

Mitigating Spectre-PHT using Speculation Barriers in Linux BPF

Luis Gerhorst

Friedrich-Alexander-Universität
Erlangen-Nürnberg (FAU)
Germany
gerhorst@cs.fau.de

Henriette Herzog

Ruhr-Universität Bochum (RUB)
Germany
henriette.hofmeier@rub.de

Peter Wägemann

Friedrich-Alexander-Universität
Erlangen-Nürnberg (FAU)
Germany
waegemann@cs.fau.de

Maximilian Ott

Friedrich-Alexander-Universität
Erlangen-Nürnberg (FAU)
Germany
ott@cs.fau.de

Rüdiger Kapitza

Friedrich-Alexander-Universität
Erlangen-Nürnberg (FAU)
Germany
ruediger.kapitza@fau.de

Timo Hönig

Ruhr-Universität Bochum (RUB)
Germany
timo.hoenig@rub.de

ABSTRACT

High-performance IO demands low-overhead communication between user- and kernel space. This demand can no longer be fulfilled by traditional system calls. Linux’s extended Berkeley Packet Filter (BPF) avoids user-/kernel transitions by just-in-time compiling user-provided bytecode and executing it in kernel mode with near-native speed. To still isolate BPF programs from the kernel, they are statically analyzed for memory- and type-safety, which imposes some restrictions but allows for good expressiveness and high performance. However, to mitigate the Spectre vulnerabilities disclosed in 2018, defenses which reject potentially-dangerous programs had to be deployed. We find that this affects 24 % to 54 % of programs in a dataset with 844 real-world BPF programs from popular open-source projects. To solve this, users are forced to disable the defenses to continue using the programs, which puts the entire system at risk.

To enable *secure and expressive* untrusted Linux kernel extensions, we propose BERRIFY, an enhancement to the kernel’s Spectre defenses that reduces the number of BPF application programs rejected from 54 % to zero. We measure BERRIFY’s overhead for all mainstream performance-sensitive applications of BPF (i.e., event tracing, profiling, and packet processing) and find that it improves significantly upon the status-quo where affected BPF programs are either unusable or enable transient execution attacks on the kernel.

CCS CONCEPTS

• **Security and privacy** → **Operating systems security**; **Side-channel analysis and countermeasures**; • **Software and its engineering** → *Automated static analysis*; Software safety; **Just-in-time compilers**.

KEYWORDS

eBPF, Transient Execution Attacks, High-Performance Networking

1 INTRODUCTION

Operating systems rely on system calls to allow the controlled communication of isolated processes with the kernel and other processes. Every system call includes a processor mode switch from the unprivileged user mode to the privileged kernel mode. Although

processor mode switches are the essential isolation mechanism to guarantee the system’s integrity, they induce direct and indirect performance costs as they invalidate parts of the processor state [1]. In recent years, high-performance network [2] and storage hardware [3] has made the user/kernel transition overhead the bottleneck for IO-heavy applications. To make matters worse, security vulnerabilities in modern processors [4], such as Meltdown [5], have prompted kernel mitigations that further increase the transition overhead.

To overcome this performance barrier, modern operating systems implement multiple alternatives to traditional blocking system calls. The three main approaches are *asynchronous IO* (io_uring [6]), *kernel-bypass* (DPDK [2]), and *safe kernel extensions* (BPF [7]). This work focuses on safe kernel extensions because they offer lower latency than asynchronous IO [8] and, unlike kernel-bypass, integrate well with the existing OS networking stack [7].

Safe kernel extensions offer very low-overhead interaction between user and kernel space. Their applications include network IO [7, 9, 10], memory optimization [11], threat detection [12], isolation [13], tracing [14], scheduling [15], and storage [16]. Specifically, Linux’s extended Berkeley Packet Filter (BPF) allows unprivileged user processes to load safety-checked bytecode into the kernel, which executes at near-native speed. Users typically develop BPF programs in a high-level programming language (e.g. C or Rust), which is compiled into BPF bytecode. This bytecode is verified in regard to its safety before being just-in-time compiled. Invoking BPF programs in the kernel and calling kernel functions from within BPF is much faster than the respective switch to/from user context [1]. However, this performance advantage comes at a cost: Executing untrusted code in the kernel address space requires precise static analysis of the untrusted program to maintain isolation.

While Meltdown and similar domain-bypass transient execution attacks (e.g., where the hardware crosses address-space boundaries [4]) can be efficiently prevented in hardware, the cross/in-domain transient execution attacks discovered in 2018 [17] still require software defenses on all high-performance central processing units (CPUs) [18, 19] to date. Two cross/in-domain transient execution attacks are a particular challenge for BPF: Spectre-PHT exploits the Pattern History Table’s potential for branch misprediction, and Spectre-STL exploits speculative store bypass (via the CPU’s Store To Load buffer). Branches and stores are common

This is the author’s version of the work. It is posted here for your personal use. Not for redistribution.

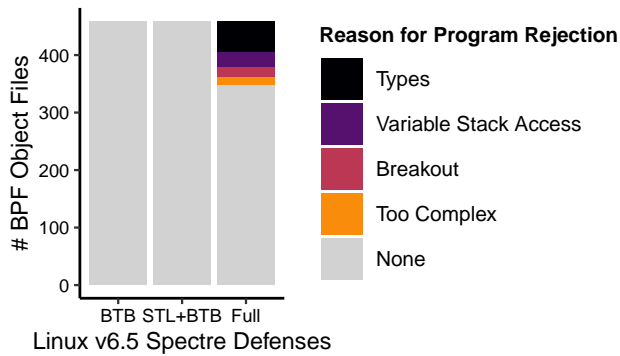


Figure 1: Reasons for program rejection in Linux when enabling defenses against transient execution attacks for 459 real-world BPF test, example, and application object files. While the Spectre-BTB and -STL defenses do not trigger program rejections, enabling Spectre-PHT defenses (Full, i.e., PHT+STL+BTB) prevents 24% of all objects (and even 54% of application objects) from being used.

in BPF bytecode, but malicious programs can use these very instructions to break out of their sandbox using transient execution. For example, following a mispredicted branch, the program can perform an unchecked out-of-bounds (OOB) access and leak the resulting kernel secret using a covert channel. Preventing these attacks entirely in hardware requires far-reaching changes to processor designs with significant performance impacts [20] (e.g., 9% to 15% for [21]), that high-performance CPU vendors to date have not implemented [22, 23, 24].

However, developing practical Spectre-resilient sandboxes is an open research problem. To the best of our knowledge, BPF is the only widely deployed sandbox that attempts to implement full software-defenses against Spectre-PHT, Spectre-BTB, and Spectre-STL. All other widely deployed sandboxes such as Java [25], WebAssembly (Wasm) [26, 27], and JavaScript [19, 28] instead rely on process isolation. This is the approach recommended by Intel [18] and V8 developers [29], but it is not applicable to BPF.

To implement the sandbox, the kernel statically analyzes the control and data flow of the BPF programs before allowing them to execute in its address space. Traditionally, the kernel only considered BPF program paths that could actually execute architecturally during this static analysis. In response to the Spectre-PHT and Spectre-STL vulnerabilities disclosed in 2018, kernel developers have extended the BPF verifier [30] to prevent leaks from transiently-executed BPF program paths. This includes (a) inserting instructions to make speculation safe (e.g., index masking) [31, 32, 33], (b) inserting instructions to prevent speculation (e.g., x86-64 Ifence) [34], or (c) statically analyzing the behavior on speculative code paths to ensure they are safe, in turn rejecting the program entirely if it exhibits any unsafe behavior [35]. The latter is required to prevent Spectre-PHT, but it severely limits BPF because it prevents the user from loading their extension into the kernel altogether.

To analyze the Spectre-PHT defenses' impact on BPF, we collect 459 BPF object files (each containing one or multiple BPF programs)

from open-source projects and enable defenses. The results, shown in Figure 1, confirm that the number of rejections is significant and are further discussed in Section 5.3. Restructuring the program to not exhibit unsafe transient behavior is tedious, as users are usually not familiar with transient execution vulnerabilities and only indirectly control the compiler-generated bytecode in their high-level sources.

To solve this problem of pessimistically rejected BPF programs, our work extends the upstream BPF verifier [30] by implementing BERRIFY, an improved defense approach for Spectre. The core benefit of BERRIFY is that it reduces the number of programs that cannot be automatically mitigated and are, thus, rejected by the state-of-the-art BPF verifier.

Our main contributions are fivefold:

- **Domain Analysis:** To inform our design decisions, we statically analyze 844 real-world BPF programs from six popular software projects regarding their code size and complexity. This domain analysis gives us a detailed picture of the BPF program landscape.
- **Formal Security Notions:** Real-world Spectre defenses lack behind the theoretical foundations [36]. We analyze BPF's security in the light of transient execution attacks using formal security notions from the literature.
- **Formal-Model-Based Analysis:** Based on formal security notions, we design the BERRIFY defense to eliminate the most common causes of program rejection using precise defenses with a minimum expected performance impact.
- **Evaluation:** We evaluate the performance impact of BERRIFY on the three most common use cases for BPF: event tracing, continuous profiling, and network load balancing. We find that BERRIFY is lightweight to applications as a whole.
- **Implementation:** We publish BERRIFY for the v6.5 Linux kernel. Our patches soundly combine speculation barriers and static analysis to not increase the kernel's attack surface.

Besides our five main contributions, we discover multiple Linux kernel bugs during our analyses. Based on our findings, we contribute proof-of-concept exploits [37] and fixes [38, 39] that have been accepted into upstream Linux.

2 BACKGROUND

This section discusses the transient execution vulnerabilities relevant to BPF and presents the respective defenses in Linux v6.5. We focus on cross/in-domain transient execution attacks [4] but not domain-bypass transient execution attacks (including Meltdown / Spectre v3 [5], MDS [40], L1TF [41]) because they can be efficiently addressed in hardware and are not specific to BPF. We focus on the **Spectre-PHT (v1)**, **-STL (v4)**, and **-BTB (v2)** attacks because these are the longest-standing vulnerabilities to which new processors are still vulnerable [22, 23, 24].

2.1 Transient Execution Vulnerabilities

While regular (architectural) timing side-channel attacks can happen only when the program computes explicitly on sensitive data (e.g., cryptographic keys, or also kernel pointers with kernel address space layout randomization, KASLR), timing side-channel attacks based on transient execution can also target victim code

that does not explicitly work with sensitive data. They become possible whenever a victim program encodes secrets into side channels (e.g., the cache) during transient execution.

2.1.1 Speculation Triggers. Transient execution CPU vulnerabilities are commonly grouped by the microarchitectural component that speculates. For BPF, Spectre-PHT, -STL, and -BTB are of particular relevance [42, 43]. Spectre-PHT [17] includes all transient execution vulnerabilities based on conditional branches as they use the PHT for target prediction. Spectre-STL exploits the fact that stores may not always become visible to subsequent loads via the STL buffer. While indirect branches (which enable Spectre-BTB attacks [17]) also occur in BPF, they are not as common and the program cannot use them directly. Because of this, the *retpoline*-based defense against Spectre-BTB is by default always enabled for privileged and unprivileged BPF programs, as it is also the case for the rest of the kernel [44, 45]. We therefore stick to the default and always keep them enabled.

2.1.2 Side Channels. Transient execution attacks can use a variety of shared hardware components to communicate sensitive data to attackers. Notably, this does not only include caches but also microarchitectural execution *ports* with simultaneous multithreading [46]. In this work, we assume the established **constant time (CT) leakage model** to capture all these incidental channels [43] by assuming that all accessed addresses are potentially leaked. This includes both instruction- and data addresses, but not the respective values at these addresses. From this, it follows that one must not branch based on secret data nor use secrets as memory offsets.

2.1.3 Unsafe Information-Flow. The CPU may even load additional sensitive data (erroneously) during speculative execution. For example, *speculative type confusion* [47] can cause the CPU to load sensitive data from an attacker-controller address and subsequently leak the data through a covert channel. To prevent any such unsafe information flow, Speculative Constant-Time (SCT) [48] and Speculative Non-Interference (SNI) [49] are notable formal notions that capture the information-flow properties of the victim program that enable transient execution attacks [43]. They are commonly referred to as *speculative security properties*.

SCT extends the constant-time programming paradigm, which modern cryptographic programs commonly use, to prevent Spectre gadgets by ensuring that the code executed speculatively does not leak sensitive data [43, 48]. A program that satisfies SCT must only operate on sensitive data using CT processor instructions [50, 51, 52, 53]. SCT can be efficiently enforced using type systems that model the sensitive data and its permitted operations [54, 55].

SNI formalizes the intuition that speculation can only leak data, which the program already leaks during architectural execution [43, 49, 56, 57]. For example, if some scalar is already leaked into a side channel architecturally, there is no point in protecting it from speculative leakage. To enforce SNI for a sandboxed program, [49] uses symbolic execution and a satisfiability modulo theories (SMT) solver. Like SCT, SNI not only protects the sandbox runtime (e.g., the kernel) but also the sandboxed program itself (e.g., a BPF program) from Spectre attacks (i.e., they prevent *poisoning attacks*) [58].

While enforcing either SCT or SNI prevents Spectre gadgets, SCT suffers from false positives while SNI cannot be efficiently

enforced for arbitrary programs. Applying them to BPF is therefore not necessarily useful. Neither is applying them trivial as the verifier was not developed with SCT or SNI in mind. However, in Section 3, we retrospectively analyze whether the BPF verifier enforces SCT or SNI for any of BPF's data types.

2.2 Spectre Defenses of Linux BPF

In this section, we present the first thorough analysis of BPF's Spectre defenses. We build upon previous talks and reports on the topic [44, 45, 59] but add more details and, in the following section, analyze whether BPF enforces any formal speculative security properties. We enable future work as BPF's defenses to date are undocumented [30]. Also, the number of contributors is low [31, 33, 34, 35] and both our work [37, 38, 39, 60] and others [32, 61] have repeatedly discovered bugs and misleading code.

2.2.1 Attacker Model. We consider both unprivileged attackers that can load arbitrary unprivileged BPF programs into the kernel, and attackers that can coerce a privileged user into loading architecturally-safe privileged BPF programs into the kernel on their behalf. As of Linux v6.5, the kernel can only effectively protect itself from the former.

For programs from unprivileged users, the kernel must verify both speculative and architectural security. For instance, the verifier must prevent unprivileged BPF programs from printing pointer values in architectural execution, as well as prevent them from leaking pointers into covert channels during speculative execution.

The second type of attacker we consider coerces a privileged user into loading an architecturally safe BPF program into the kernel. While the kernel heavily limits the operations permitted in unprivileged BPF, privileged BPF users are only restricted regarding CPU time (bounded loops) and memory usage (memory/type-safety). Consequently, BPF's Spectre-PHT and Spectre-STL defenses are not active by default as the user has to trust the program anyway. The kernel instead assumes that privileged users manually check their programs for Spectre gadgets to prevent attacks on the BPF program and on the whole kernel. We argue this assumption is unrealistic, as most developers are not familiar with transient execution attacks. Further, when writing BPF programs in high-level languages such as C, Rust, or Python, gadgets are not directly visible and may only occur because of compiler optimizations. Also, BPF programs are frequently generated ad-hoc [62] and thus never undergo code review for Spectre gadgets. From this, it follows that activating defenses for privileged BPF is desirable, but currently impossible without restricting all privileged users.

2.2.2 Architectural Safety. To understand the existing Spectre defenses, which ensure speculative safety, we first briefly present the design of the verifier and its mechanisms to ensure architectural safety. The main goal of the verifier is to limit the damage malicious or buggy BPF programs can cause. For this, it verifies a bounded execution time and memory safety but not functional correctness. Bounded execution time is mainly useful for allowing BPF use in interrupt contexts, while memory safety prevents memory leaks and program bugs that would easily enable kernel exploits.

To restrict data flow into the BPF program's registers and stack, the verifier enumerates all possible paths through the BPF program

and simulates each path's execution. To verify memory and type safety, the kernel analyzes the types (mainly different classes of pointers and scalars [63]) and the value ranges of the scalars (for which it uses *tristate numbers* [64]). This allows the verifier to ensure that all scalars used as pointer-offsets only point to locations owned-by or borrowed-to¹ the BPF program. For example, this prevents OOB accesses to the BPF stack and to network packets (the BPF program can manipulate packets directly using its context pointer [60]). Further, accessing uninitialized stack slots and registers is also prohibited. In summary, the BPF program can only access memory locations (i.e., load their value into registers) to which the kernel grants explicit access.

The kernel not only restricts how data flows into BPF program registers and stack but also limits how the program uses the data in these locations. For this, the verifier distinguishes between pointers and scalars. This allows the verifier to limit how the BPF program processes data and ensure that only safe operations are executed. For example, programs can only dereference pointers (at valid offsets) and not cast them into scalars. Further, they can only use scalars in arbitrary ALU operations and conditionally branch based on their value [60]. Therefore, the BPF program can only use kernel pointers in CT operations (except for dereferencing them), while scalars are only restricted so that the program cannot cast them into pointers.

2.2.3 Design Goals. Extending the design used to ensure architectural safety, the kernel's Spectre defenses for BPF ensure that the same restrictions also apply in speculative execution. The goal in Linux v6.5 is only to protect the kernel from BPF programs and associated user applications. The BPF program itself is not protected against other user applications that might exploit Spectre gadgets to retrieve scalars (including cryptographic keys) the program processes. Further, the existing defenses against Spectre-PHT and Spectre-STL try to be transparent to users as much as possible, as the kernel applies them to the bytecode during verification time. The kernel does not require the user or source compiler to insert speculation barriers into the program. However, the defenses fail to be entirely transparent, as they can impact performance and prevent programs from passing verification.

2.2.4 Spectre-STL. To defend against Spectre-STL (v4), the kernel inserts speculation barriers after *critical* stores to the BPF stack [34]. A store is *critical* if speculatively bypassing it would make otherwise inaccessible data (or operations upon data) available to the program. For example, initializing a stack slot or overwriting a scalar value with a pointer are critical stores as the kernel prohibits reading uninitialized stack slots and dereferencing scalars. The verifier can only skip the insertion of a barrier if a scalar is overwritten with another scalar [66]. Confusing one scalar value with another scalar cannot lead to OOB memory accesses since the verifier enforces pointer limits using branchless logic (e.g., masking). This is further discussed in the following paragraph, as it is also required to defend against Spectre-PHT.

2.2.5 Spectre-PHT. While Spectre-STL effectively causes the CPU to speculatively bypass a store, Spectre-PHT (v1) causes the CPU

¹For example, this includes pointers from the `bpf_ringbuf_reserve()` helper which must be either submitted or discarded subsequently [65].

to bypass evaluation of a branch condition. Therefore, conditional branches can no longer be reliably used to ensure memory and type safety. To defend against this, the verifier prevents OOB accesses using branchless logic [31, 32, 33] and type confusion by verifying architecturally-impossible speculative execution paths [35].

Branchless Bounds Enforcement. This includes simple masking but also more complex instruction sequences when required. First, for arrays (i.e., BPF maps) the kernel can simply ceil the size to a power of two and apply the respective index-mask before the access [31]. Second, for pointers, the kernel deducts the bounds from the conditional branches that lead to the pointer dereference. Then the kernel also enforces them directly before the access using a special sequence of ALU operations [32] (additionally, the verifier has to check a remaining corner case using verification of an architecturally impossible speculative code path). Third, for the stack, the kernel enforces that all offsets are constant to simplify the implementation [33]. Finally, there is one exception where the kernel allows the program to access OOB memory speculatively for practical reasons [67]. We have verified that the specific memory layout in this case does not expose any secrets [37].

Verification of Speculative Execution Paths. While branchless bounds enforcement prevents the BPF program from accessing forbidden data, it cannot prevent the program from using kernel pointers or scalars in an unsafe manner (e.g., dereferencing a scalar). To ensure speculative type safety, the kernel also simulates and checks the execution paths that include mispredicted branches [35]. The verifier only allows transient behavior that is also permitted architecturally. Even though this is a sound approach, it suffers from false positives because not every architecturally unsafe operation enables a Spectre attack [49] (and the architectural verification logic itself also already suffers from false positives). Using separate simulation and verification logic for the transient domain could resolve this, but it would increase the verifier's attack surface even further.

In summary, branchless bounds enforcement and verification of mispredicted branches enable reliable defense against Spectre-PHT.

3 SECURITY ANALYSIS

In this section, we analyze the foundations of BPF's Spectre defenses. We do this both from a hardware and a formal perspective. Regarding the hardware, we summarize the required instruction-set properties for the defenses (i.e., the hardware-software contract). Further, we retrospectively analyze the speculative security properties the verifier enforces.

3.1 Hardware-Software Contract

To terminate speculative execution after an STL misprediction, the verifier relies on speculation barriers. For this, it uses the undocumented `nospec` BPF bytecode instruction that is not available to user space. Notably, on x86-64, the just-in-time (JIT) implements these using the `lfence` instruction, which is in line with Intel's recommendation for Spectre-STL and -PHT [18, 68]. Independent research has confirmed that `lfence` terminates speculation [69], therefore, we deem them reliable. On ARM64, the JIT compiler lowers `nospec` to a `no-op` because the hardware already defends

against Spectre-STL in firmware [70]. To perform safe pointer arithmetic, the verifier further relies on ALU instructions that have data-independent timing [50, 52, 53].

3.2 Formal Security Notions

From a formal perspective, the verifier enforces a mix of speculative security properties for BPF programs. Overall, the verifier attempts to prevent the BPF program from speculatively breaking out of its sandbox and is successful in that (to the best of our knowledge), except for the exception discussed in Section 2.2.5 [37]. *Speculative-breakout attacks* are therefore prevented. Regarding the exception, we are also unable to construct an exploit that leaks *sensitive* kernel data [37]. However, while this in summary protects the kernel, the verifier does not enforce SNI [43] as processed scalars that developers were careful to never leaked architecturally, can still leak after a speculative scalar confusion due to Spectre-STL. Regarding SCT, the verifier enforces it only for pointers, but with the exception that the program can dereference the pointer itself. Overall, retrofitting a single formal speculative security property to BPF does not appear viable. However, we still find SNI and SCT to be of high value as we discover multiple kernel bugs by analyzing BPF’s defenses with them in mind [38, 39, 60].

4 PROBLEM STATEMENT

In this section, we summarize the limitations of the existing BPF Spectre defenses that motivate our work. We identify three key problems:

Problem # 1 *Defenses do not protect user secrets.* The verifier does not ensure SNI, and there is no support for BPF programs in protecting cryptographic keys they process. While users cannot insert speculation barriers manually in Linux v6.5, the kernel could offer basic support for this by exposing the internal BPF bytecode instruction to insert barriers (nospec) to users.

Problem # 2 *Spectre-STL defenses reduce performance.* The verifier inserts speculation barriers that reduce instruction-level parallelism (even if there is no misprediction) and, thereby, potentially limit performance.

Problem # 3 *Spectre-PHT defenses lead to program rejections.* If the verifier finds a speculative execution path (following a simulated misprediction) that performs prohibited operations, it rejects the whole program [35]. This forces users to resort to tedious restructuring of the source code or abandon BPF completely (which has a much higher performance impact than the speculation barriers [1]). Frequently, this development-overhead motivates users to simply disable defenses altogether [71], which enables BPF-based exploits as shown in the original Spectre paper [17].

Out of these three, we deem *Problem # 3* to be the most limiting to users and, therefore, focus on it in this work. Nevertheless, our evaluation also addresses *Problem # 2* by including the Spectre-STL defenses explicitly to quantify their overhead. Further, we propose a portable and performant solution to *Problem # 1* in Section 7.

5 DOMAIN ANALYSIS

As discussed in the previous section, Linux v6.5’s Spectre-PHT defenses cause BPF program rejections when the verifier cannot prove some transient execution is safe. In this case it prevents unprivileged users from using BPF altogether (*Problem # 3*). At the same time, the Spectre-STL defenses can negatively impact the execution time as they insert speculation barriers (*Problem # 2*). However, to-date it is unclear to which extent these defenses are actually triggered in real-world BPF programs.

In this section, we present a domain analysis of the BPF program landscape to analyze how BPF programs are affected by the kernel’s Spectre defenses. We collect 459 BPF object files containing 844 individual BPF programs from six popular open-source projects and analyze the number of objects rejected and speculation barriers inserted. To support future research, we publish our tools for collecting and analyzing the programs.²

5.1 Dataset

We include a diverse set of programs taken from the Linux kernel selftests and BPF samples [72], libbpf examples [73], Prevail [74], BCC [75], the Cilium Kubernetes Container Networking Interface (CNI) [76], the Parca Continuous Profiler [77], and the Loxilb Network Loadbalancer [78]³. We observe that programs from the Linux kernel selftests are often designed to only test a specific kernel interface and therefore very small. This motivates us to further group the programs into *test or example* and *application* programs (171 programs from 50 object files). The *application* group only includes programs from the BCC, Loxilb, Parca, and Cilium projects, as well as 2 specific object files from the Linux selftests with programs adapted from real-world applications.

By default, the kernel only applies Spectre-BTB defenses to privileged BPF programs to avoid possible performance overheads and rejections due to Spectre-PHT and -STL defenses. Also, because of this potential rejection by the verifier, few BPF programs are designed for unprivileged use. To ensure a reasonably-sized dataset, we modify the kernel to support activation of Spectre-PHT and -STL defenses at runtime. We expect the results to carry over to unprivileged BPF programs because their code is generally similar (only differing in the interfaces used) and generated by the same compiler. Upstream kernel developers have also used this technique to analyze the impact of the Spectre-PHT defenses on pointer arithmetic [80].

To inform our analysis, we first measure the number of BPF bytecode instructions per program. We observe that the mean number of instructions per program is only 40 overall. For the programs classified as applications, the mean is 46 instructions per program, while the average is 559 since there are some programs that are close to the verifier’s complexity limit of 1×10^6 instructions (e.g., Parca’s stack sampler). Overall, the low mean number of instructions is expected as most BPF programs only implement a fast path or policy decision in the kernel.

²<https://gitlab.cs.fau.de/un65esoq/bpf-spectre>

³Despite our best efforts, we are unable to use most of the BPF objects from the Prevail paper because they are not compatible with Linux v6.5. For the Cilium project, only the `bpf_socket` program is compatible with our `bpf_tool`-based toolchain [79].

5.2 Speculation Barriers

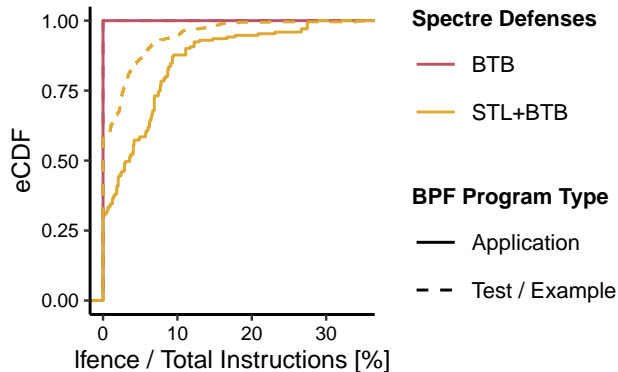


Figure 2: Percentage of speculation barriers in 844 BPF programs as a fraction of the total number of machine instructions with Spectre-STL defenses active.

We first analyze how many speculation barriers the verifier inserts to defend against Spectre-STL. Figure 2 shows the percentage of barriers per BPF program for all collected object files that are compatible with our kernel version. Comparing the programs classified as test or example programs with those from application projects, we find that application programs require more barriers per instruction on average (2.2 % instead of 1.0 % overall). This is likely because they are usually more complex and thus cannot work with registers exclusively. Simply counting the number of speculation barriers, of course, only indirectly relates to a real-world performance overhead, which heavily depends on the exact location of the barrier (e.g., in a tight loop or only at the beginning of the program during initialization). In Section 8, we will, therefore, analyze the performance overhead of the Spectre-STL defenses in real-world execution-time benchmarks. Still, we expect the overall overhead to be much lower compared to the performance overhead if BPF were not used at all [1].

5.3 Program Rejections

While the Spectre-STL defenses of Linux v6.5 only slow down the BPF program, the Spectre-PHT defenses can lead to the entire program being rejected. This forces users to resort to tedious restructuring of the source code or abandon BPF completely. In this section, we analyze how many programs from real-world projects the verifier actually rejects.

Figure 1 from Section 1 shows the number and cause of program rejections with only Spectre-BTB defenses active (*BTB*), Spectre-STL and -BTB defenses active (*STL+BTB*), and with additional Spectre-PHT defenses active (*Full*, i.e., *PHT+STL+BTB*). As expected, the Spectre-STL defenses cause no rejections. However, with the Spectre-PHT defenses active, the verifier rejects 24 % of the 459 BPF object files we have collected. Excluding test or example programs, we even find that 54 % of the remaining 50 application object files are rejected (likely because they are usually more complex and, therefore, more likely to contain any unsafe speculative behavior). We group the causes for rejections into four categories:

Table 1: Number of BPF programs rejected with the existing Linux v6.5 defenses for different software projects.

Project	# BPF Object Files	# Rejected
Linux Kernel Selftests	355	80
BCC	39	19
Linux BPF Samples	35	5
Prevail	14	1
libbpf Examples	7	1
Loxilb	4	3
Parca	4	3
Cilium	1	0

- **Types (12 %):** These are type errors on speculative program paths. For example, speculatively casting as scalar into a pointer and dereferencing it.
- **Variable Stack Access (5 %):** The program contains a speculative or architectural stack access for which the verifier cannot statically compute the offset. To simplify the Spectre-PHT defenses, kernel developers have disallowed variable-stack accesses [33].
- **Breakout (4 %):** Unsafe speculative execution where the BPF program accesses locations (memory or uninitialized registers) owned by the kernel.
- **Too Complex (3 %):** Because the verifier has to check the additional speculative paths, some programs exceed the verifier’s complexity limits (e.g., a maximum path length of 1×10^6 instructions).

In summary, the verifier rejects BPF programs due to a diverse set of unsafe behaviors that can lead to transient execution attacks on the kernel if the program leaks the resulting secret into a side channel subsequently.

In Table 1, we further analyze the number of rejections per software project. As expected, there is no notable difference in the extent to which the rejections affect the different projects since most of the bytecode-level properties that lead to unsafe speculation do not directly map to high-level constructs in the source code. The only exceptions are architectural variable-stack accesses [33], which users can avoid by not allocating arrays on the stack. In summary, all but one of the projects we have analyzed (that is Cilium, for which we only have one program compatible with our toolchain) are negatively affected by the Spectre-PHT defenses.

To conclude, we find that the Spectre-PHT defenses in Linux v6.5 are significantly more limiting to users than the Spectre-STL defenses, showcasing the significance of *Problem # 3*. The defenses cause users to resort to disabling defenses altogether (e.g., [71]), which opens the door to introducing dangerous arbitrary-read gadgets into the kernel. In the following section, we present **BERRIFY**, which defends against Spectre-PHT without rejecting the whole BPF programs.

6 DESIGN

As outlined in the previous sections, the current BPF verifier pessimistically rejects programs that pose a potential security threat. While this is preferable to leaving the system vulnerable to transient

```

I: // reg = is_ptr ? public_ptr : scalar;
A: if (!is_ptr) goto C; // mispredict
B: value = *reg;
   covert_channel[value];
C: exit();

```

Figure 3: Example from a BPF program that contains a speculative type confusion gadget and is rejected by the Linux v6.5 BPF verifier without BERRYFY.

execution attacks, it also limits BPF’s usability. These rejections seem unnecessary, especially considering that effective defense mechanisms exist, which could be included in the program instead.

In this section, we present the idea and design of BERRYFY, not only a solution to *Problem #3* but also a generic technique to build Spectre-resistant software sandboxes based on verification. BERRYFY optimistically attempts to verify all speculative execution paths and only falls back to speculation barriers when unsafe behavior is detected. Importantly, we fully reuse existing verification logic for this to balance verifier complexity (and thereby potential for kernel bugs) with BPF execution-time overheads. Surprisingly, our evaluation shows that this approach results in low application overheads because invocation-latency is more critical to BPF’s performance than code execution-time. While we only implement BERRYFY’s approach for Spectre-PHT, it can also be applied to *Problem #2* to reduce the number of speculation barriers required for Spectre-STL defenses [34, 66], and to Spectre-BTB with Intel’s Control Flow Enforcement Technology (CET) [81]. We first present our idea by example and then discuss its security and performance.

6.1 Fence or Verify

Based on our analysis of the BPF program landscape, we have identified the existing Spectre-PHT defenses to be the main real-world issue because they prevent applications from using BPF altogether. In this section, we analyze a minimal example, shown in Figure 3, of a BPF program where BERRYFY successfully prevents only the critical transient execution while the Linux v6.5 BPF verifier rejects the entire program [35]. The program executes either with `is_ptr` true or false, and we assume logic in block I implements `reg = is_ptr ? public_ptr : scalar`. We exclude the corresponding branch in block I in this analysis to give a concise example. Figure 4 illustrates all execution paths through the BPF program. Path 1 and Path 2 are architecturally possible paths through the program. By simulating the execution of each basic block, the verifier computes the invariants that will hold after executing the block based on the inputs, the block’s instructions, and the condition of the final branch that ends the basic block (e.g., `reg` will be equal to `public_ptr` and `is_ptr` must be true, when we go to block B after A). We list these invariants on the edges connecting the basic blocks in Figure 4.

Only considering the architectural execution paths, the program could be accepted because the verifier understands that block B will only execute if `reg` contains a valid pointer (in Path 1). Thus, the program does not exhibit any architecturally unsafe behavior and was therefore accepted by the kernel before 2018. However, with the defenses active (specifically [35]), the program is rejected as the verifier also simulates branch mispredictions, which leads to

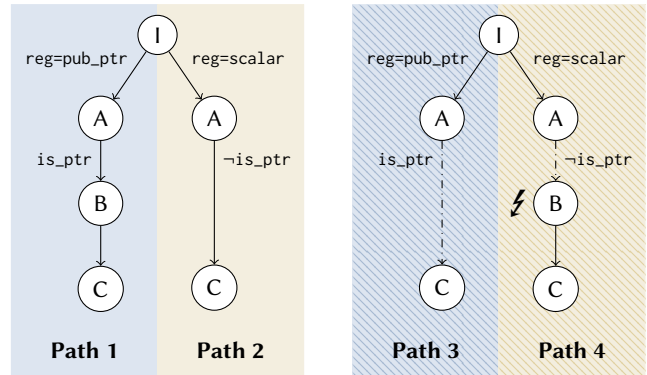


Figure 4: Execution paths through the BPF program from Figure 3. Paths marked with \rightarrow denote architectural execution, while \dashrightarrow indicates speculative execution. ! indicates unsafe behavior the verifier must prevent.

the speculative execution paths (Path 3 and 4) shown in Figure 4. With Spectre defenses, the verifier rightly finds the program to be unsafe because it performs a forbidden pointer dereference on Path 4. Specifically, when block B executes with a scalar in `reg`, the program breaks out of its sandbox by accessing an arbitrary address (! in Figure 4). Subsequently, the resulting kernel secret could be leaked into a covert channel (denoted as `covert_channel[value]` in the code example). In the Linux v6.5 kernel, this unsafe behavior causes the whole BPF program to be rejected, which either prevents unprivileged users from using the program altogether or forces privileged users to disable the Spectre defenses.

Our solution, BERRYFY, solves this by dynamically falling back to inserting a speculation barrier whenever the verifier detects unsafe program behavior after a simulated misprediction. By utilizing speculation barriers instead of program rejection to protect against unsafe program behavior, BERRYFY successfully addresses *Problem #3*. Further, whenever the verifier encounters a barrier while verifying a speculative path, it can stop verifying this path. Both of these modifications are based upon the insight that speculation barriers behave like program exits when the execution is transient. In the example from Figure 3, BERRYFY triggers the insertion of a barrier at the beginning of block B before `reg` is dereferenced. This still allows the CPU to mispredict the branch but terminates the speculative execution path when block B is reached. In contrast, the speculative execution in Path 3 is not problematic because no unsafe operation is performed in block C, even though the CPU mispredicted the branch at the end of block A. BERRYFY rightly detects this and does not insert a speculation barrier at the beginning of block C.

6.2 Security

BERRYFY only allows transient behavior that was already allowed architecturally and prevents all other execution paths using speculation barriers. We can therefore reduce its security to the security of the existing BPF verifier. Importantly, this allows BERRYFY to benefit from formal verification efforts that improve the security of the verifier [64, 82], even if they do not explicitly take Spectre into account. In summary, BERRYFY uses a reliable and easy-to-implement

method to prevent unsafe speculative behavior (and thereby the use of secret-leaking side- and covert channels).

Our extension does not change the speculative security properties enforced by the BPF verifier, because we do not permit any new speculative behavior. With the assumptions noted in Section 3, our modified verifier is therefore still *secure*, assuming the CT leakage model, an unprivileged local attacker, and the secrecy policy identified in Section 2.2.2.

While BERRIFY does not affect security against unprivileged attackers, it improves security in the light of privileged users that only manually verify architectural safety. With BERRIFY, these users no longer accidentally load BPF programs into the kernel that unexpectedly leak sensitive kernel or user data. For example, an admin unfamiliar with transient execution vulnerabilities might have loaded the code from Figure 3 into their kernel. This effectively places an arbitrary-read gadget in the kernel that can be used by anyone (even unprivileged or remote users [83]) who controls the invocation of the BPF program, its input parameter `scalar`, and the covert channel used (`covert_channel[value]`). With BERRIFY enabled for privileged users (which becomes practical because of the improved expressiveness), no arbitrary data will be read into `value` and therefore the kernel remains secure.

In summary, BERRIFY soundly prevents unsafe transient execution and directly benefits from work that discovers architectural verifier bugs. The main expected downside of BERRIFY’s approach, in comparison to a vulnerable system, is the execution-time overhead at runtime due to the speculation barriers and the added verification-time overhead. Both are the topic of the following section and are further analyzed in the evaluation.

6.3 Performance

In this section, we analyze BERRIFY’s impact on performance. In particular, BERRIFY has an impact on the execution time of the BPF program and on the verification time when loading a BPF program.

BERRIFY inserts speculation barriers into the BPF program, which reduce the instruction-level parallelism inside the BPF program. Adding to the speculation barriers introduced as Spectre-STL defenses, these insertions could further add to *Problem #2*. However, the user and kernel code that calls the BPF programs, as well as the kernel helper functions invoked by the BPF programs, are unaffected by this change and therefore still execute at maximum performance. As most applications only spend a small part of the CPU time executing BPF code, we expect the real-world overhead to be small, which is also supported by our evaluation. In any case, BERRIFY enables unprivileged users to use BPF at all, as their programs were rejected prior to BERRIFY. This saves them from having to implement their logic in user space, which would then require very expensive user/kernel switches [1]. We therefore deem BERRIFY’s negative impact on *Problem #2* the more favorable trade-off.

Further, the number of barriers inserted by BERRIFY is reduced in comparison to more naive approaches, as barriers are only inserted when the verifier actually detects unsafe behavior that could enable a transient execution attack on the kernel. In addition, verification of a speculative path is cut short when a barrier is already present. This not only reduces the verification time but also helps to keep the number of barriers inserted small. At the same time, this approach

avoids complex compile-time analysis, which would not be practical for BPF. Our evaluation supports that this approach is precise even though it simple – we find that BERRIFY inserts a lot fewer barriers than the Spectre-STL defenses already present in Linux v6.5.

Aside from the execution-time overheads, BERRIFY only impacts the verification time (in comparison to a vulnerable system) for the programs that would otherwise have been rejected. For these programs, BERRIFY continues to explore the remaining architectural and speculative paths after discovering unsafe behavior. We do not make any assumptions about the size of the CPU’s speculation window, following the reasoning from [43]. However, if the hardware were to expose this information reliably, BERRIFY could use it to significantly reduce the verification time because the verifier indeed knows the exact microarchitecture on which the code will run. To limit verification time in our prototype, we selectively cut explored speculative paths short by inserting a barrier prematurely when we approach the verifier’s configurable complexity limits (e.g., number of instructions simulated, number of branches followed⁴). Using this simple heuristic, we can even verify the largest BPF application programs accepted by the Linux v6.5 verifier without defenses. In summary, there are multiple approaches that would allow us to reduce the number of barriers even further without impacting verification complexity. Because verification is usually not part of the user application’s hot path [84] we focus on the execution-time overhead in our evaluation.

In summary, BERRIFY improves the performance of unprivileged user applications by allowing them to use BPF at all. At the same time, our approach has no impact on privileged users, who can still dynamically disable BERRIFY per-program at execution time if they choose to manually check for gadgets.

7 IMPLEMENTATION

We implement BERRIFY for the Linux [72] BPF verifier and publish our patches⁵ consisting of less than 1000 source line changes under an open-source license. The majority of changes merely restructure the existing BPF verifier to support our design. While we use Linux v6.5, our approach is not fixed to only apply to this specific Linux version. Instead, it even benefits from future improvements to the precision of the verifier’s architectural analysis (which also reduce false-positives for BERRIFY) and JIT compiler. From a practical perspective, we are not aware of any fundamental reasons for which the patches could not be merged into upstream Linux.

To make our approach portable, we introduce a distinction between speculation barriers (verifier-internal BPF bytecode instructions) against Spectre-STL, and BPF speculation barriers against Spectre-PHT into the kernel (`nospec_v4` and `nospec_v1` respectively). The verifier inserts both barriers into the bytecode, and the JIT compiler backends then either drop or lower the instructions based on the architecture’s configuration. For example, ARM64 does not require speculation barriers to defend against Spectre-STL due to its firmware defenses [70] while x86-64 does require 1 fence instructions for both Spectre-STL (unless Speculative Store Bypass Disable, SSB, is active) and Spectre-PHT. This approach is also

⁴For our benchmarks, we increase the existing limits by a factor of 4 to successfully verify all real-world BPF programs (designed for verification without Spectre) with full defenses. For Parca, we increase the limit by a factor of 32.

⁵<https://gitlab.cs.fau.de/un6esoq/linux/-/tree/bpf-spectre>

compatible with research that proposes address-specific speculation barriers like `protect` from [55]. In summary, vulnerability- and even address-specific speculation barriers allow the verifier to remain architecture-agnostic while not impacting performance on architectures that are not affected.

8 EVALUATION

We evaluate BERRIFY using both static analysis and real-world execution time benchmarks. Regarding the former, we analyze whether BERRIFY successfully protects the programs rejected by the upstream Linux v6.5 verifier and how many speculation barriers it inserts. Regarding the latter, we analyze the performance impact of the speculation barriers on three popular performance-critical applications of BPF.

8.1 Static Analysis

First, we apply BERRIFY to all BPF programs from our motivating analysis of the BPF program landscape from Section 5. BERRIFY successfully applies defenses to all real-world application programs in our dataset and the number of barriers it inserts is insignificant in comparison to the upstream Spectre-STL defense.

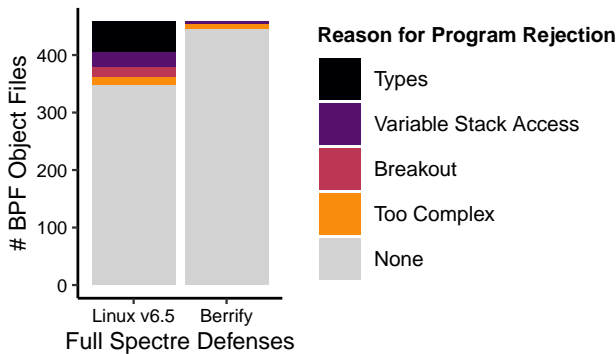


Figure 5: BERRIFY successfully applies to all real-world application programs without modifications. Analyzing all 459 object files, defenses only fail to apply to 3 % which is a significant improvement over the upstream verifier’s 24 %. BERRIFY only rejects 15 test programs from the Linux selftests which exhibit unverifiable architectural behavior that could easily be avoided in real applications.

8.1.1 Program Rejections. Figure 5 compares the number and type of verification errors with BERRIFY to the upstream Linux defenses (in both cases the Spectre-STL defenses are also active). As expected, BERRIFY is able to successfully protect all BPF programs from our *application* group and almost all programs from the full dataset. The remaining programs are all from the Linux kernel selftests which BERRIFY rejects for one of two reasons:

- **Variable Stack-Access (5/459):** The verifier currently does not implement defenses against Spectre-PHT in the light of variable stack accesses to simplify the implementation [33], our modified verifier therefore still rejects them. We deem an extension to resolve this out of scope as one can easily

avoid variable stack accesses in real-world BPF programs. Our evaluation confirms that rejection due to variable stack accesses is not a problem for any of the application’s BPF programs.

- **Too Complex (10/459):** The Linux kernel selftests contain object files with very large programs to test the verifier’s complexity limits. Because BERRIFY has to verify additional speculative program paths, verification can fail if the existing program was already designed to be as large as possible. If users were to encounter this in the real-world, they could easily circumvent it by splitting their program into multiple smaller programs that call each other using BPF tail calls.

In summary, **BERRIFY successfully verifies all BPF programs from real applications, thereby solving Problem # 3.** Further, the remaining theoretical reasons for failed verification are easy to understand and circumvent by BPF program developers.

8.1.2 Speculation Barriers. While it is BERRIFY’s desired outcome to enable more programs to be successfully loaded into the kernel, it also has the potential downside of slowing down those programs because of the speculation barriers it inserts for defense, thus, adding to *Problem # 2*. In this section and the following real-world performance evaluation, we analyze the extent to which these barriers affect performance.

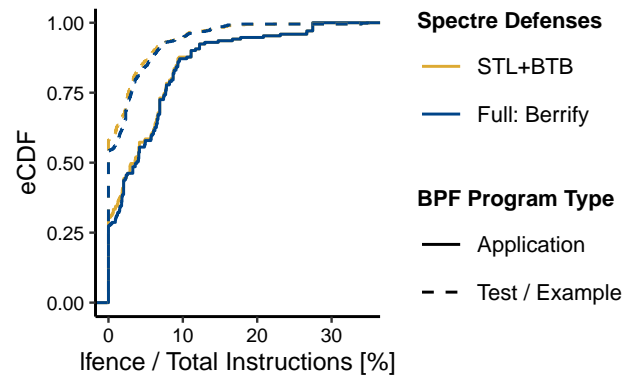


Figure 6: BERRIFY inserts only insignificantly more speculation barriers than the Spectre-STL defenses present in the upstream Linux kernel.

First, we repeat the analysis from Section 5.2 for BERRIFY, analyzing the number of speculation barriers it inserts in comparison to the Spectre-STL defenses. Figure 6 shows the fraction of BPF programs with less than X percent of speculation barriers for BERRIFY and Spectre-STL defenses (excluding the 15 rejected programs from the previous section). BERRIFY only inserts barriers when the program would otherwise have been rejected by the verifier, therefore the overall number of barriers it inserts is small.

In summary, we find that BERRIFY has a low expected performance overhead for a diverse set of applications. In the following section, we will back up this claim using performance benchmarks for three real-world applications using BPF.

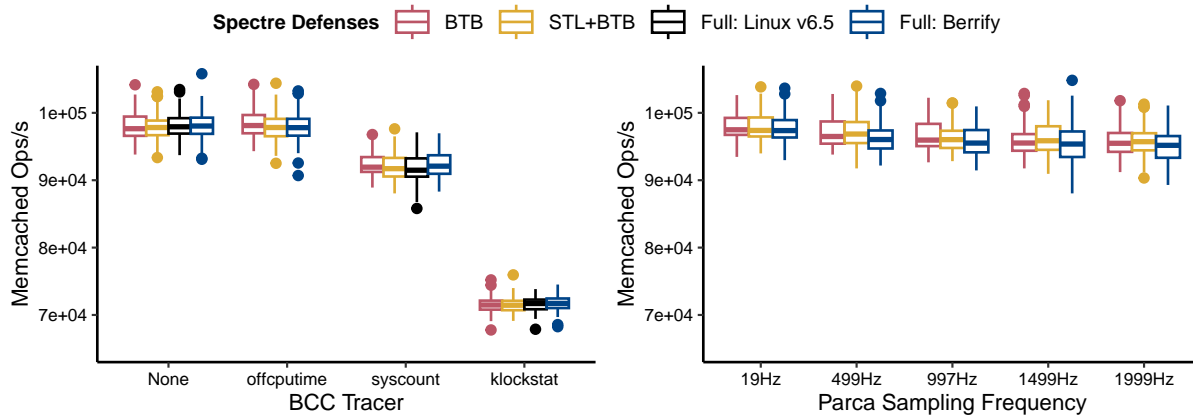


Figure 7: BERRIFY does not affect Memcached performance when using any of the 43 analyzed BCC tracers (of which the three most CPU-intensive are displayed in the left plot) or the Parca Continuous Profiler (right plot) to monitor the workload.

8.2 Real-World Performance Evaluation

In this section, we analyze the real-world application overhead of BERRIFY for over 50 distinct BPF programs. We focus on event tracing, continuous profiling, and packet processing, as these are three of the most popular performance-critical applications of BPF. As more complex programs are more likely to exhibit unsafe transient behavior, we select the Parca Continuous Profiler [77] and the Loxilb Network Load Balancer [78] to analyze BERRIFY’s overhead in extreme cases.

We run all our execution-time benchmarks in a Debian 11 GNU/Linux system with a v6.5.11 kernel on a 6-core 2.8 GHz Intel CPU from 2017 (Intel Core i5-8400). We disable dynamic voltage and frequency scaling (DVFS) to ensure reproducibility. In our graphs, we use standard boxplots [85] showing the first, second (i.e., median), and third quartile.

8.2.1 Event Tracing. Tracing is one of the most popular applications of BPF useable both for performance debugging and continuous monitoring in production. It allows users to gain valuable insights into kernel execution with little to no impact on production workloads. Still, their use poses security risks as, without Spectre defenses, BPF programs can easily introduce arbitrary-read gadgets into the kernel by mistake, therefore jeopardizing the security of the whole system. In this section, we analyze the execution-time overhead BERRIFY has for BCC’s libbpf-based tracers when monitoring a system that runs Memcached [86] together with the memtier_benchmark [87] load generator. Client and server each use three threads and communicate using Memcached’s binary protocol. We perform 15 000 requests taking 18 s on average and repeat each test 100 times.

We analyze 43 BCC tracers⁶. Out of these tracers, 21 require BERRIFY to be successfully used with Spectre-PHT defenses. When running the tracers alongside the workload, we find that 22 of the tracers record at least 10 events per second and 9 at least 1000. We conclude that most of the tracers are invoked frequently and collect

⁶This includes the 39 BPF objects from Section 5 but the number is higher here because some BPF objects are not designed to be loaded using bpf tool.

a reasonable amount of information in this setup. However, the CPU time the tracers spend execution the BPF programs themselves is small. Out of the analyzed tracers, 9 tracers spend more than 0.1 %, and 3 tracers spend more than 1 % of CPU time executing BPF code.

The left plot of Figure 7 shows the impact that the three most CPU-intensive tracers have on Memcached’s performance in four different system configurations. Without BERRIFY, the verifier cannot successfully apply Spectre defenses to offputime’s BPF program, thus the baseline system configuration (*Full: Linux v6.5*) is missing for this tracer. Defenses do successfully apply to syscount and klockstat without BERRIFY, but we still measure the overhead with BERRIFY for completeness. Overall, BERRIFY does not have any measurable impact. We also analyze the CPU time spent executing BPF code and find no measurable difference between BERRIFY and the other configurations even when we run each test for a total duration of 30 min. This is expected as tracers are designed to not impact production workloads by consuming as few resources as possible. They can therefore benefit from the security benefits BERRIFY offers without experiencing increased overheads.

8.2.2 Continuous Profiling. The second real-world BPF application we analyze is continuous profiling. Here, a BPF program that records the current user and kernel stack trace is invoked by a timer interrupt at a configurable frequency. The data can later be analyzed to find performance bugs, for example using flame graphs [88]. Traces are either collected continuously in production (sampling frequencies below 150 Hz) or on-demand by developers (usually 500 Hz to 2000 Hz).

For this benchmark, we run the Parca Continuous Profiler alongside the Memcached workload from the previous section. The results are displayed in Figure 7 on the right. As the BPF program that Parca uses to collect the stack samples is relatively complex, it cannot be successfully defended against Spectre-PHT without BERRIFY. Therefore, the baseline (*Full: Linux v6.5*) is missing from the figure. Analyzing the amount of CPU time spent in BPF code for this benchmark, we find that it is only 1.0 % even when the maximum sampling frequency we deem reasonable is used. The main overhead of the Parca tracer therefore originates from the

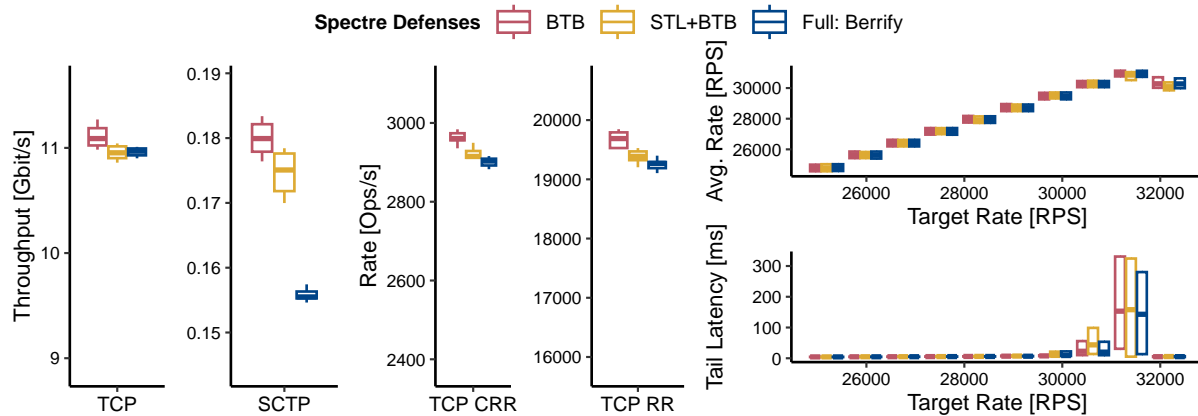


Figure 8: Transmission Control Protocol (TCP) and Stream Control Transmission Protocol (SCTP) throughput, TCP Request/Response- (RR) and Connect/Request/Response (CRR) rate, and nginx HTTP tail latency achievable over a BPF-based Loxilb load balancer. The upstream Linux v6.5 verifier cannot successfully apply Spectre defenses to Loxilb’s BPF program, therefore the baseline (Full: Linux v6.5) is missing.

code that post-processes the samples in user space. From this, it follows that both the Spectre-STL defenses and the Spectre-PHT defenses using BERRIFY do not impact Memcached’s performance even though they increase the execution time of the BPF programs by 16 % and 62 % respectively. In comparison to a vulnerable system, BERRIFY only reduces Memcached’s throughput by 0.8 % at 1999 Hz. At frequencies below 150 Hz, the overhead is no longer measurable (0.1 % for 19 Hz in this particular run). In summary, BERRIFY can be applied to Parca’s BPF program without overheads to the workload.

8.2.3 Network Load Balancing. High-performance packet processing is the application BPF was originally developed for and is still one of the most popular use cases [7, 10, 78, 89, 90, 91, 92].

In this benchmark, we replicate a real-world scenario where multiple Docker containers communicate with each other through the BPF-based Loxilb network load balancer [78]. To stress the load balancer, we extend the benchmarks the Loxilb upstream project includes. To achieve the maximum possible CPU utilization, we do not run the containers on different physical machines. Therefore, the overhead we measure represents the upper bound as the performance is not limited by the speed of the networking interface but only by the CPU’s speed.

The four left plots of Figure 8 show the maximum throughput and request rate for Transmission Control Protocol (TCP) and Stream Control Transmission Protocol (SCTP). We normalize each scale to display between 80 % and 105 % of the respective performance achievable with only the Spectre-BTB defenses active (the default in Linux v6.5). For netperf [93] and iperf3 [94], we use one client thread but confirm that multiple threads do not impact our conclusions in separate experiments (e.g., using iperf v2.0.14a). We run each benchmark for 2 min and repeat the measurement 10 times.

We observe that BERRIFY only has a small impact on TCP throughput and average latency. Compared to a vulnerable system, both BERRIFY and STL+BTB defenses reduce the TCP throughput by 1.1 % to 1.2 %. As expected, the Request/Response (RR) rate is affected the most by BERRIFY with 2.1 % overhead while the Connect/Request/

Response (CRR) rate only decreases by 1.8 %. In summary, the main overhead stems from the Spectre-STL defenses (Problem # 2), while BERRIFY only has a negligible impact.

SCTP is a message-based protocol that still ensures reliable transportation. In comparison to TCP and User Datagram Protocol (UDP), it offers native support for multihoming. However, because it is not as popular, the Linux kernel does not offer a BPF helper function to recompute the checksum for packets that are redirected. Therefore, Loxilb has to recompute the SCTP checksum in BPF itself. Packet redirection thus consumes much more CPU time and the throughput is reduced by over 60× in comparison to TCP. This also causes increased relative overheads when we activate the Spectre-STL defenses (2.7 %) and BERRIFY (14 %). To improve SCTP’s performance, it would be the most promising to either disable checksum calculation [95], or create a kernel helper function. Compared to the potential speedup by avoiding checksum recalculation in BPF altogether, the speedup achieved by disabling BERRIFY is insignificant for SCTP.

In summary, BERRIFY does not bottleneck the TCP and SCTP performance achievable over the load balancer. This is the case especially because the overhead would be even lower if the containers were not colocated on the same machine.

For our final benchmark, we run two nginx servers serving a 1 KiB payload and connect to them through a single Loxilb load balancer instance. We use the wrk2 load generator [96] with two threads to target a specific HTTP request rate and analyze the resulting tail latency (99th percentile) and average rate. We use wrk2 because it does not suffer from coordinated omissions. With two servers, the maximum achievable rate is approx. 31 kRPS. Based on this, we scale the target load from 80 % to 105 % in increments of 0.1 % and then group the data into bins of 2.5 %. We run each test for 1 min and repeat the measurement 10 times, therefore each bin contains 250 data points. We omit outliers from this plot.

The right plots of Figure 8 show the achieved average rate and tail latency for different target rates. As expected, BERRIFY does

not cause increased tail latencies. This is because BERRIFY merely increases the time needed to process each request by a small amount, but does not cause any unpredictable spikes in the processing time.

8.2.4 Summary. BERRIFY does not affect most of the applications we analyze in a measurable way. Event tracers appear unaffected, which is particularly important as developers often write small tracing programs ad-hoc [62] and therefore are unlikely to exhaustively review their programs for gadgets. For more complex BPF programs as used by the Parca Continuous Profiler and the Loxilb Network Load Balancer, we find that both the Spectre-STL defenses and BERRIFY can affect the CPU time spent executing BPF code, however, **this usually does not impact application performance, showing that Problem # 2 should only be of secondary concern.** Instead, it is much more important to be able to use BPF in the first place (i.e., Problem # 3, which BERRIFY solves). In any case, users can always disable BERRIFY per-BPF-program at runtime to trade security for performance.

9 RELATED WORK

In this section, we discuss alternatives to BPF and then relevant alternative techniques to defend against transient execution attacks.

9.1 High-Performance IO

BPF implements safe kernel extensions [97, 98, 99] for Linux. While it has numerous applications outside of high-performance IO [11, 12, 13, 14, 15], asynchronous IO and kernel-bypass can replace it in some cases.

9.1.1 Asynchronous IO. To improve performance over traditional system calls, asynchronous IO [100] is implemented by aio [101] and more recently by io_uring [6] in Linux [102]. Being closely related to system-call batching [1, 103], it amortizes mode switches over multiple operations. However, both do not improve upon the latency of traditional system-calls because data still has to pass through the OS network or storage stack to reach the application. In contrast, BPF allows applications to process (or discard) data directly on the CPU where it is first retrieved [8].

9.1.2 Kernel-Bypass. The most prominent implementations of kernel-bypass are DPDK [2, 104] and SPDK [3, 105]. By giving user applications direct access to the hardware, kernel-bypass can reduce both IO latency and throughput. However, in contrast to BPF, kernel-bypass does not integrate cooperatively with the existing networking stack and requires dedicating full CPU cores to busy-looping for low-latency packet processing (hurting power-proportionality). Further, when kernel-bypass and regular applications are colocated on a server, kernel-bypass reduces performance for regular applications because packets have to be re-injected into the kernel networking stack to reach them [7]. For this reason, Meta reportedly uses a BPF-based load balancer (similar to Loxilb from our evaluation) instead of kernel-bypass [90].

9.2 Transient-Execution-Attack Defenses

To defend against transient execution attacks, BERRIFY relies on compiler-based defenses. Alternative approaches partition resources (OS-based) or attempt to implement side-channel-resistant transient execution that is still low-overhead (hardware-based).

9.2.1 Compiler-based Defenses. To the best of our knowledge, Linux’s BPF verifier is the only widely deployed sandbox that is fully hardened against the known Spectre vulnerabilities without relying on process isolation. All production Java [25] and Wasm [27, 28] runtimes rely on process isolation, as recommended by Intel [18]. However, the kernel cannot apply process isolation to BPF without reducing its performance significantly [1]. The Windows BPF verifier (Prevail [106]) does not implement Spectre defenses as of April 2024 [107].

There exist theoretical works that can potentially detect Spectre gadgets more precisely than the BPF verifier (even with our extension) [43, 57]. However, these cannot be readily applied to BPF because they make simplifying assumptions in their implementation [49, 108] or are too complex to implement [55, 109] (and would, therefore, further increase the risk of security-critical bugs). The common C/C++ toolchains do not offer reliable, architecture-agnostic Spectre defenses that are transparent to the users [18, 110, 111, 112, 113, 114, 115], therefore the BPF verifier must apply defenses to the programs. Most related works on BPF focus on architectural security but do not consider transient execution attacks [116, 117, 118, 119]. In future work, the number of barriers BERRIFY inserts could be further reduced by using ALU instructions to detect misprediction and erase sensitive data (e.g., using Speculative Load Hardening, SLH [57, 120, 121, 122, 123] and similar techniques [80, 124, 125]).

While browsers to date rely on process isolation to prevent Wasm programs from exploiting Spectre [28, 125], there is some work on compiler-based defenses in specific runtimes. However, we find that they are either incomplete (i.e., Wasmtime [27], V8 [19, 29]) or specific to Wasm/user space (i.e., Swivel [126]) and do not apply to BPF in the kernel.

9.2.2 OS-based Defenses. To defend against transient execution attacks, most works rely on coarse-grained partitioning of resources [127, 128] with support from the OS (e.g., address-space isolation [129] and core scheduling [130]). However, address-space isolation cannot be applied to BPF without limiting its performance even further [1]. Applying Memory Protection Keys (MPKs) to BPF against Spectre appears to be a promising direction for future research [4, 18]. Existing work on the topic does not take transient execution attacks into account [119] but assumes the BPF verifier (e.g., our work) implements defenses.

9.2.3 Hardware-based Defenses. As of April 2024, there exists no practical high-performance processor implementation that is not vulnerable to Spectre [22, 23, 24]. Instead, vendors continue to recommend compiler-based defenses against all cross- and in-domain transient execution attacks [4, 131]. To date, complete protection from microarchitectural timing side-channels is not possible without significant changes to the instruction set architecture (ISA) and processor design [20, 21, 132, 133].

10 CONCLUSION

This work has presented BERRIFY for Linux BPF, a practical and easy-to-apply enhancement to the only software-based sandbox resilient against Spectre. BERRIFY is *sound*, *precise*, and *lightweight*. First, it *soundly* combines speculation barriers and static analysis

in a way that does not increase the kernel’s attack surface. Second, it *precisely* prevents unsafe transient behavior thereby reducing the number of rejected programs from 24 % to 3 % (with only test programs from the Linux selftests remaining). Finally, as our real-world performance evaluation demonstrates, it is *lightweight* because it does not impact BPF invocation latency, which is of particular importance to BPF as it complements user space by offering minimum-overhead transitions.

REFERENCES

- [1] Luis Gerhorst et al. “AnyCall: Fast and Flexible System-Call Aggregation”. In: *Proceedings of the 11th Workshop on Programming Languages and Operating Systems (PLOS’21)*. ACM, 2021, pp. 1–8. doi: [10.1145/3477113.3487267](https://doi.org/10.1145/3477113.3487267).
- [2] Thomas Monjalon et al. *DPDK: Data Plane Development Kit – v24.03*. 2024. URL: <http://git.dpdk.org/dpdk/tree/?h=v24.03&id=a9778aad> (visited on 04/12/2024).
- [3] Ziye Yang et al. “SPDK: A Development Kit to Build High Performance Storage Applications”. In: *Proceedings of the 9th International Conference on Cloud Computing Technology and Science (CloudCom’17)*. IEEE, 2017, pp. 154–161. doi: [10.1109/CloudCom.2017.14](https://doi.org/10.1109/CloudCom.2017.14).
- [4] *Refined Speculative Execution Terminology*. Mar. 11, 2022. URL: <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/best-practices/refined-speculative-execution-terminology.html> (visited on 04/04/2024).
- [5] Moritz Lipp et al. “Meltdown: Reading Kernel Memory from User Space”. In: *Proceedings of the 27th USENIX Security Symposium (USENIX Security ’18)*. USENIX, 2018, pp. 973–990. doi: [10.1145/3357033](https://doi.org/10.1145/3357033).
- [6] Jens Axboe. *[PATCHSET v5] io_uring IO interface (Mail)*. Jan. 16, 2019. URL: <https://lore.kernel.org/linux-block/20190116175003.17880-1-axboe@kernel.dk/> (visited on 04/08/2024).
- [7] Toke Høiland-Jørgensen 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 (CoNEXT ’18)*. ACM, 2018, pp. 54–66. doi: [10.1145/3281411.3281443](https://doi.org/10.1145/3281411.3281443).
- [8] Yuhong Zhong et al. “XRP: In-Kernel Storage Functions with eBPF”. In: *Proceedings of the 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI ’22)*. Carlsbad, CA: USENIX Association, July 2022, pp. 375–393. ISBN: 978-1-939133-28-1. URL: <https://www.usenix.org/conference/osdi22/presentation/zhong>.
- [9] Tomasz Osiński. *p4c-ubpf: a New Back-end for the P4 Compiler*. June 1, 2020. URL: <https://opennetworking.org/news-and-events/blog/p4c-ubpf-a-new-back-end-for-the-p4-compiler/> (visited on 02/17/2023).
- [10] Marcos A. M. Vieira et al. “Fast Packet Processing with eBPF and XDP: Concepts, Code, Challenges, and Applications”. In: *ACM Computing Surveys* 53.1 (2020), 16:1–16:36. doi: [10.1145/3371038](https://doi.org/10.1145/3371038).
- [11] Zhilu Lian et al. *eBPF-based Working Set Size Estimation in Memory Management*. 2022. doi: [10.1109/ICSS55994.2022.00036](https://doi.org/10.1109/ICSS55994.2022.00036). arXiv: [2303.05919](https://arxiv.org/abs/2303.05919).
- [12] Daniel Kim and Robert Prast. “Triaging Real-time Security Threats with eBPF-powered Observability”. SREcon22 Americas. 2022. URL: <https://www.usenix.org/conference/srecon22americas/presentation/kim> (visited on 04/12/2024).
- [13] Jinghao Jia et al. *Programmable System Call Security with eBPF*. 2023. doi: [10.48550/arXiv.2302.10366](https://doi.org/10.48550/arXiv.2302.10366). arXiv: [2302.10366](https://arxiv.org/abs/2302.10366).
- [14] Junqing Yang, Lijun Chen, and Jiaqing Bai. “Redis automatic performance tuning based on eBPF”. In: *Proceedings of the 14th International Conference on Measuring Technology and Mechatronics Automation (ICMTMA’22)*. IEEE, 2022, pp. 671–676. doi: [10.1109/ICMTMA54903.2022.00139](https://doi.org/10.1109/ICMTMA54903.2022.00139).
- [15] Tejun Heo. *[PATCHSET v3] sched: Implement BPF extensible scheduler class*. 2023. URL: <https://lore.kernel.org/all/20230317213333.2174969-1-tj@kernel.org/> (visited on 03/18/2023).
- [16] Ioannis Zarkadas et al. *BPF-oF: Storage Function Pushdown Over the Network*. 2023. doi: [10.48550/arXiv.2312.06808](https://doi.org/10.48550/arXiv.2312.06808). arXiv: [2312.06808](https://arxiv.org/abs/2312.06808).
- [17] Paul Kocher et al. “Spectre Attacks: Exploiting Speculative Execution”. In: *Proceedings of the 40th Symposium on Security and Privacy (SP’19)*. IEEE, 2019, pp. 1–19. doi: [10.1109/SP.2019.00002](https://doi.org/10.1109/SP.2019.00002).
- [18] Intel. *Managed Runtime Speculative Execution Side Channel Mitigations*. Jan. 3, 2018. URL: <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/technical-documentation/runtime-speculative-side-channel-mitigations.html> (visited on 04/04/2024).
- [19] Ross McIlroy et al. *Spectre is here to stay: An analysis of side-channels and speculative execution*. 2019. doi: [10.48550/arXiv.1902.05178](https://doi.org/10.48550/arXiv.1902.05178). arXiv: [1902.05178](https://arxiv.org/abs/1902.05178).
- [20] Wenjie Xiong and Jakub Szefer. “Survey of Transient Execution Attacks and Their Mitigations”. In: *ACM Comput. Surv.* 54.3 (May 2021). ISSN: 0360-0300. doi: [10.1145/3442479](https://doi.org/10.1145/3442479).
- [21] Jiyong Yu et al. “Speculative Taint Tracking (STT): A Comprehensive Protection for Speculatively Accessed Data”. In: *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*. MICRO ’52. Columbus, OH, USA: Association for Computing Machinery, 2019, pp. 954–968. ISBN: 9781450369381. doi: [10.1145/3352460.3358274](https://doi.org/10.1145/3352460.3358274).
- [22] *Software Techniques for Managing Speculation on AMD Processors - Revision 5.09.23*. Sept. 5, 2023. URL: <https://web.archive.org/web/20240313145139/https://www.amd.com/content/dam/amd/en/documents/epyc-technical-docs/tuning-guides/software-techniques-for-managing-speculation.pdf> (visited on 04/11/2024).
- [23] *Affected Processors: Guidance for Security Issues on Intel® Processors*. Apr. 11, 2024. URL: <https://web.archive.org/web/20240411104155/https://www.intel.com/content/www/us/en/developer/topic-technology/software-security-guidance/processors-affected-consolidated-product-cpu-model.html> (visited on 04/11/2024).

- [24] *Speculative Processor Vulnerability (Updated December 20, 2023) - Arm Security Center*. Dec. 20, 2023. URL: <https://web.archive.org/web/20240324073901/https://developer.arm.com/Arm%20Security%20Center/Speculative%20Processor%20Vulnerability> (visited on 04/11/2024).
- [25] Amir Naseredini et al. *Systematic Analysis of Programming Languages and Their Execution Environments for Spectre Attacks*. 2021. DOI: 10.48550/arXiv.2111.12528. arXiv: 2111.12528.
- [26] Jules Dejaeghere et al. “Comparing Security in eBPF and WebAssembly”. In: *Proceedings of the 1st Workshop on eBPF and Kernel Extensions (eBPF’23)*. ACM, 2023, pp. 35–41. DOI: 10.1145/3609021.3609306.
- [27] Wasmtime. *Wasmtime Docs – Security*. Feb. 15, 2024. URL: <https://docs.wasmtime.dev/security.html#spectre> (visited on 02/15/2024).
- [28] Charles Reis, Alexander Moshchuk, and Nasko Oskov. “Site Isolation: Process Separation for Web Sites within the Browser”. In: *Proceedings of the 28th USENIX Security Symposium (USENIX Security ’19)*. USENIX, 2019, pp. 1661–1678. URL: <https://www.usenix.org/conference/usenixsecurity19/presentation/reis>.
- [29] Mathias Bynens. *V8 Docs – Untrusted code mitigations*. 2018. URL: <https://v8.dev/docs/untrusted-code-mitigations#sandbox-untrusted-execution-in-a-separate-process> (visited on 04/10/2024).
- [30] *eBPF verifier – The Linux Kernel documentation (v6.5)*. We use BPF to refer to the current extended Berkeley Packet Filter facility. Nov. 8, 2023. URL: <https://www.kernel.org/doc/html/v6.5/bpf/verifier.html> (visited on 03/18/2024).
- [31] Alexei Starovoitov, John Fastabend, and Daniel Borkmann. *Linux Kernel Source Tree – bpf: prevent out-of-bounds speculation (commit #b2157399cc98)*. Jan. 7, 2018. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit?id=b2157399cc98> (visited on 04/02/2024).
- [32] Daniel Borkmann et al. *Linux Kernel Source Tree – bpf: Fix leakage of uninitialized bpf stack under speculation (commit #801c60)*. Apr. 29, 2021. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit?id=801c6058> (visited on 02/14/2024).
- [33] Andrei Matei and Alexei Starovoitov. *Linux Kernel Source Tree – bpf: Allow variable-offset stack access (commit #01f810)*. Feb. 6, 2021. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit?id=01f810ac> (visited on 03/18/2024).
- [34] Daniel Borkmann et al. *Linux Kernel Source Tree – bpf: Fix leakage due to insufficient speculative store bypass mitigation (commit #2039f2)*. July 13, 2021. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit?id=2039f26f> (visited on 02/14/2024).
- [35] Daniel Borkmann et al. *Linux Kernel Source Tree – bpf: Fix leakage under speculation on mispredicted branches (commit #918367)*. May 28, 2021. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit?id=9183671a> (visited on 02/14/2024).
- [36] Jason Kim et al. “iLeakage: Browser-based Timerless Speculative Execution Attacks on Apple Devices”. In: *Proceedings of the 30th SIGSAC Conference on Computer and Communications Security (CCS’23)*. ACM, 2023, pp. 2038–2052. DOI: 10.1145/3576915.3616611.
- [37] Luis Gerhorst. *Re: [PATCH 2/3] Revert "bpf: Fix issue in verifying allow_ptr_leaks"*. Sept. 28, 2023. URL: <https://lore.kernel.org/bpf/20230928110927.115238-1-gerhorst@amazon.de/> (visited on 02/12/2024).
- [38] Luis Gerhorst, Henriette Hofmeier, and Daniel Borkmann. *bpf: Fix pointer-leak due to insufficient speculative store bypass mitigation*. Jan. 9, 2023. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/bpf/bpf.git/commit/?id=e4f4db47> (visited on 01/22/2023).
- [39] Luis Gerhorst and Daniel Borkmann. *#082cdc - bpf: Remove misleading spec_v1 check on var-offset stack read - kernel/git/torvalds/linux.git - Linux kernel source tree*. Mar. 15, 2023. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=082cdc69> (visited on 02/08/2024).
- [40] Michael Schwarz et al. “ZombieLoad: Cross-Privilege-Boundary Data Sampling”. In: *Proceedings of the 26th ACM SIGSAC Conference on Computer and Communications Security (CCS ’19)*. ACM, 2019, pp. 753–768. DOI: 10.1145/3319535.3354252.
- [41] Jo Van Bulck et al. “Foreshadow: Extracting the Keys to the Intel SGX Kingdom with Transient Out-of-Order Execution”. In: *Proceedings of the 27th USENIX Security Symposium (USENIX Security ’18)*. USENIX, 2018, p. 991. URL: <https://www.usenix.org/conference/usenixsecurity18/presentation/bulck> (visited on 04/12/2024).
- [42] Claudio Canella et al. “A systematic evaluation of transient execution attacks and defenses”. In: *Proceedings of the 28th USENIX Security Symposium (USENIX Security ’19)*. USENIX, 2019, pp. 249–266. URL: <https://www.usenix.org/conference/usenixsecurity19/presentation/canella>.
- [43] Sunjay Cauligi et al. “SoK: Practical Foundations for Software Spectre Defenses”. In: *Proceedings of the 43rd Symposium on Security and Privacy (SP’22)*. IEEE, 2022, pp. 666–680. DOI: 10.1109/SP46214.2022.9833707.
- [44] Daniel Borkmann. *BPF and Spectre: Mitigating transient execution attacks*. eBPF Summit 2021. 2021. URL: <https://www.youtube.com/watch?v=6N30Yp5f9c4> (visited on 01/22/2023).
- [45] Piotr Krysiuk, Benedict Schlüter, and Daniel Borkmann. *BPF and Spectre: Mitigating transient execution attacks*. Keynote at PriSC’22. 2022. URL: <https://popl22.sigplan.org/details/prisc-2022-papers/11/BPF-and-Spectre-Mitigating-transient-execution-attacks> (visited on 01/25/2023).
- [46] Atri Bhattacharyya et al. “SMoTherSpectre: Exploiting Speculative Execution through Port Contention”. In: *Proceedings of the 26th SIGSAC Conference on Computer and Communications Security (CCS’19)*. ACM, 2019, pp. 785–800. DOI: 10.1145/3319535.3363194.
- [47] Ofek Kirzner and Adam Morrison. “An Analysis of Speculative Type Confusion Vulnerabilities in the Wild”. In: *Proceedings of the 30th USENIX Security Symposium (USENIX Security ’21)*. USENIX, 2021, pp. 2399–2416.

- [48] Sunjay Cauligi et al. “Constant-time foundations for the new spectre era”. In: *Proceedings of the 41st SIGPLAN Conference on Programming Language Design and Implementation (PLDI’20)*. ACM, 2020, pp. 913–926. DOI: [10.1145/3385412.3385970](https://doi.org/10.1145/3385412.3385970).
- [49] Marco Guarnieri et al. “Spectector: Principled Detection of Speculative Information Flows”. In: *Proceedings of the 41st Symposium on Security and Privacy (SP’20)*. IEEE, 2020, pp. 1–19. DOI: [10.1109/SP40000.2020.00011](https://doi.org/10.1109/SP40000.2020.00011).
- [50] Intel. *Data Operand Independent Timing Instructions*. Feb. 28, 2023. URL: <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/resources/data-operand-independent-timing-instructions.html> (visited on 02/05/2023).
- [51] Intel. *Data Operand Independent Timing ISA Guidance*. Feb. 13, 2023. URL: <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/best-practices/data-operand-independent-timing-isa-guidance.html> (visited on 02/05/2023).
- [52] Jonathan Corbet. *Constant-time instructions and processor optimizations [LWN.net]*. Feb. 3, 2023. URL: <https://lwn.net/Articles/921511/> (visited on 02/03/2023).
- [53] ARM. *ARM Documentation – DIT, Data Independent Timing – Arm Armv8-A Architecture Registers*. June 2021. URL: <https://developer.arm.com/documentation/ddi0595/2021-06/AArch64-Registers/DIT--Data-Independent-Timing> (visited on 02/05/2023).
- [54] Basavesh Ammanaghatta Shivakumar et al. “Typing High-Speed Cryptography against Spectre v1”. In: *Proceedings of the 43rd Symposium on Security and Privacy (SP’23)*. IEEE, 2023, pp. 1094–1111. DOI: [10.1109/SP46215.2023.10179418](https://doi.org/10.1109/SP46215.2023.10179418).
- [55] Marco Vassena et al. “Automatically Eliminating Speculative Leaks from Cryptographic Code with Blade”. In: *Proceedings of the ACM on Programming Languages (PACMPL)* 5 (POPL 2021), pp. 1–30. DOI: [10.1145/3434330](https://doi.org/10.1145/3434330).
- [56] Xaver Fabian, Marco Guarnieri, and Marco Patrignani. *Automatic Detection of Speculative Execution Combinations*. 2022. DOI: [10.48550/arXiv.2209.01179](https://doi.org/10.48550/arXiv.2209.01179). arXiv: [2209.01179](https://arxiv.org/abs/2209.01179).
- [57] Marco Patrignani and Marco Guarnieri. “Exorcising spectres with secure compilers”. In: *Proceedings of the 28th SIGSAC Conference on Computer and Communications Security (CCS’21)*. ACM, 2021, pp. 445–461. DOI: [10.1145/3460120.3484534](https://doi.org/10.1145/3460120.3484534).
- [58] Sunjay Cauligi et al. *A Turning Point for Verified Spectre Sandboxing*. 2022. DOI: [10.48550/arXiv.2208.01548](https://doi.org/10.48550/arXiv.2208.01548). arXiv: [2208.01548](https://arxiv.org/abs/2208.01548).
- [59] Benedict Schlüter. “Security Analysis of eBPF”. Bachelor’s Thesis. Ruhr-Universität Bochum (RUB), July 22, 2021.
- [60] Luis Gerhorst. *[PATCH 2/3] Revert “bpf: Fix issue in verifying allow_ptr_leaks”*. Sept. 13, 2023. URL: <https://lore.kernel.org/bpf/20230913122827.91591-1-gerhorst@amazon.de/> (visited on 02/07/2024).
- [61] Daniel Borkmann et al. *Linux Kernel Source Tree – bpf: Tighten speculative pointer arithmetic mask (commit #7fedb6)*. Mar. 24, 2021. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=7fedb63a> (visited on 02/08/2024).
- [62] Daniel Xu et al. *bpftrace – v0.20.3*. Sept. 27, 2022. URL: <https://github.com/bpftrace/bpftrace/tree/v0.20.3> (visited on 04/12/2024).
- [63] Alexei Starovoitov et al. *kernel/bpf/verifier.c – Line 644 – Linux Kernel Stable Tree – v6.5.11*. Nov. 8, 2023. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/stable/linux.git/tree/kernel/bpf/verifier.c?id=79944183#n644> (visited on 03/18/2024).
- [64] Harishankar Vishwanathan et al. “Sound, precise, and fast abstract interpretation with tristate numbers”. In: *Proceedings of the 20th International Symposium on Code Generation and Optimization (CGO’22)*. IEEE, 2022, pp. 254–265. DOI: [10.1109/CGO53902.2022.9741267](https://doi.org/10.1109/CGO53902.2022.9741267).
- [65] *Debian Manpages – bpf(2)*. Feb. 5, 2023. URL: <https://manpages.debian.org/bookworm/manpages-dev/bpf.2.en.html> (visited on 04/08/2024).
- [66] Luis Gerhorst, Henriette Hofmeier, and Daniel Borkmann. *Linux Kernel Source Tree – bpf: Fix pointer-leak due to insufficient speculative store bypass mitigation (commit #e4f4db47)*. Jan. 9, 2023. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/bpf/bpf.git/commit/?id=e4f4db47> (visited on 01/22/2023).
- [67] Alexei Starovoitov. *Re: [PATCH 2/3] Revert “bpf: Fix issue in verifying allow_ptr_leaks”*. Sept. 19, 2023. URL: <https://patchwork.kernel.org/project/linux-kselftest/patch/20230913122827.91591-1-gerhorst@amazon.de/#25516346> (visited on 02/07/2024).
- [68] Intel. *Speculative Store Bypass / CVE-2018-3639 / INTEL-SA-00115*. May 21, 2018. URL: <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/advisory-guidance/speculative-store-bypass.html> (visited on 03/02/2023).
- [69] Andrea Mambretti et al. “Speculator: a tool to analyze speculative execution attacks and mitigations”. In: *Proceedings of the 35th Annual Computer Security Applications Conference (ACSAC’19)*. ACM, 2019, pp. 747–761. DOI: [10.1145/3359789.3359837](https://doi.org/10.1145/3359789.3359837).
- [70] ARM. *Firmware interfaces for mitigating cache speculation vulnerabilities – System Software on Arm Specification (Version 1.3)*. May 21, 2018. URL: <https://developer.arm.com/cache-speculation-vulnerability-firmware-specification>.
- [71] Yafang Shao. *Re: [PATCH 2/3] Revert “bpf: Fix issue in verifying allow_ptr_leaks”*. Sept. 19, 2023. URL: <https://patchwork.kernel.org/project/linux-kselftest/patch/20230913122827.91591-1-gerhorst@amazon.de/#25515964> (visited on 02/08/2024).
- [72] Greg Kroah-Hartman et al. *Linux Kernel Stable Tree – v6.5.11. BPF Samples from samples/bpf and Kernel Selftests from tools/testing/selftests/bpf*. Nov. 8, 2023. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/stable/linux.git/tree/?id=79944183> (visited on 03/18/2024).
- [73] Andrii Nakryiko et al. *libbpf/libbpf-bootstrap (GitHub, commit #a7c0f7)*. Sept. 27, 2022. URL: <https://github.com/libbpf/libbpf-bootstrap/tree/a7c0f7e4> (visited on 04/09/2024).
- [74] Elazar Gershuni et al. “Simple and precise static analysis of untrusted Linux kernel extensions”. In: *Proceedings of the 40th SIGPLAN Conference on Programming Language Design*

- and Implementation (PLDI'19). ACM, 2019, pp. 1069–1084. DOI: [10.1145/3314221.3314590](https://doi.org/10.1145/3314221.3314590).
- [75] Yonghong Song et al. *BCC – v0.27.0*. We use the libbpf-based variants from the libbpf-tools folder. Apr. 3, 2023. URL: <https://github.com/iovisor/bcc/tree/v0.27.0> (visited on 04/12/2024).
- [76] Joe Stringer et al. *cilium/cilium (GitHub) – v1.12.5*. Dec. 2022. URL: <https://github.com/cilium/cilium/tree/v1.12.5> (visited on 04/09/2024).
- [77] Javier Honduvilla Coto. *parca-dev/parca-agent (GitHub) – v0.27.0*. Nov. 9, 2023. URL: <https://github.com/parca-dev/parca-agent> (visited on 04/09/2024).
- [78] PacketCrunch (GitHub User) et al. *loxilb: eBPF based cloud-native load-balancer (GitHub) – v0.9.0*. Mar. 21, 2024. URL: <https://github.com/loxilb-io/loxilb> (visited on 03/21/2024).
- [79] Linux Torvalds et al. *bpftool – v5.18 (Linux kernel release)*. We use an older version because the latest version does not support legacy map definitions that are still used by many projects. This does not impact the in-kernel verification and JIT compilation (and thereby our results). May 22, 2022. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/stable/linux.git/tree/tools/bpf/bpftool?h=v5.18> (visited on 04/12/2024).
- [80] Daniel Borkmann, Jann Horn, and Alexei Starovoitov. *Linux Kernel Source Tree – bpf: prevent out of bounds speculation on pointer arithmetic (commit #979d63d5)*. Jan. 3, 2019. URL: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=979d63d5> (visited on 02/12/2024).
- [81] Vedvyas Shanbhogue, Deepak Gupta, and Ravi Sahita. “Security Analysis of Processor Instruction Set Architecture for Enforcing Control-Flow Integrity”. In: *Proceedings of the 8th International Workshop on Hardware and Architectural Support for Security and Privacy*. HASP '19. Phoenix, AZ, USA: Association for Computing Machinery, 2019. ISBN: 9781450372268. DOI: [10.1145/3337167.3337175](https://doi.org/10.1145/3337167.3337175).
- [82] Luke Nelson et al. “Specification and verification in the field: Applying formal methods to BPF just-in-time compilers in the Linux kernel”. In: *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association, Nov. 2020, pp. 41–61. ISBN: 978-1-939133-19-9. URL: <https://www.usenix.org/conference/osdi20/presentation/nelson>.
- [83] Michael Schwarz et al. *NetSpectre: Read arbitrary memory over network*. eprint: 1807.10535. 2018. DOI: [10.48550/arXiv.1807.10535](https://doi.org/10.48550/arXiv.1807.10535). arXiv: 1807.10535.
- [84] Milo Craun, Adam Oswald, and Dan Williams. “Enabling eBPF on Embedded Systems Through Decoupled Verification”. In: *Proceedings of the 1st Workshop on eBPF and Kernel Extensions (eBPF'23)*. ACM, 2023, pp. 63–69. DOI: [10.1145/3609021.3609299](https://doi.org/10.1145/3609021.3609299).
- [85] Hadley Wickham et al. *ggplot2 3.5.0: A box and whiskers plot (in the style of Tukey) – geom_boxplot*. Feb. 23, 2024. URL: https://ggplot2.tidyverse.org/reference/geom_boxplot.html (visited on 04/08/2024).
- [86] Chris Lamb et al. *memcached – v1.6.9*. Nov. 25, 2020. URL: <https://packages.debian.org/bullseye/web/memcached> (visited on 04/09/2024).
- [87] Yossi Gottlieb et al. *memtier_benchmark – v1.4.0 – NoSQL Redis and Memcache traffic generation and benchmarking tool*. Aug. 8, 2022. URL: https://github.com/RedisLabs/memtier_benchmark/tree/29f51a82 (visited on 04/12/2024).
- [88] Brendan Gregg. “The flame graph”. In: *Communications of the ACM* 59.6 (May 23, 2016), pp. 48–57. DOI: [10.1145/2909476](https://doi.org/10.1145/2909476).
- [89] Rui Yang and Marios Kogias. “HEELS: A Host-Enabled eBPF-Based Load Balancing Scheme”. In: *Proceedings of the 1st Workshop on eBPF and Kernel Extensions (eBPF'23)*. ACM, Sept. 10, 2023, pp. 77–83. DOI: [10.1145/3609021.3609307](https://doi.org/10.1145/3609021.3609307).
- [90] Nikita Shirokov Dasineni Ranjeeth. *Engineering at Meta – Open-sourcing Katran, a scalable network load balancer*. May 22, 2018. URL: <https://engineering.fb.com/2018/05/22/open-source/open-sourcing-katran-a-scalable-network-load-balancer/> (visited on 02/16/2023).
- [91] Shaun Crampton. *Introducing the Calico eBPF dataplane*. Feb. 25, 2020. URL: <https://www.tigera.io/blog/introducing-the-calico-ebpf-dataplane/> (visited on 03/28/2023).
- [92] Thomas Graf et al. *CNI Performance Benchmark – Cilium documentation*. May 11, 2021. URL: <https://docs.cilium.io/en/v1.12/operations/performance/benchmark/> (visited on 03/28/2023).
- [93] Erik Wenzel et al. *netperf – v2.7.0*. Nov. 15, 2020. URL: <https://packages.debian.org/bullseye/netperf>.
- [94] Bruce A. Mah et al. *iperf3 – v3.9 (cJSON 1.7.13)*. Aug. 15, 2020. URL: <https://packages.debian.org/bullseye/iperf3>.
- [95] nik-netlox (GitHub User). *SCTP: >10x BPF program runtime compared to TCP (GitHub, Issue #447, loxilb)*. Nov. 29, 2023. URL: <https://github.com/loxilb-io/loxilb/issues/447#issuecomment-1832999598> (visited on 03/21/2024).
- [96] Ragnar Rova et al. *wrk2 (GitHub, commit #920ce1e)*. Nov. 24, 2023. URL: <https://github.com/rrva/wrk2/tree/920ce1e> (visited on 04/09/2024).
- [97] Manuel Fähndrich et al. “Language support for fast and reliable message-based communication in singularity OS”. In: *Proceedings of the 1st European Conference on Computer Systems (EuroSys'06)*. ACM, 2006, pp. 177–190. DOI: [10.1145/1217935.1217953](https://doi.org/10.1145/1217935.1217953).
- [98] Galen C. Hunt and James R. Larus. “Singularity: rethinking the software stack”. In: *ACM SIGOPS Operating Systems Review* 41.2 (2007), pp. 37–49. DOI: [10.1145/1243418.1243424](https://doi.org/10.1145/1243418.1243424).
- [99] Zhe Zhou et al. “Userspace Bypass: Accelerating Syscall-intensive Applications”. In: *Proceedings of the 17th USENIX Symposium on Operating Systems Design and Implementation (OSDI'23)*. USENIX, 2023, pp. 33–49.
- [100] Alexander van der Grinten. *The Managarm Project*. Mar. 14, 2024. URL: <https://managarm.org/> (visited on 03/14/2024).
- [101] Suparna Bhattacharya et al. “Asynchronous I/O Support in Linux 2.5”. In: (2003). URL: <https://www.kernel.org/doc/ols/2003/ols2003-pages-351-366.pdf>.
- [102] Diego Didona et al. “Understanding modern storage APIs: a systematic study of libaio, SPDK, and io_uring”. In: *Proceedings of the 15th ACM International Conference on Systems and Storage (SYSTOR'22)*. ACM, 2022, pp. 120–127. DOI: [10.1145/3534056.3534945](https://doi.org/10.1145/3534056.3534945).

- [103] E. Zadok et al. “Efficient and safe execution of user-level code in the kernel”. In: *Proceedings of the 19th International Parallel and Distributed Processing Symposium (IPDPS’05)*. IEEE, 2005, 8 pp.--. DOI: 10.1109/IPDPS.2005.189.
- [104] Rami Rosen. *Network acceleration with DPDK [LWN.net]*. July 5, 2017. URL: <https://lwn.net/Articles/725254/> (visited on 03/08/2024).
- [105] *SPDK: Storage Performance Development Kit*. 2024. URL: <https://spdk.io/> (visited on 03/14/2024).
- [106] Alan Jowett. (*PREVAIL*) *Spectre Mitigations Issue (GitHub, Issue #229, vbpf/ebpf-verifier)*. 2021. URL: <https://github.com/vbpf/ebpf-verifier/issues/229>.
- [107] Elazar Gershuni. *Use verifier to identify dependent reads that require speculative load hardening (GitHub, Issue #229, vbpf/ebpf-verifier)*. May 16, 2022. URL: <https://github.com/vbpf/ebpf-verifier/issues/229> (visited on 02/07/2024).
- [108] Georg Julius Liefke. “Protecting the Linux Kernel from Transient-Execution Vulnerabilities in eBPF Bytecode”. Bachelor’s Thesis. Ruhr-Universität Bochum (RUB), July 19, 2022.
- [109] Guanhua Wang et al. “KLEESpectre: Detecting information leakage through speculative cache attacks via symbolic execution”. In: *Transactions on Software Engineering and Methodology* 29.3 (2020), pp. 1–31. DOI: 10.1145/3385897.
- [110] *Visual Studio 2022: C++ Developer Guidance for Speculative Execution Side Channels*. 2021. URL: <https://learn.microsoft.com/en-us/cpp/security/developer-guidance-speculative-execution?view=msvc-170>.
- [111] Paul Kocher. *Spectre Mitigations in Microsoft’s C/C++ Compiler*. Feb. 13, 2018. URL: <https://www.paulkocher.com/doc/MicrosoftCompilerSpectreMitigation.html>.
- [112] Kristof Beyls. *D49073 from the LLVM Phabricator archive: Introducing __builtin_speculation_safe_value*. July 9, 2018. URL: <https://reviews.llvm.org/D49073>.
- [113] Kristof Beyls. *D49070: Introduce llvm.speculation_safe_value intrinsic*. July 9, 2018. URL: <https://reviews.llvm.org/D49070>.
- [114] Kristof Beyls. *D41760: Introduce __builtin_load_no_speculate*. Jan. 5, 2018. URL: <https://reviews.llvm.org/D41760>.
- [115] Devin Jeanpierre and Chandler Carruth. *Mitigating Spectre v1 Attacks in C++*. Jan. 10, 2020. URL: <https://www.openstd.org/jtc1/sc22/wg21/docs/papers/2020/p0928r1.pdf>.
- [116] Sai Veerya Mahadevan, Yuuki Takano, and Atsuko Miyaji. “PRSafe: Primitive Recursive Function based Domain Specific Language using LLVM”. In: *Proceedings of the 20th International Conference on Electronics, Information, and Communication (ICEIC’21)*. IEEE, 2021, pp. 1–4. DOI: 10.1109/ICEIC51217.2021.9369763.
- [117] Luke Nelson, Xi Wang, and Emina Torlak. “A proof-carrying approach to building correct and flexible in-kernel verifiers”. *Linux Plumbers Conference (LPC’21)*. 2021. URL: <https://homes.cs.washington.edu/~lukenels/slides/2021-09-23-lpc21.pdf>.
- [118] Luke Nelson and Xi Wang. *Exoverifier (GitHub, commit #a0166f6)*. July 29, 2022. URL: <https://github.com/uw-unsat/exoverifier/tree/a0166f6> (visited on 02/26/2023).
- [119] Hongyi Lu et al. *MOAT: Towards Safe BPF Kernel Extension*. 2023. DOI: 10.48550/arXiv.2301.13421. arXiv: 2301.13421v2 [cs.CR].
- [120] Chandler Carruth. [*llvm-dev*] *RFC: Speculative Load Hardening (a Spectre variant #1 mitigation)*. Mar. 23, 2018. URL: <https://lists.llvm.org/pipermail/llvm-dev/2018-March/122085.html> (visited on 04/04/2024).
- [121] Zhiyuan Zhang et al. *Breaking and Fixing Speculative Load Hardening*. 2022. URL: <https://www.cryptojedi.org/papers/uslh-20220605.pdf>.
- [122] Basavesh Ammanaghatta Shivakumar et al. “Spectre Declassified: Reading from the Right Place at the Wrong Time”. In: *Proceedings of the 44th Symposium on Security and Privacy (SP’23)*. IEEE, 2023, pp. 1753–1770. DOI: 10.1109/SP46215.2023.10179355.
- [123] Chandler Carruth. *Speculative Load Hardening — LLVM 16.0.0git documentation*. Feb. 12, 2024. URL: <https://llvm.org/docs/SpeculativeLoadHardening.html>.
- [124] Oleksii Oleksenko et al. *You Shall Not Bypass: Employing data dependencies to prevent Bounds Check Bypass*. 2018. DOI: 10.48550/arXiv.1805.08506v3. arXiv: 1805.08506v3.
- [125] Ejobagom John Ojogbo, Mithuna Thottethodi, and T N Vijaykumar. “Secure automatic bounds checking: prevention is simpler than cure”. In: *Proceedings of the 18th International Symposium on Code Generation and Optimization (CGO’20)*. ACM, 2020, pp. 42–55. DOI: 10.1145/3368826.3377921.
- [126] Narayan, Disselkoe, Moghimi, et al. “Swivel: Hardening WebAssembly against Spectre”. In: *Proceeding of the 30th USENIX Security Symposium (USENIX Security ’21)*. USENIX, 2021, pp. 1433–1450. ISBN: 978-1-939133-24-3. URL: <https://www.usenix.org/conference/usenixsecurity21/presentation/narayan>.
- [127] Mathé Hertogh et al. “Quarantine: Mitigating Transient Execution Attacks with Physical Domain Isolation”. In: *Proceedings of the 26th International Symposium on Research in Attacks, Intrusions and Defenses (RAID’23)*. ACM, 2023, pp. 207–221. DOI: 10.1145/3607199.3607248.
- [128] Henriette Hofmeier. “Dynamic Reconfiguration of Hardware-Vulnerability Mitigations in the Linux Kernel”. Master’s Thesis. Friedrich-Alexander Universität Erlangen-Nürnberg (FAU), June 24, 2022. URL: https://sys.cs.fau.de/publications/2022/hofmeier_22_thesis.pdf.
- [129] Jonathan Behrens et al. “Efficiently mitigating transient execution attacks using the unmapped speculation contract”. In: *Proceedings of the 14th Symposium on Operating Systems Design and Implementation (OSDI’20)*. USENIX, 2020, pp. 1139–1154. URL: <https://www.usenix.org/conference/osdi20/presentation/behrens>.
- [130] *Core Scheduling – The Linux Kernel documentation (v6.5)*. Nov. 8, 2023. URL: <https://www.kernel.org/doc/html/v6.5/admin-guide/hw-vuln/core-scheduling.html> (visited on 04/04/2024).
- [131] Intel. *An optimized mitigation approach for Load Value Injection*. Mar. 9, 2020. URL: <https://www.intel.com/content/www/us/en/develop/articles/software-security-guidance/best-practices/optimized-mitigation-approach-load-value-injection.html>.
- [132] Gernot Heiser, Toby Murray, and Gerwin Klein. “Towards Provable Timing-Channel Prevention”. In: *ACM SIGOPS*

Operating Systems Review 54.1 (2020), pp. 1–7. DOI: [10.1145/3421473.3421475](https://doi.org/10.1145/3421473.3421475).

[133] Scott Buckley et al. *Proving the Absence of Microarchitectural Timing Channels*. 2023. DOI: [10.48550/arXiv.2310.17046](https://doi.org/10.48550/arXiv.2310.17046). arXiv: [2310.17046](https://arxiv.org/abs/2310.17046).