

# COZ: Finding Code that Counts with Causal Profiling

Charlie Curtsinger    Emery D. Berger

School of Computer Science  
University of Massachusetts Amherst  
{charlie, emery}@cs.umass.edu

## Abstract

Improving performance is a central concern for software developers. To locate optimization opportunities, developers rely on software profilers. However, these profilers only report where programs spent their time: optimizing that code may have no impact on performance. Past profilers thus both waste developer time and make it difficult for them to uncover significant optimization opportunities.

This paper introduces *causal profiling*. Unlike past profiling approaches, causal profiling indicates exactly where programmers should focus their optimization efforts, and quantifies their potential impact. Causal profiling works by running *performance experiments* during program execution. Each experiment calculates the impact of any potential optimization by *virtually speeding up* code: inserting pauses that slow down all other code running concurrently. The key insight is that this slowdown has the same *relative* effect as running that line faster, thus “virtually” speeding it up.

We present COZ, a causal profiler, and evaluate it on a range of highly-tuned applications: Memcached, SQLite, and the PARSEC benchmark suite. COZ identifies previously-unknown optimization opportunities that are both significant and targeted. Guided by COZ, we improve the performance of Memcached by 9%, SQLite by 25%, and accelerate six PARSEC applications by as much as 68%; in most cases, these optimizations involve modifying under 10 lines of code.

## 1. Introduction

Improving performance is a central concern for software developers. While compiler optimizations are of some assistance, they often do not have enough of an impact on performance to meet programmers’ demands [11]. Programmers seeking to increase the throughput or responsiveness of their applications thus must resort to manual performance tuning.

Since manually inspecting an entire program to find optimization opportunities is impractical, developers use profilers. Conventional profilers rank code by its contribution to total execution time. Prominent examples include oprofile, perf, and gprof [19, 28, 30].

Unfortunately, even when a profiler accurately reports where a program is spending its time, this information can

lead programmers astray. Code that runs for a long time is not necessarily a good choice for optimization. For example, optimizing code that draws a loading animation during a file download will not make the program run any faster, even though this code runs just as long as the file download.

This phenomenon is not limited to I/O operations. Figure 1 shows a simple program that illustrates the shortcomings of existing profilers, along with its gprof profile in Figure 2a. This program spawns two threads, which invoke functions *a* and *b* respectively. Most profilers will report that these functions comprise roughly half of the total execution time. Other profilers may report that the *a* function is on the critical path, or that the main thread spends roughly equal time waiting for *a\_thread* and *b\_thread* [24]. While accurate, all of this information is potentially misleading. Optimizing *a* away entirely will only speed up the program by 4.5% because *b* becomes the new critical path.

Existing profilers do not report the potential impact of optimizations; developers are left to make these predictions given their understanding of the program. While these predictions may be easy for programs as simple as the one in Figure 1, accurately predicting the performance impact of a proposed optimization is nearly impossible for programmers attempting to optimize large applications.

This paper introduces *causal profiling*, an approach that accurately and precisely indicates where programmers should

example.cpp

```
1 void a() { // ~6.7 seconds
2   for(volatile size_t x=0; x<2000000000; x++) {}
3 }
4 void b() { // ~6.4 seconds
5   for(volatile size_t y=0; y<1900000000; y++) {}
6 }
7 int main() {
8   // Spawn both threads and wait for them.
9   thread a_thread(a), b_thread(b);
10  a_thread.join(); b_thread.join();
11 }
```

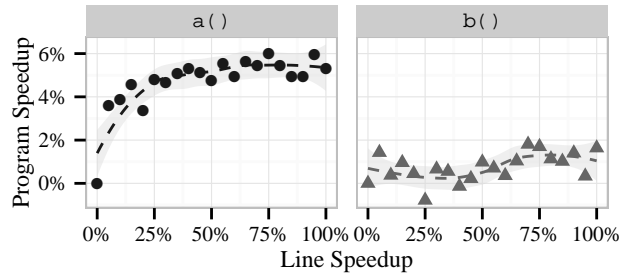
**Figure 1:** A simple multithreaded program that illustrates the shortcomings of existing profilers. Optimizing *a* will improve performance by no more than 4.5%, while optimizing *b* would have no effect on performance.

**gprof Profile For example.cpp**

% time	cumulative seconds	self seconds	calls	self Ts/call	total Ts/call	name
55.20	7.20	7.20	1			a ()
45.19	13.09	5.89	1			b ()
% time	self	children	called	name		
55.0	7.20	0.00		<spontaneous> a ()		
45.0	5.89	0.00		<spontaneous> b ()		

(a) gprof profile for example.cpp

**Causal Profile For example.cpp**



(b) Causal profile for example.cpp

**Figure 2:** The gprof and causal profiles for the code in Figure 1. In the causal profile, the y-axis shows the program speedup that would be achieved by speeding up each line of code by the percentage on the x-axis. The gray area shows standard error. While gprof reports that a and b comprise similar fractions of total runtime, it does not indicate that optimizing a will improve performance by at most 4.5%, and optimizing b would have no effect. The causal profile predicts both outcomes within 0.5%.

focus their optimization efforts, and quantifies their potential impact. Figure 2b shows the results of running COZ, our prototype causal profiler. This profile plots the hypothetical speedup of a line of code (x-axis) versus its impact on execution time (y-axis). The graph correctly shows that optimizing either a or b in isolation would have little impact on execution time.

A causal profiler conducts a series of *performance experiments* to empirically observe the impact of a potential optimization. Of course, it is not possible to automatically speedup any line of code by an arbitrary amount. Instead, during a performance experiment, the causal profiler uses the novel technique of *virtual speedups* to mimic the effect of optimizing a specific line of code by a specific amount.

Virtual speedup works by inserting pauses that slow down all code running at the same time as the line under examination. The key insight is that this slowdown has the same *relative* effect as running that line faster, thus “virtually” speeding it up. Figure 3 illustrates the relative equivalence between actual and virtual speedups: after accounting for delays, both have the same impact.

Each performance experiment measures the impact of some amount of virtual speedup to a single line. By sampling over the range of virtual speedup from between 0% (no change) and 100% (the line is completely eliminated), causal profiling can calculate the impact of *any* potential optimization on overall performance.

Causal profiling further departs from traditional profiling by making it possible to view the effect of optimizations on *throughput* and *latency*. To profile throughput, developers specify a *progress point*, indicating a line in the code that corresponds to the end of a unit of work. For example, a progress point could be the point at which a transaction concludes, when a web page finishes rendering, or when a query completes. A causal profiler then measures the rate of visits to each progress point to determine any potential optimization’s effect on throughput.

To profile latency, programmers place *two* progress points that correspond to the start and end of an event of interest, such as when a transaction begins and completes. A causal profiler then reports the effect of potential optimizations on the average latency between those two progress points.

We demonstrate causal profiling with COZ, a prototype causal profiler that works with Linux x86-64 binaries. We show that COZ imposes low execution time overhead (mean: 17%, min: 0.1%, max: 65%), making it substantially faster than gprof (up to 6× overhead).

We show that causal profiling accurately predicts optimization opportunities, and that it is effective at guiding optimization efforts. We apply COZ to Memcached, SQLite, and the extensively-studied PARSEC benchmark suite. Guided by COZ’s output, we optimized the performance of Memcached by 9%, SQLite by 25%, and six PARSEC applications by as much as 68%. These optimizations typically involved modifying under 10 lines of code. When possible to accurately measure the size of our optimizations on the line(s) identified by COZ, we compare the observed performance improvements to COZ’s predictions: in each case, we find that the real effect of our optimization matched COZ’s prediction.

## Contributions

This paper makes the following contributions:

1. It presents **causal profiling**, which identifies code where optimizations will have the largest impact. Using *virtual speedups* and *progress points*, causal profiling directly measures the effect of potential optimizations on both throughput and latency (§2).
2. It presents **COZ**, a causal profiler that works on unmodified Linux binaries. It describes COZ’s implementation (§3), and demonstrates its efficiency and effectiveness at identifying optimization opportunities (§4).

## 2. Causal Profiling Overview

Causal profiling relies on several key ideas to provide developers with actionable profiles. *Virtual speedups* let a causal profiler automatically create the effect of optimizing any fragment of code. *Progress points* let the profiler measure a program’s performance repeatedly during one run. *Performance experiments* apply a virtual speedup and measure the resulting effect on performance. Repeated performance experiments enable a causal profiler to identify fragments of code where optimizations will have the greatest impact. This section provides a detailed description of these key concepts, and describes the workflow of COZ, our prototype causal profiler.

**Virtual speedups.** A virtual speedup uses delays to create the *effect* of optimizing a fragment of code. Each time a selected fragment is executed, all other threads are briefly paused. The longer the pause, the larger the relative speedup. At the end of an execution, causal profiling subtracts the total pause time from runtime to determine the effective execution time. This technique is illustrated in Figure 3.

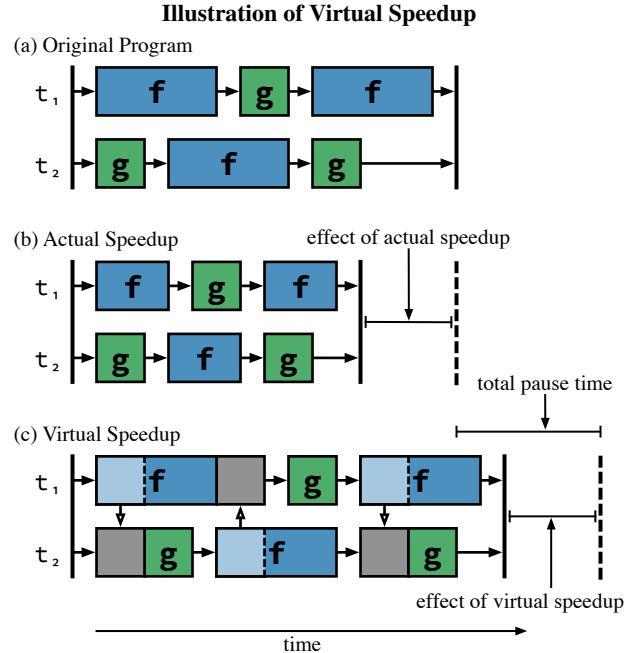
**Progress points.** A causal profiler uses progress points to measure program performance during execution. Developers must place progress points at a source location where some useful work has been completed. These points let a causal profiler conduct many performance experiments during a single run. Additionally, progress points enable measurement of both latency and throughput, and enable profiling of long-running applications where end-to-end execution time is meaningless.

**Performance experiments.** A causal profiler runs many performance experiments during a program’s execution. For each experiment, the profiler randomly selects a fragment of code to virtually speed up for the duration of the experiment. Meanwhile, the profiler measures the rate of visits to one or more progress points. Each performance experiment establishes the impact of optimizing a particular code fragment by a specific amount. Given a sufficient number of experiments, the profiler can identify which fragments will yield the largest performance gains if optimized.

A causal profiler can also identify contention, which appears as a downward sloping line on a causal profile graph. A negative slope indicates that optimizing the code fragment will hurt application performance. We find and address several instances of contention in our case studies in Section 4.

### 2.1 Causal Profiling Workflow

To demonstrate the effectiveness of causal profiling, we have implemented COZ, a prototype causal profiler. COZ implements all of the key components of a causal profiler: virtual speedups, progress points, and performance experiments. COZ identifies optimization opportunities at the granularity of source lines, but our technique can easily support any type of code fragment. We describe COZ’s profiling workflow in detail below.



**Figure 3:** An illustration of virtual speedup: (a) shows the original execution of two threads running functions  $f$  and  $g$ ; (b) shows the effect of a *actually* speeding up  $f$  by 40%; (c) shows the effect of *virtually* speeding up  $f$  by 40% by pausing the other thread each time  $f$  runs. Each inserted pause (dark gray) is equal to the size of the speedup—40% of  $f$ ’s execution time (light blue). The runtime of (c) is longer than (b) by the total pause time. Adjusting the baseline runtime from (a) by the total delay time lets us measure the virtual speedup size, which matches the effect of the actual speedup.

**Profiler startup.** A user invokes COZ using a command of the form `coz run --- <program> <args>`. At the beginning of the program’s execution, COZ collects debug information for the executable and all loaded libraries. Users may specify file and binary scope, which restricts COZ’s experiments to speedups in the specified files. By default, COZ will consider speedups in any source file from the main executable. COZ builds a map from instructions to source lines using the program’s debug information and the specified scope. Once the source map is constructed, COZ creates a profiler thread and resumes normal execution.

**Experiment initialization.** COZ’s profiler thread begins an experiment by selecting a line to virtually speed up, and a randomly-chosen percent speedup. Both parameters must be selected randomly; any systematic method of exploring lines or speedups could lead to systematic bias in profile results. Once a line and speedup have been selected, the profiler thread saves the number of visits to each progress point and begins the experiment.

**Applying a virtual speedup.** Every time the profiled program creates a thread, COZ begins sampling the instruction pointer from this thread. COZ processes samples within each thread to implement a sampling version of virtual speedups.

In Section 3.4, we show the equivalence between the virtual speedup mechanism described above and the sampling approach implemented in COZ. Every time a sample is available, a thread checks whether the sample falls in the line of code selected for virtual speedup. If so, it forces other threads to pause. This process continues until the profiler thread indicates that the experiment has completed.

**Ending an experiment.** COZ ends the experiment after a pre-determined time has elapsed. If there were too few visits to progress points during the experiment—five is the default minimum—COZ doubles the experiment time for the rest of the execution. Once the experiment has completed, the profiler thread logs the results of the experiment, including the effective duration of the experiment (runtime minus the total inserted delay), the selected line and speedup, and the number of visits to all progress points. Before beginning the next experiment, COZ will pause for a brief cooldown period to allow any remaining samples to be processed.

### 3. Implementation

The current implementation of COZ profiles Linux x86-64 executable binaries. To map program addresses to source lines, COZ uses DWARF debugging information. As long as debug information is available in a separate file, COZ can profile optimized and stripped executables. Sampling is implemented using the `perf_event` API.

#### 3.1 Profiler Startup

The COZ profiling code is inserted into a process using the `LD_PRELOAD` environment variable. This allows COZ to intercept library calls from the program, including the `libc_start_main` function, which runs before `main` and all global constructors. Before the program's normal execution begins, COZ collects the names and locations of all loaded executables by reading `/proc/self/maps`. COZ records the loaded address and path to each in-scope executable for later processing.

For all in-scope executables and libraries, COZ locates DWARF debug information for the program's main executable and libraries [15]. By default, the scope includes all source files from the main executable, but alternate source locations and libraries can be specified on the command line. If any debug information has been stripped, COZ uses the same procedure as Gdb to search standard system paths for separate debugging information [16]. Note that debug information is available even for optimized code, and most Linux distributions offer packages that include this information for common libraries.

COZ uses DWARF line tables to build a map from instruction pointer ranges to source lines. The DWARF format also includes both caller and callee information for inlined procedures. Special handling is required when an in-scope callsite is replaced by an inlined function that is *not* in scope. The inlined function's address range is assigned to the caller's

source location in the source map. This approach mirrors the process by which COZ attributes out-of-scope samples to callsites during execution (see the discussion of sample attribution, below).

**Enabling Sampling.** Before calling the program's `main` function, COZ opens a `perf_event` file to start sampling in the main thread. COZ invokes the `perf_event_open` system call to track high precision timer events via a memory-mapped file. COZ samples each thread individually using the high precision timer event, and collects instruction pointers and the user-space `callchain` in each sample.

**Sample Attribution.** Samples are attributed to source lines using the source map constructed at startup. When a sample does not fall in any in-scope source line, the profiler walks the sampled `callchain` to find the first in-scope address. This process has the effect of attributing all out-of-scope execution to the last in-scope callsite responsible. For example, a program may call `printf`, which calls `vfprintf`, which in turn calls `strlen`. Any samples collected during this chain of calls will be attributed to the source line that issues the original `printf` call.

#### 3.2 Experiment Initialization

A single profiler thread, created during program initialization, coordinates performance experiments. Before a performance experiment can begin, a line must be selected for virtual speedup. When an experiment is not running, each program thread will set the `next_line` atomic variable to its most recent sample. The profiler thread spins until this variable contains a non-null value.

Once the profiler receives a valid line from one of the program's threads, it chooses a random virtual speedup between 0% and 100%, in multiples of 5%. For any given virtual speedup, the effect on program performance is  $1 - \frac{p_s}{p_0}$ , where  $p_0$  is the period between progress point visits with no virtual speedup, and  $p_s$  is the same period measured with some virtual speedup  $s$ . Because  $p_0$  is required to compute program speedup for every  $p_s$ , a virtual speedup of 0 is selected with 50% probability. The remaining 50% is distributed evenly over the other virtual speedup amounts.

Virtual speedups must be selected randomly to prevent bias in the results of performance experiments. A seemingly reasonable (but invalid) approach would be to begin conducting performance experiments with small virtual speedups, gradually increasing the speedup until it no longer has an effect on program performance. However, this approach may both over- and under-state the impact of optimizing a particular line if its impact varies over time.

For example, a line that has no performance impact during a program's initialization would not be measured later in execution, when optimizing it could have significant performance benefit. Conversely, a line that only affects performance during initialization would have exaggerated performance impact unless future experiments re-evaluate virtual speedup values

for this line during normal execution. Any systematic approach to exploring the space of virtual speedup values could potentially lead to systematic bias in the profile output.

Once a line and virtual speedup have been selected, COZ saves the current values of all progress point counters and begins the performance experiment.

### 3.3 Running a Performance Experiment

Once a performance experiment has started, each of the program’s threads processes samples and inserts delays to perform virtual speedups. After the pre-determined experiment time has elapsed, the profiler thread logs the end of the experiment, including the current time, the number and size of delays inserted for virtual speedup, the running count of samples in the selected line, and the values for all progress point counters. After a performance experiment has finished, COZ waits at least 10ms before starting another experiment. This pause ensures that delays and samples processed by threads around the end of the experiment are not accidentally attributed to the next experiment, which would bias results.

### 3.4 Virtual Speedups

COZ uses delays to create the effect of optimizing the selected line. Every time one thread executes this line, all other threads must pause. The length of the pause determines the amount of virtual speedup; pausing other threads for half the selected line’s runtime has the effect of optimizing the line by 50%.

**Implementing Virtual Speedup.** Tracking every visit to the selected line would incur significant performance overhead and distort the program’s execution. Instead, COZ uses sampling to implement virtual speedups accurately and efficiently, delaying proportionally to the time spent in the selected line. This lets COZ virtually speed up the line by a specific percent, even though the number of visits to the line is unknown.

The expected number of samples in the selected line,  $s$ , is

$$\mathbb{E}[s] = \frac{n \cdot t}{P} \quad (1)$$

where  $P$  is the period of time between samples,  $t$  is the time required to run the selected line once, and  $n$  is the number of times the selected line is executed.

In our original model of virtual speedups, delaying other threads by time  $d$  each time the selected line is executed has the effect of shortening this line’s runtime by  $d$ . With sampling, only some executions of the selected line will result in delays. The effective runtime of the selected line *when sampled* is  $t - d$ , while executions of the selected line that are not sampled simply take time  $t$ . The average effective time to run the selected line is

$$t' = \frac{(n - s) \cdot t + s \cdot (t - d)}{n}.$$

Using (1), this reduces to

$$t' = \frac{n \cdot t \cdot (1 - \frac{t}{P}) + \frac{n \cdot t}{P} \cdot (t - d)}{n} = t \cdot (1 - \frac{d}{P}) \quad (2)$$

The percent difference between  $t$  and  $t'$ , the amount of virtual speedup, is simply

$$\Delta t = 1 - \frac{t'}{t} = \frac{d}{P}.$$

This result lets COZ virtually speed up selected lines by a specific amount without instrumentation. Inserting a delay that is half the sampling period will virtually speed up the selected line by 50%.

**Pausing Other Threads.** When one thread receives a sample in the line selected for virtual speedup, all other threads must pause. COZ triggers these pauses using two counters: a shared global counter, and per-thread local counters. These counters are used to pause threads without using expensive POSIX signals. The global counter stores the number of pauses each thread should execute, while per-thread local counters track the number of pauses each thread has executed so far. To pause all other threads, a thread increments both counters. Every thread checks the counters after each sample. If a thread’s local delay count is less than the global delay count, it must pause and increment its local counter. Each thread checks its counter against the global count and inserts any required delays immediately after processing samples.

**Ensuring Accurate Timing.** COZ uses the `nanosleep` POSIX function to insert delays. This function only guarantees that the thread will pause for *at least* the requested time, but the pause may be longer than requested. COZ tracks any excess pause time, which is subtracted from future pauses.

**Thread Creation.** To start sampling and adjust delays, COZ interposes on the `pthread_create` function. COZ first initiates `perf_event` sampling in the new thread. It then copies the parent thread’s local delay count, propagating any delays: any previously inserted delays to the parent thread also delayed the creation of the new thread.

**Thread Sampling and Delay Accounting.** COZ only interrupts a thread to process samples if the thread is running. If the thread is blocked on I/O, sample processing and delays will be performed after the blocking call returns. For blocking I/O, this is the desired behavior—inserting pauses during a file read would have no effect on the time it takes to complete the read. However, threads can also block on other threads, which complicates delay insertion.

Consider a program with two threads: thread A is currently holding a mutex, and thread B is waiting to acquire the mutex. If thread B is spinning on the mutex, delaying that thread will not necessarily have any effect on how long it waits. Unlike with blocking I/O, this is actually the desired behavior: thread A will have inserted these delays, which delays the time that thread A unlocks the mutex and B can proceed. But, if thread B is suspended while waiting for the mutex, these delays would be inserted when the thread wakes. Any delays required while the thread is blocked could be inserted twice:

Potentially <i>blocking</i> calls	
<code>pthread_mutex_lock</code>	lock a mutex
<code>pthread_cond_wait</code>	wait on a condition variable
<code>pthread_barrier_wait</code>	wait at a barrier
<code>pthread_join</code>	wait for a thread to complete
<code>sigwait</code>	wait for a signal
<code>sigwaitinfo</code>	wait for a signal
<code>sigtimedwait</code>	wait for a signal (with timeout)
<code>sigsuspend</code>	wait for a signal

**Table 1:** COZ intercepts POSIX functions that could block waiting for a thread, instrumenting them to update delay counts before and after blocking.

once by thread A before unlocking the mutex, and then again in thread B after acquiring the mutex.

To correct this behavior, blocked threads must inherit the delay count from the thread that unblocks them. This causal propagation ensures that any delays inserted before unblocking the thread would not be inserted again in the waking thread. For simplicity, COZ forces threads to execute all required delays before performing an operation that could wake a blocked thread. These operations include the POSIX calls given in Table 2.

When a thread is unblocked by one of the listed functions, COZ guarantees that all required delays have been inserted. The thread can simply skip any delays that were incurred while it was blocked. Before executing a function that may block on thread communication, a thread saves both the local and global delay counts. When the thread wakes, it sets its local delay count to the saved delay count, plus any global delays incurred since the call. This accounting is correct whether the thread was suspended or simply spun on the synchronization primitive. Table 1 lists the functions that require this additional handling.

### Optimization: Minimizing Delays

If every thread executes the selected line, forcing each thread to delay `num_threads - 1` times unnecessarily slows execution. If all but one thread executes the selected line, only that thread needs to pause. The invariant that must be preserved is the following: for each thread, the number of pauses plus the number of samples in the selected line must equal the global delay count. When a sample falls in the selected line, COZ increments only the local delay count. If the local delay count is still less than the global delay count after processing all available samples, COZ inserts pauses. If the local delay count is larger than global delay count, the thread increases the global delay count.

### 3.5 Progress Points

COZ supports three different mechanisms for progress points: *source-level*, *breakpoint*, and *sampled*.

**Source-Level Progress Points.** Source-level progress points are the only progress points that require program modification. To indicate a source-level progress point, a developer simply

Potentially <i>unblocking</i> calls	
<code>pthread_mutex_unlock</code>	unlock a mutex
<code>pthread_cond_signal</code>	wake one waiter on a c.v.
<code>pthread_cond_broadcast</code>	wake all waiters on c.v.
<code>pthread_barrier_wait</code>	wait at a barrier
<code>pthread_kill</code>	send signal to a thread
<code>pthread_exit</code>	terminate this thread

**Table 2:** COZ intercepts POSIX functions that could wake a blocked thread. To ensure correctness of virtual speedups, COZ forces threads to execute any unconsumed delays before invoking any of these functions and potentially waking another thread.

inserts the `CAUSAL_PROGRESS` macro in the program’s source code at the appropriate location.

**Breakpoint Progress Points.** Breakpoint progress points are specified at the command line. COZ uses the `perf_event` API to set a breakpoint at the first instruction in a line specified in the profiler arguments.

**Sampled Progress Points.** Like breakpoint progress points, sampled progress points are specified at the command line. However, unlike source-level and breakpoint progress points, sampled progress points do not keep a count of the number of visits to the progress point. Instead, sampled progress points count the number of samples that fall within the specified line. As with virtual speedups, the percent change in visits to a sampled progress point can be computed even when the raw counts are unknown.

**Measuring Latency.** Source-level and breakpoint progress points can also be used to measure the impact of an optimization on latency rather than throughput. To measure latency, a developer must specify two progress points: one at the start of some operation, and the other at the end. The rate of visits to the starting progress point measures the arrival rate, and the difference between the counts at the start and end points tells us how many requests are currently in progress. By denoting  $L$  as the number of requests in progress and  $\lambda$  as the arrival rate, we can solve for the average latency  $W$  via Little’s Law, which holds for nearly any queuing system:  $L = \lambda W$  [31]. Rewriting Little’s Law, we then compute the average latency as  $L/\lambda$ .

Little’s Law holds under a wide variety of circumstances, and is independent of the distributions of the arrival rate and service time. The key requirement is that Little’s Law only holds when the system is *stable*: the arrival rate cannot exceed the service rate. Note that all usable systems are stable: if a system is unstable, its latency will grow without bound since the system will not be able to keep up with arrivals.

### 3.6 Adjusting for Phases

COZ randomly selects a recently-executed line of code at the start of each performance experiment. This increases the likelihood that experiments will yield useful information—a virtual speedup would have no effect on lines that never run—but could bias results for programs with phases.

If a program runs in phases, optimizing a line will not have any effect on progress rate during periods when the line is not being run. However, COZ will not run performance experiments for the line during these periods because only currently-executing lines are selected. If left uncorrected, this bias would lead COZ to overstate the effect of optimizing lines that run in phases.

To eliminate this bias, we break the program’s execution into two logical phases: phase A, during which the selected line runs, and phase B, when it does not. These phases need not be contiguous. The total runtime  $T = t_A + t_B$  is the sum of the durations of the two phases. The average progress rate during the entire execution is:

$$P = \frac{T}{N} = \frac{t_A + t_B}{N}. \quad (3)$$

COZ collects samples during the entire execution, recording the number of samples in each line. We define  $s$  to be the number of samples in the selected line, of which  $s_{obs}$  occur during a performance experiment with duration  $t_{obs}$ . The expected number of samples during the experiment is:

$$\mathbb{E}[s_{obs}] = s \cdot \frac{t_{obs}}{t_A}, \quad \text{therefore} \quad t_A \approx s \cdot \frac{t_{obs}}{s_{obs}}. \quad (4)$$

COZ measures the effect of a virtual speedup during phase A,

$$\Delta p_A = \frac{p_A - p_A'}{p_A}$$

where  $p_A'$  and  $p_A$  are the average progress periods with and without a virtual speedup; this can be rewritten as:

$$\Delta p_A = \frac{\frac{t_A}{n_A} - \frac{t_A'}{n_A'}}{\frac{t_A}{n_A}} = \frac{t_A - t_A'}{t_A} \quad (5)$$

where  $n_A$  is the number of progress point visits during phase A. Using (3), the new value for  $P$  with the virtual speedup is

$$P' = \frac{t_A' + t_B}{N}$$

and the percent change in  $P$  is

$$\Delta P = \frac{P - P'}{P} = \frac{\frac{t_A + t_B}{N} - \frac{t_A' + t_B}{N}}{\frac{t_A + t_B}{N}} = \frac{t_A - t_A'}{T}.$$

Finally, using (4) and (5),

$$\Delta P = \Delta p_A \frac{t_A}{T} \approx \Delta p_A \cdot \frac{t_{obs}}{s_{obs}} \cdot \frac{s}{T}. \quad (6)$$

COZ multiplies all measured speedups,  $\Delta p_A$ , by the correction factor  $\frac{t_{obs}}{s_{obs}} \cdot \frac{s}{T}$  in its final report.

## 4. Evaluation

Our evaluation answers the following questions: (1) Does causal profiling enable effective performance tuning? (2) Are COZ’s performance predictions accurate? (3) Is COZ’s overhead low enough to be practical?

**Summary of Optimization Results**

Application	Speedup	Diff Size	LOC
blackscholes	2.56% ± 0.41%	-61, +4	342
dedup	8.95% ± 0.27%	-3, +3	2,570
ferret	21.27% ± 0.17%	-4, +4	5,937
fluidanimate	37.5% ± 0.56%	-1, +0	1,015
streamcluster	68.4% ± 1.12%	-1, +0	1,779
swaptions	15.8% ± 1.10%	-10, +16	970
Memcached	9.39% ± 0.95%	-6, +2	10,475
SQLite	25.60% ± 1.00%	-7, +7	92,635

**Table 3:** All benchmarks were run ten times before and after optimization. Standard error for speedup was computed using Efron’s bootstrap method, where speedup is defined as  $\frac{t_0 - t_{opt}}{t_0}$ . All speedups are statistically significant at the 99.9% confidence level ( $\alpha = 0.001$ ) using the one-tailed Mann-Whitney U test, which does not rely on any assumptions about the distribution of execution times. Lines of code do not include blank or comment-only lines.

### 4.1 Experimental Setup

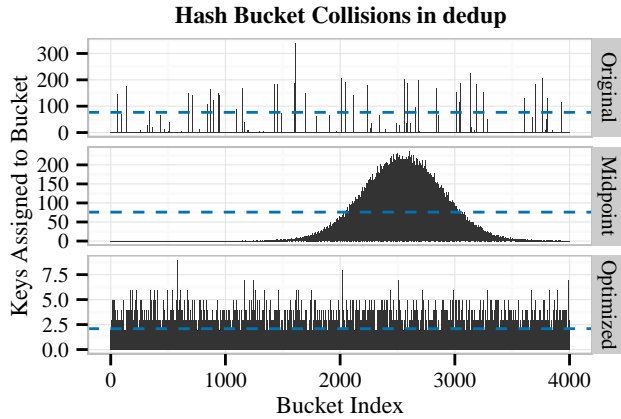
We perform all experiments on a 64 core, four socket AMD Opteron machine with 60GB of memory, running Linux 3.14 with no modifications. All applications are compiled using GCC version 4.9.1 at the `-O3` optimization level and debug information generated with `-g`. We disable frame pointer elimination with the `-fno-omit-frame-pointer` flag so the Linux can collect accurate call stacks with each sample. COZ is run with the default sampling period of 1ms, and a sample batch size of ten. Each performance experiment runs for a minimum of 100ms with a cooloff period of 10ms after each experiment. Due to space limitations, we only profile throughput (and not latency) in this evaluation.

### 4.2 Effectiveness

We demonstrate causal profiling’s effectiveness through case studies. Using COZ, we collect causal profiles for Memcached, SQLite, and the PARSEC benchmark suite. Using these causal profiles, we were able to make small changes to two of the real applications and six PARSEC benchmarks, resulting in performance improvements as large as 68%. Table 3 summarizes the results of our optimization efforts. We describe our experience using COZ with each application below.

#### 4.2.1 Case Study: blackscholes

The blackscholes benchmark, provided by Intel, solves the Black–Scholes differential equation to price a portfolio of stock options. We placed a progress point after each thread completes one round of the iterative approximation to the differential equation (`blackscholes.c:259`). COZ identifies many lines in the `CNDF` and `BlkSchlsEqEuroNoDiv` functions that would have a small impact if optimized. This same code was identified as a bottleneck by ParaShares [27]; this is the only optimization we describe here that was previously reported. This block of code performs the main numerical work of the program, and uses many temporary



**Figure 4:** In the dedup benchmark, COZ identified hash bucket traversal as a bottleneck. Collisions per-bucket for the first 4000 buckets before, midway through, and after optimization of the dedup benchmark (note different y-axes). The dashed horizontal line shows average collisions per-utilized bucket for each version. Replacing dedup’s hash function improved performance by 8%.

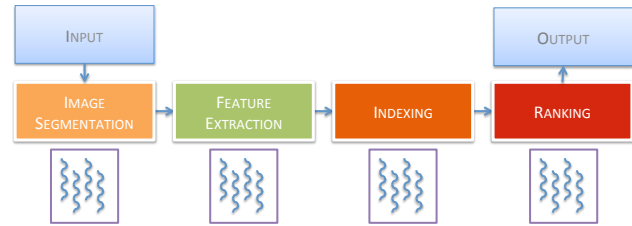
variables to break apart the complex computation. Manually eliminating common subexpressions and combining 61 piecewise calculations into 4 larger expressions resulted in a  $2.56\% \pm 0.41\%$  program speedup.

#### 4.2.2 Case Study: dedup

The dedup application performs parallel file compression via deduplication. This process is divided into three main stages: fine-grained fragmentation, hash computation, and compression. We placed a progress point immediately after dedup completes compression of a single block of data (`encoder.c:189`).

COZ identifies the source line `hashtable.c:217` as the best opportunity for optimization. This code is the top of the `while` loop in `hashtable_search` that traverses the linked list of entries that have been assigned to the same hash bucket. This suggests that dedup’s shared hash table has a significant number of collisions. Increasing the hash table size had no effect on performance. This led us to examine dedup’s hash function, which could also be responsible for the large number of hash table collisions. We discovered that dedup’s hash function maps keys to just 2.3% of the available buckets; over 97% of buckets were never used during the entire execution.

The original hash function adds characters of the hash table key, which leads to virtually no high order bits being set. The resulting hash output is then passed to a bit shifting procedure intended to compensate for poor hash functions. We removed the bit shifting step, which increased hash table utilization to 54.4%. We then changed the hash function to bitwise XOR 32 bit chunks of the key. This increased hash table utilization to 82.0% and resulted in an  $8.95\% \pm 0.27\%$  performance improvement. Figure 4 shows the rate of bucket collisions of the original hash function, the same hash function without



**Figure 5:** Ferret’s pipeline. The middle four stages each have an associated thread pool; the input and output stages each consist of one thread. The colors represent the impact on throughput of each stage, as identified by COZ: green is low impact, orange is medium impact, and red is high impact.

the bit shifting “improvement”, and our final hash function. The entire optimization required changing just three lines of code. As with ferret, this result was achieved by one graduate student who was initially unfamiliar with the code; the entire profiling and tuning effort took just two hours.

**Comparison with gprof.** We ran both the original and optimized versions of dedup with gprof. As with ferret, the optimization opportunities identified by COZ were not obvious in gprof’s output. Overall, `hashtable_search` had the largest share of highest execution time at 14.38%, but calls to `hashtable_search` from the hash computation stage accounted for just 0.48% of execution time; Gprof’s call graph actually obscured the importance of this code. After optimization, `hashtable_search`’s share of execution time reduced to 1.1%.

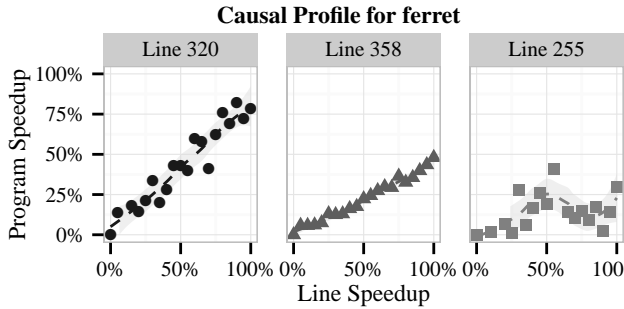
#### 4.2.3 Case Study: ferret

The ferret benchmark performs a content-based image similarity search. Ferret consists of a pipeline with six stages: the first and the last stages are for input and output. The middle four stages perform image segmentation, feature extraction, indexing, and ranking. Ferret takes two arguments: an input file and a desired number of threads, which are divided equally across the four middle stages. We first inserted a progress point in the final stage of the image search pipeline to measure throughput (`ferret-parallel.c:398`). We then ran COZ with the `--source-scope` argument to limit our attention to the `ferret-parallel.c` file, rather than across the entire ferret toolkit.

Figure 6 shows the top three lines identified by COZ, using its default ranking metric. Lines 320 and 358 are calls to `cass_table_query` from the indexing and ranking stages. Line 255 is a call to `image_segment` in the segmentation stage. Figure 5 depicts ferret’s pipeline with the associated thread pools (colors indicate COZ’s computed impact on throughput of optimizing these stages).

Because each important line falls in a different pipeline stage, and because COZ did not find any important lines in the queues shared by adjacent stages, we can easily “optimize” a specific line by shifting threads to that stage. We modified





**Figure 6:** COZ output for the unmodified ferret application. The x-axis shows the amount of virtual speedup applied to each line, versus the resulting change in throughput on the y-axis. The top two lines are executed by the indexing and ranking stages; the third line is executed during image segmentation.

ferret to let us specify the number of threads assigned to each stage separately, a four-line change.

COZ did not find any important lines in the feature extraction stage, so we shifted threads from this stage to the three other main stages. After three rounds of profiling and adjusting thread assignments, we arrived at a final thread allocation of 20, 1, 22, and 21 to segmentation, feature extraction, indexing, and ranking respectively. The reallocation of threads led to a  $21.27\% \pm 0.17\%$  speedup over the original configuration, using the same number of threads.

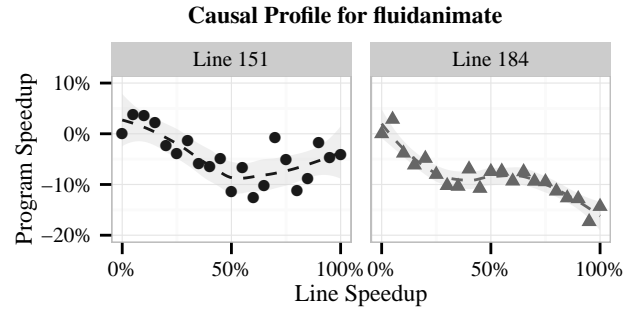
**Comparison with gprof.** We also ran ferret with gprof in both the initial and final configurations. Optimization opportunities are not immediately obvious from that profile. For example, in the flat profile, the function `cass_table_query` appears near the bottom of the ranking, and is tied with 56 other functions for most cumulative time.

Gprof also offers little guidance for optimizing ferret. In fact, its output was virtually unchanged before and after our optimization, despite a large performance change.

#### 4.2.4 Case Study: fluidanimate

The fluidanimate benchmark, also provided by Intel, is a physical simulation of an incompressible fluid for animation. The application spawns worker threads that execute in eight concurrent phases, separated by a barrier. We placed a progress point immediately after the barrier, so it executes each time all threads complete a phase of the computation.

COZ identifies a single modest potential speedup in the thread creation code, but there was no obvious way to speed up this code. However, COZ also identified two significant points of contention, indicated by a downward sloping causal profile. Figure 7 shows COZ’s output for these two lines. This result tells us that optimizing the indicated line of code would actually *slow down* the program, rather than speed it up. Both lines COZ identifies are in a custom barrier implementation, immediately before entering a loop that repeatedly calls `pthread_mutex_trylock`. Removing this spinning from the barrier would reduce the contention, but it was simpler to replace the custom barrier with the



**Figure 7:** COZ output for fluidanimate, prior to optimization. COZ finds evidence of contention in two lines in `parsec_barrier.cpp`, the custom barrier implementation used by both fluidanimate and streamcluster. This causal profile reports that optimizing either line will slow down the application, not speed it up. These lines precede calls to `pthread_mutex_trylock` on a contended mutex. Optimizing this code would increase contention on the mutex and interfere with the application’s progress. Replacing this inefficient barrier implementation sped up fluidanimate and streamcluster by 37.5% and 68.4% respectively.

default `pthread_barrier` implementation. This one line change led to a  $37.5\% \pm 0.56\%$  speedup.

#### 4.2.5 Case Study: streamcluster

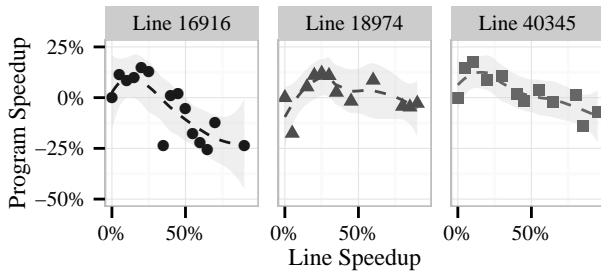
The streamcluster benchmark performs online clustering of streaming data. As with fluidanimate, worker threads execute in concurrent phases separated by a custom barrier, where we placed a progress point. COZ identified a call to a random number generator as a potential line for optimization. Replacing this call with a lightweight random number generator had a modest effect on performance ( $\sim 2\%$  speedup). As with fluidanimate, COZ highlighted the custom barrier implementation as a major source of contention. Replacing this barrier with the default `pthread_barrier` led to a  $68.4\% \pm 1.12\%$  speedup.

#### 4.2.6 Case Study: swaptions

The swaptions benchmark is a Monte Carlo pricing algorithm for swaptions, a type of financial derivative. Like blackscholes and fluidanimate, this program was developed by Intel. We placed a progress point after each iteration of the worker threads’ main loop (`HJM_Securities.cpp:99`).

COZ identified three significant optimization opportunities, all inside nested loops over a large multidimensional array. One of these loops just zeroed out consecutive values, so we replaced all but the outermost loop with a call to `memset`. A second loop filled part of the same large array with values from a distribution function, with no obvious opportunities for optimization. The third nested loop iterated over the same array again, but traversed the dimensions in an irregular order. We reordered the loops to traverse dimensions from left to right whenever possible in order to improve the locality of the loop body. This change, along with the call to `memset`, sped execution by  $15.8\% \pm 1.10\%$ .

Causal Profile for SQLite



**Figure 8:** COZ output for SQLite before optimizations. The three lines correspond to entry points for `sqlite3MemSize`, `pthreadMutexLeave`, and `pcache1Fetch`. A small optimization to each of these lines will improve program performance, but beyond about a 25% speedup, COZ predicts that the optimization would actually lead to a slowdown (because of contention). Changing indirect calls into direct calls for these functions improved performance by  $25.6\% \pm 1.0\%$ .

#### 4.2.7 Case Study: Memcached

Memcached is a widely-used in-memory caching system. To evaluate cache performance, we ran a benchmark ported from the Redis performance benchmark. This program spawns 50 parallel clients that collectively issue 100,000 SET and GET requests for randomly chosen keys. We placed a progress point at the end of the `process_command` function, which handles each client request.

Most of the lines COZ identifies are cases of contention, with a characteristic downward-sloping causal profile plot. One such line is at the start of `item_remove`, which locks an item in the cache and then decrements its reference count, freeing it if the count goes to zero. To reduce lock initialization overhead, Memcached uses a static array of locks to protect items, where each item selects its lock using a hash of its key. Consequently, locking any one item can potentially contend with independent accesses to other items whose keys happen to hash to the same lock index. Because reference counts are updated atomically, we can safely remove the lock from this function, which resulted in a  $9.39\% \pm 0.95\%$  speedup.

#### 4.2.8 Case Study: SQLite

The SQLite database library is widely used by many applications to store relational data. The embedded database, which can be included as a single large C file, is used for many applications including Firefox, Chrome, Safari, Opera, Skype, iTunes, and is a standard component of Android, iOS, Blackberry 10 OS, and Windows Phone 8. We evaluated SQLite performance using a write-intensive parallel workload, where each thread rapidly inserts rows to its own private table. While this benchmark is synthetic, it exposes any scalability bottlenecks in the database engine itself because all threads should theoretically operate independently. We placed a progress point in the benchmark itself (which is linked with the database), which executes after each insertion.

Results for Unoptimized Applications

Benchmark	Progress Point	Top Optimization
bodytrack	TicketDispenser.h:106	ParticleFilter.h:262
canneal	annealer_thread.cpp:87	netlist_elem.cpp:82
facesim	taskQDistCommon.c:109	MATRIX_3X3.h:136
freqmine	fp_tree.cpp:383	fp_tree.cpp:301
vips	threadgroup.c:360	im_Lab2LabQ.c:98
x264	encoder.c:1165	common.c:687

**Table 4:** The locations of inserted progress points for the remaining PARSEC benchmarks, and the top optimization opportunities that COZ identifies. Note that we exclude one PARSEC benchmark, raytrace, due to time constraints.

COZ identified three important optimization opportunities, shown in Figure 8. At startup, SQLite populates a large number of structs with function pointers to implementation-specific functions, but most of these functions are only ever given a default value. The three functions COZ identified unlock a standard pthread mutex, retrieve the next item from a shared page cache, and get the size of an allocated object. These simple functions do very little work, so the overhead of the indirect function call is relatively high. Replacing these indirect calls with direct calls resulted in a  $25.60\% \pm 1.00\%$  speedup.

**Comparison with conventional profilers.** Unfortunately, running a version of SQLite compiled to use gprof segfaults immediately. The application does run with the Linux perf tool, which reports that the three functions COZ identified account for a total of just 0.15% of total runtime. Using perf, a developer would be misled into thinking that optimizing these functions would be a waste of time. COZ accurately shows that the opposite is true: optimizing these functions has a dramatic impact on performance.

#### Effectiveness Summary

Our case studies confirm that COZ is effective at identifying optimization opportunities and guiding performance tuning. In every case, the information COZ provided led us directly to the optimization we implemented. COZ identified optimization opportunities in all of the PARSEC benchmarks, but some required more invasive changes that are out of scope for this paper. Table 4 summarizes our findings for the remaining PARSEC benchmarks. We have submitted patches to the developers of all the applications we optimized.

#### 4.3 Accuracy

For most of the optimizations described above, it is not possible to quantify the effect our optimization had on the specific lines that COZ identified. However, for two of our case studies—ferret and dedup—we can directly compute the effect our optimization had on the line COZ identified and compare the resulting speedup to COZ’s predictions. Our results show that COZ’s predictions are highly accurate.

To optimize ferret, we increased the number of threads for the indexing stage from 16 to 22, which increases the

throughput of line 320 by 27%. COZ predicted that this improvement would result in a 21.4% program speedup, which is nearly the same as the 21.2% we observe.

For dedup, COZ identified the top of the `while` loop that traverses a hash bucket’s linked list. By replacing the degenerate hash function, we reduced the average number of elements in each hash bucket from 76.7 to just 2.09. This change reduces the number of iterations from 77.7 to 3.09 (accounting for the final trip through the loop). This reduction corresponds to a speedup of the line COZ identified by 96%. For this speedup, COZ predicted a performance improvement of 9%, very close to our observed speedup of 8.95%.

#### 4.4 Efficiency

We measure COZ’s profiling overhead on the PARSEC benchmarks running with the native inputs. The sole exception is streamcluster, where we use the test inputs, because execution time was excessive with the native inputs.

Figure 9 breaks down the total overhead of running COZ on each of the PARSEC benchmarks by category. The average overall overhead is 17%.

The primary contributor to COZ’s overhead is the introduction of delays for virtual speedup. This source of overhead can be reduced by performing fewer performance experiments during a program’s run, in exchange for increasing the execution time required to collect useful causal profiles.

The second greatest contributor to COZ’s overhead is sampling overhead: the cost of collecting samples, processing those samples, and producing profile output. The primary cost is due to initiating sampling with the `perf` API for every new thread. In addition, sampling is disabled during introduced delays, which requires two system calls (one before the delay, and one after).

Finally, startup overhead is due to COZ’s initial processing of debugging information for the profiled application. Because the benchmarks are sufficiently long running (mean: 103s) to amortize startup time, this effect is minimal.

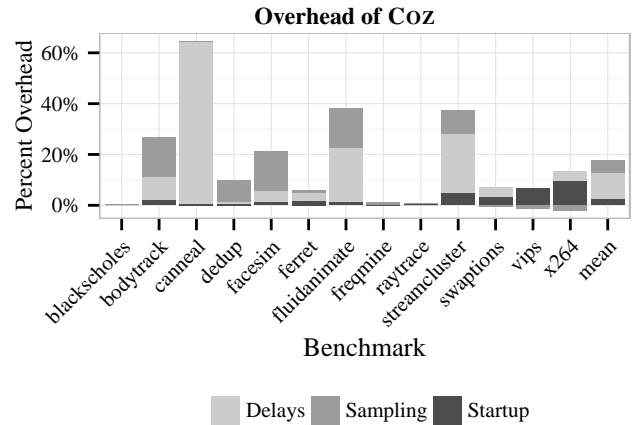
**Efficiency Summary.** COZ’s profiling overhead is on average 17% (minimum: 0.1%, maximum: 65%). For all but three of the benchmarks, its overhead was under 30%. Given that the widely used `gprof` profiler can impose much higher overhead (e.g., 6× for `ferret`, versus 6% with COZ), these results confirm that COZ has sufficiently low overhead to be used in practice.

## 5. Related Work

Causal profiling identifies and quantifies optimization opportunities, while most past work on profilers has focused on collecting detailed (though not necessarily actionable) information with low overhead.

### 5.1 General-Purpose Profilers

General-purpose profilers are typically implemented using instrumentation, sampling, or both. Systems based on sam-



**Figure 9:** Percent overhead for each of COZ’s possible sources of overhead. *Delays* are the overhead due to adding delays for virtual speedups, *Sampling* is the cost of collecting and processing samples, and *Startup* is the initial cost of processing debugging information. Note that sampling results in slight performance *improvements* for *swaptions*, *vips*, and *x264*.

pling (including causal profiling) can arbitrarily reduce *probe effect*, although sampling must be unbiased [36].

The UNIX `prof` tool and `oprofile` both use sampling exclusively [30, 43]. `Oprofile` can sample using a variety of hardware performance counters, which can be used to identify cache-hostile code, poorly predicted branches, and other hardware bottlenecks. `Gprof` combines instrumentation and sampling to measure execution time [19]. `Gprof` produces a call graph profile, which counts invocations of functions segregated by caller. Cho, Moseley, et al. reduce the overhead of `Gprof`’s call-graph profiling by interleaving instrumented and un-instrumented execution [10]. Path profilers add further detail, counting executions of each path through a procedure, or across procedures [2, 6].

### 5.2 Parallel Profilers

Past work on parallel profiling has focused on identifying the critical path or bottlenecks, although optimizing the critical path or removing the bottleneck may not significantly improve program performance.

**Critical Path Profiling.** IPS uses traces from message-passing programs to identify the critical path, and reports the amount of time each procedure contributes to the critical path [35]. IPS-2 extends this approach with limited support for shared memory parallelism [34, 45]. Other critical path profilers rely on languages with first-class threads and synchronization to identify the critical path [22, 38, 41]. Identifying the critical path helps developers find code where optimizations will have some impact, but these approaches do not give developers any information about how much performance gain is possible before the critical path changes. Hollingsworth and Miller introduce two new metrics to approximate optimization potential: *slack*, how much a proce-

duration can be improved before the critical path changes; and *logical zeroing*, the reduction in critical path length when a procedure is completely removed [23]. These metrics are similar to the optimization potential measured by a causal profiler, but can only be computed with a complete program activity graph. Collection of a program activity graph is costly, and could introduce significant probe effect.

**Bottleneck Identification.** Several approaches have used hardware performance counters to identify hardware-level performance bottlenecks [9, 13, 33]. Techniques based on binary instrumentation can identify cache and heap performance issues, contended locks, and other program hotspots [5, 32, 37]. ParaShares and Harmony identify basic blocks that run during periods with little or no parallelism [26, 27]. Code identified by these tools is a good candidate for parallelization or classic serial optimizations. Bottlenecks, a profile analysis tool, uses heuristics to identify bottlenecks using call-tree profiles [3]. Given call-tree profiles for different executions, Bottlenecks can pinpoint which procedures are responsible for the difference in performance. The FreeLunch profiler and Visual Studio’s contention profiler identify locks that are responsible for significant thread blocking time [12, 18]. BIS uses similar techniques to identify highly contended critical sections on asymmetric multiprocessors, and automatically migrates performance-critical code to faster cores [25]. Bottleneck graphs present thread execution time and parallelism in a visual format that highlights program bottlenecks [14]. Unlike causal profiling, these tools do not predict the performance impact of removing bottlenecks. All these systems can only identify bottlenecks that arise from explicit thread communication, while causal profiling can measure parallel performance problems from any source, including cache coherence protocols, scheduling dependencies, and I/O.

**Profiling for Parallelization and Scalability.** Several systems have been developed to measure potential parallelism in serial programs [17, 44, 46]. Like causal profiling, these systems identify code that will benefit from developer time. Unlike causal profiling, these tools are not aimed at diagnosing performance issues in code that has already been parallelized.

Kulkarni, Pai, and Schuff present general metrics for available parallelism and scalability [29]. The Cilkview scalability analyzer uses performance models for Cilk’s constrained parallelism to estimate the performance effect of adding additional hardware threads [21]. Causal profiling can detect performance problems that result from poor scaling on the current hardware platform.

**Time Attribution Profilers.** Time attribution profilers assign “blame” to concurrently executing code based on what other threads are doing. Quartz introduces the notion of “normalized processor time,” which assigns high cost to code that runs while a large fraction of other threads are blocked [4]. CPPROFJ extends this approach to Java programs with aspects [20]. CPPROFJ uses finer categories for time: running,

blocked for a higher-priority thread, waiting on a monitor, and blocked on other events. Tallent and Mellor-Crummey extend this approach further to support Cilk programs, with an added category for time spent managing parallelism [42]. The WAIT tool adds fine-grained categorization to identify bottlenecks in large-scale production Java systems [1]. Unlike causal profiling, these profilers can only capture interference between threads that directly affects their scheduler state.

### 5.3 Performance Guidance and Experimentation

Several systems have employed delays to extract information about program execution times. Mytkowicz et al. use inserted delays to validate the output of profilers on single-threaded Java programs [36]. Snelick, Jájá et al. use delays to profile parallel programs [39]. This approach measures the impact of slowdowns in combination, which is impractical because it requires a complete execution of the program for each of an exponential number of configurations. Active Dependence Discovery (ADD) introduces performance perturbations to distributed systems and measures their impact on response time [8]. ADD requires a complete enumeration of system components, and requires developers to insert performance perturbations manually. Song and Lu use machine learning to identify performance anti-patterns in source code [40]. None of these approaches quantify the effect of potential optimizations, which causal profiling measures directly.

## 6. Conclusion

Profilers are the primary tool in the programmer’s toolbox for identifying performance tuning opportunities. Previous profilers only observe actual executions and correlate code with execution time or performance counters. This information can be of limited use because the amount of time spent does not necessarily correspond to where programmers should focus their optimization efforts. Past profilers are also limited to reporting end-to-end execution time, an unimportant quantity for servers and interactive applications whose key metrics of interest are throughput and latency. Causal profiling is a new, experiment-based approach that establishes causal relationships between hypothetical optimizations and their effects. By virtually speeding up lines of code, causal profiling identifies and quantifies the impact on either throughput or latency of any degree of optimization to any line of code. Our prototype causal profiler, COZ, is efficient, accurate, and effective at guiding optimization efforts.

## Acknowledgments

*Omitted for double-blind reviewing.* This material is based upon work supported by the National Science Foundation under Grants No. CCF-1012195 and CCF-1439008. Charlie Curtsinger was supported by a Google PhD Research Fellowship. The authors thank Dan Barowy, Emma Tosch, and John Vilks for their feedback and helpful comments.

## References

- [1] E. R. Altman, M. Arnold, S. Fink, and N. Mitchell. Performance analysis of idle programs. In *OOPSLA*, pages 739–753. ACM, 2010.
- [2] G. Ammons, T. Ball, and J. R. Larus. Exploiting hardware performance counters with flow and context sensitive profiling. In *PLDI*, pages 85–96. ACM, 1997.
- [3] G. Ammons, J.-D. Choi, M. Gupta, and N. Swamy. Finding and removing performance bottlenecks in large systems. In *ECOOP*, volume 3086 of *Lecture Notes in Computer Science*, pages 170–194. Springer, 2004.
- [4] T. E. Anderson and E. D. Lazowska. Quartz: A tool for tuning parallel program performance. In *SIGMETRICS*, pages 115–125, 1990.
- [5] M. M. Bach, M. Charney, R. Cohn, E. Demikhovskiy, T. Devor, K. Hazelwood, A. Jaleel, C.-K. Luk, G. Lyons, H. Patil, and A. Tal. Analyzing parallel programs with Pin. *Computer*, 43(3):34–41, Mar. 2010.
- [6] T. Ball and J. R. Larus. Efficient path profiling. In *MICRO*, pages 46–57, 1996.
- [7] A. P. Black and T. D. Millstein, editors. *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications, OOPSLA 2014, part of SPLASH 2014, Portland, OR, USA, October 20-24, 2014*. ACM, 2014.
- [8] A. B. Brown, G. Kar, and A. Keller. An active approach to characterizing dynamic dependencies for problem determination in a distributed environment. In *Integrated Network Management*, pages 377–390. IEEE, 2001.
- [9] M. Burtscher, B.-D. Kim, J. R. Diamond, J. D. McCalpin, L. Koesterke, and J. C. Browne. PerfExpert: An easy-to-use performance diagnosis tool for HPC applications. In *SC*, pages 1–11. IEEE, 2010.
- [10] H. K. Cho, T. Moseley, R. E. Hank, D. Bruening, and S. A. Mahlke. Instant profiling: Instrumentation sampling for profiling datacenter applications. In *CGO*, pages 1–10. IEEE Computer Society, 2013.
- [11] C. Curtsinger and E. D. Berger. STABILIZER: Statistically sound performance evaluation. In *Proceedings of the seventeenth international conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS '13*, New York, NY, USA, 2013. ACM.
- [12] F. David, G. Thomas, J. Lawall, and G. Muller. Continuously measuring critical section pressure with the free-lunch profiler. In Black and Millstein [7], pages 291–307.
- [13] J. R. Diamond, M. Burtscher, J. D. McCalpin, B.-D. Kim, S. W. Keckler, and J. C. Browne. Evaluation and optimization of multicore performance bottlenecks in supercomputing applications. In *ISPASS*, pages 32–43. IEEE Computer Society, 2011.
- [14] K. Du Bois, J. B. Sartor, S. Eyerhan, and L. Eeckhout. Bottle graphs: Visualizing scalability bottlenecks in multi-threaded applications. In *OOPSLA*, pages 355–372, 2013.
- [15] DWARF Debugging Information Format Committee. *DWARF Debugging Information Format, Version 4*, 2010.
- [16] Free Software Foundation. *Debugging with GDB*, tenth edition.
- [17] S. Garcia, D. Jeon, C. M. Louie, and M. B. Taylor. Kremlin: rethinking and rebooting gprof for the multicore age. In *PLDI*, pages 458–469. ACM, 2011.
- [18] M. Goldin. Thread performance: Resource contention concurrency profiling in visual studio 2010. *MSDN magazine*, page 38, 2010.
- [19] S. L. Graham, P. B. Kessler, and M. K. McKusick. gprof: a call graph execution profiler. In *SIGPLAN Symposium on Compiler Construction*, pages 120–126. ACM, 1982.
- [20] R. J. Hall. CPPROFJ: Aspect-Capable Call Path Profiling of Multi-Threaded Java Applications. In *ASE*, pages 107–116. IEEE Computer Society, 2002.
- [21] Y. He, C. E. Leiserson, and W. M. Leiserson. The Cilkview scalability analyzer. In *SPAA*, pages 145–156. ACM, 2010.
- [22] J. M. D. Hill, S. A. Jarvis, C. J. Siniolakis, and V. P. Vasilev. Portable and architecture independent parallel performance tuning using a call-graph profiling tool. In *PDP*, pages 286–294, 1998.
- [23] J. K. Hollingsworth and B. P. Miller. Slack: a new performance metric for parallel programs. *University of Maryland and University of Wisconsin-Madison, Tech. Rep*, 1994.
- [24] Intel. *Intel VTune Amplifier 2015*, 2014.
- [25] J. A. Joao, M. A. Suleman, O. Mutlu, and Y. N. Patt. Bottleneck identification and scheduling in multithreaded applications. In *ASPLOS*, pages 223–234. ACM, 2012.
- [26] M. Kambadur, K. Tang, and M. A. Kim. Harmony: Collection and analysis of parallel block vectors. In *ISCA*, pages 452–463. IEEE Computer Society, 2012.
- [27] M. Kambadur, K. Tang, and M. A. Kim. Parashares: Finding the important basic blocks in multithreaded programs. In *Euro-Par*, Lecture Notes in Computer Science, pages 75–86, 2014.
- [28] kernel.org. *perf: Linux profiling with performance counters*, 2014.
- [29] M. Kulkarni, V. S. Pai, and D. L. Schuff. Towards architecture independent metrics for multicore performance analysis. *SIGMETRICS Performance Evaluation Review*, 38(3):10–14, 2010.
- [30] J. Levon and P. Elie. Oprofile: A system profiler for Linux, 2004.
- [31] J. D. Little. OR FORUM: Little’s Law as Viewed on Its 50th Anniversary. *Operations Research*, 59(3):536–549, 2011.
- [32] C.-K. Luk, R. S. Cohn, R. Muth, H. Patil, A. Klauser, P. G. Lowney, S. Wallace, V. J. Reddi, and K. M. Hazelwood. Pin: building customized program analysis tools with dynamic instrumentation. In V. Sarkar and M. W. Hall, editors, *PLDI*, pages 190–200. ACM, 2005.
- [33] B. P. Miller, M. D. Callaghan, J. M. Cargille, J. K. Hollingsworth, R. B. Irvin, K. L. Karavanic, K. Kunchithapadam, and T. Newhall. The paradyn parallel performance measurement tool. *IEEE Computer*, 28(11):37–46, 1995.
- [34] B. P. Miller, M. Clark, J. K. Hollingsworth, S. Kierstead, S.-S. Lim, and T. Torzewski. IPS-2: The second generation of a parallel program measurement system. *IEEE Trans. Parallel Distrib. Syst.*, 1(2):206–217, 1990.

- [35] B. P. Miller and C.-Q. Yang. IPS: An interactive and automatic performance measurement tool for parallel and distributed programs. In *ICDCS*, pages 482–489, 1987.
- [36] T. Mytkowicz, A. Diwan, M. Hauswirth, and P. F. Sweeney. Evaluating the accuracy of Java profilers. In *PLDI*, pages 187–197. ACM, 2010.
- [37] N. Nethercote and J. Seward. Valgrind: a framework for heavyweight dynamic binary instrumentation. In *PLDI*, pages 89–100. ACM, 2007.
- [38] Y. Oyama, K. Taura, and A. Yonezawa. Online computation of critical paths for multithreaded languages. In *IPDPS Workshops*, volume 1800 of *Lecture Notes in Computer Science*, pages 301–313. Springer, 2000.
- [39] R. Snelick, J. JáJá, R. Kacker, and G. Lyon. Synthetic-perturbation techniques for screening shared memory programs. *Software Practice & Experience*, 24(8):679–701, 1994.
- [40] L. Song and S. Lu. Statistical debugging for real-world performance problems. In Black and Millstein [7], pages 561–578.
- [41] Z. Szebenyi, F. Wolf, and B. J. N. Wylie. Space-efficient time-series call-path profiling of parallel applications. In *SC*. ACM, 2009.
- [42] N. R. Tallent and J. M. Mellor-Crummey. Effective performance measurement and analysis of multithreaded applications. In *PPOPP*, pages 229–240. ACM, 2009.
- [43] K. Thompson and D. M. Ritchie. *UNIX Programmer’s Manual*. Bell Telephone Laboratories, 1975.
- [44] C. von Praun, R. Bordawekar, and C. Cascaval. Modeling optimistic concurrency using quantitative dependence analysis. In *PPOPP*, pages 185–196. ACM, 2008.
- [45] C.-Q. Yang and B. P. Miller. Performance measurement for parallel and distributed programs: A structured and automatic approach. *IEEE Trans. Software Eng.*, 15(12):1615–1629, 1989.
- [46] X. Zhang, A. Navabi, and S. Jagannathan. Alchemist: A transparent dependence distance profiling infrastructure. In *CGO*, pages 47–58. IEEE Computer Society, 2009.