# Engineering Record And Replay For Deployability
## Extended Technical Report

Robert O'Callahan[*]          Chris Jones[*]          Nathan Froyd          Kyle Huey[*]
*Mozilla Corporation*

Albert Noll[*]          Nimrod Partush[*]
*Swisscom AG*          *Technion*

## Abstract

The ability to record and replay program executions with low overhead enables many applications, such as reverse-execution debugging, debugging of hard-to-reproduce test failures, and "black box" forensic analysis of failures in deployed systems. Existing record-and-replay approaches limit deployability by recording an entire virtual machine (heavyweight), modifying the OS kernel (adding deployment and maintenance costs), requiring pervasive code instrumentation (imposing significant performance and complexity overhead), or modifying compilers and runtime systems (limiting generality). We investigated whether it is possible to build a practical record-and-replay system avoiding all these issues. The answer turns out to be yes — if the CPU and operating system meet certain non-obvious constraints. Fortunately modern Intel CPUs, Linux kernels and user-space frameworks do meet these constraints, although this has only become true recently. With some novel optimizations, our system RR records and replays real-world low-parallelism workloads with low overhead, with an entirely user-space implementation, using stock hardware, compilers, runtimes and operating systems. RR forms the basis of an open-source reverse-execution debugger seeing significant use in practice. We present the design and implementation of RR, describe its performance on a variety of workloads, and identify constraints on hardware and operating system design required to support our approach.

## 1 Introduction

The ability to record a program execution with low overhead and play it back precisely has many applications [17, 18, 22] and has received significant attention in the research community. It has even been implemented in products such as VMware Workstation [31], Simics [23],

UndoDB [3] and TotalView [25]. Unfortunately, deployment of these techniques has been limited, for various reasons. Some approaches [20, 23, 31] require recording and replaying an entire virtual machine, which is heavyweight — users must set up and manage the virtual machine, and the state that must be recorded and replayed encompasses an entire operating system, when the user often only cares about some specific process(es). This is especially problematic in record-and-replay applications such as reverse-execution debugging, where checkpoints of the replayed state are frequently created and resumed. Other approaches [9, 17, 29, 32] require running a modified OS kernel, hindering deployment and adding security and stability risk to the system. Requiring compiler and language runtime changes [32] also hinders deployment, especially when applications include their own JIT compilers. Some approaches [27, 33, 37] require custom hardware not yet available. Many approaches [3, 10, 25, 36] require pervasive instrumentation of code, which adds complexity and overhead, especially for self-modifying code (commonly used in polymorphic inline caching [28] and other implementation techniques in modern just-in-time compilers). A performant dynamic code instrumentation engine is also expensive to build and maintain.

Mozilla developers observed that developers spend a lot of time debugging, particularly on deterministic bugs. Record-and-replay debugging seemed like a promising technology but none of the existing systems were easy to deploy, for the above reasons. We initiated the RR project to try to build a system suitable for Mozilla's use — or else discover reasons why such a system could not be built (with moderate implementation effort).

Therefore we set out to build a system that maximizes deployability: to record and replay unmodified user-space applications with stock Linux kernels, compilers, language runtimes, and x86/x86-64 CPUs, with a fully user-space implementation running without special privileges, and without using pervasive code instrumen-

---

tation. Given our limited development resources (a couple of person-years for a prototype able to debug Firefox, perhaps five person-years total to date), we avoid approaches that demand huge engineering effort. We assume RR should run unmodified applications, and they will have bugs (including data races) that we wish to faithfully record and replay, but these applications will not maliciously try to subvert recording or replay. We combine techniques already known, but not previously demonstrated working together in a practical system: primarily, using `ptrace` to record and replay system call results and signals, avoiding non-deterministic data races by running only one thread at a time, and using CPU hardware performance counters to measure application progress so asynchronous signal and context-switch events are delivered at the right moment [35]. Section 2 describes our approach in more detail.

With that in place, we discovered the main performance bottleneck for low-parallelism workloads (in particular, Firefox running the Firefox test suite) was context switching induced by using `ptrace` to monitor system calls. We implemented a novel *in-process system-call interception* technique to eliminate those context switches, dramatically reducing recording and replay overhead on important real-world workloads. This optimization relies on modern Linux kernel features: `seccomp-bpf` to selectively suppress `ptrace` traps for certain system calls, and `perf` context-switch events to detect recorded threads blocking in the kernel. Section 3 describes this work, and Section 4 gives some performance results, showing that on important application workloads RR recording and replay slowdown is less than a factor of two.

We rely on hardware and OS features designed for other goals, so it is surprising that RR works. In fact, it skirts the edge of feasibility; in particular it cannot be implemented on ARM CPUs. Section 5 summarizes RR's hardware and software requirements, which we hope will influence system designers.

RR is in daily use by many developers as the foundation of an efficient reverse-execution debugger that works on complex applications such as Samba, Firefox, Chromium, QEMU, LibreOffice and Wine. It is free software, available at `https://github.com/mozilla/rr`. This paper makes the following research contributions:

- We show that record and replay of user-space processes on modern, stock hardware and software without pervasive code instrumentation is possible and practical.

- We introduce an *in-process system-call interception* technique and show it dramatically reduces overhead.

- We show that for low-parallelism workloads, RR recording and replay overhead is reasonably low, lower than other approaches with comparable deployability.

- We identify hardware and operating system design constraints required to support our approach.

This extended technical report elaborates on published papers with additional technical details and, in Section 6, a reflection on the design, usage and evolution of RR.

## 2 Design

### 2.1 Summary

Most low-overhead record-and-replay systems depend on the observation that CPUs are mostly deterministic. We identify a boundary around state and computation, record all sources of nondeterminism within the boundary and all inputs crossing into the boundary, and reexecute the computation within the boundary by replaying the nondeterminism and inputs. If all inputs and nondeterminism have truly been captured, the state and computation within the boundary during replay will match that during recording.

To enable record and replay of arbitrary Linux applications, without requiring kernel modifications or a virtual machine, RR records and replays the user-space execution of a group of processes. To simplify invariants, and to make replay as faithful as possible, replay preserves almost every detail of user-space execution. In particular, user-space memory and register values are preserved exactly, with a few exceptions noted later in the paper. This implies CPU-level control flow is identical between recording and replay, as is memory layout.

While replay preserves user-space state and execution, only a minimal amount of kernel state is reproduced during replay. For example, file descriptors are not opened, signal handlers are not installed, and filesystem operations are not performed. Instead the recorded user-space-visible effects of those operations, and future related operations, are replayed. We do create one replay thread per recorded thread (not strictly necessary), and we create one replay address space (i.e. process) per recorded address space, along with matching memory mappings.

With this design, our recording boundary is the interface between user-space and the kernel. The inputs and sources of nondeterminism are mainly the results of system calls, and the timing of asynchronous events.

### 2.2 Avoiding Data Races

With threads running on multiple cores, racing read-write or write-write accesses to the same memory lo-

cation by different threads would be a source of non-determinism. Therefore we take the common approach [20, 31, 3, 18] running only one thread at a time. RR preemptively schedules these threads, so context switch timing is nondeterminism that must be recorded. Data race bugs can still be observed if a context switch occurs at the right point in the execution (though bugs due to weak memory models cannot be observed).

This approach is much simpler and more deployable than alternatives [10, 21, 36, 41, 32], avoids assuming programs are race-free [17, 32], and is efficient for low-parallelism workloads. There is a large slowdown for workloads with a consistently high degree of parallelism; however, even for applications which are potentially highly parallel, users often apply RR to test workloads with relatively small datasets and hence limited parallelism.

The one-thread-at-a-time restriction is implemented by supervising all tracee processes/threads via the `ptrace` system call. Whenever a thread enters a system call, RR receives a "ptrace stop" event from the kernel and decides whether to allow a context switch or continue running the current thread. (At any potentially blocking system call, we must allow a context switch.) When a context switch is allowed, RR uses PTRACE_SYSCALL to let the thread continue running to system call exit, then selects another thread to run and resumes it with a PTRACE_CONT (or similar) `ptrace` operation. When the new thread in turn reaches a reschedule point, RR queries the state of all threads using `waitpid` to determine which, if any, are ready to exit blocking system calls and are therefore candidates for scheduling.

RR performs preemptive context switching by periodically interrupting the running thread with an asynchronous signal (see Section 2.4). The normal scheduler honors Linux thread priorities strictly: higher-priority threads always run in preference to lower-priority threads. Threads with equal priority run in a round-robin manner.

## 2.3 System Calls

System calls return data to user-space by modifying registers and memory, and these changes must be recorded. The `ptrace` system call allows a process to be synchronously notified when a tracee thread enters or exits a system call. When a tracee thread enters the kernel for a system call, it is suspended and RR is notified. When RR chooses to run that thread again, the system call will complete, notifying RR again, giving it a chance to record the system call results. RR contains a model of most Linux system calls describing the user-space memory they can modify, given the system call input parame-
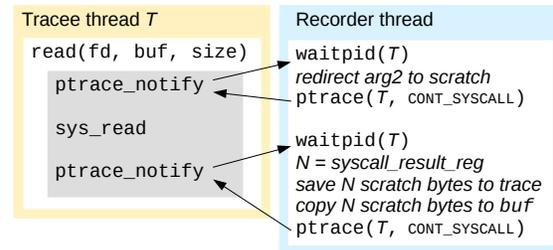


Figure 1: Recording simple system call

ters and result.

### 2.3.1 Scratch Buffers

As noted above, RR normally avoids races by scheduling only one thread at a time. However, if a system call blocks in the kernel, RR must try to schedule another application thread to run while the blocking system call completes. It's possible (albeit unlikely) that the running thread could access the system call's output buffer and race with the kernel's writes to that buffer. To avoid this, we redirect system call output buffers to per-thread temporary "scratch memory" which is otherwise unused by the application. When we get a `ptrace` event for a blocked system call completing, RR copies scratch buffer contents to the real user-space destination(s) while no other threads are running, eliminating the race. (We actually have no evidence that the races prevented by scratch buffers occur in practice, and it might be worth trying to eliminate scratch buffers.)

Figure 1 illustrates recording a simple `read` system call. The gray box represents kernel code.

### 2.3.2 `ptrace` **Emulation**

Apart from redirection through scratch buffers and recording memory effects, most system calls execute normally during recording. One important exception is `ptrace`. `ptrace` is not only used by debuggers but also by some sophisticated applications, e.g. to collect data for crash reporting in Firefox, or to provide emulation facilities in Wine [6]. Linux only allows a thread to have a single `ptrace` supervisor, so if one tracee tried to `ptrace` another directly, that would fail because the latter already has RR as a `ptrace` supervisor. Instead RR emulates all tracee `ptrace` operations and relationships. `ptrace` semantics are complicated, so this `ptrace` emulation is also necessarily rather complicated to implement.

### 2.3.3 Performing System Calls In Tracee Context

In a variety of situations (e.g. `mmap` handling, see below) RR must perform arbitrary system calls in the context of a tracee thread. RR locates a system-call instruction in the tracee's address space and uses `ptrace` to temporarily alter the tracee execution context to execute that instruction with the desired parameters in registers. Sometimes memory parameters are also required, in which case readable memory is located on the stack (or *in extremis* anywhere in the address space), temporarily overwritten with the desired parameters, and then reset to its original contents after the system call has completed.

### 2.3.4 Kernel Namespaces

Linux kernel namespaces are a powerful feature used to implement "containers", and are also used to implement sandboxing in sophisticated applications such as Firefox and Chromium. They mostly require no special handling, except that it is sometimes necessary for a sandboxed process to communicate with RR, e.g. by opening shared files or sockets. To make this possible, we open a hole in the sandbox by allocating a reserved file descriptor in all tracees that refers to the root directory as seen by RR; with this descriptor, a tracee can open temporary files or sockets created by RR even if its mount namespace or root directory have been changed.

A malicious but sandboxed tracee aware of this hole could exploit it to escape from the sandbox, which is one reason we assume tracees are not malicious.

### 2.3.5 Seccomp

Another Linux sandboxing feature is `seccomp` filtering. This installs a `seccomp-bpf` filter into the kernel, bytecode which evaluates the register parameters of each system call and determines whether the system call can proceed normally or should produce an error, signal or `ptrace` trap. Sandboxes typically use this feature to whitelist permitted system calls and trigger emulation of others. This sandboxing can disrupt RR by interfering with the system calls RR performs in the tracee context.

Immediately after each `execve` we map a page of memory at a fixed address, the "RR page", containing a "privileged" system-call instruction. We detect the system calls that install `seccomp-bpf` filters, and patch the filter bytecode with a prologue exiting early with "allow" if the program counter is at the "privileged" instruction address. We use that "privileged" instruction when performing RR system calls in tracee context.

We also record the effects of `seccomp-bpf` filters on regular system calls and replay them faithfully. We do this just by observing kernel behavior, without having to actually implement a BPF interpreter in RR.

### 2.3.6 System Call Modeling In Practice

RR assumes it is possible to efficiently determine from user-space state before and after a system call (plus, possibly, observations of previous system calls) reasonable bounds on the memory locations modified by the system call. For example, we would not be able to handle a system call that writes to random locations in memory without somehow indicating which locations were modified. This assumption holds adequately in practice, but it is a lot of work to create and maintain the model that describes how to determine the memory modified by each system call.

The biggest problem is the `ioctl` system call, which provides per-device (therefore extensible) system call namespaces. There is a convention for encoding the address and length of written memory in the ioctl parameters, but it is often not followed. RR builds in knowledge of many ioctls, but this is necessarily incomplete and needs to be extended from time to time to handle new applications.

This issue is not unique to RR — for example, Valgrind Memcheck also needs this information. It would be useful to have a shared repository of system call descriptions from which our models could be generated.

An alternative approach to modeling all system calls would be to use the Linux "soft-dirty" tracking API [4] to identify pages modified by infrequently used or unknown system calls. We haven't tried this because although it might ultimately reduce work, it adds complexity and would not be compatible with applications that use soft-dirty for their own needs. Also, requiring RR to be aware of the semantics of all system calls means that we detect new system calls that require special handling by RR.

### 2.3.7 Replay

During replay, when the next event to be replayed is an intercepted system call, we set a temporary breakpoint at the address of the system call instruction (recorded in the trace). We use `ptrace` to run the tracee thread until it hits the breakpoint, remove the breakpoint, advance the program counter past the system call instruction, and apply the recorded register and memory changes. This approach is simple and more efficient than using `PTRACE_SYSEMU` to enter and suppress execution of system calls, because it reduces the number of `ptrace`-stop notifications from two (on entry to and exit from the system call) to one. (These notifications are relatively expensive because each one requires context switching from the tracee to RR and back.)

In some cases this technique may not work reliably. For example, the tracee may write a system-call instruction to memory and then jump to it, in which case a soft-

ware breakpoint set at that location would be overwritten. We could use hardware breakpoints but that could conflict with a tracee's own use of hardware breakpoints, since x86 CPUs support a limited number of them (typically four). Therefore if the system call instruction to replay is in shared memory or writeable memory, we use PTRACE_SYSEMU to advance to the execution of the next system call and incur the cost of two ptrace-stop notifications.

### 2.3.8 Replaying Complex System Calls

Some system calls manipulate threads or address spaces and require special handling during replay. For example a recorded file mmap saves the original contents of the mapped file to the trace (often in an optimized zero-copy way; see Section 2.7). Replay maps that copy with the MAP_FIXED option to ensure the mapping is created at the correct address.

The execve system call is complicated to record and replay. After a recorded execve has completed, RR reads /proc/.../maps to discover the memory mappings in the new address space and record them. In replay, RR execve's a small stub executable for the correct architecture (32-bit or 64-bit), removes all the memory mappings in the address space (since they may be completely different from those created during recording, e.g. due to ASLR), then creates new mappings to match those recorded.

### 2.3.9 Signal Handlers

When a synchronous or asynchronous signal triggers a signal handler during recording, the kernel sets up a user-space stack frame for the handler and transfers control to the address of the handler. ptrace reports the delivery of a signal as via a ptrace-stop notification, then RR triggers the transition to the signal handler by issuing a PTRACE_SINGLESTEP request. (RR must track internally which signals have handlers so it knows in advance which signals will trigger a handler. It also needs to be aware of quirky edge cases that suppress entry into a signal handler; for example, when code triggerings a SIGILL or SIGSEGV, but the tracee has used sigprocmask to block the signal, the signal is delivered but handlers do not run (usually with fatal results).)

When RR has transitioned into the signal handler, it records the contents of the stack frame and registers set up by the kernel. During replay, no signal handlers are set up and no real signals are delivered. To replay the entry to a signal handler, we just write the recorded signal handler frame into memory and set the registers to the recorded values. By not delivering real signals during replay, we avoid having to set up signal-handler state in the replaying process.

### 2.3.10 Signal Handling And Interrupted System Calls

During normal non-recorded execution, when a system call is interrupted by a signal, the kernel adjusts the user-space registers with the program counter immediately after the system-call instruction and a special -ERESTARTSYS (or similar) value in the system-call result register. If a signal handler needs to run, the signal frame is built and the signal handler runs normally. When execution returns to the pre-handler state via sigreturn, or a signal handler did not run but the thread is resumed by some other mechanism, then just before returning to user-space the kernel's system-call restart mechanism kicks in: when the kernel returns to userspace, if the system-call result register contains -ERESTARTSYS and certain other conditions are met, the kernel automatically backs up the program counter to before the system call instruction and restores the system-call-number register. This causes the system call to reexecute.

RR avoids interfering with this restart mechanism and relies on the kernel restarting interrupted system calls as it normally would. However, because we set up state on system-call entry that expects to be processed on system-call exit, we have to detect when a ptrace system-call notification corresponds to a restarted system call instead of a new system call. There is no foolproof way to do this, but we have a heuristic that works well. We maintain a stack of pending interrupted system calls and their saved parameter registers. On every system call entry, if the parameter registers match the parameters of the most recent pending interrupted system call, we assume that system call is being restarted and pop it off the stack. If a signal handler interrupted a system call and then returns, but the interrupted system call is not immediately restarted, we assume it has been abandoned and pop it off the stack.

## 2.4 Asynchronous Events

We need to support two kinds of asynchronous events: preemptive context switches and signals. We treat the former as a special case of the latter, forcing a context switch by sending a signal to a running tracee thread. We need to ensure that during replay, a signal is delivered when the program is in exactly the same state as it was when the signal was delivered during recording.

As in previous work [20, 35, 13] we measure application progress using CPU hardware performance counters. Ideally we would count retired instructions leading up to an asynchronous event during recording, and dur-

ing replay program the CPU to fire an interrupt after that many instructions have been retired — but this approach needs modifications to work in practice.

### 2.4.1  Nondeterministic Performance Counters

We require that every execution of a given sequence of user-space instructions changes the counter value by an amount that depends only on the instruction sequence, not system state invisible to user space (e.g. the contents of caches, the state of page tables, or speculative CPU state). (Such effects manifest as noise making the counter unreliable from RR's point of view.) This property (commonly described as "determinism" [42]) does not hold for most CPU performance counters in practice [20, 42]. For example it does not hold for any "instructions retired" counter on any known x86 CPU model (e.g. because an instruction triggering a page fault is restarted and counted twice).

Fortunately, modern Intel CPUs have exactly one deterministic performance counter: "retired conditional branches" ("RCB"), so we use that. We cannot just count the number of RCBs during recording and deliver the signal after we have executed that number of RCBs during replay, because the RCB count does not uniquely determine the execution point to deliver the signal at. Therefore we pair the RCB count with the complete state of general-purpose registers (including the program counter) to identify an execution point.

In general that *still* does not uniquely identify an execution point (e.g. consider the infinite loop `label: inc [global_var]; jmp label;`). However, in practice we have found it works reliably; code that returns to the same instruction with no intervening conditional branch must be very rare, and it only matters to RR for regular replay if an asynchronous event occurs at such an instruction — in which case replay would probably diverge and fail.

When RR is used for more than just plain record-and-replay, in particular to implement reverse-execution debugging, it becomes considerably more sensitive to situations where the general-purpose registers and the RCB count do not uniquely identify a program state, because our reverse-execution implementation sometimes performs very frequent stops (e.g. while single-stepping) and many comparisons of one stopping point to another. We encountered a few practical examples where distinct program states with the same general-purpose registers and RCB count would cause reverse-execution debugging to fail. These were cases where a function A contained multiple calls to the same function B, no conditional branches occurred in A or B, and due to optimization, at some instruction in B general-purpose registers ended up with the same values in both invocations. We

addressed this issue by collecting some data from the stack at each stopping point and adding it to the program-state identifier. This now seems to be robust; we have not encountered any new issues in this area for a couple of years.

### 2.4.2  Alternative Approaches to Nondeterministic Counters

Other work [20, 13] avoids those uniqueness issues by counting the number of branches retired; an instruction at a given address can't be reached more than once without an intervening branch (except for repeating string instructions, which require the `ECX`/`RCX` register to be added to the program-state identifier). Unfortunately the branch counter is nondeterministic; for example a transition out a hardware interrupt is counted against the user process, but invisible to RR. Even worse, transitions out of "Systems Management Mode" interrupts are counted against the user process and aren't even visible to the host kernel. Working around these issues at user-level seems harder than the approach RR has been using.

Another interesting issue is that cloud virtualization services have been reluctant to enable virtualization of hardware performance counters, partly because counter nondeterminism due to low-level CPU state often represents information leakage from one guest VM to another. By definition, a deterministic counter is unaffected by information not already visible in its process, and we hope therefore more likely to be safely virtualizable. (At time of writing, most cloud providers do not virtualize any performance counters and thus RR does not work on them, except for Digital Ocean, where RR does work.)

The other obvious approach to measuring progress is to eschew hardware performance counters and use binary instrumentation to count instructions instead [36, 10, 3]. We believe the overhead and implementation complexity of binary instrumentation still make it less desirable than using hardware counters, for our purposes.

### 2.4.3  Late Interrupt Firing

Another problem with hardware performance counters is that, although CPUs can be programmed to fire an interrupt after a specified number of performance events have been observed, the interrupt does not fire immediately. In practice we often observe it firing after dozens more instructions have retired. To compensate for this, during replay, we program the interrupt to trigger some number of events earlier than the actual RCB count we are expecting. Then we set a temporary breakpoint at the program counter value for the state we're trying to reach, and repeatedly run to the breakpoint until the RCB count and the general-purpose register values match their recorded

values.

## 2.5 Shared Memory

By scheduling only one thread at a time, RR avoids issues with races on shared memory as long as that memory is written only by tracee threads. It is possible for recorded processes to share memory with other processes, and even kernel device drivers, where that non-recorded code can perform writes that race with accesses by tracee threads. Fortunately, this is rare for applications running in common Linux desktop environments, occurring in only four common cases: applications sharing memory with the PulseAudio daemon, applications sharing memory with the X server, applications sharing memory with kernel graphics drivers and GPUs, and vdso syscalls. We avoid the first three problems by automatically disabling use of shared memory with PulseAudio and X (falling back to a socket transport in both cases), and disabling direct access to the GPU from applications.

vdso syscalls are a Linux optimization that implements some common read-only system calls (e.g. gettimeofday) entirely in user space, partly by reading memory shared with the kernel and updated asynchronously by the kernel. We disable vdso syscalls by patching their user-space implementations to perform the equivalent real system call instead.

Applications could still share memory with non-recorded processes in problematic ways, though this is rare in practice and can often be solved just by enlarging the scope of the group of processes recorded by RR.

## 2.6 Nondeterministic Instructions

Almost all CPU instructions are deterministic, but some are not. One common nondeterministic x86 instruction is RDTSC, which reads a time-stamp counter. This particular instruction is easy to handle, since the CPU can be configured to trap on an RDTSC and Linux exposes this via a prctl API, so we can trap, emulate and record each RDTSC.

RDRAND generates random numbers and hopefully is not deterministic. We have only encountered it being used in one place in GNU libstdc++, so RR patches that explicitly.

The CPUID instruction is mostly deterministic, but one of its features returns the index of the running core, which affects behavior deep in glibc and can change as the kernel migrates a process between cores. We use the Linux sched_setaffinity API to force all tracee threads to run on a particular fixed core, and also force them to run on that core during replay.

We could easily avoid most of these issues in well-behaved programs if we could just trap-and-emulate the

CPUID instruction from user-space, since then we could mask off the feature bits indicating support for RDRAND, hardware transactions, etc, as well as record-and-replay the exact values returned by the recording CPU. Modern Intel CPUs support this ("CPUID faulting"); we are in the process of adding a user-accessible API for this to Linux. (In fact it just landed in Linux 4.12.)

### 2.6.1 Transactional Instructions

Modern Intel CPUs support two forms of transactional extensions. "Restricted Transactional Memory" ("RTM") exposes transactions as an architectural feature, using the XBEGIN and XEND instructions to explicitly start and end transactions. "Hardware Lock Elision" ("HLE") supports implicit transactions by adding the XACQUIRE and XRELEASE prefixes to the instructions that acquire and release a lock. The latter is a pure backwards-compatible performance hint that does not change visible semantics — failed transactions are automatically retried non-speculatively. The former has visible side effects due to explicitly exposing transaction failures.

RTM is nondeterministic from the point of view of user-space, since a hardware transaction can succeed or fail depending on CPU cache state (and probably the occurrence of hardware interrupts). Fortunately so far we have only found these being used by the system pthreads library, and we dynamically apply custom patches to that library to disable use of hardware transactions. In the future we will use mask off the RTM feature bit in CPUID to stop applications using RTM.

HLE does not have user-space-visible side effects, but the regular RCB counter counts conditional branches retired in a failed HLE transaction, making that counter nondeterministic when HLE is used. We can't disable HLE via CPUID, because HLE being backwards-compatible means code is free to use the HLE prefixes without checking CPUID first. Intel provides an option (IN_TXCP) to have hardware performance counters only count events in committed transactions, so we can use that to make the RCB counter deterministic again. Unfortunately on Linux an IN_TXCP counter can't be configured to interrupt after a certain number of events have been observed, because this could cause infinite loops in some cases. (The scenario is: counter $C$ is configured to interrupt after $N$ events, a transaction starts and executes $N$ events, the interrupt fires and causes the transaction to be rolled back, restoring the counter to its original value, the kernel resumes the process, repeat...) Our solution is to use a non-IN_TXCP counter to trigger interrupts and a IN_TXCP counter to count elapsed events. This means an interrupt could fire early, which is not a problem since we already program it to fire early.

## 2.7 Reducing Trace Sizes

For many applications the bulk of their input is memory-mapped files, mainly executable code. Copying all executables and libraries to the recorded trace on every execution would impose significant time and space overhead. RR creates hard links to memory-mapped executable files instead of copying them; as long as a system update or recompile replaces executables with new files, instead of writing to the existing files, the links retain the old file data. This works well in practice.

Even better, modern filesystems such as XFS and Btrfs offer copy-on-write logical copies of files (and even block ranges within files), ideal for our purposes. When a mapped file is on the same filesystem as the recorded trace, and the filesystem supports cloning, RR clones mapped files into the trace. These clone operations are essentially free in time and space, until/unless the original file is modified or deleted.

RR compresses all trace data, other than cloned files and blocks, with the `zlib` "deflate" method. Significantly better compression algorithms exist, and at some point we should revisit this choice.

Nevertheless, with these optimizations, in practice trace storage is a non-issue. Section 4.4 presents some results.

## 3 In-process System-call Interception

The approach described in the previous section works, but overhead is disappointingly high (see Figure 5 below). The core problem is that for every tracee system call, as shown in Figure 1 the tracee performs four context switches: two blocking `ptrace` notifications, each requiring a context switch from the tracee to RR and back. For common system calls such as `gettimeofday` or `read` from cached files, the cost of even a single context switch dwarfs the cost of the system call itself. To significantly reduce overhead, we must avoid context-switches to RR when processing these common system calls.

Therefore, we inject into the recorded process a library that intercepts common system calls, performs the system call without triggering a `ptrace` trap, and records the results to a dedicated buffer shared with RR. RR flushes the buffer to its trace every time it receives a synchronous stop notification. The concept is simple but there are problems to overcome.

### 3.1 Intercepting System Calls

A common technique for intercepting system calls in-process is to use dynamic linking to interpose wrapper functions over the C library functions that make system calls. In practice, we have found that method to be insufficient, due to applications making direct system calls, and fragile, due to variations in C libraries, and applications that require their own preloading [39, 6]).

Instead, when the tracee makes a system call, RR is notified via a `ptrace` trap and it tries to rewrite the system-call instruction to call into our interception library. This is tricky because on x86 a system call instruction is two bytes long, but we need to replace it with a five-byte `call` instruction. (On x86-64, to ensure we can call from anywhere in the address space to the interception library, we also need to also allocate trampolines within 2GB of the patched code.) In practice, frequently executed system call instructions are followed by a few known, fixed instruction sequences; for example, many system call instructions are followed by a `cmpl $0xfffff001,%eax` instruction testing the syscall result. We added five hand-written stubs to our interception library that execute post-system-call instructions before returning to the patched code. On receipt of a `ptrace` system-call notification, RR replaces the system call instruction and its following instruction with a call to the corresponding stub.

We (try to) redirect all system call instructions to the interception library, but for simplicity it only contains wrappers for the most common system calls, and for others it falls back to doing a regular `ptrace`-trapping system call.

### 3.2 Selectively Trapping System Calls

`ptrace` system-call monitoring triggers traps for all system calls, but our interception library needs to avoid traps for selected system calls. Fortunately, modern Linux kernels support selectively generating `ptrace` traps: `seccomp-bpf`. `seccomp-bpf` was designed primarily for sandboxing. A process can apply a `seccomp-bpf` filter function, expressed in bytecode, to another process; then, for every system call performed by the target process, the kernel runs the filter, passing in incoming user-space register values, including the program counter. The filter's result directs the kernel to either allow the system call, fail with a given `errno`, kill the target process, or trigger a `ptrace` trap. Overhead of filter execution is negligible since filters run directly in the kernel and are compiled to native code on most architectures.

Figure 2 illustrates recording a simple `read` system call with in-process system-call interception. The solid-border box represents code in the interception library and the grey box represents kernel code.

As mentioned above, RR injects a "RR page" into every tracee process at a fixed address. That page contains a special system call instruction — the "untraced
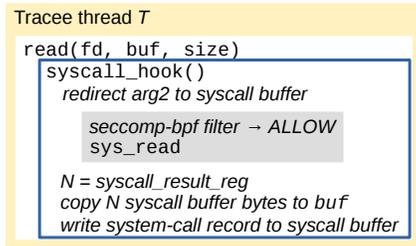
Figure 2: Recording with system-call interception

instruction". RR applies a `seccomp-bpf` filter to each recorded process that triggers a `ptrace` trap for every system call — except when the program counter is at the untraced instruction, in which case the call is allowed. Whenever the interception library needs to make an untraced system call, it uses that instruction.

## 3.3 Detecting Blocked System Calls

Some common system calls sometimes block (e.g. `read` on an empty pipe). Because RR runs tracee threads one at a time, if a thread enters an untraced blocking system call without notifying RR, it will hang and could cause the entire recording to deadlock (e.g. if another tracee thread is about to `write` to the pipe). We need the kernel to notify RR and suspend the tracee thread whenever an untraced system call blocks, to ensure we can schedule a different tracee thread.

We do this using the Linux `perf` event system to monitor `PERF_COUNT_SW_CONTEXT_SWITCHES`. The kernel raises one of these events every time it deschedules a thread from a CPU core. The interception library monitors these events for each thread and requests that the kernel send a signal to the blocked thread every time the event occurs. These signals trigger `ptrace` notifications to RR while preventing the thread from executing further. To avoid spurious signals (e.g. when the thread is descheduled due to normal timeslice expiration), the event is normally disabled and explicitly enabled during an untraced system call that might block. Still, spurious `SWITCHES` can occur at any point between enabling and disabling the event; we handle these edge cases with careful inspection of the tracee state.

Figure 3 illustrates recording a blocking `read` system call with system-call interception. The kernel deschedules the thread, triggering a `perf` event which sends a signal to the thread, *rescheduling* it, interrupting the system call, and sending a `ptrace` notification to the recorder. The recorder does bookkeeping to note that an intercepted system call was interrupted in thread *T*, then checks whether any tracee threads in blocking system calls have progressed to a system-call exit and gen-
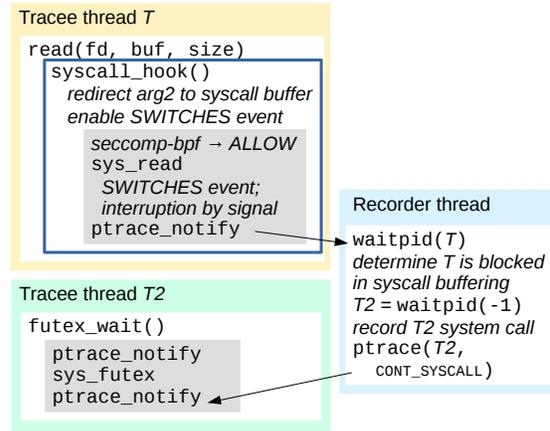


Figure 3: Recording a blocking system call

erated a `ptrace` notification. In this example *T2* has completed a (not intercepted) blocking `futex` system call, so we resume executing *T2*.

Meanwhile the `read` system call in *T* is restarted and treated much like a regular non-intercepted system call. Its parameters are rewritten to use scratch buffers. *T* is resumed using `PTRACE_SYSCALL` so that the thread will stop when the system call exits. After the system call exits, *T* will eventually be rescheduled. At that point the system-call exit will be recorded. Because the system call was recorded normally, it should not also be recorded in the interception library's trace buffer. Therefore RR sets an "abort commit" flag in the interception library's memory before resuming.

## 3.4 Signal Handling

If interception code was interrupted by a context switch or asynchronous signal partway through recording the result of a system call, that would be disastrous. When `ptrace` reports a signal arriving at a thread in interception code, RR stashes the signal away and sets a flag to indicate that the interception code must perform a traced "dummy" system call just before returning to application code, then resumes the tracee without delivering the signal. When the interception code performs a traced system call (the "dummy" system call, or some other system call for which we don't have a fast self-recording path), RR arranges to drain stashed signals before the traced system call is allowed to proceed (by setting a breakpoint at the traced system call entry point while there are stashed signals).

A difficult edge case involves untraced `sigprocmask` system calls. The kernel may try to deliver a signal to a thread in the interception library just before `sigprocmask` blocks the signal. If RR

9

were to stash the signal and try to inject it after the `sigprocmask`, the signal might not be deliverable. Therefore, while there are stashed signals, RR sets a breakpoint on the untraced system call instruction. If a signal arrived before an untraced `sigprocmask`, the breakpoint is hit and RR emulates the `sigprocmask` returning `EAGAIN`, causing the interception code to automatically retry it using a traced system call, which gives RR an opportunity to inject all stashed signals (as noted above) and always succeeds.

Another difficult edge case involves multiple signals. Normal signal delivery can automatically block future signals, sometimes permanently (e.g. using the `sigaction sa_mask` feature). For example suppose signals $S_1$ and then $S_2$ are sent to a process with two threads $T_1$ and $T_2$, and the handler for $S_1$ is configured to block $S_2$ and just exits the thread. Then suppose under RR the kernel delivers $S_1$ to $T_1$ while $T_1$ is in interception code; RR stashes $S_1$ and resumes $T_1$. If the kernel were to then deliver $S_2$ to $T_1$ we would have a problem, because $S_2$ will never be handled by any thread, whereas during non-recorded execution the kernel would have delivered it to $T_2$ instead. There is not enough information available from user-space to determine where the kernel would have delivered $S_2$ if it had been blocked in $T_1$ (in particular we can't tell the difference between a signal dispatched to a specific thread and a signal dispatched to a process).

Therefore it is important to ensure the kernel does not deliver a signal to a recorded thread unless that thread will definitely be able to handle the signal. While there is a stashed signal for a thread, we block any more signals from being delivered to the thread. This means there can never be more than one stashed signal for a thread.

## 3.5 Interrupted System Calls

An untraced system call can be interrupted by a signal causing a signal handler to run, after which the untraced system call may be resumed. This is similar to handling a blocked untraced system call as described above, and we handle it much the same way. The main extra complication is that the signal handler can run application code while interception library code is on the stack. To prevent reentry into the interception library code, we set a thread-local "locked" flag to indicate that any intercepted system calls should avoid trying to take the untraced fast path until the signal handler has returned.

## 3.6 Thread-Local Storage

The interception library code needs access to thread-local data. Unfortunately in some rare cases the standard Linux thread-local data machinery is not set up correctly,

e.g. when certain applications (Chromium!) use the raw `clone` system call to create a thread. To work around this, RR allocates a special "thread locals" memory page in each process at a fixed address, and on every context switch copies out the old thread's values and copies in the new thread's values for that page.

This approach also simplifies handling of `fork`. Unlike regular Linux thread-local storage, we make our thread-locals start off initialized to zero after a `fork`. Thus when calling into the interception library after `fork`, it can reinitialize itself.

## 3.7 Stack Handling

On x86-64, leaf functions are allowed to use a 128-byte "red zone" of memory below the stack pointer to avoid the cost of adjusting the stack pointer. The interception library needs to avoid modifying that memory. Using memory below the red zone is dangerous because it might cause a spurious stack overflow. (Go programs with many small stacks are particularly vulnerable to this.) Therefore the trampolines that enter the interception library also switch stacks to a temporary stack allocated as part of the thread's scratch buffer.

In theory this stack switching could cause problems for application signal handlers that expect to be called on the regular stack but instead see our switched stack. In practice this has not been a problem yet. If it was a problem, we would probably have to drastically modify our strategy for application signal handling so that we completely unwind the interception library stack frames and run the application signal handler as if the original system call instruction had been interrupted.

## 3.8 Handling Replay

Conceptually, during recording we need to copy system call output buffers to a trace buffer, and during replay we need to copy results from the trace buffer to system call output buffers. This is a problem because the interception library is part of the recording and replay and therefore should execute the same code in both cases. (Previous work with user-level system call interception [14, 38, 24, 32] avoided these problems by having less strict goals for replay fidelity.)

For this reason (and to avoid races of the sort discussed in Section 2.3), the interception library redirects system call outputs to write directly to the trace buffer. After the system call completes, the interception library copies the output data from the trace buffer to the original output buffer(s). During replay the untraced system call instruction is replaced with a no-op, so the system call does not occur; the results are already present in the trace buffer

so the post-system-call copy from the trace buffer to the output buffer(s) does what we need.

During recording, each untraced system call sets a result register and the interception library writes it to the trace buffer. Replay must read the result register from the trace buffer instead. We use a conditional move instruction so that control flow is perfectly consistent between recording and replay. The condition is loaded from an `is_replay` global variable, so the register holding the condition is different over a very short span of instructions (and explicity cleared afterwards).

Handling "in-out" system call memory parameters is tricky. During recording we copy the input buffer to the trace buffer, pass the system call a pointer to the trace buffer, then copy the trace buffer contents back to the input buffer. Performing that first copy during replay would overwrite the trace buffer values holding the system call results, so during replay we turn that copy into a no-op using a conditional move to set the source address copy to the destination address.

We could allow replay of the interception library to diverge further from its recorded behavior, but that would have to be done very carefully. We'd have to ensure the RCB count was identical along both paths, and that register values were consistent whenever we exit the interception library or trap to RR within the interception library. It's simplest to minimize the divergence.

Replay of the interception library could be simplified a little by taking a different approach in which untraced system calls are turned into traced system calls during replay. However, it is important for replay performance that replay of untraced system calls avoid context switching to the RR supervisor process, and for our purposes replay performance is important.

## 3.9   Optimizing Reads With Block Cloning

When an input file is on the same filesystem as the recorded trace and the filesystem supports copy-on-write cloning of file blocks, for large block-aligned `reads` the system call interception code clones the data to a per-thread "cloned-data" trace file, bypassing the normal system-call recording logic. This greatly reduces space and time overhead for file-read-intensive workloads; see the next section.

This optimization works by cloning the input blocks and then reading the input data from the original input file. This opens up a possible race: between the clone and the read, another process could overwrite the input file data, in which case the data read during replay would differ from the data read during recording, causing replay to fail. However, when a file read races with a write under Linux, the reader can receive an arbitrary mix of old and new data, so such behavior would almost certainly be a severe bug, and in practice such bugs do not seem to be common. The race could be avoided by reading from the cloned-data file instead of the original input file, but that performs very poorly because it defeats Linux's readahead optimizations (since the data in the cloned-data file is never available until just before it's needed).

## 3.10   Retrospective

In-process system-call interception looks reasonably simple in principle, but in practice it has been a large maintenance burden. It has to work unchanged during recording and replay, it has to behave like a regular system call as much as possible (e.g. being atomic with respect to signals), and it has to be protected from various sorts of application misbehavior. It cannot be protected from malicious applications. On the other hand, without it we could not meet our performance targets.

## 4   Results

## 4.1   Workloads

Benchmarks were chosen to illuminate RR's strengths and weaknesses, while also containing representatives of real-world usage. They were tuned to fit in system memory (to minimize the impact of I/O on test results), to run for about 30 seconds each (except for *cp* where a 30s run time would require it to not fit in memory).

*cp* duplicates a `git` checkout of `glibc` (revision 2d02fd07) using `cp -a` (15200 files constituting 732MB of data, according to `du -h`). `cp` is single-threaded, making intensive use of synchronous reads and a variety of other filesystem-related system calls.

*make* builds DynamoRio [11] (version 6.1.0) with `make -j8` (`-j8` omitted when restricting to a single core). This tests potentially-parallel execution of many short-lived processes.

*octane* runs the Google Octane benchmark under the Mozilla Spidermonkey Javascript engine (Mercurial revision 9bd900888753). This illustrates performance on CPU-intensive code in a complex language runtime.

*htmltest* runs the Mozilla Firefox HTML forms tests (Mercurial revision 9bd900888753). The harness is excluded from recording (using `mach mochitest -f plain --debugger RR dom/html/test/forms`). This is an example from real-world usage. About 30% of user-space CPU time is in the harness.

*sambatest* runs a Samba (git revision 9ee4678b) UDP echo test via `make test TESTS=samba4.echo.udp`. This is an example from real-world usage.

All tests run on a Dell XPS15 laptop with a quad-core Intel Skylake CPU (8 SMT threads), 16GB RAM and a 512GB SSD using Btrfs in Fedora Core 23 Linux.

## 4.2  Overhead

Table 1 shows the wall-clock run time of various configurations, normalized to the run time of the baseline configuration. *octane* is designed to run for a fixed length of time and report a score, so we report the ratio of the baseline score to the configuration-under-test score — except for replay tests, where the reported score will necessarily be the same as the score during recording. For *octane* replay tests we report the ratio of the baseline score to the recorded score, multiplied by the ratio of replay run time to recording run time. Each test was run six times, discarding the first result and reporting the geometric mean of the other five results. Thus the results represent warm-cache performance.

"Single core" reports the overhead of just restricting all threads to a single core using Linux `taskset`.

"Record no-intercept" and "Replay no-intercept" report overhead with in-process system-call interception disabled (which also disables block cloning). "Record no-cloning" reports overhead with just block cloning disabled.

"DynamoRio-null" reports the overhead of running the tests under the DynamoRio [11] (version 6.1.0) "null tool", to estimate a lower bound for the overhead of using dynamic code instrumentation as an implementation technique. (DynamoRio is reported to be among the fastest dynamic code instrumentation engines.)

## 4.3  Observations

Overhead on *make* is significantly higher than for the other workloads. Forcing *make* onto a single core imposes major slowdown. Also, *make* forks and execs 2430 processes, mostly short-lived. (The next most prolific workload is *sambatest* with 89.) In-process system-call interception only starts working in a process once the interception library has been loaded, but at least 80 system calls are performed before that completes, so its effectiveness is limited for short-lived processes.

Figure 4 shows the overall recording and replay overhead for workloads other than *make*. Error bars in figures show 95% confidence intervals; these results are highly stable across runs.

Excluding *make*, RR's recording slowdown is less than a factor of two. Excluding *make*, RR's replay overhead is lower than its recording overhead. Replay can even be faster than normal execution, in *cp* because system calls do less work. For interactive applications, not represented here, replay can take much less time than the
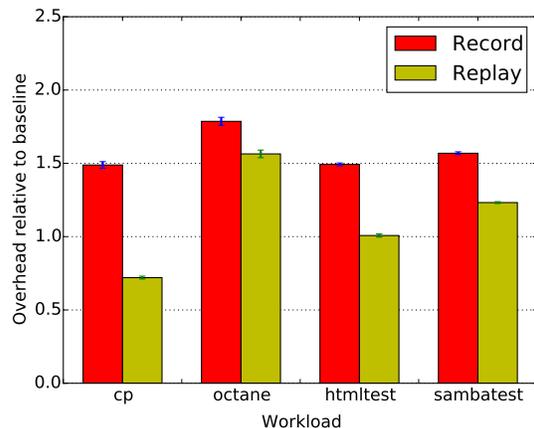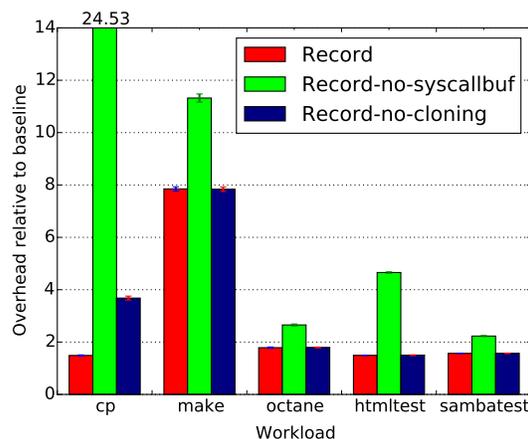


Figure 4: Run-time overhead excluding *make*



Figure 5: Impact of optimizations

original execution because idle periods are eliminated.

*octane* is the only workload here other than *make* making significant use of multiple cores, and this accounts for the majority of RR's overhead on *octane*.

Figure 5 shows the impact of system-call interception and blocking cloning on recording. The system-call interception optimization produces a large reduction in recording (and replay) overhead. Cloning file data blocks is a major improvement for *cp* recording but has essentially no effect on the other workloads.

Figure 6 compares RR recording overhead with DynamoRio's "null tool", which runs all code through the DynamoRio instrumentation engine but does not modify the code beyond whatever is necessary to maintain supervised execution; this represents a minimal-overhead code instrumentation configuration. DynamoRio crashed on *octane* [1]. *cp* executes very little user-space code and

---

[1]We reported DynamoRio's crash on our "octane" work-

12

| Workload | Baseline duration | Record | Replay | Single core | Record no-intercept | Replay no-intercept | Record no-cloning | DynamoRio-null |
|---|---|---|---|---|---|---|---|---|
| cp | 1.04s | 1.49× | 0.72× | 0.98× | 24.53× | 15.39× | 3.68× | 1.24× |
| make | 20.99s | 7.85× | 11.93× | 3.36× | 11.32× | 14.36× | 7.84× | 10.97× |
| octane | 32.73s | 1.79× | 1.56× | 1.36× | 2.65× | 2.43× | 1.80× | crash |
| htmltest | 23.74s | 1.49× | 1.01× | 1.07× | 4.66× | 3.43× | 1.50× | 14.03× |
| sambatest | 31.75s | 1.57× | 1.23× | 0.95× | 2.23× | 1.74× | 1.57× | 1.43× |

Table 1: Run-time overhead



Figure 6: Comparison with DynamoRio-null

| Workload | Compressed MB/s | `deflate` ratio | Cloned blocks MB/s |
|---|---|---|---|
| cp | 19.03 | 4.87× | 586.14 |
| make | 15.82 | 8.32× | 5.50 |
| octane | 0.08 | 8.33× | 0.00 |
| htmltest | 0.79 | 5.94× | 0.00 |
| sambatest | 6.85 | 21.87× | 0.00 |

Table 2: Storage space usage

| Workload | Baseline | Record | Replay | Single core |
|---|---|---|---|---|
| cp | 0.51 | 34.54 | 9.11 | 0.51 |
| make | 510.51 | 327.19 | 314.16 | 288.65 |
| octane | 419.47 | 610.48 | 588.01 | 392.95 |
| htmltest | 690.81 | 692.06 | 324.71 | 689.75 |
| sambatest | 298.68 | 400.49 | 428.79 | 303.03 |

Table 3: Memory usage (peak PSS MB)

DynamoRio's overhead is low on that workload. On *make* and *sambatest* DynamoRio overhead is similar to RR recording, even though on *make* DynamoRio can utilize multiple cores. On *htmltest* DynamoRio's overhead is very high, possibly because that test runs a lot of Javascript with dynamically generated and modified machine code. Implementing record-and-replay on top of dynamic instrumentation would incur significant additional overhead, so we would expect the resulting system to have significantly higher overhead than RR.

### 4.4 Storage Space Usage

RR traces contain three kinds of data: cloned (or hard-linked) files used for memory-map operations, cloned file blocks, and all other trace data, especially event metadata and the results of general system calls.

Memory-mapped files are mostly the executables and libraries loaded by tracees. While the original files are not changed or removed, which is usually true in practice, their clones take no additional space and require no data writes. RR makes no attempt to consolidate duplicate file clones, so most traces contain many duplicates and reporting meaningful space usage for these files is

load to the developers at `https://github.com/DynamoRIO/dynamorio/issues/1930`.

both difficult and unimportant in practice. The same is true for cloned file blocks.

Table 2 shows the storage usage of each workload, in MB/s, for general trace data and cloned file blocks. We compute the geometric mean of the data usage for each trace and divide by the run-time of the workload baseline configuration. Space consumption shows very little variation between runs.

Different workloads have highly varying space consumption rates, but several MB/s is easy for modern systems to handle. In real-world usage, trace storage has not been a concern.

### 4.5 Memory Usage

Table 3 shows the memory usage of each workload. Every 10ms we sum the proportional-set-size ("PSS") values of all workload processes (including RR if running); we determine the peak values for each run and take their geometric mean. In Linux, each page of memory mapped into a process's address space contributes $1/n$ pages to that process's PSS, where $n$ is the number of processes mapping the page; thus it is meaningful to sum PSS values over processes which share memory. The same data are shown in Figure 7. In the figure, the fraction of PSS
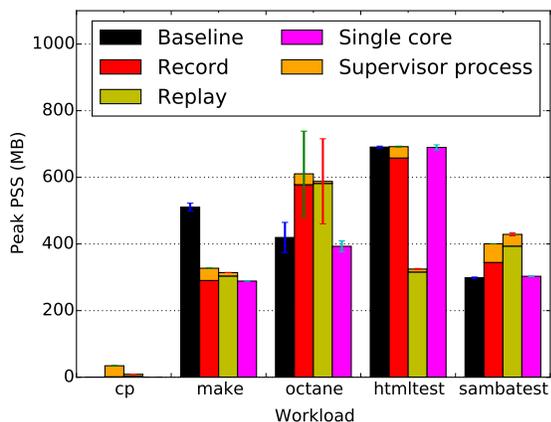
Figure 7: Memory usage

used by the RR process is shown in orange. Memory usage data was gathered in separate runs from the timing data shown above, to ensure the overhead of gathering memory statistics did not impact those results.

Given these experiments ran on an otherwise unloaded machine with 16GB RAM and all data fits in cache, none of these workloads experienced any memory pressure. *cp* uses almost no memory. In *make*, just running on a single core reduces peak PSS significantly because not as many processes run simultaneously. In *octane* memory usage is volatile (highly sensitive to small changes in GC behavior) but recording significantly increases application memory usage; recording also increases application memory usage a small amount in *sambatest* but slightly decreases it in *htmltest*. (We expect a small increase in application memory usage due to system-call interception and scratch buffers.) These effects are difficult to explain due to the complexity of the applications, but could be due to changes in timing and/or effects on application or operating system memory management heuristics.

Replay memory usage is similar to recording except in *htmltest*, where it's dramatically lower because we're not replaying the test harness.

RR's memory overhead is not an issue in practice.

## 5 Hardware/Software Design Constraints

We summarize the hardware and software features RR depends on, for system designers who may be interested in supporting RR-like record-and-replay.

### 5.1 Hardware

As discussed in Section 2.4.1, RR requires a "determinstic" hardware performance counter to measure application progress. The ideal performance counter for

our purposes would count the exact number of instructions retired as observed in user-space (e.g., counting an interrupted-and-restarted instruction once). Virtual machines should support reliable performance-counter virtualization. Currently RR works under KVM and VMware, but VMware's *VM exit clustering* optimization [7], as implemented, breaks the determinism of the RCB counter and must be manually disabled.

Some x86 CPU instructions are nondeterministic. Section 2.6 discusses our current workarounds for this. Exposing hardware and OS support for trapping CPUID is important for long-term control over these instructions.

We would like to support record-and-replay of programs using hardware transactional memory (XBEGIN/XEND). It would suffice if hardware and the OS could be configured to raise a signal on any failed transaction.

Trapping on all other nondetermnistic instructions (e.g. RDRAND) would be useful.

Porting RR to ARM failed because all ARM atomic memory operations use the "load-linked/store-conditional" approach, which is inherently nondeterminstic. The conditional store can fail because of non-user-space-observable activity, e.g. hardware interrupts, so counts of retired instructions or conditional branches for code performing atomic memory operations are nondeterminstic. These operations are inlined into very many code locations, so it appears patching them is not feasible except via pervasive code instrumentation or compiler changes. On x86(-64), atomic operations (e.g. compare-and-swap) are deterministic in terms of user-space state, so there is no such problem.

### 5.2 Software

As noted in Section 2.5, RR depends on configuring applications to avoid sharing memory with non-recorded processes.

We described how RR performance depends on modern Linux features: seccomp-bpf to selectively trap system calls, PERF_COUNT_SW_CONTEXT_SWITCHES performance events to handling blocking system calls, and copy-on-write file and block cloning APIs to reduce I/O overhead.

Efficient record-and-replay depends on clearly identifying a boundary within which code is replayed deterministically, and recording and replaying the timing and contents of all inputs into that boundary. In RR, that boundary is mostly the interface between the kernel and user-space. This suits Linux: most of the Linux user/kernel interface is stable across OS versions, relatively simple and well-documented, and it's easy to count hardware performance events occurring within the boundary (i.e. all user-space events for a specific pro-

cess). This is less true in other operating systems. For example, in Windows, the user/kernel interface is not publicly documented, and is apparently more complex and less stable than in Linux. Implementing and maintaining the RR approach for Windows would be considerably more challenging than for Linux, at least for anyone other than the OS vendor.

RR relies on models that specify for each system call, how to efficiently determine from the before-and-after user-space state which memory locations were modified by the system call. The smaller the number of different system-call interfaces (e.g. using APIs based on `read` and `write` through file descriptors instead of arbitrary `ioctl`s), the less work that is.

## 6 Lessons Learned

### 6.1 Deployability and Usage

We set out to build a record-and-replay system optimized for deployability, and we seem to have achieved that. Some Linux distributions package RR, and it's easy to download or build for those that don't. On many distributions users must tweak the `perf_events_paranoid` kernel setting to enable non-root access to `PERF_COUNT_SW_CONTEXT_SWITCHES`, but otherwise RR "just works" for unprivileged users running most applications in non-virtualized environments. This low barrier to entry is extremely important, since encouraging users to use new tools is difficult under any circumstances.

We try hard to avoid requiring manual configuration. For example, RR automatically detects kernel and virtualization bugs affecting it by running quick tests on every launch; if found, it prints warnings with advice about how to fix the environment (e.g. by upgrading software components) but automatically activates workarounds for the bugs.

Our goal was always to work very well for a significant set of users. We chose good performance for low-parallelism applications (and high slowdown for highly parallel applications) over mediocre performance for all applications because we felt the former would please more users, especially at Mozilla. It still appears we made the right choice.

RR has been mostly feature-complete for at least two years. Some of the authors have been developing new tools on top of RR, but for the core open source project we have been focusing on reliability and maintainance to ensure RR keeps working well as hardware, Linux and applications evolve. Avoiding scope creep has helped keep the complexity of RR under control, and therefore contributed to reliability. Reliability matters because it's very important to users that RR not fall over during a long debugging session. It's a lot more important that RR work 100% of the time on 90% of applications than 90% of the time on 100% of applications; the former pleases 90% of users, but the latter pleases no-one.

Space and time overhead of RR, at least for low-parallelism workloads, seems to be mostly a non-issue for our users. Reported performance problems tend to be due to system calls not having a fast-path in the interception library, or just silly RR bugs, though sometimes denying direct access to the GPU causes performance problems. For many users the overhead of RR is less than the difference between an optimized build and a debug build of their application, and therefore not very important.

Some record-and-replay techniques [8] reduce recording overhead in exchange for slower replay, but for our users replay performance is very important — often more important than recording performance (e.g. when debugging an easily-reproduced bug, recording executes code once but replay with reverse-execution may execute it many times.)

Some record-and-replay applications need to frequently create and resume checkpoints of replayed execution state (e.g. to simulate reverse execution). Making checkpoints cheap in time and memory is an often-overlooked design constraint. Anecdotally, creating and restoring checkpoints in VMWare's record-and-replay debugger was rather slow, perhaps because the entire OS state must be checkpointed. RR's process-level design made it easy for us to implement cheap checkpoints: we use `fork` to copy address spaces and delay creating the non-main threads until the checkpoint is resumed (and most checkpoints are never resumed). `fork` is (mostly) copy-on-write and is very well optimized on Linux, so creating a checkpoint typically takes less than ten milliseconds. Although checkpointing and restoring a general Linux process from user-space is rather complicated [1] because a lot of kernel state may need to be saved and reestablished, checkpointing and restoring *for replay only* is relatively simple, because only a minimal amount of kernel state is reproduced during replay.

### 6.2 Engineering RR

RR is complex and relies on poorly documented hardware and OS features. Over time we've encountered many hardware, kernel and hypervisor bugs affecting RR and had to work around them (while doing our best to get the underlying bugs fixed, too). Those bugs, and many bugs in RR itself, are often difficult to diagnose and fix. In particular, small bugs can cause small changes in replay state that don't cause observable divergence until long after the root cause. To some extent RR makes debugging easier for many at the cost of very difficult de-

bugging for its developers!

Therefore it is important to detect issues as soon after the root cause as possible. RR has a high density of consistency assertions, most of which are checked even in release builds, and many of which produce custom diagnostic messages for easier understanding in the field. For example, unsupported system calls or ioctls produce a message clearly identifying the problem and the offending system call. In general we try to make our invariants as strict as possible, even stricter than may seem necessary, the better to catch bugs and unexpected situations.

RR has a lot of logging code, which is built into debug and release builds so that end users can enable it just by setting an environment variable. This makes it easier to get meaningful reports from users, and we have been able to diagnose many bugs just by asking users to collect the right logs.

Being free software has helped. A small community has developed around RR and some users have been able to diagnose and fix the bugs they've found, particularly expanding system call coverage.

Inspecting the state of tracees during recording and replay is very important for debugging RR. Therefore RR offers an "emergency debugger" feature, which can be manually triggered or triggered automatically when an error occurs. The emergency debugger activates a `gdb` server (the same code RR uses to provide replay debugging) to which `gdb` can attach to inspect the register and memory state of the tracee, with full support for symbols, source mapping, etc.

To track down subtle errors in memory state during replay, RR supports taking checksums and full checkpoints of memory at selected points during recording and comparing them with the replay. We also support single-stepping tracees during recording and replay with logging of all register states, to help narrow down divergences (though that can cause bugs to disappear).

To understand obscure details of Linux kernel behavior and bugs, we frequently have had to resort to reading kernel source code. We have also used the kernel's `ftrace` [2] feature several times. Implementing something like RR on a system without access to full OS source code would be daunting.

RR has a fairly good suite of automated tests that run on every commit. As of May 15, 2017, RR comprises 75028 lines of C and C++ code, of which 21403 lines belong to the test suite. We have spent a lot of effort tracking down intermittent failures in the test suite; this usually reveals bugs in the tests, but sometimes has revealed bugs in RR or the kernel that would have been difficult to diagnose in the field. Since users can easily run the tests, that helps us to identify issues in a user's environment causing problems for RR.

RR is an excellent debugging tool, so we would like to be able to debug RR using RR itself. This is hard because RR uses kernel features like `ptrace` very intensively and supporting those features well in RR is difficult. Furthermore, debugging RR recording using RR usually wouldn't make sense because if recording workload $X$ doesn't work, recording workload "RR + $X$" is even less likely to work. However, RR replay is simpler than recording at the kernel level (since most tracee system calls don't run), and we support RR recording and replaying of RR replay. This has been useful, although keeping the multiple levels of RR straight in one's head can be confusing!

Anecdotally we have heard that RR has been used to help debug QEMU-based virtual machine record and replay [18]. In turn, virtual-machine-based record and replay with a good debugging experience could make debugging RR a lot easier. Because these approaches are quite different, they could be very helpful for debugging each other.

# 7 Related Work

## 7.1 Whole-System Replay

ReVirt [20] was an early project that recorded and replayed the execution of an entire virtual machine. VMware [31] used the same approach to support record-and-replay debugging in VMware Workstation, for a time, but discontinued the feature. The full-system simulator Simics supports reverse-execution debugging via deterministic reexecution [23]. There have been efforts to add some record-and-replay support to QEMU [18, 19, 40] and Xen [21, 13]. Whole-system record-and-replay can be useful, but it is often inconvenient to hoist the application into a virtual machine. Many applications of record-and-replay require cheap checkpointing, and checkpointing a VM image is generally more expensive than checkpointing one or a few processes.

## 7.2 Replaying User-Space With Kernel Support

Scribe [29], dOS [9] and Arnold [17] replay a process or group of processes by extending the OS kernel with record-and-replay functionality. Kernel changes make maintenance and deployment more difficult — unless record-and-replay is integrated into the base OS. But adding invasive new features to the kernel has risks, so if record-and-replay can be well implemented outside the kernel, moving it into the kernel may not be desirable.

16

### 7.3 Pure User-Space Replay

Pure user-space record-and-replay systems have existed since at least MEC [14], and later Jockey [38] and liblog [24]. Those systems did not handle asynchronous event timing and other OS features. PinPlay [36], iDNA [10], UndoDB [3] and TotalView ReplayEngine [25] use code instrumentation to record and replay asynchronous event timing. Unlike UndoDB and RR, PinPlay and iDNA instrument all loads, thus supporting parallel recording in the presence of data races and avoiding having to compute the effects of system calls, but this gives them higher overhead than the other systems. Compared to the other systems that support asynchronous events, RR achieves lower overhead by avoiding code instrumentation.

### 7.4 Higher-Level Replay

Record-and-replay features have been integrated into language-level virtual machines. DejaVu [15] added record-and-replay capabilities to the Jalapeño Java VM. Microsoft IntelliTrace [5] instruments CLR bytecode to record high-level events and the parameters and results of function calls; it does not produce a full replay. Systems such as Chronon [16] for Java instrument bytecode to collect enough data to provide the appearance of replaying execution, without actually doing a replay. Dolos [12] provides record-and-replay for JS applications in Webkit by recording and replaying nondeterministic inputs to the browser. R2 [26] provides record-and-replay by instrumenting library interfaces; handling data races or asynchronous events requires user effort to isolate the nondeterminism. Such systems are all significantly narrower in scope than the ability to replay general user-space execution.

### 7.5 Parallel Replay

Recording application threads running concurrently on multiple cores, with the possibility of data races, with low overhead, is extremely challenging. PinPlay [36] and iDNA/Nirvana [10] instrument shared-memory loads and report high overhead. SMP-ReVirt [21] tracks page ownership using hardware page protection and reports high overhead on benchmarks with a lot of sharing. DoublePlay [41] runs two instances of the application and thus has high overhead when the application alone could saturate available cores. ODR [8] has low recording overhead but replay can be extremely expensive and is not guaranteed to reproduce the same program states. Castor [32] instruments synchronization code by modifying compilers and runtime systems, which creates barriers to easy deployment, and cannot replay reliably in the presence of data races.

The best hope for general, low-overhead parallel recording seems to be hardware support. Projects such as FDR [43], BugNet [34], Rerun [27], DeLorean [33] and QuickRec [37] have explored low-overhead parallel recording hardware.

## 8 Future Work

Probably the main issue holding back RR deployment currently is the lack of support for virtualized hardware performance counters in the major cloud providers. This is not really a technical issue; those providers need to be convinced that the benefits of virtualizing counters (or at least those used by RR) outweigh the costs and risks. Expanding the usage of RR is part of our effort to address this problem.

Some use-cases for record-and-replay involve transporting recordings from one machine to another, which RR does not really support yet. Once `CPUID` faulting support is widely available, that will be easier to implement.

RR perturbs execution, especially by forcing all threads onto a single core, and therefore can fail to reproduce bugs that manifest outside RR. We have addressed this problem by introducing a "chaos mode" that intelligently adds randomness to scheduling decisions, enabling us to reproduce many more bugs, but that work is beyond the scope of this paper. There are many more opportunities to enhance the recorder to find more bugs.

Putting record-and-replay support in the kernel would have performance benefits, e.g. reducing the cost of recording context switches. We may be able to find reusable primitives that can be added to kernels to improve the performance of user-space record-and-replay while being less invasive than a full kernel implementation. It would be particularly interesting to find small kernel extensions that could eliminate the need for in-process system call interception. Some way to record the data passing through `copy_to_user` would be a good start.

Recording multiple processes running in parallel on multiple cores seems feasible if they do not share memory — or, if they share memory, techniques inspired by SMP-ReVirt [21], dthreads [30] or Castor [32] may work for some workloads on existing hardware. We hope that popularizing practical use of record-and-replay will help support the economic argument for adding hardware support for data-race recording [37].

The applications of record-and-replay are perhaps more interesting and important than the base technology. For example, one can perform high-overhead dynamic analysis during replay [17, 18, 36], potentially parallelized over multiple segments of the execution. With

RR's no-instrumentation approach, one could collect performance data such as sampled stacks and performance counter values during recording, and correlate that data with rich analysis generated during replay (e.g. cache simulation). Always-on record-and-replay would make finding and fixing bugs in the field much easier. Demonstrating compelling applications for record-and-replay will build the case for building support into commodity hardware and software.

## 9   Conclusions

The current state of Linux on commodity x86 CPUs enables single-core user-space record-and-replay with low overhead, without pervasive code instrumentation — but only just. This is fortuitous; we use software and hardware features for purposes they were not designed to serve. It is also a recent development; five years ago `seccomp-bpf` and the Linux file cloning APIs did not exist, and commodity architectures with a deterministic hardware performance counter usable from user-space had only just appeared (Intel Westmere)[2]. By identifying the utility of these features for record-and-replay, we hope that they will be supported by an increasingly broad range of future systems. By providing an open-source, easy-to-deploy, production-ready record-and-replay framework we hope to enable more compelling applications of this technology.

## References

[1] Checkpoint/restore in userspace. `https://criu.org/Main_Page`. Accessed: 2017-05-16.

[2] Debugging the kernel using ftrace. `https://lwn.net/Articles/365835/`. Accessed: 2017-05-16.

[3] Reversible debugging tools for C/C++ on Linux & Android. `http://undo-software.com`. Accessed: 2016-04-16.

[4] Soft-dirty ptes. `https://www.kernel.org/doc/Documentation/vm/soft-dirty.txt`. Accessed: 2017-05-15.

[5] Understanding IntelliTrace part I: What the @#$% is IntelliTrace? `https://blogs.msdn.microsoft.com/zainnab/2013/02/12/understanding-intellitrace-part-`

i-what-the-is-intellitrace. Accessed: 2016-04-16.

[6] Wine windows-on-posix framework. `https://www.winehq.org`. Accessed: 2016-09-20.

[7] O. Agesen, J. Mattson, R. Rugina, and J. Sheldon. Software techniques for avoiding hardware virtualization exits. In *Proceedings of the 2012 USENIX Annual Technical Conference*, June 2012.

[8] G. Altekar and I. Stoica. ODR: Output-deterministic replay for multicore debugging. In *Proceedings of the ACM SIGOPS 22nd Symposium on Operating Systems Principles*, October 2009.

[9] T. Bergan, N. Hunt, L. Ceze, and S. D. Gribble. Deterministic process groups in dOS. In *Proceedings of the 9th USENIX Symposium on Operating Systems Design and Implementation*, October 2010.

[10] S. Bhansali, W.-K. Chen, S. de Jong, A. Edwards, R. Murray, M. Drinić, D. Mihočka, and J. Chau. Framework for instruction-level tracing and analysis of program executions. In *Proceedings of the 2nd International Conference on Virtual Execution Environments*, June 2006.

[11] D. Bruening, Q. Zhao, and S. Amarasinghe. Transparent dynamic instrumentation. In *Proceedings of the 8th International Conference on Virtual Execution Environments*, March 2012.

[12] B. Burg, R. Bailey, A. Ko, and M. Ernst. Interactive record/replay for web application debugging. In *Proceedings of the 26th ACM Symposium on User Interface Software and Technology*, October 2013.

[13] A. Burtsev, D. Johnson, M. Hibler, E. Eide, and J. Regehr. Abstractions for practical virtual machine replay. In *Proceedings of the 12th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments*, April 2016.

[14] M. E. Chastain. `https://lwn.net/1999/0121/a/mec.html`, January 1999. Accessed: 2016-04-16.

[15] J.-D. Choi, B. Alpern, T. Ngo, M. Sridharan, and J. Vlissides. A perturbation-free replay platform for cross-optimized multithreaded applications. In *Proceedings of the 15th International Parallel and Distributed Processing Symposium*, April 2001.

[16] P. Deva. `http://chrononsystems.com/blog/design-and-architecture-of-the-chronon-record-0`, December 2010. Accessed: 2016-04-16.

---

[2]Performance counters have been usable for kernel-implemented replay [20, 35] for longer, because kernel code can observe and compensate for events such as interrupts and page faults.

[17] D. Devecsery, M. Chow, X. Dou, J. Flinn, and P. M. Chen. Eidetic systems. In *Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation*, October 2014.

[18] B. Dolan-Gavitt, J. Hodosh, P. Hulin, T. Leek, and R. Whelan. Repeatable reverse engineering with PANDA. In *Proceedings of the 5th Program Protection and Reverse Engineering Workshop*, December 2015.

[19] P. Dovgalyuk. Deterministic replay of systems execution with multi-target QEMU simulator for dynamic analysis and reverse debugging. 2012.

[20] G. Dunlap, S. King, S. Cinar, M. A. Basrai, and P. M. Chen. ReVirt: Enabling intrusion analysis through virtual-machine logging and replay. In *Proceedings of the 5th USENIX Symposium on Operating Systems Design and Implementation*, December 2002.

[21] G. W. Dunlap, D. G. Lucchetti, M. A. Fetterman, and P. M. Chen. Execution replay of multiprocessor virtual machines. In *Proceedings of the 4th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments*, March 2008.

[22] J. Engblom. A review of reverse debugging. In *System, Software, SoC and Silicon Debug Conference*, September 2012.

[23] J. Engblom, D. Aarno, and B. Werner. Full-system simulation from embedded to high-performance systems. In *Processor and System-on-Chip Simulation*, 2010.

[24] D. Geels, G. Altekar, S. Shenker, and I. Stoica. Replay debugging for distributed applications. In *Proceedings of the 2006 USENIX Annual Technical Conference*, June 2006.

[25] C. Gottbrath. Reverse debugging with the TotalView debugger. In *Cray User Group Conference*, May 2008.

[26] Z. Guo, X. Wang, J. Tang, X. Liu, Z. Xu, M. Wu, M. F. Kaashoek, and Z. Zhang. R2: An application-level kernel for record and replay. In *Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation*, December 2008.

[27] D. Hower and M. Hill. Rerun: Exploiting episodes for lightweight memory race recording. In *Proceedings of the 35th Annual International Symposium on Computer Architecture*, June 2008.

[28] U. Hlzle, C. Chambers, and D. Ungar. Optimizing dynamically-typed object-oriented languages with polymorphic inline caches. In *Proceedings of the 1991 European Conference on Object-Oriented Programming*, July 1991.

[29] O. Laadan, N. Viennot, and J. Nieh. Transparent, lightweight application execution replay on commodity multiprocessor operating systems. In *Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems*, June 2010.

[30] T. Liu, C. Curtsinger, and E. Berger. Dthreads: Efficient deterministic multithreading. In *Proceedings of the ACM SIGOPS 23rd Symposium on Operating Systems Principles*, October 2011.

[31] V. Malyugin, J. Sheldon, G. Venkitachalam, B. Weissman, and M. Xu. ReTrace: Collecting execution trace with virtual machine deterministic replay. In *Proceedings of the Workshop on Modeling, Benchmarking and Simulation*, June 2007.

[32] A. J. Mashtizadeh, T. Garfinkel, D. Terei, D. Mazières, and M. Rosenblum. Towards practical default-on multi-core record/replay. In *Proceedings of the 22nd International Conference on Architectural Support for Programming Languages and Operating Systems (to appear)*, April 2017.

[33] P. Montesinos, L. Ceze, and J. Torrellas. DeLorean: Recording and deterministically replaying shared-memory multiprocessor execution efficiently. In *Proceedings of the 35th Annual International Symposium on Computer Architecture*, June 2008.

[34] S. Narayanasamy, G. Pokam, and B. Calder. Bugnet: Continuously recording program execution for deterministic replay debugging. In *Proceedings of the 32nd Annual International Symposium on Computer Architecture*, June 2005.

[35] M. Olszewski, J. Ansel, and S. Amarasinghe. Kendo: Efficient deterministic multithreading in software. In *Proceedings of the 14th International Conference on Architectural Support for Programming Languages and Operating Systems*, March 2009.

[36] H. Patil, C. Pereira, M. Stallcup, G. Lueck, and J. Cownie. PinPlay: A framework for deterministic replay and reproducible analysis of parallel programs. In *Proceedings of the 8th Annual IEEE/ACM International Symposium on Code Generation and Optimization*, April 2010.

[37] G. Pokam, K. Danne, C. Pereira, R. Kassa, T. Kranich, S. Hu, J. Gottschlich, N. Honarmand, N. Dautenhahn, S. King, and J. Torrellas. Quick-Rec: Prototyping an intel architecture extension for record and replay of multithreaded programs. In *Proceedings of the 40th Annual International Symposium on Computer Architecture*, June 2013.

[38] Y. Saito. Jockey: A user-space library for record-replay debugging. In *Proceedings of the 6th International Symposium on Automated Analysis-driven Debugging*, September 2005.

[39] K. Serebryany, D. Bruening, A. Potapenko, and D. Vyukov. Addresssanitizer: A fast address sanity checker. In *Proceedings of the 2012 USENIX Annual Technical Conference*, June 2012.

[40] D. Srinivasan and X. Jiang. Time-traveling forensic analysis of VM-based high-interaction honeypots. In *Security and Privacy in Communication Networks: 7th International ICST Conference*, September 2011.

[41] K. Veeraraghavan, D. Lee, B. Wester, J. Ouyang, P. M. Chen, J. Flinn, and S. Narayanasamy. Doubleplay: Parallelizing sequential logging and replay. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems*, March 2011.

[42] V. Weaver, D. Terpstra, and S. Moore. Non-determinism and overcount on modern hardware performance counter implementations. In *Proceedings of the IEEE International Symposium on Performance Analysis of Systems and Software*, April 2013.

[43] M. Xu, R. Bodik, and M. D. Hill. A "flight data recorder" for enabling full-system multiprocessor deterministic replay. In *Proceedings of the 30th Annual International Symposium on Computer Architecture*, June 2003.