in the detection verdict. We also believe crowdsourcing is more applicable today than

ever: collectively, software users have overwhelmingly more hardware than any single software development organization, so leveraging end-users for race detection is particularly advantageous. Furthermore, such distribution helps reduce the per-user overhead to barely noticeable levels.

Validation is distributed across a population of users, with each user receiving only a small number of races to validate. The hive distributes validation tasks, which contain the locations in the program of two memory accesses suspected to be racing, along with a particular order of these accesses. Completing the validation task consists of confirming, in the end-user's instance of the program, whether the indicated ordering of the memory accesses is possible. If a user site receives more than one race to validate, it will first validate the race whose racing instruction is first reached during execution.

There exists a wide variety of possible assignment policies that can be implemented in RaceMob. By default, if there are more users than races, RaceMob initially randomly assigns a single race to each user for validation. Assigned validation tasks that fail to complete within a time bound are then distributed to additional users as well, in order to increase the probability of completing them. Such multiple assignment could be done from the very beginning, in order to reach a verdict sooner. The number of users asked to validate a race r could be based, for example, on the expected severity of r , as inferred based on heuristics or the static analysis phase. Conversely, in the unlikely case that there are more data races to validate than users, the default policy is to initially distribute a single validation task to each user, thereby leaving a subset of the races unassigned. As users complete their validation tasks, RaceMob assigns new tasks from among the unassigned races. Once a data race is confirmed as a true race, it is removed from the list of data races being validated, for all users.

During schedule steering, whenever a race candidate is "stubborn" and does not exercise the desired order despite the wait introduced by the instrumentation, a "Timeout" is sent to the hive. The hive then instructs the site to increase the timeout τ , to more aggressively favor the alternate order; every subsequent timeout triggers a new "Timeout" response to the hive and a further increase in τ . Once τ reaches a configured upper bound, the hive instructs the site to abandon the validation task. At this point, or even in parallel with increasing τ , the hive could assign the same task to additional users.

There are two important steps in achieving low overhead at user sites. First, the timeout τ must be kept low. For example, to preserve the low latency of interactive applications, RaceMob uses an upper bound $\tau \leq 200$ ms; for I/O bound applications, higher timeouts can be configured. Second, the instrumentation at the user site disables schedule steering for a given race after a first steering attempt for a given race, regardless of whether it succeeded or not; this is par-

ticularly useful when the racing accesses are in a tight loop. Steering is resumed when a new value of τ comes in from the hive. It is certainly possible that a later dynamic instance of the potentially racing instruction might be able to exercise the desired order, had steering not been disabled; nevertheless, in our evaluation we found that RaceMob achieves higher accuracy than state-of-the art data race detectors.

RaceMob monitors each user’s validation workload and aims to balance the global load across the user population. Rebalancing does not require users to download a new version of the program, but rather the hive simply toggles validation on/off at the relevant user sites. In other words, each instance of the instrumented program P' is capable of validating, on demand, any of the races found during the static phase—the hive instructs it which race(s) is/are of interest to that instance of P' .

If an instance of P' spends more time doing data race validation than the overall average, then the hive redistributes some of that instance’s validation tasks to other instances. Additionally, RaceMob reshuffles tasks whenever one program experiences more timeouts than the average. In this way, we reduce the number of outliers, in terms of runtime overhead, during the dynamic phase.

Crowdsourcing offers RaceMob the opportunity to tap into a large number of executions, which makes it possible to only perform a small amount of monitoring per user site without harming the precision of detection. This in turn reduces RaceMob’s runtime overhead, making it more palatable to users and easier to adopt.

3.5 REACHING A VERDICT

TRUE RACE RaceMob decides a race candidate is a true race whenever both the primary and the alternate orders are executed at a user site, or when either of the orders is executed with no intervening happens-before relationship between the corresponding memory accesses. Among the true races, some can be specification-violating races in the sense of [53] (e.g., that cause crash or deadlock). In the case of a crash, the RaceMob instrumentation catches the SIGSEGV signal and submits a crash report to the hive. In the case of an unhandled SIGINT (e.g., the user pressed Ctrl-C), RaceMob prompts the user with a dialog asking whether the program has failed to meet expectations. If yes, the hive is informed that the enforced schedule leads to a specification violation. Of course, users who find this kind of “consequence reporting” too intrusive can disable schedule steering altogether.

LIKELY FALSE POSITIVE RaceMob concludes that a potential race is likely a false positive if at least one user site reported a NoRace re-

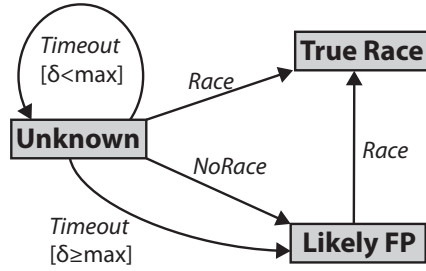


Figure 5 – The state machine used by the hive to reach verdicts based on reports from program instances. Transition edges are labeled with validation results that arrive from instrumented program instances; states are labeled with RaceMob’s verdict.

sult to the hive (i.e., on-demand race detection discovered a happens-before relationship between the accesses in the primary or alternate). RaceMob cannot provide a definitive verdict on whether the race is a false positive or not, because there might exist some other execution in which the purported false positive proves to be a real race (e.g., due to an unobserved input dependency). The “Likely False Positive” verdict, especially if augmented with the number of corresponding NoRace results received at the hive, can help developers decide whether to prioritize for fixing this particular race over others. RaceMob continues validation for “Likely False Positive” data races for as long as the developers wishes.

UNKNOWN As long as the hive receives no results from the validation of a potential race r , the hive keeps the status of the race “Unknown”. Similarly, if none of the program instances report that they reached the maximum timeout value, r ’s status remains “Unknown”. However, if at least one instance reaches the maximum timeout value for r , the corresponding report is promoted to “Likely False Positive”.

The “True Race” verdict is definite: RaceMob has proof of the race occurring in a real execution of the program. The “Likely False Positive” verdict is probabilistic: the more NoRace or Timeout reports are received by the hive as a result of distinct executions, the higher the probability that a race report is indeed a false positive, even though there is no precise probability value that RaceMob assigns to this outcome.

3.6 RACEMOB'S IMPLEMENTATION DETAILS

RaceMob can use any data race detector that outputs data race candidates; preferably it should be complete (i.e., not miss data races). We use RELAY, which analyzes programs that are turned into CIL, an intermediate language for C programs [78]. The instrumentation engine at the hive is based on [59]. We wrote a 500-LOC plugin that

converts RELAY reports to the format required by our instrumentation engine.

The instrumentation engine is an LLVM static analysis pass. It avoids instrumenting empty loop bodies that have a data race on a variable in the loop condition (e.g., of the form `while(notDone){}`). These loops occur often in ad-hoc synchronization [110]. Not instrumenting such loops avoids excessive overhead that results from running the instrumentation frequently. When such loops involve a data race candidate, they are reported by the hive directly to developers. We encountered this situation in two of the programs we evaluated, and both cases were true data races (thus validating prior work that advocates against ad-hoc synchronization [110]), so this optimization did not effect RaceMob’s accuracy.

Whereas Fig. ?? indicates three possible results from user sites (Race, NoRace, and Timeout), our prototype also implements a fourth one (NotSeen), to indicate that a user site has not witnessed the race it was expected to monitor. Technically, NotSeen can be inferred by the hive from the absence of any other results. However, for efficiency purposes, we have a hook at the exit of `main`, as well as in the signal handlers, that send a NotSeen message to the hive whenever the program terminates without having made progress on the validation task.

RaceMob uses compiler-based instrumentation, but other approaches are also possible. For example, we plan to use dynamic binary rewriting, which would also allow us to dynamically remove instrumentation for data races that are not enabled in a given instance of the program. The instrumentation is for the most part inactive; the hive activates part of the instrumentation on-demand, when the program runs.

ROOT CAUSE DIAGNOSIS OF IN-PRODUCTION FAILURES USING FAILURE SKETCHES

In this chapter, we present failure sketching, a technique that automatically produces a high level execution trace called the *failure sketch* that includes the statements that lead to a failure and the differences between the properties of failing and successful program executions.

We show in our evaluation that the differences between failing and successful executions which are displayed on the failure sketch, point to root causes [63, 88, 116], of failures the bugs we evaluated (§??). In the context of our work, we are talking about a statistical definition of a root cause: we consider events that are primarily correlated with a failure as the root causes of that failure. g Root cause diagnosis allows detecting bugs as well as providing an explanation of how the failure associated with the bug occurred: which branches were taken, which data values were computed, which thread schedule led to the failure, etc. After all, understanding how a failure occurred requires detecting what the failure is. Therefore, in the rest of this chapter, when we refer to root cause diagnosis, we imply both detection of a bug and an explanation of how a given bug led to a failure.

Failure sketching is a general technique for root cause diagnosis of concurrency bugs as well as some sequential bugs due to failing input values or rare paths that a program takes. Despite this generality, in the context of this dissertation, and especially when evaluating failure sketching, we focus on concurrency bugs.

Fig. 6 shows an example failure sketch for a bug in Pbzip2, a multithreaded compression tool [33]. Time flows downward, and execution steps are enumerated along the flow of time. The failure sketch shows the statements (in this case statements from two threads) that affect the failure and their order of execution with respect to the enumerated steps (i.e., the control flow). The arrow between the two statements in dotted rectangles indicates the difference between failing executions and successful ones. In particular, in failing executions, the statement `f->mut` from T_1 is executed before the statement `mutex_unlock(f->mut)` in T_2 . In non-failing executions, `cons` returns before `main` sets `f->mut` to `NULL`. The failure sketch also shows the value of the variable `f->mut` (i.e., the data flow) in a dotted rectangle in step 7, indicating that this value is 0 in step 7 only in failing runs. A developer can use the failure sketch to fix the bug by introducing proper synchronization that eliminates the offending thread schedule. This is exactly how pbzip2 developers fixed this bug, albeit four months after it was reported [33].

Root cause diagnosis of a failure subsumes the detection of the bug causing the failure.

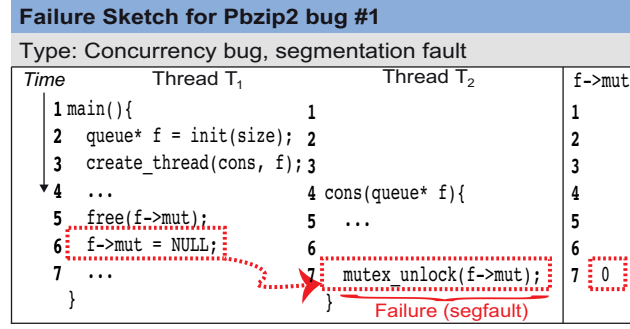


Figure 6 – The failure sketch of pbzip2 bug.

The insight behind the work presented in the chapter is that failure sketches can be built automatically, by using a combination of static program analysis and cooperative dynamic analysis. The use of a brand new hardware feature in Intel CPUs helps keep runtime performance overhead low (3.74% in our evaluation).

We built a prototype, Gist, that automatically generates a failure sketch for a given failure. Given a failure, Gist first statically computes a program slice that contains program statements that can potentially affect the program statement where the failure manifests itself. Then, Gist performs data and control flow tracking in a cooperative setup by targeting either multiple executions of the same program in a data center or users who execute the same program. Gist uses hardware watchpoints to track the values of data items in the slice, and uses Intel Processor Trace [42] to trace the control flow.

Although Gist relies on Intel Processor Trace and hardware watchpoints for practical and low-overhead control and data flow tracking, Gist’s novelty is in the combination of static program analysis and dynamic runtime tracing to judiciously select how and when to trace program execution in order to best extract the information for failure sketches.

In the rest of this chapter, we describe the overview of Gist’s design (§4.1). We then explain each component of Gist’s design in detail: we first describe static slicing in Gist (§4.4), we then explain how Gist performs slice refinement (§4.3), and we finally describe how Gist identifies the root cause of a failure (§4.4).

4.1 GIST’S DESIGN OVERVIEW

Gist, our system for building failure sketches has three main components: the static analyzer, the failure sketch computation engine, and the runtime that tracks production runs. The static analyzer and the failure sketch computation engine constitutes the server side of Gist, and they can be centralized or distributed, as needed. The runtime constitutes the client-side of Gist, and it runs in a cooperative

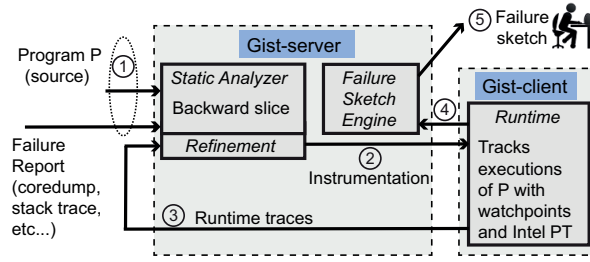


Figure 7 – The architecture of Gist

setting such as in a data center or in multiple users' machines, similar to RaceMob [54].

The usage model of Gist is shown in Fig. 13. Gist takes as input a program (source code and binary) and a failure report ① (e.g., stack trace, the statement where the failure manifests itself, etc). Gist, being a developer tool has access to the source code. Using these inputs, Gist computes a backward slice [104] by computing the set of program statements that potentially affect the statement where the failure occurs. Gist uses an interprocedural, path-insensitive and flow-sensitive backward slicing algorithm. Then, Gist instructs its runtime, running in a data center or at user endpoints, ② to instrument the programs and gather more traces (e.g., branches taken and values computed at runtime). Gist then uses these traces to refine the slice ③: refinement removes from the slice the statements that don't get executed during the executions that Gist monitors, and it adds to the slice statements that were not identified as being part of the slice initially. Refinement also determines the inter-thread execution order of statements accessing shared variables and the values that the program computes. Refinement is done using hardware watchpoints for data flow and Intel Processor Trace (Intel PT) for control flow. Gist's failure sketch engine gathers execution information from failing and successful runs ④. Then, Gist determines the differences between failing and successful runs and builds a failure sketch ⑤ for the developer to use.

Gist operates in a cooperative setting [54, 63] where multiple instances of the same software execute in a data center or in multiple users' machines. Gist's server side (e.g., a master node) performs offline analysis and distributes instrumentation to its client side (e.g., a node in a data center). Gist incurs low overhead, so it can be kept always-on and does not need to resort to sampling an execution [63], thus avoiding missing information that can increase root cause diagnosis latency.

Gist operates iteratively: the instrumentation and refinement continues as long as developers see fit, continuously improving the accuracy of failure sketches. Gist generates a failure sketch after a first failure using static slicing. Our evaluation shows that, in some cases,

this initial sketch is sufficient for root cause diagnosis, whereas in other cases refinement is necessary (§??).

We now describe how each component of Gist works and explain how they solve the challenges presented in the previous section (§1.2). We first describe how Gist computes the static slice followed by slice refinement via adaptive tracking of control and data flow information. Then we describe how Gist identifies the root cause of a failure using multiple failing and successful runs.

4.2 STATIC SLICE COMPUTATION

Gist uses an interprocedural, path-insensitive and flow-sensitive backward slicing algorithm to identify the program statements that may affect the statement where the failure manifests itself at runtime. We chose to make Gist’s slicing algorithm interprocedural because failure sketches can span the boundaries of functions. The algorithm is path-insensitive in order to avoid the cost of path-sensitive analyses that do not scale well [ammons:hotpath, bodik:valueflow]. However, this is not a shortcoming, because Gist can recover precise path information at runtime using low-cost control flow tracking (§??). Finally, Gist’s slicing algorithm is flow-sensitive because it traverses statements in a specific order (backwards) from the failure location. Flow-sensitivity generates static slices with statements in the order they appear in the program text (except some out-of-order statements due to path-insensitivity, which are fixed using runtime tracking), thereby helping the developer to understand the flow of statements that lead to a failure.

Algorithm 1 describes Gist’s static backward slicing: it takes as input a failure report (e.g., a coredump, a stack trace) and the program’s source code, and it outputs a static backward slice. For clarity, we define several terms we use in the algorithm. CFG refers to the control flow graph of the program (Gist computes a whole-program CFG as we explain shortly). An item (line 7) is an arbitrary program element. A source (line 8, 16) is an item that is either a global variable, a function argument, a call, or a memory access. Items that are not sources are compiler intrinsics, debug information, and inline assembly. The definitions for a source and an item are specific to LLVM [59], which is what we use for the prototype (§??). The function `getItems` (line 1) returns all the items in a given statement (e.g., the operands of an arithmetic operation). The function `getRetValues` (line 11) performs an intraprocedural analysis to compute and return the set of items that can be returned from a given function call. The function `getArgValues` (line 14) computes and returns the set of arguments that can be used when calling a given function. The function `getReadOperand` (line 20) returns the item that is read, and the function `getWrittenOperand` (line 23) returns the item that is written.

Input : Failure report *report*, program source code *program*
Output : Static backward slice *slice*

```

1 workSet  $\leftarrow$  getItems(failingStmt)
2 function init ()
3   failingStmt  $\leftarrow$  extractFailingStatement(report)
4 function computeBackwardSlice (failingStmt, program)
5   cfg  $\leftarrow$  extractCFG(program)
6   while !workSet.empty() do
7     item  $\leftarrow$  workSet.pop()
8     if isSource(item) then
9       slice.push(item)
10      if isCall(item) then
11        retValues  $\leftarrow$  getRetValues(item, cfg)
12        workSet  $\leftarrow$  workSet  $\cup$  retValues
13      else if isArgument(item) then
14        argValues  $\leftarrow$  getArgValues(item, cfg)
15        workSet  $\leftarrow$  workSet  $\cup$  argValues
16 function isSource (item)
17 if item is (global||argument||call||memory access) then
18   return true
19 else if item is read then
20   workSet  $\leftarrow$  workSet  $\cup$  item.getReadOperand()
21   return true
22 else if item is write then
23   workSet  $\leftarrow$  workSet  $\cup$  item.getWrittenOperand()
24   return true
25 return false

```

Algorithmus 1 : Backward slice computation (Simplified)

Gist's static slicing algorithm differs from classic static slicing [104] in two key ways:

First, Gist addresses a challenge that arises for multithreaded programs because of the implicit control flow edges that get created due to thread creations and joins. For this, Gist uses a compiler pass to build the *thread interprocedural control flow graph* (TICFG) of the program [108]. An *interprocedural control flow graph* (ICFG) of a program connects each function's CFG with function call and return edges. TICFG then augments ICFG to contain edges that represent thread creation and join statements (e.g., a thread creation edge is akin to a callsite with the thread start routine as the target function). TICFG represents an overapproximation of all the possible dynamic control flow behaviors that a program can exhibit at runtime. TICFG is useful for Gist to track control flow that is implicitly created via thread creation and join operations (§??).

Second, unlike other slicing algorithms [llvmslicer], Gist does not use static alias analysis. Alias analysis could determine an overapproximate set of program statements that may affect the computation

of a given value and augment the slice with this information. Gist does not employ static alias analysis because, in practice, it can be over 50% inaccurate [54], which would increase the static slice size that Gist would have to monitor at runtime, thereby increasing its performance overhead. Gist compensates for the lack of alias analysis with runtime data flow tracking, which adds the statements that Gist misses to the static slice (§??).

The static slice that Gist computes has some extraneous items that do not pertain to the failure, because the slicing algorithm lacks actual execution information. Gist weeds out this information using accurate control flow tracking at runtime (§??).

4.3 SLICE REFINEMENT

Slice refinement removes the extraneous statements from the slice and adds to the slice the statements that could not be statically identified. Together with root cause identification (§??), the goal of slice refinement is to build what we call the ideal failure sketch.

We define an *ideal failure sketch* to be one that: 1) contains only statements that have data and/or control dependencies to the statement where the failure occurs; 2) shows the failure predicting events that have the highest positive correlation with the occurrence of failures.

Different developers may have different standards as to what is the “necessary” information for root cause diagnosis; nevertheless, we believe that including all the statements that are related to a failure and identifying the failure predicting events constitute a reasonable and practical set of requirements for root cause diagnosis. Failure predictors are identified by determining the difference of key properties between failing and successful runs.

For example, failure sketches display the partial order of statements involved in data races and atomicity violations, however certain developers may want to know the total order of all the statements in an ideal failure sketch. In our experience, focusing on the partial order of statements that matter from the point of view of root cause diagnosis is more useful than having a total order of all statements. Moreover, obtaining the total order of all the statements in a failure sketch would be difficult without undue runtime performance overhead using today’s technology.

We now describe Gist’s slice refinement strategy in detail. We first describe adaptive tracking of a static slice to reduce the overhead of refinement (§??), then we describe how Gist tracks the control flow (§??) and the data flow (§??) to 1) add to the slice statements that get executed in production but are missing from the slice, and 2) remove from the slice statements that don’t get executed in production.

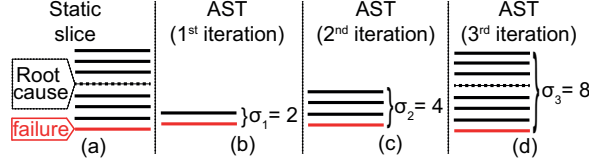


Figure 8 – Adaptive slice tracking in Gist

4.3.1 Adaptive Slice Tracking

Gist employs Adaptive Slice-Tracking (AsT) to track increasingly larger portions of the static slice, until it builds a failure sketch that contains the root cause of the failure that it targets. Gist performs AsT by dynamically tracking control and data flow while the software runs in production. AsT does not track all the control and data elements in the static slice at once in order to avoid introducing performance overhead.

It is challenging to pick the size of the slice, σ , to monitor at runtime, because 1) a too large σ would cause Gist to do excessive runtime tracking and increase overhead; 2) a too small σ may cause Gist to track too many runs before identifying the root cause, and so increase the latency of root cause diagnosis.

Based on previous observations that root causes of most bugs are close to the failure locations [gu:characterization, 84, 117], Gist initially enables runtime tracking for a small number of statements ($\sigma = 2$ in our experiments) backward from the failure point. We use this heuristic because even a simple concurrency bug is likely to be caused by two statements from different threads. This also allows Gist to avoid excessive runtime tracking if the root cause is close to the failure (i.e., the common case). Nonetheless, to reduce the latency of root cause diagnosis, Gist employs a multiplicative increase strategy for further tracking the slice in other production runs. More specifically, Gist doubles σ for subsequent AsT iterations, until a developer decides that the failure sketch contains the root cause and instructs Gist to stop AsT.

Consider the static slice for a hypothetical program in Fig. 8.(a), which displays the failure point (bottom-most solid line) and the root cause (dashed line). In the first iteration (Fig. 8.(b)), AsT starts tracking $\sigma_1 = 2$ statements back from the failure location. Gist cannot build a failure sketch that contains the root cause of this failure by tracking 2 statements, as the root cause lies further backwards in the slice. Therefore, in the second and third iterations (Fig. 8.(c-d)), AsT tracks $\sigma_2 = 4$ and $\sigma_3 = 8$ statements, respectively. Gist can build a failure sketch by tracking 8 statements.

In summary, AsT is a heuristic to resolve the tension between performance overhead, root cause diagnosis latency, and failure sketch accuracy. We elaborate on this tension in our evaluation (§??). AsT

does not limit Gist’s ability to track larger slices and build failure sketches for bugs with greater root-cause-to-failure distances, although it may increase the latency of root cause diagnosis.

4.3.2 *Tracking Control Flow*

Gist tracks control flow to increase the accuracy of failure sketches by identifying which statements from the slice get executed during the monitored production runs. Static slicing lacks real execution information such as dynamically computed call targets, therefore tracking the dynamic control flow is necessary for high accuracy failure sketches.

Static slicing and control flow tracking jointly improve the accuracy of failure sketches: control flow traces identify statements that get executed during production runs that Gist monitors, whereas static slicing identifies statements that have a control or data dependency to the failure. The intersection of these statements represents the statements that relate to the failure and that actually get executed in production runs. Gist statically determines the locations where control flow tracking should start and stop at runtime in order to identify which statements from the slice get executed.

Although control flow can be tracked in a relatively straightforward manner using software instrumentation [67], hardware facilities offer an opportunity for a design with lower overhead. Our design employs Intel PT, a set of new hardware monitoring features for debugging. In particular, Intel PT records the execution flow of a program and outputs a highly-compressed trace (~0.5 bits per retired assembly instruction) that describes the outcome of all branches executed by a program. Intel PT can be programmed to trace only user-level code and can be restricted to certain address spaces. Additionally, with the appropriate software support, Intel PT can be turned on and off by writing to processor-specific registers. Intel PT is currently available in Broadwell processors, and we control it using our custom kernel driver (§??). Future families of Intel processors are also expected to provide Intel PT functionality.

We explain how Gist tracks the statements that get executed via control flow tracking using the example shown in Fig. 9(a). The example shows a static slice composed of three statements (stmt₁, stmt₂, stmt₃). The failure point is stmt₃. Let us assume that, as part of AsT, Gist tracks these three statements. At the high level, Gist identifies all entry points and exit points to each statement and starts and stops control-flow tracking at each entry point and at each exit point, respectively. Tracing is started to capture control flow if the statements in the static slice get executed, and is stopped once those statements complete execution. We use postdominator analysis to optimize out unnecessary tracking.

In this example, Gist starts its analysis with stmt_1 . Gist converts the branch decision information to statement execution information using the technique shown in box I in Fig. 9.(a). It first determines bb_1 , the basic block in which stmt_1 resides, and determines the predecessor basic blocks $p_{11} \dots p_{1n}$ of bb_1 . The predecessor basic blocks of bb_1 are blocks from which control can flow to bb_1 via branches. As a result, Gist starts control flow tracking in each predecessor basic block $p_{11} \dots p_{1n}$ (i.e., entry points). If Gist's control flow tracking determines at runtime that any of the branches from these predecessor blocks to bb_1 was taken, Gist deduces that stmt_1 was executed.

Gist uses an optimization when a statement it already processed strictly dominates the next statement in the static slice. A statement d strictly dominates a statement n (written $d \text{ sdom } n$) if every path from the entry node of the control flow graph to n goes through d , and $d \neq n$. In our example, $\text{stmt}_1 \text{ sdom } \text{stmt}_2$, therefore, Gist will have already started control flow tracking for stmt_1 when the execution reaches stmt_2 , and so it won't need special handling to start control flow tracking for stmt_2 .

However, if a statement that Gist processed does not strictly dominate the next statement in the slice, Gist stops control flow tracking. In our example, after executing stmt_2 , since the execution may never reach stmt_3 , Gist stops control flow tracking after stmt_2 gets executed. Otherwise tracking could continue indefinitely and impose unnecessary overhead. Intuitively, Gist stops control flow tracking *right after* stmt_2 gets executed as shown in box II of Fig. 9.(a). More precisely, Gist stops control flow tracking after stmt_2 and before stmt_2 's immediate postdominator. A node p is said to strictly postdominate a node n if all the paths from n to the exit node of the control flow graph pass through p , and $n \neq p$. The immediate postdominator of a node n ($\text{ipdom}(n)$) is a unique node that strictly postdominates n and does not strictly postdominate any other strict postdominators of n .

Finally, as shown in box III in Fig. 9.(a), Gist processes stmt_3 using the combination of techniques it used for stmt_1 and stmt_2 . Because control flow tracking was stopped after stmt_2 , Gist first restarts it at each predecessor basic block $p_{31} \dots p_{3n}$ of the basic block bb_3 that contains stmt_3 , then Gist stops it after the execution of stmt_3 .

4.3.3 Tracking Data Flow

Similar to control flow, data flow can also be tracked in software, however this can be prohibitively expensive [102]. Existing hardware support can be used for low overhead data flow tracking. In this section, we describe why and how Gist tracks data flow.

Determining the data flow in a program increases the accuracy of failure sketches in two ways:

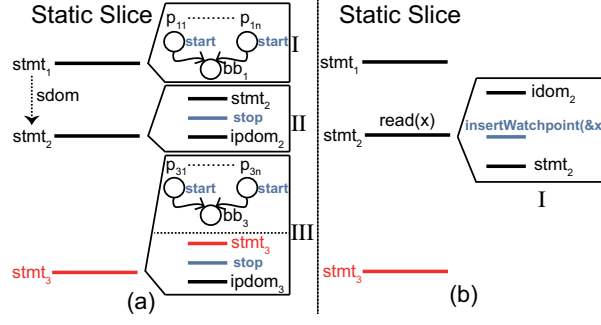


Figure 9 – Example of control (a) and data (b) flow tracking in Gist. Solid horizontal lines are program statements, circles are basic blocks.

First, Gist tracks the total order of memory accesses that it monitors to increase the accuracy of the control flow shown in the failure sketch. Tracking the total order is important mainly for shared memory accesses from different threads, for which Intel PT does not provide order information. Gist uses this order information in failure sketches to help developers reason about concurrency bugs.

Second, while tracking data flow, Gist discovers statements that access the data items in the monitored portion of the slice that were missing from that portion of the slice. Such statements exist because Gist’s static slicing does not use alias analysis (due to alias analysis’ inaccuracy) for determining all statements that can access a given data item.

Gist uses hardware watchpoints present in modern processors to track the data flow (e.g., x86 has 4 hardware watchpoints [40]). They enable tracking the values written to and read from memory locations with low runtime overhead.

For a given memory access, Gist inserts a hardware watchpoint for the address of the accessed variable at a point *right before* the access instruction. More precisely, the inserted hardware watchpoint must be located before the access and after the immediate dominator of that access. Fig. 9.(b) shows an example, where Gist places a hardware watchpoint for the address of variable x , just before stmt_2 ($\text{read}(x)$).

Gist employs several optimizations to economically use its budget of limited hardware watchpoints when tracking the data flow. First, Gist only tracks accesses to shared variables, it does not place a hardware watchpoint for the variables allocated on the stack. Gist maintains a set of active hardware watchpoints to make sure to not place a second hardware watchpoint at an address that it is already watching.

If the statements in the slice portion that AsT monitors access more memory locations than the available hardware watchpoints on a user machine, Gist uses a cooperative approach to track the memory locations across multiple production runs. In a nutshell, Gist’s collaborative approach instructs different production runs to monitor different sets of memory locations in order to monitor all the memory locations

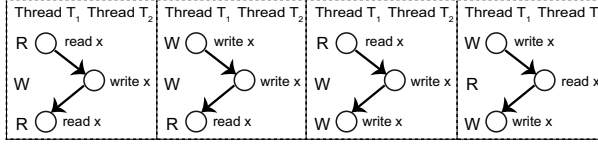


Figure 10 – Four common atomicity violation patterns (RWR, WWR, RWW, WRW). Adapted from [4].

that are in the slice portion that Gist monitors. However, in practice, we did not encounter this situation (§??).

4.4 IDENTIFYING THE ROOT CAUSE

In this section, we describe how Gist determines the differences of key execution properties (i.e., control and data flow) between failing and successful executions in order to do root cause diagnosis.

For root cause diagnosis, Gist follows a similar approach to cooperative bug isolation [4, 46, 63], which uses statistical methods to correlate failure predictors to failures in programs. A failure predictor is a predicate that, when true, predicts that a failure will occur [62]. Carefully crafted failure predictors point to failure root causes [63].

Gist-generated failure sketches contain a set of failure predictors that are both informative and good indicators of failures. A failure predictor is informative if it contains enough information regarding the failure (e.g., thread schedules, critical data values). A failure predictor is a good indicator of a failure if it has high positive correlation with the occurrence of the failure.

Gist defines failure predictors for both sequential and multithreaded programs. For sequential programs, Gist uses branches taken and data values computed as failure predictors. For multithreaded programs, Gist uses the same predicates it uses for sequential programs, as well as special combinations of memory accesses that portend concurrency failures. In particular, Gist considers the common single-variable atomicity violation patterns shown in Fig. 10 (i.e., RWR (Read, Write, Read), WWR, RWW, WRW) and data race patterns (WW, WR, RW) as concurrency failure predictors.

For both failing and successful runs, Gist logs the order of accesses and the value updates to shared variables that are part of the slice it tracks at runtime. Then, using an offline analysis, Gist searches the aforementioned failure-predicting memory access patterns in these access logs. Gist associates each match with either a successful run or a failing run. Gist is not a bug detection tool, but it can understand common failures, such as crashes, assertion violations, and hangs. Other types of failures can either be manually given to Gist, or Gist can be used in conjunction with a bug finding tool.

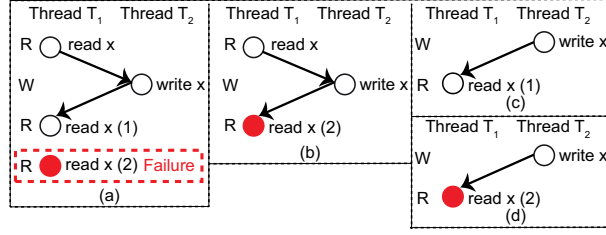


Figure 11 – A sample execution failing at the second read in T₁ (a), and three potential concurrency errors: a RWR atomicity violation (b), 2 WR data races (c-d).

Once Gist has gathered failure predictors from failing and successful runs, it uses a statistical analysis to determine the correlation of these predictors with the failures. Gist computes the precision P (how many runs fail among those that are predicted to fail by the predictor?), and the recall R (how many runs are predicted to fail by the predictor among those that fail?). Gist then ranks all the events by their F-measure, which is the weighted harmonic mean of their precision and recall $F_\beta = (1 + \beta^2) \frac{P \cdot R}{\beta^2 \cdot P + R}$ to determine the best failure predictor. Gist favors precision by setting β to 0.5 (a common strategy in information retrieval [rijsbergen:ir]), because its primary aim is to not confuse the developers with potentially erroneous failure predictors (i.e., false positives).

The failure sketch presents the developer with the highest-ranked failure predictors for each type of failure predictor (i.e., branches, data values, and statement orders). An example of this is shown in Fig. 6, where the dotted rectangles show the highest-ranked failure predictor.

As an example, consider the execution trace shown in Fig. 11.(a). Thread T₁ reads x, after which thread T₂ gets scheduled and writes to x. Then T₁ gets scheduled back and reads x twice in a row, and the program fails (e.g., the second read could be made as part of an assertion that causes the failure). This execution trace has three memory access patterns that can potentially be involved in a concurrency bug: a RWR atomicity violation in Fig. 11.(b) and two data races (or order violations) in Fig. 11.(c) and 11.(d). For this execution, Gist logs these patterns and their outcome (i.e., failure and success: 11.(b) and 11.(d) fail, whereas the pattern in 11.c succeeds. Gist keeps track of the outcome of future access patterns and computes their F-measure to identify the highest ranked failure predictors.

There are two key differences between Gist and cooperative bug isolation (CBI). First, Gist tracks all data values that are part of the slice that it monitors at runtime, allowing it to diagnose the root cause of failures caused by a certain input, as opposed to CBI, which tracks ranges of some variables. Second, Gist uses different failure predictors than CCI [46] and PBI [4], which allow developers to distinguish

between different kinds of concurrency bugs, whereas PBI and CCI use the same predictors for failures with different root causes (e.g., invalid MESI ([80] state for all of RWR, WWR, RWW atomicity violations)).

4.5 GIST'S IMPLEMENTATION DETAILS

Gist's static slicing algorithm is built on the LLVM framework [59]. As part of this algorithm, Gist first augments the intraprocedural control flow graphs of each function with function call and return edges to build the interprocedural control flow graph (ICFG) of the program. Then, Gist processes thread creation and join functions (e.g., `pthread_create`, `pthread_join`) to determine which start routines the thread creation functions may call at runtime and where those routines will join back to their callers, using data structure analysis [58]. Gist augments the edges in the ICFGs of the programs using this information about thread creation/join in order to build the thread interprocedural control flow graph (TICFG) of the program. Gist uses the LLVM information flow tracker [llvm:giri] as a starting point for its slicing algorithm.

Gist currently inserts a small amount of instrumentation into the programs it runs, mainly to start/stop Intel PT tracking and place a hardware watchpoint. To distribute the instrumentation, Gist uses `bsdiff` to create a binary patch file that it ships off to user endpoints or to a data center. We plan to investigate more advanced live update systems such as POLUS [chen:polus]. Another alternative is to use binary rewriting frameworks such as DynamoRio [bruening:dynamorio] or Paradyn [miller:paradyn].

Trace collection is implemented via a Linux kernel module which we refer to as the Intel PT kernel driver. The kernel driver configures and controls the hardware using the documented MSR (Machine Specific Register) interface. The driver allows filtering of what code is traced using the privilege level (i.e. kernel vs. user-space) and CR3 values, thus allowing tracing of individual processes. The driver uses a memory buffer sized at 2 MB, which is sufficient to hold traces for all the applications we have tested. Finally, Gist-instrumented programs use an `ioctl` interface that our driver provides to turn tracing on/off.

Gist's hardware watchpoint use is based on the `ptrace` system call. Once Gist sets the desired hardware watchpoints, it detaches from the program (using the `PTRACE_DETACH`), thereby not incurring any performance overhead. Gist's instrumentation handles hardware watchpoint triggers atomically in order to maintain a total order of accesses among memory operations. Gist logs the program counter when a hardware watchpoint is hit, which it later translates into source line information at developer site. Gist does not need debug

information to do this mapping: it uses the program counter and the offset at which the program binary is loaded to compute where in the actual program this address corresponds to.

PORTEND: CLASSIFYING DATA RACES DURING DEVELOPMENT AND TESTING

Ideally, programs would have no data races at all. In this way, programs would avoid possible catastrophic effects due to data races. This either requires programs to be data race-free by design, or it requires finding and fixing all data races in a program. However, modern software still has data races either because it was written carelessly, or the complexity of the software made it very difficult to properly synchronize threads, or the benefits of fixing all data races using expensive synchronization did not justify its costs.

In order to construct programs that are free of data races by design, novel languages and language extensions that provide a deterministic programming model have been proposed [99, 11]. Deterministic programs are race-free, and therefore, their behavior is not timing dependent. Even though these models may be an appropriate solution for the long term, the majority of modern concurrent software is written in mainstream languages such as C, C++ and Java that don't provide any data race-freedom guarantees.

In order to eliminate all data races in current mainstream software, developers need to first find them. This can be achieved using data race detectors.

Even if all the real races in a program are found using race detectors, eliminating all of them still appears impractical. First, synchronizing all racing memory accesses would introduce performance overheads that may be considered unacceptable. For example, developers have not fixed a race that can lead to lost updates in memcached for a year—ultimately finding an alternate solution—because, it leads to a 7% drop in throughput [72]. Performance implications led to 23 data races in Internet Explorer and Windows Vista being purposely left unfixed [77]. Similarly, several races have been left unfixed in the Windows 7 kernel, because fixing those races did not justify the associated costs [47].

Another reason why data races go unfixed is that 76%–90% of data races are actually considered to be harmless [47, 77, 26, 103]—*harmless races* are assumed to not affect program correctness, either fortuitously or by design, while *harmful races* lead to crashes, hangs, resource leaks, even memory corruption or silent data loss. Deciding whether a race is harmful or not involves a lot of human labor (with industrial practitioners reporting that it can take days, even weeks [36]), so time-pressed developers may not even attempt this high-investment/low-return activity.

Given the large number of data race reports (e.g., Google’s Thread Sanitizer [94] reports over 1,000 unique data races in Firefox when the browser starts up and loads `http://bbc.co.uk`), we argue that data race detectors should also triage reported data races based on the consequences they could have in future executions. This way, developers are better informed and can fix the critical bugs first. A race detector should be capable of inferring the possible consequences of a reported race: is it a false positive, a harmful race, or a race that has no observable harmful effects and left in the code perhaps for performance reasons?

Alas, automated classifiers [45, 47, 77, 100] are often inaccurate (e.g., [77] reports a 74% false positive rate in classifying harmful races). To our knowledge, no data race detector/classifier can do this without false positives.

In this section, we describe Portend, a technique and tool that detects data races and, based on an analysis of the code, infers each race’s potential consequences and automatically classifies them into four categories: “specification violated”, “output differs”, “k-witness harmless” and “single ordering”. In Portend, harmless is circumstantial rather than absolute; it implies that Portend did not witness a harmful effect of the data race in question for k different executions.

For the first two categories, Portend produces a replayable trace that demonstrates the effect, making it easy on the developer to fix the race.

Portend operates on binaries, not on source code (more specifically on LLVM [59] bitcode obtained from a compiler or from a machine-code-to-LLVM translator like RevGen [19]). Therefore it can effectively classify both source-code-level races and assembly-level races that are not forbidden by any language-specific memory model (e.g., C [44] and C++ [43]).

Portend extends our earlier data race classification tool, Portend, with the support for classifying data races under different memory models using a technique called symbolic memory consistency modeling (SMCM). To evaluate this technique, we added specific support for a variant of the weak memory model [24]. We then performed data race classification for microbenchmarks for which the underlying memory model makes a difference with regards to data race classification. We show how SMCM enables more accurate classification under different memory models, and evaluate its runtime and memory overheads. In a nutshell, SMCM allows Portend to achieve 100% classification accuracy for classifying races under different memory models, while incurring low overhead.

5.1 A FINE-GRAINED WAY TO CLASSIFY DATA RACES

A simple harmless vs. harmful classification scheme is undecidable in general (as will be explained below), so prior work typically resorts to “likely harmless” and/or “likely harmful.” Alas, in practice, this is less helpful than it seems (§??). We therefore propose a new scheme that is more precise.

Note that there is a distinction between false positives and harmless races: when a purported race is not a true race, we say it is a *false positive*. When a true race’s consequences are deemed to be harmless for all the witnessed executions, we refer to it as a harmless race¹.

A false positive is harmless in an absolute sense (since it is not a race to begin with), but not the other way around, harmless races are still true races. Static [26], lockset [90], and hybrid [79] data race detectors typically report false positives.

If a data race does not have any observable effect (crash, hang, data corruption) for all the witnessed executions, we say that the data race is harmless with respect to those executions. Harmless races are still true races. Note that our definition of harmlessness is circumstantial; it is not absolute. It is entirely possible that a harmless race for some given executions can become harmful in another execution.

Our proposed scheme classifies the true races into four categories: “spec violated”, “output differs”, “k-witness harmless”, and “single ordering”. We illustrate this taxonomy in Fig. 12.

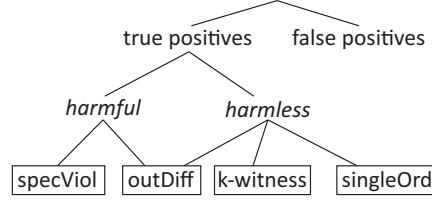


Figure 12 – Portend taxonomy of data races.

“Spec violated” corresponds to races for which at least one ordering of the racing accesses leads to a violation of the program’s specification. These are, by definition, harmful. For example, races that lead to crashes or deadlocks are generally accepted to violate the specification of any program; we refer to these as “basic” specification violations. Higher level program semantics could also be violated, such as the number of objects in a heap exceeding some bound, or a checksum being inconsistent with the checksummed data. Such semantic properties must be provided as explicit predicates to Portend, or be embedded as assert statements in the code.

“Output differs” is the set of races for which the two orderings of the racing accesses can lead to the program generating different out-

1. In the rest of the text, whenever we mention harmless data races, we refer to this definition. If we refer to another definition adopted by prior work, we make the distinction clear.

puts, thus making the output depend on scheduling that is beyond the application’s control. Such races are often considered harmful: one of those outputs is likely “the incorrect” one. However, “output differs” races can also be considered as harmless, whether intentional or not. For example, a debug statement that prints the ordering of the racing memory accesses is intentionally order-dependent, thus an intentional and harmless race. An example of an unintentional harmless race is one in which one ordering of the accesses may result in a duplicated syslog entry—while technically a violation of any reasonable logging specification, a developer may decide that such a benign consequence makes the race not worth fixing, especially if she faces the risk of introducing new bugs or degrading performance when fixing the bug.

As with all high level program semantics, automated Portends *cannot* decide on their own whether an output difference violates some non-explicit specification or not. Moreover, whether the specification has been violated or not might even be subjective, depending on which developer is asked. It is for this reason that we created the “output differs” class of races: we provide developers a clear characterization of the output difference and let them decide using the provided evidence whether that difference matters.

“*K-witness harmless*” are races for which the harmless classification is performed with some quantitative level of confidence: the higher the k , the higher the confidence. Such races are guaranteed to be harmless for at least k combinations of paths and schedules; this guarantee can be as strong as covering a virtually infinite input space (e.g., a developer may be interested in whether the race is harmless for all positive inputs, not caring about what happens for zero or negative inputs). Portend achieves this using a symbolic execution engine [17, 16] to analyze entire equivalence classes of inputs. Depending on the time and resources available, developers can choose k according to their needs—in our experiments we found $k = 5$ to be sufficient to achieve 99% accuracy (manually verified) for all the tested programs. The value of this category will become obvious after reading §??. We also evaluate the individual contributions of exploring paths versus schedules in §6.3.

“*Single ordering*” are races for which only a single ordering of the accesses is possible, typically enforced via ad-hoc synchronization [110]. In such cases, although no explicit synchronization primitives are used, the shared memory could be protected using busy-wait loops that synchronize on a flag. Considering these to be non-races is inconsistent with our definition (§??) because the ordering of the accesses is not enforced using non-ad-hoc synchronization primitives, even though it may not actually be possible to exercise both interleavings of the memory accesses (hence the name of the category). Such ad-hoc synchronization, even if bad practice, is frequent

in real-world software [110]. Previous data race detectors generally cannot tell that only a single order is possible for the memory accesses, and thus report this as an ordinary data race. Such data races can turn out to be both harmful [110] or they can be a major source of harmless data races [45, 100]. That is why we have a dedicated class for such data races.

5.2 PORTEND'S DESIGN OVERVIEW

Portend feeds the target program through its own race detector (or even a third party one, if preferred), analyzes the program and the report automatically, and determines the potential consequences of the reported data race. The report is then classified, based on these predicted consequences, into one of the four categories in Fig. 12. To achieve the classification, Portend performs targeted analysis of multiple schedules of interest, while at the same time using symbolic execution [17] to simultaneously explore multiple paths through the program; we call this technique *multi-path multi-schedule data race analysis*. Portend can thus reason about the consequences of the two orderings of racing memory accesses in a richer execution context than prior work. When comparing program states or program outputs, Portend employs *symbolic output comparison*, meaning it compares constraints on program output as well as path constraints when these outputs are made, in addition to comparing the concrete values of the output, in order to generalize the comparison to more possible inputs that would bring the program to the specific data race and to determine if the data race affects the constraints on program output. Unlike prior work, Portend can accurately classify even races that, given a fixed ordering of the original racing accesses, are harmless along some execution paths, yet harmful along others. In §?? we go over one such race (Fig. 15) and explain how Portend handles it.

Fig. 13 illustrates Portend's architecture. Portend is based on Cloud9 [16], a parallel symbolic execution engine that supports running multi-threaded C/C++ programs. Cloud9 is in turn based on KLEE [17], which is a single-threaded symbolic execution engine. Cloud9 has a number of built-in checkers for memory errors, overflows and division-by-0 errors; on top of which, Portend adds an additional deadlock detector. Portend has a built-in race detector that implements a dynamic happens-before algorithm [57]. This detector relies on a component that tracks Lamport clocks [57] at runtime (details are in §??). Portend's analysis and classification engine performs multi-path multi-schedule data race analysis and symbolic output comparison. This engine also works together with the Lamport clock tracker and the symbolic memory consistency modeling (SMCM) plugin to perform classification. The SMCM plugin defines the rules according to which a memory read operation from a shared memory location can return

previously written values to that location. SMCM plugin is crucial for classifying races under different memory consistency models.

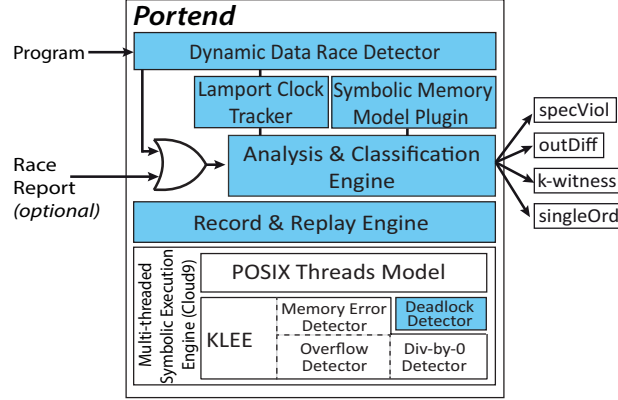


Figure 13 – High-level architecture of Portend. The six shaded boxes indicate new code written for Portend, whereas clear boxes represent reused code from KLEE [17] and Cloud9 [16].

When Portend determines that a race is of “spec violated” kind, it provides the corresponding evidence in the form of program inputs (including system call return values) and thread schedule that reproduce the harmful consequences deterministically. Developers can replay this “evidence” in a debugger, to fix the race.

We now give an overview of our approach and illustrate it with an example, describe the first step, single-path/single-schedule analysis (§??), followed by the second step, multi-path analysis and symbolic output comparison (§??) augmented with multi-schedule analysis (§??). We introduce SMCM and describe how it can be used to model various memory models while performing data race classification (§??). We describe Portend’s race classification (§??) and the generated report that helps developers debug the race (§??).

Portend’s race analysis starts by executing the target program and dynamically detecting data races (e.g., developers could run their existing test suites under Portend). Portend detects races using a dynamic happens-before algorithm [57]. Alternatively, if a third party detector is used, Portend can start from an existing execution trace; this trace must contain the thread schedule and an indication of where in the trace the suspected race occurred. We developed a plugin for Thread Sanitizer [94] to create a Portend-compatible trace; we believe such plugins can be easily developed for other dynamic race detectors [38].

Portend has a record/replay infrastructure for orchestrating the execution of a multi-threaded program; it can preempt and schedule threads before/after synchronization operations and/or racing accesses. Portend uses Cloud9 to enumerate program paths and to collect symbolic constraints.

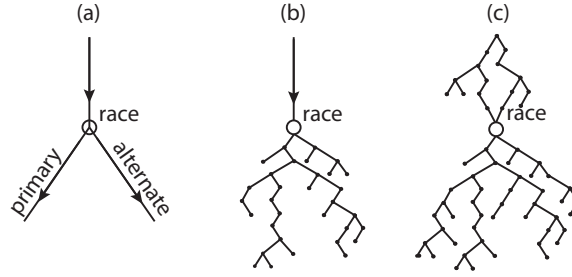


Figure 14 – Increasing levels of completeness in terms of paths and schedules: [a. single-pre/single-post] \ll [b. single-pre/multi-post] \ll [c. multi-pre/multi-post].

A trace consists of a schedule trace and a log of system call inputs. The schedule trace contains the thread id and the program counter at each preemption point. Portend treats all POSIX threads synchronization primitives as possible preemption points and uses a single-processor cooperative thread scheduler (see §?? for a discussion of the resulting advantages and limitations). Portend can also preempt threads before and after any racing memory access. We use the following notation for the trace: $(T_0 : pc_0) \rightarrow (T_1 \rightarrow RaceyAccess_{T_1} : pc_1) \rightarrow (T_2 \rightarrow RaceyAccess_{T_2} : pc_2)$ means that thread T_0 is preempted after it performs a synchronization call at program counter pc_0 ; then thread T_1 is scheduled and performs a memory access at program counter pc_1 , after which thread T_2 is scheduled and performs a memory access at pc_2 that is racing with the previous memory access of T_1 . The schedule trace also contains the absolute count of instructions executed by the program up to each preemption point. This is needed in order to perform precise replays when an instruction executes multiple times (e.g., a loop) before being involved in a race; this is not shown as part of the schedule trace, for brevity. The log of system call inputs contains the non-deterministic program inputs (e.g., *gettimeofday*).

In a first analysis step (illustrated in Fig. 14a), Portend replays the schedule in the trace up to the point where the race occurs. Then it explores two different executions: one in which the original schedule is followed (the *primary*) and one in which the alternate ordering of the racing accesses is enforced (the *alternate*). As described in §??, some classifiers compare the primary and alternate program state immediately after the race, and, if different, flag the race as potentially harmful, and, if same, flag the race as potentially harmless. Even if program outputs are compared rather than states, “single-pre/single-post” analysis (Fig. 14a) may not be accurate, as we will show below. Portend uses “single-pre/single-post” analysis mainly to determine whether the alternate schedule is possible at all. In other words, this stage identifies any ad-hoc synchronization that might prevent the alternate schedule from occurring.

If there is a difference between the primary and alternate post-race states, we do not consider the race as necessarily harmful. Instead, we allow the primary and alternate executions to run, independently of each other, and we observe the consequences. If, for instance, the alternate execution crashes, the race is harmful. Of course, even if the primary and alternate executions behave identically, it is still not certain that the race is harmless: there may be some unexplored pair of primary and alternate paths with the same pre-race prefix as the analyzed pair, but which does not behave the same. This is why single-pre/single-post analysis is insufficient, and we need to explore *multiple* post-race paths. This motivates “single-pre/multi-post” analysis (Fig. 14b), in which multiple post-race execution possibilities are explored—if any primary/alternate mismatch is found, the developer must be notified.

Even if all feasible post-race paths are explored exhaustively and no mismatch is found, one still cannot conclude that the race is harmless: it is possible that the absence of a mismatch is an artifact of the specific pre-race execution prefix, and that some different prefix would lead to a mismatch. Therefore, to achieve higher confidence in the classification, Portend explores multiple feasible paths even in the pre-race stage, not just the one path witnessed by the race detector. This is illustrated as “multi-pre/multi-post” analysis in Fig. 14c. The advantage of doing this vs. considering these as different races is the ability to systematically explore these paths.

Finally, we combine multi-path analysis with *multi-schedule* analysis, since the same path through a program may generate different outputs depending on how its execution segments from different threads are interleaved. The branches of the execution tree in the post-race execution in Fig. 14c correspond to different paths that stem from both multiple inputs and schedules, as we detail in §??.

Of course, exploring all possible paths and schedules that experience the race is impractical, because their number typically grows exponentially with the number of threads, branches, and preemption points in the program. Instead, we provide developers a “dial” to control the number k of path/schedule alternatives explored during analysis, allowing them to control the “volume” of paths and schedules in Fig. 14. If Portend classifies a race as “ k -witness harmless”, then a higher value of k offers higher confidence that the race is harmless for *all* executions (i.e., including the unexplored ones), but it entails longer analysis time. We found $k = 5$ to be sufficient for achieving 99% accuracy in our experiments in less than 5 minutes per race on average.

To illustrate the benefit of multi-path multi-schedule analysis over “single-pre/single-post” analysis, consider the code snippet in Fig. 15, adapted from a real data race bug. This code has racing accesses to the global variable *id*. Thread T_0 spawns threads T_1 and T_2 ; thread T_1

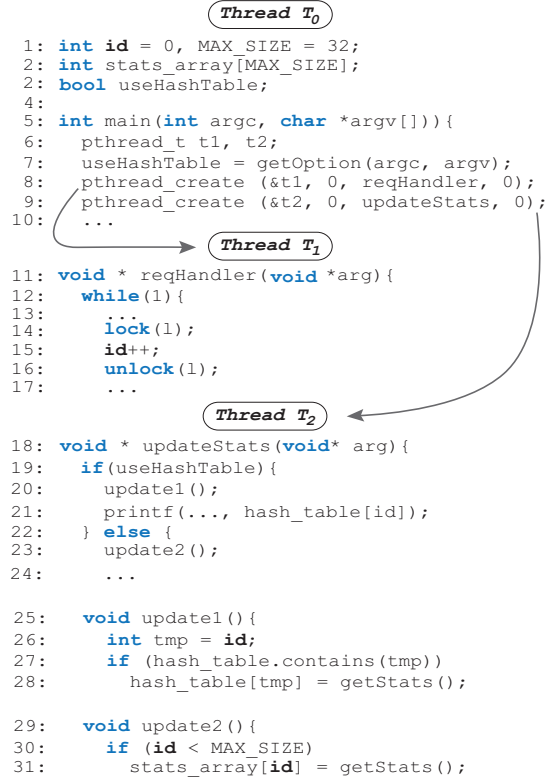


Figure 15 – Simplified example of a harmful race from Ctrace [71] that would be classified as harmless by classic race classifiers.

updates *id* (line 15) in a loop and acquires a lock each time. However, thread T_2 , which maintains statistics, reads *id* without acquiring the lock—this is because acquiring a lock at this location would hurt performance, and statistics need not be precise. Depending on program input, T_2 can update the statistics using either the *update1* or *update2* functions (lines 20-23).

Say the program runs in Portend with input *-use-hash-table*, which makes *useHashTable*=*true*. Portend records the primary trace ($T_0 : pc_9 \rightarrow \dots (T_1 \rightarrow RaceyAccess_{T_1} : pc_{15}) \rightarrow (T_2 \rightarrow RaceyAccess_{T_2} : pc_{26}) \rightarrow \dots T_0$). This trace is fed to the first analysis step, which replays the trace with the same program input, except it enforces the alternate schedule ($T_0 : pc_9 \rightarrow \dots (T_2 \rightarrow RaceyAccess_{T_2} : pc_{26}) \rightarrow (T_1 \rightarrow RaceyAccess_{T_1} : pc_{15}) \rightarrow \dots T_0$). Since the printed value of *hash_table[id]* at line 21 would be the same for the primary and alternate schedules, a “single-pre/single-post” classifier would deem the race harmless.

However, in the multi-path multi-schedule step, Portend explores additional paths through the code by marking program input as symbolic, i.e., allowing it to take on any permitted value. When the trace is replayed and Portend reaches line 19 in T_2 in the alternate schedule, *useHashTable* could be both true and false, so Portend splits the execution into two executions, one in which *useHashTable* is set to *true* and one in which it is *false*. Assume, for example, that *id* = 31 when check-

```

Input : Primary execution trace primary
Output : Classification result  $\in \{specViol, outDiff, outSame, singleOrd\}$ 
1 current  $\leftarrow$  execUntilFirstThreadRacyAccess(primary)
2 preRaceCkpt  $\leftarrow$  checkpoint(current)
3 execUntilSecondThreadRacyAccess(current)
4 postRaceCkpt  $\leftarrow$  checkpoint(current)
5 current  $\leftarrow$  preRaceCkpt
6 preemptCurrentThread(current)
7 alternate  $\leftarrow$  execWithTimeout(current)
8 if alternate.timedOut then
9   if detectInfiniteLoop(alternate) then
10    return specViol
11   else
12    return singleOrd
13 else
14   if detectDeadlock(alternate) then
15    return specViol
16 primary  $\leftarrow$  exec(postRaceCkpt)
17 if detectSpecViol(primary)  $\vee$  detectSpecViol(alternate) then
18   return specViol
19 if primary.output  $\neq$  alternate.output then
20   return outDiff
21 else
22   return outSame

```

Algorithmus 2 : Single-Pre/Single-Post Analysis (singleClassify)

ing the if condition at line 30. Due to the data race, *id* is incremented by T_1 to 32, which overflows the statically allocated buffer (line 31). Note that in this alternate path, there are two racing accesses on *id*, and we are referring to the access at line 31.

Portend detects the overflow (via Cloud9), which leads to a crashed execution, flags the race as “spec violated”, and provides the developer the execution trace in which the input is *-no-hash-table*, and the schedule is $(T_0 : pc_9) \rightarrow \dots (T_2 \rightarrow RaceyAccess_{T_2} : pc_{30}) \rightarrow (T_1 \rightarrow RaceyAccess_{T_1} : pc_{15}) \rightarrow (T_2 : pc_{31})$. The developer can replay this trace in a debugger and fix the race.

Note that this data race is harmful only if the program input is *-no-hash-table*, the given thread schedule occurs, and the value of *id* is 31; therefore the crash is likely to be missed by a traditional single-path/single-schedule data race detector.

We now describe Portend’s race analysis in detail: §??–§?? focus on the *exploration* part of the analysis, in which Portend looks for paths and schedules that reveal the nature of the race, and §?? focuses on the *classification* part.

5.3 SINGLE-PATH ANALYSIS

The goal of this first analysis step is to identify cases in which the alternate schedule of a race cannot be pursued, and to make a first classification attempt based on a single alternate execution. Algorithm 2 describes the approach.

Portend starts from a trace of an execution of the target program, containing one or more races, along with the program inputs that generated the trace. For example, in the case of the Pbzip2 file compressor used in our evaluation, Portend needs a file to compress and a trace of the thread schedule.

As mentioned earlier, such traces are obtained from running, for instance, the developers' test suites (as done in CHES [74]) with a dynamic data race detector enabled.

Portend takes the primary trace and plays it back (line 1). Note that *current* represents the system state of the current execution. Just before the first racing access, Portend takes a checkpoint of system state; we call this the *pre-race checkpoint* (line 2). The replay is then allowed to continue until immediately after the second racing access of the race we are interested in (line 3), and the primary execution is suspended in this post-race state (line 4).

Portend then primes a new execution with the pre-race checkpoint (line 5) and attempts to enforce the alternate ordering of the racing accesses. To enforce this alternate order, Portend preempts the thread that did the first racing access (T_i) in the primary execution and allows the other thread (T_j) involved in the race to be scheduled (line 6). In other words, an execution with the trace $\dots(T_i \rightarrow \text{RaceAccess}_{T_i} : pc_1) \rightarrow (T_j \rightarrow \text{RaceAccess}_{T_j} : pc_2)\dots$ is steered toward the execution $\dots(T_j \rightarrow \text{RaceAccess}_{T_j} : pc_2) \rightarrow (T_i \rightarrow \text{RaceAccess}_{T_i} : pc_1)\dots$

This attempt could fail for one of three reasons: (a) T_j gets scheduled, but T_i cannot be scheduled again; or (b) T_j gets scheduled but RaceAccess_{T_j} cannot be reached because of a complex locking scheme that requires a more sophisticated algorithm [49] than Algorithm 2 to perform careful scheduling of threads; or (c) T_j cannot be scheduled, because it is blocked by T_i . Case (a) is detected by Portend via a timeout (line 8) and is classified either as “spec violated”, corresponding to an infinite loop (i.e., a loop with a loop-invariant exit condition) in line 10 or as ad-hoc synchronization in line 12. Portend does not implement the more complex algorithm mentioned in (b), and this may cause it to have false positives in data race classification. However, we have not seen this limitation impact the accuracy of data race classification for the programs in our evaluation. Case (c) can correspond to a deadlock (line 15) and is detected by Portend by keeping track of the lock graph. Both the infinite loop and the deadlock case cause the race to be classified as “spec violated”, while the ad-hoc synchronization case classifies the race as “single ordering” (more details in

§??). While it may make sense to not stop if the alternate execution cannot be enforced, under the expectation that other paths with other inputs might permit the alternate ordering, our evaluation suggests that continuing adds little value.

If the alternate schedule succeeds, Portend executes it until it completes, and then records its outputs. Then, Portend allows the primary to continue (while replaying the input trace) and also records its outputs. During this process, Portend watches for “basic” specification violations (crashes, deadlocks, memory errors, etc.) as well as “high level” properties given to Portend as predicates—if any of these properties are violated, Portend immediately classifies (line 18) the race as “spec violated”. If the alternate execution completes with no specification violation, Portend compares the outputs of the primary and the alternate; if they differ, the race is classified as “output differs” (line 20), otherwise the analysis moves to the next step. This is in contrast to replay-based classification [77], which compares the program state immediately after the race in the primary and alternate interleavings.

5.4 MULTI-PATH ANALYSIS

The goal of this step is to explore variations of the single paths found in the previous step (i.e., the primary and the alternate) in order to expose Portend to a wider range of execution alternatives.

First, Portend finds multiple primary paths that satisfy the input trace, i.e., they (a) all experience the same thread schedule (up to the data race) as the input trace, and (b) all experience the target race condition. These paths correspond to *different* inputs from the ones in the initial race report. Second, Portend uses Cloud9 to record the “symbolic” outputs of these paths—that is, the constraints on the output, rather than the concrete output values themselves—as well as path constraints when these outputs are made, and compares them to the outputs and path constraints of the corresponding alternate paths; we explain this below. Algorithm 3 describes the functions invoked by Portend during this analysis in the following order: 1) on initialization, 2) when encountering a thread preemption, 3) on a branch that depends on symbolic data, and 4) on finishing an execution.

Unlike in the single-pre/single-post step, Portend now executes the primary *symbolically*. This means that the target program is given symbolic inputs instead of regular concrete inputs. Cloud9 relies in large part on KLEE [17] to interpret the program and propagate these symbolic values to other variables, corresponding to how they are read and operated upon. When an expression with symbolic content is involved in the condition of a branch, *both* options of the branch are explored, if they are feasible. The resulting path(s) are annotated with a constraint indicating that the branch condition holds true (re-

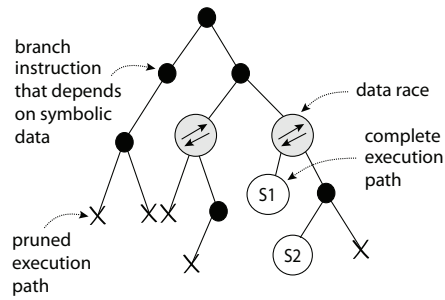


Figure 16 – Portend prunes paths during symbolic execution.

spectively false). Thus, instead of a regular single-path execution, we get a tree of execution paths, similar to the one in Fig. 16. Conceptually, at each such branch, program state is duplicated and constraints on the symbolic parameters are updated to reflect the decision taken at that branch (line 11). Describing the various techniques for performing symbolic execution efficiently [17, 16] is beyond the scope of this article.

An important concern in symbolic execution is “path explosion,” i.e., that the number of possible paths is large. Portend offers two parameters to control this growth: (a) an upper bound M_p on the number of primary paths explored; and (b) the number and size of symbolic inputs. These two parameters allow developers to trade performance vs. classification confidence. For parameter (b), the fewer inputs are symbolic, the fewer branches will depend on symbolic input, so less branching will occur in the execution tree.

Determining the optimal values for these parameters may require knowledge of the target system as well as a good sense of how much confidence is required by the system’s users. Reasonable (i.e., good but not necessarily optimal) values can be found through trial and error relatively easily—we expect development teams using Portend to converge onto values that are a good fit for their code and user community, and then make these values the defaults for their testing and triage processes. We empirically study in §?? the impact of these parameters on classification accuracy on a diverse set of programs and find that relatively small values achieve high accuracy for a broad range of programs.

During symbolic execution, Portend prunes (Fig. 16) the paths that do not obey the thread schedule in the trace (line 8), thus excluding the (many) paths that do not enable the target race. Moreover, Portend attempts to follow the original trace only until the second racing access is encountered; afterward, it allows execution to diverge from the original schedule trace. This enables Portend to find more executions that partially match the original schedule trace (e.g., cases in which the second racing access occurs at a different program counter, as in Fig. 15). Tolerating these divergences significantly increases Por-

Input : Schedule trace $trace$, initial program state S_o , set of states $S = \emptyset$, upper bound M_p on the number of primary paths

Output : Classification result $\in \{specViol, outDiff, singleOrd\}$ k -witness

```

1 function init ()
2    $S \leftarrow S \cup S_o$ 
3    $current \leftarrow S.head()$ 
4    $pathsExplored \leftarrow 0$ 
5 function onPreemption ()
6    $t_i \leftarrow \text{scheduleNextThread}(current)$ 
7   if  $t_i \neq \text{nextThreadInTrace}(trace, current)$  then
8      $S \leftarrow S.remove(current)$ 
9      $current \leftarrow S.head()$ 
10 function onSymbolicBranch ()
11    $S \leftarrow S \cup current.fork()$ 
12 function onFinish ()
13    $classification \leftarrow classification \cup \text{classify}(current)$ 
14   if  $pathsExplored < M_p$  then
15      $pathsExplored \leftarrow pathsExplored + 1$ 
16 else
17   return  $classification$ 
18 function classify (primary)
19    $result \leftarrow \text{singleClassify}(primary)$ 
20   if  $result = outSame$  then
21      $alternate \leftarrow \text{getAlternate}(primary)$ 
22     if  $\text{symbolicMatch}(primary.symState, alternate.symState)$  then
23       return  $k$ -witness
24     else
25       return  $outDiff$ 
26 else
27   return  $result$ 

```

Algorithmus 3 : Multi-path Data Race Analysis (Simplified)

tend's accuracy over the state of the art [77], as will be explained in §??.

Once the desired paths are obtained (at most M_p , line 14), the conjunction of branch constraints accumulated along each path is solved by KLEE using an SMT solver [32] in order to find concrete inputs that drive the program down the corresponding path. For example, in the case of Fig. 16, two successful leaf states S_1 and S_2 are reached, and the solver provides the inputs corresponding to the path from the root of the tree to S_1 , respectively S_2 . Thus, we now have $M_p = 2$ different primary executions that experience the data race.

5.5 SYMBOLIC OUTPUT COMPARISON

Portend now records the output of each of the M_p executions, like in the single-pre/single-post case, and it also records the path con-

straints when these outputs are made. However, this time, in addition to simply recording concrete outputs, Portend propagates the constraints on symbolic state all the way to the outputs, i.e., the outputs of each primary execution contain a mix of concrete values and symbolic constraints (i.e., symbolic formulae). Note that by output we mean all arguments passed to output system calls.

Next, for each of the M_p executions, Portend produces a corresponding alternate (analogously to the single-pre/single-post case). The alternate executions are fully concrete, but Portend records constraints on the alternate's outputs (lines 19-21) as well as the path constraints when these outputs are made. The function *singleClassify* in Algorithm 3 performs the analysis described in Algorithm 2. Portend then checks whether the constraints on outputs of each alternate and the path constraints when these outputs are made match the constraints of the corresponding primary's outputs and the path constraints when primary's outputs are made. This is what we refer to as *symbolic output comparison* (line 22). The purpose behind comparing symbolic outputs is that Portend tries to figure out if the data race caused the constraints on the output to be modified or not, and the purpose behind comparing path constraints is to be able generalize the output comparison to more possible executions with different inputs.

This symbolic comparison enables Portend's analysis to extend over more possible primary executions. To see why this is the case, consider the example in Figure 17. In this example, the Main thread reads input to the shared variable i and then spawns two threads T_1 and T_2 which perform racing writes to $globalx$. T_1 prints 10 if the input is positive. Let us assume that during the symbolic execution of the primary, the write to $globalx$ in T_1 is performed before the write to $globalx$ in T_2 . Portend records that the output at line 6 is 10 if the path constraint is $i \geq 0$. Let us further assume that Portend runs the program while enforcing the alternate schedule with input 1. The output of the program will still be 10 (since $i \geq 0$) and the path constraint when the output will be made will still be $i \geq 0$. The output and the path constraint when the output is made is therefore the same regardless of the order with which the accesses to $globalx$ are performed (i.e., the primary or the alternate order). Therefore, Portend can assert that the program output for $i \geq 0$ will be 10 regardless of the way the data race goes even though it only explored a single alternate ordering with input 1.

This comes at the price of potential false negatives because path constraints can be modified due to a different thread schedule; despite this theoretical shortcoming, we have not encountered such a case in practice, but we plan to investigate further in future work.

False negatives can also arise because determining semantic equivalence of output is undecidable, and our comparison may still wrongly

classify as “output differs” a sequence of outputs that are equivalent at some level (e.g., $\langle \text{print } ab; \text{print } c \rangle$ vs. $\langle \text{print } abc \rangle$).

```

1:  int globalx = 0;
2:  int i = 0;

Thread T1

3:  void * work0 (void *arg) {
4:      globalx = 1;
5:      if(i >= 0)
6:          printf("10\n");
7:      return 0;
8:  }

Thread T2

9:  void * work1 (void *arg) {
10:     globalx = 2;
11:     return 0;
12:  }

Main Thread

13: int main (int argc, char *argv[]){
14:     pthread_t t0, t1;
15:     int rc;
16:     i = getInput(argc, argv)
17:     rc = pthread_create (&t0, 0, work0, 0);
18:     rc = pthread_create (&t1, 0, work1, 0);
19:     pthread_join(t0, 0);
20:     pthread_join(t1, 0);
21:     return 0;
22: }
```

Figure 17 – A program to illustrate the benefits of symbolic output comparison

When executing the primaries and recording their outputs and the path constraints, Portend relies on Cloud9 to track all symbolic constraints on variables. To determine if the path constraints and constraints on outputs match for the primary and the alternates, Portend directly employs an SMT solver [32].

As will be seen in §6.3, using symbolic comparison in conjunction with multi-path multi-schedule analysis leads to substantial improvements in classification accuracy.

We do not detail here the case when the program reads input *after* the race—it is a natural extension of the algorithm above.

5.6 MULTI-SCHEDULE ANALYSIS

The goal of multi-schedule analysis is to further augment the set of analyzed executions by diversifying the thread schedule.

We mentioned earlier that, for each of the M_p primary executions, Portend obtains an alternate execution. Once the alternate ordering of the racing accesses is enforced, Portend randomizes the schedule of the *post-race* alternate execution: at every preemption point in the alternate, Portend randomly decides which of the runnable threads to schedule next. This means that every alternate execution will most likely have a different schedule from the original input trace (and thus from the primary).

Consequently, for every primary execution P_i , we obtain *multiple* alternate executions A_i^1, A_i^2, \dots by running up to M_a multiple instances of the alternate execution. Since the scheduler is random, we expect practically every alternate execution to have a schedule that differs from all others. Recently proposed techniques [73] can be used to quantify the probability of these alternate schedules discovering the harmful effects of a data race.

Portend then uses the same symbolic comparison technique as in §?? to establish equivalence between the constraint on outputs and path constraints of $A_i^1, A_i^2, \dots, A_i^{M_a}$ and the symbolic outputs and path constraints of P_i .

Schedule randomization can be employed also in the pre-race stage of the alternate-execution generation as well as in the generation of the primary executions. We did not implement these options, because the level of multiplicity we obtain with the current design appears to be sufficient in practice to achieve high accuracy. Note however that, as we show in §6.3, multi-path multi-schedule analysis is indeed crucial to attaining high classification accuracy.

In summary, multi-path multi-schedule analysis explores M_p primary executions and, for each such execution, M_a alternate executions with different schedules, for a total of $M_p \times M_a$ path-schedule combinations. For races that end up being classified as “k-witness harmless”, we say that $k = M_p \times M_a$ is the lower bound on the number of concrete path-schedule combinations under which this race is harmless.

Note that the k executions can be simultaneously explored in parallel: if a developer has p machines with q cores each, she could explore $p \times q$ parallel executions in the same amount of time as a single execution. Given that Portend is “embarrassingly parallel,” it is appealing for cluster-based automated bug triage systems.

5.7 PORTEND'S IMPLEMENTATION DETAILS

Portend works on programs compiled to LLVM [59] bitcode and can run C/C++ programs for which there exists a sufficiently complete symbolic POSIX environment [16]. We have tested it on C programs as well as C++ programs that do not link to libstdc++; this latter limitation results from the fact that an implementation of a standard C++ library for LLVM is in progress, but not yet available [22]. Portend uses Cloud9 [16] to interpret and symbolically execute LLVM bitcode; we suspect any path exploration tool will do (e.g., CUTE [93], SAGE [35]), as long as it supports multi-threaded programs.

T Portend intercepts various system calls, such as *write*, under the assumption that they are the primary means by which a program communicates changes in its state to the environment. A separate Portend module is responsible for keeping track of symbolic out-

puts in the form of constraints, as well as of concrete outputs. Portend hashes program outputs (when they are concrete) and can either maintain hashes of all concrete outputs or compute a hash chain of all outputs to derive a single hash code per execution. This way, Portend can deal with programs that have a large amount of output.

Portend uses a dynamic race detector that is based on the happens before relationship [57]. The happens before relationships are tracked using the Lamport clock tracker component. This component, together with the SMC plugin is used by the analysis and classification engine to perform classification under a particular memory model.

Portend keeps track of Lamport clocks per execution state on the fly. Note that it is essential to maintain the happens before graph per execution state because threads may get scheduled differently depending on the flow of execution in each state and therefore synchronization operations may end up being performed in a different order.

The state space exploration in Portend is exponential in the number of values that “read”s can return in a program. Therefore the implementation needs to handle bookkeeping as efficiently as possible. There are several optimizations that are in place to enable a more scalable exploration. The most important ones are: 1) Copy-on-write for keeping the happens before graph and 2) Write buffer compression.

Portend employs copy-on-write for tracking the happens before graphs in various states. Initially, there is a single happens before graph that gets constructed during program execution before any state is forked due to a read with multiple possible return values. Then, when a state is forked, the happens before graph is not duplicated. The forking state rather maintains a reference to the old graph. Then, when a new synchronization operation is recorded in either one of the forked states, this event is recorded as an incremental update to previously saved happens before graph. In this way maximum sharing of the happens before graph is achieved among forked states.

The copy-on-write scheme for states can be further improved if one checks whether two different states perform the same updates to the happens before graph. If that is the case, these updates can be merged and saved as part of the common happens before graph. This feature is not implemented in the current prototype, but it is a potential future optimization.

The second optimization is write buffer compression. This is performed whenever the same value is written to the same shared variable’s buffer and the constraints imposed by the happens before relationship allow these same values to be returned upon a read. Then, in such a case, these two writes are compressed into one, as returning two of them would be redundant from the point of view of state

exploration. For example, if a thread writes 1 to a shared variable `globalx` twice before this value is read by another thread, the write buffer will be compressed to behave as if the initial thread has written 1 once.

Portend clusters the data races it detects, in order to filter out similar races; the clustering criterion is whether the racing accesses are made to the same shared memory location by the same threads, and the stack traces of the accesses are the same. Portend provides developers with a single representative data race from each cluster.

T The timeout used in discovering ad-hoc synchronization is conservatively defined as 5 times what it took Portend to replay the primary execution, assuming that reversing the access sequence of the racing accesses should not cause the program to run for longer than that.

In order to run multi-threaded programs in Portend, we extended the POSIX threads support found in Cloud9 to cover almost the entire POSIX threads API, including barriers, mutexes and condition variables, as well as thread-local storage. Portend intercepts calls into the POSIX threads library to maintain the necessary internal data structures (e.g., to detect data races and deadlocks) and to control thread scheduling.

EVALUATION

In this section, we evaluate all the prototypes we built for all the techniques we presented in the previous four chapters. For each prototype evaluation, we first describe the experimental setup followed with a description of prototype-specific experiments. We first present the evaluation of RaceMob (§6.1). Next we present the evaluation results of Gist, followed by the evaluation results of Portend.

6.1 RACEMOB'S EVALUATION

Program	Apache	SQLite	Memcached	Fmm	Barnes	Ocean	Pbzip2	Knot	Aget	Pfscan
Size (LOC)	138,456	113,326	19,397	9,126	7,580	6,551	3,521	3,586	2,053	2,033
Race candidates	118	88	7	176	166	115	65	65	24	17
True Race	Causes hang	0	3	0	0	0	0	0	0	0
	Causes crash	0	0	0	0	0	3	0	0	0
	Both orders	0	0	1	5	10	0	2	0	0
	Single order	8	0	0	53	6	3	4	2	4
Likely FP	Not aliasing	10	31	0	33	65	13	0	18	2
	Context	61	10	2	61	28	42	21	28	10
	Synchronization	1	37	3	10	49	47	34	13	7
Unknown	38	7	1	14	8	10	1	4	1	0

Table 1 – Data race detection with RaceMob. The static phase reports *Race candidates* (row 2). The dynamic phase reports verdicts (rows 3-10). *Causes hang* and *Causes crash* are races that caused the program to hang or crash. *Single order* are true races for which either the primary or the alternate executed (but not both) with no intervening synchronization; *Both orders* are races for which both executed without intervening synchronization.

In this section, we address the following questions about RaceMob: Can it effectively detect true races in real code (§??)? Is it efficient (§??)? How does it compare to state-of-the-art data race detectors (§??) and interleaving-based concurrency testing tools (§??)? Finally, how does RaceMob scale with the number of threads (§??)?

6.1.1 Experimental Setup

We evaluated RaceMob using a mix of server, desktop and scientific software: Apache httpd is a Web server that serves around

65% of the Web [2]—we used the `mpm-worker` module of Apache to operate it in multi-threaded server mode and detected races in this specific module. SQLite [97] is an embedded database used in Firefox, iOS, Chrome, and Android, and has 100% branch coverage with developer’s tests. Memcached [28] is a distributed memory-object caching system, used by Internet services like Twitter, Flickr, and YouTube. Knot [8] is a web server. Pbzip2 [33] is a parallel implementation of the popular bzip2 file compressor. Pfscan [27] is a parallel file scanning tool that provides the combined functionality of `find`, `xargs`, and `fgrep` in a parallel way. Aget is a parallel variant of `wget`. Fmm, Ocean, and Barnes are applications from the SPLASH suite [95]. Fmm and Barnes simulate interactions of bodies, and Ocean simulates ocean movements.

Our evaluation results are obtained primarily using a test environment simulating a crowdsourced setting, and we also have a small scale, real deployment of RaceMob on our laptops. For the experiments, we use a mix of workloads derived from actual program runs, test suites, and test cases devised by us and other researchers [113]. We configured the hive to assign a single dynamic validation task per user at a time. Altogether, we have execution information from 1,754 simulated user sites. Our test bed consists of a 2.3 GHz 48-core AMD Opteron 6176 machine with 512 GB of RAM running Ubuntu Linux 11.04 and a 2 GHz 8-core Intel Xeon E5405 machine with 20 GB of RAM running Ubuntu Linux 11.10. The hive is deployed on the 8-core machine, and the simulated users on both machines. The real deployment uses ThinkPad laptops with Intel 2620M processors and 8 GB of RAM, running Ubuntu Linux 12.04.

We used C programs in our evaluation because RELAY operates on CIL, which does not support C++ code. Pbzip2 is a C++ program, but we converted it to C by replacing references to STL `vector` with an array-based implementation. We also replaced calls to `new/delete` with `malloc/free`.

6.1.2 Effectiveness

To investigate whether RaceMob is an effective way to detect data races, we look at whether RaceMob can detect true data races, and whether its false positive and false negative rates are sufficiently low.

RaceMob’s data race detection results are centralized in Table 1. RaceMob detected a total of 106 data races in ten programs. Four races in Pbzip2 caused the program to crash, three races in SQLite caused the program to hang, and one race in Aget caused a data corruption (that we confirmed manually). The other races did not lead to any observable failure. We manually confirmed that the “True Race” verdicts are correct, and that RaceMob has no false positives in our experiments.

The “Likely FP” row represents the races that RaceMob identified as likely false positives: (1) *Not aliasing* are reports with accesses that do not alias to the same memory location at runtime; (2) *Context* are reports whose accesses are only made by a single thread at runtime; (3) *Synchronization* are reports for which, the accesses are synchronized, an artifact that the static detector missed. The first two sources of likely false positives (53% of all static reports) are identified using DCI, whereas the last source (24% of all static reports) is identified using on demand race detection. In total, 77% of all statically detected races are likely false positives.

As we discussed in §??, RaceMob’s false negative rate is determined by its static data race detector. We manually verified that none of RELAY’s sources of false negatives (i.e., inline assembly and pointer arithmetic) are present in the programs in our evaluation. Furthermore, Chimera [60], a deterministic record/replay system, relies on RELAY; for deterministic record/replay to work, all data races must be detected; in Chimera’s evaluation (which included Apache, Pbzp2, Knot, Ocean, Pfscan, Aget), RELAY did not have any false negatives [60]. We therefore cautiously conclude that RaceMob’s static phase had no false negatives in our evaluation. However, this does not exclude the possibility that for other programs there do exist false negatives.

For all the programs, we initially set the timeout for schedule steering to $\tau = 1$ ms. As timeouts fired during validation, the hive increased the timeout 50 ms at a time, up to a maximum of 200 ms. Developers may choose to adapt this basic scheme depending on the characteristics of their programs. For instance, the timeout could be increased multiplicatively instead of linearly.

In principle, false negatives may also arise from τ being too low or from there being insufficient executions to prove a true race. We increased τ in our experiments by $4\times$, to check if this would alter our results, and the final verdicts were the same. After manually examining races that were not encountered during dynamic validation, we found that they were either in functions that are never called but are nonetheless linked to the programs, or they are not encountered at runtime due to the workloads used in the evaluation.

6.1.3 Efficiency

The more efficient a detector is, the less runtime overhead it introduces, i.e., the less it slows down a user’s application (as a percentage of uninstrumented execution). The static detection phase is offline, and it took less than 3 minutes for all programs, except Apache and SQLite, for which it took less than 1 hour. Therefore, in this section, we focus on the dynamic phase.

Apache	SQLite	Memcached	Fmm	Barnes	Ocean	Pbzip2	Knot	Aget	Pfscan
1.74	1.60	0.10	4.54	2.98	2.05	2.90	1.27	3.00	3.03

Table 2 – Runtime overhead of race detection as a percentage of uninstrumented execution. Average overhead is 2.32%, and maximum overhead is 4.54%.

Program	Apache	SQLite	Memcached	Fmm	Barnes	Ocean	Pbzip2	Knot	Aget
RaceMob	8	3	1	58	16	3	9	2	4
TSAN	8	3	0	58	16	3	9	2	2
RELAY	118	88	7	176	166	115	65	157	256

Table 3 – Race detection results with RaceMob, ThreadSanitizer (TSAN), and RELAY. Each cell shows the number of reported races. The data races reported by RaceMob and TSAN are all true data races. The only true data races among the ones detected by RELAY are the ones in the row “RaceMob”. To the best of our knowledge, two of the data races that cause a hang in SQLite were not previously reported.

Table 2 shows that runtime overhead of RaceMob is typically less than 3%. The static analysis used to remove instrumentation from empty loop bodies reduced our worst case overhead from 25% to 4.54%. The highest runtime overhead is 4.54%, in the case of Fmm, a memory-intensive application that performs repetitive computations, which gives the instrumentation more opportunity to introduce overhead. Our results suggest that there is no correlation between the number of race candidates (row 2 in Table 1) and the runtime overhead (Table 2)—overhead is mostly determined by the frequency of execution of the instrumentation code.

The overhead introduced by RaceMob is due to the instrumentation plus the overhead introduced by validation (DCI, on-demand detection, and schedule steering). Fig. 18 shows the breakdown of overhead for our ten target programs. We find that the runtime overhead without detection is below 1% for all cases, except the memory-intensive Fmm application, for which it is 2.51%. We conclude that, in the common case when a program is instrumented by RaceMob but no detection is performed, the runtime overhead is negligible; this property is what makes RaceMob suitable for always-on operation.

The dominant component of the overhead of race detection (the black portion of the bars in Fig. 18) is due to dynamic data race validation. The effect of DCI is negligible: it is below 0.1% for all cases; thus, we don’t show it in Fig. 18. Therefore, it is feasible to leave DCI

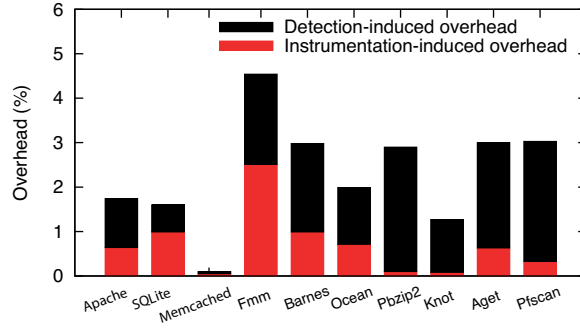


Figure 18 – Breakdown of average overhead into instrumentation-induced overhead and detection-induced overhead.

on for all executions. This can help RaceMob to promote a race from “Likely FP” to “True Race” with low overhead.

If RaceMob assigns more than one validation task at a time per user, the aggregate overhead that a user experiences will increase. In such a scenario, the user site would pick a validation candidate at runtime depending on which potentially racing access is executed. This scheme introduces a lookup overhead to determine at runtime which racing access is executed, however, it would not affect the per-race overhead, because of RaceMob’s on-demand race detection algorithm.

6.1.4 Comparison to Existing Data Race Detectors

RaceMob uses atomic operations to update internal shared structures related to dynamic data race validation and signal-wait synchronization to perform schedule steering; in this section, we analyze the effect these operations have on RaceMob’s scalability as the number of application threads increases.

We configured multiple clients to concurrently request a 10 MB file from Apache and Knot using the Apache benchmarking tool `ab`. For SQLite and Memcached, we inserted, modified, and removed 5,000 items from the database and the object cache, respectively. We used Pbzip2 to decompress a 100 MB file. For Ocean, we simulated currents in a 256×256 ocean grid. For Barnes, we simulated interactions of 16,384 bodies (default number for Barnes). We varied the number of threads from 2 – 32. For all programs, we ran the instrumented versions of the programs while performing data race detection and measured the overhead relative to uninstrumented versions on the 8-core machine.

Fig. 19 shows the results. We expected RaceMob’s overhead to become less visible after the thread count reached the core count. We wanted to verify this, and that is why we used the 8-core machine. For instance, for Apache the overhead is 1.16% for 2 threads, it slightly rises to its largest value of 2.31% for 8 threads, and then it decreases

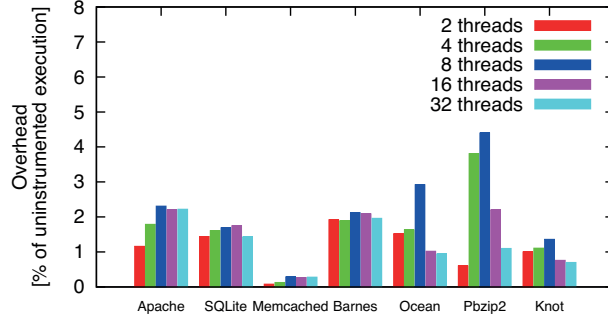


Figure 19 – RaceMob scalability: Induced overhead as a function of the number of application threads.

as the number of threads exceeds the number of cores. We observe a similar trend for all other applications. We conclude that RaceMob’s runtime overhead remains low as the number of threads in the test programs increases.

6.2 ROOT CAUSE DIAGNOSIS RESULTS

In this section we aim to answer the following questions about Gist and failure sketching: Is Gist capable of automatically computing failure sketches (§6.2.2)? Are these sketches accurate (§6.2.3)? How efficient is the computation of failure sketches in Gist (§6.2.4)?

6.2.1 Experimental Setup

To answer these questions we benchmark Gist with several real world programs: Apache httpd [2] is a popular web server. Cppcheck [cppcheck] is a C/C++ static analysis tool integrated with popular development tools such as Visual Studio, Eclipse, and Jenkins. Curl [curl] is a data transfer tool for network protocols such as FTP and HTTP, and it is part of most Linux distributions and many programs, like LibreOffice and CMake. Transmission [transmission] is the default BitTorrent client in Ubuntu and Fedora Linux, as well as Solaris. SQLite [97] is an embedded database used in Chrome, Firefox, iOS, and Android. Memcached is a distributed memory object cache system used by services such as Facebook and Twitter [28]. Pbzip2 [33] is a parallel file compression tool.

We developed an extensible framework called Bugbase [7] in order to reproduce the known bugs in the aforementioned software. Bugbase can also be used to do performance benchmarking of various bug finding tools. We used Bugbase to obtain our experimental results.

We benchmark Gist on bugs (from the corresponding bug repositories) that were used by other researchers to evaluate their bug finding and failure diagnosis tools [3, 81, 53]. Apart from bugs in Cppcheck

Failure Sketch for Curl bug #965		
Type: Sequential bug, data-related		
Time		url
1 operate(struct char* url, ...){		1
2 for(i = 0; (url = next_url(urls)); i++){		2 { }
3 }		3
4 }		4
<div style="display: flex; align-items: center;"> <div style="flex: 1;"> 5 next_url(urls* urls){ 6 len = strlen(urls->current); 7 } </div> <div style="flex: 1; border-left: 1px dotted black; padding-left: 5px; margin-left: 10px;"> horizontal line separates different functions </div> </div>		urls->current
		5
		6 0
		7
		Failure (segmentation fault)

Figure 20 – The failure sketch of Curl bug #965.

and Curl, all bugs are concurrency bugs. We use a mixture of workloads from actual program runs, test suites, test cases devised by us and other researchers [114], Apache’s benchmarking tool ab, and SQLite’s test harness. We gathered execution information from a total of 11,360 executions.

The distributed cooperative setting of our test environment is simulated, as opposed to employing real users, because CPUs with Intel PT support are still scarce, having become available only recently. In the future we plan to use a real-world deployment. Altogether we gathered execution information from 1,136 simulated user endpoints. Client-side experiments were run on a 2.4 GHz 4 core Intel i7-5500U (Broadwell) machine running a Linux kernel with an Intel PT driver [intelpt:branch]. The server side of Gist ran on a 2.9 GHz 32-core Intel Xeon E5-2690 machine with 256 GB of RAM running Linux kernel 3.13.0-44.

6.2.2 Automated Generation of Sketches

For all the failures shown in Table ??, Gist successfully computed the corresponding failure sketches after gathering execution information from 11,360 runs in roughly 35 minutes. The results are shown in the rightmost two columns. We verified that, for all sketches computed by Gist, the failure predictors with the highest F-measure indeed correspond to the root causes that developers chose to fix.

In the rest of this section, we present two failure sketches computed by Gist, to illustrate how developers can use them for root cause diagnosis and for fixing bugs. These two complement the failure sketch for the Pbzip2 bug already described in §??. Aside from some formatting, the sketches shown in this section are exactly the output of Gist. We renamed some variables and functions to save space in the figures. The statements or variable values in dotted rectangles denote failure predicting events with the highest F-measure values. We integrated Gist with KCachegrind [kcache:grind], a call graph viewer that allows easy navigation of the statements in the failure sketch.

Failure Sketch for Apache bug #21287				
Type: Concurrency bug, double-free				
Time	Thread T ₁	Thread T ₂	obj->refcnt	
1	decrement_refcount(obj){	1 decrement_refcount(obj){	1	
2	if (!obj->complete) {	2 if (!obj->complete) {	2	
3	object_t *mobj = ...	3 object_t *mobj = ...	3	
4	dec(&obj->refcnt);	4	4	1
5		5 dec(&obj->refcnt);	5	0
6		6 if (!obj->refcnt) {	6	
7		7 free(obj);	7	
8	if (!obj->refcnt) {	8 }	8	
9	free(obj);	9 }	9	
	} Failure (double free)			

Figure 21 – The failure sketch of Apache bug #21287. The grayed-out components are not part of the ideal failure sketch, but they appear in the sketch that Gist automatically computes.

Fig. 20 shows the failure sketch for Curl bug #965, a sequential bug caused by a specific program input: passing the string “{}{” (or any other string with unbalanced curly braces) to Curl causes the variable `urls->current` in function `next_url` to be NULL in step 6. The value of `url` in step 2 (“{}{”) and the value of `urls->current` in step 6 (0) are the best failure predictors. This failure sketch suggests that fixing the bug consists of either disallowing unbalanced parentheses in the input url, or not calling `strlen` when `urls->current` is NULL. Developers chose the former solution to fix this bug [curl-965].

Fig. 21 shows the failure sketch for Apache bug 21287, a concurrency bug causing a double free. The failure sketch shows two threads executing the `decrement_refcount` function with the same `obj` value. The `dec` function decrements `obj->refcount`. The call to `dec`, the `if` condition checking, namely `!obj->refcount`, and the call to `free` are not atomic, and this can cause a double free if `obj->refcount` is 0 in step 6 in T₂ and step 8 in T₁. The values of `obj->refcount` in steps 4 and 5 (1 and 0 respectively), and the double call to `free(obj)` are the best failure predictors. Developers fixed this bug by ensuring that the decrement-check-free triplet is executed atomically [apache-21287].

The grayed-out statements in the failure sketch in Fig. 21 are not part of the ideal failure sketch. The adaptive slice tracking part of Gist tracks them during slice refinement, because Gist does not know the statements in the ideal failure sketch a priori. For the Curl bug in Fig. 20, we do not show any grayed-out statements, because, adaptive slice tracking happens to track only the statements that are in the ideal failure sketch.

6.2.3 Accuracy of Failure Sketches

In this section, we measure the accuracy (A) of failure sketches computed by Gist (Φ_G), as compared to ideal failure sketches that

we computed by hand (Φ_I), according to our ideal failure sketch definition (§??). We define two components of failure sketch accuracy:

1) **Relevance** measures the extent to which a failure sketch contains all the statements from the ideal sketch and no other statements. We define relevance as the ratio of the number of LLVM instructions in $\Phi_G \cap \Phi_I$ to the number of statements in $\Phi_G \cup \Phi_I$. We compute relevance accuracy as a percentage, and define it as $A_R = 100 \cdot \frac{|\Phi_G \cap \Phi_I|}{|\Phi_G \cup \Phi_I|}$.

2) **Ordering** measures the extent to which a failure sketch correctly represents the partial order of LLVM memory access instructions in the ideal sketch. To measure the similarity in ordering between the Gist-computed failure sketches and their ideal counterparts, we use the normalized Kendall tau distance [kendalltau] τ , which measures the number of pairwise disagreements between two ordered lists. For example, for ordered lists $\langle A, B, C \rangle$ and $\langle A, C, B \rangle$, the pairs (A, B) and (A, C) have the same ordering, whereas the pair (B, C) has different orderings in the two lists, hence $\tau = 1$. We compute the ordering accuracy as a percentage defined by $A_O = 100 \cdot (1 - \frac{\tau(\Phi_G, \Phi_I)}{\# \text{ of pairs in } \Phi_G \cap \Phi_I})$. Note that # of pairs in $\Phi_G \cap \Phi_I$ can't be zero, because both failure sketches will at least contain the failing instruction as a common instruction.

We define overall accuracy as $A = \frac{A_R + A_O}{2}$, which equally favors A_O and A_R . Of course, different developers may have different subjective opinions on which one matters most.

We show Gist's accuracy results in Fig. 22. Average relevance accuracy is 92%, average ordering accuracy is 100%, and average overall accuracy is 96%, which leads us to conclude that Gist can compute failure sketches with high accuracy. The accuracy results are deterministic from one run to the next.

Note that, for all cases when relevance accuracy is below 100%, it is because Gist's failure sketches have (relative to the ideal sketches) some excess statements in the form of a prefix to the ideal failure sketch, as shown in gray in Fig. 21. We believe that developers find it significantly easier to visually discard excess statements clustered as a prefix than excess statements that are sprinkled throughout the failure sketch, so this inaccuracy is actually not of great consequence.

We show in Fig. 23 the contribution of Gist's three analysis and tracking techniques to overall sketch accuracy. To obtain these measurements, we first measured accuracy when using just static slicing, then enabled control flow tracking and re-measured, and finally enabled also data flow tracking and re-measured. While the accuracy results are consistent across runs, the individual contributions may vary if, for example, workload non-determinism causes different paths to be exercised through the program.

A small contribution of a particular technique does not necessarily mean that it does not perform well for a given program, but it means that the other techniques that Gist had enabled prior to this

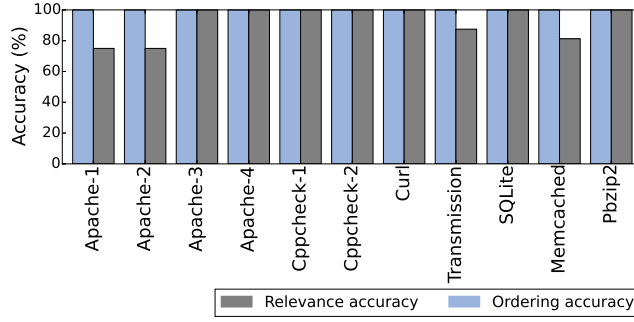


Figure 22 – Accuracy of Gist, broken down into relevance accuracy and ordering accuracy.

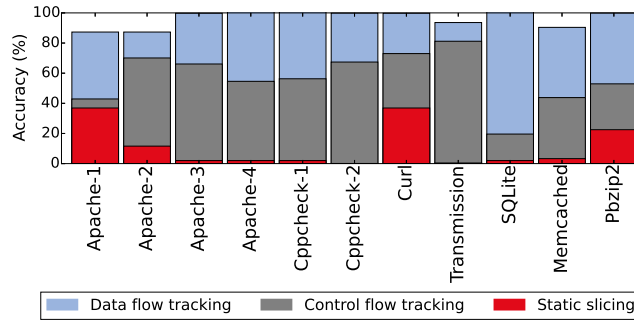


Figure 23 – Contribution of various techniques to Gist's accuracy.

technique “stole its thunder” by being sufficient to provide high accuracy. For example, in the case of Apache-1, static analysis performs well enough that control flow tracking does not need to further refine the slice. However, in some cases (e.g., for SQLite), tracking the inter-thread execution order of statements that access shared variables using hardware watchpoints is crucial for achieving high accuracy.

We observe that the amount of individual contribution varies substantially from one program to the next, which means that neither of these techniques would achieve high accuracy for all programs on its own, and so they are all necessary if we want high accuracy across a broad spectrum of software.

6.2.4 Efficiency

Now we turn our attention to the efficiency of Gist: how long does it take to compute a failure sketch, how much runtime performance overhead does it impose on clients, and how long does it take to perform its offline static analysis. We also look at how these measures vary with different parameters.

The last column of Table ?? shows Gist's failure sketch computation latency broken down into three components. We show the number of failure recurrences required to reach the best sketch that Gist can compute, and this number varies from 2 to 5 recurrences. We then

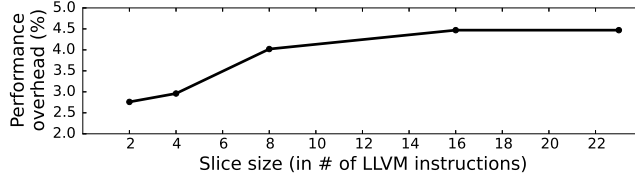


Figure 24 – Gist’s average runtime performance overhead across all runs as a function of tracked slice size.

show the total time it took in our simulated environment to find this sketch; this time is always less than 6 minutes, varying from $\langle 0m:23s \rangle$ to $\langle 5m:34s \rangle$. Not surprisingly, this time is dominated by how long it takes the target failure to recur, and in practice this depends on the number of deployed clients and the variability of execution circumstances. Nevertheless, we present the values for our simulated setup to give an idea as to how long it took to build a failure sketch for each bug in our evaluation. Finally, in parentheses we show Gist’s offline analysis time, which consists of computing the static slice plus generating instrumentation patches. This time is always less than 3 minutes, varying between $\langle 0m:2s \rangle$ and $\langle 2m:32s \rangle$. We therefore conclude that, compared to the debugging latencies experienced by developers today, Gist’s automated approach to root cause diagnosis presents a significant advantage.

In the context of adaptive slice tracking, the overhead incurred on the client side increases monotonically with the size of the tracked slice, which is not surprising. Fig. 24 confirms this experimentally. The portion of the overhead curve between the slice sizes 16 and 22 is relatively flat compared to the rest of the curve. This is because, within that interval, Gist only tracks a few control flow events for Apache-1 and Curl (these programs have no additional data flow elements in that interval), which introduces negligible overhead.

The majority of the overhead incurred on the client side stems from control flow tracking. In particular, the overhead of control flow tracking varies from a low of 2.01% to a high of 3.43%, whereas the overhead of data flow tracking varies from a low of 0.87% to a high of 1.04%.

What is perhaps not immediately obvious is the trade-off between initial slice size σ and the resulting accuracy and latency. In Fig. 25, we show the average failure sketch accuracy across all programs we measured (right y-axis) and Gist’s latency in # of recurrences (left y-axis) as a function of σ that Gist starts with (x-axis). As long as the initial slice size is less than the one for the best sketch that Gist can find, Gist’s adaptive approach is capable of guiding the developer to the highest accuracy sketch. Of course, the time it takes to find the sketch is longer the smaller the starting slice size is, because the necessary # of recurrences is higher. There is thus an incentive to

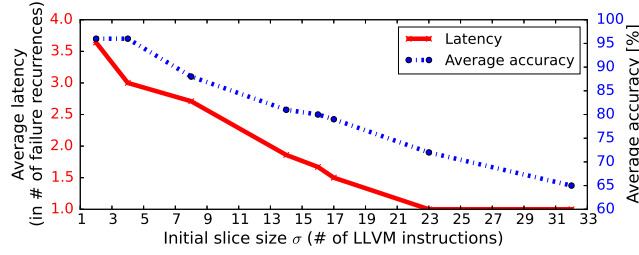


Figure 25 – Tradeoff between slice size and the resulting accuracy and latency. Accuracy is in percentage, latency is in the number of failure recurrences.

start with a larger slice size. Unfortunately, if this size overshoots the size of the highest accuracy sketch, then the accuracy of the outcome suffers, because the larger slice includes extraneous elements.

As we mentioned in §??, the extraneous statements that can lower Gist’s accuracy are clustered as a prefix to the ideal failure sketch, allowing developers to easily ignore them. Therefore, if lower root cause diagnosis latency is paramount to the developers, they are comfortable ignoring the prefix of extraneous statements, and they can tolerate the slight increase in Gist’s overhead, it is reasonable to configure Gist to start with a large σ (e.g., $\sigma = 23$ achieves a latency of *one failure recurrence* for all our benchmarks).

For the benchmarks in our evaluation, starting AsT at $\sigma = 4$ would achieve the highest average accuracy at the lowest average latency of 3, with an average overhead of 3.98%.

Finally, Fig. 26 compares Intel PT, the hardware-based control flow tracking mechanism we use in Gist, to Mozilla rr, a software-based state-of-the-art record & replay system. In particular, we compare the performance overhead imposed by the two tracking mechanisms on the client application. The two extremes are Cppcheck, where Mozilla rr is on par with Intel PT, and Transmission and SQLite, where Mozilla rr’s overhead is over many orders of magnitude higher than Intel PT’s¹. For the benchmarks in our evaluation, full tracing using Intel PT incurs an average overhead of 11%, whereas full program record & replay incurs an average runtime overhead of 984%. Unlike Intel PT, Mozilla rr also gathers data flow information, but with Gist we have shown that full program tracing is not necessary for automating root cause diagnosis.

In conclusion, our empirical evaluation shows that Gist, a failure sketching prototype, is capable of automatically computing failure sketches for failures caused by real bugs in real systems (§??), these sketches have a high accuracy of 96% on average (§??), and the average performance overhead of failure sketching is 3.74% with $\sigma = 2$

1. Full tracing overheads of Transmission and SQLite are too low to be reliably measured for Intel PT, thus they are shown as 0%, and the corresponding Mozilla rr/Intel PT overheads for these systems are shown as ∞ .

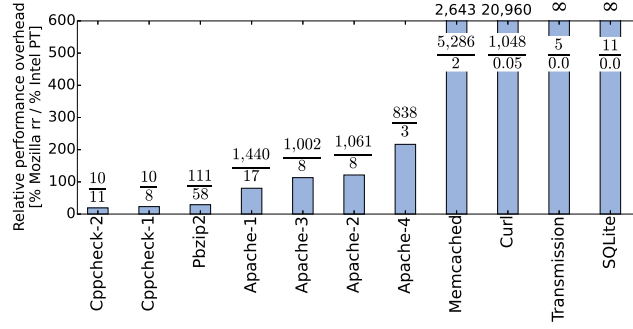


Figure 26 – Comparison of the full tracing overheads of Mozilla rr and Intel PT.

(§??). We therefore believe failure sketching to be a promising approach for helping developers debug elusive bugs that occur only in production.

6.3 DATA RACE CLASSIFICATION RESULTS

In this section, we answer the following questions: Is Portend effective in telling developers which races are true bugs and in helping them fix buggy races (§??)? How accurately does it classify race reports into the four categories of races (§6.3)? How long does classification take, and how does it scale (§??)? How does Portend compare to the state of the art in race classification (§??)? How effectively and efficiently does Portend implement symbolic memory consistency modeling and what is its memory overhead (§??, §6.3.7)?. Throughout this section, we highlight the synergy of the techniques used in Portend: in particular §?? shows how symbolic output comparison allows more accurate race classification compared to post-race state comparison, and §6.3 shows how the combination of multi-path multi-schedule analysis improves upon traditional single-path analysis.

6.3.1 Experimental Setup

We apply Portend to 7 applications: SQLite, an embedded database engine (used, for example, by Firefox, iOS, Chrome, and Android), that is considered highly reliable, with 100% branch coverage [97]; Pbzip2, a parallel implementation of the widely used bzip2 file compressor [33]; Memcached [28], a distributed memory object cache system (used, for example, by services such as Flickr, Twitter and Craigslist); Ctrace [71], a multi-threaded debug library; Bbuf [112], a shared buffer implementation with a configurable number of producers and consumers; Fmm, an n-body simulator from the popular SPLASH2 benchmark suite [107]; and Ocean, a simulator of eddy currents in oceans, from SPLASH2.

We additionally evaluate Portend on homegrown micro-benchmarks that capture most classes of races considered as harmless in the literature [94, 77]: “redundant writes” (RW), where racing threads write the same value to a shared variable, “disjoint bit manipulation” (DBM), where disjoint bits of a bit-field are modified by racing threads, “all values valid” (AVV), where the racing threads write different values that are nevertheless all valid, and “double checked locking” (DCL), a method used to reduce the locking overhead by first testing the locking criterion without actually acquiring a lock. Additionally, we have 4 other micro-benchmarks that we used to evaluate the SMCM. We detail those micro-benchmarks in §??. Table 4 summarizes the properties of our 15 experimental targets.

Program	Size (LOC)	Language	# Forked threads
SQLite 3.3.0	113,326	C	2
ocean 2.0	11,665	C	2
fmm 2.0	11,545	C	3
memcached 1.4.5	8,300	C	8
pbzip2 2.1.1	6,686	C++	4
ctrace 1.2	886	C	3
bbuf 1.0	261	C	8
AVV	49	C++	3
DCL	45	C++	5
DBM	45	C++	3
RW	42	C++	3
no-sync	45	C++	3
no-sync-bug	46	C++	3
sync	47	C++	3
sync-bug	48	C++	3

Table 4 – Programs analyzed with Portend. Source lines of code are measured with the `cloc` utility.

We ran Portend on several other systems (e.g., HawkNL, pfscan, swarm, fft), but no races were found in those programs with the test cases we ran, so we do not include them here. For all experiments, the Portend parameters were set to $M_p = 5$, $M_a = 2$, and the number of symbolic inputs to 2. We found these numbers to be sufficient to achieve high accuracy in a reasonable amount of time. To validate Portend’s results, we used manual investigation, analyzed developer change logs, and consulted with the applications’ developers when possible. All experiments were run on a 2.4 GHz Intel Core 2 Duo E6600 CPU with 4 GB of RAM running Ubuntu Linux 10.04 with kernel version 2.6.33. The reported numbers are averages over 10 experiments.

6.3.2 Effectiveness

Of the 93 distinct races detected in 7 real-world applications, Portend classified 5 as definitely harmful by watching for “basic” properties (Table 5): one hangs the program and four crash it.

Program	Total # of races	# of “Spec violated” races		
		Deadlock	Crash	Semantic
SQLite	1	1	0	0
pbzip2	31	0	3	0
ctrace	15	0	1	0
Manually inserted errors				
fmm	13	0	0	1
memcached	18	0	1	0

Table 5 – “Spec violated” races and their consequences.

To illustrate the checking for “high level” semantic properties, we instructed Portend to verify that all timestamps used in fmm are positive. This caused it to identify the 6th “harmful” race in Table 5; without this semantic check, this race turns out to be harmless, as the negative timestamp is eventually overwritten.

To illustrate a “what-if analysis” scenario, we turned an arbitrary synchronization operation in the memcached binary into a no-op, and then used Portend to explore the question of whether it is safe to remove that particular synchronization point (e.g., we may be interested in reducing lock contention). Removing this synchronization induces a race in memcached; Portend determined that the race could lead to a crash of the server for a particular interleaving, so it classified it as “spec violated”.

Portend’s main contribution is the classification of races. If one wanted to eliminate all harmful races from their code, they could use a static race detector (one that is complete, and, by necessity, prone to false positives) and then use Portend to classify these reports.

For every harmful race, Portend’s comprehensive report and replayable traces (i.e., inputs and thread schedule) allowed us to confirm the harmfulness of the races within minutes. Portend’s report includes the stack traces of the racing threads along with the address and size of the accessed memory field; in the case of a segmentation fault, the stack trace of the faulting instruction is provided as well—this information can help in automated bug clustering. According to developers’ change logs and our own manual analysis, the races in Table 5 are the only known harmful races in these applications.

Program	Number of data races						
	Distinct races	Race instances	Spec violated	Output differs	K-witness harmless		Single ordering
					states same	states differ	
SQLite	1	1	1	0	0	0	0
ocean	5	14	0	0	0	1	4
fmm	13	517	0	0	0	1	12
memcached	18	104	0	2	0	0	16
pbzip2	31	97	3	3	0	0	25
ctrace	15	19	1	10	0	4	0
bbuf	6	6	0	6	0	0	0
AVV	1	1	0	0	1	0	0
DCL	1	1	0	0	1	0	0
DBM	1	1	0	0	1	0	0
RW	1	1	0	0	1	0	0

Table 6 – Summary of Portend’s classification results. We consider two races to be distinct if they involve different accesses to shared variables; the same race may be encountered multiple times during an execution—these two different aspects are captured by the *Distinct races* and *Race instances* columns, respectively. Portend uses the stack traces and the program counters of the threads making the racing accesses to identify distinct races. The last 5 columns classify the distinct races. The *states same/differ* columns show for how many races the primary and alternate states were different after the race, as computed by the Record/Replay Analyzer [77].

6.3.3 Accuracy and Precision

To evaluate Portend’s accuracy and precision, we had it classify all 93 races in our target applications and micro-benchmarks. Table 6 summarizes the results. The first two columns show the number of distinct races and the number of respective instances, i.e., the number of times those races manifested during race detection. The “spec violated” column includes all races from Table 5 minus the semantic race in fmm and the race we introduced in memcached. In the “k-witness harmless” column we show for which races the post-race states differed vs. not.

By accuracy, we refer to the correctness of classification: the higher the accuracy, the higher the ratio of correct classification. Precision on the other hand, refers to the reproducibility of experimental results: the higher the precision, the higher the ratio with which experiments are repeated with the same results.

To determine accuracy, we manually classified each race and found that Portend had correctly classified 92 of the 93 races (99%) in our

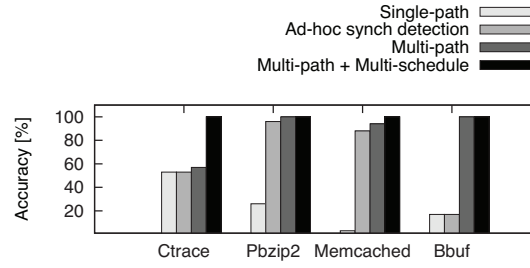


Figure 27 – Breakdown of the contribution of each technique toward Portend’s accuracy. We start from single-path analysis and enable one by one the other techniques: ad-hoc synchronization detection, multi-path analysis, and finally multi-schedule analysis.

target applications: all except one of the races classified “k-witness harmless” by Portend are indeed harmless in an absolute sense, and all “single ordering” races indeed involve ad-hoc synchronization.

To measure precision, we ran 10 times the classification for each race. Portend consistently reported the same data set shown in Table 6, which indicates that, for these races and applications, it achieves full precision.

As can be seen in the “k-witness harmless” column, for each and every one of the 7 real-world applications, a state difference (as used in [77]) does not correctly predict harmfulness, while our “k-witness harmless” analysis correctly predicts that the races are harmless with one exception.

This suggests that differencing of concrete state is a poor classification criterion for races in real-world applications with large memory states, but may be acceptable for simple benchmarks. This also supports our choice of using symbolic output comparison.

Multi-path multi-schedule exploration proved to be crucial for Portend’s accuracy. Fig. 27 shows the breakdown of the contribution of each technique used in Portend: ad-hoc synchronization detection, multi-path analysis, and multi-schedule analysis. In particular, for 16 out of 21 “output differs” races (6 in bbuf, 9 in ctrace, 1 in pbzip2) and for 1 “spec violated” race (in ctrace), single-path analysis revealed no difference in output; it was only multi-path multi-schedule exploration that revealed an output difference (9 races required multi-path analysis for classification, and 8 races required also multi-schedule analysis). Without multi-path multi-schedule analysis, it would have been impossible for Portend to accurately classify those races by just using the available test cases. Moreover, there is a high variance in the contribution of each technique for different programs, which means that none of these techniques alone would have achieved high accuracy for a broad range of programs.

We also wanted to evaluate Portend’s ability to deal with false positives, i.e., false race reports. Race detectors, especially static ones, may report false positives for a variety of reasons, depending on

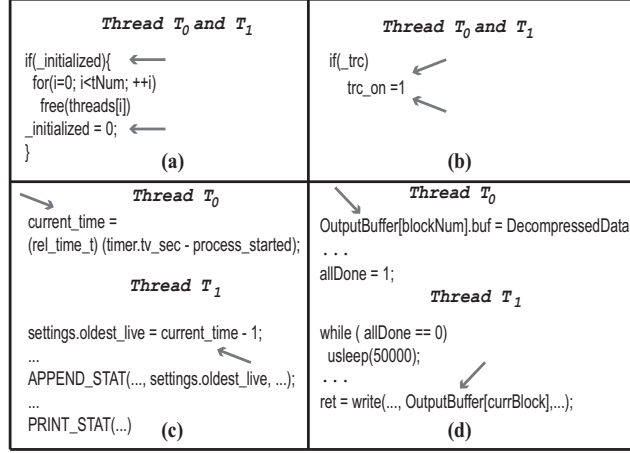


Figure 28 – Simplified examples for each race class from real systems. (a) and (b) are from ctrace, (c) is from memcached and (d) is from pbzip2. The arrows indicate the pair of racing accesses.

which technique they employ. To simulate an imperfect detector for our applications, we deliberately removed from Portend’s race detector its awareness of mutex synchronizations. We then eliminated the races in our micro-benchmarks by introducing mutex synchronizations. When we re-ran Portend with the erroneous data race detector on the micro-benchmarks, all four were falsely reported as races by the detector, but Portend ultimately classified all of them as “single ordering”. This suggests Portend is capable of properly handling false positives.

Fig. 28 shows examples of real races for each category: (a) a “spec violated” race in which resources are freed twice, (b) a “k-witness harmless” race due to redundant writes, (c) an “output differs” race in which the schedule-sensitive value of the shared variable influences the output, and (d) a “single ordering” race showing ad-hoc synchronization implemented via busy wait.

6.3.4 Efficiency

We evaluate the performance of Portend in terms of efficiency and scalability. Portend’s performance is mostly relevant if it is to be used interactively, as a developer tool, and also if used for a large scale bug triage tool, such as in Microsoft’s Windows Error Reporting system [34].

We measure the time it takes Portend to classify the 93 races; Table 7 summarizes the results. We find that Portend classifies all detected data races in a reasonable amount of time, the longest taking less than 11 minutes. For bbuf, ctrace, ocean and fmm, the slowest classification time is due to a race from the “k-witness harmless”

category, since classification into this category requires multi-path multi-schedule analysis.

The second column reports the time it took Cloud9 to interpret the programs with concrete inputs. This provides a sense of the overhead incurred by Portend compared to regular LLVM interpretation in Cloud9. Both data race detection and classification are disabled when measuring baseline interpretation time. In summary, the overhead introduced by classification ranges from $1.1\times$ to $49.9\times$ over Cloud9.

In order to get a sense of how classification time scales with program characteristics, we measured it as a function of program size, number of preemption points, number of branches that depend (directly or indirectly) on symbolic inputs, and number of threads. We found that program size plays almost no role in classification time. Instead, the other three characteristics play an important role. We show in Fig. 29 how classification time varies with the number of dependent branches and the number of preemptions in the schedule (which is roughly proportional to the number of preemption points and the number of threads). Each vertical bar corresponds to the classification time for the indicated data race. We see that, as the number of preemptions and branches increase, so does classification time.

Program	Cloud9 running time (sec)	Portend classification time (sec)		
		<i>Avg</i>	<i>Min</i>	<i>Max</i>
SQLite	3.10	4.20	4.09	4.25
ocean	19.64	60.02	19.90	207.14
fmm	24.87	64.45	65.29	72.83
memcached	73.87	645.99	619.32	730.37
pbzip2	15.30	360.72	61.36	763.43
ctrace	3.67	24.29	5.54	41.08
bbuf	1.81	4.47	4.77	5.82
AVV	0.72	0.83	0.78	1.02
DCL	0.74	0.85	0.83	0.89
DBM	0.72	0.81	0.79	0.83
RW	0.74	0.81	0.81	0.82

Table 7 – Portend’s classification time for the 93 races in Table 6.

We analyzed Portend’s accuracy with increasing values of k and found that $k = 5$ is sufficient to achieve overall 99% accuracy for all the programs in our evaluation. Fig. 30 shows the results for Ctrace, Pbzip2, Memcached, and Bbuf. We therefore conclude that it is possi-

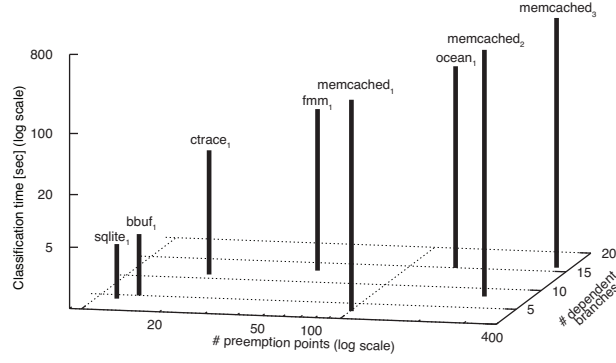


Figure 29 – Change in classification time with respect to number of preemptions and number of dependent branches for some of the races in Table 6. Each sample point is labeled with race id.

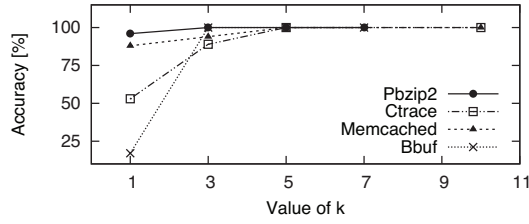


Figure 30 – Portend’s accuracy with increasing values of k.

ble to achieve high classification accuracy with relatively small values of k.

6.3.5 Comparison to Existing Data Race Detectors

T We compare Portend to the Record/Replay-Analyzer technique [77], Helgrind⁺’s technique [45], and Ad-Hoc-Detector [100] in terms of the accuracy with which races are classified. We implemented the Record/Replay-Analyzer technique in Portend and compared accuracy empirically. For the ad-hoc synchronization detection techniques, since we do not have access to the implementations, we analytically derive the expected classification based on the published algorithms. We do not compare to RACEFUZZER [92], because it is primarily a bug finding tool looking for harmful races that occur due to exceptions and memory errors; it therefore does not provide a fine-grain classification of races. Similarly, no comparison is provided to DataCollider [47], since race classification in this tool is based on heuristics that pertain to races that we rarely encountered in our evaluation.

In Table 8 we show the accuracy, relying on manual inspection as “ground truth”. Record/Replay-Analyzer does not tolerate replay failures and classifies races that exhibit a post-race state mismatch as harmful (shown as specViol), causing it to have low accuracy (10%) for that class. When comparing to Helgrind⁺ and Ad-Hoc-Detector,

we conservatively assume that these tools incur no false positives when ad-hoc synchronization is present, even though this is unlikely, given that both tools rely on heuristics. This notwithstanding, both tools are focused on weeding out races due to ad-hoc synchronization, so they cannot properly classify the other races (36 out of 93). In contrast, Portend classifies a wider range of races with high accuracy.

	specViol	k-witness	outDiff	singleOrd
Ground Truth	100%	100%	100%	100%
Record/Replay Analyzer	10%	95%	- (not-classified)	
Ad-Hoc-Detector, Helgrind⁺	- (not-classified)			100%
Portend	100%	99%	99%	100%

Table 8 – Accuracy for each approach and each classification category, applied to the 93 races in Table 6. “Not-classified” means that an approach cannot perform classification for a particular class.

The main advantage of Portend over Record/Replay-Analyzer is that it is immune to replay failures. In particular, for all the races classified by Portend as “single ordering”, there was a replay divergence (that caused replay failures in Record/Replay-Analyzer), which would cause Record/Replay-Analyzer to classify the corresponding races as harmful despite them exhibiting no apparent harmfulness; this accounts for 57 of the 84 misclassifications. Note that even if Record/Replay-Analyzer were augmented with a phase that pruned “single ordering” races (57/93), it would still diverge on 32 of the remaining 36 races and classify them as “spec violated”, whereas only 5 are actually “spec violated”. Portend, on the other hand, correctly classifies 35/36 of those remaining races. Another advantage is that Portend classifies based on symbolic output comparison, not concrete state comparison, and therefore, its classification results can apply to a range of inputs rather than a single input.

We manually verified and, when possible, checked with developers that the races in the “k-witness harmless” category are indeed harmless. Except for one race, we concluded that developers intentionally left these races in their programs because they considered them harmless. These races match known patterns [77, 47], such as redundant writes to shared variables (e.g., we found such patterns in Ctrace). However, for one race in Ocean, we confirmed that Portend did not figure out that the race belongs in the “output differs” category (the race can produce different output if a certain path in the code is followed, which depends indirectly on program input). Portend was not able to find this path even with $k = 10$ after one hour. Manual

investigation revealed that this path is hard to find because it requires a very specific and complex combination of inputs.

6.3.6 Efficiency and Effectiveness of Symbolic Memory Consistency Modeling

The previous evaluation results were using the SMCM plugin in the sequential consistency mode. The sequential memory consistency mode is the default in Cloud9 as well as in Portend. In this section, we answer the following questions while operating the SMCM plugin in Portend's weak consistency mode: (1) Is Portend effective in discovering bugs that may surface under its weak consistency model?, (2) What is Portend's efficiency and (3) memory usage while operating the SMCM plugin in Portend's weak consistency mode?

We use simple micro-benchmarks that we have constructed to test the basic functionality of SMCM. The simplified source code for these micro-benchmarks can be seen in Fig. 31. These benchmarks are:

- *no-sync*: The source code for this micro-benchmark can be seen in Fig. 31-a: A program with opportunities for write reordering. Reorderings cause the `printf` statement on line 17 to produce different program outputs.
- *no-sync-bug*: The source code for this benchmark can be seen in Fig. 31-b: A program with opportunities for write reordering. A particular write reordering causes the program to crash; however the program does not crash under sequential consistency.
- *sync*: The source code for this micro-benchmark can be seen in Fig. 31-c: A program with no opportunities for write reordering. There is a race on both `globalx` and `globaly`. Since both threads 1 and 2 write the same value 2 to `globalx`, the output of the program is the same for any execution, assuming writes are atomic².
- *sync-bug*: The source code for this micro-benchmark can be seen in Fig. 31-d: A program with opportunities for write reordering. The barrier synchronization does not prevent the write to `globalx` and `globaly` from reordering. A particular write reordering causes the program to crash; however the program does not crash under sequential consistency.

To evaluate Portend's effectiveness in finding bugs that may only arise under Portend's weak ordering, we ran the micro-benchmarks with Portend's SMCM plugin configured in two modes: sequential consistency (Portend-seq) and Portend's weak consistency (Portend-weak). We provide the number of bugs found by each configuration of Portend and also the percentage of possible execution states

2. If writes are non-atomic, even a seemingly benign race, where two threads write the same value to a shared variable, may end up producing unexpected results. Details can be found in [13]

<pre> 1: int volatile globalx = 0; 2: int volatile globaly = 0; Thread T₁ 3: void * work0 (void *arg) { 4: globalx = 2; 5: globaly = 1; 6: return 0; 7: } Thread T₂ 8: void * work1 (void *arg) { 9: globalx = 2; 10: return 0; 11: } Thread Main 12: int main (int argc, char *argv[]){ 13: pthread_t t0, t1; 14: int rc; 15: rc = pthread_create (&t0, 0, work0, 0); 16: rc = pthread_create (&t1, 0, work1, 0); 17: printf("%d,%d", globalx, globaly); 18: pthread_join(t0, 0); 19: pthread_join(t1, 0); 20: return 0; 21: } </pre> <p style="text-align: center;">(a)</p>	<pre> 1: int volatile globalx = 0; 2: int volatile globaly = 0; Thread T₁ 3: void * work0 (void *arg) { 4: globalx = 2; 5: globaly = 1; 6: return 0; 7: } Thread T₂ 8: void * work1 (void *arg) { 9: globalx = 2; 10: return 0; 11: } Thread Main 12: int main (int argc, char *argv[]){ 13: pthread_t t0, t1; 14: int rc; 15: rc = pthread_create (&t0, 0, work0, 0); 16: rc = pthread_create (&t1, 0, work1, 0); 17: if(globalx == 0 && globaly == 2) 18: ; //crash! 19: pthread_join(t0, 0); 20: pthread_join(t1, 0); 21: return 0; 22: } </pre> <p style="text-align: center;">(b)</p>
<pre> 1: int volatile globalx = 0; 2: int volatile globaly = 0; Thread T₁ 3: void * work0 (void *arg) { 4: globalx = 2; 5: pthread_barrier_wait(&barr); 6: globaly = 1; 7: return 0; 8: } Thread T₂ 9: void * work1 (void *arg) { 10: globalx = 2; 11: pthread_barrier_wait(&barr); 12: return 0; 13: } Thread Main 13: int main (int argc, char *argv[]){ 14: pthread_t t0, t1; 15: int rc; 16: pthread_barrier_init(&barr, NULL, 2); 17: rc = pthread_create (&t0, 0, work0, 0); 18: rc = pthread_create (&t1, 0, work1, 0); 19: printf("%d,%d", globalx, globaly); 20: pthread_join(t0, 0); 21: pthread_join(t1, 0); 22: return 0; 23: } </pre> <p style="text-align: center;">(c)</p>	<pre> 1: int volatile globalx = 0; 2: int volatile globaly = 0; Thread T₁ 3: void * work0 (void *arg) { 4: globalx = 2; 5: globaly = 1; 6: pthread_barrier_wait(&barr); 7: return 0; 8: } Thread T₂ 9: void * work1 (void *arg) { 10: globalx = 2; 11: pthread_barrier_wait(&barr); 12: return 0; 13: } Thread Main 13: int main (int argc, char *argv[]){ 14: pthread_t t0, t1; 15: int rc; 16: pthread_barrier_init(&barr, NULL, 2); 17: rc = pthread_create (&t0, 0, work0, 0); 18: rc = pthread_create (&t1, 0, work1, 0); 19: if(globalx == 0 && globaly == 2) 20: ; //crash! 21: pthread_join(t0, 0); 22: pthread_join(t1, 0); 23: return 0; 24: } </pre> <p style="text-align: center;">(d)</p>

Figure 31 – Micro-benchmarks used for evaluating SMCM.

that each configuration covers if Portend’s weak consistency were assumed. Note that ground truth (that is the total number of states that can be covered under Portend’s weak consistency) in this case is manually identified as the number of possible states is small. Effectively identifying this percentage for arbitrarily large programs is undecidable.

We present the results in Table 9. As it can be seen, Portend-weak discovers the bugs that can only be discovered under Portend’s weak consistency whereas Portend-seq cannot find those bugs because of sequential consistency assumptions. Similar reasoning applies to state exploration. Portend covers all the possible states that may arise from returning multiple values at “read”s whereas Cloud9 simply returns the last value that was written to a memory location and hence has lower coverage.

System	Number of bugs		State coverage (%)	
	Portend-seq	Portend-weak	Portend-seq	Portend-weak
no-sync	0/0	0/0	50	100
no-sync-bug	0/1	1/0	50	100
sync	0/0	0/0	100	100
sync-bug	0/1	1/1	50	100

Table 9 – Portend’s effectiveness in bug finding and state coverage for two memory model configurations: sequential memory consistency mode and Portend’s weak memory consistency mode.

We also evaluate the performance of Portend-weak for our micro-benchmarks and compare its running time to that of Portend-seq. The results of this comparison can be seen in Fig. 32. The running times of the benchmarks under Portend-seq essentially represent the “native” LLVM interpretation time. For the *no-sync* benchmark we can see that the running time of Portend is about 2 seconds more than that of Portend-seq. This is expected as Portend-weak covers more states compared to Portend-seq.

However, it should be noted that in the case of the *no-sync-bug* benchmark, the running times are almost the same for Portend-weak and Portend-seq (although not visible on the graph, the running time of Portend-weak is slightly larger than that of Portend-seq, on the order of a few milliseconds). This is simply due to the fact that the bug in *no-sync-bug* is immediately discovered after the exploration state has been forked in Portend. The bug is printed at the program output, and the exploration ends for that particular state. Similar reasoning applies to the other benchmark pair, namely *sync* and *sync-bug*.

6.3.7 Memory Consumption of Symbolic Memory Consistency Modeling

In this final section of the evaluation we measure the peak memory consumption of Portend-weak and Portend-seq for the micro-benchmarks we have tested. The results can be seen in Fig. 33. The memory consumption increases for all the benchmarks. This is be-

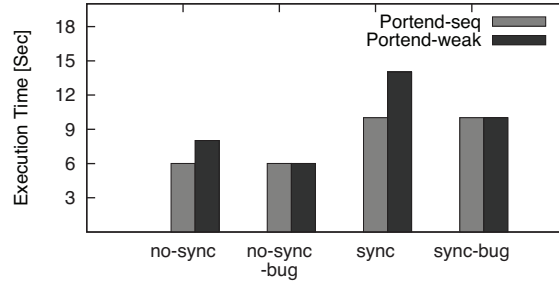


Figure 32 – Running time of Portend-weak and Portend-seq

cause for all the benchmarks, Portend-weak always forks off more states and/or performs more bookkeeping than Portend-seq, even though it does not always explore those states.

Although the memory consumption consistently increases for Portend-weak, it does not increase proportionally with the state forking. This is possible due to the copy-on-write mechanism employed for exploring states and keeping track of the happens-before graph. However, when we ran Portend-weak on real world programs, the memory consumption quickly exceeded the memory capacity of the workstation we used for our experiments. We plan on incorporating techniques like partial order reduction from model checking in order to reduce the number of states that SMCM needs to explore and improve its scalability as part of future work.

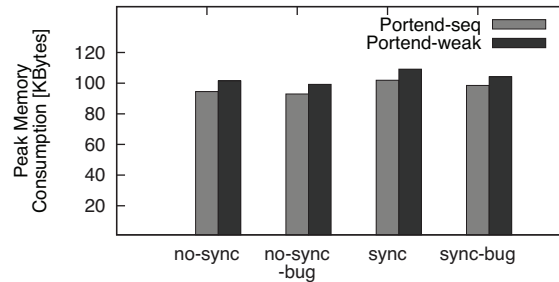


Figure 33 – Memory usage of Portend-weak and Portend-seq

In summary, Portend is able to classify with 99% accuracy and full precision all the 93 races into four data race classes defined in §5.1 in under 5 minutes per race on average. Furthermore, Portend correctly identifies 6 serious harmful races. Compared to previously published race classification methods, Portend performs more accurate classification (92 out of 93, 99%) and is able to correctly classify up to 89% more data races than existing replay-based tools (9 out of 93, 10%). Portend also correctly handles false positive race reports. Furthermore, SMCM allows Portend to accurately perform race classification under relaxed memory consistency models, with low overhead.

Part III

WRAPPING UP

In this final part, we discuss ongoing and future work and we present concluding remarks.

FUTURE WORK

Adversarial compiler, security applications of control flow tracking, process-level atomicity violations, privacy implications of in-production schemes,

CONCLUSIONS

BIBLIOGRAPHY

- [1] Rahul Agarwal et al. "Optimized Run-time Race Detection and Atomicity Checking Using Partial Discovered Types." In: *ASE*. 2005.
- [2] *Apache httpd*. <http://httpd.apache.org>. 2013.
- [3] Joy Arulraj, Guoliang Jin, and Shan Lu. "Leveraging the Short-term Memory of Hardware to Diagnose Production-run Software Failures." In: (2014).
- [4] Joy Arulraj et al. "Production-run Software Failure Diagnosis via Hardware Performance Counters." In: *ASPLOS*. 2013.
- [5] Gogul Balakrishnan and Thomas Reps. "Analyzing Memory Accesses in x86 Executables." In: *Intl. Conf. on Compiler Construction*. 2004.
- [6] Baris Kasikci. *Are "data races" and "race condition" actually the same thing in context of concurrent programming*. <http://stackoverflow.com/questions/11276259/are-data-races-and-race-condition-actually-the-same-thing-in-context-of-conc/>. 2013.
- [7] George Candea Baris Kasikci Benjamin Schubert. *Gist*. <http://dslab.epfl.ch/proj/gist/>. 2015.
- [8] Rob von Behren et al. "Capriccio: Scalable threads for Internet services." In: *Symp. on Operating Systems Principles*. 2003.
- [9] Al Bessey et al. "A Few Billion Lines of Code Later: Using Static Analysis to Find Bugs in the Real World." In: *Commun. ACM* (2010).
- [10] Swarnendu Biswas, Minjia Zhang, and Michael D. Bond. *Lightweight Data Race Detection for Production Runs*. Tech. rep. OSU-CICRC-1/15-TR01. Ohio State University, 2015.
- [11] Robert L. Bocchino Jr. et al. "A type and effect system for deterministic parallel Java." In: *OOPSLA*. 2009.
- [12] Hans Boehm. *Programming with Threads: Questions Frequently Asked by C and C++ Programmers*. <http://www.hboehm.info/c++mm/user-faq.html>.
- [13] Hans-J. Boehm. "How to miscompile programs with "benign" data races." In: *USENIX Workshop on Hot Topics in Parallelism*. 2011.

- [14] Hans-J. Boehm and Sarita V. Adve. "Foundations of the C++ concurrency memory model." In: *Proceedings of the 2008 ACM SIGPLAN conference on Programming language design and implementation*. Intl. Conf. on Programming Language Design and Implem. 2008.
- [15] Michael D. Bond, Katherine E. Coons, and Kathryn S. McKinley. "PACER: Proportional detection of data races." In: *Intl. Conf. on Programming Language Design and Implem.* 2010.
- [16] Stefan Bucur et al. "Parallel Symbolic Execution for Automated Real-World Software Testing." In: *ACM EuroSys European Conf. on Computer Systems*. 2011.
- [17] Cristian Cadar, Daniel Dunbar, and Dawson R. Engler. "KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs." In: *Symp. on Operating Sys. Design and Implem.* 2008.
- [18] Subhachandra Chandra and Peter M. Chen. "Whither Generic Recovery from Application Faults? A Fault Study Using Open-Source Software." In: *DSN*. 2000.
- [19] V. Chipounov and G. Candea. "Enabling sophisticated analyses of x86 binaries with RevGen." In: *Intl. Conf. on Dependable Systems and Networks*. 2011.
- [20] Jong-Deok Choi and Andreas Zeller. "Isolating Failure-inducing Thread Schedules." In: *ISSTA*. 2002.
- [21] Jong-Deok Choi et al. "Efficient and precise datarace detection for multithreaded object-oriented programs." In: *SIGPLAN Notices* 37.5 (2002), pp. 258–269.
- [22] Chris Lattner. *libc++*. <http://libcxx.llvm.org/>. 2012.
- [23] CVE's related to races. <http://www.cvedetails.com/vulnerability-list/cweid-362/vulnerabilities.html>.
- [24] Michel Dubois, Christoph Scheurich, and Faye Briggs. "Memory Access Buffering in Multiprocessors." In: *Proc. 13th Ann. Intl. Symp. on Computer Architecture* (1986), pp. 374–442.
- [25] Tayfun Elmas, Shaz Qadeer, and Serdar Tasiran. "Goldilocks: A race and transaction-aware Java runtime." In: *Intl. Conf. on Programming Language Design and Implem.* San Diego, California, USA, 2007.
- [26] Dawson Engler and Ken Ashcraft. "RacerX: Effective, Static Detection of Race Conditions and Deadlocks." In: *Symp. on Operating Systems Principles*. 2003.
- [27] Peter Eriksson. *Parallel File Scanner*. <http://ostatic.com/pfscan>. 2013.
- [28] Brad Fitzpatrick. *Memcached*. <http://memcached.org>. 2013.

- [29] Cormac Flanagan and Stephen N. Freund. "FastTrack: Efficient and precise dynamic race detection." In: *Intl. Conf. on Programming Language Design and Implem.* 2009.
- [30] Cormac Flanagan and Stephen N. Freund. "Type-based Race Detection for Java." In: *PLDI '00*. 2000.
- [31] Cormac Flanagan, Stephen N. Freund, and Jaeheon Yi. "Velodrome: A sound and complete dynamic atomicity checker for multithreaded programs." In: *Intl. Conf. on Programming Language Design and Implem.* 2008.
- [32] Vijay Ganesh and David L. Dill. "A decision procedure for bit-vectors and arrays." In: *Intl. Conf. on Computer Aided Verification*. 2007.
- [33] Jeff Gilchrist. *Parallel BZIP2*. <http://compression.ca/pbzip2>. 2013.
- [34] Kirk Glerum et al. "Debugging in the (very) large: ten years of implementation and experience." In: *Symp. on Operating Systems Principles*. 2009.
- [35] Patrice Godefroid, Michael Y. Levin, and David Molnar. "Automated Whitebox Fuzz Testing." In: *Network and Distributed System Security Symp.* 2008.
- [36] Patrice Godefroid and Nachiappan Nagappan. "Concurrency at Microsoft – An Exploratory Survey." In: *Intl. Conf. on Computer Aided Verification*. 2008.
- [37] *Hacking Starbucks for unlimited coffee*. <http://sakurity.com/blog/2015/05/21/starbucks.html>.
- [38] *Helgrind*. <http://valgrind.org/docs/manual/hg-manual.html>. 2012.
- [39] IEEE. "1003.1 Standard for Information Technology Portable Operating System Interface (POSIX) Rationale (Informative)." In: *IEEE Std 1003.1-2001. Rationale (Informative)* (2001).
- [40] Intel. *Intel 64 and IA-32 Architectures Software Developer's Manual*. Vol. 2. 325383-038US. 2015.
- [41] Intel Corp. *Parallel Inspector*. <http://software.intel.com/en-us/articles/intel-parallel-inspector>. 2012.
- [42] Intel Corporation. *Intel Processor Trace*. <https://software.intel.com/en-us/blogs/2013/09/18/processor-tracing>. 2013.
- [43] *ISO/IEC 14882:2011: Information technology – Programming languages – C++*. International Organization for Standardization. 2011.
- [44] *ISO/IEC 9899:2011: Information technology – Programming languages – C*. International Organization for Standardization. 2011.

- [45] Ali Jannesari and Walter F. Tichy. "Identifying Ad-hoc Synchronization for Enhanced Race Detection." In: *Intl. Parallel and Distributed Processing Symp.* 2010.
- [46] Guoliang Jin et al. "Instrumentation and sampling strategies for cooperative concurrency bug isolation." In: *SIGPLAN Not.* (2010).
- [47] Sebastian Burckhardt John Erickson Madanlal Musuvathi and Kirk Olynyk. "Effective Data-Race Detection for the Kernel." In: *Symp. on Operating Sys. Design and Implem.* 2010.
- [48] John Regehr. *Race Condition vs. Data Race*. <http://blog.regehr.org/archives/490>. 2011.
- [49] Vineet Kahlon, Franjo Ivančić, and Aarti Gupta. "Reasoning About Threads Communicating via Locks." In: *Intl. Conf. on Computer Aided Verification*. 2005.
- [50] Vineet Kahlon et al. "Fast and Accurate Static Data-race Detection for Concurrent Programs." In: *CAV*. 2007.
- [51] Vineet Kahlon et al. "Static Data Race Detection for Concurrent Programs with Asynchronous Calls." In: *FSE*. 2009.
- [52] Baris Kasikci, Cristian Zamfir, and George Candea. "Automated Classification of Data Races Under Both Strong and Weak Memory Models." In: *ACM Transactions on Programming Languages and Systems* 37.3 (2015).
- [53] Baris Kasikci, Cristian Zamfir, and George Candea. "Data Races vs. Data Race Bugs: Telling the Difference with Portend." In: *Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*. 2012.
- [54] Baris Kasikci, Cristian Zamfir, and George Candea. "RaceMob: Crowdsourced Data Race Detection." In: *Symp. on Operating Systems Principles*. 2013.
- [55] Baris Kasikci et al. "Efficient Tracing of Cold Code Via Bias-Free Sampling." In: *USENIX Annual Technical Conf.* 2014.
- [56] Baris Kasikci et al. "Failure Sketching: A Technique for Automated Root Cause Diagnosis of In-Production Failures." In: *Symp. on Operating Systems Principles*. 2015.
- [57] Leslie Lamport. "Time, clocks, and the ordering of events in a distributed system." In: *Communications of the ACM* 21.7 (1978).
- [58] Chris Lattner. "Macroscopic Data Structure Analysis and Optimization." PhD thesis. University of Illinois at Urbana-Champaign, May 2005.
- [59] Chris Lattner and Vikram Adve. "LLVM: A Compilation Framework for Lifelong Program Analysis and Transformation." In: *Intl. Symp. on Code Generation and Optimization*. 2004.

- [60] Dongyoon Lee et al. "Chimera: Hybrid program analysis for determinism." In: *Intl. Conf. on Programming Language Design and Implem.* 2012.
- [61] Nancy G. Leveson and Clark S. Turner. "An Investigation of the Therac-25 Accidents." In: *IEEE Computer* (1993).
- [62] Ben Liblit et al. "Scalable Statistical Bug Isolation." In: *Intl. Conf. on Programming Language Design and Implem.* 2005.
- [63] Benjamin Robert Liblit. "Cooperative Bug Isolation." PhD thesis. University of California, Berkeley, Dec. 2004.
- [64] Shan Lu. "Understanding, Detecting and Exposing Concurrency Bugs." PhD thesis. UIUC, 2008.
- [65] Shan Lu et al. "AVIO: Detecting Atomicity Violations via Access Interleaving Invariants." In: *ASPLOS*. 2006.
- [66] Shan Lu et al. "Learning from Mistakes – A Comprehensive Study on Real World Concurrency Bug Characteristics." In: *Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*. 2008.
- [67] Chi-Keung Luk et al. "PIN: building customized program analysis tools with dynamic instrumentation." In: *Intl. Conf. on Programming Language Design and Implem.* 2005.
- [68] Nuno Machado, Brandon Lucia, and Luís Rodrigues. "Concurrency Debugging with Differential Schedule Projections." In: *PLDI*. 2015.
- [69] Jeremy Manson, William Pugh, and Sarita V. Adve. "The Java Memory Model." In: *POPL*. 2005.
- [70] Daniel Marino, Madanlal Musuvathi, and Satish Narayanasamy. "LiteRace: Effective sampling for lightweight data-race detection." In: *Intl. Conf. on Programming Language Design and Implem.* 2009.
- [71] Cal McPherson. *Ctrace*. <http://ctrace.sourceforge.net>. 2012.
- [72] *Memcached issue 127*. <http://code.google.com/p/memcached/issues/detail?id=127>.
- [73] Madanlal Musuvathi et al. "A Randomized Scheduler with Probabilistic Guarantees of Finding Bugs." In: *Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*. 2010.
- [74] Madanlal Musuvathi et al. "Finding and Reproducing Heisenbugs in Concurrent Programs." In: *Symp. on Operating Sys. Design and Implem.* 2008.
- [75] Mayur Naik, Alex Aiken, and John Whaley. "Effective static race detection for Java." In: *Intl. Conf. on Programming Language Design and Implem.* 2006.

- [76] Mayur Naik, Alex Aiken, and John Whaley. "Effective static race detection for Java." In: *Intl. Conf. on Programming Language Design and Implem.* 2006.
- [77] Satish Narayanasamy et al. "Automatically classifying benign and harmful data races using replay analysis." In: *Intl. Conf. on Programming Language Design and Implem.* (2007).
- [78] George C. Necula et al. "CIL: Intermediate Language and Tools for Analysis and Transformation of C Programs." In: *Intl. Conf. on Compiler Construction.* 2002.
- [79] Robert O'Callahan and Jong-Deok Choi. "Hybrid dynamic data race detection." In: *Symp. on Principles and Practice of Parallel Computing.* 2003.
- [80] Mark S. Papamarcos and Janak H. Patel. "A Low-overhead Coherence Solution for Multiprocessors with Private Cache Memories." In: *ISCA.* 1984.
- [81] Soyeon Park et al. "Do You Have to Reproduce the Bug at the First Replay Attempt? – PRES: Probabilistic Replay with Execution Sketching on Multiprocessors." In: *Symp. on Operating Systems Principles.* 2009.
- [82] Eli Pozniansky and Assaf Schuster. "Efficient on-the-fly data race detection in multithreaded C++ programs." In: *Symp. on Principles and Practice of Parallel Computing.* 2003.
- [83] Polyvios Pratikakis, Jeffrey S. Foster, and Michael Hicks. "LOCK-SMITH: context-sensitive correlation analysis for race detection." In: *Intl. Conf. on Programming Language Design and Implem.* 2006.
- [84] Feng Qin et al. "Rx: Treating bugs as allergies – a safe method to survive software failures." In: *ACM Transactions on Computer Systems* 25.3 (2007).
- [85] Quora. *What is a coder's worst nightmare?* <http://www.quora.com/What-is-a-coders-worst-nightmare>.
- [86] Caitlin Sadowski and Jaeheon Yi. "How Developers Use Data Race Detection Tools." In: *Proceedings of the 5th Workshop on Evaluation and Usability of Programming Languages and Tools.* PLATEAU. 2014.
- [87] Swarup Kumar Sahoo, John Criswell, and Vikram Adve. "An Empirical Study of Reported Bugs in Server Software with Implications for Automated Bug Diagnosis." In: *ICSE.* 2010.
- [88] Swarup Kumar Sahoo et al. "Using Likely Invariants for Automated Software Fault Localization." In: (2013).
- [89] Amit Sasturkar et al. "Automated Type-based Analysis of Data Races and Atomicity." In: *PPoPP.* 2005.

- [90] Stefan Savage et al. "Eraser: A dynamic data race detector for multithreaded programs." In: *ACM Transactions on Computer Systems* 15.4 (1997).
- [91] Edith Schonberg. "On-the-fly detection of access anomalies (with retrospective)." In: *SIGPLAN Notices* 39.4 (2004).
- [92] Koushik Sen. "Race directed random testing of concurrent programs." In: *Intl. Conf. on Programming Language Design and Implem.* (2008).
- [93] Koushik Sen, Darko Marinov, and Gul Agha. "CUTE: a concolic unit testing engine for C." In: *Symp. on the Foundations of Software Eng.* 2005.
- [94] Konstantin Serebryany and Timur Iskhodzhanov. "ThreadSanitizer - Data race detection in practice." In: *Workshop on Binary Instrumentation and Applications.* 2009.
- [95] Jaswinder Pal Singh, Wolf-Dietrich Weber, and Anoop Gupta. *SPLASH: Stanford Parallel Applications for Shared Memory*. Tech. rep. CSL-TR-92-526. Stanford University Computer Systems Laboratory, 1992.
- [96] Yannis Smaragdakis et al. "Sound Predictive Race Detection in Polynomial Time." In: (2012).
- [97] *SQLite*. <http://www.sqlite.org/>. 2013.
- [98] The Associated Press. *General Electric Acknowledges Northeastern Blackout Bug*. <http://www.securityfocus.com/news/8032>. Feb. 12, 2004.
- [99] William Thies, Michal Karczmarek, and Saman P. Amarasinghe. "StreamIt: A Language for Streaming Applications." In: *CC*. 2002.
- [100] Chen Tian et al. "Dynamic recognition of synchronization operations for improved data race detection." In: *Intl. Symp. on Software Testing and Analysis*. 2008.
- [101] *TIOBE Programming Community Index*. http://www.tiobe.com/tiobe_index/. Nov. 2004.
- [102] Joseph Tucek et al. "Triage: diagnosing production run failures at the user's site." In: *Symp. on Operating Systems Principles*. 2007.
- [103] Jan Wen Voun, Ranjit Jhala, and Sorin Lerner. "RELAY: Static race detection on millions of lines of code." In: *Symp. on the Foundations of Software Eng.* 2007.
- [104] Mark Weiser. "Program slicing." In: *Intl. Conf. on Software Engineering*. 1981.
- [105] Robert P. Wilson and Monica S. Lam. "Efficient context-sensitive pointer analysis for C programs." In: *Intl. Conf. on Programming Language Design and Implem.* La Jolla, CA, 1995.

- [106] *Windows Process and Thread Functions*. [https://msdn.microsoft.com/en-us/library/windows/desktop/ms684847\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/desktop/ms684847(v=vs.85).aspx).
- [107] Steven Cameron Woo et al. "The SPLASH-2 programs: characterization and methodological considerations." In: *Intl. Symp. on Computer Architecture* (1995).
- [108] Jingyue Wu, Heming Cui, and Junfeng Yang. "Bypassing races in live applications with execution filters." In: *Symp. on Operating Sys. Design and Implem.* 2010.
- [109] Xinwei Xie and Jingling Xue. "Acculock: Accurate and Efficient Detection of Data Races." In: *CGO*. 2011.
- [110] Weiwei Xiong et al. "Ad-Hoc Synchronization Considered Harmful." In: *Symp. on Operating Sys. Design and Implem.* 2010.
- [111] Min Xu, Rastislav Bodík, and Mark D. Hill. "A Serializability Violation Detector for Shared-memory Server Programs." In: *PLDI*. 2005.
- [112] Yu Yang et al. "Distributed dynamic partial order reduction based verification of threaded software." In: *Intl. SPIN Workshop*. 2007.
- [113] Jie Yu and Satish Narayanasamy. "A case for an interleaving constrained shared-memory multi-processor." In: *Intl. Symp. on Computer Architecture*. Austin, TX, USA, 2009.
- [114] Jie Yu and Satish Narayanasamy. "A Case for an Interleaving Constrained Shared-Memory Multi-Processor." In: *Intl. Symp. on Computer Architecture*. 2009.
- [115] Misun Yu, Sang-Kyung Yoo, and Doo-Hwan Bae. "SimpleLock: Fast and Accurate Hybrid Data Race Detector." In: *PDCAT*. 2013.
- [116] Andreas Zeller and Ralf Hildebrandt. "Simplifying and Isolating Failure-Inducing Input." In: *IEEE Transactions on Software Engineering* (2002).
- [117] Wei Zhang et al. "ConSeq: Detecting Concurrency Bugs through Sequential Errors." In: *Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*. 2011.
- [118] Pin Zhou, Radu Teodorescu, and Yuanyuan Zhou. "HARD: Hardware Assisted Lockset-based Race Detection." In: *International Symposium on High-Performance Computer Architecture*. 2007.