

EXTENDED BERKELEY PACKET FILTER FOR INTRUSION DETECTION IMPLEMENTATIONS

COMP4906 HONOURS THESIS

William Findlay

March 20, 2020

Under the supervision of Dr. Anil Somayaji
Carleton University

Abstract

System introspection is becoming an increasingly attractive option for maintaining operating system stability and security. This is primarily due to the many recent advances in system introspection technology; in particular, the 2013 introduction of *Extended Berkeley Packet Filter* (*eBPF*) into the Linux Kernel [45, 46] along with the recent development of more usable interfaces such as the *BPF Compiler Collection* (*bcc*) [19] has resulted in a highly compelling, performant, and (perhaps most importantly) safe subsystem for both kernel and userland instrumentation.

The scope, safety, and performance of eBPF system introspection has potentially powerful applications in the domain of computer security. In order to demonstrate this, we present *ebpH*, an eBPF implementation of Somayaji’s [41] *Process Homeostasis* (*pH*). *ebpH* is an intrusion detection system (IDS) that uses eBPF programs to instrument system calls and establish normal behavior for processes, building a profile for each executable on the system; subsequently, *ebpH* can warn the user when it detects process behavior that violates the established profiles. Experimental results show that *ebpH* can detect anomalies in process behavior with negligible overhead. Furthermore, *ebpH*’s anomaly detection comes with zero risk to the system thanks to the safety guarantees of eBPF, rendering it an ideal solution for monitoring production systems.

This thesis will discuss the design and implementation of *ebpH* along with the technical challenges which occurred along the way. We will present experimental data and performance benchmarks that demonstrate *ebpH*’s ability to monitor process behavior with minimal overhead. Finally, we conclude with a discussion on the merits of eBPF IDS implementations and potential avenues for future work therein.

ebpH is licensed under GPLv2 and full source code is available at <https://github.com/willfindlay/ebph>.

Acknowledgments

First and foremost, I would like to thank my advisor, Anil Somayaji, for his tireless efforts to ensure the success of this project, as well as for providing the original design for pH along with invaluable advice and ideas. Implementing ebpH and writing this thesis has been a long process and not without its challenges. Dr. Somayaji's support and guidance have been quintessential to the success of this undertaking.

I would also like to thank the members and contributors of the *Iovisor Project*, especially Yonghong Song (<https://github.com/yonghong-song>) and Teng Qin (<https://github.com/palmtenor>) for their guidance and willingness to respond to issues and questions related to the *bcc* project. Brendan Gregg's (<http://www.brendangregg.com>; <https://github.com/brendangregg>) writings and talks have been a great source of inspiration, especially with respect to background research. Sasha Goldstein's (<https://github.com/goldshtn>) *syscount.py* was an invaluable basis for my earliest proof-of-concept experimentation, although none of that original code has made it into this iteration of ebpH.

For their love and tremendous support of my education, I would like to thank my parents, Mark and Terri-Lyn. Without them, I am certain that none of this would have been possible. I would additionally like to thank my mother for suffering through the first draft of this thesis, and finding the many errors that come with writing a paper this large in Vim with no grammar checker.

Finally, I want to thank my dear friend, Amanda, for all the support she has provided me throughout my university career. I couldn't have made it this far without you.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction and Motivation | 1 |
| 2 | Background | 2 |
| 2.1 | An Overview of the Linux Tracing Landscape | 2 |
| 2.1.1 | Dtrace | 4 |
| 2.2 | Classic BPF | 6 |
| 2.3 | eBPF: Linux Tracing Superpowers | 7 |
| 2.3.1 | How eBPF Works at a High Level | 8 |
| 2.3.2 | The Verifier: The Good, the Bad, and the Ugly | 10 |
| 2.4 | System Calls | 11 |
| 2.5 | Intrusion Detection | 12 |
| 2.5.1 | A Survey of Data Collection in Intrusion Detection Systems | 13 |
| 2.6 | Process Homeostasis | 15 |
| 2.6.1 | Anomaly Detection Through Lookahead Pairs | 16 |
| 2.6.2 | Homeostasis Through System Call Delays | 16 |
| 3 | Implementing ebpH | 22 |
| 3.1 | Userspace Components | 23 |
| 3.1.1 | The ebpH Daemon | 23 |
| 3.1.2 | ebph-ps | 24 |
| 3.1.3 | ebph-admin | 24 |
| 3.2 | ebpH Profiles | 24 |
| 3.2.1 | Writing Profiles to Disk and Reading Profiles from Disk | 25 |
| 3.3 | Tracing Processes | 27 |
| 3.3.1 | Profile Creation and Association | 27 |
| 3.3.2 | Profile Association and Sequence Duplication | 28 |
| 3.3.3 | Dealing with Signal Handlers and Non-Determinism | 29 |
| 3.3.4 | Reaping Processes | 30 |
| 3.4 | Training, Testing, and Anomaly Detection | 30 |
| 3.4.1 | A Simple Example of ebpH Anomaly Detection | 30 |
| 3.5 | Issuing Commands to ebpH | 33 |
| 3.5.1 | Setting Runtime Parameters | 33 |
| 3.5.2 | Examining Profiles and Processes | 33 |
| 3.5.3 | Issuing More Complex Commands | 33 |
| 4 | Technical Challenges of an eBPF Intrusion Detection System | 33 |
| 4.1 | Soothing the Verifier | 33 |
| 5 | Measuring ebpH's Overhead | 34 |
| 5.1 | Methodology | 35 |
| 5.2 | Results | 36 |
| 5.2.1 | bronte Macro-Benchmark | 37 |
| 5.2.2 | homeostasis Macro-Benchmark | 41 |

| | | |
|----------|---|-----------|
| 5.2.3 | arch-3day Macro-Benchmark: Using ebpH on a Personal Computer . | 41 |
| 5.2.4 | bronte-lmbench: Measuring System Latency with lmbench Micro-Benchmark | 41 |
| 5.3 | Comparing Results with the Original pH | 44 |
| 6 | Future Work | 44 |
| 6.1 | Security Analysis | 44 |
| 6.2 | Controlling for Further Sources of Non-Deterministic Behavior | 44 |
| 6.3 | Automating ebpH Response | 44 |
| 6.4 | General System Introspection and the Future of ebpH | 44 |
| 7 | Methodology and Future Work | 44 |
| 7.1 | Planned Testing Strategy | 44 |
| 7.1.1 | Gathering and Analyzing Profile Data | 45 |
| 7.1.2 | Gathering and Analyzing Performance Data | 45 |
| 7.2 | Potential Improvements to ebpH | 46 |
| 7.2.1 | Delaying System Calls | 46 |
| 7.2.2 | Measuring Other System Behavior | 47 |
| 7.2.3 | Overall System Visibility | 48 |
| | References | 49 |
| | Appendices | 54 |
| A | eBPF Design Patterns | 54 |
| B | bpfbench Source Code | 55 |

List of Figures

| | | |
|-----|---|----|
| 2.1 | A high level overview of the broad categories of Linux instrumentation . . . | 4 |
| 2.2 | Comparing Dtrace and eBPF functionality with respect to API design | 5 |
| 2.3 | A high level overview of various eBPF use cases | 7 |
| 2.4 | Basic topology of eBPF with respect to userland and the kernel | 9 |
| 2.5 | The set participation of valid C and eBPF programs. | 18 |
| 2.6 | Complexity and verifiability of eBPF programs. | 19 |
| 2.7 | An overview of the basic categories of IDS. | 20 |
| 2.8 | An overview of the most common data collection categories and subcategories in IDS | 21 |
| 3.1 | The dataflow between various components of ebpH. | 22 |
| 3.2 | A sample (read, close) lookahead pair in the ebpH profile for ls. | 26 |
| 3.3 | Two sample (write, write) lookahead pairs in the ebpH profile for anomaly.c. | 32 |
| 5.1 | Top 20 most frequent system calls with base times of under $3\mu s$ from the bronte-7day dataset (after discarding outliers). Smaller difference in times is better. Standard error is given as error bars. | 40 |
| 5.2 | Mean system call times from the bronte-1mbench dataset. Smaller difference in times is better. Standard deviation is given as error bars. | 42 |
| 5.3 | Mean process creation times from the bronte-1mbench dataset. Smaller differ- ence in times is better. Standard deviation is given as error bars. | 43 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | A summary of various system introspection technologies available for GNU/Linux systems. | 3 |
| 2.2 | Various map types [15, 19] available in eBPF programs, as of kernel 5.3. . . . | 10 |
| 3.1 | Main event categories in ebpH. | 24 |
| 3.2 | Important system calls in ebpH. | 28 |
| 5.1 | Systems used for the collection of ebpH benchmarking data. | 35 |
| 5.2 | ebpH macro-benchmarking datasets. | 36 |
| 5.3 | ebpH micro-benchmarking datasets. | 37 |
| 5.4 | Top 20 most frequent system calls with base times of under $3\mu s$ from the bronte-7day dataset (after discarding outliers). Smaller overhead is better. . | 38 |
| 5.5 | Results of the system call benchmarks from the bronte-1mbench dataset. Standard deviations are given in parentheses, smaller overhead is better. | 41 |
| 5.6 | Results of the process creation benchmarks from the bronte-1mbench dataset. Standard deviations are given in parentheses, smaller overhead is better. . . | 41 |
| 7.1 | A summary of the various systems that will be used to test ebpH | 44 |

List of Listings

| | | |
|-----|--|----|
| 3.1 | A simplified definition of the ebpfH profile struct. | 25 |
| 3.2 | A simplified definition of the ebpfH process struct. | 27 |
| 3.3 | A simple program to demonstrate anomaly detection in ebpfH. | 30 |
| 3.4 | The flagged anomaly in the <code>anomaly</code> binary as shown in the ebpfH logs. Note that ebpfH also logs the offending sequence, reordering it so that most recent system calls appear on the right. | 31 |
| A.1 | Handling large datatypes in eBPF programs. | 54 |
| B.1 | The eBPF component of <code>bpfbench</code> | 55 |

1 Introduction and Motivation

As our computer systems grow increasingly complex, so too does it become more difficult to gauge precisely what they are doing at any given moment. Modern computers are often running hundreds, if not thousands of processes at any given time, the vast majority of which are running silently in the background. As a result, users often have a very limited notion of what exactly is happening on their systems, especially beyond that which they can actually see on their screens. An unfortunate corollary to this observation is that users *also* have no way of knowing whether their system may be *misbehaving* at a given moment, whether due to a malicious actor, buggy software, or simply some unfortunate combination of circumstances.

Recently, a lot of work has been done to help bridge this gap between system state and visibility, particularly through the introduction of powerful new tools such as *Extended Berkeley Packet Filter* (eBPF). Introduced to the Linux Kernel in a 2013 RFC and subsequent kernel patch [45, 46], eBPF offers a promising interface for kernel introspection, particularly given its scope and unprecedented level of safety therein; although eBPF can examine any data structure or function in the kernel through the instrumentation of tracepoints, its safety is guaranteed via a bytecode verifier. What this means in practice is that we effectively have unlimited, highly performant, production-safe system introspection capabilities that can be used to monitor as much or as little system state as we desire.

Certainly, eBPF offers unprecedented system state visibility, but this is only scratching the surface of what this technology is capable of. With limitless tracing capabilities, we can construct powerful applications to enhance system security, stability, and performance. In theory, these applications can perform much of their work autonomously in the background, but are equally capable of functioning in a more interactive role, keeping the end user informed about changes in system state, particularly if these changes in state are undesired. To that end, I propose *ebpH* (a portmanteau of eBPF and pH), an intrusion detection system based entirely on eBPF that monitors process state in the form of system call sequences. By building and maintaining per-executable behavior profiles, ebpH can dynamically detect when processes are behaving outside of the status quo, and notify the user so that they can understand exactly what is going on.

A prototype of ebpH has been written using the Python interface provided by the *BPF Compiler Collection* (*bcc*) [19], and preliminary tests show that it is capable of monitoring system state under moderate to heavy workloads with negligible overhead. What's more, zero kernel panics occurred during ebpH's development and early testing, which simply would not have

been possible without the safety guarantees that eBPF provides. The rest of this proposal will cover the necessary background material required to understand ebpH, describe several aspects of its implementation, including the many findings and pitfalls encountered along the way, and discuss the planned methodology for testing and iterating on this prototype going forward.

2 Background

In the following sections, I will provide the necessary background information needed to understand ebpH; this includes an overview of system introspection and tracing techniques on Linux including eBPF itself, and some background on system calls and intrusion detection.

While my work is primarily focused on the use of eBPF for maintaining system security and stability, the working prototype for ebpH borrows heavily from Anil Somayaji’s *pH* or *Process Homeostasis* [41], an anomaly-based intrusion detection and response system written as a patch for Linux Kernel 2.2. As such, I will also provide some background on the original pH system and many of the design choices therein.

2.1 An Overview of the Linux Tracing Landscape

System introspection is hardly a novel concept; for years, developers have been thinking about the best way to solve this problem and have come up with several unique solutions, each with a variety of benefits and drawbacks. Table 2.1 presents an overview of some prominent examples relevant to GNU/Linux systems.

These technologies can, in general, be classified into a few broad categories (Figure 2.1), albeit with potential overlap depending on the tool:

- (1) Userland libraries.
- (2) Ptrace-based instrumentation.
- (3) Loadable kernel modules.
- (4) Kernel subsystems.

Applications such as strace [49, 50] which make use of the ptrace system call are certainly a viable option for limited system introspection with respect to specific processes. However, this does not represent a complete solution, as we are limited to monitoring the system calls made by a process to communicate with the kernel, its memory, and the state of its registers, rather than the underlying kernel functions themselves [36]. The scope of ptrace-based solutions is also limited by ptrace’s lack of scalability; ptrace’s API is conducive to

Table 2.1: A summary of various system introspection technologies available for GNU/Linux systems.

| Name | Interface and Implementation | Citations |
|--------------|--|------------------|
| strace | Uses the ptrace system call to trace userland processes | [49, 50] |
| ltrace | Uses the ptrace system call to trace library calls in userland processes | [5, 39] |
| SystemTap | Dynamically generates loadable kernel modules for instrumentation; newer versions can optionally use eBPF as a backend instead | [31, 37] |
| ftrace | Sysfs pseudo filesystem for tracepoint instrumentation located at <code>/sys/kernel/debug/tracing</code> | [38] |
| perf_events | Linux subsystem that collects performance events and returns them to userspace | [56] |
| LTtng | Loadable kernel modules, userland libraries | [26] |
| dtrace4linux | A Linux port of DTrace via a loadable kernel module | [10] |
| sysdig | Loadable kernel modules for system monitoring; native support for containers | [51] |
| eBPF | In-kernel virtual machine for running pre-verified bytecode | [11, 19, 45, 46] |

tracing single processes at a time rather than tracing processes system wide. Its limited scale becomes even more obvious when considering the high amount of context-switching between kernel space and user space required when tracing multiple processes or threads, especially when these processes and threads make many hundreds of system calls per second [23].

Although library call instrumentation through software such as ltrace [5, 39] does not necessarily suffer from the same performance issues as described above, it still constitutes a suboptimal solution for many use cases due to its limited scope. In order to be effective and provide a complete picture of what exactly is going on during a given process’ execution, library tracing needs to be combined with other solutions. In fact, ltrace does exactly this; when the user specifies the `-S` flag, ltrace uses the ptrace system call to provide strace-like system call tracing functionality.

LKM-based implementations such as sysdig [51] and SystemTap [37] offer an extremely deep and powerful tracing solution given their ability to instrument the entire system, including the kernel itself. Their primary detriment is a lack of safety guarantees with respect to the modules themselves. No matter how vetted or credible a piece of software might be, running it natively in the kernel always comports with an inherent level of risk; buggy code might cause system failure, loss of data, or other unintended and potentially catastrophic consequences.

Custom tracing solutions through kernel modules carry essentially the same risks. No sane

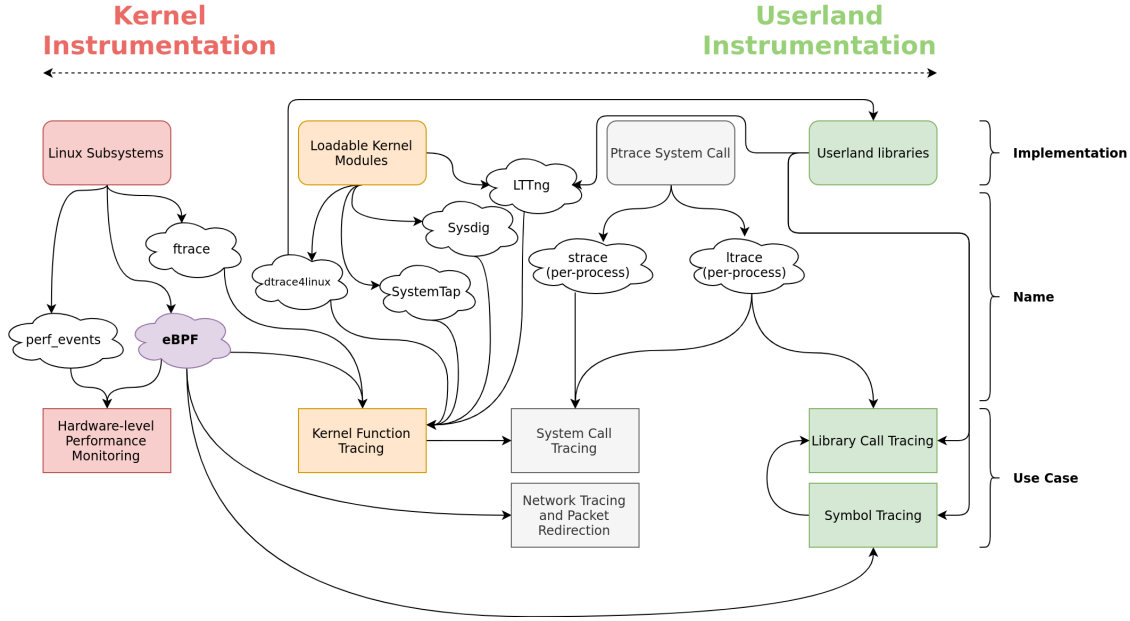


Figure 2.1: A high level overview of the broad categories of Linux instrumentation. This does not represent a complete picture of all available tools and interfaces, but instead presents many of the most popular ones. Note how eBPF covers every presented use case.

manager would consent to running untrusted, unvetted code natively in the kernel of a production system; the risks are simply too great and far outweigh the benefits. Instead, such code must be carefully examined, reviewed, and tested, a process which can potentially take months. What’s more, even allowing for a careful testing and vetting process, there is always some probability that a bug can slip through the cracks, resulting in the catastrophic consequences outlined above.

Built-in kernel subsystems for instrumentation seem to be the most desirable choice of any of the presented solutions. In fact, eBPF [45, 46] itself constitutes one such solution. However, for the time being, we will focus on a few others, namely ftrace [38] and perf_events [56] (eBPF programs actually *can* and *do* use both of these interfaces anyway). While both of these solutions are safe to use (assuming we trust the user), they suffer from limited documentation and relatively poor user interfaces. These factors in tandem mean that ftrace and perf_events, while quite useful for a variety of system introspection needs, are less extensible than other approaches.

2.1.1 Dtrace

It is worth spending a bit more time comparing eBPF with Dtrace, as both APIs are quite full-featured and designed with similar functionality in mind. The original Dtrace [4] was

designed in 2004 for Solaris and lives on to this day in several operating systems, including Solaris, FreeBSD, MacOS X [16], and Linux [10] (we will examine the dtrace4linux implementation with more scrutiny later in this section).

In general, the original Dtrace and the current version of eBPF share much of the same family and cover similar use cases [4, 45, 46]. This includes perhaps most notably dynamic instrumentation in both userspace and kernelspace, arbitrary context instrumentation (i.e. the ability to instrument essentially any aspect of the system), and guarantees of safety and data integrity. The difference between the two systems generally boils down to the following three points:

- (1) eBPF supports a superset of Dtrace’s functionality;
- (2) Dtrace provides only a high level interface, while eBPF provides both low level and high level interfaces (see Figure 2.2); and
- (3) eBPF is natively supported in Linux, while Dtrace ports are purely LKM-based.

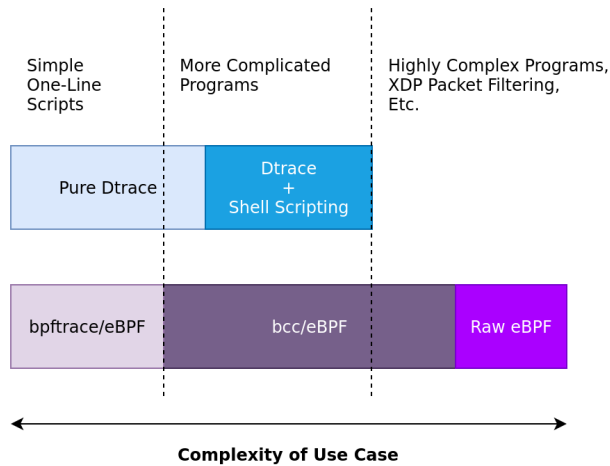


Figure 2.2: Comparing Dtrace and eBPF functionality with respect to API design (adapted from [13]). Note that eBPF covers more complex use cases and supports both low level and high level APIs. Dtrace needs to be used in tandem with shell scripting to support more complex use cases.

dtrace4linux [10] is a free and open source port of Sun’s Dtrace [4] for the Linux Kernel, implemented as a loadable kernel module (LKM). While Dtrace offers a powerful API for full-system tracing, its usefulness is, in general, eclipsed by that of eBPF [13] and requires extensive shell scripting for use cases beyond one-line tracing scripts. In contrast, with the help of powerful and easy to use front ends like bcc [19], developing complex eBPF programs for a wide variety of use cases is becoming an increasingly painless process.

Not only does eBPF cover more complex use cases than Dtrace, but it also provides support

for simple one-line programs through tools like `bpftool` [13, 20] which has been designed to provide a high-level Dtrace-like tracing language for Linux using eBPF as a backend. Although `bpftool` only provides a subset of Dtrace’s functionality [13], its feature set has been carefully curated in order to cater to the most common use cases and more functionality is being added on an as-needed basis.

Additional work is being done to fully reimplement Dtrace as a new BPF program type [54] which will further augment eBPF’s breadth and provide full backward compatibility for existing Dtrace scripts to work with eBPF. This seems to be by far the most promising avenue for Linux Dtrace support thus far, as it seeks to combine the higher level advantages of Dtrace with the existing eBPF virtual machine.

2.2 Classic BPF

In 1992, McCanne and Jacobson [27] introduced the original BPF¹ or *Berkeley Packet Filter* as a mechanism for capturing, monitoring, and filtering network traffic in the BSD kernel. Classic BPF’s primary insights were two-fold:

- (1) network traffic events are *frequent* and *fast*, and therefore an efficient filtering mechanism was needed; and
- (2) a limited, register-based bytecode being run in an in-kernel virtual machine provides precisely the mechanism described in point (1).

The virtual machine described above was used to implement the *filter* component of BPF, while in-kernel network function tracepoints implemented the *tap* component. The tap forwarded packet events to the filter, which then decided what to do with the packets according to a user-defined BPF program. McCanne and Jacobson showed that their approach was much faster than other contemporary packet filtering techniques, namely NIT [34] and CSPF [32].

While Classic BPF is certainly a powerful technique for filtering packets, Starovoitov [45, 46] realized that its tap and filter mechanism represented a desirable approach for general system introspection. Therefore, in 2013, he proposed *Extended BPF* (eBPF), a superset of Classic BPF, which vastly increased the capabilities of the original BPF virtual machine.

Since its original introduction, eBPF has offered a consistent, powerful, and production-safe mechanism for general Linux introspection and continues to improve rapidly over time. We discuss eBPF in more detail in the following section.

¹Hereafter, we will refer to the original BPF as *Classic BPF* to avoid confusion with eBPF and the BPF programming paradigm.

2.3 eBPF: Linux Tracing Superpowers

In 2016, eBPF was described by Brendan Gregg [12] as nothing short of *Linux tracing superpowers*. I echo that sentiment here, as it summarizes eBPF’s capabilities perfectly. Through eBPF programs, we can simultaneously trace userland symbols and library calls, kernel functions and data structures, and hardware performance. What’s more, through an even newer subset of eBPF, known as *XDP* or *Express Data Path* [17], we can inspect, modify, redirect, and even drop packets entirely before they even reach the main kernel network stack. Figure 2.3 provides a high level overview of these use cases and the corresponding eBPF instrumentation required.

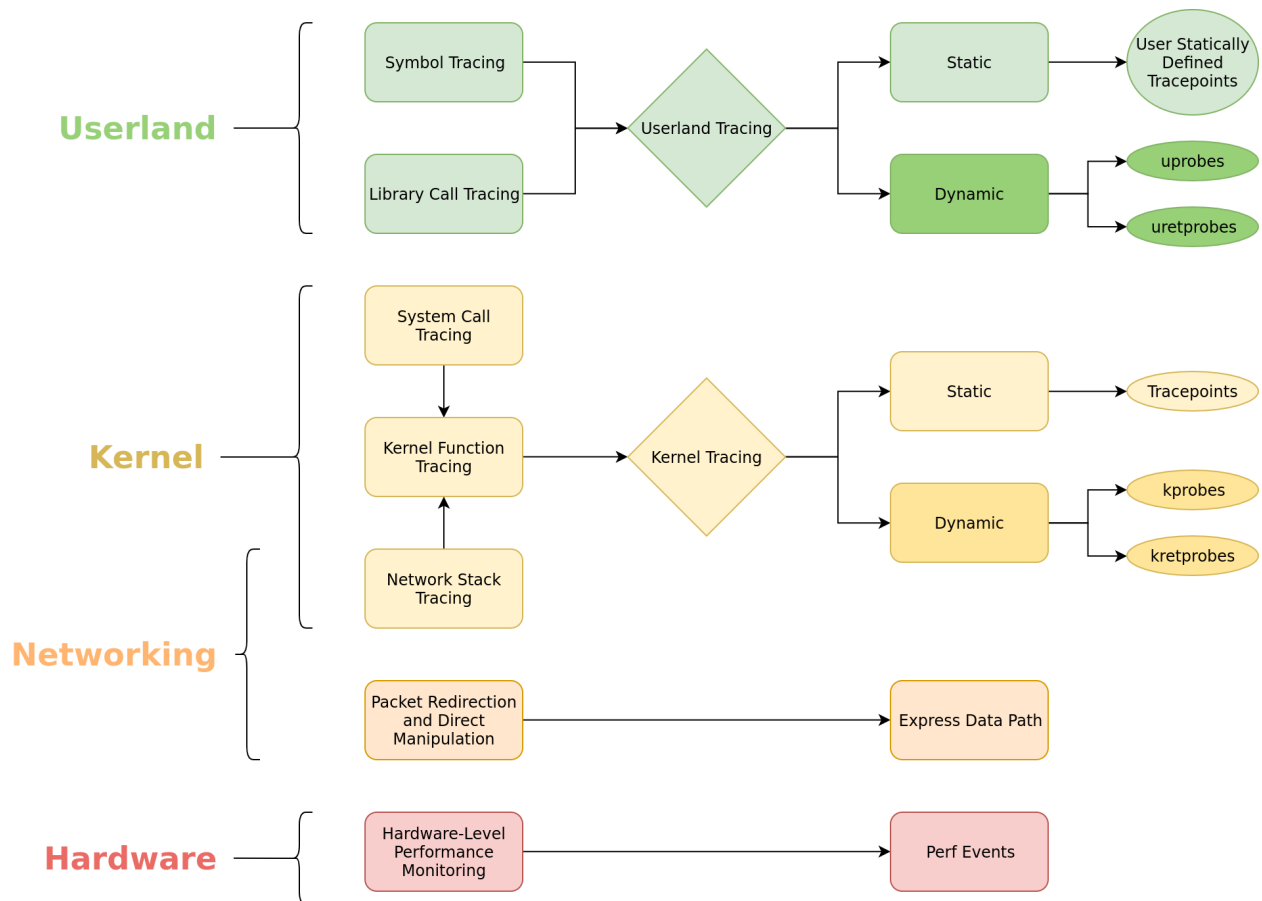


Figure 2.3: A high level overview of various eBPF use cases. Note the high level of flexibility that eBPF provides with respect to system tracing.

The advantages of eBPF extend far beyond scope of traceability; eBPF is also extremely performant, and runs with guaranteed safety. In practice, this means that eBPF is an ideal tool for use in production environments and at scale.

Safety is guaranteed with the help of an in-kernel verifier that checks all submitted bytecode

before its insertion into the BPF virtual machine. While the verifier does limit what is possible (eBPF in its current state is *not* Turing complete), it is constantly being improved; for example, a recent patch [48] that was mainlined in the Linux 5.3 kernel added support for verified bounded loops, which greatly increases the computational possibilities of eBPF. The verifier will be discussed in further detail in Section 2.3.2.

eBPF’s superior performance can be attributed to several factors. On supported architectures,² eBPF bytecode is compiled into machine code using a *just-in-time* (*JIT*) compiler; this both saves memory and reduces the amount of time it takes to insert an eBPF program into the kernel. Additionally, since eBPF runs in-kernel and communicates with userland via direct map access and perf events, the number of context switches required between userland and the kernel is greatly diminished, especially compared to approaches such as the ptrace system call.

2.3.1 How eBPF Works at a High Level

From the perspective of a user, the eBPF workflow is surprisingly simple. Users can elect to write eBPF bytecode directly (not recommended) or use one of many front ends to write in higher level languages that are then used to generate the respective bytecode. `bcc` [19] offers front ends for several languages including Python, Go, and C++; users write eBPF programs in C and interact with `bcc`’s API in order to generate eBPF bytecode and submit it to the kernel.

Figure 2.4 presents an overview of eBPF’s architecture and dataflow, including the interaction between userspace programs, eBPF programs in kernelspace, and the rest of the kernel. This interaction occurs via the `bpf(2)` system call [3] which is used to load and verify BPF programs, issue commands to BPF programs, and interact with eBPF maps. These maps are the mechanism for sending data between BPF programs and other BPF programs or BPF programs and userspace.

There are several map types available in eBPF which cover a wide variety of use cases. These map types along with a brief description are provided in Table 2.2 [3, 8, 19]. Thanks to this wide arsenal of maps, eBPF developers have a powerful set of both general-purpose and specialized data structures at their disposal; as we will see in coming sections, many of these maps are quite versatile and have use cases beyond what might initially seem pertinent. For example, the `ARRAY` map type may be used to initialize large data structures to be copied into a general purpose `HASH` map (refer to Listing A.1 in Appendix A). This can be effectively used

²x86-64, SPARC, PowerPC, ARM, arm64, MIPS, and s390 [8]

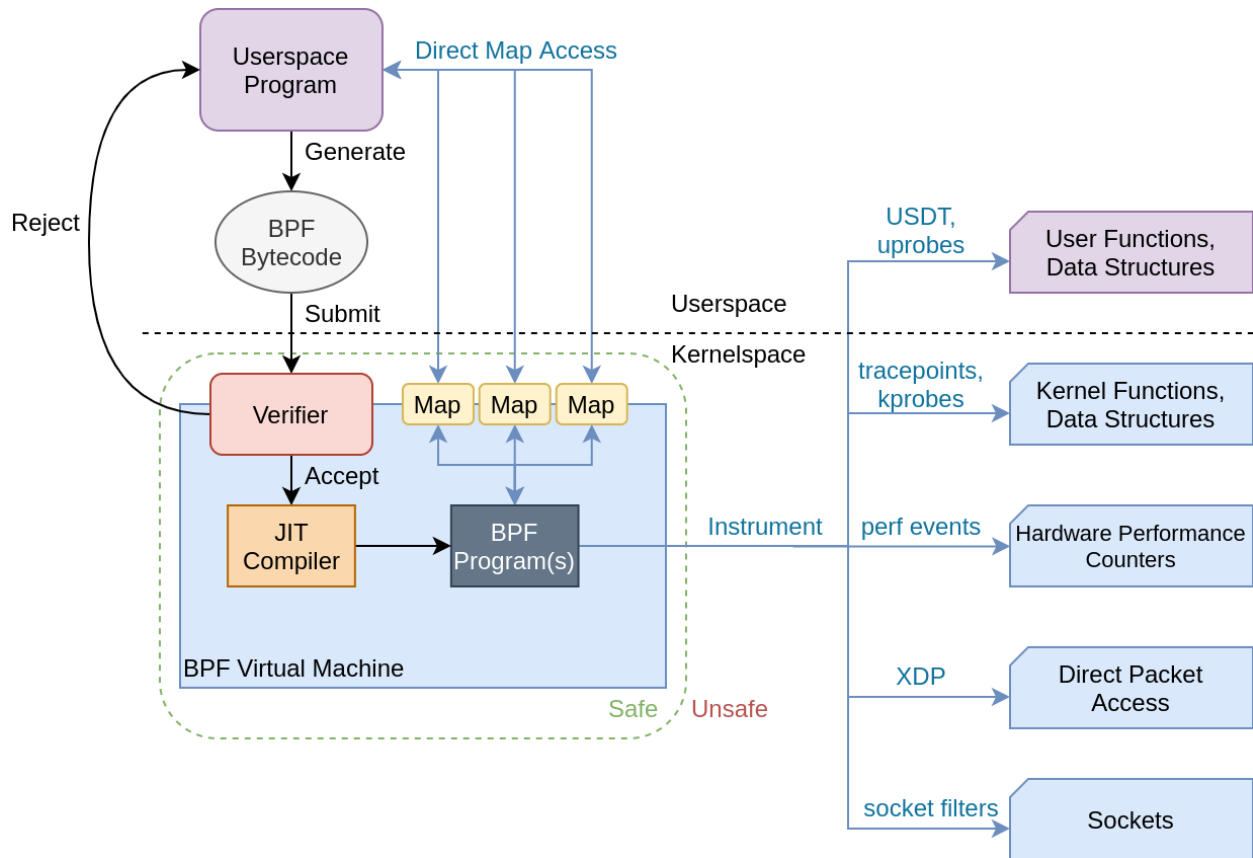


Figure 2.4: Basic topology of eBPF with respect to userland and the kernel. Note the bidirectional nature of dataflow between userspace and kernelspace using maps.

to bypass the verifier’s stack space limitations, which are discussed in detail in Section 2.3.2.

Table 2.2: Various map types [15, 19] available in eBPF programs, as of kernel 5.3.

| Map Type | Description |
|------------------|---|
| HASH | A hashtable of key-value pairs |
| ARRAY | An array indexed by integers; members are zero-initialized |
| PROG_ARRAY | A specialized array to hold file descriptors to other BPF programs; used for tail calls |
| PERF_EVENT_ARRAY | Holds perf event counters for hardware monitoring |
| PERCPU_HASH | Like HASH but stores a different copy for each CPU context |
| PERCPU_ARRAY | Like ARRAY but stores a different copy for each CPU context |
| STACK_TRACE | Stores stack traces for userspace or kernelspace functions |
| CGROUP_ARRAY | Stores pointers to cgroups |
| LRU_HASH | Like a HASH except least recently used values are removed to make space |
| LRU_PERCPU_HASH | Like LRU_HASH but stores a different copy for each CPU context |
| LPM_TRIE | A "Longest Prefix Matching" trie optimized for efficient traversal |
| ARRAY_OF_MAPS | An ARRAY of file descriptors into other maps |
| HASH_OF_MAPS | A HASH of file descriptors into other maps |
| DEVMAP | Maps the ifindex of various network devices; used in XDP programs |
| SOCKMAP | Holds references to sock structs; used for socket redirection |
| CPUMAP | Allows for redirection of packets to remote CPUs; used in XDP programs |

2.3.2 The Verifier: The Good, the Bad, and the Ugly

The verifier is responsible for eBPF's unprecedented safety, one of its most attractive qualities with respect to system tracing. While this verifier is quintessential to the safety of eBPF given its impressive scope and power, it is not without its drawbacks. In this section, we describe how the verifier works, its nuances and drawbacks, and recent work that has been done to improve the verifier's support for increasingly complex eBPF programs.

Proving the safety of arbitrary code is by definition a difficult problem. This is thanks in part to theoretical limitations on what we can actually prove; a famous example is the halting problem described by Turing circa 1937 [53]. This difficulty is further compounded by stricter requirements for safety in the context of an eBPF program; in fact, the problem that we are effectively trying to solve is one of *untrusted* code running in the kernel, an implicitly trusted environment.

To illustrate the importance of this problem of safety with respect to eBPF, let us consider a simple example. We will again consider the halting problem described above. Suppose we have two eBPF programs, program *A* and program *B*, that each hook onto a mission-critical kernel function (`schedule()`, for example). The only difference between these two programs is that program *A* always terminates, while program *B* runs forever without stopping. Program *B* effectively constitutes a denial of service attack [18] on our system, intentional or otherwise; allowing untrusted users to load this code into our kernel spells immediate disaster for our

system.

While we have established that verifying the safety of eBPF programs is an important problem to solve, the question remains as to whether it is *possible* to solve. For the reasons outlined above, this problem should intuitively seem impossible, or at least far more difficult than should be feasible. So, what can we do? The answer is to *change the rules* to make it easier. In particular, while it is difficult to prove that the set of all possible eBPF programs are safe, it is much easier³ to prove this property for a subset of all eBPF programs. Figure 2.5 depicts the relationship between potentially valid eBPF code and verifiably valid eBPF code.

The immediate exclusion of eBPF programs meeting certain criteria is the crux of eBPF’s safety guarantees. Unfortunately, it also rather intuitively limits what we are actually able to do with eBPF programs. In particular, eBPF is not a Turing complete language; it prohibits arbitrary jump instructions, cycles in execution graphs, and unverified memory access. Further, we limit stack allocations to only 512 bytes – far too small for many practical use cases. From a security perspective, these limitations are a *good thing*, because they allow us to immediately exclude eBPF programs with unverifiable safety; but from a usability standpoint, particularly that of a new eBPF developer, the trade-off is not without its drawbacks.

Fortunately, the eBPF verifier is getting better over time (Figure 2.6). When we say *better*, what we mean is that it is able to prove the safety of increasingly complex programs. Perhaps the best example of this steady improvement is a recent kernel patch [48] that added support for bounded loops in eBPF programs. With this patch, the set of viable eBPF programs was *greatly* increased; in fact, ebpfH in its current incarnation relies heavily on bounded loop support. Prior to bounded loops, eBPF programs relied on *unrolling* loops at compile time, a technique that was both slow and highly limited. This is just one example of the critical work that is being done to improve the verifier and thus improve eBPF as a whole.

2.4 System Calls

ebpfH (and the original pH system upon which it is based) works by instrumenting *system calls* in order to establish behavioral patterns for all binaries running on the system. Understanding pH and ebpfH requires a reliable mental model of what a system call is and how programs use them to communicate with the kernel.

At the time of writing this paper, the Linux Kernel [52] supports an impressive 436 distinct system calls, and this number generally grows with subsequent releases. In general, userspace

³Easier here means *computationally easier*, certainly not trivial.

libraries such as the C standard library implement a subset of these system calls, with the exact specifications varying depending on architecture. These system calls are used to request services from the operating system kernel; for example, a program that needs to write to a file would make an `open` call to receive a file descriptor into that file, followed by one or more `write` calls to write the necessary data, and finally a `close` call to clean up the file descriptor. These system calls form the basis for much of our process behavior, from I/O as seen above, to process management, memory management, and even the execution of binaries themselves.

Through the instrumentation of system calls, we can establish a clear outline of exactly how a process is behaving, the critical operations it needs to perform, and how these operations interact with one another. In fact, system call-based instrumentation forms a primary use case for several of the tracing technologies previously discussed in Section 2.1, perhaps most notably *strace*. We will discuss the behavioral implications of system calls further in Section 2.6.1.

2.5 Intrusion Detection

At a high level, intrusion detection systems (IDS) strive to monitor systems at a particular level and use observed data to make decisions about the legitimacy of these observations [22]. Intrusion detection systems can be broadly classified into several categories based on data collection, detection technique(s), and response. Figure 2.7 presents a broad and incomplete overview of these categories.

In general, intrusion detection systems can either attempt to detect anomalies (i.e. mismatches in behavior when compared to normally observed patterns) or misuse, which generally refers to matching known attack patterns to observed data [22]. In general, anomaly-based approaches cover a wider variety of attacks while misuse-based approaches tend to yield fewer false positives. Misuse-based approaches can also be further broken down into specification-based and signature-based, which deal in behavioral blacklists and whitelists respectively. A hybrid approach between any of these techniques is also possible.

Data collection is generally either host-based or network based. Network-based IDSes examine network traffic and analyze it to detect attacks or anomalies. In contrast, host-based IDSes analyze the state of the local system [22, 41].

Responses can vary significantly based on the system, but can be classified into two main categories: alerts and counter-attacks. Systems can either choose to alert an administrator about the potential issue, or choose to mount counter-measures to defeat or mitigate the

perceived attack [22]. Naturally, systems also have the option to take a hybrid approach here.

Using the above metrics, ebpH can be broadly classified as a host-based anomaly detection system that responds to anomalies by issuing alerts. This is generally quite similar to the original pH (Section 2.6) with one major exception: As we will see, the original pH also responds to anomalies by delaying system calls [41] and preventing anomalous `execves`. Implementing this functionality in ebpH is a topic for future work (refer to Section 7.2).

2.5.1 A Survey of Data Collection in Intrusion Detection Systems

We have presented the general classification of intrusion detection systems through the establishment of three core elements of an IDS and several categories therein. As it relates to eBPF, the *data collection* component of intrusion detection systems is of particular interest; what is especially exciting about eBPF is its impressive scope, safety, and performance with respect to general system introspection; this presents a perfect trifecta of traits for collecting system data. As such, it is worth examining data collection techniques from various intrusion detection systems in more detail.

We have established that data collection in intrusion detection systems can primarily be separated into two relatively broad categories:

- (1) host-based data collection which involves collecting data about the use and behavior of a local machine; and
- (2) network-based data collection which involves monitoring network traffic and looking for established patterns.

While the above two categories are generally sufficient for understanding intrusion detection at a high level, there are in fact several distinct subcategories therein. Figure 2.8 presents an overview of the most common data collection subcategories in IDSes.

Internal and External Sensors. Kerschbaum et al. [24] introduce the concept of *internal* sensors for intrusion detection and contrast them with the far more popular *external* sensors. An internal sensor by definition is included in the source code of a monitored component, while an external sensor is implemented outside of the monitored component. These two categories of sensors each present unique advantages and disadvantages [24]. In particular, external sensors are easily modifiable and extensible, although they introduce potential delays, and are generally weaker to tampering by intruders; internal sensors minimize overhead (assuming correct implementation) and are much more resistant to intruder

tampering, but suffer from reduced portability, difficulty in implementation, and may incur severe performance penalties if implemented incorrectly.

eBPF would fall under the internal sensor classification [24] due to its implementation within the Linux Kernel; however, eBPF presents a rather unique case, as it overcomes many of the disadvantages proposed by Kerschbaum et al. while maintaining the advantages. Specifically, eBPF is completely application transparent, portable to any modern version of Linux⁴, easy to update and modify, and has guaranteed performance and safety due to the verifier.

Internal Host-Based Approaches. System call pattern analysis was examined in detail by Forrest et al. [9] and culminated in the development of the original pH system [41] on which ebpH is based. Somayaji and Inoue [42] compared two methods of organizing system call data for intrusion detection (full sequences and lookahead pairs), which we will discuss further in Section 2.6.1.

Kerschbaum et al. also describe a generic method of application-specific internal sensors through the addition of audit checks in program source code [24]. However, the primary caveat here is that such checks need to be integrated into a specific application early enough such that refactoring is minimized [24, 35]. This approach is also far less generic than other internal sensor approaches described here.

Another potential internal source for data is through host-based network stack introspection. Classic BPF [27] and eBPF/XDP [17, 19, 45, 46] are quite excellent at this. Host-based network introspection allows the analysis of network traffic at various points in the kernel’s networking stack, and XDP packet redirection [17] allows fast detection and response before a packet even reaches the main networking stack.

ebpH itself constitutes an internal host-based approach; that is, it uses eBPF for in-kernel instrumentation of system calls (internal) on a given host (host-based). As we discuss in Section 7.2, a potential avenue for future research in ebpH is moving beyond system call monitoring to *general system introspection* (c.f. Figure 2.8). This is specifically a possibility due to eBPF’s unique classification as an internal sensor capable of monitoring the entire system dynamically, safely, and with minimal overhead.

External Host-Based Approaches. External host-based data collection is very popular in intrusion detection. This can be primarily attributed due to the advantages described by

⁴Although eBPF is available on all modern kernels, some of its features are specific to the very newest versions. In particular, recent verifier updates which allow for increased complexity have only been available since version 5.3. See Section 2.3 and Section 2.3.2 for more details.

Kerschbaum et al. [24], particularly with respect to ease of implementation and portability.

AAFID [43] uses a *combined* internal/external approach based on separate autonomous agents running continuously on multiple hosts. These agents make use of various data sources, such as external programs (i.e. `ps`, `netstat`, and `df`), file system state, and network interface packet capture (i.e. hooking into the host’s networking stack). Agents supplement collected data by analyzing audit logs generated by the system [24].

In 1999, Kuperman and Spafford [25] proposed the use of library interpolation for intrusion detection in dynamically linked binaries. Library interpolation is a method of interposing a custom library implementation between a dynamically linked executable and its shared objects. This effectively allows the generation of custom audit data on each library call that a process makes.

Internal and External Network-Based Approaches. Network-based approaches [44] to intrusion detection involve the inspection of network traffic en route to its destination. This typically comes in the form of inspecting packets headers, payloads, and frequency to establish patterns for analysis. Generally, network-based approaches have a choice between using either inline (internal) sensors, or passive (external) sensors for data collection [44]. An inline sensor either hooks into a network device, or is built into specialized hardware; traffic passes directly through the sensor and is analyzed directly.

In contrast, passive sensors create copies of network traffic for analysis. This approach is typically favored since it does not introduce delays into the traffic itself, instead sacrificing the ability to respond to threats before they reach their destination [44]. This result is consistent with Kerschbaum et al.’s observation that external sensor approaches tend to be favored over their internal counterparts [24].

2.6 Process Homeostasis

Anil Somayaji’s *Process Homeostasis* [41], styled as *pH*, forms the basis for ebpH’s core design; as such, it is worth exploring the original implementation, design choices, and rationale therein. Using the same IDS categories from the previous section, we can classify pH as a host-based anomaly detection system that responds by both issuing alerts *and* mounting countermeasures to reduce the impact of anomalies; in particular pH responds to anomalies by injecting delays into a process’ system calls proportionally to the number of recent anomalies that have been observed [41]. It is in this way that pH lives up to its name: these delays make process behavior *homeostatic*.

2.6.1 Anomaly Detection Through Lookahead Pairs

pH uses a technique known as *lookahead pairs* [41, 42] for detecting anomalies in system call data. This is in stark contrast to other anomaly detection systems at the time that primarily relied on *full sequence analysis*. Here we describe lookahead pairs, their use for anomaly detection, and offer a comparison with the more widely-known full sequence analysis.

In order to identify normal process behavior, profiles are constructed for each executable on the system. On calls to `execve`, pH associates the correct profile with a process and begins monitoring its system calls, modifying the lookahead pairs associated with the testing data of a profile. Once enough normal samples have been gathered and the profile has reached a specified maturity date, the process is then placed into training mode wherein sequences of system calls are compared with the existing lookahead pairs for a given profile.

Somayaji and Inoue [42] contrasted full sequence analysis with lookahead pairs and found that lookahead pairs produce fewer false positives than full sequences and maintain this property even with very long window lengths. This comes at the expense of potentially reduced sensitivity to some attacks as well as more vulnerable to mimicry attacks. However, as part of their work, Somayaji and Inoue showed that longer sequences can help mitigate these shortcomings in practice [42].

Both pH and ebpH use lookahead pair window sizes of 9, which has been shown to be effective at both minimizing false positive rates and mitigating potential mimicry attacks [41]. This window size also carries the advantage that lookahead pairs can be expressed with exactly 8 bits of information (one bit for every previous position $i \in \{1..9\}$).

2.6.2 Homeostasis Through System Call Delays

Perhaps the most unique aspect of pH’s approach is the means by which it achieves the eponymous concept of *process homeostasis*: system call delays. Inspired by the biological process of the same name, pH attempts to maintain homeostatic process behavior by injecting delays into system calls that are detected as being anomalous [41].

By scaling this response in proportion to the number of recent anomalies detected in a profile, pH is able to effectively mitigate attacks while minimizing the impact of occasional false positives. For example, a process that triggers several dozen anomalies will be slowed down to the point of effectively stopping, while a process that triggers only one or two might only be delayed by a few seconds. Admittedly, this relies upon the assumption of low burstiness for false positives. While this assumption generally holds, Somayaji acknowledges in his dissertation [41] that the possibility of attackers purposely provoking pH into causing

denial of service attacks is a potential problem. Additionally, users may become frustrated with pH's refusal to allow otherwise legitimate behavior simply due to the fact that it has not yet been observed.

In its current incarnation, ebpH does not yet delay system calls like its predecessor. The primary reason for this gap in functionality is that a solution still needs to be developed that works well with the eBPF paradigm; in particular, injecting delays via eBPF tracepoints or probes seems untenable due to the verifier's refusal to accommodate the code required for such an implementation. The addition of system call delays into ebpH is currently a topic for future work (see Section 7.2).

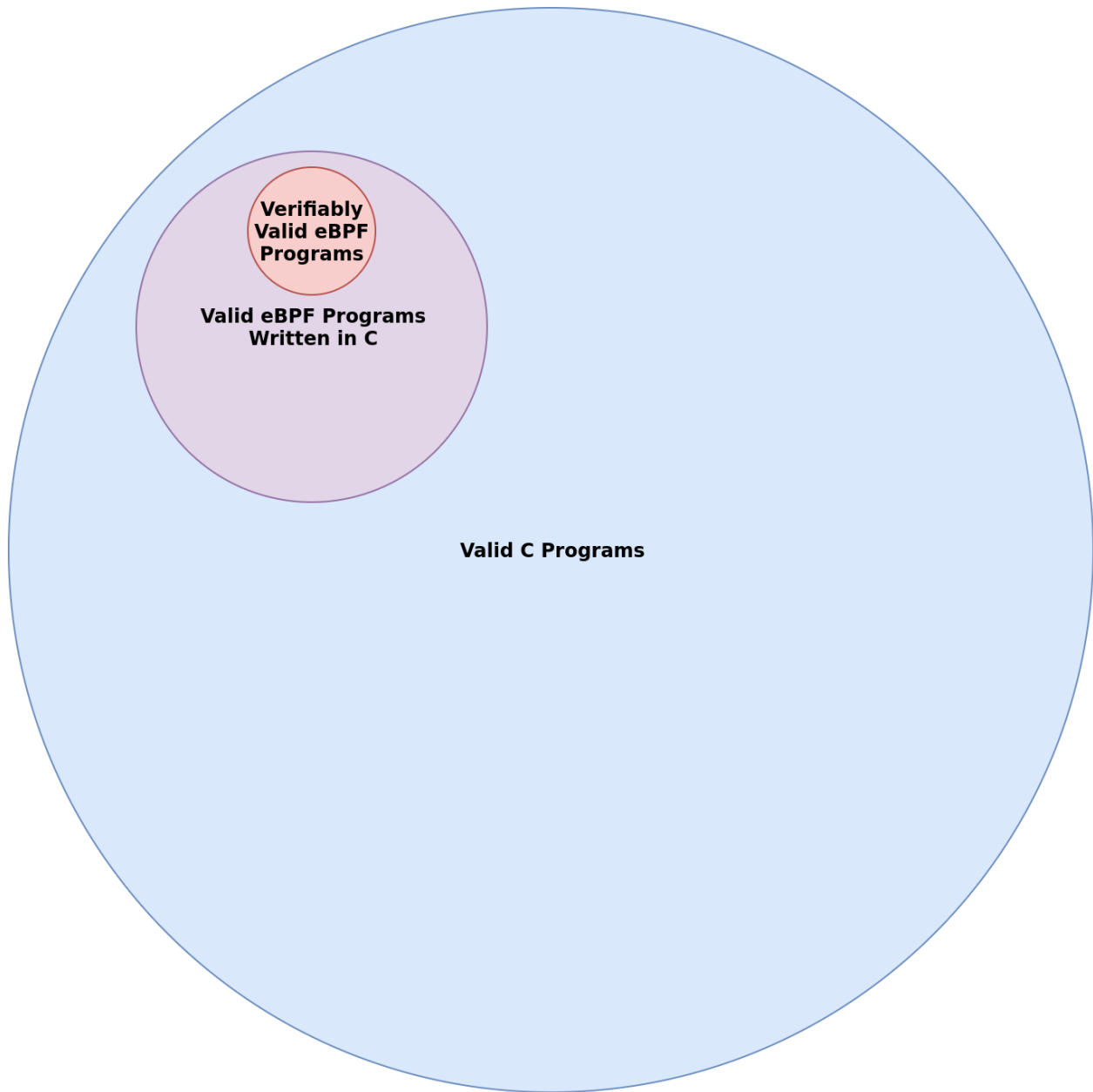


Figure 2.5: The set participation of valid C and eBPF programs. Valid eBPF programs written in C constitute a small subset of all valid C programs. Verifiably valid eBPF programs constitute an even smaller subset therein.

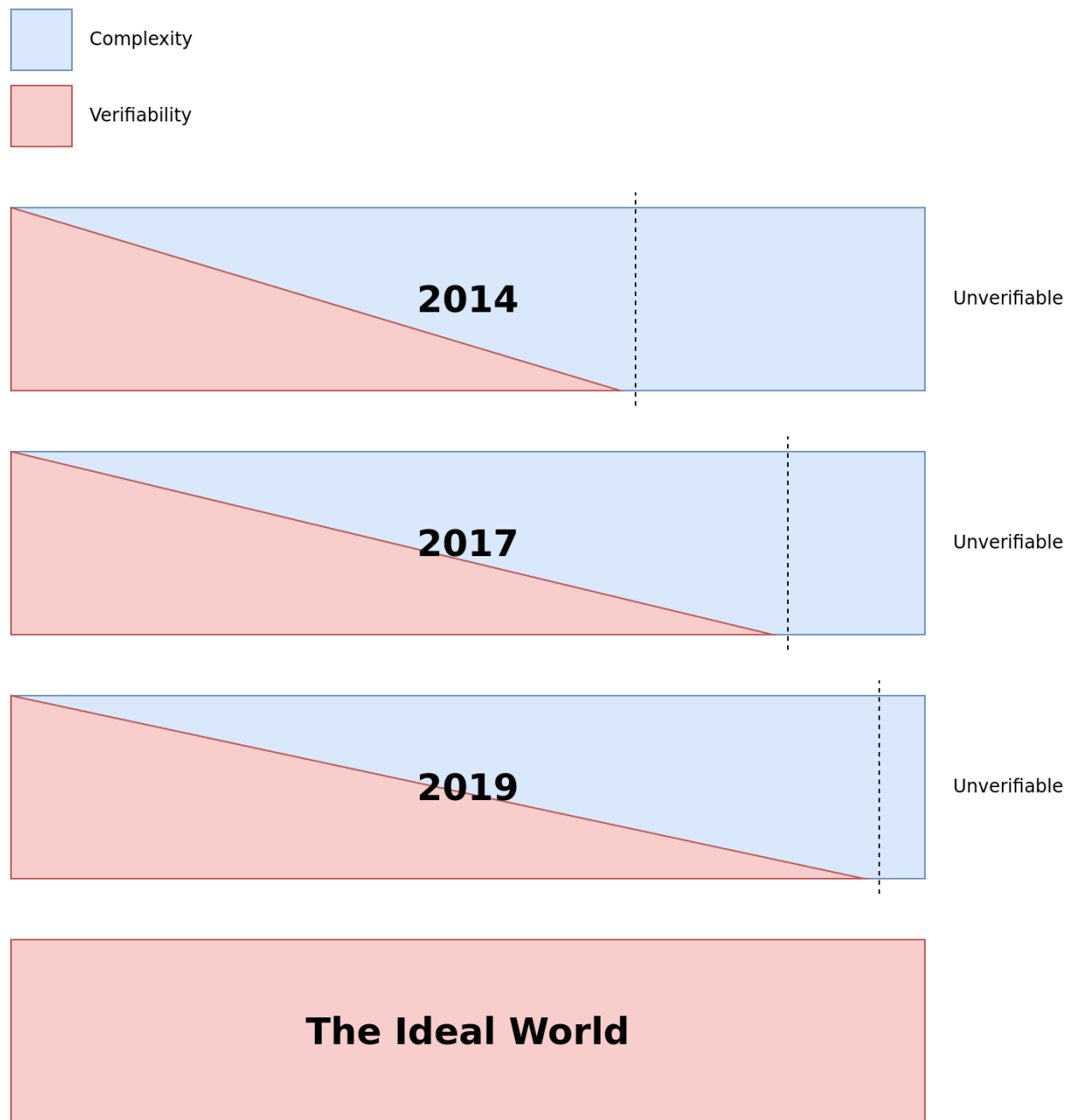


Figure 2.6: Complexity and verifiability of eBPF programs. Safety guarantees for eBPF programs rely on certain compromises. Ideally we would like to have a relationship as shown on the bottom; in practice, we have something that is getting closer over time, but is still far from the ideal.

Detection Technique

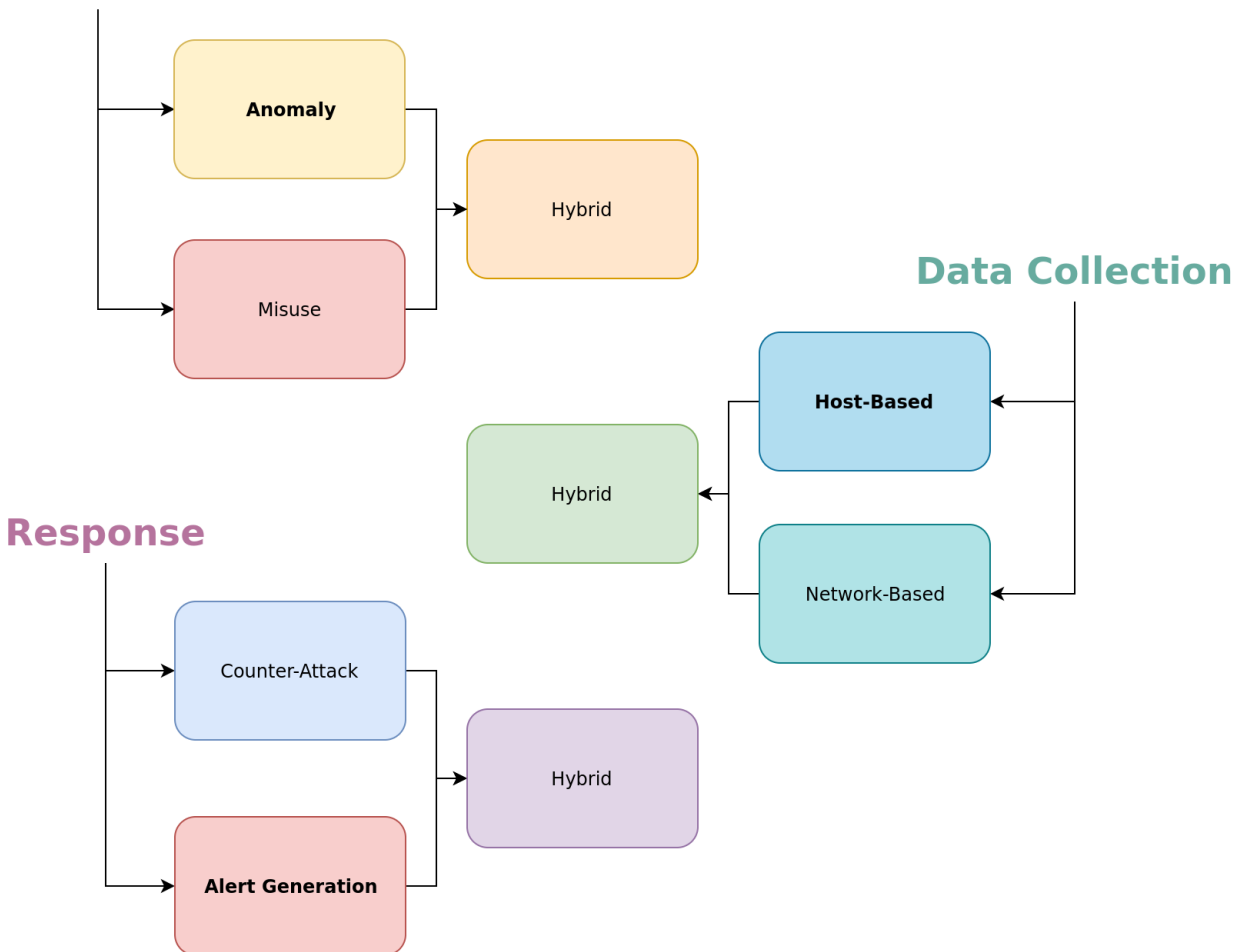


Figure 2.7: A broad overview of the basic categories of IDS. The current version of ebpH can be classified according to the categories labeled in **boldface**. Note that intrusion detection system classification can often be more nuanced than the basic overview presented here. However, this should present a good enough overview to understand IDSes in the context of ebpH.

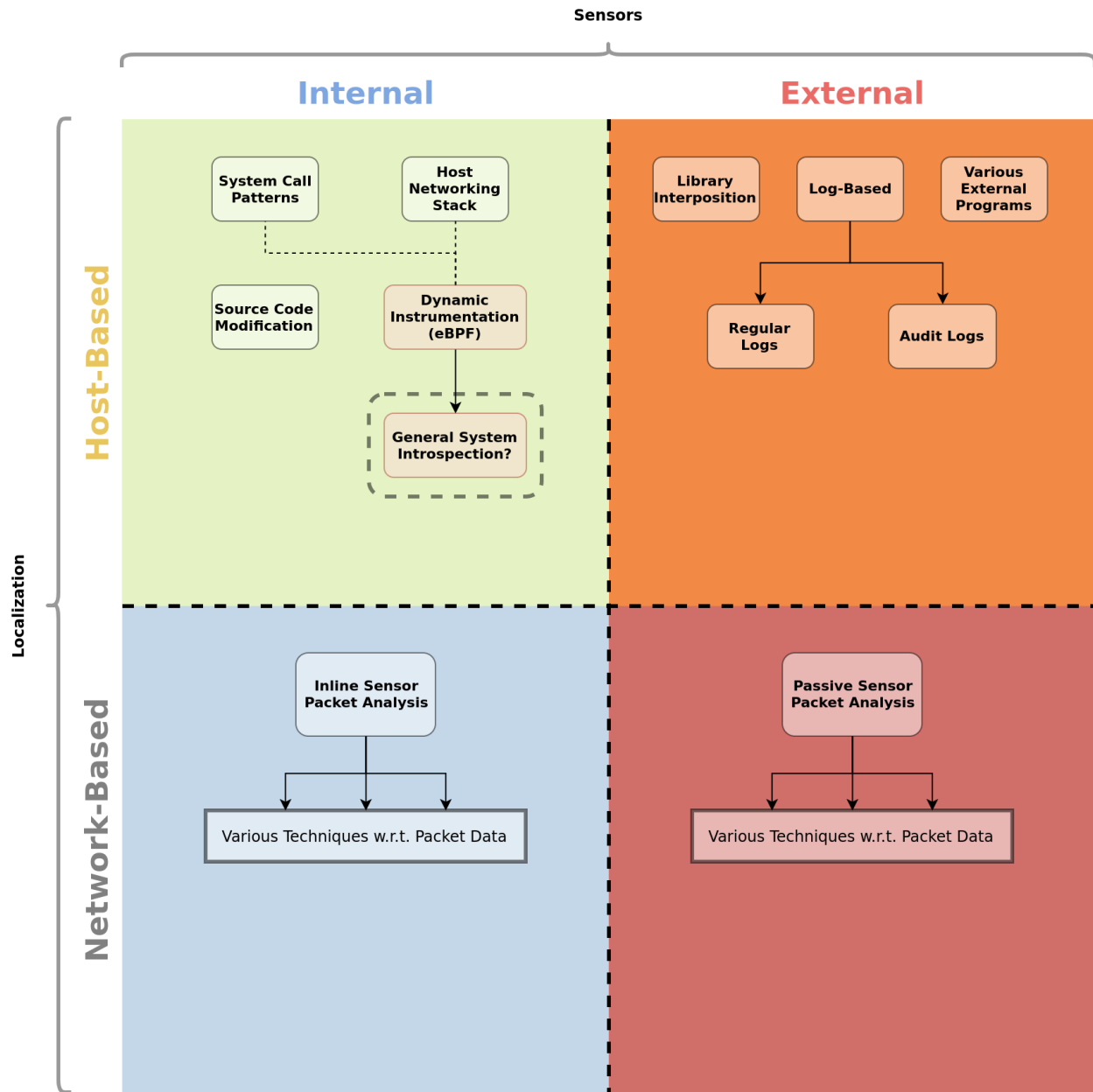


Figure 2.8: An overview of the most common data collection categories and subcategories in IDS, as well as a potentially new and promising category, *general system introspection*, thanks to eBPF. This figure primarily synthesizes the technologies presented in [24, 44].

3 Implementing ebpH

At a high level, ebpH is an intrusion detection system that profiles executable behavior by sequencing the system calls that processes make to the kernel; this essentially serves as an implementation of the original pH system described by Somayaji [41]. What makes ebpH unique is its use of an eBPF program for system call instrumentation and profiling (in contrast to the original pH which was implemented as a Linux 2.2 kernel patch).

ebpH can be thought of as a combination of several distinct components, functioning in two different parts of the system: userspace, and kernelspace (specifically within the eBPF virtual machine). In particular it includes a daemon, a CLI, and a GUI (described in Section 3.1) as well as an eBPF program (described in Section 3.2 and onwards). The dataflow between these components is depicted in Figure 3.1.

In order to implement the ebpH prototype described here, it was necessary to circumvent several challenges associated with the eBPF verifier and make several critical design choices with respect to dataflow between userspace and the eBPF virtual machine running in the kernel. This section attempts to explain these design choices, describe any specific challenges faced, and justify why eBPF was ultimately well-suited to an implementation of this nature.

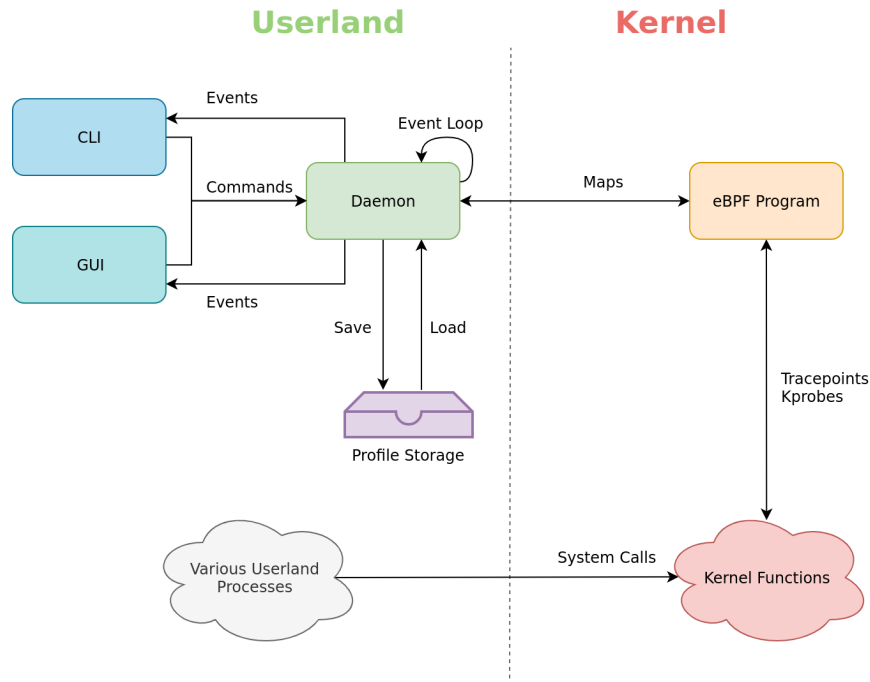


Figure 3.1: The dataflow between various components of ebpH.

3.1 Userspace Components

The userspace components of ebpb are comprised of several distinct and related programs. In particular, we can divide these programs into two sets: the *ebpb Daemon* (ebpb) and several CLI (command line interface) programs used to interact with it. The daemon is responsible for submitting BPF programs to the kernel, managing their state, and providing an API to other userspace programs. The CLI programs used to interact with the daemon include `ebpb-ps`, used to list actively traced processes, threads, and profiles, providing information about each, and `ebpb-admin`, used to issue commands to the daemon and to check the status of the BPF program. In order to issue more complex commands to the BPF program, `ebpb` leverages a userspace shared library, `libebpb.so` which provides functions that can be connected to arbitrary BPF programs via `uprobes`. Earlier versions of ebpb also included a GUI, however the GUI needs to be refactored in order to work with ebpb’s new architecture and this is currently a topic for future work.

3.1.1 The ebpb Daemon

The ebpb Daemon is implemented as a Python3 script that runs as a daemonized background process. When started, the daemon uses `bcc`’s Python front end [19] to generate the BPF bytecode responsible for tracing system calls, building profiles, and detecting anomalous behavior. It then submits this bytecode to the verifier and JIT compiler for insertion into the eBPF virtual machine.

Once the eBPF program is running in the kernel, the daemon continuously polls a set of specialized BPF maps called perf buffers which are updated on the occurrence of specific events. Table 3.1 presents an overview of the most important events we care about. As events are consumed, they are handled by the daemon and removed from the buffer to make room for new events. These buffers offer a lightweight and efficient method to transfer data from the eBPF program to userspace, particularly since buffering data in this way significantly reduces the number of required context switches between kernelspace and userspace.

In addition to perf buffers, the daemon is also able to communicate with the eBPF program through direct access to its maps. We use this direct access to issue commands to the eBPF program, check program state, and gather several statistics, such as profile count, anomaly count, and system call count. At the core of ebpb’s design philosophy is the combination of system visibility and security, and so providing as much information as possible about

⁵The majority of these values are subject to significant optimization in future iterations of ebpb. The 2^8 value is a sensible default chosen by `bcc`. In practice, many of these events are infrequent enough that smaller buffer sizes would be sufficient.

Table 3.1: Main event categories in ebp_H.

| Event | Description | Memory Overhead ⁵ |
|-------------------------------------|---|------------------------------|
| ebp _H _on_anomaly | Reports anomalies in specific processes and which profile they were associated with | 2 ⁸ pages |
| ebp _H _on_create_profile | Reports when new profiles are created | 2 ⁸ pages |
| ebp _H _on_pid_assoc | Reports new associations between PIDs and profiles | 2 ⁸ pages |
| ebp _H _error | A generic event for reporting errors to userspace | 2 ² pages |
| ebp _H _warning | A generic event for reporting warnings to userspace | 2 ² pages |
| ebp _H _debug | A generic event for reporting debug information to userspace | 2 ² pages |
| ebp _H _info | A generic event for reporting general information to userspace | 2 ² pages |

system state is of paramount importance.

The daemon also uses direct map access to save and load profiles to and from the disk. Profiles are saved automatically at regular intervals, configurable by the user, as well as any time ebp_H stops monitoring the system. These profiles are automatically loaded every time ebp_H starts.

3.1.2 ebph-ps

3.1.3 ebph-admin

3.2 ebp_H Profiles

In order to monitor process behavior, ebp_H keeps track of a unique profile (Listing 3.1) for each executable on the system. It does this by maintaining a hashmap of profiles, hashed by a unique per-executable ID; this ID is a 64-bit unsigned integer which is calculated as a unique combination of filesystem device number and inode number:

$$\text{key} = (\text{device number} \ll 32) + \text{inode number}$$

where \ll is the left bitshift operation. In other words, we take the filesystem's device ID in the upper 32 bits of our key, and the inode number in the lower 32 bits. This method provides a simple and efficient way to uniquely map keys to profiles.

The profile itself is a C data structure that keeps track of information about the executable,

Listing 3.1: A simplified definition of the ebpH profile struct.

```

1 struct ebpH_profile_data
2 {
3     u8 flags[SYSCALLS][SYSCALLS]; /* System call lookahead pairs */
4     u64 last_mod_count; /* Syscalls since profile was last modified */
5     u64 train_count;    /* Syscalls seen during training */
6 };
7
8 struct ebpH_profile
9 {
10     struct ebpH_profile_data train; /* Training data */
11     struct ebpH_profile_data test;  /* Testing data */
12     u8 frozen;                      /* Is the profile frozen? */
13     u8 normal;                      /* Is the profile normal? */
14     u64 normal_time;                /* Minimum system time required for normalcy */
15     u64 anomalies;                  /* Number of anomalies in the profile */
16     char comm[128];                 /* Name of the executable file */
17 };

```

as well as a sparse two-dimensional array of lookahead pairs [42] to keep track of system call patterns. Each entry in this array consists of an 8-bit integer, with the i^{th} bit corresponding to a previously observed distance i between the two calls. When we observe this distance, we set the corresponding bit to 1. Otherwise, it remains 0. Each profile maintains lookahead pairs for each possible pair of system calls. Figure 3.2 presents a sample (read, close) lookahead pair for the 1s binary.

Each process (Section 3.3) is associated with exactly one profile at a time. Profile association is updated whenever we observe a process making a call to `execve`. Whenever a process makes a system call, ebpH looks up its associated profile, and sets the appropriate lookahead pairs according to the process’ most recent system calls. This forms the crux of how ebpH is able to monitor process behavior.

Just like in the original pH [41], profile state is tracked using the `frozen` and `normal` fields. When a profile’s behavior has stabilized, it is marked frozen. If a profile has been frozen for one week (i.e. system time has reached `normal_time`), the profile is then marked normal. Profiles are unfrozen when new behavior is observed and anomalies are only flagged in normal profiles.

3.2.1 Writing Profiles to Disk and Reading Profiles from Disk

In order to allow profile data to persist across machine reboots, ebpH periodically writes profile data to disk, at an interval configurable the user, as well as when the BPF program

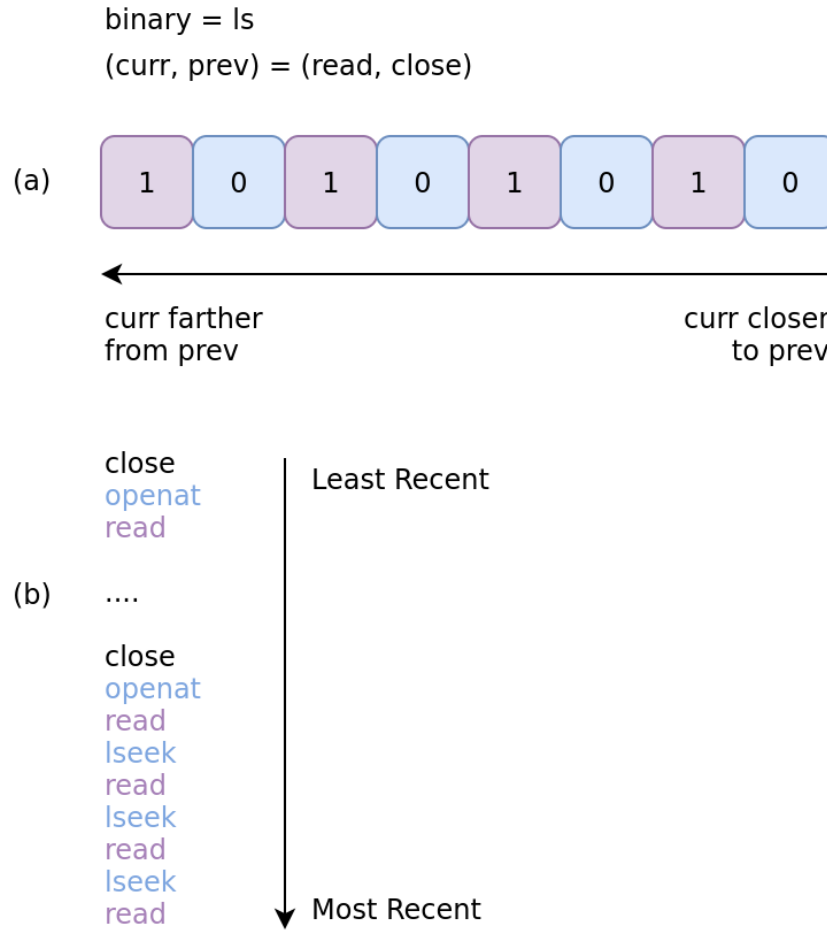


Figure 3.2: A sample (read, close) lookahead pair in the ebp_H profile for `ls`. (a) shows the lookahead pair and (b) shows two relevant system call sequences, separated by several omitted calls. Note that the first three system calls in both the first and second sequence are responsible for the two least significant bits of the lookahead pair.

is unloaded by the user. Profiles are read from disk when ebp_H first loads.

In the original pH, profile data was saved and loaded in kernelspace [41] which meant that it required kernelspace file I/O, which is often regarded as an unsafe practice. ebp_H solves this problem by moving all file I/O operations into userspace. This is made possible due to the bidirectional nature of dataflow with respect to eBPF maps. Specifically, when writing to disk, the daemon queries profile data from each entry in the profile map and writes that data to a file (`/var/lib/ebpH/<profile_key>`). When reading from disk, the daemon reads profile data from the appropriate files (`/var/lib/ebpH/<profile_key>`) and associates that data with keys in the newly created profile map.

3.3 Tracing Processes

Like profiles, process information is also tracked through a global hashmap of process structs. The process struct's primary purpose is to maintain the association between a process and its profile, maintain a sequence of system calls, and keep track of various metadata. See Listing 3.2 for a simplified definition of the ebp_H process struct.

Listing 3.2: A simplified definition of the ebp_H process struct.

```

1  struct ebpH_sequence
2  {
3      long seq[9];      /* Remember 9 most recent system calls in order */
4      u8 count;         /* How many system calls are in our sequence? */
5  };
6
7  struct ebpH_sequence_stack
8  {
9      ebpH_sequence[3]; /* Keep track of up to 3 sequences at a time */
10     int top;           /* Top of the sequence stack, values from 0-2 */
11     int should_pop;    /* Pop from the stack on next system call */
12 };
13
14 struct ebpH_process
15 {
16     struct ebpH_sequence_stack;
17     u32 pid;           /* Kernel tgid */
18     u32 tid;           /* Kernel pid */
19     u64 profile_key;   /* Associated profile key */
20     u8 in_execve;      /* Are we in the middle of an execve? */
21 };

```

ebp_H monitors process behavior by instrumenting tracepoints for both system call entry and return. The nine most recent system calls made by each process are stored in its respective process struct, and are used to set the correct lookahead pairs in the associated profile struct.

While we keep track of every system call made by a process, we pay special attention to a select few system calls which are directly related to profile creation, association, and disassociation. These system calls and their respective side effects are summarized in Table 3.2.

3.3.1 Profile Creation and Association

There are several important considerations here. First, we need a way to assign profiles to processes, which is done by instrumenting the `execve` system call using a tracepoint, as well as part of its underlying implementation via a kprobe. In particular, we hook the `do_open_execat`

Table 3.2: Important system calls in ebp_H.

| System Call | Description | ebp _H Side Effect |
|---------------------------|---|---|
| <code>execve</code> | Execute a program | (Re)associate a process with a profile, creating the profile if necessary; wipe the process' current sequence of system calls |
| <code>execveat</code> | Execute a program | |
| <code>fork</code> | Create a new process by duplicating calling process | Start tracing a process and associate with parent's profile; also copy the parent process' current sequence into the child |
| <code>vfork</code> | Create a child process and block parent | |
| <code>clone</code> | Create a new process or thread | |
| <code>rt_sigreturn</code> | Return from a signal handler | Pop the frame at the top of the process' sequence stack |

kernel function in order to access the file's inode and filesystem device number. This information is used to compute a key that uniquely maps to an individual executable on disk. Without this, we would be unable to differentiate between two paths that resolve to a binary with the same name, for example `/usr/bin/ls` and `./ls`.

The entry and exit points to the `execve` system call are used to differentiate a true `execve` call from the kernel routines responsible for loading shared libraries, which both invoke the aforementioned `do_open_execat` subroutine. When we first hit an `execve`, we set an indicator variable in the process struct to say that we are in the middle of an `execve`. Subsequent calls to `do_open_execat` are then ignored until we hit `execve`'s return tracepoint and unset the indicator variable.

In addition to associating a process with the correct profile, we also wipe the process' current sequence of system calls, to ensure that there is no carryover between two unrelated profiles when constructing their lookahead pairs.

3.3.2 Profile Association and Sequence Duplication

Another special consideration is with respect to `fork` and `clone` family system calls. A forked process should begin with the same state as its parent and should (at least initially) be associated with the same profile as its parent. A subsequent `execve` (i.e. the `fork-execve` pattern) would then overwrite this association. In order to accomplish this, we instrument tracepoints for the `fork`, `vfork`, and `clone` system calls, ensuring that we associate the child process with the parent's profile, if it exists. If ebp_H detects an `execve` as outlined above, it will simply overwrite the initial profile association provided by the `fork`. The parent's current system call sequence is also copied to the child to prevent forks from being used to

break sequences.

3.3.3 Dealing with Signal Handlers and Non-Determinism

As an anomaly-based intrusion detection system, it is critical that ebpH be able to establish normal profiles of program behavior in a timely manner. As presented in previous sections, establishing the normalcy of a profile requires that the it has been active for at least a week and that the ratio of total system calls seen during training to system calls the last time the profile was modified be sufficiently large. As a corollary, every time a process makes a system call that results in a previously unobserved sequence, this ratio becomes increasingly difficult to achieve. In practice, this means that it is much harder to normalize profiles that exhibit less deterministic behavior. As a practical example, consider the time required to stabilize a relatively simple binary, such as `ls` versus a complex web browser like `firefox`; not only does `firefox` make significantly more system calls during an average run, but it is also far more likely to produce a previously unseen sequence at any given time.

This problem of normalizing profiles is compounded by the non-deterministic behavior introduced by signals and signal handlers. This phenomenon was first noted by Amai et al. [1] in a 2005 technical paper on the original pH system. In particular, they noted that signal handlers were a significant source of non-deterministic behavior in processes that ultimately led to significantly longer wait times until profile normalcy. This effect is not difficult to see in practice, especially in the context of complex programs that run for extended periods of time, such as the above `firefox` example. Suppose that we have some sequence of system calls (A, B, C, D, E) and a signal handler that invokes system calls (F, G, H) ; depending on when this signal is caught during the initial sequence, the resulting sequence can vary significantly. For example, we might see (A, F, G, H, B, C, D, E) in one instance and (A, B, C, D, F, G, H, E) in another. This results in a significant deterioration in profile stability, and subsequently profile normalcy times.

ebpH deals with the problem of signal handlers in the same manner proposed by Amai et al. [1]. Specifically, we maintain a stack of system call sequences in each process struct; each time we receive a signal, we push a frame onto this stack, and each time we exit our signal handler, we pop the frame. This has the effect of temporarily wiping ebpH's memory of a process' current system call sequence whenever we enter a signal handler, allowing subsequent lookahead pairs to be unaffected by the execution context prior to the handler's invocation. In order to decide when to push, we instrument a new eBPF kprobe on the kernel's `do_signal` implementation; this allows us to detect when a process receives a signal that will be handled. Subsequently, we detect a return from a signal handler by checking for

the `re_sigreturn` system call; when `ebpH` detects such a return, it pops the top frame from the sequence stack.

3.3.4 Reaping Processes

`ebpH` reaps tasks from its process map whenever detects that they have exited. By reaping process structs from our map as we are finished with them we ensure that the map neither fills up, nor does it consume more memory than necessary. In order to detect when a task exits, we instrument the `sched_process_exit` tracepoint provided by the kernel's trace API. This tracepoint is triggered whenever the scheduler handles the termination of a task. We simply determine the task's PID and delete that key from our process map.

3.4 Training, Testing, and Anomaly Detection

`ebpH` profiles are tracked in two phases, *training mode* and *testing mode*. Profile data is considered training data until the profile becomes normal (as described in Section 3.2). Once a profile is in testing mode, the lookahead pairs generated by its associated processes are compared with existing data. When mismatches occur, they are flagged as anomalies which are reported to userspace via a perf event buffer. The detection of an anomaly also prompts `ebpH` to remove the profile's normal flag and return it to training mode.

3.4.1 A Simple Example of `ebpH` Anomaly Detection

As an example, consider the simple program shown in Listing 3.3. This program's normal behavior is to simply print a message to the terminal. However, when issued an extra argument (in practice, this could be a secret keyword for activating a backdoor), it prints one extra message. This will cause a noticeable change in the lookahead pairs associated with the program's profile, and this will be flagged by `ebpH` if the profile has been previously marked normal.

Listing 3.3: A simple program to demonstrate anomaly detection in `ebpH`.

```
1  /* anomaly.c */
2
3  #include <stdio.h>
4  #include <unistd.h>
5
6  int main(int argc, char **argv)
7  {
8      /* Execute this fake anomaly
9       * when we provide an argument */
10     if (argc > 1)
11         printf("Oops!\n");
```

```
12  /* Say hello */
13  printf("Hello world!\n");
14
15  return 0;
16 }
```

In order to test this, we artificially lower ebp_H's normal time to three seconds instead of one week. Then, we run our test program several times with no arguments to establish normal behavior. Once the profile has been marked as normal, we then run the same test program with an argument to produce the anomaly. ebp_H immediately detects the anomalous system calls and flags them. These anomalies are then reported to userspace via a perf buffer as shown in Listing 3.4.

Listing 3.4: The flagged anomaly in the `anomaly` binary as shown in the ebp_H logs. Note that ebp_H also logs the offending sequence, reordering it so that most recent system calls appear on the right.

```
1  WARNING: Anomalies in PID 11162 (anomaly 38803844):
2  MPROTECT, MPROTECT, MPROTECT, MUNMAP, FSTAT, BRK, BRK, WRITE, WRITE
```

From here, we can figure out exactly what went wrong by inspecting the system call sequences produced by the `anomaly` program, in both cases and comparing them with their respective lookahead pair patterns. Figure 3.3 provides an example of this comparison.

While this contrived example is useful for demonstrating ebp_H's anomaly detection, process behavior in practice is often more nuanced. ebp_H collects at least a week's worth of data about a process' system calls before marking it normal, which often corresponds with several branches of execution. In a real example, the multiple consecutive write calls might be a perfectly normal execution path for this process; by ensuring that we take our time before deciding whether a process' profile has reached acceptable maturity for testing, we decrease the probability of any false positives.

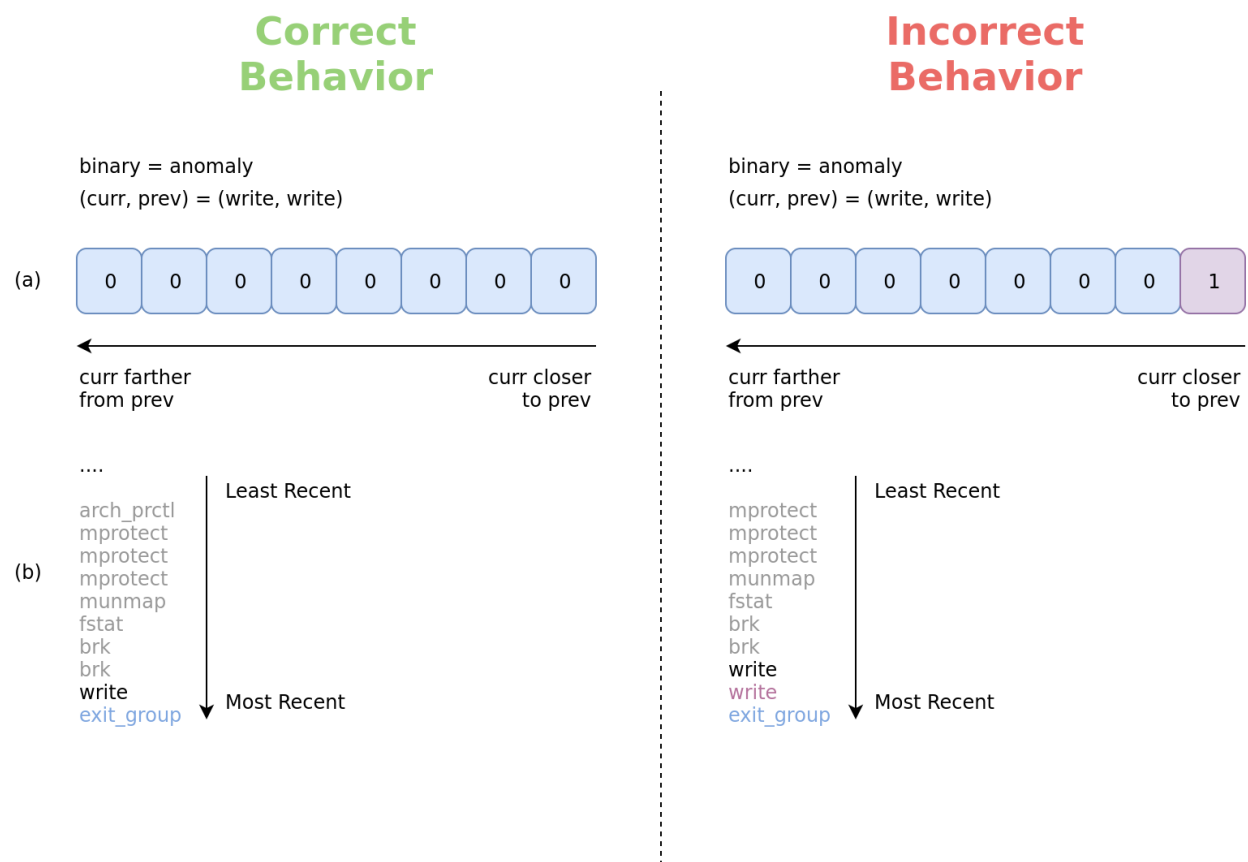


Figure 3.3: Two sample (write, write) lookahead pairs in the ebpH profile for `anomaly.c`. (a) shows the lookahead pair and (b) shows two relevant system call sequences. The left hand side depicts normal program behavior, while the right hand side depicts our artificially generated anomaly. There are several other anomalous lookahead pairs which result from this extra write call, but we focus on (write, write) for simplicity.

3.5 Issuing Commands to ebpfH

3.5.1 Setting Runtime Parameters

3.5.2 Examining Profiles and Processes

3.5.3 Issuing More Complex Commands

4 Technical Challenges of an eBPF Intrusion Detection System

4.1 Soothing the Verifier

The development of ebpfH elicited many challenges with respect to the eBPF verifier. As we have seen in Section 2.3.2, eBPF programs become more difficult to verify as they increase in complexity; as a corollary, when developing large and complex eBPF programs, a great deal of care and attention must be paid to ensure that the verifier will not reject our code.

The problem of dealing with the eBPF verifier can be expressed in the form of several subproblems as follows:

- (1) Many kernel functions and programming constructs are prohibited in eBPF;
- (2) eBPF programs come with a hard stack space limit of 512 bytes;
- (3) Dynamic memory allocation is prohibited and memory access requires strict safety checks;
- (4) Support for bounded loops is in its infancy and loops will not work without an easy proof that the induction variable will result in termination;
- (5) The verifier tends to err on the side of caution and will produce false positives with non-negligible frequency.

Subproblem (1) means that, at the moment, there is no simple means of injecting system call delay into system calls from within the eBPF program, an important part of the original pH's functionality [41]. Kernel scheduling and delay functions do not work in eBPF due to unsafe jump instructions, and so other means of delaying processes need to be explored. This is currently a topic for future work (see Section 7.2).

From subproblems (2) and (3), one immediate issue arises: with no means of explicit dynamic memory allocation and a stack space limit of 512 bytes, how do we instantiate the large structs described in previous sections? Both the `ebpfH_profile` and `ebpfH_process` structs are larger than would be allowed in the eBPF stack. Fortunately, we can creatively solve this

problem by using a `BPF_ARRAY` for initialization. Since a `BPF_ARRAY`’s entries are preinitialized with `0`, we can create an array of size 1 for each large datatype and copy its contents into the entries of a larger hashmap as needed. This technique constitutes the design pattern outlined in Listing A.1 of Appendix A.

On the topic of memory, another convenient feature of eBPF maps is the ability to flag them as being implicitly dynamically allocated. This means that the map will only use as much space as its current amount of entries requires. Memory management is handled automatically by the map. This combined with the aforementioned method of struct initialization gives us the means by which to safely and efficiently handle large data structures in eBPF programs.

From subproblem (4), we have the obvious issue that loops need to be “simple” enough for the eBPF verifier to reason about them. For example, loops that have entrypoints in the middle of iteration will potentially be flagged if the verifier is unable to correctly identify the loop structure [6]. Since the verifier relies on pattern matching in order to identify induction variables, LLVM optimizations to eBPF bytecode introduce an element of fragility to loop verification [6]. Bounded loops that perform memory access using the induction variable are also quite finicky at best; the verifier must be able to show that memory access is safe for each possible state of the loop. For these reasons, development of eBPF programs that require bounded loops is still far from perfect, but we now at least have the tools with which to implement complex eBPF programs like `ebpH`.

Subproblem (5) is perhaps the most difficult to reckon with, but is quite understandable when considering the gravity of the problem that the verifier is trying to solve. As we have already seen, guaranteeing the safety of arbitrary untrusted programs is a difficult problem and concessions need to be made in order for such guarantees to be tenable. False positives are unfortunately one of those concessions. When the verifier rejects code due to a false positive, there is simply no better solution than to try a different approach. Fortunately, well-constructed eBPF programs do not often suffer from false verifier positives, particularly as one learns the nuances of how the verifier works and how to coax it into accepting programs.

5 Measuring `ebpH`’s Overhead

One of the primary advantages of eBPF is its relatively low overhead [15, 45, 46] compared to many other system introspection solutions (c.f. Section 2.1 and Section 2.3). In order to justify this claim in the context of an eBPF intrusion detection system, we need to ascertain the overhead associated with running `ebpH` on a variety of systems under a variety

of workloads (artificial and otherwise). Here we describe the tests that were conducted in order to determine this overhead. Section 5.1 outlines the systems and tools used for testing and provides an overview of the collected datasets. The specifics of each benchmarking test along with the results are provided in Section 5.2.

5.1 Methodology

Since ebpH’s kernelspace functionality resides in system call hooks, we can get an idea of what overhead it imposes on the system by running macro-benchmarks on the time required to make system calls. Initially, I planned to use `syscount` [21] from `bcc-tools` for this purpose, however this tool currently has a race condition that may affect results due to its use of `BPF_HASH` rather than `BPF_PERCPU_ARRAY` for data storage (c.f. Table 2.2 on page 10). Instead, a custom benchmarking tool, `bpfbench`⁶, was written in eBPF for this purpose. Like `syscount`, `bpfbench` measures system call overhead by taking the difference of `ktime` (in nanoseconds) between system call entry and return; this difference along with the number of calls observed is stored in an eBPF map for later analysis. Unlike `syscount`, `bpfbench` stores this data in a per-cpu array, aggregating data at the end when necessary; this means that neither the system call count nor the system call overhead is subject to race conditions like its predecessor. See Appendix B for the BPF portion of `bpfbench`’s source code.

Macro-benchmarking data was collected on various systems, including a server used in production, a personal computer, and a CCSL (Carleton Computer Security Lab) workstation. Tests were run under a variety of workloads and benchmarking data was collected using `bpfbench`. For each dataset, the same test was conducted on the system twice: once with ebpH running, and once without. All ebpH data was collected while ebpH was monitoring the entire system (i.e. started immediately on boot via a `systemd` unit) and running with normal parameters and logging settings. Table 5.1 summarizes each of the systems used for the collection of benchmarking data and Table 5.2 provides a description of each dataset, including the system and the workload tested.

Table 5.1: Systems used for the collection of ebpH benchmarking data.

| System | Description |
|--------------------------|--|
| <code>arch</code> | Personal workstation |
| <code>bronte</code> | CCSL workstation |
| <code>homeostasis</code> | Mediawiki server for COMP3000 class wiki |

⁶Full source code available at <https://github.com/willfindlay/bpfbench>.

Table 5.2: ebpH macro-benchmarking datasets.

| Dataset | System | Workload | Description |
|------------------|-------------|------------|---|
| bronte-7day | bronte | Idle | bpfbench , 7 days with ebpH and 7 days without |
| homeostasis-3day | homeostasis | Production | bpfbench , 3 days with ebpH and 3 days without |
| arch-3day | arch | Normal use | bpfbench , 3 days with ebpH and 3 days without |

After benchmarking data was collected, overhead was calculated according to the following equation:

$$\text{Overhead}_{\text{syscall}} = \frac{T_{\text{ebpH}_{\text{syscall}}} - T_{\text{base}_{\text{syscall}}}}{T_{\text{base}_{\text{syscall}}}}$$

where,

$$T_{\text{syscall}} = \frac{\text{Total time}}{\text{Number of occurrences}}$$

as measured by **bpfbench**.

In addition to system call overhead, we are also interested in how ebpH affects overall system performance. In particular, ebpH should have negligible impact on normal system use. In order to ascertain this, micro-benchmarking data was collected for a variety of test cases. Table 5.3 provides a description of each micro-benchmark dataset, including the system and the workload tested. Additional details of each micro-benchmark test are provided in their respective results sections.

5.2 Results

In this section, we present the results of several macro and micro-benchmarks to ascertain the overhead caused by running ebpH.

Table 5.3: ebp_H micro-benchmarking datasets.

| Dataset | System | Workload | Description |
|----------------|--------|------------|--|
| bronte-process | bronte | Artificial | Three process creation microbenchmarks, 1,000,000 trials each, with and without ebp _H |
| bronte-x11perf | bronte | Artificial | x11perf [28, 33] full benchmarking suite, with and without ebp _H |
| bronte-s1bench | bronte | Artificial | s1bench [14] microbenchmark by Brendan Gregg, with and without ebp _H |

5.2.1 bronte Macro-Benchmark

A macro-benchmark of all system calls was run using bpfbench for a period of 14 days on bronte, a workstation in the CCSL lab at Carleton University. For the first seven days, the system had been running ebp_H since boot, and so all processes were being monitored. ebp_H was run with normal parameters, profile saving and loading enabled, and logging enabled. For the final seven days, the benchmark was taken without ebp_H. All benchmarks for this test were run under an idle workload. After concluding the benchmark tests, overhead was compared between ebp_H and non-ebp_H data in order to ascertain what effect ebp_H had on the system.

The tests were run as follows:

```

1 # enable ebphd systemd unit and reboot
2 sudo systemctl enable ebphd && reboot now
3 # test with ebpH for one week
4 sudo bpfbench -d 1w ebph-results.log
5 # stop ebphd
6 sudo systemctl stop ebphd
7 # test without ebpH for one week
8 sudo bpfbench -d 1w base-results.log

```

Table 5.4 presents the results of the benchmark. In order to eliminate outliers, we exclude results more than 3 standard deviations from the mean. Outliers may be caused by a number of factors, since we are measuring the entire system over an arbitrary period of time, although one common explanation would be blocking system calls taking an inordinate amount of time to finish. Indeed, of the six data points discarded that would have made it into the top 20, all six were system calls with some blocking component: futex(2), select(2), poll(2), epoll_pwait(2), pselect6(2), and ioctl(2). The first five all directly block until some event occurs, while many ioctl(2) implementations involve some degree of locking within the

Table 5.4: Top 20 most frequent system calls with base times of under $3\mu\text{s}$ from the bronte-7day dataset (after discarding outliers). Smaller overhead is better.

| System Call | $T_{\text{base}} (\mu\text{s})$ | $T_{\text{ebpH}} (\mu\text{s})$ | Diff. (μs) | % Overhead |
|------------------------|---------------------------------|---------------------------------|-------------------------|------------|
| ioperm | 1.423 | 1.433 | 0.010 | 0.70 |
| personality | 0.801 | 1.059 | 0.258 | 32.21 |
| sched_get_priority_max | 0.291 | 1.149 | 0.858 | 294.85 |
| sched_get_priority_min | 0.446 | 1.327 | 0.881 | 197.53 |
| clock_gettime | 1.316 | 1.896 | 0.580 | 44.07 |
| timer_delete | 1.477 | 2.213 | 0.736 | 49.83 |
| timer_settime | 2.706 | 3.527 | 0.821 | 30.34 |
| fchownat | 2.905 | 4.138 | 1.233 | 42.44 |
| ioprio_set | 2.189 | 3.148 | 0.959 | 43.81 |
| listen | 1.082 | 1.821 | 0.739 | 68.30 |
| shmctl | 1.148 | 2.054 | 0.906 | 78.92 |
| clock_getres | 1.343 | 2.254 | 0.911 | 67.83 |
| setregid | 1.620 | 2.728 | 1.108 | 68.40 |
| getpgrp | 0.622 | 1.595 | 0.973 | 156.43 |
| setreuid | 1.807 | 2.787 | 0.980 | 54.23 |
| llistxattr | 2.170 | 3.134 | 0.964 | 44.42 |
| fgetxattr | 2.619 | 3.513 | 0.894 | 34.14 |
| linkat | 2.334 | 2.522 | 0.188 | 8.05 |
| sched_getparam | 0.602 | 1.333 | 0.731 | 121.43 |
| flistxattr | 0.794 | 1.721 | 0.927 | 116.75 |

device driver.

Of the data presented in Table 5.4, the `bpf(2)` system call represents the most extreme overhead, with an overhead of 691% and a standard deviation of 2.48. This high deviation is explicable by the fact that many of the `bpf(2)` system calls made directly by `ebpH` involve direct map reads and writes of several megabytes at a time when saving and loading profiles to and from disk. When `ebpH` is not running, the `bpf(2)` system calls made by other BPF-related programs on the system generally do not involve such intensive operations.

Other points in the dataset shown in Table 5.4 generally do not present with overhead higher than about 130%. While 130% may seem high, it is important to remember that these overheads are on system calls that have an execution time of a few microseconds. In practice, `ebpH` is only adding a few microseconds of additional execution to these system calls, a slowdown that will likely go unnoticed by users. If we focus on system calls with a very small execution time, in the order of 7 microseconds or less base time, we see that this pattern holds, with the exception of the `bpf(2)` call described above. ?? shows the 20 most frequent system calls with base times of less than 7 microseconds. ?? and ?? show

barplots of times and overheads respectively for system calls with base times of less than 7 microseconds.

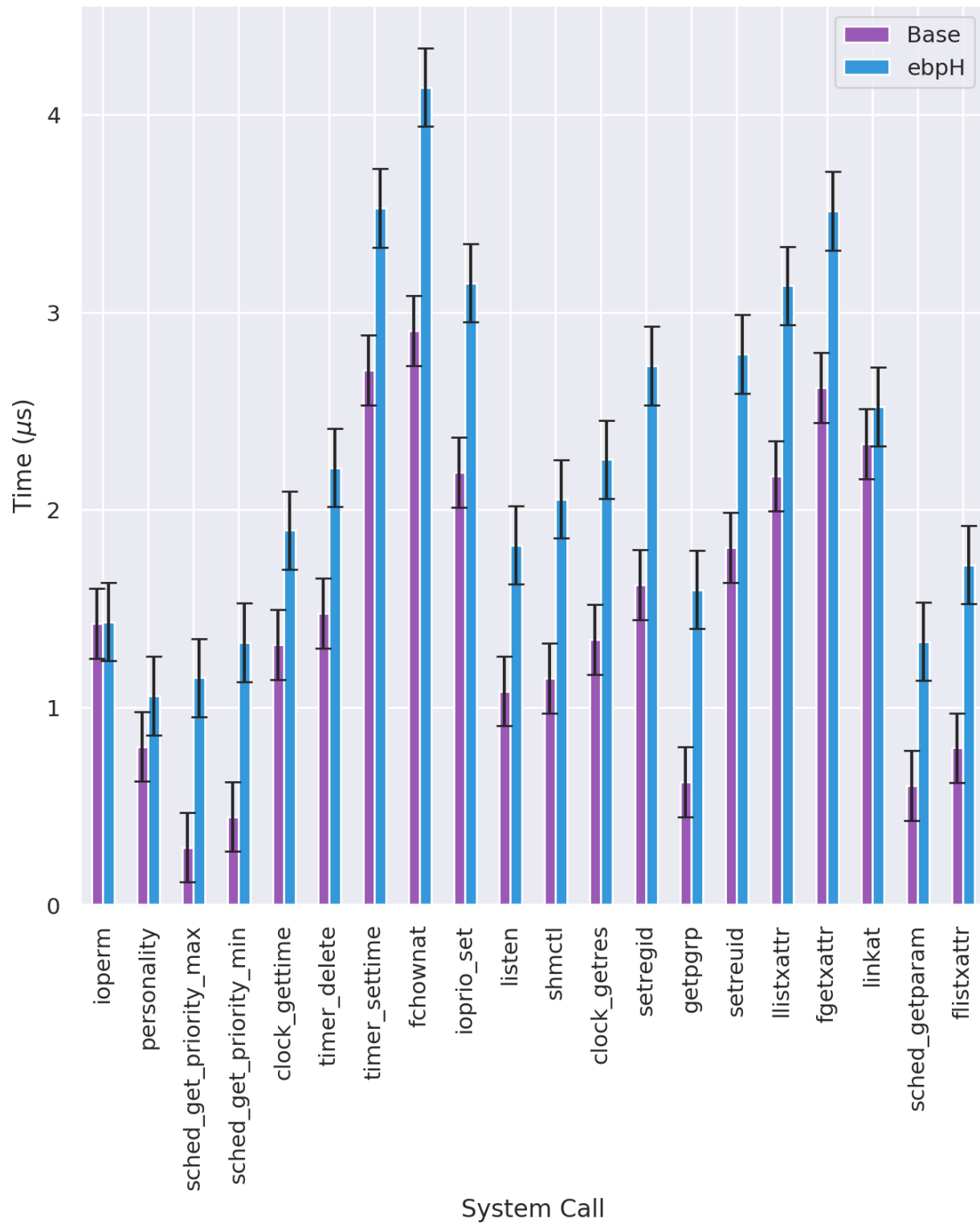


Figure 5.1: Top 20 most frequent system calls with base times of under $3\mu\text{s}$ from the bronte-7day dataset (after discarding outliers). Smaller difference in times is better. Standard error is given as error bars.

5.2.2 homeostasis Macro-Benchmark

5.2.3 arch-3day Macro-Benchmark: Using ebpH on a Personal Computer

5.2.4 bronte-lmbench: Measuring System Latency with lmbench Micro-Benchmark

McVoy’s lmbench [29, 30] is a Linux micro-benchmarking suite that has seen prominent use in academia [2, 7, 40] for establishing various performance metrics of UNIX-like systems. In particular, we are interested in the *OS*-category benchmarks provided by lmbench, especially *system call latency* and *dynamic process creation latency*. The former will give us a better idea of how much overhead ebpH imposes on simple (short execution time, and non-blocking) system calls while the latter will allow us to conjecture about ebpH’s overhead in practice.

1000 ebpH and 1000 non-ebpH trials were run on bronte and results were then averaged and compared to determine overhead. Table 5.5 and Figure 5.2 show the system call latency results and Table 5.6 and Figure 5.3 show the dynamic process creation latency results.

Table 5.5: Results of the system call benchmarks from the bronte-lmbench dataset. Standard deviations are given in parentheses, smaller overhead is better.

| System Call | $T_{\text{base}} (\mu\text{s})$ | $T_{\text{ebpH}} (\mu\text{s})$ | Diff. (μs) | % Overhead |
|-------------|---------------------------------|---------------------------------|-------------------------|------------|
| getppid | 0.055 (0.0010) | 0.390 (0.0081) | 0.335 | 611.94 |
| write | 0.103 (0.0016) | 0.436 (0.0083) | 0.333 | 321.83 |
| read | 0.174 (0.0032) | 0.510 (0.0122) | 0.336 | 193.16 |
| fstat | 0.180 (0.0018) | 0.513 (0.0081) | 0.334 | 185.32 |
| stat | 0.535 (0.0066) | 0.889 (0.0138) | 0.354 | 66.29 |
| open/close | 0.994 (0.0284) | 1.743 (0.0354) | 0.749 | 75.38 |

Table 5.6: Results of the process creation benchmarks from the bronte-lmbench dataset. Standard deviations are given in parentheses, smaller overhead is better.

| Program | $T_{\text{base}} (\mu\text{s})$ | $T_{\text{ebpH}} (\mu\text{s})$ | Diff. (μs) | % Overhead |
|-----------------|---------------------------------|---------------------------------|-------------------------|------------|
| fork+exit | 196.795 (6.6845) | 195.656 (6.1970) | -1.140 | -0.58 |
| fork+execve | 542.112 (10.8958) | 596.077 (12.8336) | 53.965 | 9.95 |
| fork+/bin/sh -c | 1488.232 (14.7249) | 1643.726 (18.1702) | 155.494 | 10.45 |

According to the results of the micro-benchmark, ebpH adds negligible overhead to the execution time of the fork+exit program. The specific overhead was measured to be -0.58% , but this perceived “performance improvement” is well within the margin of error for measurements. Actual overhead is like to be on the order of about 0 to 3%. This result is consistent with macro-benchmarking data from the bronte-7day dataset where the base and

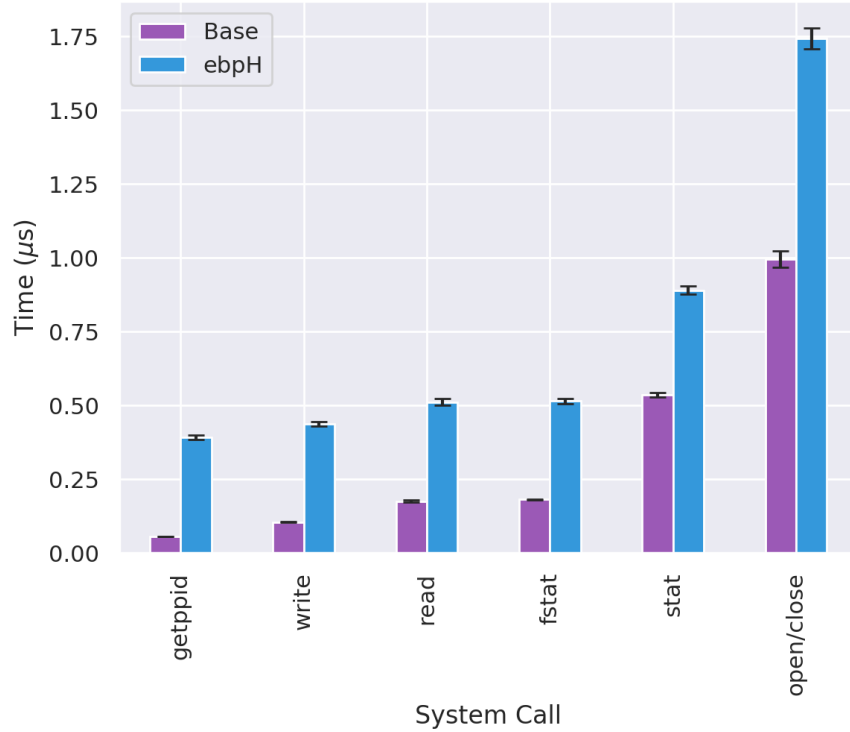


Figure 5.2: Mean system call times from the `bronte-lmbench` dataset. Smaller difference in times is better. Standard deviation is given as error bars.

ebpH times for the `clone(2)` system call were 105.553 and 106.649 respectively, a difference of just over 1 microsecond.

For the `fork+execve` program, ebpH added only about 54 microseconds of overhead, or a difference of about 9.95%. This result is *significantly* better than the results for the original pH system [41], but the difference can be explained by the fact that ebpH does not need to load profiles when programs first execute – the profiles are already resident in the profile map; the only exception to this rule is that space for additional profiles needs to be allocated as the map fills (c.f. Section 3.2) or when ebpH first loads (c.f. Section 3.2.1).

Similarly, ebpH outperforms the original pH system [41] on the `fork+/bin/sh` test. For this test, ebpH added just over 155 microseconds of overhead, a difference of about 10.45%. The overhead increase can be attributed to the fact that invoking a shell to execute a binary involves many more system calls⁷ than the traditional `fork` and `exec` pattern. More system calls implies more opportunities for ebpH’s per-system call overhead to impact performance. Even so, this 10.45% performance loss is essentially imperceptible on individual runs and therefore would not impact system behavior from the user’s perspective under normal work-

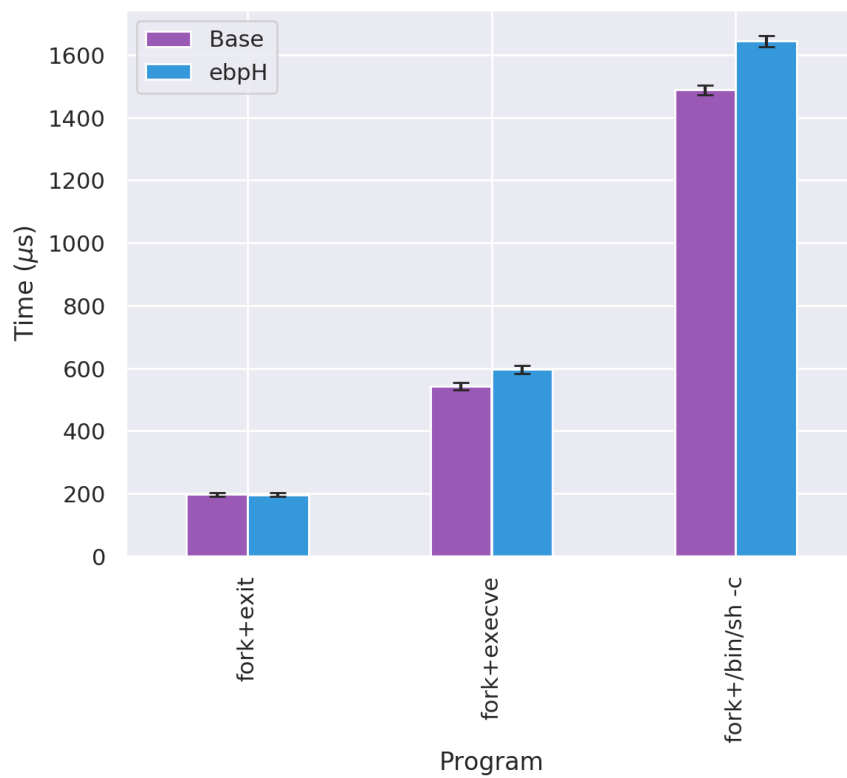


Figure 5.3: Mean process creation times from the `bronte-lmbench` dataset. Smaller difference in times is better. Standard deviation is given as error bars.

loads.

⁷111 or so additional system calls per execution (and 11 executions per run), to be more precise.

5.3 Comparing Results with the Original pH

6 Future Work

6.1 Security Analysis

6.2 Controlling for Further Sources of Non-Deterministic Behavior

6.3 Automating ebpH Response

6.4 General System Introspection and the Future of ebpH

7 Methodology and Future Work

While the ebpH prototype is certainly capable of monitoring a system for anomalies, much testing and work remains to be done in order to completely reimplement the original pH and ascertain whether eBPF is truly the best choice for implementing such an IDS. Here, we discuss the planned strategies for testing ebpH, as well as plans for iteration on the initial prototype, and future work.

7.1 Planned Testing Strategy

The ebpH prototype as well as its future iterations will be heavily tested on several machines under a variety of workloads. Table 7.1 summarizes the currently planned systems and any relevant details therein. Additional testing will be done on virtual machines to simulate multiple systems under artificially constructed workloads and attacks.

Table 7.1: A summary of the various systems that will be used to test ebpH.

| System | Description |
|---------------------------|---|
| arch | My personal desktop computer, which has been running the current ebpH prototype for one month |
| archlaptop | My personal laptop computer, which has been running the current ebpH prototype for one month |
| homeostasis | Dr. Somayaji's computer at Carleton University, which also serves the Wiki for his courses |
| CCSL Servers | The servers at the Carleton Computer Security Lab |
| Assorted Virtual Machines | Several virtual machines running a variety of test workloads |

7.1.1 Gathering and Analyzing Profile Data

ebpH will be retrofitted with the ability to generate CSV datasets and plots that will summarize profile data. Through the examination of profile data (specifically lookahead pair patterns), we can get a clear picture of ebpH’s understanding of system behavior. Results will be gathered for a wide variety of systems as depicted above in Table 7.1; this will yield the opportunity to test ebpH on production systems of various scale (i.e. *homeostasis* and the *CCSL Servers*) as well as personal computers for everyday use (i.e. *arch* and *archlaptop*).

Furthermore, the addition of several virtual machines for testing will provide the means to conduct reproducible experiments across various conditions, including measuring ebpH’s response to a variety of simulated attacks. Snapshots will ensure controlled and consistent system state between runs, and will be particularly useful in controlling for initial ebpH profile state during each round of testing.

During normal testing, we are particularly interested in the rate of false positives and false negatives observed by ebpH. A false positive will be defined as any anomaly detected by ebpH that corresponds to ordinary system behavior, while a false negative will be defined as a failure to detect the presence of an attack. In order to test for both false positive and false negative rates, we will observe ebpH on ordinary systems as well as systems under attack and compare results.

As a general-purpose anomaly-based intrusion detection system, it is important to show that ebpH is capable of detecting a wide variety of attacks. The mimicry attacks described in Wagner and Soto’s paper [55] are particularly interesting, as they were directly designed to defeat the original pH system (albeit an earlier version with much shorter lookahead pair window length) by constructing attack patterns that generate false negatives.

7.1.2 Gathering and Analyzing Performance Data

One of the primary advantages cited for using eBPF to build intrusion detection systems is lower overhead. In order to test the validity of this claim, we need reliable metrics to measure ebpH’s memory and CPU overhead under a variety of workloads and systems. A recent patch to the Linux Kernel has added the ability to measure individual eBPF program performance [47]. Additionally, we can combine this approach through hardware performance measurement with eBPF perf events. This approach should provide the combined advantage of measuring the specific overhead associated with ebpH along with its impact on the overall memory and CPU usage of its environment.

Another important consideration with respect to overhead is ebpH’s direct impact on the

user; in particular, we want to avoid annoying the user by introducing noticeable delays into their workflow. Therefore, in addition to rigorous quantitative testing, ebpH’s overhead will also be qualitatively tested for noticeable impact on system performance during everyday use.

7.2 Potential Improvements to ebpH

The system described in this proposal is a prototype, designed to implement the basic functionality of the pH intrusion detection system in eBPF, in order to ascertain whether such an implementation would be viable. While I believe I have achieved that goal, there is still plenty of room for future work on ebpH. Topics for future work include adding a mechanism for delaying system calls, using ebpH to increase overall system visibility, and the potential introduction of alternative behavioral metrics to provide a more comprehensive picture of system state and make better predictions about its validity. Furthermore, I plan to make extensive improvements to the ebpH GUI to complement the aforementioned augmentations to ebpH daemon functionality.

7.2.1 Delaying System Calls

The most obvious improvement to ebpH is the introduction of system call delays in a future iteration. This feature comprises a large part of the original pH’s response strategy and would be a vital part of a full reimplementaion. As previously discussed, this is not necessarily an easy thing to accomplish due to the eBPF verifier’s restrictions on program safety. From the perspective of what eBPF is trying to accomplish, this makes sense. In Somayaji’s dissertation [41], he discussed the potential for pH itself to cause denial of service on a system, due to intentional provocation from an attacker or simply an edge case in program behavior. eBPF is designed with safety in mind; allowing eBPF programs to cause denial of service in this way would be the antithesis of what eBPF is trying to accomplish. Therefore, another solution is needed.

A kernel-based implementation for process delays would certainly work, but would be far from ideal – this sacrifices a lot of the advantages that come with an eBPF implementation of pH in the first place, namely easy portability between Linux systems and guaranteed safety. Such an implementation would either be in the form of a kernel patch or a loadable kernel module; both of these solutions suffer from safety issues as we have discussed at length in Section 2.1. Additionally, a kernel patch in particular limits the portability of ebpH, which currently runs on any vanilla Linux Kernel above version 5.3.

We can also consider the possibility of busy waiting within the eBPF tracepoints themselves, although this also carries a few obvious drawbacks. Firstly, busy waiting means that the process that is being slowed down will continue to occupy CPU time instead of yielding it to another process by ceding time to the scheduler. Another obvious drawback is with respect to the verifier itself; verifier support for the bounded loops required for busy waiting is conditional on several factors. This may result in the rejection of busy waiting due to perceived safety violations.

A third possibility is issuing delays from userspace via the SIGSTOP and SIGCONT signals; the daemon would simply issue these signals to offending processes and space them proportionally with respect to recent anomalies produced. While this solution *could* work, it suffers from a few obvious flaws. Firstly, there is no guarantee that we can prevent another signal from waking up the process early; in fact, an attacker with the ability to send arbitrary signals has already completely circumvented this type of response. Additionally, there is no guarantee that a process will receive this signal in time to stop the offending system call(s). By the time the process receives the signal, it may already be too late to stop the attack.

None of these solutions seem ideal; it is likely that a presently unknown fourth alternative will present the best approach. Perhaps there may be an entirely eBPF-based solution on the horizon pending updates to the verifier – time will tell. For now, it may be worthwhile exploring what options are available at the present to determine if any are suitable for use in practice.

7.2.2 Measuring Other System Behavior

In his dissertation [41], Somayaji discusses the potential of having multiple homeostatic systems at work on a given machine. This approach would more closely resemble the concept of homeostasis we know in biology, wherein multiple subsystems work together to add stability to overall system state. pH was a great starting point, and pending the introduction of system call delays as discussed in the above section, ebpH will follow in its footsteps in that regard. However, much of the true power of eBPF comes from its ability to monitor *so much* system state at once; there’s no reason ebpH has to stop at system call tracing.

By using eBPF to monitor multiple facets of system state, we can get a clearer picture of normal process behavior, which could in turn yield more accurate anomaly detection results. Perhaps these metrics could include memory allocations, number of incoming network connections, socket I/O, file I/O, CPU time per process, or any number of such metrics. eBPF can measure all of that and more; and it can do so reliably, efficiently, and with guaranteed

safety.

7.2.3 Overall System Visibility

As an intrusion detection system, ebpH's role is well-defined: monitor the system, detect misbehaving processes, and report them to the user. However, there is one glaring problem with this approach, particularly as we venture into the territory of automated responses via system call delays: users do not necessarily *want* a system that chooses not to perform a requested action; they also do not necessarily *want* a system that harasses them with warnings about program behavior that they either don't care about or don't necessarily understand.

One potential solution to this problem is providing other benefits to the user through ebpH *in addition to* intrusion detection and response. For example, future versions of ebpH could include a performance analysis component, a debugger component, or any number of other metrics for increased system visibility. After all, one of the primary use cases for system introspection is precisely that: allowing a user to observe their system. By adding this extra functionality, we can provide complimentary benefits to the user that may incentivize them to run ebpH in the first place.

It should also not be overlooked that, in many cases, increased system state visibility can provide implicit security benefits to the experienced user. For example, an experienced system administrator could use a future version of ebpH to find vulnerabilities in their system before an attack even occurs.

References

- [1] W. A. Amai, E. A. Walther, and A. B. Somayaji, “An Immunological Basis for High-Reliability Systems Control,” Sandia National Laboratories, Tech. Rep., Mar. 2005.
- [2] P. Barham, B. Dragovic, K. Fraser, *et al.*, “Xen and the art of virtualization,” *SIGOPS Oper. Syst. Rev.*, vol. 37, no. 5, pp. 164–177, Oct. 2003, ISSN: 0163-5980. DOI: [10.1145/1165389.945462](https://doi-org.proxy.library.carleton.ca/10.1145/1165389.945462). [Online]. Available: <https://doi-org.proxy.library.carleton.ca/10.1145/1165389.945462>.
- [3] *bpf(2) Linux Programmer’s Manual*, Linux, Aug. 2019.
- [4] B. M. Cantrill, M. W. Shapiro, and A. H. Leventhal, “Dynamic Instrumentation of Production Systems,” in *Proceedings of the Annual Conference on USENIX Annual Technical Conference*, ser. ATEC ’04, Boston, MA: USENIX Association, 2004, pp. 2–2. [Online]. Available: https://www.usenix.org/legacy/publications/library/proceedings/usenix04/tech/general/full_papers/cantrill/cantrill.pdf.
- [5] J. Cespedes and P. Machata, *Ltrace(1) linux user’s manual*, Ltrace project, Jan. 2013.
- [6] J. Corbet, *Bounded loops in BPF programs*, Dec. 2018. [Online]. Available: <https://lwn.net/Articles/773605/>.
- [7] M. Degermark, A. Brodnik, S. Carlsson, and S. Pink, “Small forwarding tables for fast routing lookups,” *SIGCOMM Comput. Commun. Rev.*, vol. 27, no. 4, pp. 3–14, Oct. 1997, ISSN: 0146-4833. DOI: [10.1145/263109.263133](https://doi-org.proxy.library.carleton.ca/10.1145/263109.263133). [Online]. Available: <https://doi-org.proxy.library.carleton.ca/10.1145/263109.263133>.
- [8] M. Fleming, *A thorough introduction to eBPF*, Dec. 2017. [Online]. Available: <https://lwn.net/Articles/740157/>.
- [9] S. Forrest, S. A. Hofmeyr, A. Somayaji, and T. A. Longstaff, “A sense of self for unix processes,” in *Proceedings 1996 IEEE Symposium on Security and Privacy*, May 1996, pp. 120–128. DOI: [10.1109/SECPRI.1996.502675](https://doi-org.proxy.library.carleton.ca/10.1109/SECPRI.1996.502675).
- [10] P. D. Fox, *Dtrace4linux/linux*, Sep. 2019. [Online]. Available: <https://github.com/dtrace4linux/linux>.
- [11] S. Goldstein, “The Next Linux Superpower: eBPF Primer,” USENIX SRECon16 Europe, Jul. 2016. [Online]. Available: <https://www.usenix.org/conference/srecon16europe/program/presentation/goldshtein-ebpf-primer>.
- [12] B. Gregg, *Linux BPF Superpowers*, Mar. 2016. [Online]. Available: <http://www.brendangregg.com/blog/2016-03-05/linux-bpf-superpowers.html>.

- [13] B. Gregg, *bpfftrace (DTrace 2.0) for Linux 2018*, Oct. 2018. [Online]. Available: <http://www.brendangregg.com/blog/2018-10-08/dtrace-for-linux-2018.html>.
- [14] B. Gregg, *s1bench*, Feb. 2018. [Online]. Available: <https://github.com/brendangregg/Misc/blob/master/s1bench/s1bench.c>.
- [15] B. Gregg, *BPF Performance Tools*. Addison-Wesley Professional, 2019, ISBN: 0-13-655482-2.
- [16] B. Gregg, J. Mauro, and B. M. Cantrill, *DTrace: dynamic tracing in Oracle Solaris, Mac OS X and FreeBSD*. Prentice Hall, 2014.
- [17] T. Høiland-Jørgensen, J. D. Brouer, D. Borkmann, *et al.*, “The eXpress Data Path: Fast Programmable Packet Processing in the Operating System Kernel,” in *Proceedings of the 14th International Conference on Emerging Networking Experiments and Technologies*, ser. CoNEXT ’18, Heraklion, Greece: ACM, 2018, pp. 54–66, ISBN: 978-1-4503-6080-7. DOI: [10.1145/3281411.3281443](https://doi.org/10.1145/3281411.3281443). [Online]. Available: <http://doi.acm.org/10.1145/3281411.3281443>.
- [18] A. Hussain, J. Heidemann, J. Heidemann, and C. Papadopoulos, “A Framework for Classifying Denial of Service Attacks,” in *Proceedings of the 2003 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, ser. SIGCOMM ’03, Karlsruhe, Germany: ACM, 2003, pp. 99–110, ISBN: 1-58113-735-4. DOI: [10.1145/863955.863968](https://doi.org/10.1145/863955.863968). [Online]. Available: <http://doi.acm.org.proxy.library.carleton.ca/10.1145/863955.863968>.
- [19] IOVisor, *Iovisor/bcc*, Oct. 2019. [Online]. Available: <https://github.com/iovisor/bcc>.
- [20] IOVisor, *Iovisor/bpfftrace*, Nov. 2019. [Online]. Available: <https://github.com/iovisor/bpfftrace>.
- [21] IOVisor, *Syscount.py*, Oct. 2019. [Online]. Available: <https://github.com/iovisor/bcc/blob/master/tools/syscount.py>.
- [22] R. A. Kemmerer and G. Vigna, “Intrusion Detection: A Brief History and Overview,” *IEEE Security and Privacy 2002*, vol. 35, no. 4, pp. 27–30, Apr. 2002, ISSN: 1558-0814. DOI: [10.1109/MC.2002.1012428](https://doi.org/10.1109/MC.2002.1012428).
- [23] J. Keniston, A. Mavinakayanahalli, P. Panchamukhi, and V. Prasad, “Ptrace, Utrace, Uprobes: Lightweight, Dynamic Tracing of User Apps,” in *Proceedings of the Linux Symposium*, vol. 1, pp. 215–224. [Online]. Available: <https://www.linuxsecrets.com/kdocs/mirror/ols2007v1.pdf#page=41>.

- [24] F. Kerschbaum, E. Spafford, and D. Zamboni, “Using Internal Sensors and Embedded Detectors for Intrusion Detection,” *Journal of Computer Security*, vol. 10, pp. 23–70, Jan. 2002. DOI: [10.3233/JCS-2002-101-203](https://doi.org/10.3233/JCS-2002-101-203).
- [25] B. Kuperman and E. Spafford, “Generation of Application Level Audit Data via Library Interposition,” Sep. 1999.
- [26] *LTTng v2.11 - LTTng Documentation*, Oct. 2019. [Online]. Available: <https://lttng.org/docs/v2.11/>.
- [27] S. McCanne and V. Jacobson, “The BSD Packet Filter: A New Architecture for User-level Packet Capture,” *USENIX Winter*, vol. 93, 1992. [Online]. Available: <https://www.tcpdump.org/papers/bpf-usenix93.pdf>.
- [28] J. McCormack, P. Karlton, S. Angebrannt, *et al.*, *x11perf(1)*, Xorg. [Online]. Available: <https://www.x.org/releases/X11R7.7/doc/man/man1/x11perf.1.xhtml>.
- [29] L. W. McVoy, *lmbench*, Intel, Dec. 2019. [Online]. Available: <https://github.com/intel/lmbench>.
- [30] L. W. McVoy, C. Staelin, *et al.*, “lmbench: Portable tools for performance analysis,” in *USENIX annual technical conference*, USENIX, San Diego, CA, USA, 1996, pp. 279–294.
- [31] A. Merey, *Introducing stapbpf - SystemTap’s New BPF Backend*, Dec. 2017. [Online]. Available: <https://developers.redhat.com/blog/2017/12/13/introducing-stapbpf-systemtaps-new-bpf-backend/>.
- [32] J. Mogul, R. Rashid, and M. Accetta, “The Packer Filter: An Efficient Mechanism for User-level Network Code,” in *Proceedings of the Eleventh ACM Symposium on Operating Systems Principles*, ser. SOSP ’87, Austin, Texas, USA: ACM, 1987, pp. 39–51, ISBN: 0-89791-242-X. DOI: [10.1145/41457.37505](https://doi.org/10.1145/41457.37505). [Online]. Available: <http://doi.acm.org/10.1145/41457.37505>.
- [33] M. Moraes and J. McCormack, *x11perfcomp(1)*, Xorg. [Online]. Available: <https://www.x.org/releases/X11R7.7/doc/man/man1/x11perfcomp.1.xhtml>.
- [34] *NIT(4p) SunOS 4.1.1 Reference Manual*, Sun Microsystems Inc., Sep. 1990.
- [35] W. W. Peng and D. R. Wallace, *Software Error Analysis*. Silicon Press, 1995.
- [36] *ptrace(2) Linux User’s Manual*, Oct. 2019.
- [37] Red Hat, *Understanding How SystemTap Works Red Hat Enterprise Linux 5*. [Online]. Available: https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/5/html/systemtap_beginners_guide/understanding-how-systemtap-works.

- [38] S. Rostedt, *Documentation/ftrace.txt*, 2008. [Online]. Available: <https://lwn.net/Articles/290277/>.
- [39] R. Rubira Branco, “Ltrace internals,” in *Proceedings of the Linux Symposium*, vol. 1, pp. 41–52. [Online]. Available: <https://www.linuxsecrets.com/kdocs/mirror/ols2007v1.pdf#page=41>.
- [40] S. Soltesz, H. Pötzl, M. E. Fiuczynski, *et al.*, “Container-based operating system virtualization: A scalable, high-performance alternative to hypervisors,” in *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007*, ser. EuroSys ’07, Lisbon, Portugal: Association for Computing Machinery, 2007, pp. 275–287, ISBN: 9781595936363. DOI: [10.1145/1272996.1273025](https://doi.org/10.1145/1272996.1273025). [Online]. Available: <https://doi-org.proxy.library.carleton.ca/10.1145/1272996.1273025>.
- [41] A. B. Somayaji, “Operating System Stability and Security through Process Homeostasis,” PhD thesis, University of New Mexico, 2002. [Online]. Available: <https://people.scs.carleton.ca/~soma/pubs/soma-diss.pdf>.
- [42] A. B. Somayaji and H. Inoue, “Lookahead Pairs and Full Sequences: A Tale of Two Anomaly Detection Methods,” in *Proceedings of the 2nd Annual Symposium on Information Assurance Academic track of the 10th Annual 2007 NYS Cyber Security Conference*. NYS Cyber Security Conference, 2007, pp. 9–19. [Online]. Available: <http://people.scs.carleton.ca/~soma/pubs/inoue-albany2007.pdf>.
- [43] E. H. Spafford and D. Zamboni, “Intrusion Detection Using Autonomous Agents,” *Comput. Netw.*, vol. 34, no. 4, pp. 547–570, Oct. 2000, ISSN: 1389-1286. DOI: [10.1016/S1389-1286\(00\)00136-5](https://doi.org/10.1016/S1389-1286(00)00136-5). [Online]. Available: [http://dx.doi.org/10.1016/S1389-1286\(00\)00136-5](http://dx.doi.org/10.1016/S1389-1286(00)00136-5).
- [44] W. Stallings and L. Brown, *Computer security: principles and practice*. [Online]. Available: <http://www.informit.com/articles/article.aspx?p=782118>.
- [45] A. Starovoitov, “tracing filters with BPF,” The Linux Foundation, RFC Patch 0/5, Dec. 2013. [Online]. Available: <https://lkml.org/lkml/2013/12/2/1066>.
- [46] A. Starovoitov, “net: filter: rework/optimize internal BPF interpreter’s instruction set,” The Linux Foundation, Kernel Patch, Mar. 2014. [Online]. Available: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=bd4cf0ed331a275e9bf5a49e%206d0fd55dfc551b8>.
- [47] A. Starovoitov, *[v2,bpf-next,1/4] bpf: Enable program stats 1047415 diff mbox series*, Feb. 2019. [Online]. Available: <https://patchwork.ozlabs.org/patch/1047415/>.

- [48] A. Starovoitov and D. Borkmann, *bpf: introduce bounded loops*, Jun. 2019. [Online]. Available: <https://git.kernel.org/pub/scm/linux/kernel/git/davem/net-next.git/commit/?id=2589726d12a1b12eaaa93c7f1ea64287e383c7a5>.
- [49] Strace Project, *Strace*. [Online]. Available: <https://strace.io/>.
- [50] *strace(1) Linux User's Manual*, 5.3, Strace Project, Sep. 2019.
- [51] Sysdig Inc., *Draios/sysdig*, Nov. 2019. [Online]. Available: <https://github.com/draios/sysdig>.
- [52] L. Torvalds, *Linux/uapi/asm-generic/unistd.h*. [Online]. Available: <https://github.com/torvalds/linux/blob/master/include/uapi/asm-generic/unistd.h>.
- [53] A. M. Turing, "On Computable Numbers, with an Application to the Entscheidungsproblem," *Proceedings of the London Mathematical Society*, vol. s2-42, no. 1, pp. 230–265, Jan. 1937, ISSN: 0024-6115. DOI: [10.1112/plms/s2-42.1.230](https://doi.org/10.1112/plms/s2-42.1.230). eprint: <http://oup.prod.sis.lan/plms/article-pdf/s2-42/1/230/4317544/s2-42-1-230.pdf>. [Online]. Available: <https://doi.org/10.1112/plms/s2-42.1.230>.
- [54] K. Van Hees, "BPF, Trace, DTrace: DTrace BPF Program Type Implementation and Sample Use," The Linux Foundation, RFC Patch 00/11, May 2019, pp. 1–56. [Online]. Available: <https://lwn.net/Articles/788995/>.
- [55] D. Wagner and P. Soto, "Mimicry Attacks on Host-based Intrusion Detection Systems," in *Proceedings of the 9th ACM Conference on Computer and Communications Security*, ser. CCS '02, Washington, DC, USA: ACM, 2002, pp. 255–264, ISBN: 1-58113-612-9. DOI: [10.1145/586110.586145](https://doi.org/10.1145/586110.586145). [Online]. Available: <http://doi.acm.org/10.1145/586110.586145>.
- [56] V. Weaver, *Perf_event_open(2) linux user's manual*, Oct. 2019.

Appendices

A eBPF Design Patterns

Listing A.1: Handling large datatypes in eBPF programs.

```
1  /* This is way too large to fit within
2   * the eBPF stack limit of 512 bytes */
3  struct bigdata_t
4  {
5      char foo[4096];
6  };
7
8  /* We read from this array every time we want to
9   * initialize a new struct bigdata_t */
10 BPF_ARRAY(__bigdata_t_init, struct bigdata_t, 1);
11
12 /* The main hashmap used to store our data */
13 BPF_HASH(bigdata_hash, u64, struct bigdata_t);
14
15 /* Suppose this is a function where we need to use our
16 * bigdata_t struct */
17 int some_bpf_function(void)
18 {
19     /* We use this to look up from our
20      * __bigdata_t_init array */
21     int zero = 0;
22     /* A pointer to a bigdata_t */
23     struct bigdata_t *bigdata;
24     /* The key into our main hashmap
25      * Its value not important for this example */
26     u64 key = SOME_VALUE;
27
28     /* Read the zeroed struct from our array */
29     bigdata = __bigdata_t_init.lookup(&zero);
30     /* Make sure that bigdata is not NULL */
31     if (!bigdata)
32         return 0;
33     /* Copy bigdata to another map */
34     bigdata = bigdata_hash.lookup_or_try_init(&key, bigdata);
35
36     /* Perform whatever operations we want on bigdata... */
37
38     return 0;
39 }
```

B bpfbench Source Code

Listing B.1: The eBPF component of bpfbench.

```

1  /* bpfbench A better benchmarking tool written in eBPF.
2  * Copyright (C) 2020 William Findlay
3  *
4  * Heavily inspired by syscount from bcc-tools:
5  * https://github.com/iovisor/bcc/blob/master/tools/syscount.py
6  *
7  * This program is free software: you can redistribute it and/or modify
8  * it under the terms of the GNU General Public License as published by
9  * the Free Software Foundation, either version 3 of the License, or
10 * (at your option) any later version.
11 *
12 * This program is distributed in the hope that it will be useful,
13 * but WITHOUT ANY WARRANTY; without even the implied warranty of
14 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
15 * GNU General Public License for more details.
16 *
17 * You should have received a copy of the GNU General Public License
18 * along with this program. If not, see <https://www.gnu.org/licenses/>. */
19
20 #include <uapi/asm/unistd_64.h>
21 #include <linux/sched.h>
22 #include <linux/signal.h>
23
24 struct intermediate_t
25 {
26     u64 pid_tgid;
27     u64 start_time;
28 };
29
30 struct data_t
31 {
32     u64 count;
33     u64 overhead;
34 };
35
36 BPF_PERCPU_ARRAY(intermediate, struct intermediate_t, 1);
37 BPF_PERCPU_ARRAY(syscalls, struct data_t, NUM_SYSCALLS);
38
39 #ifdef FOLLOW
40 BPF_HASH(children, u32, u8);
41
42 RAW_TRACEPOINT_PROBE(sched_process_fork)
43 {
44     struct task_struct *p = (struct task_struct *)ctx->args[0];
45     struct task_struct *c = (struct task_struct *)ctx->args[1];
46
47     u32 ppid = p->tgid;

```

```

48
49     /* Filter ppid */
50     if (ppid != TRACE_PID && !children.lookup(&ppid))
51     {
52         return 0;
53     }
54
55     u32 cpid = c->tgid;
56
57     u8 zero = 0;
58
59     children.update(&cpid, &zero);
60
61     return 0;
62 }
63
64 RAW_TRACEPOINT_PROBE(sched_process_exit)
65 {
66     u32 pid = (bpf_get_current_pid_tgid() >> 32);
67
68     /* Filter ppid */
69     if (pid != TRACE_PID && !children.lookup(&pid))
70     {
71         return 0;
72     }
73
74     children.delete(&pid);
75
76     return 0;
77 }
78 #endif
79
80 TRACEPOINT_PROBE(raw_syscalls, sys_enter)
81 {
82     u64 pid_tgid = bpf_get_current_pid_tgid();
83
84     /* Maybe filter by PID */
85     #if defined(TRACE_PID) && defined(FOLLOW)
86     u32 pid = (pid_tgid >> 32);
87     if (pid != TRACE_PID && !children.lookup(&pid))
88     {
89         return 0;
90     }
91     #elif defined(TRACE_PID)
92     if (pid_tgid >> 32 != TRACE_PID)
93     {
94         return 0;
95     }
96     #endif
97
98     /* Don't trace self */
99     if (pid_tgid >> 32 == BPFBENCH_PID)
100     {
101         return 0;

```



```

102     }
103
104     int zero = 0;
105     struct intermediate_t *start = intermediate.lookup(&zero);
106     if (!start)
107     {
108         return 0;
109     }
110
111     /* Record pit_tgid of initiating process,
112      * we use this for error checking later */
113     start->pid_tgid = pid_tgid;
114     /* Record start time */
115     start->start_time = bpf_ktime_get_ns();
116
117     return 0;
118 }
119
120 TRACEPOINT_PROBE(raw_syscalls, sys_exit)
121 {
122     u64 pid_tgid = bpf_get_current_pid_tgid();
123
124     /* Maybe filter by PID */
125     #if defined	TRACE_PID) && defined(FOLLOW)
126     u32 pid = (pid_tgid >> 32);
127     if (pid != TRACE_PID && !children.lookup(&pid))
128     {
129         return 0;
130     }
131     #elif defined	TRACE_PID)
132     if (pid_tgid >> 32 != TRACE_PID)
133     {
134         return 0;
135     }
136     #endif
137
138     /* Don't trace self */
139     if (pid_tgid >> 32 == BPFBENCH_PID)
140     {
141         return 0;
142     }
143
144     int zero = 0;
145     int syscall = args->id;
146
147     /* Discard restarted syscalls due to system suspend */
148     if (args->id == __NR_restart_syscall)
149     {
150         return 0;
151     }
152
153     struct data_t *data = syscalls.lookup(&syscall);
154     struct intermediate_t *start = intermediate.lookup(&zero);
155     if (start && data)

```

```
156 {
157     /* We don't want to count twice for calls that return in two places */
158     if (pid_tgid != start->pid_tgid)
159     {
160         return 0;
161     }
162     data->count++;
163     data->overhead += bpf_ktime_get_ns() - start->start_time;
164 }
165 if (start)
166 {
167     start->pid_tgid = 0;
168     start->start_time = 0;
169 }
170
171 return 0;
172 }
```