

Towards Secure BPF Kernel Extension with Hardware-enhanced Memory Isolation

Lijian Huang*, Hongyi Lu*, Shuai Wang†, and Fengwei Zhang†

Abstract—The Linux kernel extensively uses the Berkeley Packet Filter (BPF) to allow user-written BPF applications to execute in the kernel space. The BPF employs a verifier to check the security of user-supplied BPF code statically. Recent attacks show that BPF programs can evade security checks and gain unauthorized access to kernel memory, indicating that the verification process is not flawless. In this paper, we present MOAT, a novel hardware-assisted, cross-platform isolation framework designed to protect the kernel from malicious BPF programs. MOAT introduces a two-layer memory isolation scheme that leverages hardware features such as Intel MPK and Arm Stage-2 translation to enforce isolation. Our design overcomes several key challenges, including the limited scalability of available hardware isolation mechanisms and the risk of helper function abuse. We implement MOAT for Intel x86 and Arm on Linux (ver. 6.1.38), and our evaluation shows that MOAT delivers low-cost isolation of BPF programs under mainstream use cases, such as isolating a BPF packet filter with only 3% throughput loss.

Index Terms—BPF, Linux kernel, Intel MPK, Arm Stage-2, Memory Isolation.

I. INTRODUCTION

IT is common to extend kernel functionality by allowing user applications to download code into the kernel space. In 1993, the well-known Berkeley Packet Filter (BPF) was introduced for this purpose [1]. The classic BPF is an infrastructure that inspects network packets and decides whether to forward or discard them. With the introduction of its extended version (referred to as eBPF) in the Linux kernel, BPF soon became more powerful and is now utilized in numerous real-world scenarios, such as load balancing, system tracing, and system call filtering [2–7].

To ensure security, BPF is equipped with a *verifier* [8]. The verifier performs a variety of static analyses to ensure the user-supplied code is secure. For instance, the verifier tracks the bounds of all pointers to prevent out-of-bound access. Given that BPF code runs directly within the kernel, the verifier

becomes crucial for BPF security. Nevertheless, as pointed out by recent studies [9–13], the current verifier has various limitations and is insufficient for the overall security of BPF. First, the current BPF ecosystem supports a variety of functionalities, such as packet forwarding and kernel debugging [14, 15]. Supporting all these functionalities in the verifier results in a complicated verification process. Though the verifier has been partially verified via formal methods [16], the unverified part and the gap between abstraction and implementation still result in vulnerabilities [17–24]. Second, due to the rapid expansion of BPF capabilities, the verifier is frequently updated, and it is inherently difficult to update a complex static verification tool without introducing new vulnerabilities [25]. To date, the BPF subsystem has been repeatedly exploited. For instance, two privilege-escalation vulnerabilities have been discovered in `bpf_ringbuf`, a recent BPF feature introduced in 2020 [1]. Further, the verifier’s register-value tracking is quite complex and often bypassed via corner-case operations (e.g., sign extension) [17–20].

Given the increasing security threats in BPF and the challenge of enforcing safe BPF programs with merely static verification, we seek to employ hardware extensions to sandbox untrusted BPF programs. We present MOAT, a hardware-assisted, cross-platform isolation mechanism for BPF programs. MOAT leverages two classes of hardware primitives: key-based hardwares, such as Intel Memory Protection Keys (MPK) [26] and Arm Permission Overlay Extension (POE) [27], and virtualization-based hardwares, such as Arm Stage-2 translation [27], AMD Rapid Virtualization Indexing (RVI) [28] and RISC-V H-mode [29]. Given the widespread support for Intel MPK and Arm Stage-2 translation on modern processors, we select these two primitives from each class as representative examples to illustrate MOAT’s cross-platform design. Our approach is not limited to these primitives and can be adapted to other hardware features, as discussed in Sec. IX.

To avoid confusion, we refer to these two hardware-assisted isolation schemes as MOAT-KEY (key-based, implemented with Intel MPK) and MOAT-VIR (virtualization-based, implemented with Arm Stage-2 translation), while MOAT denotes the holistic design that encompasses both approaches throughout this paper. For MOAT-KEY, we assign separate MPK keys to the kernel and BPF programs, restricting BPF’s access to kernel memory during execution. For MOAT-VIR, we construct a dedicated Stage-2 page table to isolate BPF programs from the kernel.

Despite the promising potential of these hardware features in designing MOAT, we encountered and addressed *two major technical hurdles*. First, both classes of hardware primitives of-

Lijian Huang* and Hongyi Lu* contributed equally to this work. Shuai Wang† and Fengwei Zhang† are the corresponding authors.

Lijian Huang is with Research Institute of Trustworthy Autonomous Systems, and Department of Computer Science and Engineering, Southern University of Science and Technology, China. E-mail: huanglj2023@mail.sustech.edu.cn.

Hongyi Lu is with Research Institute of Trustworthy Autonomous Systems, Southern University of Science and Technology, China, and also with Department of Computer Science and Engineering, Hong Kong University of Science and Technology, China. E-mail: luhy2017@mail.sustech.edu.cn.

Shuai Wang is with Department of Computer Science and Engineering, Hong Kong University of Science and Technology, China. E-mail: shuaiw@cse.ust.hk.

Fengwei Zhang is with Department of Computer Science and Engineering, and Research Institute of Trustworthy Autonomous Systems, Southern University of Science and Technology, China. E-mail: zhangfw@sustech.edu.cn.

fer limited support for scalable, fine-grained isolation, making it difficult to efficiently isolate a large number of concurrent BPF programs. To address this hurdle, we propose a novel two-layer isolation scheme that protects both the kernel and benign BPF programs from malicious BPF programs. Layer-I leverages the hardware isolation primitives to construct three isolation domains, preventing unauthorized kernel access from BPF programs. Layer-II enforces intra-BPF isolation within the same domain by assigning each BPF program a dedicated address space, while mitigating the Translation Lookaside Buffer (TLB) flush overhead with emerging hardware features. Second, the *helper functions* provided by the BPF subsystem may also be exploited by attackers. On the one hand, MOAT must permit benign BPF programs to invoke these helpers without restriction. On the other hand, it must ensure that the helpers cannot be abused by malicious BPF programs. However, designing individual security policies for each helper function entails significant engineering effort and risks bloating the codebase. To address this, we propose two generic defense schemes that are agnostic to the internal design of individual helper functions. We demonstrate these schemes effectively cover a wide range of helpers (see Appendix A).

This paper extends our previous work [30], published in the Proceedings of the 33rd USENIX Security Symposium 2024, with significant enhancements focused on cross-platform capability. We generalize the design of MOAT to support multiple architectures and introduce a new approach, MOAT-VIR, for platforms equipped with virtualization-based hardware primitives, where key-based primitives such as Intel MPK are unavailable. While our implementation of MOAT-VIR utilizes Arm’s Stage-2 translation which is a mature and widely supported hardware feature of Arm processes, the approach is not limited to Arm and can be adapted to other platforms supporting similar virtualization features.

We systematically examine how MOAT mitigates attacks targeting the BPF ecosystem and assess potential threats to MOAT itself. We also conduct an empirical analysis of all recent CVEs relevant to MOAT’s application domain, demonstrating that MOAT successfully mitigates each one. To assess performance impact, we evaluate MOAT across a range of micro and macro benchmark settings. In network application scenarios, MOAT-KEY incurs a maximum performance overhead of 3%, while MOAT-VIR experiences higher overhead in certain cases. Both prototypes demonstrate strong performance in various benchmarks. In the Phoronix Test Suite [31] network benchmark, MOAT-KEY and MOAT-VIR show average overheads of 2.7% and 3.5%, respectively. We further evaluate MOAT with system tracing workloads, where MOAT introduces an average overhead of 5.5% on Intel and 8.7% on Arm. These results indicate that MOAT delivers fundamental security improvements with modest performance overhead across realistic BPF use cases.

Building upon our previous work, this paper presents the final, extended version of MOAT, and the following contributions reflect both our original design and the significant cross-platform extensions introduced in this work.

- Instead of merely relying on the BPF verifier to statically validate BPF programs, this paper advocates isolating BPF

programs with hardware extensions, including key-based and virtualization-based primitives, to effectively ensure the memory safety of BPF programs.

- Technically, MOAT is designed to address domain-specific challenges, including limited hardware isolation support and the risk of helper function abuse in the BPF ecosystem. MOAT features a two-layer isolation scheme to protect both the kernel and benign BPF programs from malicious ones, and incorporates various design considerations to ensure memory safety with minimal performance overhead.
- We implemented the prototype of MOAT on Linux 6.1.38¹ and thoroughly evaluated their security over different attack scenarios (including all memory-relevant BPF CVEs in the past decade) and performance using various benchmark suite. The evaluation shows that MOAT delivers a principled security warranty with minimum overhead.

II. BACKGROUND

A. Berkeley Packet Filter (BPF)

In this section, we provide a brief overview of the BPF and its core components.

BPF Overview. BPF [1] was originally introduced to facilitate flexible network packet filtering. Instead of inspecting packets in the user space, users can provide BPF instructions specifying packet filter rules, which are directly executed in the kernel. BPF allows configurable packet filtering without costly context switching and data copying. Modern Linux kernel features extended BPF (eBPF), a Linux subsystem which supports a wide range of use cases, such as kernel profiling, load balancing, and firewalls. Popular applications such as Docker [33], Katran [14], and kernel debugging utilities like Kprobes [15] utilize or are built directly on top of BPF.

Fig. 1 depicts an overview of how BPF programs are compiled and deployed. The BPF subsystem offers ten general-purpose 64-bit registers, a stack, BPF customized data structures (often called BPF maps), and a set of BPF helper functions. To use BPF (e.g., for system tracing), users first write their own BPF programs (in C code) to specify the functionality, which, in turn, are compiled into bytecode and loaded into the kernel. Given that BPF code is written by untrusted users, the kernel employs a verifier to conduct several checks during the bytecode loading stage (see below). By default, the verified bytecode is further compiled into native code by an in-kernel Just-In-Time (JIT) compiler for better performance. Additionally, on platforms without the JIT support, the bytecode is alternatively executed by the BPF interpreter. The BPF program is then attached to certain kernel components based on its specific end goal. For instance, as shown in Fig. 1, a BPF program attaches to the kernel as the packet filter, monitoring network traffic and sending statistics back to the user space via a BPF map.

BPF Verifier. BPF programs are written in C and compiled into a RISC-like instruction set. As aforementioned, the kernel strictly verifies the BPF programs upon loading to ensure they are safe to execute. Fig. 2 illustrates the verification process

¹We release the codebase of MOAT on our site [32]. We will maintain MOAT to benefit the community and follow-up research.

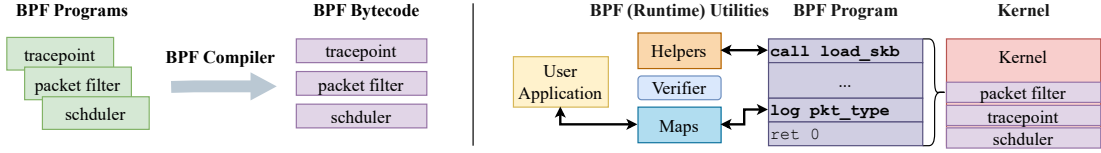


Fig. 1: BPF overview. We illustrate the BPF compilation procedure and execution context of a sample BPF packet filter.

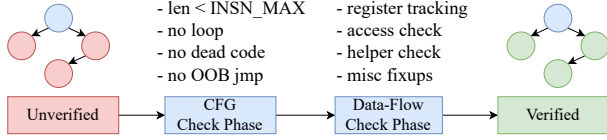


Fig. 2: BPF verification process.

in a holistic manner. First, a BPF program is parsed into a control flow graph (CFG) by the verifier, which performs a CFG check phase to ensure four key properties: 1) the program size is within a limit; 2) there are no back edges (loops) on its CFG; 3) there is no unreachable code; and 4) all jumps are direct jumps and refer to a valid destination.

The verifier then tracks the value flow of every register to deduce its value ranges conservatively. With these ranges, the verifier decides if a pointer accesses safe memory and if a parameter is valid. Since this analysis is performed statically, it is possible for a malicious BPF program to exploit a vulnerability to bypass it [17–24].

BPF Helpers. The kernel also limits the functions a BPF program may call. Those functions are dubbed BPF helpers, as shown in Fig. 1. To date, there are over 200 helpers provided by the kernel [34]. Depending on the task, a BPF program can usually call a group of relevant helpers. For example, a BPF packet filter can call `skb_load` to read packet data, but is not allowed to call any helper related to system tracing.

BPF Maps. Out of security concern, the kernel also sets a strict space limit on BPF programs. Each program, by default, can only use up to 512 bytes of stack space and 10 registers, which is far from enough for certain BPF programs. To address this problem, BPF maps can be allocated to provide additional space for BPF programs. To date, there are over 30 types of maps supported by kernel [35]. Based on the isolation requirements, they can be roughly categorized into two types. The first type is maps that own a memory region. The most commonly used maps, hash maps and array maps, belong to this category. BPF programs use them to store data and communicate with user space. Therefore, a proper access permission has to be set for these maps (see MOAT’s solution in Sec. IV-A). The second type holds references to other kernel resources (e.g., file descriptors). BPF programs are restricted to using helpers to interact with this type of map. Thus, MOAT forbids BPF programs from directly accessing them.

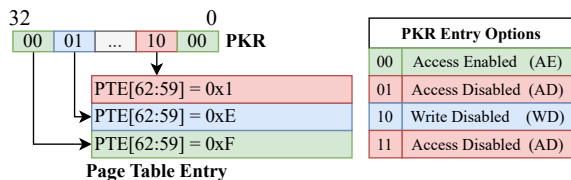


Fig. 3: Intel MPK overview.

Classic BPF (cBPF). cBPF specializes in tasks like syscall

filtering (e.g., `seccomp-BPF`) and has more restrictions than eBPF. We clarify that MOAT supports both of them. We also evaluate MOAT using `seccomp-BPF` with cBPF in Sec. VI-D. In this paper, we use BPF to refer to both cBPF and eBPF, as the kernel internally converts cBPF to eBPF.

B. Hardware Features in MOAT

In this section, we introduce the hardware features leveraged by MOAT to efficiently isolate BPF programs.

Intel MPK. Intel introduced MPK [26] to provide efficient page table permissions control. By assigning an MPK protection key to the page table entries (PTEs) of one process, users can enable intra-process isolation and confidential data access control [36–39]. As illustrated in Fig. 3, MPK uses four reserved bits [62:59] in each Page Table Entry (PTE) to indicate which protection key is attached to this page. Those three PTEs in Fig. 3 are assigned with keys `0x1`, `0xE` and `0xF`, respectively. Since there are only 4 bits involved, the maximum number of keys is 16. Then, a new 32-bit register named Protection Key Register (PKR) is introduced to specify the access permission of these protection keys. Each key occupies two bits in PKR, whose values flag the access permission of the page. In Fig. 3, the access permissions of the three pages are 01 access-disabled (AD), 10 write-disabled (WD), and 00 access-enabled (AE), respectively. By writing to certain bits in PKR, the access permission of corresponding pages can be configured efficiently without having to modify the PTEs. There are actually two versions of MPK. One applies to the user space, while the other applies to the kernel space. For brevity, we refer to these two versions in their conventional abbreviations as Protection Key Supervisor (PKS) and Protection Key User (PKU), respectively. In MOAT-KEY, we use PKS instead since our goal is to isolate in-kernel BPF program, while the corresponding PKR is a Model Specific Register (MSR) named `IA32_PKRS`.

Arm Stage-2 Translation. The Arm platform uses a two-stage memory translation process [27] that converts Virtual Address (VA) to Physical Address (PA) via an Intermediate Physical Address (IPA) at Stage-2. This design allows the operating system to manage memory independently of the actual physical memory layout. The base address of current Stage-2 page table is stored in the `VTTBR_EL2` register. Access permissions in Stage-2 translation are controlled via the Stage-2 Access Permissions (S2AP) field located in bits [7:6] of the PTEs. The S2AP encodes four levels of access control, of which we use three: 00 access-disabled (AD), 01 write-disabled (WD), and 11 access-enabled (AE). These permissions are enforced independently of the stage-1 permissions set by the guest OS. As a result, even if a guest process marks a page as writable in its own page table, the hypervisor can override this setting

and restrict access through Stage-2 controls. If a fault occurs when performing the Stage-2 translation, an ABORT exception will be triggered, and the hypervisor is expected to handle it.

Reducing TLB overhead To reduce address-space switching overhead, MOAT leverages the Address Space Identifier (ASID) on modern hardware platforms. Specifically, ASID reduces overhead by avoiding unnecessary TLB flushes between address-spaces with the same ASID. We introduce this feature here as it is crucial to our later discussion. Note that Arm and Intel refer to this feature as Virtual Machine Identifier (VMID) and Process Context Identifier (PCID), respectively. We use these names interchangeably in this paper for brevity.

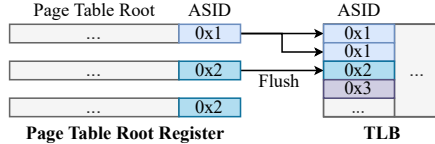


Fig. 4: ASID overview.

PCID. MOAT-KEY uses the PCID to reduce the overhead caused by address-space switching; we introduce PCID here. On the x86 platform, the CR3 register holds the page table root of the current process. Modifying the CR3 register causes a complete TLB flush and is therefore costly. Fortunately, Intel introduced PCID to address this issue. The lower 12 bits [11:0] of CR3 register are PCID, identifying the owner of the page table, while the highest bit of the new CR3 value controls the flushing behavior of TLB. If the highest bit is 1, this modification does not flush TLB at all; if the highest bit is 0, then this modification only flushes the TLB entries of the PCID in this new CR3 value. This feature enables fast address-space switch without costly TLB flush. Since there are only 12 reserved bits for PCID, it supports up to 4096 different processes with isolated TLB entries.

VMID. Similar to PCID on x86, MOAT-VIR leverages VMID to reduce the overhead associated with address-space switching on Arm platform. The VMID is associated with the Stage-2 translation base register (VTTBR_EL2), identifying the owner of the virtual address space. This mechanism allows efficient context switching between address spaces with different VMIDs without invalidating previous translation table mappings from the TLB. When switching between different Stage-2 page tables, the hardware compares the VMID of each TLB entry with the current VMID in VTTBR_EL2.VMID, bits [63:48]. If they match, the TLB entry remains valid; otherwise, it is invalidated. This selective invalidation mechanism minimizes the need for a full TLB flush during context switches. A standard 16-bit VMID supports up to 65,535 distinct BPF programs, which we consider sufficient for the deployment of MOAT-VIR.

III. MOTIVATION AND THREAT MODEL

A. Motivation

In this section, we discuss the typical threats to the BPF verifier, the restriction on unprivileged BPF brought by these threats, and lastly, the motivation for our research.

Fast Feature Evolving. As a fast-developing technology, threats may come from the inconsistency between the constantly expanding BPF capabilities and the rigorous static verification process imposed on them [21, 25]. It is a common practice to add corresponding verification procedures simultaneously when introducing new features to BPF programs. However, it is difficult to implement a verifier that supports all these features yet still does not miss any edge case, which already has over 10K LoC with various functionalities [8].

Challenging Register-value Tracking. Second type of threats originates from the complexity of the register-value tracking. Although the soundness of such a tracking mechanism is formally proved [16], there exist gaps between the actual implementation and abstraction of the register-value tracking, especially in some corner cases, such as sign extension, truncation, and bit operators [17–24].

Unprivileged BPF. BPF was originally designed as a restricted interface for *unprivileged* users to extend kernel functionality. It comes with a fine-grained privilege system [40] that allows users to tune a specific part of the kernel without root. However, numerous vulnerabilities [17–20] indicate that the verifier is not reliable, and consequently, major distributions have banned unprivileged users from loading BPF programs [41, 42]. Despite this, there is still a long-lasting desire for unprivileged BPF in the community. For example, the `seccomp-BPF` users have been asking for unprivileged BPF support for a long time [43].² Moreover, there have been continuous efforts (from 2016 to 2023) in the community to re-emerge unprivileged BPF again [45–47]. Unfortunately, these efforts fail as they only focus on enhancing the verifier itself, which is already over-complicated and error-prone.

Motivation. Overall, seeing BPF’s potential and its current restriction, we propose MOAT as an isolation scheme complementary to the BPF verifier. On the one hand, this isolation scheme shall make BPF more accessible to unprivileged users whilst maintaining security. On the other hand, even for privileged users, MOAT provides the security guarantee that the BPF programs obtained from a potentially untrusted source are isolated from the kernel. Overall, we aim to provide a more secure and accessible BPF ecosystem, thereby promoting its development and adoption in the community.

B. Threat Model

Our threat model considers a practical setting that is aligned with existing BPF vulnerabilities [17–24]. Attackers can load their prepared BPF code into the kernel space to launch exploitation. In particular, we assume attackers are *non-privileged users* with BPF access since a root user already has control over almost the entire kernel. MOAT isolates user-submitted BPF programs and prevents them from accessing kernel memory regions. As will be introduced in Sec. IV, a BPF program is given only the necessary resources and privileges to complete its task. We present the threat models of major components in our research context as follows.

²We clarify that `seccomp` already supports classic-BPF (cBPF), which lacks expressiveness and no longer updates [44]. BPF here refers to eBPF.

BPF Programs. We assume that malicious BPF programs are able to bypass checks statically performed by the verifier; they may thus behave maliciously during runtime. Our threat model deems BPF programs as *untrusted*.

BPF Helper Functions. These helpers act as the intermediate layer between the BPF subsystem and kernel. Certain malicious BPF programs can abuse these helpers to perform attacks, and therefore, we assume they are also *untrusted*. MOAT mitigates risks raised by adversarial-manipulated helper functions with practical defenses.

Out of Scope. The main objective of MOAT is to mitigate memory exploitation performed by BPF programs. Other subtle attacks (not relevant to memory exploits), such as speculation, race condition, and Denial of Service (DoS) toward the BPF subsystem [48, 49] are not considered. They do not specifically exist in BPF [50, 51], and are addressed by relevant research [52–55]. We thus treat them as orthogonal. Also, BPF subsystem comes with a set of user-space facilities such as `libbpf`; bugs in them are not considered by MOAT. Note that MOAT mitigates information leakage that is due to out-of-bounds memory access; if the leakage is due to issues like speculation [49], then it is out of the scope of MOAT. If an object is already accessible to the BPF program, MOAT does not prevent information leakage through legitimate accesses. Privileged loaders (e.g., CNIs) and vulnerabilities from other subsystems (e.g., memory management) are out of scope. Corrupted BPF return values that only affect program functionality, but do not compromise kernel security, are not the focus of this work.

We clarify that MOAT focuses on the kernel memory exploitation via BPF, its most prevalent threat. The vulnerabilities mitigated by MOAT typically receive high threat scores in vulnerability databases [17–25] with public PoC exploits [56], whereas above-precluded vulnerabilities often lack exploits [57].

IV. DESIGN

MOAT Overview. As described in Sec. III-A, the current security design against malicious BPF programs solely relies on the static analysis performed by the BPF verifier, which is seen as a weak point and exploitable by non-privileged users. MOAT instead delivers a principled isolation of BPF programs from the rest part of the kernel using hardware features and prevents bypasses.

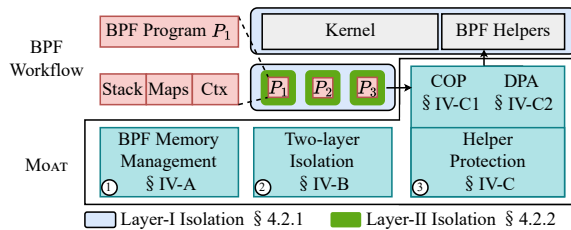


Fig. 5: MOAT overview.

Fig. 5 depicts an overview of MOAT and how it is integrated into the workflow of BPF programs. ① Given a user-submitted BPF program P , MOAT statically allocates the necessary memory regions the program needs, such as stack, maps, and context based on P 's metadata (Sec. IV-A). ② When the

kernel invokes P , MOAT isolates P from the kernel using hardware memory isolation primitives (Layer-I in Sec. IV-B1), and constrains P in its isolated address space (Layer-II in Sec. IV-B2). ③ On the occasions that P calls helpers, depending on the helper types, MOAT adjusts the involved memory region permissions (Sec. IV-C1) and validates the helper parameters (Sec. IV-C2) to prevent the helpers from being abused. For brevity, we refer to the platform-agnostic design as MOAT, while using MOAT-KEY and MOAT-VIR to denote platform-specific implementations on Intel and Arm respectively.

A. BPF Memory Management in MOAT

Further to the overview, we introduce how MOAT manages the BPF memory. The BPF memory refers to the memory regions a BPF program needs to function properly, including descriptor tables, stacks, maps, and runtime context.

Essential Regions. Each platform requires certain essential memory regions for fundamental operations such as exception handling and interrupt processing. These regions are assigned to a shared area accessible by all BPF programs. On x86 platforms, this includes the Global Descriptor Table (GDT) and Interrupt Descriptor Table (IDT); on Arm, the vector tables for each exception level; and on RISC-V, the trap handler and vector mode. To ensure correct exception and interrupt handling during BPF program execution, these regions are mapped into the BPF memory with read-only permissions.

Stack. BPF programs use a 512-byte stack space to store local variables and function frames. The verifier determines if a program makes out-of-bounds access toward the stack. Thus, if the BPF program passes the static checks, its stack is directly allocated from the kernel stack. However, as discussed in Sec. III-A, certain vulnerabilities may allow BPF programs to bypass this check. Thus, MOAT needs to allocate the stack as a part of BPF memory and swap stacks to prevent the BPF programs from tampering with the kernel stack.

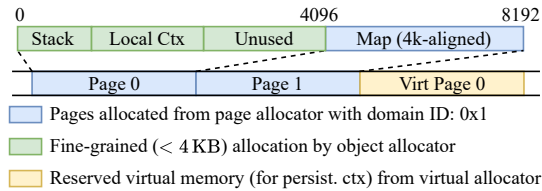
Maps. As described in Sec. II-A, the maps are utilized to store data and communicate with the user space. Linux provides a set of helper functions for BPF programs to interact with maps. For example, `bpf_map_lookup_elem` returns the pointer of an element so that the program can modify its value. This means that BPF programs must have access to these elements' memory. Thus, MOAT allocates these maps as a part of BPF memory. Note that we do not allocate the metadata of these maps inside BPF memory since they contain exploitable structures like function pointers. For example, only the key-value pair fields of the hashmap elements are allocated in the BPF memory.

Runtime Context. The context refers to BPF program parameters, which vary depending on the BPF program types. We investigated the BPF contexts of common BPF program types and summarized our findings in Table I. Most of these contexts are local objects on the kernel stack and are passed to BPF programs as parameters, such as `bpf_cgroup_dev_ctx`. For this type of BPF context, MOAT allocates them on the BPF stack instead so that the BPF programs can still access them without the permission to access the kernel stack. However,

TABLE I: BPF context of common program types.

Category	Program Type	Context Type	Persistent	Nested	Note
Network	Socket Filter	sk_buff *	Yes	Yes	Socket packet buffer
	Socket Ops	bpff_sock_ops *	No	Yes	Socket events (timeout, retransmission, ...)
	Socket Lookup	bpff_sk_lookup *	No	Yes	Packet information for socket lookup
	XDP	xdp_md *	No	Yes	Metadata of xdp_buff
Tracing	Kprobe	pt_regs *	No	No	Register status on probed location
	Tracepoints	Depending on tracepoint types	No	No	Relevant tracepoint information
	Perf Event	bpff_perf_event_data *	No	No	Perf. event (register status, sample period)
Cgroup	Cgroup Socket Filter	sk_buff *	Yes	Yes	Socket packet buffer under specific cgroup
	Cgroup Device	bpff_cgroup_dev_ctx *	No	No	Device ID, access type (read, write)

there also exist contexts that are not local objects on the stack but persistent kernel structures. For example, `sk_buff` holds the network packet received by a socket and is also passed to BPF socket filter programs as context. For this type of persistent context (denoted in the fourth column of Table I), MOAT dynamically maps the physical page of the corresponding context into the BPF memory. The reason why we choose to map instead of creating a local copy is that `sk_buff` is typically hundreds of bytes. Our preliminary experiment shows that syncing between the local copy and the actual kernel object brings non-trivial overheads. Furthermore, network-related BPF contexts (e.g., `bpff_sock_ops`) may contain nested pointers to other kernel structures (denoted in the fifth column of Table I). Including only these pointers in BPF memory triggers a false alarm, as these nested structures are not included. We clarify that BPF programs only access limited fields of these nested structures. Thus, MOAT reserves a part of BPF memory to mirror these nested fields efficiently so that they can be accessed by the BPF programs.

**Fig. 6:** BPF memory allocators.

BPF Memory Allocators. As shown in Fig. 6, MOAT provides three types of allocators to manage these BPF memory regions. When loading a BPF program, the page allocator first allocates physical pages (Page 1) for its BPF memory; these pages become a part of its BPF memory. The object allocator handles fine-grained allocations from the allocated pages (Page 0) that are less than the page size (4 KB), e.g., the BPF stacks (512 bytes). Lastly, the virtual allocator controls the virtual memory not backed with concrete physical pages, such as reserving a part of virtual BPF memory (Virt Page 0) to map persistent BPF contexts like `sk_buff`. Note that we also modify the implementation of BPF maps to use MOAT’s allocator.

B. Two-layer Isolation

Challenge. MOAT aims to use hardware features to isolate the kernel from malicious BPF programs. However, existing platforms offer limited support for scalable and flexible memory isolation. For example, Intel’s MPK supports only 16 protection keys, while other platforms currently lack an equivalent mechanism for fine-grained memory isolation.

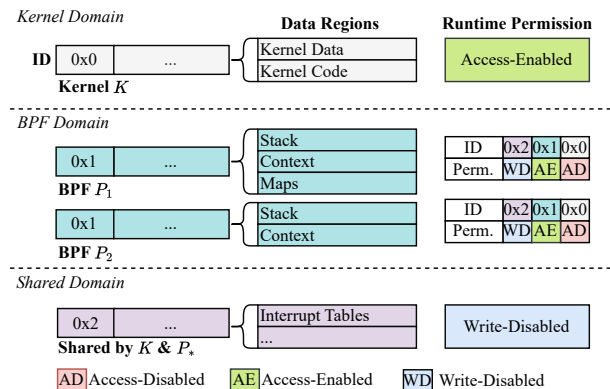
Solution. We propose a novel two-layer isolation scheme. The first layer is a lightweight isolation domain, dividing

the memory into three domains. The second layer is an isolated address space. Though isolating address spaces is less efficient, we manage to reduce its overhead to a minimum using contemporary TLB flush reduction hardware features.

1) *Layer-I: Lightweight Isolation Domain.* The main objective of MOAT is to isolate the kernel from malicious BPF programs. Thus, we design the first layer of isolation domain, using different memory isolation primitives on Intel and Arm platforms to achieve this goal. MOAT builds three kinds of isolated domains: the BPF domain, the kernel domain, and the shared domain. Specifically, MOAT-KEY leverages key-based isolation primitives such as Intel MPK, whereas MOAT-VIR employs virtualization-based mechanisms like Arm Stage-2 translation to enforce domain separation.

As depicted in Fig. 7, all BPF programs reside in the BPF domain with the domain ID `0x1`. MOAT grants a BPF program access (i.e., access-enabled; AE) to the BPF domain (`0x1`) when executing the program. The kernel domain holds all kernel pages with domain ID `0x0` and is only accessible by the kernel itself. When entering a BPF program, this kernel domain (`0x0`) becomes access-disabled (AD). However, the shared domain (`0x2`) comprises memory regions like interrupt handling tables. These regions are crucial for low-level routines. Thus, they are made write-disabled (WD) instead of access-disabled for BPF programs.

This domain design effectively mitigates malicious BPF programs targeting the kernel. For these malicious programs, a modus operandi is to introduce an unsanitized kernel pointer by exploiting a verifier vulnerability. Then, the malicious BPF program arbitrarily tampers the kernel using that pointer, leading to a full-blown exploitation. Isolating BPF programs from the kernel effectively stops such attacks, as all malicious kernel access directly from BPF programs is prevented by hardware memory isolation primitives.

**Fig. 7:** Memory domains of MOAT.

MOAT-KEY uses MPK to enforce these three domains. MOAT-KEY assigns the protection key `0x0` to the kernel

domain, 0×1 to the BPF domain, and 0×2 to the shared domain, where each key value equals to its domain's ID. As depicted in Fig. 7, MOAT-KEY then grants the access permission to each domain by setting the PKR bits of the corresponding protection key. This design only needs three (out of 16) keys from MPK.

MOAT-VIR utilizes Arm Stage-2 translation, to enforce isolation across three domains. The kernel domain is mapped within the system's original Stage-2 page table, while the BPF domain is mapped to a dedicated Stage-2 page table for BPF programs, ensuring that no unintended kernel memory mappings are present. The shared domain's memory is mapped into the BPF programs' dedicated Stage-2 page table with write-disabled (WD) permission, preventing unauthorized modifications while allowing controlled access.

2) *Layer-II: Isolated BPF Address Space:* In Sec. IV-B1, we have discussed how MOAT prevents BPF attacks targeting the kernel. However, MOAT only builds isolation between the kernel and BPF programs. All BPF programs still share the domain, allowing a malicious BPF program to tamper with the memory of benign BPF programs. Inspired by prior works in user-space isolation [58], we set up an isolated address space for each BPF program to prevent such tampering. Consequently, when a malicious BPF program tries to access the memory regions of another BPF program, a page fault occurs, and the malicious BPF program is immediately terminated.

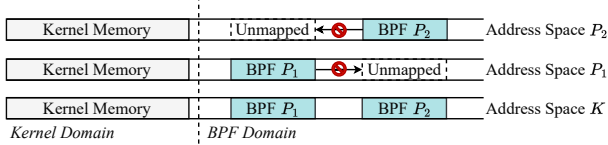


Fig. 8: Isolated address spaces of BPF programs.

Fig. 8 illustrates the isolated address spaces of two BPF programs, P_1 and P_2 . In the address space of P_1 , the mapping of P_2 does not exist. Similarly, the mapping of P_1 does not exist in the address space of P_2 either. This effectively prevents BPF programs from accessing each other and addresses the abovementioned issue. For MOAT-KEY, the isolated address space resides within the kernel space. For MOAT-VIR, isolation is enforced through Stage-2 translation. The isolated address space can be constructed by splitting the second-highest level of the BPF domain's Stage-2 page table, with each entry dedicated to a separate BPF program.

To mitigate the high TLB flush overhead (and TLB misses) from the address-space switching, we use PCID on x86 and VMID on Arm to keep the TLB entries from different BPF programs isolated. RISC-V also supports similar ASID features. MOAT allocates a non-overlapped virtual address space with a unique PCID or VMID for each of them. In rare cases where the number of concurrently executing BPF programs exceeds the available PCIDs or VMIDs, MOAT has to flush the TLB when two BPF programs that share the same identifier run consecutively. Since there are 4,096 PCIDs and up to 65,536 VMIDs, we expect such conflicts to be rare, especially considering that the kernel also periodically flushes TLB, thus clearing the conflicting entries. Even when such cases occur, MOAT only flushes the TLB entries of the conflicting PCID or VMID, preserving the integrity of other cached entries.

Why Intra-BPF Isolation. One may question the necessity for the isolation between BPF programs, as most existing BPF exploits target the kernel. However, since BPF maps are the only bridge between the BPF programs and the user space, the configurations of a BPF program have to be saved in its maps so that users can change its behavior without reloading it. This paradigm makes cross-BPF attack a noticeable threat [59]. For example, an attacker may load a malicious BPF program to change the behavior of another program by tampering with its configuration maps, disrupting resources accounting, or even nullifying security checks. Intra-BPF isolation is essential for preventing such attacks. Careful readers may note that `bpf2bpf` calls [1] appear to involve cross-BPF calls. Since the callee is a subroutine within the same BPF object and shares the caller's address space. Thus, Intra-BPF Isolation is not required for `bpf2bpf` calls.

C. Helper Security Mechanism

As mentioned in Sec. II-A, the kernel provides a set of helper functions for BPF programs. As these helpers act as the interfaces between the kernel and BPF programs, they can also be abused by malicious programs to launch attacks. MOAT needs to protect the helpers from such abuse.

Our investigation of existing BPF vulnerabilities shows that the malicious BPF programs typically abuse the BPF helpers in two ways: ① the helper contains a defect, which is exploited by the malicious programs [60]; ② the helper itself is correct, but the malicious programs pass invalid parameters to abuse it [25]. For ①, a typical case is that the helper itself contains defects such as heap overflow. These defects are leveraged by malicious programs to overwrite the sensitive fields (e.g., function pointers) of BPF-related kernel objects. For ②, since the helper by itself is correct, the malicious programs are typically restricted to leaking kernel pointers (by passing invalid parameters) and cannot conduct full-blown exploitation. Notably, in most cases, since the helper parameters are checked by the BPF verifier, the malicious programs still need to leverage the register-value-tracking vulnerabilities in the verifier (Sec. III-A) to bypass this check [17–24].

Challenge. However, protecting these BPF helpers is not trivial. First, the BPF helpers, by design, need to access BPF-related objects in the kernel memory; blindly isolating helpers leads to spurious alarms and impedes benign programs. Second, there are over 200 BPF helpers in the kernel, so MOAT's design must be generic enough to apply to most of these helpers (see Appendix A for the supported helpers).

Design Consideration and Solution. To prevent such abuse, we aim to identify and guard sensitive BPF-related objects (instead of all BPF-related objects) from the defective BPF helpers. Besides directly protecting sensitive objects, since the attackers need to deliver malformed parameters to conduct helper abuse, we also wish to ensure the validity of the helper parameters at runtime. To this end, we designed the following two defense schemes: Critical Object Protection (COP) and Dynamic Parameter Auditing (DPA). COP protects sensitive BPF-related objects from being tampered with, while DPA dynamically checks if the arguments of the helpers are within

legitimate ranges. To clarify, COP and DPA should be enabled together to deliver protection. DPA only constrains the arguments to their expected ranges. This stops most exploitation attempts [17–21] but may not ward them off completely in the presence of a buggy helper [60]; COP, in this case, prevents the buggy helper from accessing sensitive objects.

1) *Critical Object Protection (COP)*: Although BPF helpers have to access BPF-related kernel objects to complete their tasks, the sensitive BPF-related objects should still not be accessed by any helper. For example, `array_map_ops` is a function pointer in the BPF array maps that should *only* be accessed from system calls. However, `array_map_ops` is close to other helper-needed objects in the address space, making it a potential victim of the abused helpers. Based on this, we designed the COP scheme. As shown in Fig. 9, instead of treating the entire kernel domain as a whole, we divide it into a normal domain and a critical-object domain. Permissions of these critical objects are managed via an extra page and domain ID. When entering helper functions, instead of setting the entire kernel space as access-enabled (AE), those critical objects remain access-disabled (AD), preventing the helpers from accessing them. We (i) enumerate published BPF CVEs and mark the abused objects; (ii) pattern-match their structural templates (e.g., ops tables with function pointers) to surface unlabeled candidates referenced by helpers; and (iii) retain only those whose writable or leakable fields could be exploited. Although executed manually, the process can be automated: scan for pointer-dense and helper-reachable structures, then apply structural and reachability filtering. While our analysis is grounded in existing CVEs, we believe the likelihood of omitting critical objects is low, as our methodology systematically captures common patterns of exploitable objects. We identified a total of 44 critical objects (see Appendix C), including some not previously reported in existing CVEs. These objects could either leak the sensitive base address of the kernel (e.g., `iter_seq_info`) or even be tampered with to launch a full-blown exploit (e.g., `array_map_ops`). In addition to these 44 identified objects, we set MOAT itself and `cred` as critical objects. The former contains sensitive data of MOAT (e.g., the saved state of `IA32_PKRS`), while the latter tracks the privilege of a process. We believe protecting the identified critical objects provides a practical security guarantee for the BPF helpers. Nonetheless, unidentified critical objects could exist; we will discuss their potential threat in Sec. VIII. Moreover, it is always feasible to extend COP to protect other kernel objects.

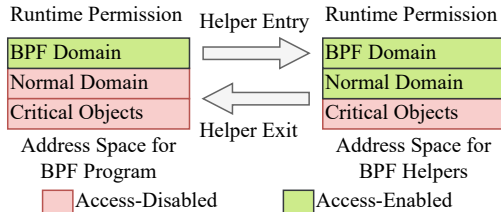


Fig. 9: Critical object protection (COP).

2) *Dynamic Parameter Auditing (DPA)*: To further regulate the helpers, we propose Dynamic Parameter Auditing (DPA), which leverages the information obtained from the BPF verifier to dynamically check if the parameters are within their

legitimate ranges and type of certain pointers. As illustrated in Fig. 10, the verifier can deduce the value range of each register via static analysis (aligned with the uncovered verifier inaccuracies [19, 22]; our DPA design tolerates even *invalidly* deduced value ranges; see clarification below). MOAT logs such value ranges and instruments the BPF programs to insert checks before helper calls. During runtime, these checks ensure that the provided parameters of the helpers are within the verifier-deduced value ranges. In our example, we can check if `r0==0x10; r1==0x11` when `BPF_HELPER` is called. If the parameter runtime values do not match with the static analysis results, the BPF program is terminated immediately.

For BPF object pointers, DPA checks if the register value falls within the valid address range determined at load time, preventing type confusion attacks (see CVE-2021-34866 in Sec. VI-A2). For kernel-shared object pointers, DPA uses verifier-recorded type info and enforces type validation before helper calls, terminating the BPF program on mismatch. Constructing a type reference map and replacing accessing with map lookups could be an alternative to refine the type checking, but it is labor-intensive and error-prone. Given the robustness of the COP scheme, we adopt the current approach for simplicity and efficiency.

Program Termination. MOAT monitors all kernel resources allocated to BPF programs (e.g., spinlocks). Upon termination, it returns a dummy value and automatically deallocates or resets associated resources, ensuring no resource leakage or residual effects remain in the system.

<code>r0 = 0x10</code>	<code>r0 = 0x10</code>	<code>r0 = 0x10</code> <code>r1 = 0xbe</code>
<code>r1 = r0 + 0x1</code>	<code>r0 = 0x10</code> <code>r1 = 0x11</code>	<code>r0 = 0x10</code> <code>r1 = 0x11</code>
<code>call BPF_HELPER</code>	<code>r0 = 0x10</code> <code>r1 = 0x11</code>	<code>r0 = 0x10</code> <code>r1 = 0x11</code>
BPF Instructions	Static Register Value Inferred by Verifier	Runtime Register Values for Each Instruction

Fig. 10: Register value tracking of the verifier.

Clarification. In this DPA strategy, one may wonder if the “value ranges” deduced by the verifier are wrong [19, 22]. To clarify this, we list possible cases of a BPF variable v ’s value range and the corresponding system states in Table II; our discussions are as follows.

TABLE II: Four cases of a BPF variable v ’s value ranges. R denotes the runtime value of v , D denotes the verifier’s deduced value of v , E denotes verifier’s *expected* legitimate value range of v , while T denotes the *ground truth* legitimate value range of v . The last column denotes this case is safe (✓), mitigated by verifier (✓_v), mitigated by MOAT (✓_M), or unsafe (✗)

	R	D	E	T	State
1	0x10	0x10	[0, 0x20]	[0, 0x20]	✓
2	0xba	0xba	[0, 0x20]	[0, 0x20]	✓ _v
3	0xba	0x10	[0, 0x20]	[0, 0x20]	✓ _M
4	0xba	0xba	[0, 0xba]	[0, 0x20]	✗

Case 1 illustrates the value range of a variable v in a benign BPF program. The runtime value aligns with verifier’s deduction which further falls within the *expected* and *true* legitimate value ranges simultaneously ($R = D \in E = T$, see the caption of Table II). Case 2 demonstrates the value range of variable v in a malformed BPF program. The runtime value 0xba is out-of-bounds, and this invalid value is detected by the verifier through static analysis. Therefore, this program is rejected by the verifier, and the system remains safe ($R = D \notin E = T$). Case 3 shows the value range

of v in a malicious BPF program. The runtime value `0xba` is out-of-bounds. However, due to the incomplete analysis caused by vulnerabilities, the verifier deduces that v 's value is `0x10`, which is within the verifier's expectation. Since DPA operates in the runtime and checks whether the runtime value actually matches the verifier's deduction, this mismatch is then detected by DPA, and this malicious BPF program is terminated ($R \neq D \in E = T$).

While the above three cases cannot be exploited, Case 4 implies a scenario where DPA fails and the helper is abused. The verifier's *expected* value range differs from the *ground truth*, legitimate value range. This discrepancy allows an out-of-bounds value `0xba` to be passed as an argument to a helper for exploitation. For this to occur, the following conditions must be satisfied simultaneously: ① The verifier has an *incorrect* expectation (i.e., $E \neq T$). ② The incorrect expectation E is *unsafe* (i.e., $T - E$ overlaps an exploitable structure). ③ The BPF program is carefully tweaked to be aligned with D and evade DPA (i.e., $R = D$). For BPF programs, it is usually straightforward for the verifier to obtain E statically (e.g., E encodes the array size). It is thus hard to satisfy ① and ② simultaneously. For today's known BPF exploits (all of which fall into Case 3), the verifier has the correct expectation $E = T$ but makes the incomplete deduction $R \neq D$; therefore, the discrepancy $E \neq T$ is never encountered in practice.

V. IMPLEMENTATION

MOAT is implemented on Linux 6.1.38. MOAT-KEY adds 2,911 LoC for Intel MPK, and MOAT-VIR adds 2,707 LoC for Arm Stage-2 translation, with about 1,000 LoC shared. The patch is modular and minimally intrusive, follows Linux multi-arch conventions (arch-agnostic core ~1K LoC; per-arch code under `/arch` with unified interfaces), easing upstreaming and deployment on mixed x86/Arm clusters. We focus on these two platforms and explain key points below.

A. Kernel Interrupt Handling

MOAT has to cooperate with many low-level routines inside the kernel. For instance, an interrupt might occur and take over the control flow during BPF programs. We introduce how MOAT handles kernel interrupt in this section.

Intel x86. Note that most interrupt handlers require access to kernel memory, and as a result, the MPK would presumably raise spurious alerts. Thus, we need to temporarily disable MPK inside these handlers and re-enable it once the handlers finish. To avoid the overhead when there is no BPF program, we use a per-CPU variable `in_bpf` to identify whether the processor is executing a BPF program. Since BPF programs only occupy a tiny fraction of kernel execution time, we observe little performance loss due to this, even under cases where interrupt frequently occurs (e.g., intensive network activity in Sec. VI-D).

Arm. MOAT-VIR, like MOAT-KEY, must temporarily disable isolation (i.e., restore the original Stage-2 page table) for interrupt handling, using a per-CPU variable. MOAT-VIR relies on Stage-2 translation, which is managed in the hypervisor-space. To efficiently switch the Stage-2 page table, MOAT-VIR sets

the IMO and FMO bits in the `HCR_EL2` register, causing the hypervisor-space to intercept physical interrupts. It then restores the original Stage-2 page table to give the handler full kernel access, and revert to the dedicated one once execution returns to the BPF program. These bits are set only when a BPF program begins execution, ensuring that intervention occurs only when necessary.

B. Granularity of MOAT

As both Intel MPK and Arm Stage-2 are based on the page table, MOAT only protects memory in the granularity of a page (i.e., 4 KB). However, the objects used by BPF programs may not be aligned to 4 KB, which means they could interleave with critical kernel structures. Therefore, granting BPF programs access to these objects also enables access to those kernel structures and leads to exploitation. To prevent this, we have modified BPF-related objects (e.g., maps) so that they are page-aligned and not interleaved with other structures.

C. Instrumentation

DPA Check Generation. To deploy DPA from Sec. IV-C2, we modify the BPF JIT compiler (`bpf_jit_comp.c` with about 2,500 LoC) to instrument BPF programs. As shown in Fig. 11, our modified JIT compiler receives a set of expected ranges types of certain pointers from the verifier. Then, for each parameter, the JIT compiler emits assembly instructions to check whether the parameter is within the expected range mismatches the type. If not, we terminate the program (bad label in Fig. 11). This prevents malicious programs from passing invalid parameters to abuse BPF helpers.

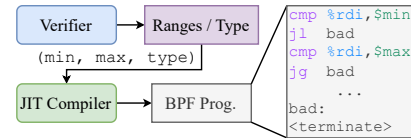


Fig. 11: DPA Check Generation.

Secure Call Gate. To enforce isolation, MOAT relies on hardware isolation primitives when executing BPF programs. However, BPF programs may invoke helper functions and tailcalls to other BPF programs, requiring a control flow transition. To ensure the security of this transition, MOAT first checks whether the called helper is allowed based on the program type and whether the target tailcall is within the BPF domain. Then, MOAT instruments BPF programs using the JIT compiler. Specifically, MOAT-KEY implements a secure call gate by inserting MPK setting and address space switching instructions before and after each helper function call. Additionally, MOAT-VIR must switch the Stage-2 page table to enforce isolation. However, since the base address of the Stage-2 translation is stored in the `VTTBR_EL2` register that is accessible only in hypervisor-space, this transition requires privilege escalation. To achieve this, MOAT-VIR instruments the HVC instruction, enabling a controlled trap to hypervisor-space.

VI. EVALUATION

To evaluate MOAT, we first analyze how MOAT mitigates various attack interfaces and then benchmark its CVE detectability in Sec. VI-A. We then assess the performance of MOAT under different BPF program setups in Sec. VI-C and Sec. VI-D. As we implement MOAT with Intel MPK and Arm Stage-2 translation, our evaluation focuses on these platforms.

A. Security Evaluation

1) *Analysis of Attack Mitigation:* We systematically analyze how MOAT mitigates the representative attacks on the BPF ecosystem as well as the potential threats to MOAT itself. Our analyses are illustrated in Fig. 12.

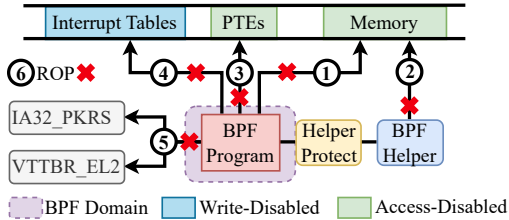


Fig. 12: Analysis of attack mitigation.

① **Arbitrary Kernel Access.** Currently, the most prevalent threat to the BPF ecosystem is the ability of malicious BPF programs to arbitrarily modify kernel memory. In order to accomplish this, these BPF programs typically employ corner-case operations to deceive the verifier during the loading phase and to behave maliciously during runtime. This type of attack is effectively mitigated due to the fact that MOAT derives the necessary memory regions of each BPF program and uses hardware features to prevent any runtime access beyond this domain (Sec. IV-B), mitigating such illegal access.

② **Helper Function Abuse.** Apart from launching an attack directly from BPF programs, a malicious BPF program may carefully prepare parameter values and pass them to abuse certain helpers. To prevent such abuse, MOAT deploys security enforcement schemes (Sec. IV-C) to dynamically audit helper parameters and also protect critical kernel objects during the execution of these helpers. Thus, the attacker can no longer take advantage of these helpers.

③ **PTE Corruption.** On Intel platforms, a page's MPK region is configured via its PTE. Consequently, a malicious BPF program may attempt to tamper with the PTEs to disable MOAT-KEY. However, this is impossible since MOAT-KEY sets these PTEs as access-disabled; they are thus protected by MPK like other kernel resources. On Arm platforms, the isolation is enforced through dedicated Stage-2 page tables. A malicious BPF program may attempt to modify the PTEs of the kernel memory to bypass isolation. However, the Stage-2 page tables are allocated within the hypervisor-space, making them inaccessible from kernel-space, where BPF programs execute.

④ **Interrupt Tables Tampering.** Interrupt tables like Intel GDT/IDT and Arm vector tables are essential for segmentation and interrupt handling. Therefore, blindly setting them as access-disabled would cause system crashes. However, since these tables are only accessed in a read-only manner, MOAT sets them as write-disabled, thus preventing malicious BPF programs from using them to compromise the kernel.

⑤ **Hardware Register Tampering.** Besides memory-based attacks, on Intel platforms attackers may also directly disable MPK through hardware configurations. As described in Sec. II-B, `IA32_PKRS` is a critical register for configuring MPK. One may disable MPK by modifying `IA32_PKRS`. However, this register can only be modified via special instructions, and BPF instruction sets do not include any of these. Thus, a BPF program with these instructions is rejected immediately. On Arm platforms, the base address of the Stage-2 page table is stored in the `VTTBR_EL2` register. However, this register is only accessible from the hypervisor-space; from the kernel-space, it can be modified only by invoking a hypercall via the `HVC` instruction. Since the BPF programs are set to $W \oplus X$ (meaning write and executable permissions are not simultaneously enabled), adding these instructions via self-modification is also impossible. Additionally, MOAT-VIR ensures that the hypercall handler code is not mapped in the BPF domain's Stage-2 page table. As a result, any unauthorized hypercall attempt triggers a Stage-2 instruction ABORT.

⑥ **Return-Oriented Programming.** Two properties of the BPF instruction set prevent potential control-flow hijacking attacks like return-oriented programming (ROP). First, BPF only supports jump instructions with *constant and instruction-level* offsets. This means the destinations of jumps are trivially known during the compile time, and there are no *unintended ROP gadgets* (jumps between instructions) like x86 [61]. Secondly, as a specialized instruction set, BPF does not include any instructions that may modify hardware configurations. These two properties allow MOAT to reliably detect invalid instructions and prevent BPF programs from tampering with hardware settings.

2) *Real-world CVE Evaluation:* We surveyed the BPF CVEs in the past ten years. A total of 26 CVEs are memory exploits (Appendix B) and thus fall within the scope of MOAT. We tested MOAT's effectiveness on all of these CVEs. For CVEs with publicly available PoC, we ported and ran the PoC on MOAT-enabled kernel. For CVEs without PoC, we studied the fixes and ensure that MOAT mitigates them. In sum, we report that MOAT successfully mitigates *all* of them. We now present the following case studies.

CVE Case Study. To better explain how MOAT mitigates these CVEs, we elaborate on the exploit paths for three of them, 2022-23222, 2020-27194, and 2021-34866.

CVE-2022-23222 is a pointer mismatch vulnerability introduced via a rather new BPF feature, `bpf_ringbuf`. This new feature was brought to BPF in 2020, along with a new pointer type named `PTR_TO_MEM_OR_NULL`. However, the verifier had not been updated to track the bounds of this new type, resulting in this vulnerability. As shown in Fig. 13a, the malicious payload first retrieves a `nullptr` via `ringbuf_reserve` (line 1), which returns this newly added pointer type named `PTR_TO_MEM_OR_NULL`. Since this new type is not tracked by the verifier, the payload can bypass pointer checks by tricking the verifier that `r1` is `0x0` when it is `0x1` (line 3). `r1` can then be multiplied with any offset to perform arbitrary kernel access (line 6). However, such access violates the Layer-I isolation by MPK or causes a Stage-2 data ABORT, and is terminated by MOAT (line 7).

CVE-2020-27194 is a vulnerability due to incorrect truncation. As in Fig. 13b, the user first inputs an arbitrary value in the range of $[0, 0 \times 600000001]$ (line 1). Then, the conditional clause helps the verifier to determine its value range (line 3). However, when tracking the BPF_OR operator, the verifier performs a wrong truncation on its upper bound. After the truncation, the user-controlled `r5` is viewed by the verifier as a legitimate constant 0×1 (line 5), which is later used as the offset to perform arbitrary access to the kernel (line 6). Similarly, such access is stopped by MOAT.

CVE-2021-34866 is a helper-abuse vulnerability. As shown in Fig. 13c, the malicious payload tries to pass an invalid map to the `ringbuf_reserve` to cause heap overflow (line 3). However, since the runtime value of `r5` does not match the argument of `ringbuf_reserve`, DPA prevents such a mismatched helper call (line 2). Moreover, supposing that the DPA is not enabled, and the helper tries to tamper with exploitable kernel objects (e.g., `array_map_ops`). COP protects these objects and thus prevents the helper from accessing them (line 3). Lastly, even if neither COP nor DPA is enabled, and the abused helper manages to return a leaked kernel pointer, accessing the leaked pointer violates the Layer-I isolation, and the malicious program is terminated by MOAT (line 4).

```

1 r0 = ringbuf_reserve(fd, INT_MAX, 0)
2 r1 = r0 + 1 // R:r0=0;r1=1 V:r0=r1=?
3 if (r0 != nullptr) // R:r0=0;r1=1 V:r0=r1=?
4   exit(1)
5 off = <bad off> // R:r0=0;r1=1 V:r0=r1=0
6 off = off * r1 // R:off=<bad off> V:off=0
7 *(ptr+off) = 0xbad // Domain violation

```

(a) Code snippet of CVE-2022-23222

```

1 r5 = <bad addr>
2 r6 = 0x600000002
3 if (r5>=r6||r5<=0) // R&V:0x1<=r5<=0x600000001
4   exit(1)
5 r5 = r5 | 0 // R:r5=<bad addr> V: r5=0x1
6 *(ptr+r5)=0xbad // Domain violation

```

(b) Code snippet of CVE-2020-27194

```

1 r5 = <bad map fd>
2 <DPA checks> // DPA violation
3 r0=ringbuf_reserve(r5, INT_MAX, 0) // COP-guarded
4 *(r0+ptr_off) = 0xbad // Domain violation

```

(c) Code snippet of CVE-2021-34866

Fig. 13: CVE case study. R denotes variable runtime statuses. V denotes verifier-deduced values of variables.

B. Performance Evaluation Setup

We evaluate MOAT-KEY on a 5-core Intel 8505 processor with MPK support, and MOAT-VIR on a Raspberry Pi 4B with a 4-core Cortex-A72 processor, on Linux v6.1.38. To reduce variance, hyper-threading, turbo-boost, and frequency scaling are disabled. All evaluated BPF programs are executed in JIT mode, given that BPF JIT is enabled by default on all supported platforms. Moreover, both COP and DPA (Sec. IV-C) are enabled; COP is configured to protect the critical objects identified in Sec. IV-C1. We manually inspected the CPU

utilization to ensure it is close to 100%, and that the overhead is not hidden by the increased CPU load.

C. Micro Benchmark

For the micro benchmark, we measure the CPU cycles of four key operations in MOAT on both Intel and Arm platforms, as outlined in Table III. On the Intel platform, `set_pkrs` changes region permissions by modifying `IA32_PKRS` via `WRMSR`. `get_pkrs` reads the current permission configuration in `IA32_PKRS` via `RDMSR`. On the Arm platform, `hypercall` traps to hypervisor-space. `stage2_switch` is the cost of switching the Stage-2 page table, including a Stage-2 TLB flush. On both platforms, `bpf_{entry/exit}` is the total cost of entering/exiting a BPF program, which includes operations like swapping stack, managing BPF context, and configuring region permissions with `set_pkrs` or using `hypercall` to perform `stage2_switch`. `dpa_check_args` is the cost of checking helper parameters. Each operation is measured by averaging ten runs of one million invocations to eliminate randomness. Given Intel’s “performance core” and “efficient core”, we measure their cycles independently.

As shown in Table III, on the tested Intel platform, the overall switching cost of MOAT-KEY is under 200 cycles, which is negligible for most BPF programs (see Sec. VI-D for details). Notably, setting and getting the region permissions (`set_pkrs/get_pkrs`) in PKS is more expensive than its user-space variant in `libmpk` [39]. We presume that this is because, in PKU, the region permission is controlled via a dedicated register named `PKRU` with two special instructions `RDPKRU`/`WRPKRU` (user-space `RDPKRU`, `WRPKRU` take 0.5 and 23.3 cycles, respectively), whereas MOAT-KEY’s PKS region permission is stored in an MSR named `IA32_PKRS` without any special instruction. To configure the permission in `IA32_PKRS`, one has to use the `RDMSR`/`WRMSR` instructions with the MSR ID `0x6E1`. We observe that the operations of MOAT-KEY are not substantially affected by the difference between performance and efficient cores. On the tested Arm platform, a `hypercall` costs approximately 900 cycles, while a `stage2_switch` operation takes around 600 cycles. Since `bpf_{entry/exit}` involves two `hypercall` invocations and two `stage2_switch` operations, the total overhead is approximately 2,300 cycles. In real-world use cases (see Sec. VI-D), this overhead is negligible for most BPF programs. Lastly, while the cost of `dpa_check_args` varies depending on the checked range type (e.g., a value point $[0 \times 1, 0 \times 1]$ costs less than a value range $[0 \times 1, 0 \times 10]$), we report that these costs are all under ten cycles.

TABLE III: Micro benchmark results comparing Intel and Arm. We use P to denote the cycles of performance cores and E for efficient cores (all cores of the test Arm device are the same).

Intel Operation	Intel #Cycles		Arm Operation	Arm #Cycles
	P	E		
<code>get_pkrs/RDMSR</code>	36	43	<code>hypercall</code>	868
<code>set_pkrs/WRMSR</code>	111	112	<code>stage2_switch</code>	578
<code>bpf_{entry/exit}</code>	154	173	<code>bpf_{entry/exit}</code>	2311
<code>dpa_check_args</code>	≤ 10	≤ 10	<code>dpa_check_args</code>	≤ 10

MOAT’s Overhead vs. #BPF Programs. To show MOAT supports numerous BPF programs, we prepare the following

experiments. We attach 1, 10, 32, 64, and 128 BPF programs to trace `execve`³, run a program that continuously creates processes for one minute, and measure the number of processes created. In this setting, each invocation of `execve` stresses MOAT to constantly switch between the BPF programs. Moreover, we craft each BPF program as succinct (programs that directly return) so that MOAT’s relative overhead is not “dominated” by the overly lengthy BPF programs.

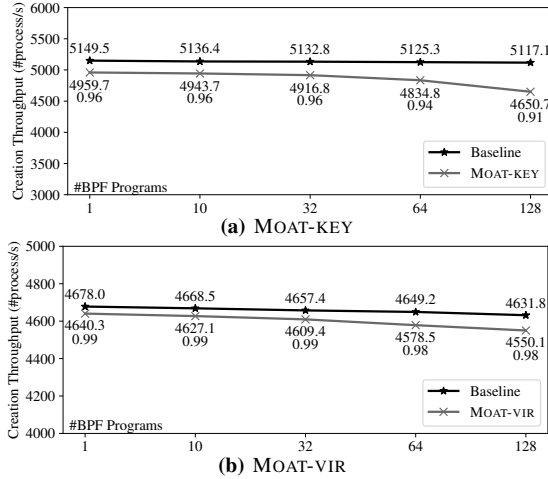


Fig. 14: MOAT’s overhead with respect to #BPF programs.

We report our results in Fig. 14. Since we use simple BPF programs, the baseline performance (without MOAT) is not observably affected by the number of BPF programs; we expect that in real-world cases, the baseline performance would also drop as the number of BPF programs increases. MOAT’s overhead stays largely the same (4% on MOAT-KEY and 1% on MOAT-VIR) when there are 1, 10, or 32 BPF programs and then slowly increases with the number of BPF programs. It eventually incurs 9% overhead with 128 BPF programs on MOAT-KEY and 2% on MOAT-VIR. With further inspection, we find that over 64 BPF programs on MOAT-KEY with isolated TLB entries pose heavy stress to TLB, resulting in increased overhead. However, having over 64 BPF programs attached to the same place is extremely rare (if it occurs at all). Nonetheless, even in such cases, MOAT-KEY incurs a performance penalty of less than 10%. And on MOAT-VIR, the VMID helps to reduce TLB flushing stress, succinct BPF programs introduce negligible overhead.

To complement the above experiment, we prepare the following experiment where BPF programs are attached to different kernel locations. There are, in total, 685 BPF tracepoints of system calls in Linux [63]. Following a similar setting as above, we attach each of these tracepoints with a simple BPF program and run UnixBench [64] to measure the overall system performance. In sum, there are nearly 700 BPF programs in the kernel with diverse invocation patterns. Thus, this setting stresses the system in a manner distinct from the previous experiment. We report that in such settings, the incurred overhead of both MOAT-KEY and MOAT-VIR are about 5%.

³We clarify that one BPF tracepoint only supports up to 64 programs [62], so we attach to both entry and exit tracepoints of `execve` in these experiments.

From these two experiments, we interpret that the overhead of MOAT is slightly affected by the number of BPF programs due to the TLB stress. Nevertheless, MOAT’s overhead falls in a reasonable and promising range (4%~9% on MOAT-KEY and 1%~2% on MOAT-VIR), even for cases that are much heavier and rarely seen in real-world ones.

D. Macro Benchmark

For the macro benchmark, we set up three mainstream BPF use cases: network, system tracing, and system-call filtering. On the network cases (i.e., Socket and XDP), we use the smallest applicable packet size because BPF operates on each packet. Thus, under the same bandwidth, smaller-sized packets will incur higher throughput (i.e., more packets) and lead to more invocations of BPF programs, eventually putting more stress on MOAT.

Network — Socket. To evaluate MOAT’s overhead on the network applications, we simulate a traffic monitoring scenario. A traffic generator sends UDP packets for one minute, with a packet size of 16 bytes, to our tested device. A server on the tested device receives these packets. Both the sender and receiver have a 1 GbE network interface controller. As for the tested BPF programs, we use five socket filtering programs from the Linux source tree, similar to previous works [65]:

- `drop`: directly ignores the packet.
- `byte`: monitors the traffic in bytes from each protocol.
- `pkt`: monitors the traffic in packets from each protocol.
- `trim`: only keeps the packet header to the socket.
- `flow`: monitors the network traffic based on protocol, interface, source, destination, and port.

We attach a socket to the receiver’s network interface and set up these BPF programs to monitor the network traffic over the socket. In addition to the evaluation of each program, we also conduct a full-on experiment where five BPF programs are attached simultaneously to stress MOAT.

TABLE IV: MOAT’s traffic monitoring performance in Thousand Packets per Second (TPPS). The full-on experiment is denoted as “all”. The “vanilla” throughput without BPF program is 596.3 (Intel) / 244.19 (Arm) TPPS; the relative throughput is denoted in parentheses, e.g., (99.73%).

Throughput (TPPS)	drop	byte	pkt	trim	flow	all
Baseline-Intel	594.39 (99.70%)	594.67 (99.73%)	594.26 (99.66%)	594.74 (99.73%)	594.39 (99.68%)	587.22 (98.47%)
MOAT-KEY	593.10 (99.46%)	594.31 (99.66%)	594.43 (99.68%)	594.69 (99.73%)	593.10 (99.46%)	575.33 (96.48%)
Baseline-Arm	224.29 (91.85%)	221.40 (90.67%)	220.69 (90.38%)	212.71 (87.11%)	211.33 (86.54%)	170.80 (69.95%)
MOAT-VIR	219.36 (89.83%)	178.47 (73.09%)	178.49 (73.10%)	175.46 (71.86%)	108.45 (44.41%)	74.40 (30.47%)

As shown in Table IV, MOAT-KEY incurs negligible overhead (<1%) for all BPF programs if they are solely executing in the kernel. Even in the full-on experiment, which forces MOAT-KEY to constantly switch between these BPF programs, MOAT-KEY brings only a very small throughput drop of 2%. However, on the tested Arm platform, MOAT-VIR’s overhead scales with the baseline performance cost of standard BPF programs. This scaling is due to the limited resources (e.g., CPU) of the Raspberry Pi 4B, which amplifies the relative overhead of frequent interrupts. Higher-end Arm devices are expected to exhibit lower overhead. In practical scenarios such as Apache and Nginx web servers (see Table IX), this

overhead remains moderate, typically below 5%, ensuring minimal impact on real-world deployments.

Network — XDP. Besides processing packets from a socket buffer, BPF provides a direct way to control the network — eXpress Data Path (XDP). XDP processes packets at an early stage in the network stack to achieve fast packet processing. Following the settings in the socket experiment, we simulate a packet processing scenario. Similar to prior works [65], we run five XDP programs from the Linux source tree:

- `xdp1`: parses the IP header, keeps packets count in a BPF map, and drops the packets.
- `xdp2`: same as `xdp1`, but re-sends the packets.
- `adj`: trims the packets into ICMP packets, sends them back, and keeps packet count in a BPF map.
- `rxq1`: counts and drops the packets in each receive queue.
- `rxq2`: same as `rxq1`, but re-sends the packets.

Unlike socket filters, XDP programs require packets above a certain size, so we tune our traffic generator to send packets that go through the maximum possible execution path of the tested programs. We send packets of 64 bytes for `xdp1` and `xdp2` and packets of 100 bytes for `adj`, `rxq1`, and `rxq2`.

TABLE V: MOAT’s XDP performance in TPPS. The “vanilla” throughput without XDP program is 532.9 (Intel) / 213.4 (Arm) TPPS with 100-byte packets, and 561.5 (Intel) / 222.5 (Arm) TPPS with 64-byte packets; the relative throughput is denoted in parentheses, e.g., (99.55%).

Throughput (TPPS)	<code>xdp1</code>	<code>xdp2</code>	<code>adj</code>	<code>rxq1</code>	<code>rxq2</code>
Baseline-Intel	560.58 (99.84%)	557.78 (99.34%)	531.11 (99.66%)	528.36 (99.15%)	530.52 (99.55%)
MOAT-KEY	560.15 (99.76%)	557.76 (99.33%)	530.65 (99.58%)	527.57 (99.00%)	527.66 (99.05%)
Baseline-Arm	218.80 (98.33%)	179.52 (80.68%)	181.52 (85.06%)	207.56 (97.26%)	187.60 (87.91%)
MOAT-VIR	218.20 (98.06%)	168.45 (75.70%)	137.67 (64.51%)	205.65 (96.36%)	125.37 (58.75%)

As illustrated in Table V, MOAT-KEY incurs negligible performance penalties (<1%) when executing XDP programs. MOAT-VIR exhibits similar scaling behavior in the XDP test, where its overhead remains proportional to the baseline overhead of standard BPF programs.

System Tracing. System tracing is another mainstream BPF use case. To evaluate MOAT’s overhead on system tracing, we prepare 11 BPF programs to trace frequent system events like process creation, context switch, and file operations. These programs collect relevant system statistics for user-space analysis. Then, we run UnixBench [64] to measure the overall system performance. UnixBench includes the following tests: ① `execl` throughput, ② file copy, ③ pipe throughput, ④ pipe-based context switching, ⑤ process creation, ⑥ shell scripts, and ⑦ system call.

We report the results in Fig. 15. We find that both MOAT-KEY and MOAT-VIR impose a small slowdown ($\leq 8\%$) for most UnixBench tests. The maximum performance loss brought by MOAT-KEY is 13% in test ⑦ and 23% by MOAT-VIR in test `fc ②`. Such overhead seems moderate. However, note that the BPF programs without MOAT already bring a non-trivial performance penalty (e.g., 63% slowdown in test ③ of MOAT-KEY). Therefore, the performance loss brought by MOAT is reasonably low for system tracing.

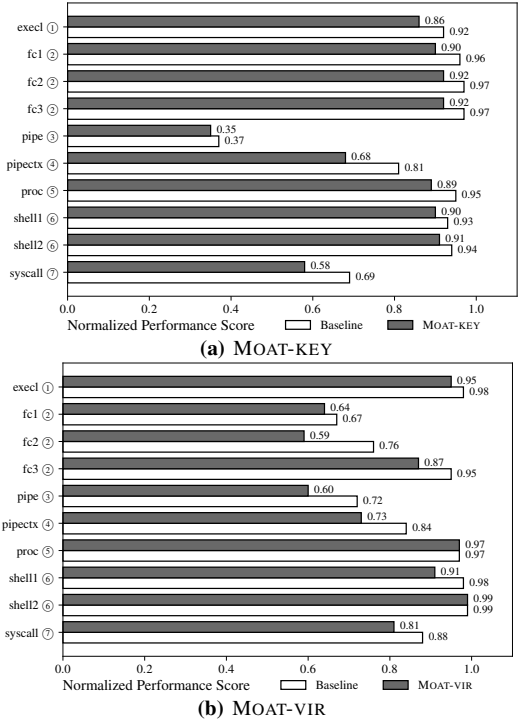


Fig. 15: UnixBench normalized scores with respect to the “vanilla” scores without BPF tracepoints. The “`fc*` ②” and “`shell*` ⑥” are file copy tests and shell tests with different settings.

Syscall Filtering. BPF is also used to enhance software security [66–68]. `seccomp`-BPF allows filtering the system calls of a process with BPF. `sysfilter` [67] is an automated tool that analyzes a program, creates the set of system calls the program needs, and restricts the program using `seccomp`-BPF. Note that `sysfilter` is only available on x86-64 platforms, so we only evaluate MOAT-KEY here. We use `seccomp`-BPF and `sysfilter` to evaluate MOAT-KEY’s overhead in such a use case. We apply the MOAT-hardened filter to Nginx and benchmark it using `wrk` [69] with one, two, and three client processes; each sends requests for one minute with 128 connections. Nginx is configured with the same number of worker processes. To clarify this setting: the number of processes for Nginx typically shall not exceed the number of cores, and adding more processes does not increase the throughput due to the context-switch cost. All requests are sent over the loopback (`10`) to minimize network interference.

TABLE VI: Nginx throughput in Thousands of Request per Second (Treq/s). The relative throughput is in parentheses, e.g., (95.6%).

Throughput (Treq/s)	1 worker	2 worker	3 worker
Vanilla	148.1 (100%)	179.5 (100%)	165.2 (100%)
(no <code>seccomp</code> -BPF)	± 12.81	± 8.35	± 4.72
Baseline	147.2 (99.4%)	171.3 (95.4%)	160.5 (97.2%)
	± 9.56	± 8.08	± 5.28
MOAT-KEY	142.3 (96.1%)	166.3 (92.6%)	158.0 (95.6%)
	± 8.77	± 6.70	± 4.48

As shown in Table VI, MOAT-KEY incurs 3% additional throughput drop. Moreover, the standard deviations of throughput are within the normal range. Therefore, MOAT-KEY does not introduce fluctuation to Nginx throughput.

E. Additional Evaluation

In addition to the above experiments, we also evaluate MOAT’s memory footprints, instrumentation cost from DPA, and compare MOAT’s performance with a prior work [70].

Memory Footprint. As mentioned in Sec. V, MOAT aligns BPF-related objects (e.g., maps) to 4 KB to ensure that they do not interleave with other kernel structures, introducing extra memory footprints. We provide a detailed breakdown of MOAT’s memory footprints in Table VII. Specifically, MOAT uses four pages to set up the page table of isolated address-spaces, one page for the stack and one for the context. As for critical objects identified in Sec. IV-C1, MOAT uses an extra page to toggle their permissions independently. Additionally, the memory used by the BPF program binaries and BPF maps is aligned up to a multiple of page size.

TABLE VII: Breakdown of MOAT’s memory footprint. AS: Address Space; ST: Stack; Ctx: Context; CO: Critical Objects. P and M denote #pages of program and map, respectively.

Type	AS	ST	CO	Ctx	Prog	Map
#Page	4	1	1	1	P	M

Though MOAT’s memory footprints seem non-trivial from Table VII, we clarify that they are static and thus do not grow dynamically during the runtime. Even if there are thousands of BPF programs, MOAT’s memory footprints are merely a few megabytes, which is negligible for modern systems.

Helper Instrumentation Cost. As described in Sec. IV-C2, MOAT instruments BPF programs to insert DPA checks. Table VIII shows the number of helper calls and memory access made by the BPF programs⁴ in Sec. VI-D. The former reflects the number of DPA checks inserted by MOAT, while the latter reflects the number of checks from SFI-based solutions. We report that, in most cases, the number of DPA checks is smaller than that of SFI-based solutions, let alone that SFI only offers memory safety, which is offered by Layer-I isolation in MOAT.

TABLE VIII: The number of checks of DPA and SFI-based isolation. Note that the tracepoint (marked with *) consists of multiple BPF programs, and we report their average number of inserted checks.

Name	drop	byte	pkt	trim	flow	xdp1	xdp2	adj	rxq	tracepoint*
#Helper	0	1	1	0	2	3	3	5	3	6.4
#Mem.	0	4	3	1	61	16	29	53	37	6.9

Comparison with SandBPF. In this section, we compare MOAT with SandBPF [70]. SandBPF enforces isolation by inserting software runtime checks into the memory access of BPF programs, which generally leads to higher overhead compared to hardware-based methods. To validate this, we conduct the following direct comparative study. The authors of SandBPF conducted an evaluation with Phoronix Test Suite [31]; we reproduce the same experiments on MOAT. Note that the source code of SandBPF is not public at the time of writing, so we directly refer to their data in Table IX.

As shown in Table IX, MOAT’s overhead is lower than SandBPF in most testcases (Rel. v.s. Ref.). Again, we interpret the advantage of MOAT is attributed to the reduced number of inserted checks compared to SFI-based approaches.

VII. RELATED WORK

In this section, we compare MOAT’s design with other works in kernel isolation. This helps highlight the contribution of MOAT, in comparison to previous research.

⁴The BPF programs in syscall filtering scenario are cBPF programs and thus do not support helper functions; we do not include them here.

TABLE IX: Comparison with SandBPF [70]. We provide MOAT’s relative overhead (Rel.) and SandBPF’s overhead (Ref.).

(a) MOAT-KEY										
Test #Conn (req./s)	XDP				Socket Filter					
	Base	MOAT-KEY	Rel.	Ref.	Base	MOAT-KEY	Rel.	Ref.		
Apache 20	34,303	33,689	2%	0%	40,666	40,286	1%	4%		
Apache 100	31,929	30,726	4%	8%	37,998	36,546	4%	4%		
Apache 200	27,751	26,657	4%	5%	32,652	31,344	4%	3%		
Apache 500	24,786	24,439	1%	7%	30,262	29,423	3%	7%		
Apache 1000	24,597	24,470	1%	6%	29,545	28,961	2%	7%		
Nginx 20	22,688	21,892	3%	7%	23,359	23,530	0%	10%		
Nginx 100	21,492	20,689	4%	7%	22,870	22,482	2%	8%		
Nginx 200	19,972	19,216	4%	6%	21,562	20,984	3%	8%		
Nginx 500	18,470	17,814	4%	6%	19,421	18,713	4%	7%		
Nginx 1000	17,024	16,735	2%	3%	17,392	17,098	2%	6%		

(b) MOAT-VIR										
Test #Conn (req./s)	XDP				Socket Filter					
	Base	MOAT-VIR	Rel.	Ref.	Base	MOAT-VIR	Rel.	Ref.		
Apache 20	2,185	2,024	7%	0%	2,855	2,669	6%	4%		
Apache 100	2,284	2,131	7%	8%	3,027	2,834	4%	4%		
Apache 200	2,249	2,134	5%	5%	3,220	3,151	2%	3%		
Apache 500	2,279	2,209	3%	7%	3,214	3,071	4%	7%		
Apache 1000	2,372	2,298	1%	6%	3,279	3,126	5%	7%		
Nginx 20	426.0	409.6	4%	7%	430.4	430.3	0%	10%		
Nginx 100	418.5	410.6	2%	7%	430.3	428.5	0%	8%		
Nginx 200	405.2	395.5	2%	6%	411.2	407.6	1%	8%		
Nginx 500	376.8	370.9	2%	6%	388.4	374.3	4%	7%		
Nginx 1000	316.8	308.4	3%	3%	328.7	315.8	4%	6%		

MPK-based Isolation. Prior to PKS, Intel first announced its user-space variant PKU. Consequently, most existing works [39, 58, 71] using MPK focus on user-space isolation. To better utilize PKU as an isolation primitive, Park et al. [39] proposed libmpk, which resolves the semantic discrepancies between PKU and conventional mprotect. VDom and EPK [58, 71] aim to provide unlimited keys in the *user space* via key virtualization. Despite the similarity, we clarify that isolating BPF programs in the *kernel* is a distinct scenario and comes with its own challenges. This is why we have proposed the lightweight two-layer design to efficiently isolate BPF programs. Apart from using PKU to isolate user applications, efforts are made to isolate trusted applications in SGX via PKU [37, 38]. SGXLock [37] establishes mutual distrust between the kernel and the trusted SGX applications, while EnclaveDom [38] enables intra-isolation within one enclave. PKU has also been used for kernel security. IskiOS [72] applies PKU to kernel pages by marking them as user-owned.

Virtualization-based Isolation. There is a line of research works on isolating kernel components via virtualization [73–75]. However, among these prior works, lightweight solutions like SKEE [75] are not compatible with the low-level routines (e.g., interrupt) in Linux. SKEE disables the interrupt upon entry, but disabling the interrupt will significantly impede BPF’s network performance. To incorporate with low-level routines in the kernel, non-trivial modification to the system is often needed. For example, LVD and LXDs [73, 74] require a hypervisor and an implanted micro-kernel to manage the isolated components. On the one hand, MOAT-VIR utilizes Arm interrupt route to effectively switches domains upon interrupts, ensuring minimal overhead in interrupt handling. On the other hand, MOAT-VIR only implements a tiny codebase to switch Stage-2 page tables, and achieves intra-BPF isolation by simply splitting the BPF domain’s Stage-2 page table, maintaining isolation with minimal complexity.

SFI-based Isolation. Prior works also proposed Software Fault Injection (SFI) for kernel security [70, 76]. SFI inserts checks before memory access to ensure that they fall into valid ranges. However, inserting checks for memory access

often brings higher overheads (see Sec. VI-E for an empirical comparison with MOAT). MOAT uses hardware isolation primitives to ensure memory safety and only inserts checks before helpers. Since the number of helper calls is much smaller than that of memory access, this design largely reduces the overhead of MOAT.

BPF Security. There are prior works [9–13] on securing the BPF ecosystem. Most of them use formal methods to enhance the following BPF components: the verifier, the JIT compiler, or the BPF program. Huang et al. [77] systematically analyze memory safety risks in the eBPF ecosystem, highlighting limitations of the verifier and current kernel defenses. To enhance the BPF verifier, Gershuni et al. [11] propose PREVAIL based on abstract interpretation, which supports more program structures (e.g., loops) and is more efficient than the standard verifier. PRSafe [12], on the other hand, designs a new domain-specific language, whose ultimate goal is to build a mathematically verifiable compiler for BPF programs. Jitk [10] offers a JIT compiler whose correctness is proven manually, and Nelson et al. [9] generate automated proof for real-world BPF JIT compilers. Lastly, Luke Nelson [13] builds proof-carrying BPF programs, requiring developers to provide a correctness proof with the program before loading it into the kernel. HIVE [78] is another BPF isolation solution on Arm. It treats BPF programs as kernel-mode applications and establishes an isolated environment for them. HIVE uses load/store unprivileged (LSU) instructions to de-privilege memory accesses in BPF programs, and pointer authentication (PAC) for access control of kernel objects. By leveraging Arm-specific hardware features, HIVE achieves fine-grained isolation with low overhead. In contrast, MOAT provides a cross-platform isolation scheme for BPF programs, utilizing the mature hardware features available on each platform.

VIII. LIMITATIONS

In this section, we discuss the limitations of MOAT, including the unidentified critical objects, issues brought by the granularity of MPK and Stage-2 page table, and potential barriers to deploying MOAT.

Unidentified Critical Objects. In Sec. IV-C1, we manually identified the critical objects in the BPF subsystem that have been exploited in the wild or similar to the ones that have been exploited. There might still exist unidentified critical objects, however, to exploit these unidentified objects, one still has to find a BPF helper that can be abused to access these critical objects and bypass the parameter checks (i.e., DPA) enforced by MOAT. Thus, even if few unidentified critical objects exist, they would be hard, if at all possible, to exploit due to DPA.

Page-size Granularity. MOAT leverages MPK and Arm Stage-2 translation to enforce isolation, which only supports 4 KB granularity. Though we carefully adjust BPF-related objects so that they are page-aligned, there still exist a few corner-cases where MOAT does not apply. In particular, BPF programs shall only access certain fields of the context `sk_buff`, which is enforced via the static checks by the verifier. However, applying to these fields would significantly bloat `sk_buff` due to the granularity. MOAT thus cannot

constrain the access of BPF programs. As a result, `sk_buff` might have some bits leaked to BPF programs if the verifier's static checks are bypassed. Fortunately, BPF programs only receive a copy of `sk_buff` and cannot tamper with the original structure. Therefore, the consequence of such leakage, per our observation, seems trivial.

Barrier to Deploying MOAT. As we mentioned in Sec. VI-E, MOAT introduces additional memory footprints. Though the introduced footprints are only a few pages and negligible for modern systems, it might still create obstacles for some embedded systems, where memory is a scarce resource. Third, MOAT comprises approximately 3,000 LoC for MOAT-KEY using Intel MPK, and around 2,700 LoC for MOAT-VIR using Arm Stage-2 translation. About 1,000 LoC are architecture-agnostic, requiring moderate engineering effort to port MOAT to other platforms. We implemented MOAT under Linux multi-architecture convention to facilitate porting.

IX. DISCUSSION

Adapting to New BPF Features. BPF frequently introduces new helpers, map types, and program types. Since new program types follow the standard BPF framework, MOAT can support them with minimal changes. For new map types, MOAT's extensible design allows allocating the data section of map elements (i.e., excluding metadata) in the BPF domain without major modifications. To support new helper functions, MOAT can extend the COP configuration to cover new critical objects, following the methodology in Sec. IV-C1.

Cross-Platform Portability. MOAT is designed for portability across platforms by leveraging two classes of hardware isolation primitives: key-based and virtualization-based, corresponding to MOAT-KEY and MOAT-VIR. For MOAT-KEY, key-based primitives are used; on Intel, MPK enforces isolation, while on Arm, POE provides similar isolation, though POE is limited to Armv9.4 and newer. For MOAT-VIR, virtualization-based primitives are utilized and are broadly available: on Arm, Stage-2 translation isolates memory; on x86, Intel's Extended Page Table (EPT), AMD's RVI, and the H-mode in RISC-V can be used. Adapting MOAT to other architectures may require additional engineering effort, as similar hardware features across different platforms often have distinct management interfaces and configuration requirements.

X. CONCLUSION

Despite using BPF to extend kernel functionality, malicious BPF applications can bypass static security checks and conduct unauthorized kernel accesses. We present MOAT to isolate potentially malicious BPF applications from the kernel. MOAT delivers practical and extensible protection with a low cost, in compensation to contemporary BPF verifiers.

ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers and COMPASS members for their insightful comments. This work is partly supported by the National Natural Science Foundation of China under Grant No.U2541211, No.62372218, No.U24A6009. The HKUST authors are supported in part by a RGC CRF grant under the contract C6015-23G.

REFERENCES

- [1] *BPF Documentation — The Linux Kernel Documentation*, 2022. [Online]. Available: <https://docs.kernel.org/bpf/index.html>
- [2] I. V. Project, “BPF Compiler Collection,” <https://github.com/iovisor/bcc>, 2022.
- [3] Y. Zhong, H. Li, Y. J. Wu, I. Zarkadas, J. Tao, E. Mesterhazy, M. Makris, J. Yang, A. Tai, R. Stutsman, and A. Cidon, “XRP: In-Kernel storage functions with eBPF,” in *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. Carlsbad, CA: USENIX Association, Jul. 2022, pp. 375–393.
- [4] H. Zhou, S. Wu, X. Luo, T. Wang, Y. Zhou, C. Zhang, and H. Cai, “NCScope: hardware-assisted analyzer for native code in android apps,” in *ISSTA 22: 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, Virtual Event, South Korea, July 18 - 22, 2022*, S. Ryu and Y. Smaragdakis, Eds. ACM, 2022, pp. 629–641.
- [5] Cilium, “Cilium,” <https://github.com/cilium/cilium>, 2022.
- [6] P. Enberg, A. Rao, and S. Tarkoma, “Partition-aware packet steering using XDP and eBPF for improving application-level parallelism,” in *Proceedings of the 1st ACM CoNEXT Workshop on Emerging in-Network Computing Paradigms, ENCP@CoNEXT 2019, Orlando, FL, USA, December 9, 2019*. ACM, 2019, pp. 27–33.
- [7] T. Høiland-Jørgensen, J. D. Brouer, D. Borkmann, J. Fastabend, T. Herbert, D. Ahern, and D. Miller, “The EXpress data path: Fast programmable packet processing in the operating system kernel,” in *Proceedings of the 14th International Conference on Emerging Networking EXperiments and Technologies*, ser. CoNEXT ’18. New York, NY, USA: Association for Computing Machinery, 2018, p. 54–66.
- [8] *eBPF Verifier — The Linux Kernel Documentation*, 2022. [Online]. Available: <https://docs.kernel.org/bpf/verifier.html>
- [9] L. Nelson, J. V. Geffen, E. Torlak, and X. Wang, “Specification and verification in the field: Applying formal methods to BPF just-in-time compilers in the linux kernel,” in *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association, Nov. 2020, pp. 41–61.
- [10] X. Wang, D. Lazar, N. Zeldovich, A. Chlipala, and Z. Tatlock, “Jitk: A trustworthy In-Kernel interpreter infrastructure,” in *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. Broomfield, CO: USENIX Association, Oct. 2014, pp. 33–47.
- [11] E. Gershuni, N. Amit, A. Gurfinkel, N. Narodytska, J. A. Navas, N. Rinetzky, L. Ryzhyk, and M. Sagiv, “Simple and precise static analysis of untrusted linux kernel extensions,” in *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI 2019. New York, NY, USA: Association for Computing Machinery, 2019, p. 1069–1084.
- [12] S. V. Mahadevan, Y. Takano, and A. Miyaji, “PRSafe: Primitive recursive function based domain specific language using llvm,” in *2021 International Conference on Electronics, Information, and Communication (ICEIC)*, 2021, pp. 1–4.
- [13] E. T. Luke Nelson, Xi Wang, “A proof-carrying approach to building correct and flexible in-kernel verifiers.” Linux Plumbers Conference, 2021.
- [14] Facebook, “Katan: A high performance layer 4 load balancer,” Oct 2023. [Online]. Available: <https://github.com/facebookincubator/katan/tree/main>
- [15] *Kprobes Documentation — The Linux Kernel Documentation*, 2022. [Online]. Available: <https://docs.kernel.org/trace/kprobes.html>
- [16] H. Vishwanathan, M. Shachnai, S. Narayana, and S. Narakatte, “Sound, precise, and fast abstract interpretation with tristate numbers,” in *Proceedings of the 20th IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO ’22. IEEE Press, 2022, p. 254–265.
- [17] MITRE, “CVE-2020-27194.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2020-27194>.
- [18] MITRE, “CVE-2020-8835.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2020-8835>.
- [19] MITRE, “CVE-2021-31440.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-31440>.
- [20] MITRE, “CVE-2021-33200.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-33200>.
- [21] MITRE, “CVE-2021-3444.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-3444>.
- [22] MITRE, “CVE-2021-3490.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-3490>.
- [23] MITRE, “CVE-2021-45402.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-45402>.
- [24] MITRE, “CVE-2022-2785.” <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-CVE-2022-2785>.
- [25] MITRE, “CVE-2022-23222.” <http://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-23222>.
- [26] *Intel 64 and IA-32 Architectures Software Developer Manuals*, 2022. [Online]. Available: <https://www.intel.com/content/www/us/en/developer/articles/technical/intel-sdm.html>
- [27] *ARM Architecture Reference Manual*, 2024. [Online]. Available: <https://developer.arm.com/documentation/ddi0487/la/>
- [28] *AMD64 Architecture Programmer’s Manual*, 2024. [Online]. Available: https://docs.amd.com/v/u/en-US/40332-PUB_4.08
- [29] *RISC-V Technical Specifications*, 2025. [Online]. Available: <https://lf-riscv.atlassian.net/wiki/x/kYD2>
- [30] H. Lu, S. Wang, Y. Wu, W. He, and F. Zhang, “MOAT: Towards safe BPF kernel extension,” in *33rd USENIX Security Symposium (USENIX Security 24)*. Philadelphia, PA: USENIX Association, Aug. 2024, pp. 1153–1170.
- [31] M. Larabel, “Phoronix test suite,” Mar 2024. [Online]. Available: <https://www.phoronix-test-suite.com/>
- [32] “MOAT website,” Oct 2023. [Online]. Available:

- <https://sites.google.com/view/safe-bpf/>
- [33] D. Merkel, “Docker: lightweight linux containers for consistent development and deployment,” *Linux journal*, vol. 2014, no. 239, p. 2, 2014.
 - [34] *BPF-Helpers(7) - Linux Manual Page*, 2021. [Online]. Available: <https://man7.org/linux/man-pages/man7/bpf-helpers.7.html>
 - [35] *eBPF Maps — The Linux Kernel Documentation*, 2022. [Online]. Available: <https://docs.kernel.org/bpf/maps.html>
 - [36] A. Vahldiek-Oberwagner, E. Elnikety, N. O. Duarte, M. Sammler, P. Druschel, and D. Garg, “ERIM: Secure, efficient in-process isolation with protection keys MPK,” in *28th USENIX Security Symposium (USENIX Security 19)*. Santa Clara, CA: USENIX Association, Aug. 2019, pp. 1221–1238.
 - [37] Y. Chen, J. Li, G. Xu, Y. Zhou, Z. Wang, C. Wang, and K. Ren, “SGXLock: Towards efficiently establishing mutual distrust between host application and enclave for SGX,” in *31st USENIX Security Symposium (USENIX Security 22)*. Boston, MA: USENIX Association, Aug. 2022, pp. 4129–4146.
 - [38] M. S. Melara, M. J. Freedman, and M. Bowman, “EnclaveDom: Privilege separation for large-tcb applications in trusted execution environments,” 2019.
 - [39] S. Park, S. Lee, W. Xu, H. Moon, and T. Kim, “libmpk: Software abstraction for intel memory protection keys (Intel MPK),” in *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. Renton, WA: USENIX Association, Jul. 2019, pp. 241–254.
 - [40] *Capabilities - Overview of Linux capabilities — The Linux manual page*, 2022. [Online]. Available: <https://man7.org/linux/man-pages/man7/capabilities.7.html>
 - [41] S. Support, “Security hardening: Use of ebpf by unprivileged users has been disabled by default,” Jan 2022. [Online]. Available: <https://www.suse.com/support/kb/doc/?id=000020545>
 - [42] C. Inc, “Unprivileged ebpf disabled by default for ubuntu 20.04 lts, 18.04 lts, 16.04 esm,” Mar 2022. [Online]. Available: <https://discourse.ubuntu.com/t/unprivileged-ebpf-disabled-by-default-for-ubuntu-20-04-lts-18-04-lts-16-04-esm/27047>
 - [43] “eBPF seccomp() filters,” May 2021. [Online]. Available: <https://lwn.net/Articles/857228/>
 - [44] “cBPF is frozen without fix/extension,” May 2019. [Online]. Available: <https://lwn.net/ml/netdev/20190509044720.fxlcldi74atev5id@ast-mbp/>
 - [45] “Unprivileged bpf(),” Oct 2015. [Online]. Available: <https://lwn.net/Articles/660331/>
 - [46] “Reconsidering unprivileged BPF,” Aug 2019. [Online]. Available: <https://lwn.net/Articles/796328/>
 - [47] “Unprivileged BPF and authoritative security hooks,” Apr 2023. [Online]. Available: <https://lwn.net/Articles/929746/>
 - [48] MITRE, “CVE-2021-4001.” <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-4001>.
 - [49] MITRE, “CVE-2021-29155.” <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-29155>.
 - [50] G. Chen, S. Chen, Y. Xiao, Y. Zhang, Z. Lin, and T. Lai, “SgxPectre attacks: Leaking enclave secrets via speculative execution,” *CoRR*, vol. abs/1802.09085, 2018.
 - [51] J. Cui, J. Z. Yu, S. Shinde, P. Saxena, and Z. Cai, “SmashEx: Smashing SGX enclaves using exceptions,” in *CCS ’21: 2021 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, Republic of Korea, November 15 - 19, 2021*, Y. Kim, J. Kim, G. Vigna, and E. Shi, Eds. ACM, 2021, pp. 779–793.
 - [52] H. M. Demoulin, I. Pedisich, N. Vasilakis, V. Liu, B. T. Loo, and L. T. X. Phan, “Detecting asymmetric application-layer Denial-of-Service attacks In-Flight with FineLame,” in *Proceedings of the 2019 USENIX Annual Technical Conference (USENIX ATC 19)*. Renton, WA: USENIX Association, 2019, pp. 693–708.
 - [53] S. Cauligi, C. Disselkoben, D. Moghimi, G. Barthe, and D. Stefan, “Sok: Practical foundations for software spectre defenses,” in *2022 IEEE Symposium on Security and Privacy (SP)*, 2022, pp. 666–680.
 - [54] J. Zhipeng, C. Hua, F. Jingyi, K. Xiaoyun, Y. Yiwei, L. Haoyuan, and F. Limin, “A combined countermeasure against side-channel and fault attack with threshold implementation technique,” *Chinese Journal of Electronics*, vol. 32, no. 2, pp. 199–208, 2023.
 - [55] W. Yichuan, G. Wen, H. Xinhong, and D. Yanning, “Method and practice of trusted embedded computing and data transmission protection architecture based on android,” *Chinese Journal of Electronics*, vol. 33, no. 3, pp. 623–634, 2024.
 - [56] tr3e, “CVE-2022-23222: Linux Kernel eBPF Local Privilege Escalation,” Jun 2022. [Online]. Available: <https://github.com/tr3ee/CVE-2022-23222>
 - [57] CVEDetails, “CVE-2020-27171.” <https://www.cvedetails.com/cve/CVE-2020-27171>.
 - [58] Z. Yuan, S. Hong, R. Chang, Y. Zhou, W. Shen, and K. Ren, “VDom: Fast and unlimited virtual domains on multiple architectures,” in *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, ser. ASPLOS 2023. New York, NY, USA: Association for Computing Machinery, 2023, p. 905–919.
 - [59] Bootlin, “BPF Configuration Map — config_map,” https://elixir.bootlin.com/linux/v5.10/source/samples/bpf/xdp_rxq_info_kern.c#L32, 2023.
 - [60] MITRE, “CVE-2021-34866.” <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-34866>.
 - [61] N. Carlini and D. Wagner, “ROP is still dangerous: Breaking modern defenses,” in *Proceedings of the 23rd USENIX Conference on Security Symposium*, ser. SEC’14. USA: USENIX Association, 2014, p. 385–399.
 - [62] Bootlin, “Linux — BPF_TRACE_MAX_PROGS,” https://elixir.bootlin.com/linux/latest/source/kernel/trace/bpf_trace.c#L2115, 2023.
 - [63] “Analysing behaviour using events and tracepoints,” 2023. [Online]. Available: <https://www.kernel.org/doc/Documentation/trace/tracepoint-analysis.txt>
 - [64] K. Lucas, “UnixBench: the original BYTE UNIX benchmark suite, updated and revised by many people

- over the years.” Apr 2023. [Online]. Available: <https://github.com/kdlucas/byte-unixbench>
- [65] D. Jin, V. Atlidakis, and V. P. Kemerlis, “EPF: Evil packet filter,” in *2023 USENIX Annual Technical Conference (USENIX ATC 23)*. Boston, MA, USA: USENIX Association, 2023, pp. 735–751.
- [66] *Seccomp BPF (SECure COMputing with filters)*, 2022. [Online]. Available: https://www.kernel.org/doc/html/latest/userspace-api/seccomp_filter.html
- [67] N. DeMarinis, K. Williams-King, D. Jin, R. Fonseca, and V. P. Kemerlis, “sysfilter: Automated system call filtering for commodity software,” in *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*. San Sebastian: USENIX Association, Oct. 2020, pp. 459–474.
- [68] W. He, H. Lu, F. Zhang, and S. Wang, “Ringuard: Guard io_uring with eBPF,” in *Proceedings of the 1st Workshop on eBPF and Kernel Extensions*, ser. eBPF ’23. New York, NY, USA: Association for Computing Machinery, 2023, p. 56–62.
- [69] W. Glozer, “wrk - a http benchmarking tool,” Feb 2021. [Online]. Available: <https://github.com/wg/wrk>
- [70] S. Y. Lim, X. Han, and T. Pasquier, “Unleashing unprivileged ebpf potential with dynamic sandboxing,” in *Proceedings of the 1st Workshop on EBPF and Kernel Extensions*, ser. eBPF ’23. New York, NY, USA: Association for Computing Machinery, 2023, p. 42–48.
- [71] J. Gu, H. Li, W. Li, Y. Xia, and H. Chen, “EPK: Scalable and efficient memory protection keys,” in *2022 USENIX Annual Technical Conference (USENIX ATC 22)*. Carlsbad, CA: USENIX Association, Jul. 2022, pp. 609–624.
- [72] S. Gravani, M. Hedayati, J. Criswell, and M. L. Scott, “Fast intra-kernel isolation and security with IskiOS,” in *24th International Symposium on Research in Attacks, Intrusions and Defenses*, ser. RAID ’21. New York, NY, USA: Association for Computing Machinery, 2021, p. 119–134.
- [73] V. Narayanan, A. Balasubramanian, C. Jacobsen, S. Spall, S. Bauer, M. Quigley, A. Hussain, A. Younis, J. Shen, M. Bhattacharyya, and A. Burtsev, “LXD: Towards isolation of kernel subsystems,” in *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. Renton, WA: USENIX Association, Jul. 2019, pp. 269–284.
- [74] V. Narayanan, Y. Huang, G. Tan, T. Jaeger, and A. Burtsev, “Lightweight kernel isolation with virtualization and vm functions,” in *Proceedings of the 16th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments*, ser. VEE ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 157–171.
- [75] A. M. Azab, K. Swidowski, R. Bhutkar, J. Ma, W. Shen, R. Wang, and P. Ning, “Skee: A lightweight secure kernel-level execution environment for arm,” in *Network and Distributed System Security Symposium*, 2016.
- [76] M. Castro, M. Costa, J. Martin, M. Peinado, P. Akritidis, A. Donnelly, P. Barham, and R. Black, “Fast byte-granularity software fault isolation,” in *Proceedings of the 22nd ACM Symposium on Operating Systems Principles 2009, SOSP 2009, Big Sky, Montana, USA, October 11-14, 2009*. ACM, 2009, pp. 45–58.
- [77] K. Huang, M. Payer, Z. Qian, J. Sampson, G. Tan, and T. Jaeger, “Sok: Challenges and paths toward memory safety for ebpf,” in *2025 IEEE Symposium on Security and Privacy (SP)*, 2025, pp. 848–866.
- [78] P. Zhang, C. Wu, X. Meng, Y. Zhang, M. Peng, S. Zhang, B. Hu, M. Xie, Y. Lai, Y. Kang, and Z. Wang, “HIVE: A hardware-assisted isolated execution environment for eBPF on AArch64,” in *33rd USENIX Security Symposium (USENIX Security 24)*. Philadelphia, PA: USENIX Association, Aug. 2024, pp. 163–180.



Lijian Huang is working on his Master’s degree from Southern University of Science and Technology (SUSTech). He received his Bachelor’s degree in Cyberspace Security from Hangzhou Dianzi University in 2023. His research interests include system security and virtualization on the Arm architecture.



Hongyi LU is a Ph.D. candidate in CSE at HKUST and SUSTech, jointly supervised by Professor Shuai Wang and Professor Fengwei Zhang. His research focuses on areas of system security, including kernel hardening, hardware-software co-design security mechanisms and automated software vulnerability detection.



Shuai Wang is an Associate Professor at CSE, HKUST. He received his Ph.D. from Penn State University, and B.S. from Peking University. He is broadly interested in computer security and specializes in binary code analysis, software and system security, and low-level security techniques.



Fengwei Zhang is an Associate Professor in Department of Computer Science and Engineering at Southern University of Science and Technology (SUSTech). His primary research interests are in the areas of systems security, with a focus on trustworthy execution, hardware-assisted security, debugging transparency, and plausible deniability encryption. Before joining SUSTech, he spent four years as an Assistant Professor at Department of Computer Science at Wayne State University.

APPENDIX A SUPPORTED BPF HELPERS

We list the helper functions that are tested on MOAT in Table X. To test these helper functions, we adapt the BPF programs included in the Linux kernel tree to invoke these helpers. We also check their results to ensure these helpers are executing correctly on a MOAT-enabled system.

TABLE X: MOAT-supported helpers.

Type	Supported BPF Helpers
Map	bpff_map_lookup_elem, bpff_map_update_elem, bpff_map_delete_elem, bpff_map_push_elem, bpff_map_pop_elem, bpff_map_peek_elem
String	bpff_strtol, bpff_strtoul, bpff_strncmp
Utilities	bpff_trace_vprintk, bpff_get_retval, bpff_set_retval, bpff_user_rnd_u32, bpff_get_raw_cpu_id, bpff_get_smp_processor_id, bpff_ktime_get_ns, bpff_ktime_get_boot_ns, bpff_ktime_get_coarse_ns, bpff_get_current_pid_tgid, bpff_get_current_uid_gid, bpff_jiffies64, bpff_get_attach_cookie
Cgroup	bpff_get_current_cgroup_id, bpff_get_cgroup_classid_curr, bpff_get_cgroup_classid, bpff_skb_cgroup_id, bpff_sk_ancestor_cgroup_id, bpff_skb_cgroup_classid, bpff_sk_cgroup_id, bpff_skb_ancestor_cgroup_id, bpff_get_current_ancestor_cgroup_id
Tracing	bpff_probe_read_compat_str, bpff_probe_read_compat, bpff_probe_read_kernel_str, bpff_probe_read_kernel, bpff_get_current_task, bpff_get_func_ip_tracing, bpff_task_pt_regs, bpff_perf_event_read, bpff_perf_event_read_value, bpff_perf_event_output, bpff_get_func_ret, bpff_get_func_arg, bpff_get_func_arg_cnt, bpff_get_func_ip, bpff_get_ns_current_pid_tgid
Ringbuf	bpff_ringbuf_discard, bpff_ringbuf_query, bpff_ringbuf_submit, bpff_ringbuf_reserve, bpff_ringbuf_output
XDP	bpff_xdp_fib_lookup, bpff_xdp_load_bytes, bpff_xdp_store_bytes, bpff_xdp_adjust_head, bpff_xdp_adjust_meta, bpff_xdp_adjust_tail, bpff_xdp_get_buff_len
Socket	bpff_get_listener_sock, bpff_skb_get_pay_offset, bpff_skc_to_mptcp_sock, bpff_skc_to_tcp6_sock, bpff_skc_to_tcp_request_sock, bpff_skc_to_tcp_sock, bpff_skc_to_tcp_timewait_sock, bpff_skc_to_udp6_sock, bpff_skc_to_unix_sock, bpff_sk_fullsock, bpff_sk_release, bpff_tcp_sock, bpff_skb_load_helper_8_no_cache, bpff_skb_load_helper_16_no_cache, bpff_skb_load_helper_32_no_cache, bpff_sock_ops_cb_flags_set, bpff_task_storage_delete, bpff_skb_load_bytes, bpff_skb_load_helper_16, bpff_skb_load_helper_32, bpff_skb_load_helper_8

APPENDIX B BPF CVE LIST

We provide the list of evaluated BPF CVEs in Table XI.

TABLE XI: Evaluated BPF CVEs.

CVE ID
2016-2383, 2017-16995, 2017-16996, 2017-17852, 2017-17853, 2017-17854, 2017-17855, 2017-17856, 2017-17857, 2017-17862, 2017-17863, 2017-17864, 2018-18445, 2020-8835, 2020-27194, 2021-34866, 2021-3489, 2021-3490, 2021-20268, 2021-3444, 2021-33200, 2021-45402, 2022-2785, 2022-23222, 2023-39191, 2023-2163

APPENDIX C CRITICAL OBJECTS

We list the critical objects identified by us in Table XII. We categorize the objects in Table XII based on their locations. Despite different location (i.e., map and iterator), both of them are used to implement dynamic dispatching in the kernel.

TABLE XII: Critical objects in the BPF subsystem.

Location	Critical Object
Map	array_map_ops, percpu_array_map_ops, prog_array_map_ops, sock_map_ops, cgroup_storage_map_ops, htab_map_ops, htab_percpu_map_ops, htab_lru_map_ops, htab_lru_percpu_map_ops, trie_map_ops, task_storage_map_ops, dev_map_ops, sk_storage_map_ops, cpu_map_ops, xsk_map_ops, perf_event_array_map_ops, queue_map_ops, stack_map_ops, bpff_struct_ops_map_ops, ringbuf_map_ops, bloom_filter_map_ops, cgroup_storage_map_ops, cgroup_array_map_ops, array_of_maps_map_ops, stack_trace_map_ops, htab_of_maps_map_ops, user_ringbuf_map_ops, inode_storage_map_ops
Iterator	cgroup_iter_seq_info, sock_map_iter_seq_info, sock_hash_iter_seq_info, ksym_iter_seq_info, bpff_link_seq_info, bpff_map_seq_info, sock_hash_iter_seq_info, sock_map_iter_seq_info, ipv6_route_seq_info, iter_seq_info, ksym_iter_seq_info, netlink_seq_info, tcp_seq_info, udp_seq_info, unix_seq_info, bpff_prog_seq_info

APPENDIX D GLOSSARY

We provide a glossary of platform-specific jargon used in the paper in Table XIII.

TABLE XIII: Glossary of platform-specific jargon

Term	Description
MPK	Memory Protection Keys, a hardware feature in the Intel x86 architecture that allows processes to set memory access permissions on a per-page basis.
Stage-2 Translation	The second stage of address translation in the Arm architecture, which translates the intermediate physical address (IPA) to the final physical address (PA).
PCID	Process Context Identifier, a unique identifier to mark a process address space in the x86 architecture.
VMID	Virtual Machine Identifier, a unique identifier to mark IPA in the Arm architecture.
Stage-2 TLB flush	Flushing the translation lookaside buffer entries for the Arm Stage-2 page table in the Arm architecture. On Arm, TLB entries tagged with VMID map VA to PA, rather than IPA.