

# The Benefit of Hindsight: Tracing Edge-Cases in Distributed Systems

Lei Zhang  
Emory University

Zhiqiang Xie  
Max Planck Institute for Software Systems

Jonathan Mace  
Max Planck Institute for Software Systems

Vaastav Anand  
Max Planck Institute for Software Systems

Ymir Vigfusson  
Emory University

## Abstract

Today’s distributed tracing frameworks only trace a small fraction of all requests. For application developers troubleshooting rare edge-cases, the tracing framework is unlikely to capture a relevant trace at all, because it cannot know which requests will be problematic until after-the-fact. Application developers thus heavily depend on luck.

In this paper, we remove the dependence on luck for any edge-case where symptoms can be programmatically detected, such as high tail latency, errors, and bottlenecked queues. We propose a lightweight and always-on distributed tracing system, Hindsight, where each constituent node acts analogously to a car dash-cam that, upon detecting a sudden jolt in momentum, persists the last hour of footage. Hindsight implements a *retroactive sampling* abstraction: when the symptoms of a problem are detected, Hindsight retrieves and persists coherent trace data from all relevant nodes that serviced the request. Developers using Hindsight receive the exact edge-case traces they desire; by comparison existing sampling-based tracing systems depend wholly on serendipity.

Our experimental evaluation shows that Hindsight successfully collects edge-case symptomatic requests in real-world use cases. Hindsight adds only nanosecond-level overhead to generate trace data, can handle GB/s of data per node, transparently integrates with existing distributed tracing systems, and persists full, detailed traces when an edge-case problem is detected.

## 1 Introduction

Distributed tracing frameworks record detailed, end-to-end traces of requests executing in distributed systems, and are helpful for a wide range of troubleshooting use cases [54, 57]. Particularly vexing for application developers are *edge cases*: performance problems that have adverse effects but are too rare or complex to ascribe obvious root causes.

Consider, as running examples, the problems of diagnosing uncommon errors, mitigating tail-latency, and conducting temporal provenance. Although high latency is easy to detect and immediately experienced by the end user [21], its root causes may be scattered across previously traversed machines and layers galore, with an exponential number of configuration parameters or hardware components that may

have contributed to the delay. Similarly, high-level API exceptions may culminate from prior interactions with other requests, system misconfiguration, hardware failures, transient faults [34], elusive programming mistakes, or even a significantly hard-to-reproduce combination of issues. Third, a request exhibiting issues because of an overfull queue somewhere in the distributed system is merely the casualty of an earlier problem [63]; troubleshooting requires systematically investigating these other, unrelated requests.

In typical production environments, tracing every request—including transmitting, processing, and storing comprehensive telemetry—requires enormous backend infrastructure and storage that is unacceptable to infrastructure operators. Tracing frameworks manage this overhead by collecting a small sample (0.001%) of traces [14, 33, 36, 59]. But since an edge case is rare, by definition, the framework will most likely *not yield edge-case traces for the developer to analyze*.

In this paper, we resolve the problem of tracing edge-case requests in production environments. We focus our attention on *symptomatic* edge cases, where the performance effects of the problem manifest shortly after its causes and where the impacts can be observed programmatically—a broad family of problems that includes the three classes of use cases above.

We built Hindsight—an always-on, lightweight distributed tracing system that is compatible with existing tracing APIs—as a practical tool for edge-case analysis. Hindsight offers *retroactive sampling*, inspired by network provenance [17, 34, 70, 72], to collect telemetry data back in time from the present moment of detection, across all machines that handled the request. Under retroactive sampling, all trace data is recorded locally but only reported when a symptom is detected, allowing applications to generate copious trace data in case they are needed without encumbering the tracing system’s backend collection infrastructure. Retroactive sampling ultimately reports the same volume of trace data as other sampling methods, but ensures that edge-case traces are not missed. To provide efficient and coherent retroactive sampling, Hindsight’s design carefully separates its dataplane – *e.g.*, generating trace data into fast local memory – from control logic – *e.g.*, for indexing metadata, coordinating among machines, and triggering collection for symptomatic requests on demand.

As demonstration, we apply Hindsight on three use cases corresponding to our running examples. We run experiments on the DeathStar Microservices Benchmark [25], the Hadoop Distributed File System [58], and on several micro- and macro-benchmarks. We have also integrated Hindsight as a replacement collection component for X-Trace [24] and OpenTracing [51]. Our experimental results show that Hindsight imposes nanosecond-level overhead when generating trace data, can scale to GB/s of data per node, rapidly reconstructs traces when triggered, and effectively captures problematic traces, as well as related lateral traces, within tens of milliseconds of identifying a symptom.

In summary, our paper makes the following contributions.

- We describe the retroactive sampling abstraction for capturing traces of symptomatic edge-cases.
- We present the design of Hindsight, a distributed tracing system that implements retroactive sampling. Hindsight is compatible with existing tracing APIs and can be transparently integrated with existing applications.
- We apply Hindsight on real-world use cases and show that efficiently collecting edge-case requests is practical.
- We evaluate Hindsight on multiple benchmarks and real systems, showing that it can achieve nanosecond-level overhead on trace data generation and handle GB/s data per node while collecting coherent traces.
- We illustrate that Hindsight is compatible and performs better than state-of-the-art tracing systems (X-Trace and Jaeger) with more efficient trace-data generation and lower overhead, while providing edge-case tracing.

## 2 Motivation

Distributed tracing systems are in widespread use in both open-source [30, 48, 73] and major internet companies [33, 53, 59]. These tools record traces of *end-to-end requests*: a trace contains logs, metrics, and other event data, along with timing and ordering, generated from every machine visited by the request. End-to-end request traces have proved to be especially useful for troubleshooting distributed systems since they explicitly tie together the individual slices of work performed across different machines, enabling an operator to observe how the work done by one machine influences, and is influenced by, work done on others [24, 53, 59]. Prior research has demonstrated a range of troubleshooting tasks using request traces, including common-case analyses centered on aggregate system behavior, distributions over data, and relationships between system components [54, 55, 57]; and edge-case analyses such as error diagnosis [40, 66, 67] and tail-latency troubleshooting [21, 38, 47, 60, 68]. In this paper, we consider general edge-case analysis use cases, so we begin with some examples.

**Error diagnosis (UC1).** Hardware failures, component errors, exceptions, and programming mistakes abound in large distributed systems. Despite advances in testing [15] and verification [27, 28, 37], new or complex problems slip through

the cracks and may wreak havoc if not promptly addressed. Application developers thus often play the role of detective, to identify root causes of errors. Yet without a trace of a problematic request, the developer will face a daunting task: each request traverses many different processes and machines and its outcome is influenced by every machine it visits. The symptoms of a problem often manifest far from the root causes [23, 40, 44], and the potential root causes are manifold, perhaps a combination of software or hardware problems on multiple nodes or network links [34].

**Tail-latency troubleshooting (UC2).** Distributed systems track a wide range of high-level health metrics, such as API distributions, latency percentiles, resource utilization, and many others [32, 33]. An operator may observe an unusual metric jump, say the 99.9<sup>th</sup> percentile latency has spiked for some important API. However, knowing about the spike is not enough; the application developer wants to identify the specific service, code paths, or conditions that contribute to the peak, so that any underlying problems can be addressed [21, 38, 47].

**Temporal provenance (UC3).** Many modern distributed systems respond to requests through an architecture of loosely coupled microservices [57]. Application developers need tools for tracking queuing issues [8–11, 13] when the number of components in a distributed system is large, since a request  $r$  exhibiting symptoms (*e.g.*, prolonged queuing time) may not be the true culprit for the backlogged queue. Rather, the developer wants to follow the *temporal provenance* of  $r$  to determine *lateral traces* of other related requests with which  $r$  interacted through shared components and queues [63].

### 2.1 Trace Collection Infrastructure

Distributed tracing frameworks require a supporting backend *trace collection infrastructure* that is distinct from the traced applications. Applications continually generate trace data using a client library, which internally serializes and transmits the data over the network. The trace collection infrastructure, or *backend* for short, receives trace data, and thereafter processes it in various way to combine data from different machines, construct trace objects, apply user-defined functions, and ultimately store the trace in a database [33].

Production applications cannot trace *every* request, as doing so would generate far too much data and impose an enormous burden on the backend [33, 39, 59]. At Facebook, for example, several MBs of tracing data are generated per traced request [33]; at Google, traces are typically more detailed than debug-level logging [59]. Thus a key design concern for current tracing frameworks is to reduce the volume of data reported to the backend.

**Head-based sampling.** The de facto approach of all existing distributed tracing tools is to simply capture fewer traces. Instead of tracing every request, the application will only record traces for a small number of requests, chosen by *ran-*

*dom sampling* in the client library [57]. Sampling decisions occur at the *beginning* of a request when it enters the system and before it starts executing. Trace data is only recorded if the sampling decision is successful. Head-based sampling satisfies an all-or-nothing property: if a request is sampled, then applications will record an *entire* trace of the request, including all data it produces across all machines it visits. Otherwise, the application records nothing about the request on *any* machine. Coherent traces are essential for distributed tracing. A partial or fragmented trace has limited value in diagnosis [24, 59] because it sacrifices the end-to-end visibility that makes the trace useful in the first place [53].

Overall, head-based sampling reduces the volume of trace data sent to the trace collection infrastructure, with the sampling probability determining how much trace data is collected. In production, the sampling probability is typically very low: Jaeger’s default is 0.1% [31]; some systems sample as few as 0.001% [33, 59].

## 2.2 Sampling Derails Edge-Case Analysis

Head-based sampling is indiscriminate: the fate of a trace is determined a priori, upon arrival to the system. With low sampling probabilities, the vast majority of requests are not traced. For edge-case analysis especially, requests of interest are inherently rare and thus unlikely to be traced at all.

The outcome for the application developer investigating edge cases is thus disappointing: problems arise, but traces do not exist. For example, the developer may have reports that errors took place (UC1) yet the corresponding ‘rare’ requests were not sampled when those requests began. The developer therefore lacks the detailed cross-machine data necessary for finding the error’s root cause. Likewise, the application’s high-level metric monitoring may indicate a spike in end-to-end tail-latency (UC2); a developer is thus aware that these high-latency outliers exist, yet without a trace, the developer lacks the ability to localize the problem to a particular component or request class. The situation is even more problematic when investigating bottlenecked queues via temporal provenance (UC3): the tracing system will have only a vanishing probability that traces of *all* relevant requests in the queue were captured, since each request was sampled independently.

Practitioners have long pointed out a discord between what traces are interesting and what traces gets sampled [12, 14, 16, 49, 50]. *Tail-based sampling* describes how the backend trace collection infrastructure could identify ‘interesting’ traces and only persist those to storage [35, 36, 55], yet such approaches elide the basic scalability challenge faced by the trace collection infrastructure in the first place: we cannot collect all traces. To date, the desire for a low-overhead sampling system that would support edge-case analysis remains unfulfilled.

## 3 Challenges

The crucible of edge-case analysis under prevailing sampling approaches is that developers miss out on *relevant* traces. Our

goal is to remedy this situation: to capture all relevant request traces without relying on serendipity or imposing additional overhead on the tracing backend. We first address a series of motivating questions.

**What traces are relevant?** Typically, strong signals of an interesting request come towards the end of its execution when more information is known about its fate [10, 14, 16, 49, 50, 55]. Moreover, looking to our motivating use cases, we observe that many issues manifest a known set of symptoms that are cheaply detectable soon after they occur, like error codes (UC1), high end-to-end response time (UC2), or an overfull queue (UC3). In each scenario, the symptoms can be detected programmatically even though the underlying causes remain unknown. We call such issues *symptomatic edge-cases*.

### Why not just collect traces after we see the symptoms?

Paradoxically, since the request has already executed, enabling tracing at a late point of the request means we have already missed the events that led to the anomaly. The sure-fire way of obtaining coherent traces for any edge-case is to trace from the very beginning of the request.

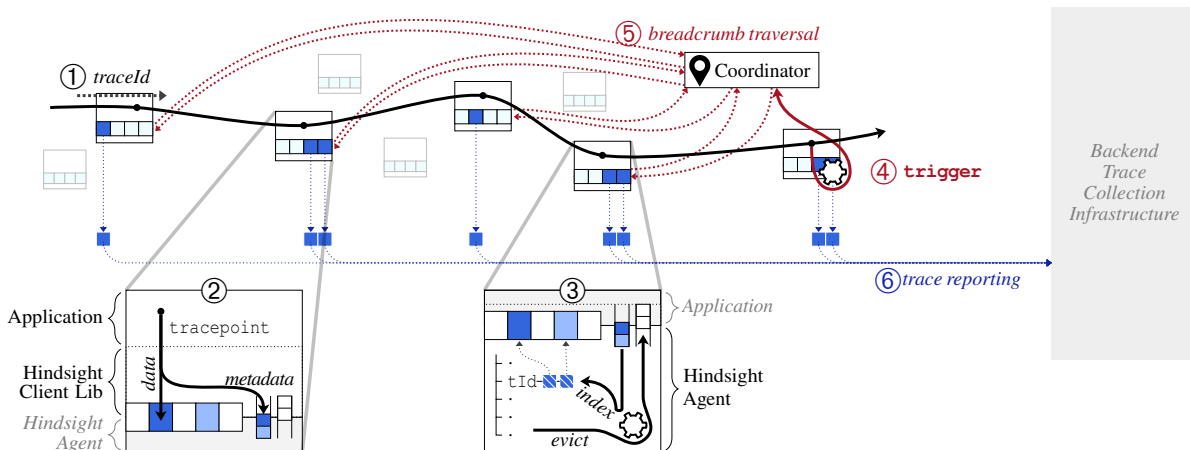
### Why don’t tracing frameworks use ring-buffers?

Machine-level logging and telemetry systems often make use of in-memory ring-buffers for temporarily persisting telemetry in case it is needed in the near-future [17, 20, 69]. While conceptually appealing, the mismatch in data granularity makes them a poor fit for distributed tracing. Machines execute many requests concurrently and the generated trace data on any given machine is interleaved with that of other requests, making it difficult to extract any one request’s specific slice of data. Conversely, each request executes across many machines, so trace data must be identified and extracted from all relevant machines.

## 4 Retroactive Sampling

Building on these answers, we now detail our approach to the problem, using an abstraction we call *retroactive sampling*. Whereas traditional distributed tracing systems based on sampling obtain edge-case traces only by luck, retroactive sampling enables applications to obtain them explicitly. Inspired by work on network provenance [17, 18, 34, 70, 72], retroactive sampling endows the application with the benefit of hindsight: when the application detects something wrong or anomalous about a request, the tracing system can reach back into the recent past to retrieve its trace data as follows.

**Nodes generate, but do not collect, all trace data.** Our tracing system is not clairvoyant: it lacks information about whether a request will be impacted by some sort of problem (and thereby be interesting) until after the symptoms manifest themselves [14, 16, 49, 50]. As argued, the only conclusion is to record *all* trace data. Recording the data, however, need not burden the backend infrastructure until the traces are reported, giving applications leeway to generate fine-detail trace data



**Fig. 1: The end-to-end lifecycle of a trace in Hindsight (§5.1).** A request (①, solid black line) traverses the processes of the system, depositing trace data data using Hindsight’s client library (②). A Hindsight agent on each machine locally indexes data (③). If the application detects outlier symptoms, it can trigger trace collection (④). Hindsight’s coordinator traverses breadcrumbs left behind on each machine (⑤), instructing each agent to set aside and report this request’s trace data. Subsequently, agents will report the triggered trace data to the existing backend trace collection infrastructure (⑥).

locally in case it will be beneficial after a problem occurs.

**A trigger signals that a trace is relevant.** Under retroactive sampling, a tracing system programmatically signals after-the-fact that a request was indeed relevant. We call this a *trigger* that may be fired at any point during or even after a request is served, since symptoms (*e.g.*, high request latency) may only be measurable after the fact.

**Traces are collected across machines.** So far, traces are only recorded locally into memory without any forwarding or coordination among the nodes—we deliberately refrain from doing more complex processing for every trace. Thus for each individual request, its trace data will be dispersed across memory of the multiple machines that it visited. Yet when a machine fires a trigger to collect a trace, it is effectively requesting *all* corresponding trace data from *all* machines visited by the request. In response to the trigger mechanism, systems implementing retroactive sampling must therefore notify pertinent machines that they must *report* the relevant trace data to the backend.

**Trace data expires.** The tracing system deposits trace data in fast local memory on the machine where the data is generated. Lazy collection via triggers implies that, eventually, we will exhaust the memory available to record traces; naturally, we should overwrite old trace data to make way for traces of new requests. We call the implicit time duration between generating data and overwriting it the *event horizon*. After the event horizon, the trace data has been overwritten and is gone forever. With large and detailed traces, or highly constrained memory, the event horizon will be short—tens of seconds.

**Triggers are automatic.** To guarantee trace collection of

any traces of interest within the event horizon, the developer needs *automatic* triggers. Symptomatic edge cases, by definition, can be detected programmatically within a short interval of their root causes. An automatic trigger for tail-latency (UC2), for example, could simply check for a tardy response.

## 5 Design

### 5.1 Overview

Hindsight is a distributed tracing framework that implements retroactive sampling. We begin with a high-level overview of Hindsight’s main components before describing our key design choices. Fig. 1 depicts the end-to-end flow of a request executing in a distributed system, interacting with Hindsight’s tracing API along the way.

- ① Upon request arrival, Hindsight generates a random unique traceId for this request and thereafter propagates it alongside the request to any machine visited [44].
- ② While the request executes, it generates trace data (*e.g.*, logs, tracepoints, spans) that the application reports using Hindsight’s tracepoint client API. When the request completes, its trace data is left scattered across the machines visited.
- ③ Internally, a Hindsight *agent* running on each machine receives and manages trace data. Hindsight agents do not eagerly report trace data to the backend collection infrastructure – instead, agents index and organize metadata in-memory and await an explicit trigger. For most traces, there is no trigger, the trace is not reported, and the data residing on the local machine is overwritten by new data.
- ④ The application may detect an outlier symptom (*e.g.*, an erroneous response value, high latency, or a bottlenecked queue) during the request and thus want to persist its trace. The application invokes Hindsight’s trigger API on the machine

where the symptom is detected, passing the `traceId`.

- ⑤ The local agent immediately initiates trace collection by sending `traceId` to Hindsight’s logically centralized *coordinator* service. The coordinator’s role is to recursively contact all machines that the request visited by following *breadcrumbs* deposited by the request at each machine it visited; a breadcrumb is a pointer to another machine involved in the request (*e.g.*, to the RPC caller or callee).
- ⑥ When instructed by the coordinator, an agent will set aside relevant data for this `traceId` and report it to the backend.

## 5.2 API Compatibility

Hindsight centers on redesigning the *internals* of a distributed tracing system to support retroactive sampling while remaining compatible with other use cases, such as tracing common-case requests that establish aggregate system behavior. Hindsight’s design is encompassing: rare triggers for edge-case requests can intermingle with frequent triggers for baseline common-case requests. Moreover, Hindsight does not redesign or replace the backend trace collection components.

Hindsight is compatible with existing distributed tracing APIs and can be transparently integrated into an application’s existing tracing instrumentation (*e.g.*, OpenTracing [51]). In our evaluation (*cf.* §7), we integrated Hindsight as the collection component for both X-Trace [24] and OpenTracing [51]. Existing tracing instrumentation API calls relay directly to Hindsight. For example, API calls that create and annotate spans are proxied as `tracepoint` calls. Hindsight seamlessly supports any existing head-based sampling policy: `trigger` is invoked immediately for a positive sampling decision or if the `sampled` flag is set within a received OpenTracing context. Lastly, Hindsight piggybacks its `breadcrumbs` within OpenTracing’s context propagation; transparently enabling Hindsight’s trace coordination procedure.

For a developer wanting to benefit from retroactive sampling, the only additional concern is when to invoke `trigger`. For example, this may entail adding a `trigger` call within a service’s exception handler, or after checking for outlier latency upon a request’s completion. A developer can explicitly decide the conditions for triggering a request, but Hindsight also provides a library of automatic triggers based on metric percentiles, categorical features, and exceptions. All our use cases (UC1–UC3) can be implemented using Hindsight’s `autotrigger`s (*cf.* §6.2).

## 5.3 Data Coherence

A trace is only useful if it has coherent data: the trace must contain *all* information about the request across *all* machines visited. If a trace is patchy (*i.e.*, missing data from one or more machines) then the trace’s value as a whole is compromised.

Hindsight encounters potential trace incoherence in several places: at ③ when an agent chooses which old data to evict; at ⑤ if the coordinator is too slow to contact all machines that handled a request; and at ⑥ if more traces have been triggered

than can be reported. The threat is exacerbated because it only takes one agent dropping its slice of a trace to render the remaining data on other agents effectively worthless.

We ensure coherence as follows. Within a single agent (③), a trace is an atomic entity; when an agent evicts a trace, the trace’s data are evicted in its entirety. Agents carefully organize and index metadata to account for a trace being non-contiguous and fragmented in memory, and enable efficient insertion and eviction. After a trace is triggered (⑤), the coordinator service rapidly disseminates the trigger, recursively following `breadcrumbs`, to ensure no agent inadvertently evicts its slice of the trace. `Breadcrumbs` are scalable, as the coordinator only needs to contact services explicitly known to be part of a request’s execution, and trivially scales by sharding on the triggered `traceId`. Triggers are eagerly disseminated; subsequently agents report the actual trace data to the backend collection infrastructure asynchronously in the background. Hindsight enforces a user-configured rate-limit on background trace reporting.

Unlike head-based sampling that traces a fixed percentage of requests, retroactive sampling may fire triggers arbitrarily often. A `trigger-happy` application may cause a backlog of unreported traces within the agent (⑥). If the agent passes a threshold of triggered data, it will begin to drop data for triggered traces. A backlog of data on one agent strongly implies a backlog on others, thus agents making different choices about which data to drop risk trace incoherence. Hindsight uses consistent hashing of `TraceIds` in several places to decide priority both for reporting and evicting trace data.

Finally, a developer might invoke `trigger` at several different places in their application, for a diverse range of symptoms; developers can distinguish different triggers by providing a `triggerId`. Hindsight can prevent a profuse `trigger` from impacting trace collection of other, low-frequency triggers: agents implement weighted fair sharing over triggers for reporting and evicting trace data, and users can configure weights and rate-limits for each `triggerId`.

## 5.4 Efficient Data Management

From the perspective of the backend trace collection infrastructure, an application using Hindsight will report a similar volume of traces and have similar demands to conventional head-based sampling. For an individual machine (③), however, retroactive sampling entails the application generate a substantially higher volume of trace data into local memory. For perspective, traces in production systems can be verbose, containing tens of thousands of events [33], and often incorporating detailed debug information; *e.g.*, our Hadoop experiments, (*cf.* §7.1), sustained 30 MB/s per process. Even when `tracepoints` involve string formatting and argument packing, prior work on efficient logging has attained lower bounds of a few nanoseconds per logging point [64, 65]. More generally, Hindsight’s `tracepoint` API can be a sink for arbitrary telemetry sources such as function tracing or new technol-

ogy like Intel® Processor Trace (PT) [29, Ch. 36], which can reach 100–200 MB/s of data per core at 5–15% runtime overhead. Handling this data at large volumes requires a design centered on performance.

The most sensitive performance bottleneck for Hindsight occurs between the client application generating data via `tracepoint`, and the local Hindsight agent that manages trace data. Agents must efficiently receive new trace data for *every* request, evict old data to be overwritten, and set aside triggered data for reporting, while managing trace data coherently with respect to other agents. Simple approaches fail to meet these needs; ring-buffers, *e.g.*, are effective across many systems designs but deteriorate in our setting by incoherently overwriting trace data and posing significant costs to extract triggered data via expensive scanning. Instead, our design establishes a clear split between *control* and *data* activities, which congregates general-purpose data and efficiency in the data plane, and embeds all logic in the control plane.

Hindsight’s data plane is concerned with efficiently writing trace data from client applications. Hindsight writes trace data to a large shared memory pool, divided up into *buffers*. A client invoking `tracepoint` writes to a buffer, and each buffer only belongs to one trace at a time. Full buffers are continually circulated to the agent via shared memory queues; the agent likewise continually frees up old buffers for reuse.

Hindsight’s agent process encapsulates control activities: it receives buffer metadata, indexes them according to `traceId`, and coherently evicts old buffers. Agents receive triggers and communicate with Hindsight’s coordinator, manage breadcrumbs linking the trace data that is strewn across many agents, extract triggered trace data, and report data asynchronously to the backend trace collection infrastructure. Hindsight’s control and data distinction yields an efficient agent implementation, and confers additional desirable properties, such as making it easy to change and update control logic like eviction and fair sharing policies.

In our description thus far, we assume that an application will trace *all* requests into local memory. Such completeness is not a strict requirement for Hindsight; an application may opt to trace a smaller fraction of requests (*e.g.*, 10%), for instance, if generating data is expensive or unoptimized, or if the application is highly performance-sensitive. Hindsight remains effective in such a scenario since we still capture significantly more traces to local memory than can be reported to tracing backends (*e.g.*, 0.1% and lower). The scale-back is supported via a `trace percentage` option (default 100%). Hindsight enforces scale-back coherently across agents through consistent `traceId` hashing (cf. §5.3), and will maintain breadcrumbs and support triggers for all requests, traced or not.

## 5.5 Divorcing Triggers from Traces

Applications initiate retroactive sampling via Hindsight’s trigger API (⊕). Triggers are orthogonal to traces in Hindsight for several reasons.

**Efficiency.** In principle, triggers could be calculated directly atop trace data. For example, end-to-end latency can be calculated directly from span start and end times recorded in a trace. Yet applying user-defined functions in-dataplane poses an infeasible performance challenge. Moreover, traces are brittle data structures: in practice, as engineers continually modify and update instrumentation, the traces may sometimes ‘break’ [1–7, 41], leading to incorrect or degenerate derivative features and metrics [33, 43, 59].

**Integration.** In the common case, symptoms are easy to detect and localize: top-level error codes; high latency; increased queue time. Such symptoms can be readily recognized and cheaply computed without the trigger mechanism needing the trace data itself. It further leads to a straightforward integration of triggers into existing metric monitoring systems regardless of their architecture.

**Lateral traces.** Outlier behavior may not map directly to a single request; instead there may be several other related *lateral* requests. For example, to diagnose a bottlenecked queue (UC3), a trigger needs to capture traces for the previous  $N$  requests to understand what led to queue buildup [63]; to diagnose a write-ahead log, we desire all requests blocking on a log sync [9, 13]; to diagnose resource contention we require all requests contending for a slow disk or network [8, 10, 11]. Hindsight enables an application to atomically trigger a group of related lateral `traceIds`; internally Hindsight will ensure that the group as a whole is coherently collected.

## 6 Implementation

We have implemented Hindsight in  $\approx 4$ KLOC in C for the client library,  $\approx 3$ KLOC in Go for the agent and coordinator, and  $\approx 300$ LOC for a JNI-based Java client library for integrating with Hadoop in our experiments (cf. §7). We chose C for dataplane efficiency and Go for its ease of use for the more complex control plane logic. In this section, we elaborate how we meet the Hindsight design goals (cf. §5).

### 6.1 Agent Data Management

Hindsight pre-allocates a fixed-size *buffer pool* in shared memory for storing trace data. We choose a fixed-size pool to bound memory overheads, a desirable property for telemetry systems [62]. Hindsight subdivides the buffer pool into fixed-size buffers, by default 32 kB. Client processes write trace data to buffers via Hindsight’s client API. The agent process does not touch data in the buffer pool except when reporting triggered traces. At each point in time, a buffer can only contain trace data of a single request; no two different requests will write trace data to the same buffer at the same time. A single trace will typically comprise (1) multiple non-contiguous buffers and (2) many buffers scattered across numerous machines. Buffers are the granularity of data management within Hindsight. Within clients and agents, a buffer is addressable by its `bufferId`—its offset into the buffer pool.

<code>begin(traceId)</code>	Record that request with id <code>traceId</code> is processing in the current thread.
<code>tracepoint({payload})</code>	Record data for the current trace; <code>payload</code> is arbitrary size.
<code>breadcrumb(address)</code>	Adds a breadcrumb to the current trace on the local node pointing to some other node address.
<code>serialize()</code>	Obtain the current <code>traceId</code> and a <code>breadcrumb</code> to the current node.
<code>end()</code>	Request finished processing in current thread; flush and remove buffers.
<code>trigger(traceId,trigger lateralTraceIds...)</code>	Tell Hindsight to trigger <code>traceId</code> for collection; <code>triggerId</code> is an arbitrary identifier used for rate-limiting.

**Table 1: The Hindsight client API** can be called by applications directly, or by other distributed tracing tools (e.g., X-Trace [24], OpenTracing [51]) as the trace collection component (cf. §7).

## 6.2 Client Library

**Writing trace data.** Table 1 outlines Hindsight’s client API. When a request begins executing in a thread, it must call `begin`; subsequently it may call `tracepoint` an arbitrary number of times; and finally when it completes executing in a thread, it must call `end`. As noted in §5.2, this aligns with existing tracing instrumentation calls. Hindsight internally maintains thread-local state including the current `traceId` and a pointer to a buffer. Calls to `tracepoint` write directly to the thread-local buffer without needing any synchronization. Synchronization is only required when acquiring a new buffer or returning a buffer to the agent; these operations touch shared-memory queues but are infrequent. A buffer must be acquired during `begin`, returned during `end`, and is otherwise only acquired or flushed when it becomes full (i.e., from successive `tracepoint` calls).

**Communicating with agents.** To acquire a buffer, the client library polls a shared-memory *available queue*. The queue returns an integer `bufferId` that points into the buffer pool. Clients do not block if the available queue is empty—clients immediately return, and instead of writing to the buffer pool, write to a ‘null buffer’, which is discarded afterwards. When the client fills a buffer, it writes its `traceId` and the `bufferId` to a shared-memory *complete queue*. The agent continually drains the complete queue, and likewise continually returns fresh buffers to the available queue. Shared memory queues are lock-free and support batch operations; using batch operations, agents are robust to queue contention from multiple client writer threads.

This paired channel design forms a natural separator between control and data with two desirable properties: (1) queues only communicate metadata—they avoid data copying and use a single integer `bufferId` to represent, by default, a 32 kB buffer; (2) communication is infrequent, occurring only when buffers are filled or a thread switches over to execute a different request, thereby minimizing synchronization.

From the client library’s perspective, it cheaply and blindly writes trace data into shared memory and forwards only the control metadata to agents; conversely agents are agnostic to buffer contents—they do not inspect data in the shared memory pool and use only the metadata communicated via the complete queue.

**Depositing breadcrumbs.** Hindsight clients deposit *breadcrumbs* on every machine visited by a request. Each breadcrumb addresses another Hindsight agent that handled the request. When a request arrives at a node, it is carrying the breadcrumb of the previous node. Hindsight’s breadcrumb API is called during trace context deserialization, which establishes within the local agent a pointer to the previous node visited for this `traceId`. Similarly, during trace context serialization prior to communication with another node or an RPC response, Hindsight will insert the current node’s breadcrumb transparently as a key-value baggage field of the trace context [43, 51]. For synchronous RPCs, breadcrumbs are sufficient for reconstructing full traces triggered by any node, including requests with arbitrary concurrency and fan-out. To handle asynchronous settings, clients can also use the breadcrumb API to eagerly establish *forward-breadcrumbs* to a named destination node prior to communication. Internally, the breadcrumb API call reports the `traceId` and breadcrumb to a shared memory *breadcrumb queue*. Agents poll this queue and index breadcrumbs alongside the buffer metadata. However, agents do not forward or act upon breadcrumbs until a trace is explicitly collected with a trigger.

**Triggering trace collection.** Applications initiate trace collection by invoking `trigger`, passing a `traceId` and a `triggerId`; internally this writes to a shared-memory *trigger queue*. The `trigger` API allows any condition that can be programmatically detected to initiate trace collection in Hindsight—e.g., a trigger detecting latency outliers; a trigger for erroneous return values; or a trigger for high queue latency. A trigger can specify `traceIds` for a group of related lateral traces (cf. §5.5). A developer can implement custom outlier detection and invoke `trigger` directly, or they can make use of Hindsight’s autotrigger library, a separate collection of triggers that track simple conditions over time and automatically invoke `trigger` when a condition is met:

<code>PercentileTrigger(p)</code>	Clients call <code>addSample(traceID, measurement)</code> . Trigger fires for measurements above percentile $p$ . (e.g., high latency or resource consumption)
<code>CategoryTrigger(f)</code>	Clients call <code>addSample(traceID, label)</code> . Trigger on rare categorical data that is less frequent than threshold $f$ (e.g., rare API calls or odd attributes)
<code>ExceptionTrigger</code>	Trigger on an exception or error code
<code>TriggerSet(T,N)</code>	Trigger all lateral traces in a sliding window of length $N$ when trigger $T$ fires (e.g., gather the last $N$ traceIDs in a bottlenecked queue)

`TriggerSet` is noteworthy as a building block for lateral tracing; it captures recent traces related to a request that exhibited symptoms.

### 6.3 Agent

Agents maintain metadata for each trace in a map keyed by `traceId`. The metadata for each `traceId` includes a list of `bufferIds` and a list of breadcrumbs. Agents also maintain metadata for each trigger that has fired, including the `traceId`, `triggerId`, and zero or more lateral `traceIds`.

**Indexing and reusing buffers.** Agents poll the complete queue, each time reading the `traceId` and `bufferId` of a full buffer, then adding this `bufferId` to the trace’s metadata. The agent then updates an LRU of `traceIds`. The agent performs eviction when 80% or more of the buffer pool is indexed, by removing the least recently used `traceId` and returning all its `bufferIds` to the available queue. Agents never touch data within buffers; all management is done using `bufferIds`.

**Local triggers.** Agents poll the trigger queue, each time draining the trigger metadata. The agent immediately forwards the trigger metadata and breadcrumbs to the coordinator, which begins disseminating the trigger to other related agents by recursively following breadcrumbs. Meanwhile the agent schedules the trigger for collection. In the case of a spammy trigger, an agent may decide to immediately discard the trigger instead of forwarding and scheduling it—this is implemented using a per-`triggerId` token bucket. Doing so prevents flooding other agents with triggers.

**Remote triggers.** Agents receive remote triggers fired by other machines via the coordinator. To facilitate rapid trigger dissemination, the agent immediately responds with any breadcrumbs it has accumulated for the specified `traceId`. Agents do not rate-limit remote triggers (unlike local triggers) since it risks incoherent traces. Instead, all remote triggers are scheduled for collection.

**Reporting traces.** When a trigger is scheduled for collection, its `traceId` and lateral `traceIds` are removed from the agent’s LRU and can no longer be evicted by the regular buffer eviction cycle. The trigger metadata is then inserted into a per-`triggerId` priority queue. In the normal case when an agent is not overloaded, the queue will be empty. The agent asynchronously pulls trigger metadata from the queues; fetches all buffers for all `traceIds` from the buffer pool; reports the buffer *contents* to the backend collection infrastructure; and finally frees the `bufferIds` by returning them to the available queue. If configured, the agent will apply a rate limit to pace its data reporting to the collection infrastructure. New data for any triggered `traceIds` can continue to arrive via the available queue, possibly rescheduling the trigger if existing data is already reported.

**Ignoring triggers during overload.** If trace collection becomes overloaded, the priority queues will begin to fill up. The agent tracks the proportion of buffer pool pending for collection, and when this exceeds 40%, the agent will pick a `triggerId` and abandon the lowest-priority trigger from

its queue. Abandoning a trigger entails removing it from the priority queue and returning all buffers of all `traceIds` to the available queue. If a `traceId` happens to be included in *multiple* triggers, its buffers are only returned after all triggers are abandoned. The agent provides a fair allocation of buffer pool to each `triggerId` and selects whichever `triggerId` most exceeds its weighted max-min fair share. In addition to memory pressure, a trigger will be automatically abandoned if it has not been reported after a configurable delay (5 minutes by default).

**Reporting traces during overload.** During overload, the agent continues to report traces as described above for the normal case. The agent implements weighted fair queueing to select the next `triggerId` queue, and will adhere to any configured per-`triggerId` rate limits. From the selected queue, the agent dequeues the highest-priority trigger, and reports its data as described above for the normal case.

**Coherence during overload.** The per-`triggerId` priority queues uses consistent hashing of the `traceId` to determine priority. This priority is consistent across all agents that may have relevant data and results in a given trace enjoying the same priority across all agents. Consequently, even during overload, agents will consistently report the highest-priority traces and evict the lowest-priority traces. The agent’s fair allocation of buffer pool to each `triggerId` and fair sharing of reporting bandwidth means a low-throughput `triggerId` is not impacted by a spammy trigger.

## 7 Evaluation

We now evaluate how effectively Hindsight overcomes the fundamental problem of head-based tracing methods in examples (UC1)–(UC3) and meets the goals of retroactive sampling to provide lightweight and effective request tracing.

**Tracing integration.** We compare Hindsight to two existing distributed tracing systems, X-Trace [24] and Jaeger [30]. We also integrate Hindsight with OpenTracing [51] and X-Trace [24], and, in our experiments, we evaluate Hindsight as a replacement backend to X-Trace [24].

**Systems.** We have integrated Hindsight with two distributed systems: the Hadoop Distributed File System (HDFS) [58], the DeathStar Microservices Benchmark (DSB) [25]. We also benchmark Hindsight’s single-node characteristics and cross-node trace retrieval performance.

**Summary.** Our experiments demonstrate the following.

- Hindsight effectively provides retroactive sampling and collects relevant edge-case traces across real use cases.
- Hindsight is lightweight and not a bottleneck for client applications. Hindsight can achieve < 5 ns tracepoint latency and tolerate write throughput up to 55 GB/s. Hindsight’s control/data split lets Hindsight’s agent match client data generation.
- With large trace data volumes, Hindsight’s ‘event horizon’



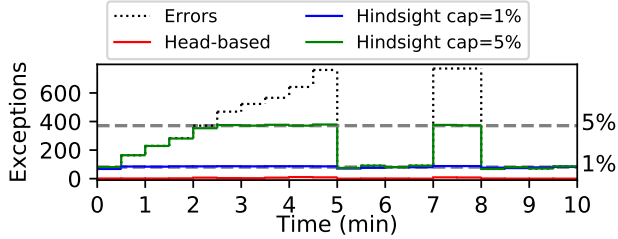


Fig. 2: Case Study: Error diagnosis (UC1). Exceptions captured by different sampling strategies as the error rate varies.

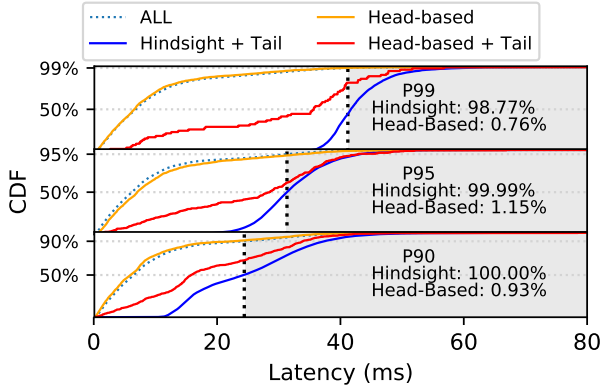


Fig. 3: Case Study: Tail-latency troubleshooting (UC2). Latency of requests captured through different sampling strategies with different tail-latency triggers (top to bottom).

extends tens of seconds into the past.

- Hindsight has substantially lower overheads than X-Trace and Jaeger when generating trace data.

## 7.1 Case Studies

We first evaluate Hindsight’s tracing libraries and trigger mechanisms on collecting edge-case requests on our motivating use cases.

**Error diagnosis (UC1).** We deploy Hindsight on DSB, a microservice system with 12 microservices and 17 backends [25]. We add an `ExceptionTrigger` from Hindsight’s autotrigger library to the `ComposePostService` module in DSB’s Social Network Benchmark. We run DSB’s default workload with 300 req/s. We randomly inject exceptions in the `ComposePostService` module with error rates ranging from 1% to 10%, varied after each 30 s; Fig. 2 shows the collected exceptions for each 30 s time window. We limit Hindsight’s agent to trigger on only 1% and 5% exceptions among all requests. Fig. 2 shows Hindsight is able to collect nearly all exceptions except when they surpass its collection limits. We add a burst of 10% error rates at 7 seconds and Hindsight rapidly meets its collection limit. Unlike Hindsight, a head-based sampling methods (at 1% sampling rate) can only collect 1% of the exceptions at random (see Fig. 2).

**Tail-latency troubleshooting (UC2).** We add a `Per-`  
`centileTrigger` from Hindsight’s autotrigger library to the

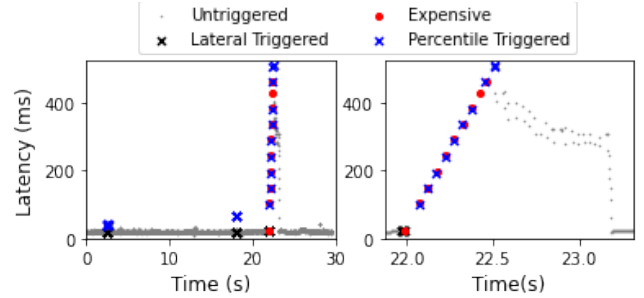


Fig. 4: Case Study: Temporal provenance (UC3). Lateral requests gathered (blue) after triggering on a high-latency request (red) due to an overfull queue in Hadoop’s HDFS.

`ComposePostService` module in the same setting as above, invoking `addSample` at the end of each `ComposePost` RPC call and providing the measured RPC duration. We set  $p$  to 0.99, 0.95, and 0.9, as different thresholds for tail latency. We inject 10% requests at random with 20–30 ms latency. We also compare with head-based sampling methods with 1% sampling rate. Fig. 3 shows that Hindsight predominantly collects requests near the threshold, whereas head-based sampling methods cannot. The vertical dotted line marks the tail-latency percentile threshold. The figure also shows that Hindsight collected 98.77%–100% of requests above the latency threshold; head-based sampling could gather only  $\approx 1\%$ .

**Temporal provenance (UC3).** We add a `QueueTrigger` from Hindsight’s autotrigger library to the HDFS NameNode queue — the `QueueTrigger` combines a `TriggerSet` with a `PercentileTrigger`, parameterized to capture  $N = 10$  most recently dequeued lateral requests when 99.99<sup>th</sup> percentile queueing latency is observed. We deploy HDFS on 10 machines (8 DataNodes, 1 NameNode, and 1 client) and run a Hindsight agent on each machine. We run a closed-loop workload of random 8 kB reads with 10 concurrent requests.

Fig. 4 (left) shows NameNode queue latency over time. We inject a burst of 10 expensive `createfile` requests 21 second into the trace that briefly saturate the queue—Fig. 4 (right) zooms in on this time window. The figure shows high-latency requests ( $\bullet$ ), requests that fire the autotrigger ( $\times$ ), and the additional lateral requests that were triggered to Hindsight ( $\times$ ). The first expensive request occurred at 22 seconds, followed by a pause while it was executed. Upon dequeuing the subsequent `read8k` request, `QueueTrigger` fired due to high queue latency, and Hindsight retroactively sampled the 10 prior traces leading up to the trigger. The sample included the culprit expensive request. Overall, all 10 expensive requests were sampled, 8 unrelated requests prior to the first expensive request, and 9 additional `read8k` requests. Moreover, several intermittent latency spikes occurred unrelated to the experiment (Fig. 4, left), which Hindsight also captured; upon investigation, these were due to garbage collection.

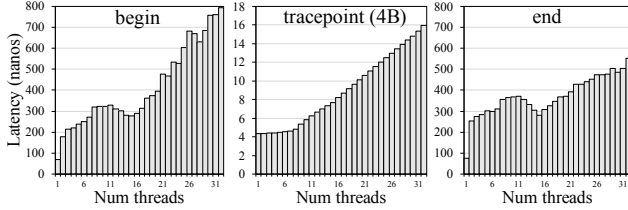


Fig. 5: Client API latency.

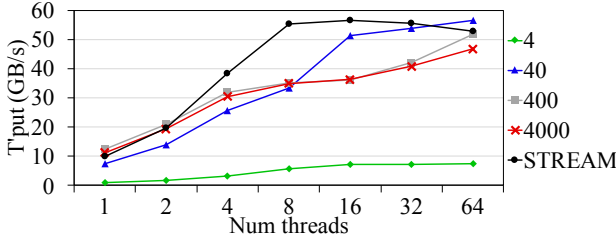


Fig. 6: Client Throughput.

## 7.2 Hindsight Tracing Performance

Hindsight’s design emphasizes low-overhead data ingestion. To evaluate this, we consider (a) the latency of client API operations; (b) achievable client data throughput; and (c) agent buffer indexing throughput. We run these experiments on an Intel Xeon W-2245 3.9 GHz workstation with 8 physical cores (16 hyperthreads) and 128 GB RAM. We run a benchmarking program that uses the Hindsight client library along with Hindsight’s agent sidecar process.

**Client API.** In this experiment, we run a client application that executes a loop comprising begin, 100 tracepoint, and end calls. We configure Hindsight with a 1 GB buffer pool divided into 32 kB buffers. Each tracepoint writes only a 4-byte payload—this use case is common in statically-preprocessed call site logging (*e.g.*, wherever a program logs only static strings [64]). We vary the number of threads from 1 to 32. Fig. 5 plots latency of begin, end, and tracepoint calls, as we vary the number of threads. For up to 8 threads, tracepoint latency ranges from 4.3–4.8 ns/call, benefiting from each thread writing to its own thread-local buffer. Beyond 8 cores, tracepoint latency increases linearly as we begin to multiplex threads on cores. By comparison, begin and end are more expensive and variable since they contend for shared-memory queues to acquire and return buffers; fortunately, these concurrent operations are uncommon by design. Latency is between 200–400 ns with at most 8 threads, and much lower for a single thread with no concurrent interference.

**Client throughput.** We repeat the experiment, keeping a buffer size of 32 kB. We now vary the size of the payload used in tracepoint calls, using 4, 40, 400, and 4,000 bytes per call. Fig. 6 plots the throughput achieved in GB/s. As expected, small payloads of 4 bytes fail to fully saturate mem-

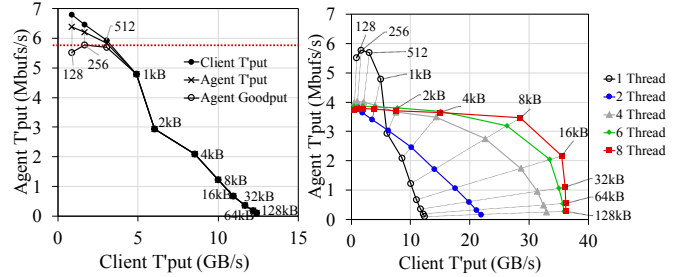


Fig. 7: Agent Throughput.

ory bandwidth, achieving only 887 MB/s with one thread and peaking at 7.55 GB/s with 64 threads. By contrast, even a modest increase in payload size to 40 bytes is enough to nearly saturate memory bandwidth; with 400 byte payloads, we achieve throughput of 12.5 GB/s on a single core. We include in Fig. 6 measurements of peak memory bandwidth from the STREAM benchmark [46]. These throughputs occur because in the common case, tracepoint is a memory copy to a thread-local buffer, and thus applications will never be bottlenecked by Hindsight’s client library.

**Agent throughput.** Client data throughput is moot if the agent cannot keep up; recall the agent must continually index buffers and cycle them back to the available queue for clients—if no buffer is available, the client will write to a ‘null buffer’ and the new data will be lost. This is a last-resort mechanism in Hindsight that fails to respect coherence, and thus must be avoided at all costs.

In this experiment we run one thread, fix the tracepoint payload size to 1 kB, and vary the buffer size from 128 B to 128 kB (Hindsight will fragment payloads across multiple buffers when necessary). Fig. 7 (left) shows the client-side data throughput along agent-side buffer throughput achieved by one thread. The data points are annotated with corresponding buffer size. We see, for example, that large buffer sizes (128 kB) can achieve peak client data throughput (12.1 GB/s) while requiring little of the agent. Conversely, tiny buffer sizes (128 B) stress the agent buffer throughput since we more frequently cycle buffers through the queues. Fig. 7 (left) plots three lines and indicates two important phenomena. The client throughput line plots the rate at which the client writes buffers, whereas the agent throughput line plots the rate at which the agent cycles buffers; the delta in-between are ‘null buffers’, written by the client because the available queue is empty, *i.e.*, the agent cannot keep up. Writing to null buffers means lost trace data; the third line, agent goodput, counts only the agent-side throughput of buffers that are useful, *i.e.*, buffers of traces that did not lose data. We observe that the goodput with 128 B buffers is lower than with 256 B buffers due to greater loss. In general, with  $\geq 1$  kB buffers, the agent is able to consistently keep up without losing data.

Fig. 7 repeats this experiment with varying numbers of threads, and plots client-side data throughput and agent-side

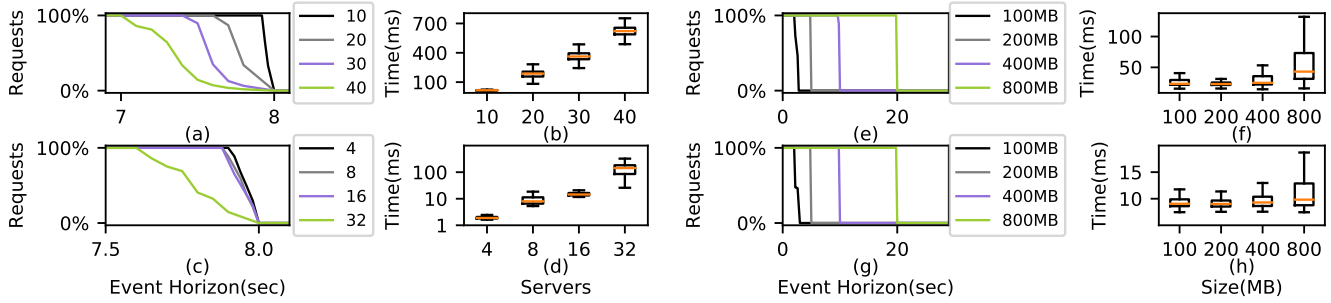


Fig. 8: Hindsight’s **event horizon** and **trace reconstruction** on gRPC-Chain and gRPC-Branch as we vary the number of servers ((a)–(d)) and the buffer pool size ((e)–(h)).

buffer goodput. Surprisingly, peak buffer throughput is only achievable with one thread; beyond that, client-side contention on the shared queues reduces the achievable client-side throughput. The same effect appears in Fig. 5 in increased begin and end latency. We also observe that buffer sizes of 16 kB and higher are sufficient for reaching peak write throughput while remaining comfortably within agent throughput limits; by default, we select 32 kB for Hindsight.

### 7.3 Retroactive Sampling

Hindsight is only useful if trace data is still in-memory on all agents when a trigger fires. Hindsight’s *event horizon* describes the time window between a request completing and its data being overwritten: triggers must fire within this window, or else it will be too late and the data will be overwritten.

We evaluate Hindsight’s event horizon on two distributed benchmarks: gRPC-Branch is a binary tree of  $N$  gRPC servers; each server concurrently invokes its 2 child servers and calls `tracepoint` with a configurable payload. gRPC-Chain is a configurable chain of  $N$  gRPC servers that invoke each other in sequence and calls `tracepoint`. We deployed all benchmarks on 12 CloudLab c8220 nodes [22], each with two Intel E5-2660 v2 10-core CPUs@2.2 GHz, and ran a Hindsight agent per process.

**Event horizon.** In this set of experiments, we introduce a delay between requests completing and subsequently firing a trigger. We randomly trigger 1% of requests, adding variable delay to the trigger, and measure how many coherent traces are ultimately received by the collection backend.

We vary the number of nodes between 4 and 40 (Fig. 8 (a)–(d)). We configure a 40 MB buffer pool with 4 kB buffers. The request rate is 1,000 r/s and each request writes 4 kB trace payload per server. We intentionally select a small buffer pool to more readily illustrate the event horizon. Fig. 8 (a,c) shows the percentage of coherent traces collected for each configured trigger delay. They illustrate how, 8 seconds after a request is complete, the trace data is no longer available for collection. Nevertheless, Hindsight successfully reports 100% of triggered traces after 7 seconds with 40 servers in the gRPC-Chain benchmark, where breadcrumb coordination

spans all 40 nodes. Fig. 8(b,d) shows the trace reporting time in milliseconds—the time between issuing a trigger and data received by the tracing backends. Trace reporting typically takes 100–400 ms, and even with 40 nodes to traverse, completes in less than one second.

**Extending the event horizon.** Hindsight’s event horizon is proportional to the size of the buffer pool configured for agents. To illustrate, we vary the buffer pool size (Fig. 8 (e)–(h)) from 100 MB to 800 MB. We configure the number of servers as 10 and 8, and use 4 kB buffers. The request rate is 2,000 r/s and each request writes 16 kB trace payload per server. Fig. 8 (e,g) illustrates how increasing the buffer pool size effectively extends the event horizon possible. Fig. 8 (f,h) shows that the delay between triggering a trace and receiving its data increases marginally with larger buffer pool sizes.

*Discussion.* For Hindsight to coherently capture a trace, its trigger must both fire, and be fully disseminated, within the event horizon – otherwise an agent may have already evicted relevant data. We note that once an agent learns of a trigger, the trace data is safe from eviction and can be reported asynchronously in the background. In tandem with buffer pool size, Hindsight’s event horizon is influenced by generated trace data volume. Lastly, for use cases extending even further into the past, Hindsight’s optional trace percentage (cf. §5.4) offers flexibility orthogonal to buffer pool size.

### 7.4 Comparison with the State-of-the-Art

We finally evaluate how Hindsight compares with existing state-of-the-art tracing systems.

**Comparison with X-Trace.** To provide more context, we integrate Hindsight with X-Trace [24]; all data generation is done within X-Trace, but instead of reporting data in-band, we reroute it to Hindsight. To a user, this integration is transparent. We run a benchmark that uses X-Trace to create traces and log 1,000 X-Trace events. With 8 threads and 32 kB buffers, the client generates 1.33 million events/s, totaling 259.6 MB/s and consuming 66,458 buffers/s—well below Hindsight’s peak ingestion throughput. By comparison, X-Trace’s client library experiences local performance bottlenecks and drops

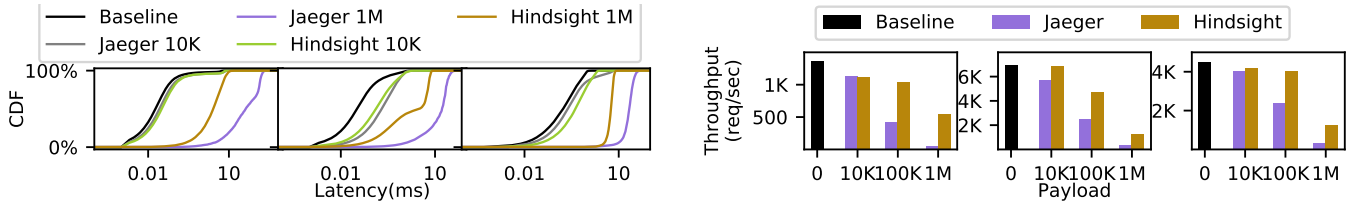


Fig. 9: **Comparison with Jaeger.** End-to-end latency (left) and throughput (right) with three workloads (compose-post, read-home-timeline, read-user-timeline).

data from overfull queues, only reporting data at 5.9 MB/s. This underscores how a dedicated dataplane design can yield substantial performance improvements.

We also offer numbers from HDFS [58] (cf. §7.1). We route Hadoop’s debug-level logging into the Hindsight-modified X-Trace. A sustained request workload generates  $\approx 30$  MB/s from Hadoop’s NameNode and 5 MB/s from each of Hadoop’s DataNodes. In our 10-node experimental setup (§7.1), this results in a total of 70 MB/s.

**Comparison with Jaeger.** In the final experiment, we deploy Hindsight in DSB on the cluster using `docker swarm`, interposing a Hindsight agent within each of the 12 DSB microservices. We augment the existing Jaeger instrumentation of DSB with breadcrumbs and masquerade Jaeger spans as tracepoint payloads. We use the built-in `wrk2` workload generator with `compose-post`, `read-home-timeline`, and `read-user-timeline`.

We compare Hindsight’s performance to Jaeger by measuring the tracing overhead on DSB, as well as a baseline of DSB stripped of all tracing instrumentation. To saturate the systems, we use a 100% trigger rate for Hindsight and a 100% sampling rate for Jaeger. We repeat the experiment with two different per-tracepoint payload sizes of 10 kB and 1 MB, to compensate for DSB’s short request length and minimal instrumentation. For apples-to-apples comparison, we add the payload to both Hindsight and Jaeger. Fig. 9 shows the end-to-end latency and throughput comparison. Hindsight displays significantly better performance than Jaeger at larger payloads. At 1 MB payloads, Hindsight boasts average latency of 3.2 ms and throughput of 549 req/sec compared to Jaeger’s average latency of 66.7 ms and 56 req/sec.

## 8 Related Work

**Distributed tracing.** Numerous prior works identify end-to-end requests as a useful granularity for slicing telemetry data and troubleshooting distributed systems. Example use cases include detecting anomalous request structures [36, 59, 63], diagnosing changes in the steady-state [19, 52, 56], modeling workloads [45, 61], and identifying resource and queue contention [26, 42, 63]. Distributed tracing systems have been presented in industry [33, 59], as open-source tools [30, 48, 51, 73], and in academia [24, 44]. We described how head-based sam-

pling is widespread in practice (cf. §2), and despite a desire for the ability to discriminate towards edge-cases, no techniques to do so have emerged to date. Tail-based sampling shares similar aspirations [35, 36]; its goal is to identify, from a collection of traces, which are most ‘interesting’ [12, 14, 16, 49, 50]. Yet tail-based sampling—perhaps a misnomer—does *not* capture fewer traces or reduce overheads on tracing backends; instead it is useful in tracing backends after collection and processing, for deciding whether traces can be discarded.

**Network provenance.** Hindsight is similar in spirit to network packet provenance systems that chronicle the history of network state, enabling use cases such as tracking the origin or path traversed by a packet across the network. Earlier systems, like ExSPAN [72] and SNP [70], adopt this abstraction; more recent works like SyNDB [34] and SPP [17] apply network provenance for packet-level root-cause analysis on Internet scale. Packet provenance systems primarily trace only packet metadata, which is well-structured and can be summarized in-band; these systems tackle additional trust challenges outside of Hindsight’s purview. By contrast, handling metadata to reconstruct the path of a trace is but one concern for Hindsight; Hindsight is focused on handling arbitrary payloads (*i.e.*, trace data), and the resulting performance, coherence, and fairness challenges. Hindsight also draws inspiration from works focused on temporal provenance [71] and packet reputation [18] in distributed systems, although Hindsight’s tracing abstractions operate entirely at the application level.

## 9 Conclusion

Hindsight circumvents the false dilemma between performance and usefulness for diagnosing symptomatic edge cases by providing developers detailed traces from the recent past when they encounter symptoms of failures, performance regressions, SLO violations, or other issues among the requests that pass through their systems, while making use of standard instrumentation for distributed tracing and without burdening the backend. We believe the retroactive sampling abstraction, and our Hindsight implementation of it, can shift the conversation around tracing away from mechanism (how to collect traces) to a question of policy (what traces should be collected), and allow distributed tracing systems to support edge-cases analysis: a key use case for which they were originally conceived.

## References

- [1] ACCUMULO-3725: Majc trace tacked onto minc trace. Retrieved May 2021 from <https://issues.apache.org/jira/browse/ACCUMULO-3725>.
- [2] CASSANDRA-7644: Tracing does not log commit-log/memtable ops when the coordinator is a replica. Retrieved May 2021 from <https://issues.apache.org/jira/browse/CASSANDRA-7655>.
- [3] CASSANDRA-7657: Tracing doesn't finalize under load when it should. Retrieved May 2021 from <https://issues.apache.org/jira/browse/CASSANDRA-7657>.
- [4] HBASE-13077: BoundedCompletionService doesn't pass trace info to server. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HBASE-13077>.
- [5] HBASE-15880: RpcClientImpl#tracedWriteRequest incorrectly closes HTrace span. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HBASE-15880>.
- [6] HTRACE-330: Add to Tracer, TRACE-level logging of push and pop of contexts to aid debugging "Can't close TraceScope..". Retrieved May 2021 from <https://issues.apache.org/jira/browse/HTRACE-330>.
- [7] Spring Cloud Sleuth 410: Trace ID problem when using Spring ThreadPoolTaskExecutor. Retrieved May 2021 from <https://github.com/spring-cloud/spring-cloud-sleuth/issues/410>.
- [8] HDFS-3751: DN should log warnings for lengthy disk IOs. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HDFS-3751>, 2014.
- [9] HBASE-8228: Investigate time taken to snapshot memstore. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HDFS-8228>, 2015.
- [10] HBASE-8744: Enable HBase to log the entire latency profile for HDFS packets resulting in slow writes. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HDFS-8744>, 2016.
- [11] HDFS-11461: DataNode Disk Outlier Detection. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HDFS-11461>, 2017.
- [12] Jaeger Issue 365: Adaptive Sampling. Retrieved May 2021 from <https://github.com/jaegertracing/jaeger/issues/365>, 2017.
- [13] HDFS-6110: adding more slow action log in critical write path. Retrieved May 2021 from <https://issues.apache.org/jira/browse/HDFS-6110>, 2018.
- [14] Jaeger Issue 1861: Delayed Sampling. Retrieved May 2021 from <https://github.com/jaegertracing/jaeger/issues/1861>, 2019.
- [15] Peter Alvaro, Joshua Rosen, and Joseph M Hellerstein. Lineage-driven fault injection. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 331–346, 2015.
- [16] Narayanan Arunachalam. Zipkin Secondary Sampling. Retrieved May 2021 from <https://github.com/openzipkin-contrib/zipkin-secondary-sampling>, 2019.
- [17] Ang Chen, Andreas Haeberlen, Wenchao Zhou, and Boon Thau Loo. One primitive to diagnose them all: Architectural support for internet diagnostics. In *Proceedings of the Twelfth European Conference on Computer Systems*, pages 374–388, 2017.
- [18] Ang Chen, Yang Wu, Andreas Haeberlen, Wenchao Zhou, and Boon Thau Loo. The good, the bad, and the differences: Better network diagnostics with differential provenance. In *Proceedings of the 2016 ACM SIGCOMM Conference, SIGCOMM '16*, page 115–128, New York, NY, USA, 2016. Association for Computing Machinery.
- [19] Mike Y Chen, Anthony Accardi, Emre Kiciman, Jim Lloyd, Dave Patterson, Armando Fox, and Eric Brewer. Path-based failure and evolution management. In *1st USENIX Symposium on Networked Systems Design & Implementation (NSDI'04)*, pages 23–23, 2004.
- [20] Weidong Cui, Xinyang Ge, Baris Kasikci, Ben Niu, Upamanyu Sharma, Ruoyu Wang, and Insu Yun. REPT: Reverse debugging of failures in deployed software. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI'18)*, pages 17–32, 2018.
- [21] Jeffrey Dean and Luiz André Barroso. The Tail at Scale. *Communications of the ACM*, 56(2):74–80, 2013.
- [22] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, et al. The design and operation of CloudLab. In *2019 USENIX Annual Technical Conference (USENIX ATC'19)*, pages 1–14, 2019.
- [23] Úlfar Erlingsson, Marcus Peinado, Simon Peter, Mihai Budiu, and Gloria Mainar-Ruiz. Fay: Extensible distributed tracing from kernels to clusters. *ACM Transactions on Computer Systems (TOCS)*, 30(4):1–35, 2012.
- [24] Rodrigo Fonseca, George Porter, Randy H Katz, and Scott Shenker. X-trace: A pervasive network tracing framework. In *4th USENIX Symposium on Networked Systems Design & Implementation (NSDI'07)*, 2007.

- [25] Yu Gan, Yanqi Zhang, Dailun Cheng, Ankitha Shetty, Priyal Rathi, Nayan Katarki, Ariana Bruno, Justin Hu, Brian Ritchken, Brendon Jackson, et al. An open-source benchmark suite for microservices and their hardware-software implications for cloud & edge systems. In *Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'19)*, pages 3–18, 2019.
- [26] Yu Gan, Yanqi Zhang, Kelvin Hu, Dailun Cheng, Yuan He, Meghna Pancholi, and Christina Delimitrou. Seer: Leveraging big data to navigate the complexity of performance debugging in cloud microservices. In *Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'19)*, pages 19–33, 2019.
- [27] Dennis Geels, Gautam Altekar, Petros Maniatis, Timothy Roscoe, and Ion Stoica. Friday: Global comprehension for distributed replay. In *NSDI*, volume 7, pages 285–298, 2007.
- [28] Chris Hawblitzel, Jon Howell, Manos Kapritsos, Jacob R Lorch, Bryan Parno, Michael L Roberts, Srinath Setty, and Brian Zill. IronFleet: proving practical distributed systems correct. In *Proceedings of the 25th Symposium on Operating Systems Principles (SOSP'15)*, pages 1–17, 2015.
- [29] Intel Corporation. *Intel 64 and IA-32 architectures software developer's manual*, volume 3 (3A, 3B, 3C & 3D): System Programming Guide. Intel, 2016.
- [30] Jaeger: Open Source, End-to-End Distributed Tracing. Retrieved May 2021 from <https://www.jaegertracing.io/>.
- [31] Jaeger. Jaeger 1.26 Documentation: Sampling. Retrieved in Sep 2021 from <https://www.jaegertracing.io/docs/1.26/sampling>, 2021.
- [32] Chris Jones, John Wilkes, Niall Murphy, and Cody Smith. *Site Reliability Engineering: How Google Runs Production Systems*. O'Reilly Media, 2016. <https://landing.google.com/sre/sre-book/chapters/service-level-objectives/>.
- [33] Jonathan Kaldor, Jonathan Mace, Michał Bejda, Edison Gao, Wiktor Kuropatwa, Joe O'Neill, Kian Win Ong, Bill Schaller, Pingjia Shan, Brendan Viscomi, et al. Canopy: An end-to-end performance tracing and analysis system. In *Proceedings of the 26th ACM Symposium on Operating Systems Principles (SOSP'17)*, pages 34–50, 2017.
- [34] Pravein Govindan Kannan, Nishant Budhdev, Raj Joshi, and Mun Choon Chan. Debugging transient faults in data centers using synchronized network-wide packet histories. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*, pages 253–268. USENIX Association, April 2021.
- [35] Pedro Las-Casas, Jonathan Mace, Dorgival Guedes, and Rodrigo Fonseca. Weighted Sampling of Execution Traces: Capturing More Needles and Less Hay. In *9th ACM Symposium on Cloud Computing (SOCC '18)*, 2018.
- [36] Pedro Las-Casas, Giorgi Papakerashvili, Vaastav Anand, and Jonathan Mace. Sifter: Scalable sampling for distributed traces, without feature engineering. In *Proceedings of the ACM Symposium on Cloud Computing (SOCC'19)*, pages 312–324, 2019.
- [37] K Rustan M Leino. Dafny: An automatic program verifier for functional correctness. In *International Conference on Logic for Programming Artificial Intelligence and Reasoning*, pages 348–370. Springer, 2010.
- [38] Jialin Li, Naveen Kr Sharma, Dan RK Ports, and Steven D Gribble. Tales of the Tail: Hardware, OS, and Application-Level Sources of Tail Latency. In *Proceedings of the 5th ACM Symposium on Cloud Computing (SoCC)*, 2014.
- [39] Lightstep. OpenTelemetry Issue #407: Tail-Based Sampling Scalability Issues. Retrieved May 2021 from <https://github.com/open-telemetry/opentelemetry-collector/issues/407>, 2020.
- [40] Liang Luo, Suman Nath, Lenin Ravindranath Sivalingam, Madan Musuvathi, and Luis Ceze. Troubleshooting transiently-recurring errors in production systems with blame-proportional logging. In *2018 USENIX Annual Technical Conference (USENIX ATC'18)*, pages 321–334, 2018.
- [41] Dan Luu. A simple way to get more value from tracing. Retrieved May 2021 from <https://danluu.com/tracing-analytics/>, 2020.
- [42] Jonathan Mace, Peter Bodik, Rodrigo Fonseca, and Madanlal Musuvathi. Retro: Targeted Resource Management in Multi-Tenant Distributed Systems. In *Proceedings of the 12th USENIX Symposium on Networked Systems Design and Implementation (NSDI '15)*, 2015.
- [43] Jonathan Mace and Rodrigo Fonseca. Universal context propagation for distributed system instrumentation. In *Proceedings of the 13th European Conference on Computer Systems (EuroSys'18)*, pages 1–18, 2018.

- [44] Jonathan Mace, Ryan Roelke, and Rodrigo Fonseca. Pivot Tracing: Dynamic Causal Monitoring for Distributed Systems. In *25th ACM Symposium on Operating Systems Principles (SOSP'15)*, 2015.
- [45] Gideon Mann, Mark Sandler, Darja Krushevska, Sudipto Guha, and Eyal Even-Dar. Modeling the parallel execution of black-box services. In *Proceedings of USENIX Workshop on Hot Topics in Cloud Computing (HotCloud'11)*, 2011.
- [46] John D. McCalpin. Memory bandwidth and machine balance in current high performance computers. *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter*, pages 19–25, December 1995.
- [47] Pulkit A Misra, María F Borge, Íñigo Goiri, Alvin R Lebeck, Willy Zwaenepoel, and Ricardo Bianchini. Managing tail latency in datacenter-scale file systems under production constraints. In *Proceedings of the Fourteenth EuroSys Conference 2019*, pages 1–15, 2019.
- [48] OpenTelemetry: An Observability Framework for Cloud-Native Software. Retrieved May 2021 from <http://opentelemetry.io/>.
- [49] OpenTelemetry Specification Issue 307: Allow samplers to be called during different moments in the Span lifetime. Retrieved May 2021 from <https://github.com/open-telemetry/opentelemetry-specification/issues/307>, 2019.
- [50] OpenTelemetry Enhancement Proposal 115: Allow Additional Sampling Hooks. Retrieved May 2021 from <https://github.com/open-telemetry/oteps/pull/115>, 2020.
- [51] OpenTracing: Vendor-Neutral APIs and Instrumentation for Distributed Tracing. Retrieved May 2021 from <http://opentracing.io/>.
- [52] Krzysztof Ostrowski, Gideon Mann, and Mark Sandler. Diagnosing latency in multi-tier black-box services. In *4th International Workshop on Large-Scale Distributed Systems and Middleware (LADIS'11)*, 2011.
- [53] Maulik Pandey. Building Netflix's Distributed Tracing Infrastructure. Retrieved May 2021 from <https://netflixtechblog.com/building-netflixs-distributed-tracing-infrastructure-bb856c319304>, 2019.
- [54] Austin Parker, Daniel Spoonhower, Jonathan Mace, Ben Sigelman, and Rebecca Isaacs. *Distributed Tracing in Practice: Instrumenting, Analyzing, and Debugging Microservices*. O'Reilly Media, 2020.
- [55] Raja R Sambasivan, Ilari Shafer, Jonathan Mace, Benjamin H Sigelman, Rodrigo Fonseca, and Gregory R Ganger. Principled workflow-centric tracing of distributed systems. In *Proceedings of the Seventh ACM Symposium on Cloud Computing*, pages 401–414, 2016.
- [56] Raja R Sambasivan, Alice X Zheng, Michael De Rosa, Elie Krevat, Spencer Whitman, Michael Stroucken, William Wang, Lianghong Xu, and Gregory R Ganger. Diagnosing performance changes by comparing request flows. In *8th USENIX Symposium on Networked Systems Design & Implementation (NSDI'11)*, volume 5, pages 1–1, 2011.
- [57] Yuri Shkuro. *Mastering Distributed Tracing*. Packt Publishing, Feb 2019.
- [58] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The Hadoop distributed file system. In *2010 IEEE 26th symposium on mass storage systems and technologies (MSST)*, pages 1–10. Ieee, 2010.
- [59] Benjamin H. Sigelman, Luiz André Barroso, Mike Burrows, Pat Stephenson, Manoj Plakal, Donald Beaver, Saul Jaskan, and Chandan Shanbhag. Dapper, a large-scale distributed systems tracing infrastructure. Technical report, Google, Inc., 2010.
- [60] Kun Suo, Jia Rao, Luwei Cheng, and Francis CM Lau. Time capsule: Tracing packet latency across different layers in virtualized systems. In *Proceedings of the 7th ACM SIGOPS Asia-Pacific Workshop on Systems*, pages 1–9, 2016.
- [61] Eno Thereska, Brandon Salmon, John Strunk, Matthew Wachs, Michael Abd-El-Malek, Julio Lopez, and Gregory R Ganger. Stardust: Tracking Activity in a Distributed Storage System. In *2006 ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '06)*, 2006.
- [62] Abhishek Verma, Luis Pedrosa, Madhukar R. Korupolu, David Oppenheimer, Eric Tune, and John Wilkes. Large-scale cluster management at Google with Borg. In *Proceedings of the 10th European Conference on Computer Systems (EuroSys'15)*, Bordeaux, France, 2015.
- [63] Yang Wu, Ang Chen, and Linh Thi Xuan Phan. Zeno: diagnosing performance problems with temporal provenance. In *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI'19)*, pages 395–420, 2019.
- [64] Stephen Yang, Seo Jin Park, and John Ousterhout. NanoLog: a nanosecond scale logging system. In *2018 USENIX Annual Technical Conference (ATC'18)*, pages 335–350, 2018.

- [65] Kundi Yao, Guilherme B. de Pádua, Weiyi Shang, Steve Sporea, Andrei Toma, and Sarah Sajedi. Log4perf: Suggesting logging locations for web-based systems' performance monitoring. In *Proceedings of the 2018 ACM/SPEC International Conference on Performance Engineering (ICPE'18)*, pages 127–138, 2018.
- [66] Zuoning Yin, Xiao Ma, Jing Zheng, Yuanyuan Zhou, Lakshmi N Bairavasundaram, and Shankar Pasupathy. An empirical study on configuration errors in commercial and open source systems. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP'11)*, pages 159–172, 2011.
- [67] Ding Yuan, Jing Zheng, Soyeon Park, Yuanyuan Zhou, and Stefan Savage. Improving software diagnosability via log enhancement. In *16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'11)*, 2011.
- [68] Yunqi Zhang, David Meisner, Jason Mars, and Lingjia Tang. Treadmill: Attributing the source of tail latency through precise load testing and statistical inference. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pages 456–468. IEEE, 2016.
- [69] Xu Zhao, Kirk Rodrigues, Yu Luo, Michael Stumm, Ding Yuan, and Yuanyuan Zhou. Log20: Fully automated optimal placement of log printing statements under specified overhead threshold. In *Proceedings of the 26th ACM Symposium on Operating Systems Principles (SOSP'17)*, pages 565–581, 2017.
- [70] Wenchao Zhou, Qiong Fei, Arjun Narayan, Andreas Haeberlen, Boon Thau Loo, and Micah Sherr. Secure network provenance. In *Proceedings of the twenty-third ACM symposium on operating systems principles*, pages 295–310, 2011.
- [71] Wenchao Zhou, Suyog Mapara, Yiqing Ren, Yang Li, Andreas Haeberlen, Zachary Ives, Boon Thau Loo, and Micah Sherr. Distributed time-aware provenance. *Proceedings of the VLDB Endowment*, 6(2):49–60, 2012.
- [72] Wenchao Zhou, Micah Sherr, Tao Tao, Xiaozhou Li, Boon Thau Loo, and Yun Mao. Efficient querying and maintenance of network provenance at internet-scale. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pages 615–626, 2010.
- [73] Zipkin: A Distributed Tracing System. Retrieved May 2021 from <http://zipkin.io/>.