

bpfbbox: Simple Precise Process Confinement with eBPF

William Findlay

Anil Somayaji

David Barrera

Carleton University

will@ccsl.carleton.ca

November 5, 2020



Outline of this Talk

1. The Process Confinement Problem
2. The Status Quo
3. eBPF 101
4. bpfbox Overview
5. bpfbox Design and Implementation
6. bpfbox Performance Evaluation
7. The Future of eBPF and Security
8. Conclusion

The Process Confinement Problem

What is Process Confinement?

We want to be able to **confine** our **processes**.

- ▶ Also known as *sandboxing*

Why do we want to do this?

- ▶ Default protection mechanisms are too:
 - ▶ Coarse-grained
 - ▶ User-centric
 - ▶ Discretionary
- ▶ DAC can be **overridden**
 - ▶ Superuser (*nix)
 - ▶ Administrator (Windows)

How do we protect the user from *themselves* and *their own processes*?

An Interesting Question

How many processes do you think are running on your computer **right now**?

- ▶ Probably a lot more than you think
- ▶ You probably didn't start most of them yourself
- ▶ You might not even know what some of them are for
- ▶ **Do you trust them?**

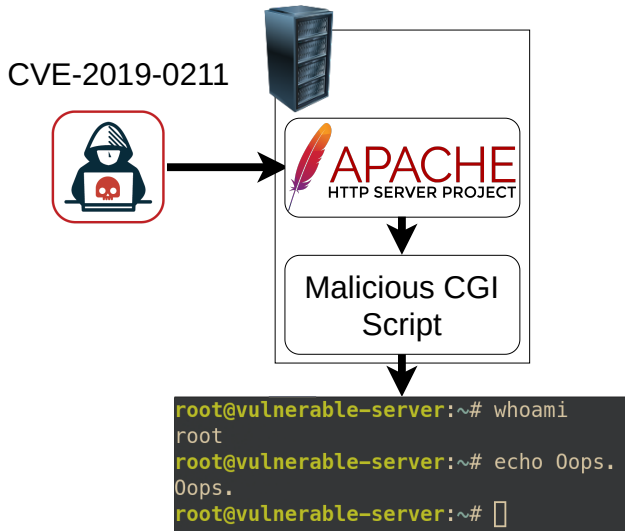
Threat Model

Compromised Processes

- ▶ Web servers
- ▶ Daemons
- ▶ Chat applications
- ▶ etc.

The Morris Worm

- ▶ Backdoor in sendmail daemon
- ▶ Buffer overflow in fingerd
- ▶ Estimated damage:
\$100,000–\$10,000,000

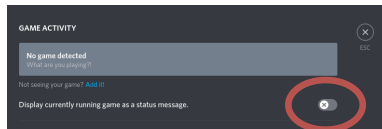


Threat Model

Semi-Honest Software

- ▶ Software that does its job...
- ▶ But also performs potentially **unwanted actions**

```
[16:13] hh@arch:~ | strace -fqo /tmp/discord.log discord-canary
```

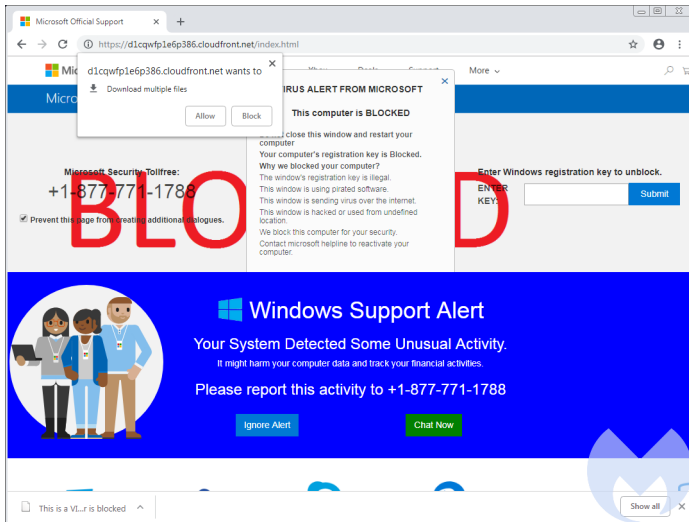


```
912433 readlink("/proc/506824/exe", "/opt/draw.io/drawio", 1023) = 19
912433 readlink("/proc/634535/exe", "/usr/lib/firefox/firefox", 1023) = 24
912433 readlink("/proc/909407/exe", "/usr/bin/alacrity", 1023) = 18
912433 readlink("/proc/3205375/exe", "/usr/bin/tmux", 1023) = 13
912433 readlink("/proc/3244098/exe", "/usr/bin/zsh", 1023) = 12
912433 readlink("/proc/3387/exe", "/usr/bin/picom", 1023) = 14
912433 readlink("/proc/3412/exe", "/usr/bin/xbindkeys", 1023) = 18
912433 readlink("/proc/3421/exe", "/usr/bin/volnoti", 1023) = 16
912433 readlink("/proc/3449/exe", "/usr/bin/pulseaudio", 1023) = 19
```

Threat Model

Malicious Software

- ▶ Viruses
- ▶ Trojans
- ▶ Ransomware
- ▶ Spyware
- ▶ etc.



Threat Model

Attack goals?

- ▶ Installing backdoors/rootkits
- ▶ Information leakage
- ▶ Denial of service
- ▶ Data ransom
- ▶ Setting up a botnet

Process confinement
reduces the attack
surface.

The Process Confinement Problem

- ▶ “A Note on the Confinement Problem” (Lampson, 1973)

Operating
Systems

C. Weissman
Editor

A Note on the Confinement Problem

Butler W. Lampson
Xerox Palo Alto Research Center

- ▶ A (mostly) open problem for nearly **five decades**

The Status Quo

Unix DAC

- ▶ **User-centric** permissions
- ▶ Permission bits and ACLs

User			Group			Other		
1	1	1	1	0	1	1	0	1
r	w	x	r	-	x	r	-	x

Problems?

- ▶ My own processes can still access all my files
- ▶ Abuse of DAC for “confinement”
- ▶ root-owned processes ignore everything
- ▶ `chmod -R 777 .` (I've seen too many COMP3000 students do this)

POSIX Capabilities

- ▶ All-or-nothing superuser privileges are a problem
- ▶ Split them up into capabilities
 - ▶ `CAP_DAC_OVERRIDE` (override DAC)
 - ▶ `CAP_CHOWN` (change file owners)
 - ▶ `CAP_NET_BIND_SERVICE` (bind to privileged ports)
 - ▶ etc.

Pick your poison:

- ▶ Replace SUID binaries with capabilities as xattrs
- ▶ Drop capabilities with `prctl(2)`

POSIX Capabilities

Problems?

- ▶ Dropping capabilities requires **modifying applications**
- ▶ **Complicates** the Linux permission model (40 new permission bits and counting)
- ▶ Doesn't really **solve** process confinement
 - ▶ Unprivileged processes are still a problem
 - ▶ A process under your own UID can still access your files, do networking, read keyboard input, etc.

Namespaces and Cgroups

Namespaces

- ▶ Virtualize enumerable resources
- ▶ Give a process group a **private view** of the resources
- ▶ PID namespace, UID namespace, Mount namespace, etc.

Cgroups

- ▶ Limit availability of system resources
- ▶ Memory, CPU, etc.

Namespaces and Cgroups

Problems?

- ▶ Not so easy for end users to configure
- ▶ Not application transparent on their own
- ▶ Not a full confinement implementation
 - ▶ We have **virtualization**, but we don't have **least-privilege**
 - ▶ Needs to be combined with *something else*

System Call Interposition

Linux seccomp-bpf (**NOT** eBPF)

- ▶ Processes call `seccomp(2)` to enter a secure mode
- ▶ All system calls are denied except:
 - ▶ `read(2)`
 - ▶ `write(2)`
 - ▶ `sigreturn(2)`
 - ▶ `exit(2)`
- ▶ Applications can write classic BPF programs to allow additional system calls

Problems?

- ▶ Not application-transparent
- ▶ Classic BPF is **arcane**
- ▶ Equivalent system calls (e.g. `open(2)` and `openat(2)`)
- ▶ Do users understand system call semantics?

System Call Interposition

OpenBSD `pledge`

- ▶ Group system calls into higher-level, meaningful categories
- ▶ For example, `stdio` includes `read(2)` and `write(2)`
- ▶ A process commits to what it wants to use before calling `pledge(2)`

Problems?

- ▶ Not application-transparent
- ▶ Too **coarse-grained**
- ▶ A process can escape the sandbox with `execve(2)`

System Call Interposition

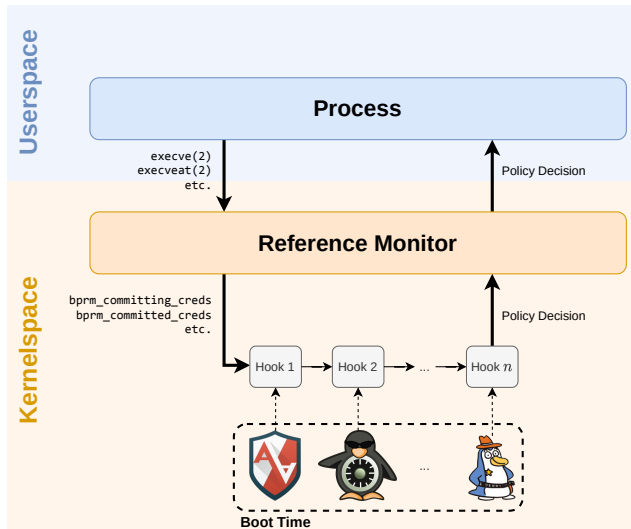
FreeBSD capsicum

- ▶ Implements *capabilities* (**NOT** POSIX capabilities)
- ▶ A process receives capabilities from the OS as file descriptors
- ▶ These capabilities restrict access to global namespaces
- ▶ The process calls `cap_enter` to enter *capability mode*

Problems?

- ▶ Not application-transparent
- ▶ More usable than `seccomp-bpf`, but still meant for security experts

Linux MAC



Linux MAC

SELinux sucks.

apache.te

```
apache_content_template(user)
ubac_constrained(httpd_user_script_t)
userdom_user_home_content(httpd_user_content_t)
userdom_user_home_content(httpd_user_htaccess_t)
userdom_user_home_content(httpd_user_script_exec_t)
userdom_user_home_content(httpd_user_ra_content_t)
userdom_user_home_content(httpd_user_rw_content_t)
```

apache.fc

```
/etc/apache(2)?(/.*)? gen_context(system_u:object_r:httpd_config_t,s0)
/etc/apache-ssl(2)?(/.*)? gen_context(system_u:object_r:httpd_config_t,s0)
/etc/cherokee(/.*)? gen_context(system_u:object_r:httpd_config_t,s0)
/etc/nginx(/.*)? gen_context(system_u:object_r:httpd_config_t,s0)
/etc/hiawatha(/.*)? gen_context(system_u:object_r:httpd_config_t,s0)
/etc/httpd(/.*)? gen_context(system_u:object_r:httpd_config_t,s0)
```

apache.if

```
can_exec(httpd_$1_script_t, httpd_$1_script_exec_t)

allow httpd_$1_script_t httpd_$1_ra_content_t:dir { list_dir_perms add_ent
allow httpd_$1_script_t httpd_$1_ra_content_t:file { append_file_perms rea
setattr_file_perms };
allow httpd_$1_script_t httpd_$1_ra_content_t:lnk_file read_lnk_file_perms

allow httpd_$1_script_t { httpd_$1_content_t httpd_$1_script_exec_t }:dir
allow httpd_$1_script_t httpd_$1_content_t:file read_file_perms;
allow httpd_$1_script_t { httpd_$1_content_t httpd_$1_script_exec_t }:lnk_
```

- ▶ Each file on the left is actually *thousands* of lines of the same nonsense
- ▶ This problem is generalizable across MAC implementations
- ▶ Policy is designed to be written by **security experts**
- ▶ Not suitable for ad-hoc confinement

Containers / Containerized Package Management

High level policy.

- ▶ Package maintainers write coarse-grained package manifests
- ▶ Users supply command line arguments

Complex enforcement.

- ▶ Virtualization with namespaces, cgroups, filesystem mounts
- ▶ Least-privilege with seccomp-bpf, SELinux, AppArmor
- ▶ **Whole userlands** need to be secured for **each application**

Containers / Containerized Package Management

Problems?

- ▶ Overpermission:
 - ▶ Permissions are **overly-generalized**, not application specific
- ▶ Auditability:
 - ▶ **Four** line package manifest
 - ▶ **Thousands** of lines of AppArmor/seccomp-bpf
- ▶ Usability:
 - ▶ What command line arguments do I use for application *x*?
 - ▶ What if I want to write my own policy?

Main Takeaways

1. Process confinement is **hard to get right**.
2. Trade-off between **usability** and **security**.
3. Trade-off between **terseness** and **expressiveness**.

Process confinement is *not* a solved problem.

eBPF 101

eBPF in the Beginning

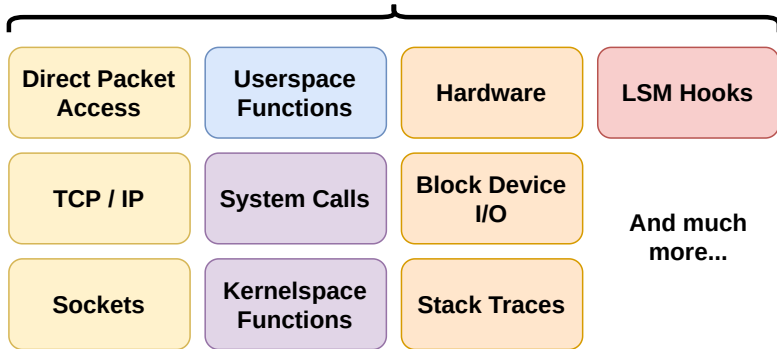
eBPF \equiv **E**xtended **B**erkley **P**acket **F**ilter

- ▶ But it has little to do with Berkley, packets, or filtering nowadays
- ▶ The name BPF is preserved for historical reasons

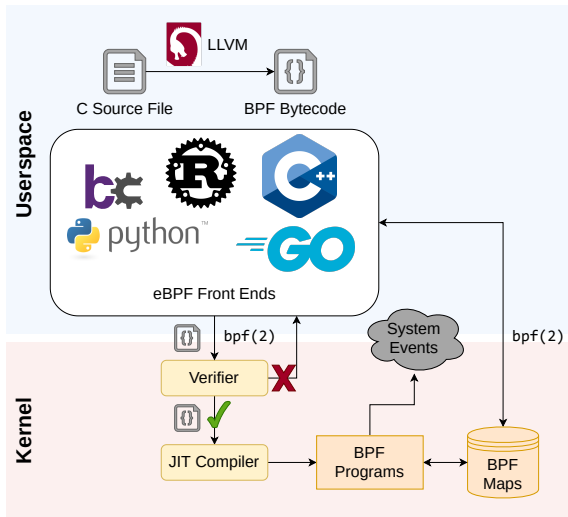
So then **what is eBPF?**

- ▶ A major re-write of the Linux BPF engine
 - ▶ Alexei Starovoitov and Daniel Borkmann
- ▶ Merged into the Linux kernel in 2014
- ▶ The original goal was fine-grained, cross-layer **system introspection**

What Can eBPF Do?



How eBPF Works



Verifiably Safe Programs

Limited instruction set.

- ▶ 11 registers (10 general purpose)
- ▶ 114 instructions (vs 2000+ in x86)
- ▶ Access to a limited set of **kernel helpers** with `call` instruction

Restricted execution context.

- ▶ 512 byte stack limit
- ▶ Memory access must be bounds-checked
- ▶ No unbounded loops
- ▶ No back-edges in control flow
- ▶ **eBPF is not Turing complete**

eBPF in 2020

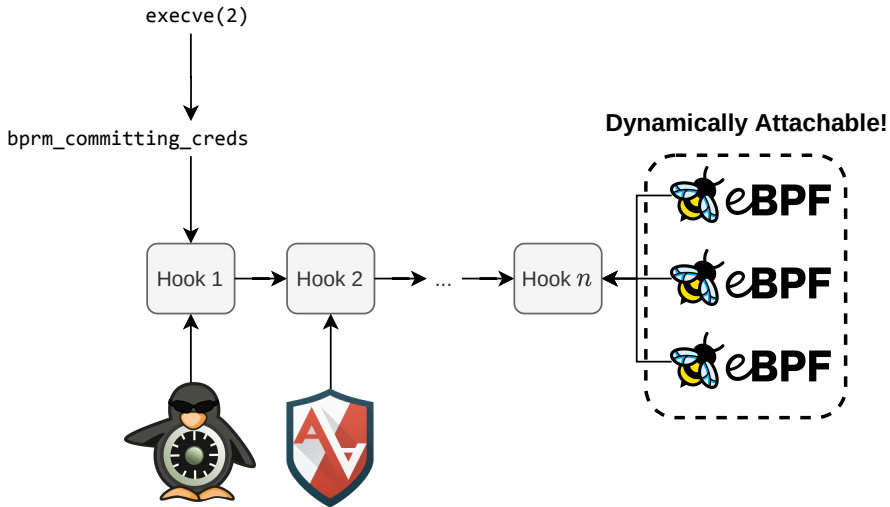
eBPF is now **more than just an observability tool**.

- ▶ eBPF provides a **safe, efficient, and flexible** way for privileged users to extend the kernel
- ▶ eBPF turns Linux into a **programmable kernel**

Linux 5.7 → KRSI (**K**ernel **R**untime **S**ecurity **I**nstrumentation)

- ▶ Attach BPF programs to LSM hooks
- ▶ Make security decisions and generate audit logs with eBPF

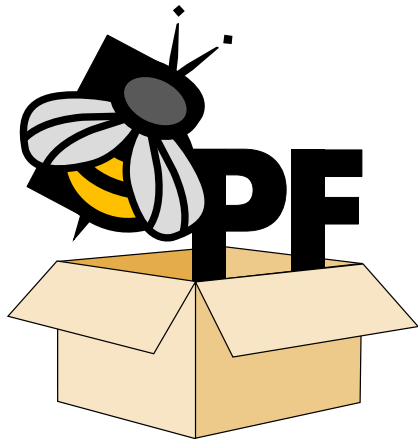
KRSI: BPF LSM Programs



bpfbox Overview

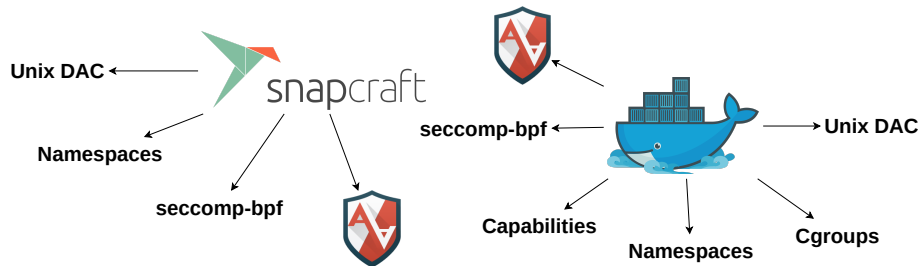
bpfbbox at a Glance

- ▶ bpfbbox is a novel **process confinement mechanism** for Linux
 - ▶ Using a new Linux technology called eBPF
- ▶ Users write per-application policy in a **simple policy language**
- ▶ Policy is enforced by attaching **eBPF programs** to **LSM hooks**
 - ▶ Integrates cross-layer state into policy decisions



Motivation

- ▶ Existing process confinement mechanisms are **complex**



- ▶ Existing process confinement mechanisms are **difficult to use**



SELinux



AppArmor



TOMOYO

- ▶ Can we do any better?

eBPF Changes the Game

eBPF enables:

- ▶ Fine-grained system introspection
- ▶ Integration of **cross-layer state** with policy enforcement
- ▶ Rapid prototyping
- ▶ Safe production deployment of new security solutions

We have an opportunity to **rethink process confinement** from the ground up.

bpfbbox Design and Implementation

Policy Design Goals

1. Simplicity

- ▶ Policy should be simple enough for ad hoc confinement

2. Application transparency

- ▶ Policy should not require changes to the confined application

3. Flexibility

- ▶ Policy should offer optional layers of granularity

4. Security

- ▶ Policy should follow the principle of least privilege
- ▶ It should be difficult to write an insecure policy

Rules and Directives

Rules specify access to system objects:

- ▶ `fs(file, access)`
- ▶ `net(socket, access)`
- ▶ `signal(prog, sig)`
- ▶ etc.

Directives augment blocks of rules:

- ▶ `#[directive]` syntax
- ▶ Specify **actions to be taken** on a block of rules
- ▶ Add **additional context** to a block of rules

Taints and Transitions

- ▶ `#[taint]` → Start confinement
- ▶ `#[transition]` → Switch profiles on `execve`

```
#![profile "/bin/mywebdaemon"]

#[taint] {
    net(inet, any)
    net(inet6, any)
}

/* ... */

#[transition] {
    fs("/bin/myhelper", getattr|read|exec)
}
```

Policy at the Function Call Level

- ▶ `#[func "foo"]` → Apply rules only within a call to `foo()`
- ▶ `#[kfunc "foo"]` → Same thing, but for kernel functions

```
#[profile "/sbin/mylogin"]

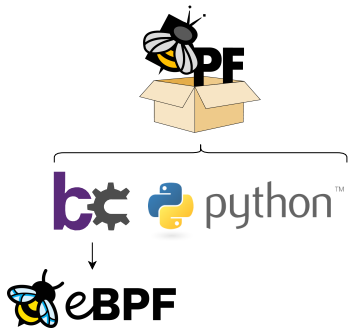
#[func "check_password"]
#[allow] {
    fs("/etc/passwd", read)
    fs("/etc/shadow", read)
}

#[func "add_user"]
#[allow] {
    fs("/etc/passwd", read|append)
    fs("/etc/shadow", read|append)
}

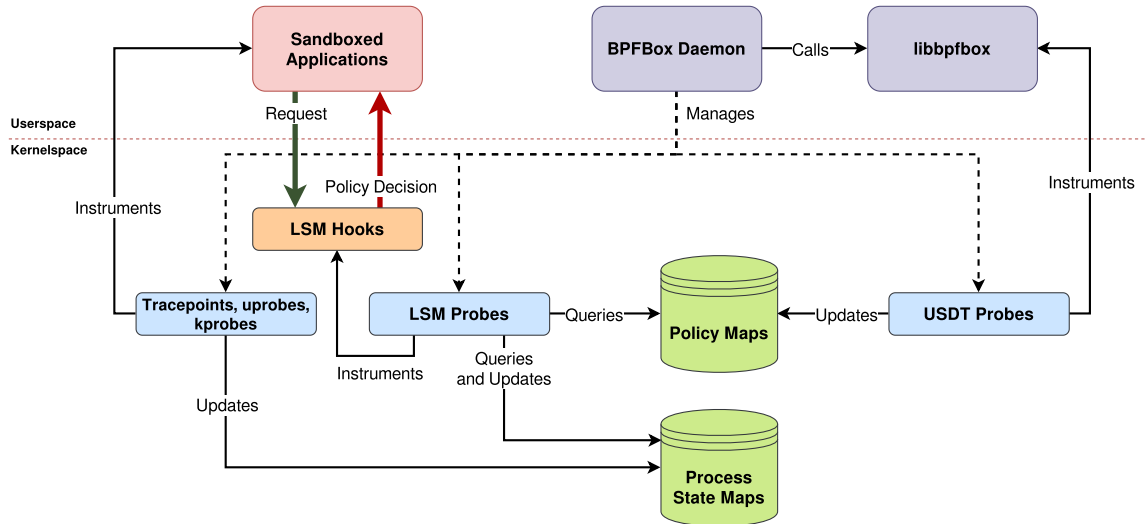
/* ... */
```


bpffbox Implementation

- ▶ Userspace daemon using the Python3 bcc framework
- ▶ Kernelspace components are all written in eBPF
 - ▶ LSM probes (KRSI), kprobes, uprobes
 - ▶ Under 2000 source lines of kernelspace code
- ▶ Thanks to eBPF, bpffbox is **light-weight**, **flexible**, and **production-safe**
 - ▶ Works out of the box on any vanilla Linux kernel ≥ 5.8



bpfbox Architecture



bpfbox Performance Evaluation

Methodology

- ▶ Phoronix Test Suite OSBench
 - ▶ Measures basic OS functionality
 - ▶ (spawning processes, memory allocations, etc.)
- ▶ Phoronix Test Suite Apache
 - ▶ Benchmark Apache `httpd` packets per second
- ▶ Kernel compilation benchmarks
 - ▶ Measure Linux kernel compilation performance
 - ▶ Heavy workload, spawning lots of processes

Methodology

Two modes of operation for each test.

- ▶ Passive mode
 - ▶ bpfbox and AppArmor instrument hooks, but do not enforce or audit
 - ▶ Test lowest possible overhead
- ▶ Complaining mode
 - ▶ bpfbox and AppArmor complain about (log) every security-sensitive operation
 - ▶ Test worst case overhead

Results

- ▶ Phoronix OSBench
 - ▶ Passive: bpfbox is **roughly equivalent** to AppArmor
 - ▶ Complaining: bpfbox performs **significantly better** than AppArmor
- ▶ Phoronix Apache
 - ▶ bpfbox and AppArmor are **roughly equivalent**
- ▶ Kernel compilation
 - ▶ Passive: bpfbox is **roughly equivalent** to AppArmor
 - ▶ Complaining: bpfbox performs **better in kernelspace** overhead and **worse in userspace** overhead

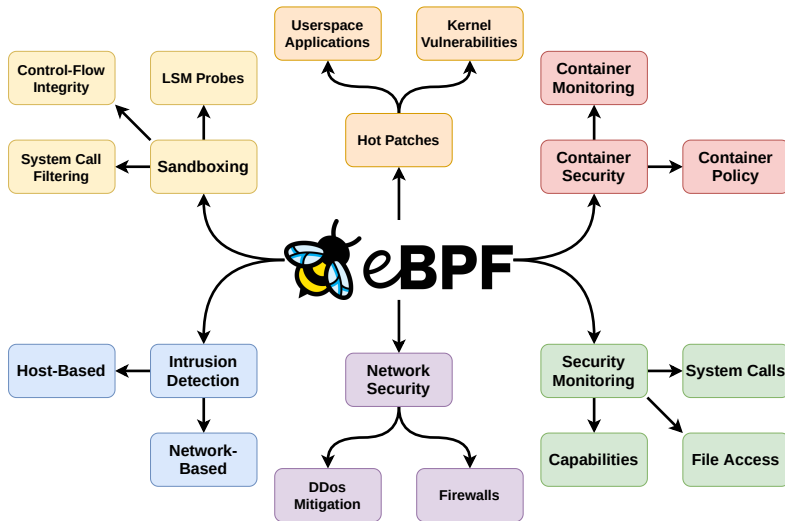
Results

- Results of the benchmarks
- Percent differences from the base are given in parentheses

	Base	Passive				Permissive			
		bpfbox		AppArmor		bpfbox		AppArmor	
Phoronix OSBench (lower is better):									
Create Files (μ s)	27.86	28.94	(3.81%)	28.01	(0.55%)	32.31	(14.80%)	96.56	(110.44%)
Create Threads (μ s)	25.96	26.90	(3.56%)	26.28	(1.24%)	27.67	(6.39%)	26.09	(0.51%)
Launch Programs (μ s)	75.12	78.02	(3.79%)	77.64	(3.30%)	87.31	(15.01%)	102.43	(30.76%)
Create Processes (μ s)	51.32	52.53	(2.34%)	51.61	(0.57%)	51.85	(1.04%)	52.11	(1.54%)
Memory Allocations (ns)	113.98	112.33	(-1.45%)	112.29	(-1.50%)	112.75	(-1.09%)	112.74	(-1.09%)
Kernel Compilation (lower is better):									
User (s)	14457.01	14564.80	(0.74%)	14711.42	(1.74%)	14829.11	(2.54%)	14432.09	(-0.17%)
System (s)	1712.59	1760.02	(2.73%)	1765.69	(3.05%)	1804.10	(5.20%)	2544.72	(39.09%)
Elapsed (s)	2086.92	2114.83	(1.33%)	2130.38	(2.06%)	2397.48	(13.85%)	2261.09	(8.01%)
Phoronix Apache (higher is better):									
Requests Per Second (r/s)	14686.95	13887.59	(-5.59%)	13743.88	(-6.63%)	13504.23	(-8.39%)	13431.34	(-8.93%)

The Future of eBPF and Security

Security Applications of eBPF



New Directions

Userspace LSM (Self-Confinement)

- ▶ Attach uprobes to a shared library
- ▶ Userspace applications make calls to the library to declare privileges
- ▶ uprobes update a policy map in kernelspace

Dynamic Capabilities

- ▶ Users define custom capabilities
- ▶ Enforced in kernelspace with dynamic LSM probes
- ▶ E.g. CAP_ACCESS_PHOTOS to grant access to ~/pictures

New Directions

Hot Patches (Userspace)

- ▶ Patch vulnerabilities before security updates are available
- ▶ uprobes to hook into functions
- ▶ `bpf_probe_write_user()` to replace userspace memory

Hot Patches (Kernel)

- ▶ Replace vulnerable kernel functions with BPF programs
- ▶ Alter/drop malicious packets before they reach the networking stack
- ▶ E.g. patch packet-of-death vulnerability with an XDP program

Conclusion

bpfbox Future Work

- ▶ Consider alternative policy languages
 - ▶ yaml?
 - ▶ rego?
- ▶ Incorporate new kernel features
 - ▶ `task_local_storage`, `inode_local_storage` (Linux 5.10)
 - ▶ Boot-time loading BPF programs (Linux 5.9)
- ▶ Container integration?
 - ▶ Enforce policy at the container level

Acknowledgements

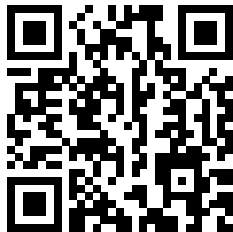
Special thanks to:

- ▶ **Alexei Starovoitov** and **Daniel Borkmann** (creators of eBPF)
- ▶ **K.P. Singh** (creator of KRSI)
- ▶ Fellow **bcc contributors** (an awesome eBPF framework)
- ▶ Anonymous **CCSW'2020 reviewers** (valuable feedback)

This work was supported by NSERC through a Discovery Grant.

Contributions

- ▶ First **policy enforcement engine** written in **eBPF**
- ▶ Integration of **userspace** and **kernelspace** state with **LSM layer enforcement**
- ▶ A simple policy language for **ad hoc process confinement**
 - ▶ But with optional complexity for **fine-grained protection**



github.com/willfindlay/bpfbbox

Check out the project on GitHub!