

Jockey: A User-space Library for Record-replay Debugging

Yasushi Saito

Hewlett-Packard Laboratories¹
Palo Alto, CA, USA

Google, Inc.
Mountain View, CA, USA

yasushi.saito@gmail.com

ABSTRACT

Jockey is an execution record/replay tool for debugging Linux programs. It records invocations of system calls and CPU instructions with timing-dependent effects and later replays them deterministically. It supports process checkpointing to diagnose long-running programs efficiently. Jockey is implemented as a shared-object file that runs as a part of the target process. While this design is the key for achieving Jockey's goal of safety and ease of use, it also poses challenges. This paper discusses some of the practical issues we needed to overcome in such environments, including low-overhead system-call interception, techniques for segregating resource usage between Jockey and the target process, and an interface for fine-grain control of Jockey's behavior.

Categories and Subject Descriptors

D.2.5 [Software]: Software Engineering—*Testing and Debugging*;
D.4.9 [Operating Systems]: Systems Programs and Utilities

General Terms

Reliability, Experimentation, Languages, Verification

Keywords

Debugging, Execution record and replay, Checkpointing, x86, Jockey, Linux

1. INTRODUCTION

Jockey is a record/replay tool for Linux. It logs the execution of an ordinary program and replays deterministically later. Jockey is designed to help debug interactive or distributed programs that communicate with the operating system or other computers in a complex fashion. We plan to make Jockey publicly available via <http://www.freshmeat.net>.

¹This work was done while the author was with HP Labs. The author is now employed by Google.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

AADEBUG'05, September 19–21, 2005, Monterey, California, USA.
Copyright 2005 ACM 1-59593-050-7/05/0009 ...\$5.00.

Jockey was originally developed as a debugging aid for FAB (Federated Array of Bricks) [23]. FAB is a high-availability disk array built on a cluster of commodity servers. It provides accesses to logical volumes to iSCSI clients using complex peer-to-peer-style replication and erasure-coding protocols.

Traditional debuggers, such as gdb, provide comprehensive support for debugging single-node, sequential programs. They are, however, not as useful for interactive or distributed programs such as FAB [10, 4]. We identify three key problems and discuss how Jockey alleviates them.

First, the execution of such programs is inherently nondeterministic. The behavior of a process will diverge, depending on interactions with the OS, the user, or other processes. Jockey helps debug such programs by recording every nondeterministic choice the process makes, and replaying the execution as many times as the developer wishes. Thus, debugging for a nondeterministic program is reduced to that for a sequential, repeatable program.

Second, these programs often run for a long period of time, either because they need lots of resources (e.g., scientific computation), or they are server programs (e.g., FAB and distributed hash tables), or they need substantial user interactions (e.g., spreadsheet). Simply reproducing the bug often tests a developer's patience. Jockey alleviates this problem by transparently checkpointing the process state during execution. The developer can start replaying from any checkpoint and easily “time-travel” through the history of execution to diagnose the problem. Checkpointing also bounds the log-space overhead, as log records older than the checkpoint can be discarded.

Third, running a distributed system such as FAB requires starting processes on multiple computers, which is cumbersome and increases the turn-around time for program development. Jockey alleviates this problem by recording and replaying each process independently—after recording the execution of the whole system, the developer can replay each process under a traditional debugger. This can also be a limitation; Jockey could be less useful when one wants to investigate the execution of the whole system at once. We discuss our experience in Section 5.2.

The remainder of this section overviews Jockey's design and discusses its benefits and challenges.

1.1 Goals and approaches

Jockey is designed with two pragmatic goals in mind. First is *ease of use*: Jockey must be easy and safe to deploy. It should work without requiring changes to the target program, the operating system, or the debugger. Second is *generality*. Jockey should be able to handle generic Linux programs, not just those written in a particular programming language or API, such as MPI or CORBA [11].

We achieve the first goal by implementing Jockey as a user-space library that runs as a part of the target process. In contrast to kernel-

```
// test1.c
int main() {
    FILE *f = fopen("/dev/random", "r");
    printf("%x\n", getc(f));
}
```

Figure 1: A simple program, `test1.c`, that reads and displays a random number every time it is run.

```
% cc -o test1 test1.c
% LD_PRELOAD=libjockey.so \
  JOCKEYRC="replay=0" ./test1 # recording
38
% LD_PRELOAD=libjockey.so \
  JOCKEYRC="replay=1" ./test1 # replaying
38
```

Figure 2: The most basic use of Jockey. The program outputs the same number even though it is reading from `/dev/random`. Setting environment variable `LD_PRELOAD` causes the dynamic linker to load `libjockey.so` before other object files. Environment variable `JOCKEYRC` passes parameters to `libjockey.so`.

based approaches [26], Jockey can be used by anyone without administrative privilege or a patched kernel. Developers can continue using their favorite debuggers without change. In addition, this design allows the target program to control or extend Jockey easily, as we discuss in Section 4. Our second goal is achieved by recording and replaying events at a fairly low level—system calls and CPU instructions.

1.2 Non-goals

We do not try to make program execution under Jockey identical to native execution without Jockey. Because Jockey internally needs to perform mmap and file accesses, the mmap addresses and file descriptors allocated to the target process may differ between a native and a recording run. This usually does not cause an additional problem, because programs targeted by Jockey are usually nondeterministic to begin with.

Also, performance is a secondary goal. Jockey is used only during testing and debugging. Some slowdown under Jockey should be acceptable, as long as it does not change the target program’s behavior qualitatively. In practice, as we show in Section 5, Jockey’s overhead is at most 30% for I/O intensive programs, more often close to zero—well within our limit of tolerance.

1.3 Example

The meat of Jockey is `libjockey.so`, an x86 shared-object file. Figure 1 shows a simple program that reads from `/dev/random`, Linux’s random number device. Figure 2 shows the most basic use of Jockey. Recording or replaying program execution requires no change to the source code or the executable file. Simply loading `libjockey.so` on startup causes Jockey to take control of the process. In this example, Jockey intercepts the call to the `read` system call made via `getc`. It logs the value read during the recording phase. When replaying, it reads the value from the log without actually reading from the random device. Jockey can also be invoked in several different ways, as shown in Figure 3.

```
% LD_PRELOAD=libjockey.so \
./test1 --jockey=replay=0 # recording
82
% LD_PRELOAD=libjockey.so \
./test1 --jockey=replay=1 # replaying
82
```

(a) One can pass a command-line parameter `--jockey=` to the target process to control Jockey. This parameter is parsed by `libjockey.so`. For this method to work, the target program must be designed to ignore a command-line string that starts with `--jockey=`.

```
% jockey --replay=0 ./test1 # recording
a9
% jockey --replay=1 ./test1 # replaying
a9
```

(b) One can also start the test program from the `jockey` frontend. `jockey` is a small script just sets the environment variables and executes the target program.

```
% cc test.c -ljockey
% ./test1 --jockey=replay=0 # recording
c1
% ./test1 --jockey=replay=1 # replaying
c1
```

(c) Jockey can also be linked manually to the target program, as shown above. This method is convenient when one frequently runs the program under a debugger.

Figure 3: Alternative ways of running a program under Jockey.

1.4 Challenges and limitations

Our decision to co-locate Jockey with the target process poses challenges and limitations. First, Jockey could be compromised by a seriously buggy or malicious target program—if it wishes, for example, the target program can destroy a memory region used internally by Jockey. We address this problem by segregating the use of resources as much as possible between the target program and Jockey. As a result, Jockey has been able to record and replay most common memory bugs, including accessing free’d memory blocks and off-by-one array accesses. Resource segregation is discussed further in Section 3.2.

The second challenge is recording and replaying events that are not directly initiated by system calls. We describe our solutions to two such types of events, signals and memory-mapped file I/Os in Sections 3.4 and 3.6. There are, however, events that are fundamentally impossible to capture. For example, memory access races that happen with kernel-based pthreads cannot be replayed, because thread context switches are out of Jockey’s control. For this reason, Jockey does not support kernel multi-threading. Similarly, it does not support any program or API that interacts with other processes (or devices) via shared memory or files—e.g., uDAPL [1] for memory-mapped network I/Os. Note that Jockey does support user-space threads, such as Capriccio [30]—in fact, FAB is built on a similar package.

2. RELATED WORK

Execution record/replay has long been advocated as an effective debugging method [6, 22, 21]. This section reviews prior approaches to record/replay debugging and relates them to Jockey.

2.1 Record/replay debugging for a single process

Bugnet was one of the earliest deterministic record/replay tools [32]. It intercepted I/O activities by the processes and took checkpoints of the system periodically. Bugnet, however, supported only programs written to a special API, unlike Jockey that supports generic Linux programs. Flashback [26] is the most recent work along this line. It offers a similar set of functionalities as Jockey—recording and replaying system calls and fork-based checkpointing—but Flashback is offered as a kernel patch. As such, it is less easy and safe to use than Jockey.

In a slightly different approach, some systems record and replay individual memory accesses [19, 16, 25, 7]. They have several advantages over event-based approaches like Bugnet, Flashback, or Jockey. First, some of them enable *reverse execution*—stepping CPU instructions literally backward [19, 7]. Second, they could be more generic because they need not know deeply about the semantics of system calls and other interactions with OS. However, they require a special compiler and have a large logging overhead. Even with sophisticated optimizations, these systems generate logs at the rate of multiple megabytes per second for CPU-intensive programs [16, 25]. Jockey, in contrast, generates only a few hundred bytes per second for such programs, as we show in Section 5.1.

2.2 Record/replay using virtual machines

Revirt is a virtual machine that records and replays low-level interrupts and device activities [2]. It has been proved to be useful for network intrusion detection and diagnosing kernel bugs [5]. Several other papers also propose distributed-system emulation using virtual machines [3, 17]. While these systems are powerful, they are also cumbersome to use—for example, one needs to create a full file system tree for each virtual machine. They are overkill when one is interested only in debugging user-space programs. Jockey is designed to be simpler and easier to use than these systems.

2.3 Record/replay for parallel and distributed programs

Deterministic record/replay has been most effective in parallel and distributed environments [21]. Indeed, earliest tools specifically targeted such environments [32, 19]. Since then, many theoretical improvements have been proposed for both shared-memory parallel programs [13] and message-passing programs [14, 15]. Jockey does not yet support deterministic replay of a distributed system—it can only replay processes within the system independently. As we discuss in Section 5.2, we found this limitation not to be a serious obstacle so far.

3. IMPLEMENTATION OF JOCKEY

We have implemented Jockey on Linux in C++. The dynamic linker (`ld.so`) invokes Jockey's initialization routine immediately after `libjockey.so` is loaded, before the target program starts execution. The initialization routine performs the following tasks.

- (1) For each system call in `libc` with timing- or context- dependent effects, Jockey rewrites its first few instructions and intercepts calls to it. Jockey currently intercepts 80 Linux system calls, including `time`, `recvfrom`, and `select`. Jockey

logs the values generated by these calls during recording, and reads the value from the log during replay.

- (2) Jockey does the same for CPU instructions with nondeterministic effects. It currently patches only `rdtsc`, the x86 instruction for reading the CPU's timestamp counter. It is used, for example, as a pseudo random-number generator in `libc`.
- (3) Jockey checkpoints the process state just before returning control to the target program. In the replay mode, Jockey simply loads the checkpoint. Checkpointing is needed to ensure that the target sees the same set of environment variables and command line parameters during both record and replay. We discuss checkpointing in more detail in Section 3.3.
- (4) Jockey transfers control to the target program. From this moment, Jockey becomes active only when the target executes a system call or a nondeterministic CPU instruction.

The next section describes the first two steps in more detail. Section 3.2 discusses Jockey's efforts to segregate itself from the target program to avoid unnecessary interference. Section 3.3 describes Jockey's checkpointing function (step (3)), along with the challenges we had to overcome.

3.1 Instruction patching

As an example, Figure 4 shows how the `time` system call is recorded and replayed. (a) shows the first few CPU instructions of `time` in `libc.so`. When Jockey starts, it writes a `jmp` instruction in the first 5 bytes of the procedure, as shown in (b). If the 5th byte is in the middle of another CPU instruction, as is the case with `time`, Jockey overwrites up to the next instruction boundary (and fills the memory with `nop` as needed). In (c), Jockey also copies the original first 5 bytes (6 bytes for `time`) of the function to a newly allocated memory region so that Jockey can run the old implementation if necessary. (d) shows the pseudocode of the entry point for the new implementation of `time`. As we discuss further in Section 3.2, this code is dynamically generated so that Jockey can intercept system calls on a separate stack and avoid corrupting target memory. Finally, (e) shows `newtime`, Jockey's implementation of `time`. While recording, `newtime` calls the original implementation (c) and logs the returned value. While replaying, it simply supplies the value from the log without actually executing the system call.

One might wonder why `libjockey.so` does not just provide a new implementation of a system call with the same name—in fact, `LD_PRELOAD` is often used for that purpose. The reason is that doing so will miss system calls made inside `libc` or the dynamic linker—for example, a call to `read` made by `getc`. These internal calls are pre-resolved by the static linker (`ld`), and they cannot be overridden merely by redefining using `LD_PRELOAD`.

For task (2), Jockey rewrites all offending CPU instructions found in the target process. This is done in two ways, *slow* mode and *cached* mode. In *slow* mode, Jockey first reads the special file `/proc/N/maps` (N is the target process ID) that shows the virtual-memory mappings of the target process. It then reads the header of each mapped shared object file, discovers the locations of the text sections, and scans each text section. Jockey finds nondeterministic CPU instructions in the section (if any), and patches them. It also intercepts invocations of the `mmap` system call and does the same when a shared object file is newly loaded onto the process's address space.

Jockey needs to parse CPU instructions during steps (1) and (2), not a trivial task given x86's complex instruction encoding. It uses

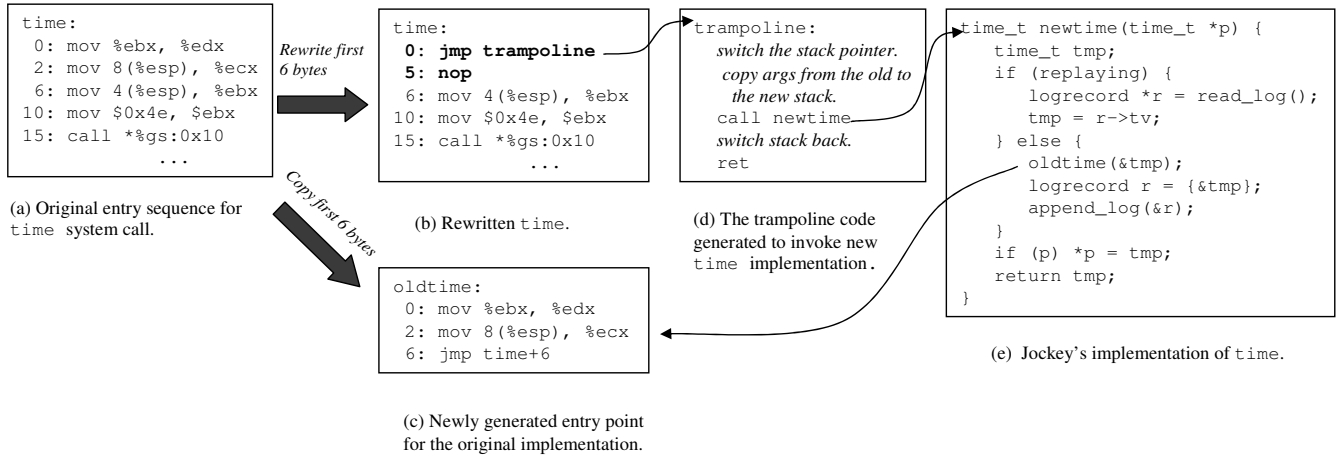


Figure 4: Recording and replaying `time`.

a pidgin table-based parser for common instructions and consults *libdisasm* [8], an open-source x86 disassembler library, for uncommon cases. A few tables that map opcodes and operands to their instruction length let us quickly parse more than 80% of all instruction occurrences.

Even using this technique, however, parsing all CPU instructions in a typical Linux program takes about 350 milliseconds on a 1.5GHz Pentium-M processor, which may be too slow for some users. To reduce the startup latency further, Jockey also employs cached-mode instruction patching. Here, after finishing the slow mode, Jockey writes the locations of nondeterministic instructions found for each shared object in file `~/.jockey-sig`. When the program starts the next time, it just reads `~/.jockey-sig` without scanning the process's virtual memory, unless the timestamp of the object file has changed.

Jockey's instruction-patching approach is simpler and faster than full-program binary translation, employed by ATOM [27] or Valgrind [24]. It needs to patch only a few bytes at the beginning of system calls and nondeterministic CPU instructions; the rest of the target program executes natively. Indeed, as we show in Section 5, Jockey's performance overhead is negligible for CPU-intensive programs.

3.2 Segregating resource usage

Jockey and the target application run as part of the same process and share all resources. Jockey must segregate the use of resources to prevent Jockey from unnecessarily changing the target's behavior, and to minimize the chance of a misbehaving target program breaking Jockey. This section discusses Jockey's treatment of three types of shared resources: heap, stack, and file descriptors.

3.2.1 Heap

Jockey cannot use standard libc functions, such as `malloc` or `sbrk`, to manage its internal data. Doing so increases the likelihood of a misbehaving target program breaking Jockey. Moreover, it changes the memory layout of the target process between record and replay. It would thus become impossible to replay invalid memory accesses correctly, e.g., accessing freed memory, which is one of the common programming errors.

Instead, Jockey stores all its internal data in a mmaped region at a fixed virtual address that is unlikely to be accessed accidentally by the target program. This address is by default set to `0x63000000`, but it could be changed via an environment variable if the tar-

get needs to access this address legitimately. We use an internal malloc-like library to carve the memory out to individual data structures and build a custom C++ STL memory allocator on top of it. Thus, the Jockey code has full access to STL features, including maps and dynamic vectors. This design has considerably simplified the development of Jockey.

One restriction is that Jockey cannot make internal calls to libc functions that use `malloc`. Examples include high-level I/O functions (`fopen`, `std::fstream`) and DNS resolvers (`gethostbyname`).

3.2.2 Stack

Jockey also segregates the use of stack. This is necessary to replay a program that improperly accesses data beyond the stack pointer (e.g., accessing an on-stack array with a negative index). Figure 4 (d) illustrates how this is done. In the first few instructions after it intercepts the call to `time`, Jockey saves the stack pointer to an internal variable, switches the stack to an internal buffer, copies the parameters to `time` (a 4-byte pointer) from the old to the new stack, and calls `newtime`. Once the new implementation returns, Jockey restores the stack pointer. This allows for deterministic replay of even a buggy program because Jockey never uses the target's stack.

This stack-switching process must be done without touching any CPU register other than the stack pointer. For this purpose, all the data structures involved here are allocated statically. This makes Jockey nonreentrant, but it is not an issue because Jockey does not support multi-threading.

3.2.3 File descriptors

Jockey must perform its own file accesses occasionally, for example, when opening a log file or dumping a checkpoint. Because Jockey and the target process share the same file-descriptor table, Jockey must ensure that its file operations do not alter the descriptor allocation scheme seen by the target. To achieve this goal, Jockey moves file descriptors it internally opens to a fixed range not likely to be used by the target (430~439).

An alternative approach would be to create an indirection table that maps file descriptors between the target program and the operating system kernel. Jockey would then intercept every system call that takes a file descriptor and translates it using the table. We chose our approach for two reasons. First, the former approach simplifies implementation, especially for system calls such as `select`. Second, the absence of descriptor indirection allows the user to inspect

```
% jockey --checkpointfrequency=30 \
  --retaincheckpoints=5 \
  -- httpd -X
... later ...
% jockey --restore=log/checkpoint-3 httpd
```

Figure 5: Taking automatic checkpoints of `httpd` (Apache) every 30 seconds. The `-X` option runs Apache in foreground. Option `--retaincheckpoints=5` causes only the last five checkpoints to be retained. The last line replays `httpd` from the third checkpoint.

the process’s state more transparently, for example, by using the system call tracer, such as `strace`.

Gdb (debugger) poses another problem. When starting the target process, gdb opens a few extra file descriptors in addition to the usual `stdin`, `stdout`, and `stderr`. Thus, if an execution is recorded under a normal shell and then replayed under gdb, the files opened by the target processes will be assigned different descriptors, which make replaying divergent.² We solve this problem by having Jockey open dummy files for descriptors 0 to 9 before starting the target program (it leaves descriptors inherited from the parent process untouched). Assuming that gdb opens at most 10 descriptors when it starts the target, we can ensure that the target has the same set of files opened upon record and replay.

3.3 Checkpointing

Jockey allows process state to be checkpointed automatically. Figure 5 shows an example. Checkpointing serves two purposes. First, it allows the developer to time-travel through the history of execution quickly. Second, it bounds log-space consumption, because log records older than the oldest checkpoint can be deleted from disk.

Following the technique pioneered by `libckpt` [20] and `Flashback` [26], Jockey first forks the target process. It then dumps the state of the child, while letting the parent continue running. Jockey reads the file `/proc/N/maps` (N is the process ID) to obtain the virtual memory mappings of the process and dumps only those sections that are mapped privately and read-write. To restore a checkpoint, for each section recorded in the checkpoint file, Jockey unmaps the memory region if it is already occupied, and either restores the contents from the checkpoint file or remaps the file.

We discuss two particular problems we faced, both related to dynamic linking.

3.3.1 Preventing brain damage to the dynamic linker

One of the challenges of checkpoint restoration is that Jockey needs to overwrite the memory that is potentially used by the restoration code itself. The process would crash if restoration is done naively. Here, two types of memory regions need to be taken care of: Jockey’s internal heap (Section 3.2) and the heap used by the dynamic linker. For example, Jockey must execute the `read` system call to load checkpoint contents. If the call to `read` happens to be the first ever made by the target application or Jockey, then the dynamic linker is invoked to resolve the symbol “`read`”, which involves modifying the linker’s heap.

Jockey handles its internal heap by excluding it from checkpointing, but the dynamic linker poses a particular challenge—we cannot know a priori where the memory used by the dynamic linker is

²In fact, this problem is not just specific to gdb. It happens whenever the target process inherits more than the standard number of file descriptors from the parent.

(the linker performs an anonymous `mmap` of its heap memory; all anonymous-mmapped sections look the same to Jockey). We resolve this issue by eagerly linking all libc functions that are called during snapshot restoration, by making dummy calls to functions such as `open` and `read` before it restores any checkpoint.

3.3.2 Exec shield

Exec-shield is a facility found in some Linux kernels (e.g., Red Hat, Fedora Core) to thwart buffer-overflow attacks [12]. One of its features is randomization of the loading addresses of shared-object files. This feature breaks Jockey because Jockey needs to keep data structures that are specific to the process’s memory layout. We currently require that this feature be disabled by doing the below on machine boot.

```
echo 0 >/proc/sys/kernel/exec-shield
```

3.4 Handling signals

Signals, especially those that happen asynchronously (for example, `SIGALRM`, `SIGINT`) present a special challenge, because they need to be delivered at exactly the same point in the execution during record and replay. We handle them in a way similar to [28]. Each signal delivery is first intercepted by Jockey. Jockey’s signal handler simply records the parameters to the signal (signal number and the CPU register values) and finishes. At the end of the Jockey’s handler for a system call or `rdtsc` CPU instruction, Jockey checks if a signal was intercepted in the past. If so, it logs the signal (so that it can be replayed) and calls the target-defined signal handler. This way, we convert asynchronous signals to synchronous upcalls that only happen immediately after a system call.

This technique may distort program behavior when the target program runs without issuing a system call (or executing non-deterministic CPU instructions) for a long period and receives signals in the meantime. However, our primary targets, I/O-oriented programs, usually do not suffer from this problem.

3.5 Reducing logging overhead for I/O system calls

Jockey employs two different types of logging techniques, depending on the types of system calls, to reduce the log-space overhead.

- For requests to regular files or directories, Jockey performs “undo” logging [9]. That is, for system calls that update a file, Jockey logs enough information to restore its contents *before* the modification. For example, when a `write` system call overwrites the mid-section of a file, Jockey logs the offset and the old contents of the section. Or, when `write` appends to the end of the file, Jockey just logs the old size of the file. In the replay mode, Jockey scans the log from the end to the start and restores the file contents. Read-only system calls (e.g., `read`) to regular files are simply executed directly on the file.
- For all other types of events—I/Os to sockets, pipes, fifos, devices, or `select`, `time`, or `rdtsc`—Jockey performs “redo” logging. Jockey logs the value produced by the event during recording, as illustrated in Figure 4 (e). During replay, Jockey just reads the values from the log without executing the actual system call.

System calls such as `read` and `write` can operate on both types of files. We intercept calls to functions that create file descriptors—e.g., `open`, `socket`, and `accept`—remember the type of each descriptor, and dispatch based on the descriptor type. File descrip-

tors inherited from the parent process (e.g., `stdin`) are always redologged.

Various studies have shown that majority of I/Os to regular files are reads, and that most of the write traffic is actually appends [18, 29]. For these common cases, our design allows Jockey to only log the type and the offset of the requests, not the actual contents. Thus, it drastically reduces the logging overhead for file I/O system calls.

The downside of the undo-based logging is that the user cannot modify the files accessed by the target program between record and replay. So far, we have not found this to be a significant burden.

3.6 Handling memory-mapped I/Os

Updates to memory mapped files are handled using user-space memory-protection mechanisms. Jockey intercepts calls to the `mmap` system call. For each file requested to be mapped read-write in a shared mode (i.e., `MAP_SHARED`),³ Jockey makes the mapped region read-only, and takes a page fault (`SIGSEGV` signal) after the first write access to each page in the region. In the `SIGSEGV` handler, Jockey logs the current page contents (Section 3.5), makes the page writable, and returns the control to the target. A similar approach is adopted by Flashback [26], albeit using a kernel extension.

These memory pages are made read-only again just before checkpointing, so that Jockey can restore the contents of the file at the moment of each checkpoint.

4. CONTROLLING JOCKEY

Jockey is designed to replay executions without requiring modification to the target source code. Sometimes, however, allowing the target program to change the behavior of Jockey could enable more efficient program execution or debugging. Jockey’s library-based design makes it easy to offer such control knobs for the target. This section introduces some of them.

4.1 Controlling the behavior of fork

By default, upon `fork`, Jockey continues recording only the parent and disables tracing the child. Procedure `jockey_set_fork_trace_mode(mode)` can be called by the target program to record only the child, or both (the behavior of `fork` can also be controlled via `JOCKEYRC` environment variable.) It can be used, for example, for daemon-type programs that fork to detach themselves from the parent process.

4.2 Target-specific function call interception

Procedure `jockey_redirect_calls(name, newproc, size)` is used to transfer the control to `newproc` whenever function `name` is called. Parameter `argsize` is the size of the on-stack parameters to the function. This function is implemented using instruction-patching service discussed in Section 3.1. This feature can be used, for example, to provide record/replay functionality for obscure `ioctl` commands.

4.3 User-defined invariant checker

Jockey allows an arbitrary object file to be linked into the target program during replay. Figure 6 shows an example. Let us assume that we ran `test2.c` under Jockey and found that procedure `bar` behaves anomalously when `i == 95999`. We could diagnose the bug by setting a breakpoint on `bar` in a debugger and waiting until it hits 95999 times, but Jockey offers a quicker alternative, as shown in Figure 7.

³Accesses to a private mapping (`MAP_PRIVATE`) need not be intercepted, because private mapping is essentially a heap memory with particular initial contents.

```
// test2.c
void bar(int i) { ... do something complex ... }
void main() {
    for (int i = 0; i < 100000; i++) bar(i);
}
```

Figure 6: A small program that executes procedure `bar` many times.

```
// check.c
#include <jockey/jockey.h>
void check_bar(int i) {
    if (i == 95999)
        jockey_breakpoint();
}
void init() {
    jockey_interpose_calls("bar",
                          check_bar, 4);
}
```

(a) A user-defined checker code that sets a breakpoint when procedure `bar` is called 95999 times.

```
% cc -c check.c -o check.o
% gdb test2
(gdb) b jockey_breakpoint
(gdb) run --jockey=replay=1;checker=check.o
```

(b) Running the user-defined checker.

Figure 7: Using a user-defined invariant checker.

The developer writes `check.c` in Figure 7 (a) to diagnose the problem. Procedure `init` is called automatically by Jockey when the object file is loaded into memory. `jockey_interpose_calls` is similar to `jockey_redirect_calls` (Section 4.2), but it returns the control to the original procedure after the callback returns. In this example, it will cause `check_bar` to be called just before `bar` is called.

The callback can be an arbitrary procedure as far as it does not modify the state of the program. It can set a conditional breakpoint as shown in this example, or it can check if some application-specific invariant holds. User-defined checkers offers several advantages over similar features offered by traditional debuggers, such as conditional breakpoints and watchpoints. First, it is more flexible because the checker can evaluate arbitrary application-specific expressions. They are also faster because they run at the native CPU speed.

The implementation of this feature is tricky, because we cannot use the dynamic linker to load the checker object file into the target process—doing so would alter the target’s memory usage (Section 3.3.1), which would cause the program execution to diverge between record and replay.

Jockey instead invokes the static linker, `ld`, to create a binary image at runtime. When Jockey tries to load a checker object, say `check.o`, it first discovers the memory addresses of all public symbols in the target process by invoking the `nm` command for each loaded shared object. Jockey then invokes the `ld` command on `check.o`. It passes the addresses of discovered public symbols and instructs `ld` to resolve symbols in `check.o` starting from a fixed virtual address unlikely to be accessed by the target program

Name	Run time			Log size	
	Native	Record	Replay	#bytes	#records
g++	1.33	1.51	1.49	73KB	80
xclock	N/A	180	0.4	80KB	4639
Emacs	N/A	210	5.81	1.4MB	20769
httpd	16.7	17.5	9.5	2.0MB	140180
FAB	33.7	44.1	31.1	34MB	887000

Table 1: The performance and log-space overheads of Jockey. Run times are in seconds. “Native” is the run-time without Jockey. “Record” and “Replay” show the runtime during recording and replaying, respectively.

(0x62000000 by default). The binary image created by `ld` is then read directly into memory at address 0x62000000 and executed.

5. EVALUATION

This section reports performance and space overheads of Jockey and discusses our experiences using Jockey to debug real-world programs.

5.1 Performance and log-space overheads

The evaluation was performed on a Fedora Core 3 Linux machine with a 1.5GHz Pentium-M CPU, 512MB of memory, and a 7200 rpm ATA disk drive. We ran a variety of programs under Jockey, as listed below. Stock binary executable files from the Fedora Core distribution were used, except for FAB.

g++: gcc 3.4.2 compiling a small C++ program that uses an STL map. The result shows the sum of the frontend (g++), backend (cc1plus), assembler (as), and linker (ld).

xclock: a digital clock for the X window system with a screen update every second.

Emacs: Emacs 21.3 running a program-development session, involving active typing, file reading, and saving.

httpd: Apache 2.0.52 serving 100000 HTTP GET requests for a static 0.5KB file. It was configured to run as a single, non-threaded process.

FAB: a four-process FAB cluster [23] serving 80000 random 1KB read and write iSCSI requests.

g++ is an example of a short-running, CPU-intensive program, which is not among Jockey’s primary targets. This example still shows that Jockey has a very low log-space overhead compared to approaches that involve memory-access logging [16], which could consume up to a few megabytes per second for logging. For g++, most of the slowdown is due to checkpointing that happens at the beginning of the execution (Section 3).

Xclock and Emacs are examples of interactive applications. Jockey exhibits reasonable log-space overheads for them. It is able to replay their execution extremely fast, because they need not wait for timeouts or user inputs during replay. This translates to more efficient debugging sessions.

Apache and FAB are examples of server programs. FAB represents the worst case for Jockey. Not only does FAB perform large amount of network I/O, it also overwrites existing files repeatedly, resulting in a large amount of logging traffic (Section 3.5). In comparison, Apache has a lower logging overhead because it only reads from HTML files and appends to access-log files.

5.2 Experiences

We have used Jockey regularly for FAB development. Our experiences have overall been positive. Jockey has been most useful when diagnosing bugs that happen after long stress or regression tests. Before Jockey, we were forced to recompile and reboot the system many times, each time with a slightly different set of “printf” statements, hoping that we would eventually reproduce and catch the error. Jockey allows us to reproduce the bug reliably as often as we wish. Fixing such bugs, however, is still difficult even with Jockey. The real cause of the bug often happens minutes before the bug exhibits, often on a different machine. The programmer needs to replay the execution of multiple processes repeatedly to locate the cause.

On the other hand, we have also found Jockey to be surprisingly effective in diagnosing bugs that exhibit quickly, e.g., while processing the first request from the client (indeed, most real-world bugs are of this type). Jockey cuts the debugging turn-around time by allowing the developer to replay a single process quickly instead of restarting the entire cluster.

Our experiences so far suggest that deterministic distributed replay system (Section 2.3) is not worth the complexity, at least for a system like FAB. The most important feature of a record/replay tool is the ability to replay quickly and reduce developers’ turn-around time. The whole-system replay does not improve this issue; it may actually increase the replay latency.

There are a few Jockey features that sound useful in theory, but have turned out to be not quite so in practice. First is user-defined invariant checking (Section 4.3). Debugging is an ad-hoc activity—writing and compiling a program every time one wants to debug is awkward. A debugger support, such as transparently compiling and loading a user-defined watchpoint to the program, would help. Another problem is the checker can only do only limited things—for example, it cannot intercept calls in the middle of function execution, nor can it inspect on-stack variables in the call chain.

Second, the concept of “time travel” using periodic automatic checkpoints (Section 3.3) has turned out to be powerful but cumbersome to use. The developer must manually restart the process every time he or she wants to switch to a different checkpoint. The developer can easily lose track of which part of the execution he or she is replaying. An extension to debuggers, such as automatic checkpoint scanning for detecting invariant violation [5, 31], would go a long way toward making this feature useful.

6. CONCLUSION

This paper described Jockey, a Linux tool for deterministic record/replay debugging. To achieve Jockey’s goals of safety and easy of use, it is implemented as a user-space library that runs as a part of the target process. It intercepts calls to nondeterministic system calls and CPU instructions, logs the effects of these operations during recording, and replays them from the log during replay. Jockey has a small performance and log-space overhead. Jockey has been extensively used to develop FAB.

7. REFERENCES

- [1] DAT collaborative. User-level direct access transport APIs (uDAPL), 2004. <http://www.datcollaborative.org/udapl.html>.
- [2] George W. Dunlap, Samuel T. King, Sukru Cinar, Murtaza Basrai, and Peter M. Chen. Revirt: Enabling intrusion analysis through virtual-machine logging and replay. In *5th Symp. on Op. Sys. Design and Impl. (OSDI)*, Boston, MA, USA, December 2002.

- [3] Timothy L. Harris. Dependable software needs pervasive debugging. In *10th ACM SIGOPS European Workshop*, Saint Emilion, France, September 2002.
- [4] Joel Huselius. Debugging parallel systems: A state of the art report. Technical Report 63, Dept. of CSE, Malardalen University, September 2002.
- [5] Samuel T. King, George W. Dunlap, and Peter M. Chen. Debugging operating systems with time-traveling virtual machines. In *USENIX Annual Tech. Conf.*, Anaheim, CA, USA, April 2005.
- [6] Lap Chung Lam. A survey of data breakpoint and reverse execution. SUNY Stony Brook RPE report, <http://www.ecsl.cs.sunysb.edu/tr/rpe12.ps.gz>, September 2001.
- [7] Bill Lewis. Debugging backwards in time. In *5th Workshop on Automated and Algorithmic Debugging (AADEBUG)*, Ghent, Belgium, September 2003.
- [8] libdisasm. Libdisasm: x86 disassembler library, 2004. <http://bastard.sourceforge.net/libdisasm.html>.
- [9] David E. Lowell and Peter M. Chen. Discount checking: Transparent, low-overhead recovery for general applications. Technical Report CSE-TR-410-99, University of Michigan, November 1998.
- [10] Charles E. McDowell and David P. Helmbold. Debugging concurrent programs. *ACM Computing Surveys*, 21(4):593–622, December 1989.
- [11] Michael S. Meier, Kevan L. Miller, Donald P. Pazel, Josyula R. Rao, and James R. Russell. Experiences with building distributed debuggers. In *SIGMETRICS Symposium on Parallel and Distributed Tools (SPDT)*, pages 70–79, Philadelphia, PA, USA, May 1996.
- [12] Ingo Molner. Exec shield, new Linux security feature. <http://people.redhat.com/mingo/exec-shield/ANNOUNCE-exec-shield>, 2004.
- [13] Robert H. B. Netzer. Optimal tracing and replay for debugging shared-memory parallel programs. In *ACM workshop on parallel and distributed debugging*, San Diego, CA, USA, May 1993.
- [14] Robert H. B. Netzer and Barton P. Miller. Optimal tracing and replay for debugging message-passing parallel programs. In *Supercomputing*, Mineapolis, MN, USA, November 1992.
- [15] Robert H. B. Netzer, Sairam Subramanian, and Jian Xu. Critical-path-based message logging for incremental replay of message-passing programs. In *14th Int. Conf. on Dist. Comp. Sys. (ICDCS)*, pages 404–413, Poznan, Poland, June 1994.
- [16] Robert H. B. Netzer and Mark H. Weaver. Optimal tracing and incremental reexecution for debugging long-running programs. In *SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, Orlando, FL, USA, June 1994. Also available as Brown University Technical Report CS-94-11.
- [17] Oliver Oppitz. A particular bug trap: Execution replay using virtual machines. In *5th Workshop on Automated and Algorithmic Debugging (AADEBUG)*, Ghent, Belgium, September 2003.
- [18] John K. Ousterhout, Herv Da Costa, David Harrison, John A. Kunze, Michael D. Kupfer, and James G. Thompson. A trace-driven analysis of the UNIX 4.2 BSD file system. In *10th Symp. on Op. Sys. Principles (SOSP)*, pages 15–24, Orcas Island, WA, USA, December 1985.
- [19] Douglas Z. Pan and Mark A. Linton. Supporting reverse execution of parallel programs. In *ACM workshop on parallel and distributed debugging*, Madison, WI, USA, May 1988.
- [20] James S. Plank, Micah Beck, Gerry Kingsley, and Kai Li. Libckpt: Transparent checkpointing under UNIX. In *USENIX Winter Tech. Conf.*, New Orleans, LA, USA, January 1995.
- [21] Michiel Ronsse, Koen De Bosschere, Mark Christiaens, Jacques Chassin de Kergommeaux, and Dieter Kranzlmüller. Record/replay for non-deterministic program executions. *Comm. of the ACM (CACM)*, 46(9), September 2003.
- [22] Michiel Ronsse, Koen De Bosschere, and Jacques Chassin de Kergommeaux. Execution replay and debugging. In *4th Workshop on Automated and Algorithmic Debugging (AADEBUG)*, Munich, Germany, August 2000.
- [23] Yasushi Saito, Svend Frølund, Alistair Veitch, Arif Merchant, and Susan Spence. FAB: Building distributed enterprise disk arrays from commodity components. In *11th Int. Conf. on Arch. Support for Prog. Lang. and Op. Sys. (ASPLOS-XI)*, Boston, MA, USA, October 2004.
- [24] Julian Seward et al. Valgrind: A GPL'd system for debugging and profiling x86-linux programs. <http://valgrind.kde.org/>, 2004.
- [25] Michael W. Shapiro. RDB: A system for incremental replay debugging. Master's thesis, Dept of. Computer Science, Brown University, 1997.
- [26] Sudarshan M. Srinivasan, Srikanth Kandula, Christopher R. Andrews, and Yuanyuan Zhou. Flashback: A lightweight extension for rollback and deterministic replay for software debugging. In *USENIX Annual Tech. Conf.*, Boston, MA, USA, June 2004.
- [27] Amitabh Srivastava and Alan Eustace. ATOM: a system for building customized program analysis tools. In *SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 196–205, Orlando, FL, USA, June 1994.
- [28] Daniel Stodolsky, Brian N. Bershad, and J. Bradley Chen. Fast Interrupt Priority Management in Operating System Kernels. *Usenix Workshop on Microkernels*, pages 105–110, September 1993.
- [29] Werner Vogels. File system usage in Windows NT 4.0. In *17th Symp. on Op. Sys. Principles (SOSP)*, pages 93–109, Kiawah Island, SC, USA, December 1999.
- [30] Rob von Behren, Jeremy Condit, Feng Zhou, George C. Necula, and Eric Brewer. Capriccio: Scalable threads for Internet services. In *19th Symp. on Op. Sys. Principles (SOSP)*, Bolton Landing, NY, USA, October 2003.
- [31] Andrew Whitaker, Richard S. Cox, and Steven D. Gribble. Configuration debugging as search: Finding the needle in the haystack. In *6th Symp. on Op. Sys. Design and Impl. (OSDI)*, San Francisco, CA, USA, December 2004.
- [32] Larry D. Wittie. Debugging distributed C programs by real time replay. In *ACM workshop on parallel and distributed debugging*, pages 57–67, Madison, WI, USA, May 1988.