

# Extending Applications Safely and Efficiently

**Yusheng Zheng**<sup>1</sup> • Tong Yu<sup>2</sup> • Yiwei Yang<sup>1</sup> • Yanpeng Hu<sup>3</sup> Xiaozheng Lai<sup>4</sup> • Dan Williams<sup>5</sup> • Andi Quinn<sup>1</sup>

<sup>1</sup>UC Santa Cruz <sup>2</sup>eunomia-bpf Community <sup>3</sup>ShanghaiTech University <sup>4</sup>South China University of Technology <sup>5</sup>Virginia Tech

### Extensions are everywhere



### What are extensions?

• Customize software without modifying source code

### Why do we need them?

• Different deployments, different needs

## Nginx firewall example

### **Before deployment, user:**

- Writes firewall using nginx APIs
- Associates firewall with request processing extension entry.

### **During runtime, Nginx:**

- Jumps to firewall when reaching request processing entry.
- Executes firewall in the extension runtime execution context.





## Extension Problems & requirements

# Real-world safety violations:

• Bilibili CDN outage, Apache buffer overflow, Redis RCE

### **Performance penalty:**

 WebAssembly/Lua impose 10-15% overhead

### **Requirements:**

- Fine-grained safety and interconnectedness trade-offs
- Isolation
- Efficiency



### State-of-the-Art Falls Short

- Dynamic loading: efficiency but no isolation or finegrained safety-interconnectedness policies (LD\_PRELOAD, DBI tools)
- Software Fault Isolation: safety with 10–15 % performance penalty (XFI [OSDI 06], NaCL [SOSP 09], RL-Box [USENIX Security 20], Wasm and Lua)
- Subprocess: strong isolation but high IPC overhead (Wedge [NSDI 08], Shreds [IEEE SP 16], IwC [OSDI 16], and Orbit [OSDI 22])
- Kernel eBPF uprobes: isolation at micro second-level trap cost, low efficiency

### Contributions

# Extension Interface model (EIM)

Navigate fine-grained safety/interconectedness trade-offs for extensions

#### **Bpftime runtime**

Efficient support for EIM and isolation through userspace eBPF runtime

• Up to 6x less overhead than current state-of-the-art!

## Outline

- Background & motivation: Extensions
- → Extension Interface Model (EIM): Fine-grained Interface
- bpftime Runtime: safety & performance
- Evaluation

### **EIM: Extension Interface Model**

- Goal: enable fine-grained safety/interconnectedness trade-offs
- Challenge: supporting per deployment tradeoffs
- Solution:

 $_{\odot}$  Two-phase specification (Development Time and Deployment Time)  $_{\odot}$  Model all resources as capabilities



1. During Development

2. Before Deployment

At Deployment/Runtime

## **EIM: Development Time Specification**

- Developers annotate code for capabilities
- Automatically extracted into capability manifest



1. During Development

```
EIM_STATE_DEFINE(readPid, read, ngx_pid);
EIM_HFUNC_DEFINE_WITH_CONSTRAINTS(
    nginxTime,
    HF_RET_POSITIVE);
EIM_EXTENSION_ENTRY_DEFINE(
    processBegin,
    ngx_http_process_request,
    int,
    struct Request *);
```

# **EIM: Deployment Time Specification**

- YAML policies specify safety/interconnectedness tra deoffs
- Compact policies (avg of 30 lines in evaluation).

# 1 Extension\_Class( 2 name = "observeProcessBegin", 3 extension\_entry = "processBegin",

4 allowed = {instructions < inf, nginxTime, readPid, read(r)})

#### 5 Extension\_Class(

```
6 name = "updateResponse",
```

- 7 extension\_entry = "updateResponseContent"



## Outline

- → Background & motivation: Extensions
- Extension Interface Model (EIM): Fine-grained Interface
- → **bpftime Runtime**: safety & performance
- Evaluation

## bpftime: userspace eBPF extension framework

- Goal: efficiently support EIM and isolation
- Challenge: Existing extension runtimes use heavyweight safety & isolation techniques
- Solution:
  - build new design that exploits eBPF-style verification, binary rewriting, and hardware features to enable efficient intra-process extensions



## Outline

- Background & motivation: Extensions
- Extension Interface Model (EIM): Fine-grained Interface
- bpftime Runtime: safety & performance
- → Evaluation

### Six Real-World Use Cases

<b>bpftime</b> Public	
Userspace eBPF runtime for Observability, Network, GPU & General Extensions Framework	
● C++ 🏠 1k 😵 99	

GitHub: <u>https://github.com/eunomia-bpf/bpftime</u>

### Customization

- Nginx Firewall
- Redis Durability
- FUSE Metadata Cache

### Observability

- DeepFlow
- Syscount
- Sslsniff

### **Customization:** Nginx firewall

 5× to 6× less overhead than lua or WebAssembly



### Observability: sslsniff





### Contributions

### Questions?

**Extension Interface model (EIM)** 

Navigate fine-grained safety/interconectedness trade-offs for extensions **Bpftime runtime** 

Efficient support for EIM and isolation through userspace eBPF runtime

Up to 6x less overhead than current state-of-the-art!



bpftime load ./example/malloc/malloc
bpftime start nginx -c ./nginx.conf

### Backup

## Customization: Nginx firewall

- 5× to 6× improvement
- Less is better



### Observability: sslsniff

 21% less overhead than kernel eBPF



### Four Roles in an Extension Ecosystem



### Micro-Benchmark

Compare with eBPF:

- **Uprobe Dispatch**: 2.56  $\mu$ s  $\rightarrow$  190 ns (14× faster)
- Syscall Tracepoint: 151 ns  $\rightarrow$  232 ns (1.5 × slower)
- **Memory access** (Table 3): user-space read/write 2 ns vs 20 ns (10× faster)
- **Overall**: average 1.5× faster than ubpf/rbpf (Figure 11)

### Extensions have issues

• Example issues caused by extension safety violations

Bug	Software	Summary	
Bilibili [73]	Nginx	Livelock (infinite loop) in an ex-	
		tension caused production out-	
		age.	
CVE-2021-44790 [47]	Apache	Buffer overflow in httpd's lua	
		module causes application to	
		crash.	
CVE-2024-31449 [42]	Redis	Stack overflow in Lua script	
		leads to arbitrary remote code execution.	

• The performance penalty of existing approaches





To get started, you can build and run a libbpf based eBPF program starts with bpftime cli:

٢Ų

# Get started

make -C example/malloc # Build the eBPF program example
export PATH=\$PATH:~/.bpftime/
bpftime load ./example/malloc/malloc

- Use uprobe to monitor userspace malloc function
- Try eBPF in GitHub codespace!(Unprivilidge d container)

In another shell, Run the target program with eBPF inside:

\$ bpftime start ./example/malloc/victim
Hello malloc!
malloc called from pid 250215
continue malloc...
malloc called from pid 250215

### Loader & Runtime Workflow

- Intercept standard eBPF syscalls from libbpf/bcc.
- Parse EIM manifests and DWARF/BTF to generate constraints.
- Verify byte-code via kernel's eBPF verifier with added assertions.
- JIT-Compile verified byte-code into native x86.
- Inject user-runtime via ptrace + Frida + Capstone trampolines.
- Execute extension: flip MPK key → jump to code → flip back → resume.

## Efficient Safety & Isolation

### • The eBPF compatibility challenge:

Linux eBPF has tightly coupled components (compilers, runtime, kernel)
 Prior user eBPF failed by re-implementing entire stack
 **bpftime solution**: Interpose on eBPF syscalls only

### • Key design principles:

Lightweight EIM enforcement

 $\circ$  Concealed extension entries: 10× faster uprobe

### Contribution

- Extension Interface Model (EIM): Fine-grained capability control
- **bpftime Runtime**: Kernel-grade safety with library-grade performance

### State-of-the-Art Falls Short

Approach	Safety	Isolation	Efficiency	Fine-Grained Control
Dynamic Loading	X	X	<ul> <li>✓</li> </ul>	X
SFI (Wasm, Lua)	Limited	$\checkmark$	X (10-15% overhead)	X
Subprocess	~	$\checkmark$	X (context switches)	Limited
eBPF uprobes	$\checkmark$	$\checkmark$	X (kernel traps)	Limited

• No single framework satisfies all requirements

# Summary of EIM

- Existing frameworks → no control OR coarse-grained bundles
- Treats safety and interconnectedness as independent dimensions
- Example policies:

 $\odot$  Monitoring extension: read-only access to specific variables  $\odot$  Firewall extension: read/write for response modification

## bpftime - Why We Need a New Runtime

- Can't existing frameworks enforce EIM efficiently?
  - WebAssembly/SFI: 10-15% overhead, Subprocess isolation: Expensive switches, Kernel eBPF uprobes: Kernel traps
- A userspace extension framework in eBPF
  - $_{\odot}$  Compatibility and Work together with kernel eBPF extensions
  - verification for safety
  - $\odot$  Conceal for efficient
  - $\odot\, {\rm Mpk}$  for isolation

















## Nginx firewall example

User wants to have a firewall to block malicious requests

User write custom firewall logic using nginx helper functions Load their extension at an extension entry for request pr

# **Extension execution model:** Thread → Extension entry → Jump to extension → Execute by extension runtime → Return to host



	•
	•
Extension Duntime	•
	•

40

### **EIM: Extension Interface Model**

- Solution to nav fine-grained safety-interconnectedness trade-offs
- Two-Phase Specification
  - Development-Time (by Developer)
  - Deployment-Time (by Manager)
- Capabilities as Resources



### **EIM:** Development-Time Specification

- Developers annotate code for capabilities
- Automatically extracted into capability manifest



```
EIM_STATE_DEFINE(readPid, read, ngx_pid);
EIM_HFUNC_DEFINE_WITH_CONSTRAINTS(
    nginxTime,
    HF_RET_POSITIVE);
EIM_EXTENSION_ENTRY_DEFINE(
    processBegin,
    ngx_http_process_request,
    int,
    struct Request *);
```

### **EIM: Extension Interface Model**



### bpftime: userspace eBPF extension framework

- Challenge for compatibility and efficiency:
  - eBPF: tightly coupled components
  - Bpftime: Intercept syscalls & Share memory maps



















## bpftime: userspace eBPF extension framework

Provide efficient

solution to enforce EIM and isolation

- Verification: EIM: no runtime cost
- Execution extension
   runtime in thesame process
   forefficiency
- Conceal extension entry for efficiency: using binary rewriting to remove unused extension entries
- Hardware features for efficient isolation
- o Compatibility
- . . . .





## **EIM:** Development-Time Specification

- Developers annotate code for capabilities
- Automatically extracted into capability manifest



```
EIM_STATE_DEFINE(readPid, read, ngx_pid);
EIM_HFUNC_DEFINE_WITH_CONSTRAINTS(
    nginxTime,
    HF_RET_POSITIVE);
EIM_EXTENSION_ENTRY_DEFINE(
    processBegin,
    ngx_http_process_request,
    int,
    struct Request *);
```

## **EIM: Deployment-Time Specification**

 Extension Manager write simple YAML policies to explore interconnectedness/safety trade-offs without recompiling





### contribution

- Background & motivation: Extensions
- Extension Interface Model (EIM): Fine-grained Interface
- bpftime Runtime: safety & performance
- Evaluation
- ..... add more and make it a contribution

# Q & A?

- GitHub repo:
- Get started:

bpftime load ./example/malloc/malloc
bpftime start nginx -c ./nginx.conf



### Contribution

• Extension Interface Model (EIM):

Solution to navigate fine-grained safety-interconnectedness tradeoff & Two-Phase Specification: Development-Time (by Developer) Deployment-Time (by Manager) • **bpftime Runtime**: An userspace eBPF runtime implemented EIM with isolation and efficiency

• Evaluation: 6 usecases and Up to 6x less overhead

# Nginx firewall example

Offline:

- Writes custom firewall logic using nginx helper functions
- Loads their extension at an extension entry for request processing

At runtime

- Nginx jumps to the extension runtime when reaches the extension entry
- Extension runtime execute the extension entry and return to Nginx



## Six Real-World Use Cases

- Nginx Firewall
- Redis Durability
- FUSE Metadata Cache
- DeepFlow
- Syscount
- Sslsniff

(grounp and figture?)

GitHub: <u>https://github.com/eunomia-bpf/bpftime</u>



### **Extension Problems**

#### • Real-world safety violations:

Bilibili CDN outage, Apache buffer overflow, Redis RCE

 Performance penalty: WebAssembly/Lua impose 10-15% overhead



### contribution

- Background & motivation: Extensions
- Extension Interface Model (EIM): Fine-grained Interface
- bpftime Runtime: safety & performance
- Evaluation
- ..... add more and make it a contribution

# Q & A?

- GitHub repo:
- Get started:

bpftime load ./example/malloc/malloc
bpftime start nginx -c ./nginx.conf



### **EIM: Extension Interface Model**

- Goal: enable fine-grained safety/interconnectedness trade-offs
- Challenge: supporting per deployment tradeoffs
- Solution:

 $_{\odot}$  Two-Phase Specification (Development-Time and deployment-Time)  $_{\odot}$  Model all resources as capabilities



During Development

Before Deployment

At Deployment/Runtime

### Contributions

• Extension Interface Model (EIM)

EIM

Solution to navigate fine-grained safety-interconnectedness trade-off



bpftime



Up to 6x less overhead than current state-of-the-art!

### Contributions



Two phase Specification

Up to 6x less overhead than current state-of-the-art!

### **EIM: Extension Interface Model**

- Goal: enable fine-grained safety/interconnectedness trade-offs
- Challenge: supporting per deployment tradeoffs
- Solution:

 $_{\odot}$  Two-Phase Specification (Development-Time and deployment-Time)  $_{\odot}$  Model all resources as capabilities



### **Extension Requirements**

- Fine-grained safety and interconnected ness trade-offs
- Isolation:
- Efficiency:



### **EIM: Extension Interface Model**

- Goal: enable fine-grained safety/interconnectedness trade-offs
- Challenge: supporting per deployment tradeoffs
- Solution:

 $_{\odot}$  Two-Phase Specification (Development-Time and deployment-Time)  $_{\odot}$  Model all resources as capabilities



1. During Development

2. Before Deployment

3. At Deployment/Runtime

## **Extension Problems & requirements**

# Real-world safety violations:

• Bilibili CDN outage, Apache buffer overflow, Redis RCE

### **Performance penalty:**

 WebAssembly/Lua impose 10-15% overhead

### **Requirements:**

- Fine-grained safety and interconnectedness trade-offs
- Isolation
- Efficiency

