

SoK: Challenges and Paths Toward Memory Safety for eBPF

Kaiming Huang^{*}, Mathias Payer[†], Zhiyun Qian[‡], Jack Sampson^{*}, Gang Tan^{*}, Trent Jaeger[‡]

^{*} Penn State [†] EPFL [‡] UC Riverside

eBPF Has Been Drawing More Attention

- Networking

BMC: Accelerating Memcached using Safe In-kernel Caching and Pre-stack Processing
Yoann Ghigoff, Orange Labs, Sorbonne Université, Inria, LIP6; Julien Sopena, Sorbonne Université, LIP6; Kahina Lazri, Orange Labs; Antoine Blin, Gandi; Gilles Muller, Inria

State-Compute Replication: Parallelizing High-Speed Stateful Packet Processing Taking 5G RAN Analytics and Control to a New Level
Xenofon Foukas, Bozidar Radunovic, Matthew Balkwill, Zhihua Lai, Microsoft, Cambridge, United Kingdom
kz_computer@email.ru.nl

Program Warping
Inella^{1,3}, Giuseppe

Surfing on the Edge of the Cloud
Thomas Wirtgen, Tom Rousseau, Quentin De Coninck, and Nicolas Rühwiler

Surfing on the Edge of the Cloud
Xuan Zhang^Δ, London, UK

- Optimization

λ-IO: A Unified IO Stack for Computational Storage
MERLIN: Multi-tier Optimization of eBPF Code for Data Planes
XRP: In-Kernel Storage Functions with eBPF
SPRIGHT: Extracting the Server from Serverless Computing!

Extension Framework for File Systems in User space
Ashish Bijlani, Georgia Institute of Technology
Umakishore Ramachandran, Georgia Institute of Technology

- Security

HIVE: A Hardware-assisted Isolated Execution Environment for eBPF on AArch64
Daihua Zhang^{1,2}, Changqian Wu^{1,2,3}, Yixiong Mao^{4†}, Xingyu Zhang⁵, Mingfan Deng^{1,2}, Shiyang Wang^{1,2,3*}

Falco
Di Jin, Vagdeli, https://github.com/falco-go

An Analysis of eBPF Security
Ofek Kirzner and Adam Morrison, Tel Aviv University
George Mason University, Tsinghua University

Xijia Che
University and BNRist

- Tracing

Network-Centric Distributed Tracing with DeepFlow:
The Design and Implementation of Hyperupcalls

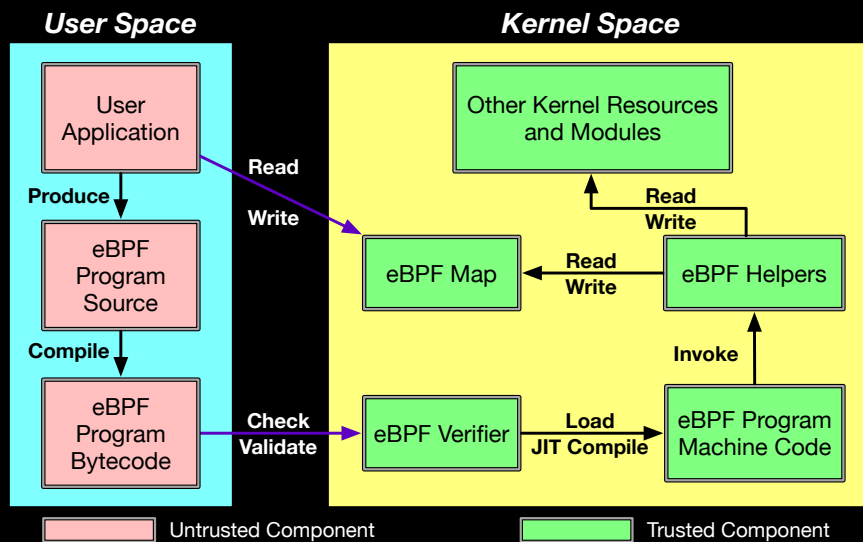
Eliminating eBPF Tracing Overhead in User Processes
bpftrace
Dynamic Tracing for Linux
bpftrace is a high-level tracing language for Linux and provides a quick and easy way for people to write observability-based eBPF programs, especially those unfamiliar with the complexities of eBPF.

Hussain
Virginia Tech
Blacksburg, VA, USA
hussain@vt.edu

Uddhav Gautam
Virginia Tech
Blacksburg, VA, USA
upgautam@vt.edu

Dan Williams
Virginia Tech
Blacksburg, VA, USA
dwillia@vt.edu

eBPF Workflow and Its Trust Model



Do they work as expected?

- eBPF programs should not perform unsafe memory accesses.
- eBPF helper functions are trusted kernel APIs but do not have any validations.
- The verifier must ensure that accesses to the kernel data do not populate memory errors.
- The verifier must be free of implementation bugs, as any bug can be exploited to load unsafe programs.
- The eBPF trust model relies critically on the eBPF verifier to enforce memory safety.

Memory Safety Issue in eBPF Verifier

- eBPF verifier has been becoming a significant source of bugs
 - 46 CVEs in eBPF verifier in 2024
 - 325 Syzbot-reported bug related to eBPF submodule in Linux Kernel
- Checks are unsound and incomplete
 - Bugs left unchecked amid removal of safety checks by optimizations
 - Checks are incomplete for ensuring full memory safety
- Checks are limited in scope in terms of complete workflow
 - Checks of the verifier are limited to the eBPF bytecode

Memory Safety Issue in eBPF Verifier

- Checks are limited in scope in terms of complete workflow
 - Checks of the verifier are limited to the eBPF bytecode

```
1 static void *__dev_map_lookup_elem(  
2     struct bpf_map *map, u32 key){  
3     struct bpf_dtab *dtab =  
4         container_of(map, struct bpf_dtab, map);  
5     struct bpf_dtab_netdev *obj;  
6     if (key >= map->max_entries)  
7         return NULL;  
8     obj = rcu_dereference_check(dtab->netdev_map[key],  
9         rcu_read_lock_bh_held());  
10    return obj;  
11 }
```

No checks in `map_lookup_elem` to ensure the initialization of `obj`

```
1 SEC("classifier")  
2 int example_prog(struct __sk_buff *skb) {  
3     int index = 0; // Key for accessing dev_map  
4     int *dev_ifindex;  
5     // Use dev_map_lookup_elem to retrieve the interface  
6     dev_ifindex = dev_map_lookup_elem(&dev_map, &index);  
7     if (!dev_ifindex) {  
8         return TC_ACT_SHOT; // Drop packet if fails  
9     }  
10    // Uninitialized memory access  
11    *dev_ifindex += 1; // KMSAN uninit warning  
12    // Final decision to accept or drop the packet  
13    return TC_ACT_OK;  
14 }
```

Attacker can easily forge a malicious eBPF program to exploit UBI

Kernel Defenses

Category	Kernel Defensive Features	Description
Required Defense	eBPF Verifier	Validates security of eBPF programs.
Optional Defense	Capability CAP_BPF BPF LSM (Linux Security Modules) BPF Type Format (BTF) and CO-RE	Permits only privileged users to attach eBPF programs. Enforces access control over eBPF programs Validates data type and version compatibility.
General Defense	CFI and Execute-Only Memory (XOM) Memory Tagging Shadow Stacks kASAN kASLR SMAP and SMEP	Prevents control flow hijacking and code reuse attacks. Prevents pointers from being tampered and forged. Protects return addresses. Detects memory errors at runtime. Randomizes memory layout. Prevents unauthorized user-space memory access in kernel mode.

- eBPF-specific defenses are limited by optional settings and leave room for attacks with limited privilege.
- General defenses fail to fully block eBPF-based attacks.

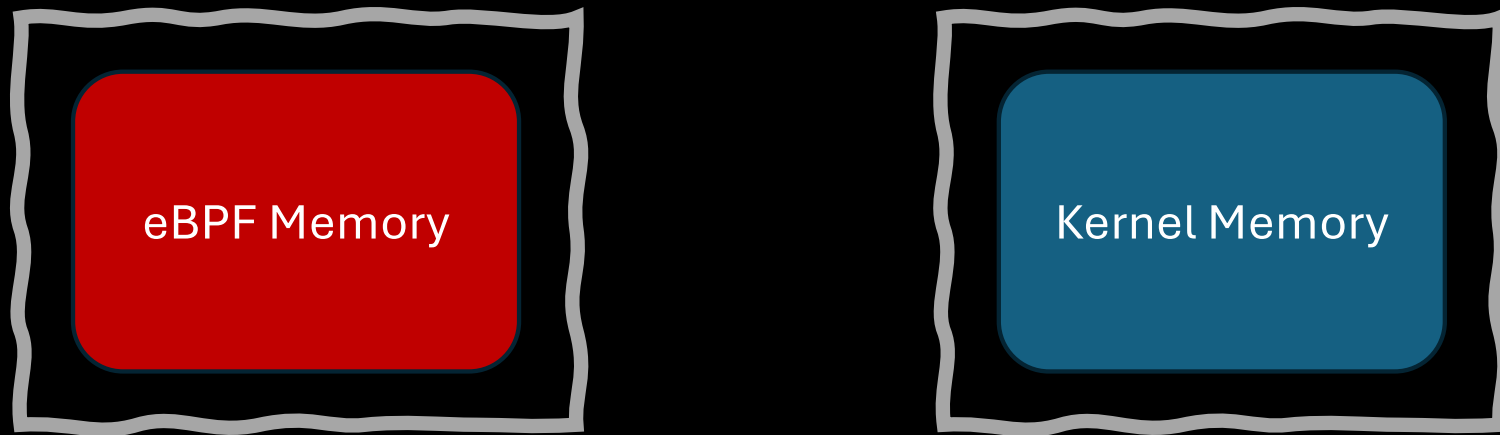
Take Capability CAP_BPF as an Example

- Introduced in **Linux 5.8 (Aug 2020)**
 - Designed to **restrict unprivileged users** from attaching **eBPF** programs
- CAP_BPF is not a hard restriction
 - Users can **opt out** and still attach eBPF programs
- Privileged enforcement reduces flexibility
 - **Vendors such as Cilium** rely on **unprivileged eBPF**
- CAP_BPF illustrates the **tension between security and usability**
 - **Unprivileged execution** remains common in practice

Research Directions

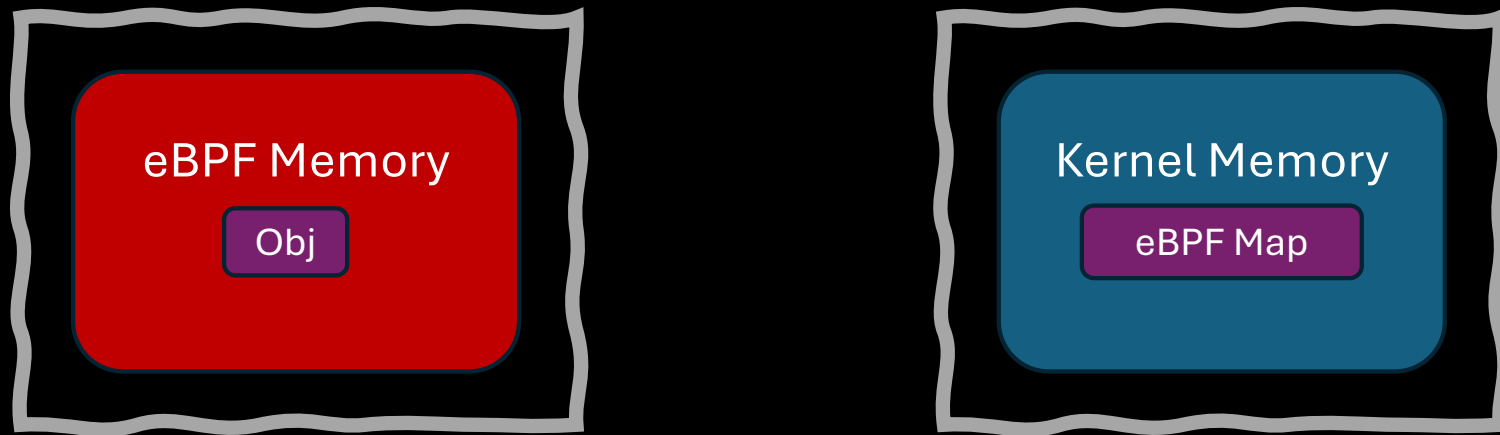
- Fuzzing
 - Inherently incomplete
 - Hard to generate eBPF programs that both pass verifier and trigger bugs
- Isolation
 - Aims to restrict access to eBPF and shared (eBPF map) memory
 - Does not address risks from **indirect kernel access**
- Runtime Checks
 - Limited by the resource constraints and instruction limits.
- Static Validation
 - Existing approaches are either unsound or incomplete.

Limitation – Isolation as Example



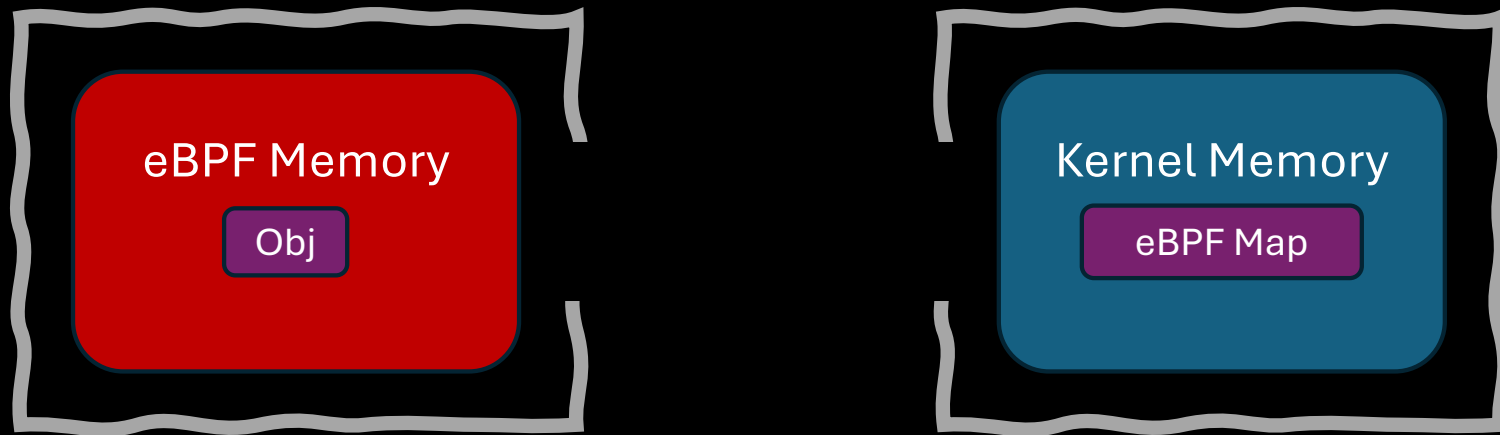
- Isolation separates eBPF memory and Kernel memory
- Unauthorized memory accesses are prevented at isolation boundary

Limitation – Isolation as Example



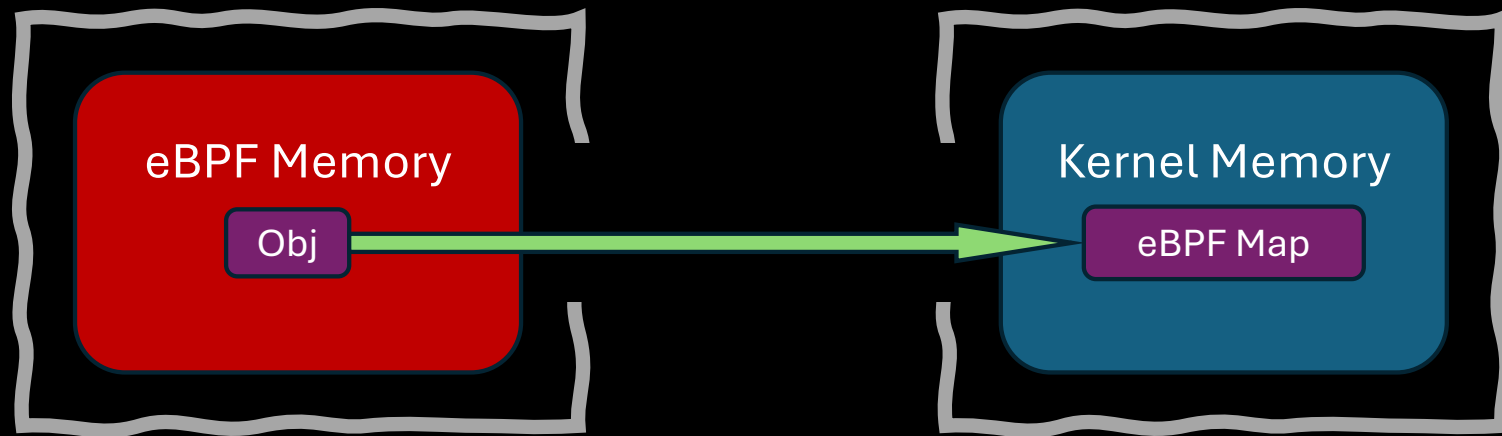
- However, objects in eBPF program can be transmitted to be saved in kernel data (e.g., eBPF map) via eBPF helpers.

Limitation – Isolation as Example



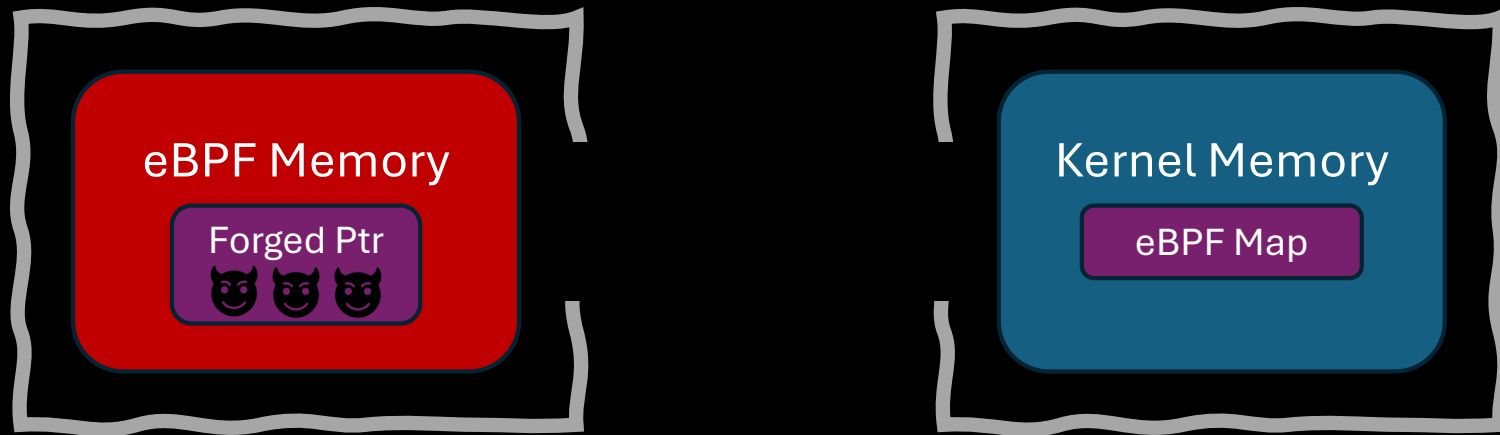
- However, objects in eBPF program can be transmitted to be saved in kernel data (e.g., eBPF map) via eBPF helpers.
- The access to kernel data through such pointers need to be preserved.

Limitation – Isolation as Example



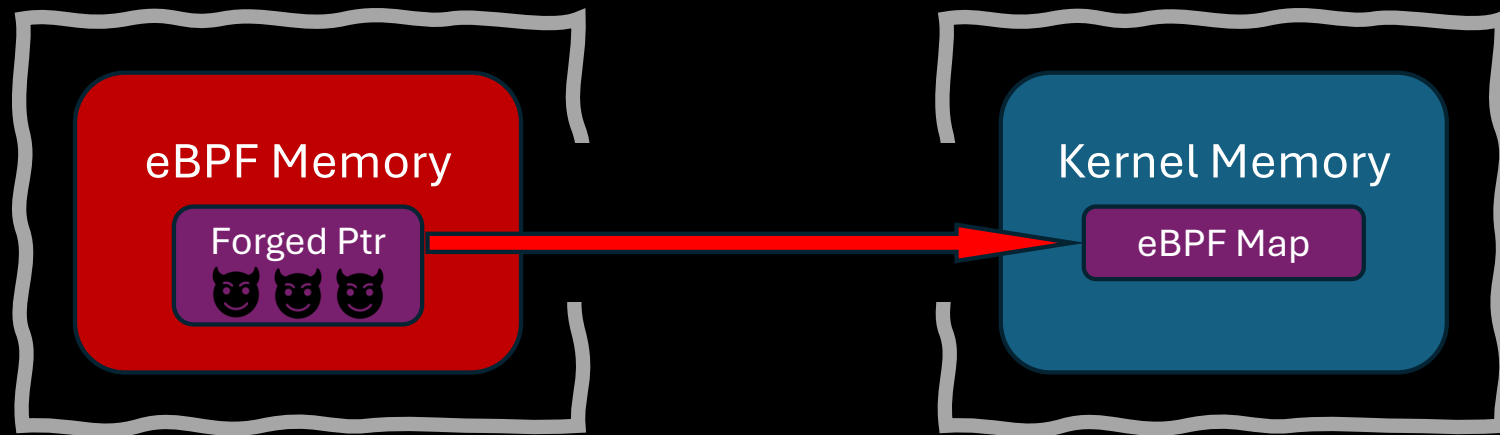
- However, objects in eBPF program can be transmitted to be saved in kernel data (e.g., eBPF map) via eBPF helpers.
- The access to kernel data through such pointers need to be preserved.

Limitation – Isolation as Example



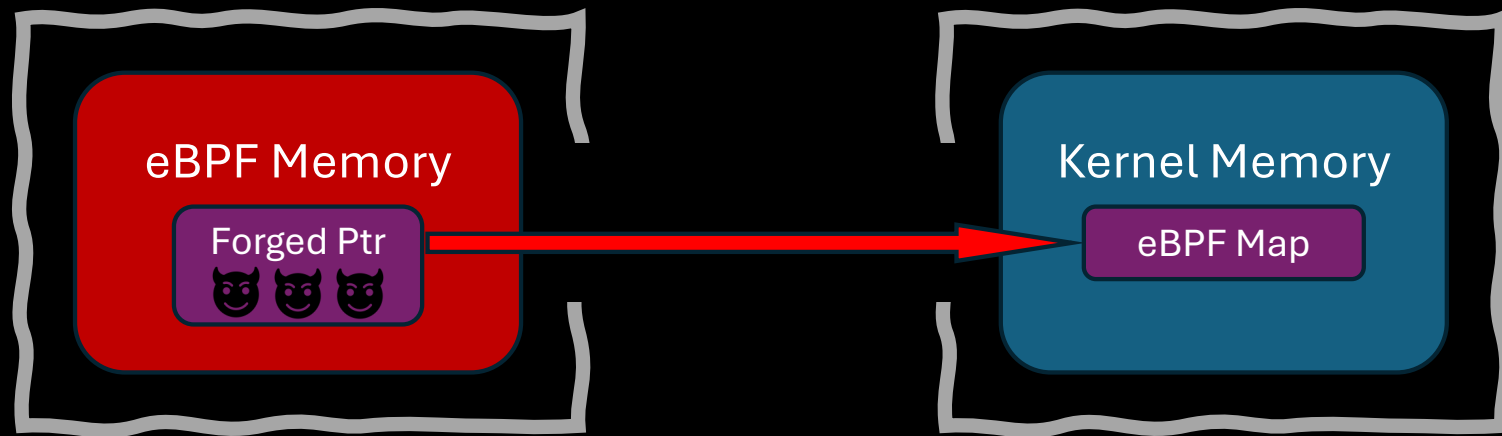
- Such opened-gates allows attacker to
 - exploit the memory errors in eBPF program
 - forge a pointer arbitrarily

Limitation – Isolation as Example



- Such opened-gates allows attacker to
 - exploit the memory errors in eBPF program
 - forge a pointer arbitrarily
 - Pass the pointer to the kernel (through helpers) for unauthorized accesses
 - Known as Cross-boundary Interface Vulnerabilities (CIVs).

Limitation – Isolation as Example



- Attacker can use the forged pointer to escalate exploitability.
 - Examined by EPF and Interp-flow Hijacking attacks.
- *Linux eBPF - new privilege escalation techniques* – Pentera Labs

Protection Scope of Existing Defenses

		eBPF-Only			Shared Objs		
		Spatial	Type	Temp	Spatial	Type	Temp
eBPF Verifier [30]	V	●	●	●	●	●	○
HyperBee [47]	V	●	●	○	○	○	○
KFuse [49]	V	●	●	○	●	●	○
PREVAIL [112]	V	●	●	●	●	●	○
SandBPF [36]	II	●	○	○	●	○	○
SafeBPF [37]	II	●	○	○	●	○	○
HIVE [39]	II	●	○	○	●	○	○
MOAT [38]	II	●	○	○	●	○	○
Prevail2Radius [107]	T	●	●	○	●	●	○
Seccomp-eBPF [131]	T	●	○	○	●	○	○
TnumArith [43]	T	●	○	○	○	○	○
RangeAnalysis [44]	T	●	○	○	○	○	○

None of the defenses, whether currently deployed or proposed in research, fully or soundly cover any category of unsafe ops.

Memory Safety in eBPF Context

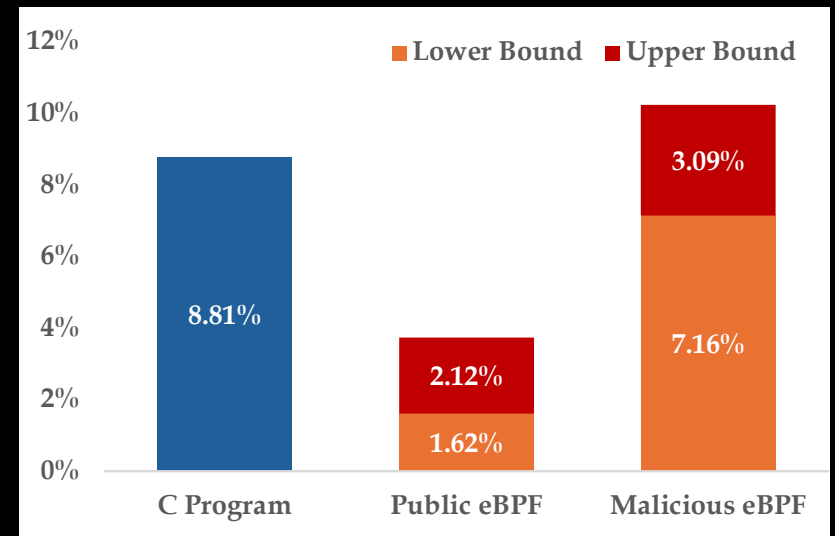
- **Goal:** Memory safe accesses in eBPF programs and kernel use of eBPF-generated pointers
- **Spatial Safety**
 - All accesses of a memory object must only access memory within the object's allocated region.
- **Type Safety**
 - All accesses of a memory object must only access the same data types for each offset and each field.
- **Temporal Safety**
 - All accesses of a memory object must not access the object's allocated region before allocation nor after the object's deallocation.

Identify Unsafe Memory Operation in eBPF

- **Hypothesis:** eBPF should be close to memory safe in terms of low fraction of unsafe operations, but how to identify them?
- Approach – DataGuard (NDSS 2022) and Uriah (CCS 2024)
- Dataset
 - Public eBPF programs – Linux Kernel and BCC
 - Malicious eBPF programs – CVE PoCs and Syzbot reproducers
 - General C programs – evaluated by DataGuard and Uriah

Fraction of Unsafe Memory Operations

- General C Program
 - Fraction similar to Malicious eBPF programs but far more in absolute number of unsafe ops.
- Malicious – higher fractions of unsafe ops
 - **7.16 (lower bound) to 10.25% (upper bound)**
 - Despite being crafted to exploit bugs, the verifier still limits unsafe memory use.
- Public - significantly lower fractions
 - **1.62% (lower bound) to 3.74% (upper bound)**
 - This gap is due to missing kernel-specific constraint information in static analysis.
 - **Upper bound reduced to 1.74%** with updated static analyses for kernel constraints extraction.



How Far are We toward Memory-safe eBPF?

- **Insight 1:** eBPF's linear design makes the fraction of unsafe ops low.
 - Good start for full memory safety validation and enforcement
- **Insight 2:** Must ensure that all eBPF operations cannot exploit victim objects to prevent illicit modification of shared data/pointers
 - Ideally, make eBPF programs run in a memory safe manner
 - E.g., Easier to validate temporal safety statically
- **Insight 3:** Must ensure that all pointer values shared with the kernel are memory safe

Future Directions

- Enhancing static memory safety validation
 - Extract and apply kernel-specific constraints
 - Adopt compiler-informed techniques, e.g., Rust, WASM
 - Incorporate syntactic annotations, e.g., checked-c
- Advancing finer-grained isolation
 - Pointer forging for indirect corruption
 - Cross-boundary interface vulnerabilities
- Migrating to memory-safe languages