



Red Hat



Carnegie  
Mellon  
University



# Revisiting eBPF Seccomp Filters

Jinghao Jia

University of Illinois at Urbana-Champaign

LPC '22 (Networking & BPF Track)



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

with YiFei Zhu, Andrea Arcangeli, Hubertus Franke, Tobin Feldman-Fitzthum, Claudio Canella, Dimitrios Skarlatos, Daniel Gruss, Dan Williams, Tianyin Xu

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Seccomp

- The syscall interposition mechanism for Linux
  - Restricts the set of syscalls for an application
  - First implemented by Andrea Arcangeli and merged in kernel 2.6.12



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Seccomp

- The syscall interposition mechanism for Linux
  - Restricts the set of syscalls for an application
  - First implemented by Andrea Arcangeli and merged in kernel 2.6.12
- Two modes



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Seccomp

- The syscall interposition mechanism for Linux
  - Restricts the set of syscalls for an application
  - First implemented by Andrea Arcangeli and merged in kernel 2.6.12
- Two modes
  - *Restricted mode*: only allows `read()`, `write()`, `exit()`, and `sigreturn()`
    - Kills the process if other syscalls are used



A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Seccomp

- The syscall interposition mechanism for Linux
  - Restricts the set of syscalls for an application
  - First implemented by Andrea Arcangeli and merged in kernel 2.6.12
- Two modes
  - *Restricted mode*: only allows `read()`, `write()`, `exit()`, and `sigreturn()`
    - Kills the process if other syscalls are used
  - *Filter mode*: allows custom policies implemented in **cBPF** (classic BPF)
    - Custom set of allowed/disallowed syscalls and non-pointer arguments
    - Custom actions (e.g., allow, return error code, kill process)



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative green pipe graphic runs along the top and left edges of the slide. It features various fittings, elbows, and a valve. The pipe is rendered in a bright green color with a slight shadow effect.

# cBPF is simple

- Small, RISC-like instruction set
- Only allows forward jumps
- Limited to 4096 instructions in length



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, running along the top, left, and bottom edges of the slide. The pipes are connected at various points, forming a complex but clean layout.

# cBPF is simple

- Small, RISC-like instruction set
- Only allows forward jumps
- Limited to 4096 instructions in length

Easy to verify safety properties



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the central text and logo.

Seccomp-cBPF is widely used

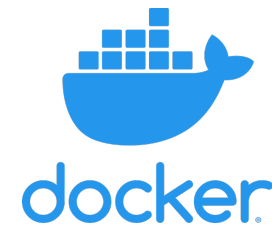


Linux  
Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**



# Seccomp-cBPF is widely used

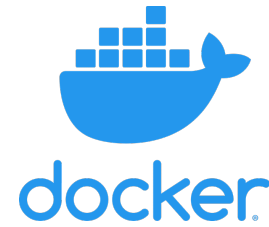
- Virtualization/Containers
  - Docker, Rkt, gVisor, Kubernetes, etc



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Seccomp-cBPF is widely used

- Virtualization/Containers
  - Docker, Rkt, gVisor, Kubernetes, etc
- Sandboxing
  - Systemd, Android



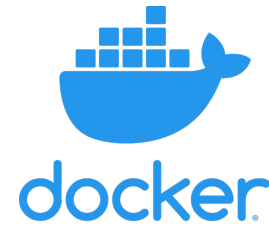
[ ● ◀ ] systemd



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Seccomp-cBPF is widely used

- Virtualization/Containers
  - Docker, Rkt, gVisor, Kubernetes, etc
- Sandboxing
  - Systemd, Android
- Various other applications
  - OpenSSH, Tor, Firefox, ...



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide. The pipes are rendered in a bright green color with a slight shadow effect.

# Limitation of cBPF

 **Linux**  
**Plumbers Conference** | Dublin, Ireland **Sept. 12-14, 2022**

A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide.

# Limitation of cBPF

- **Stateless**: only support **static** allow/deny lists



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, running along the top, left, and bottom edges of the slide.

## Limitation of cBPF

- **Stateless**: only support **static** allow/deny lists
- **Size limit**: complex policies cannot be encoded in one filter [1]



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, running along the left and top edges of the slide.

## Limitation of cBPF

- **Stateless**: only support **static** allow/deny lists
- **Size limit**: complex policies cannot be encoded in one filter [1]
- **Simple instruction set**: cannot use advanced utilities



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide.

## Limitation of cBPF

- **Stateless**: only support **static** allow/deny lists
- **Size limit**: complex policies cannot be encoded in one filter [1]
- **Simple instruction set**: cannot use advanced utilities
- **Programmability**: cannot use modern eBPF tools



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide.

# Limitation of cBPF

- **Stateless**: only support **static** allow/deny lists
- **Size limit**: complex policies cannot be encoded in one filter [1]
- **Simple instruction set**: cannot use advanced utilities
- **Programmability**: cannot use modern eBPF tools
  - Requires writing cBPF instructions by hand (mitigated by libseccomp)



A decorative graphic of a green pipe network with various fittings, elbows, and valves, running along the top, left, and bottom edges of the slide.

## Limitation of cBPF

- **Stateless**: only support **static** allow/deny lists
- **Size limit**: complex policies cannot be encoded in one filter [1]
- **Simple instruction set**: cannot use advanced utilities
- **Programmability**: cannot use modern eBPF tools
  - Requires writing cBPF instructions by hand (mitigated by libseccomp)

cBPF cannot express complex syscall policies,  
e.g., state-based filtering



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the central text. The pipes are a vibrant green color and are set against a light gray shadow.

# Revisiting eBPF Seccomp filters

to improve *programmability* and *extensibility* of system call security.



Linux  
Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**

A decorative graphic of green pipes with various fittings, elbows, and valves, running vertically and horizontally around the slide content.

# Advances of eBPF

A new level of programmability



Linux

Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**



# Advances of eBPF

A new level of programmability

- **Richer instruction set and more flexible control flow**



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Advances of eBPF

A new level of programmability

- **Richer instruction set and more flexible control flow**
- **eBPF map structures** for efficient storage



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022





# Advances of eBPF

A new level of programmability

- **Richer instruction set and more flexible control flow**
- **eBPF map structures** for efficient storage
- **Helper interface** for rich interaction with kernel



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Advances of eBPF

A new level of programmability

- **Richer instruction set and more flexible control flow**
- **eBPF map structures** for efficient storage
- **Helper interface** for rich interaction with kernel
- **Toolchain support:** write C, compile to eBPF



Linux  
Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content. The pipes are bright green with a slight shadow effect.

# Implications on Seccomp



Linux  
Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**

A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the slide content.

# Implications on Seccomp

- Advanced checking logic from new instruction set



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the slide content.

# Implications on Seccomp

- Advanced checking logic from new instruction set
- Stateful filters via map storage



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative green pipe graphic runs vertically along the left side of the slide, featuring various fittings, valves, and elbows. It starts at the top left, goes down, then right, then down again, and finally right at the bottom. 

# Implications on Seccomp

- Advanced checking logic from new instruction set
- Stateful filters via map storage
- Access to more kernel context using kernel helper functions



A decorative graphic of green pipes with various fittings, valves, and elbows, running along the left and top edges of the slide. The pipes are rendered in a 3D style with shadows.

# Implications on Seccomp

- Advanced checking logic from new instruction set
- Stateful filters via map storage
- Access to more kernel context using kernel helper functions
- Mature toolchains



Linux

Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**



# Previous efforts

- Seccomp-eBPF patch from Dhillon (Feb. 2018)
  - <https://lwn.net/Articles/747229/>
  - + A new Seccomp program type for eBPF
  - + A new eBPF flag in Seccomp with code that allocates eBPF filters
  - No eBPF map support nor security model
  - Only supports 4 helpers
- Seccomp-eBPF patch from Hromatka (Feb. 2018)
  - <https://groups.google.com/g/libseccomp/c/pX6QkVFoF74?pli=1>
  - + A new eBPF filter mode in Seccomp that loads and allocates the filter.
  - Incomplete -- no Seccomp-eBPF program type/whitelisted helpers



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Some common opinions



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, running along the top, left, and bottom edges of the slide.

# Some common opinions

- “What's the reason for adding eBPF support? seccomp shouldn't need it, and it only makes the code more complex. I'd rather stick with cBPF until we have an **overwhelmingly good reason** to use eBPF as a "native" seccomp filter language.”

Lacking use cases



# Some common opinions

- “What's the reason for adding eBPF support? seccomp shouldn't need it, and it only makes the code more complex. I'd rather stick with cBPF until we have an **overwhelmingly good reason** to use eBPF as a "native" seccomp filter language.”

Lacking use cases

- This is the blocker as far as I'm concerned: there is no story for unprivileged eBPF. And even IF there was a story there, I find **the rate of security-related flaws in eBPF to be way too high for a sandboxing primitive to depend on...** I just can't bring myself to accept that level of risk for seccomp.

Security as a concern



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe system with various fittings, valves, and elbows, running along the top and left sides of the slide.

# Many use cases of Seccomp-eBPF

- Syscall count limiting
- Enhanced temporal syscall specialization
- Syscall-flow-integrity protection
- Checking syscall sequence for intrusion detection systems
- Accelerating security checks
- Many others!



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe system with various fittings, elbows, and valves, running along the top and left edges of the slide.

# Syscall count limiting

Real world example in runc



Linux

Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**



# Syscall count limiting

Real world example in runc

```
if l.cfg.Config.Seccomp != nil &&  
    !l.cfg.NoNewPrivileges {  
    seccomp.InitSeccomp(...)  
}
```

Sets up Seccomp filters





# Syscall count limiting

## Real world example in runc

```
if l.cfg.Config.Seccomp != nil &&  
    !l.cfg.NoNewPrivileges {  
    seccomp.InitSeccomp(...)  
}  
finalizeNamespace(...)  
...  
if unix.Getppid() != l.parentPid { ... }  
...  
unix.Write(fd, []byte("0"))
```

Sets up Seccomp filters

namespace & privilege setup, etc



# Syscall count limiting

## Real world example in runc

```
if l.cfg.Config.Seccomp != nil &&  
    !l.cfg.NoNewPrivileges {  
    seccomp.InitSeccomp(...)  
}
```

Sets up Seccomp filters

```
finalizeNamespace(...)
```

```
...  
if unix.Getppid() != l.parentPid { ... }
```

namespace & privilege setup, etc

```
...  
unix.Write(fd, []byte("0"))  
system.Exec(name, l.cfg.Args[0:], os.Env())
```

exec into target application



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Syscall count limiting

## Real world example in runc

```
if l.cfg.Config.Seccomp != nil &&  
    !l.cfg.NoNewPrivileges {  
    seccomp.InitSeccomp(...)  
}  
finalizeNamespace(...)  
...  
if unix.Getppid() != l.parentPid { ... }  
...  
unix.Write(fd, []byte("0"))  
system.Exec(name, l.cfg.Args[0:], os.Env())
```

seccomp w/o NO\_NEW\_PRIVS is privileged!  
Needed before dropping privileges



# Syscall count limiting

## Real world example in runc

```
if l.cfg.Config.Seccomp != nil &&
    !l.cfg.NoNewPrivileges {
    seccomp.InitSeccomp(...)
}
finalizeNamespace(...)
...
if unix.Getppid() != l.parentPid { ... }
...
unix.Write(fd, []byte("0"))
system.Exec(name, l.cfg.Args[0:], os.Env())
```

syscalls for later setups need to be allowed  
e.g. execve, prctl, capset



# Syscall count limiting

## Real world example in runc

```
if l.cfg.Config.Seccomp != nil &&
    !l.cfg.NoNewPrivileges {
    seccomp.InitSeccomp(...)
}
finalizeNamespace(...)
...
if unix.Getppid() != l.parentPid { ... }
...
unix.Write(fd, []byte("0"))
system.Exec(name, l.cfg.Args[0:], os.Env())
```

syscalls for later setups need to be allowed  
e.g. execve, prctl, capset

**cBPF filters cannot block setup syscalls during application run.**

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content. The pipes are a vibrant green color and are set against a light gray shadow.

# Syscall count limiting

- Solution: Implement a state-based policy in eBPF



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Syscall count limiting

- Solution: Implement a state-based policy in eBPF

```
unsigned int no_exec;
SEC("seccomp")
int exec_once(struct seccomp_data *ctx)
{
    if (ctx->nr == __NR_execve) {
        if (!no_exec) {
            no_exec = 1;
            return SECCOMP_RET_ALLOW;
        }
        return SECCOMP_RET_ERRNO | EPERM;
    }
    return SECCOMP_RET_ALLOW;
}
```





# Syscall count limiting

- Solution: Implement a state-based policy in eBPF

```
unsigned int no_exec;
SEC("seccomp")
int exec_once(struct seccomp_data *ctx)
{
    if (ctx->nr == __NR_execve) {
        if (!no_exec) {
            no_exec = 1;
            return SECCOMP_RET_ALLOW;
        }
        return SECCOMP_RET_ERRNO | EPERM;
    }
    return SECCOMP_RET_ALLOW;
}
```

eBPF map to store states, i.e., whether  
execve has been called



# Syscall count limiting

- Solution: Implement a state-based policy in eBPF

```
unsigned int no_exec;
SEC("seccomp")
int exec_once(struct seccomp_data *ctx)
{
    if (ctx->nr == __NR_execve) {
        if (!no_exec) {
            no_exec = 1;
            return SECCOMP_RET_ALLOW;
        }
        return SECCOMP_RET_ERRNO | EPERM;
    }
    return SECCOMP_RET_ALLOW;
}
```

If execve hasn't been called,  
allow and mark it as called



# Syscall count limiting

- Solution: Implement a state-based policy in eBPF

```
unsigned int no_exec;
SEC("seccomp")
int exec_once(struct seccomp_data *ctx)
{
    if (ctx->nr == __NR_execve) {
        if (!no_exec) {
            no_exec = 1;
            return SECCOMP_RET_ALLOW;
        }
        return SECCOMP_RET_ERRNO | EPERM;
    }
    return SECCOMP_RET_ALLOW;
}
```

Otherwise reject the syscall

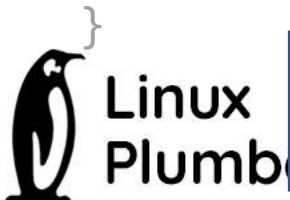


# Syscall count limiting

- Solution: Implement a state-based policy in eBPF

```
unsigned int no_exec;
SEC("seccomp")
int exec_once(struct seccomp_data *ctx)
{
    if (ctx->nr == __NR_execve) {
        if (!no_exec) {
            no_exec = 1;
            return SECCOMP_RET_ALLOW;
        }
        return SECCOMP_RET_ERRNO | EPERM;
    }
    return SECCOMP_RET_ALLOW;
}
```

Otherwise reject the syscall



Additional attack surface eliminated



# Temporal syscall specialization [1]

- Restrict different syscalls at different phases of an application



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Temporal syscall specialization [1]

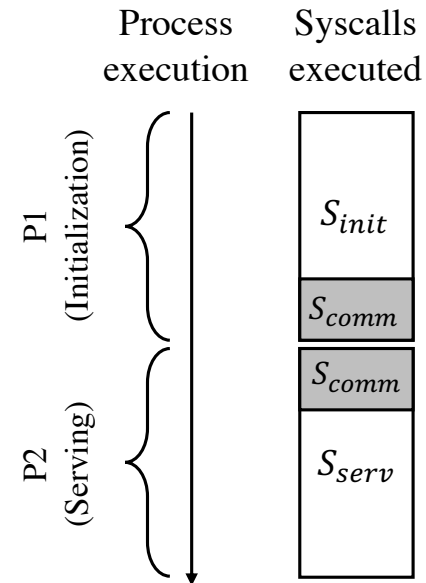
- Restrict different syscalls at different phases of an application
  - E.g., webservers



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Temporal syscall specialization [1]

- Restrict different syscalls at different phases of an application
  - E.g., webserver

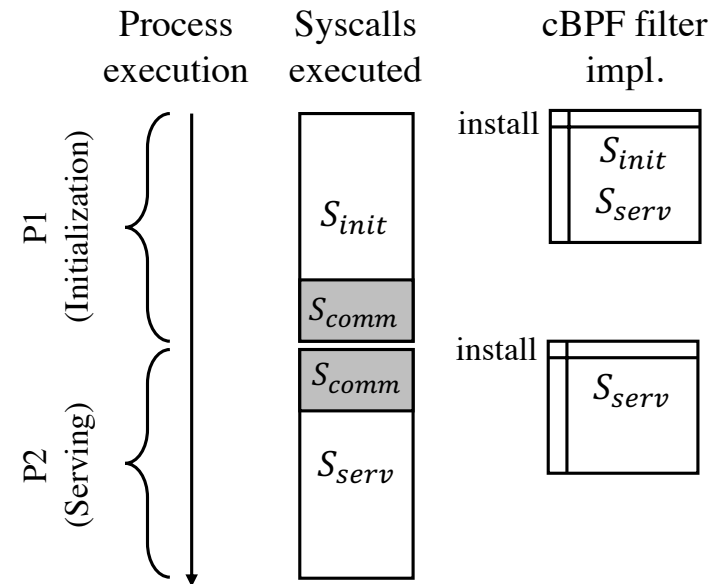


Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Temporal syscall specialization [1]

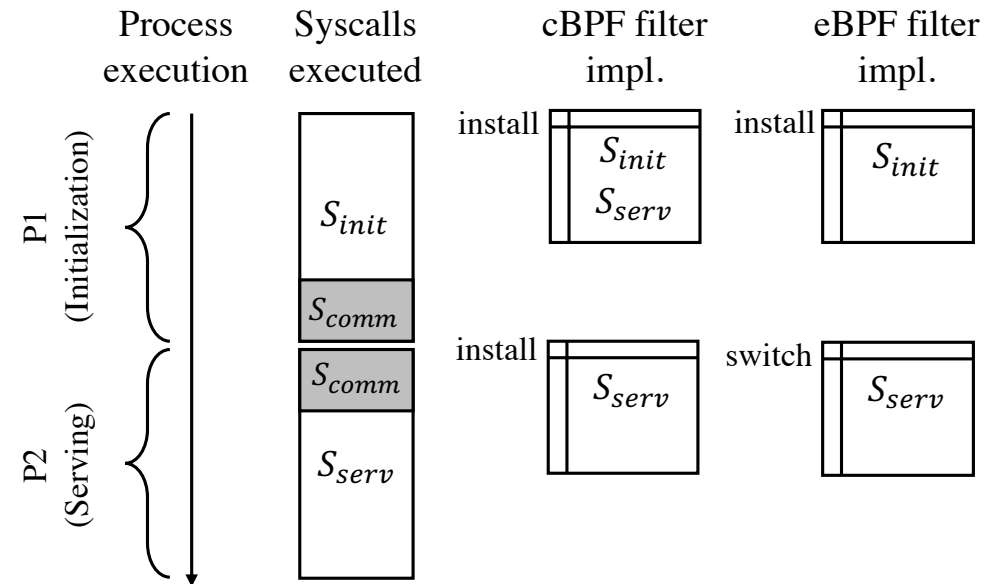
- Restrict different syscalls at different phases of an application
  - E.g., webserver
- cBPF cannot enforce the policy **precisely**



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Temporal syscall specialization [1]

- Restrict different syscalls at different phases of an application
  - E.g., webservers
- cBPF cannot enforce the policy **precisely**
- eBPF map can track application phases
  - Different policies for different phases



A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the top and left sides of the slide.

# Syscall-flow-integrity protection [1]



**Linux**  
**Plumbers Conference** | Dublin, Ireland **Sept. 12-14, 2022**

[1] SFIP: Coarse-Grained Syscall-Flow-Integrity Protection in Modern Systems, arXiv:2202.13716, 2020.

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Syscall-flow-integrity protection [1]

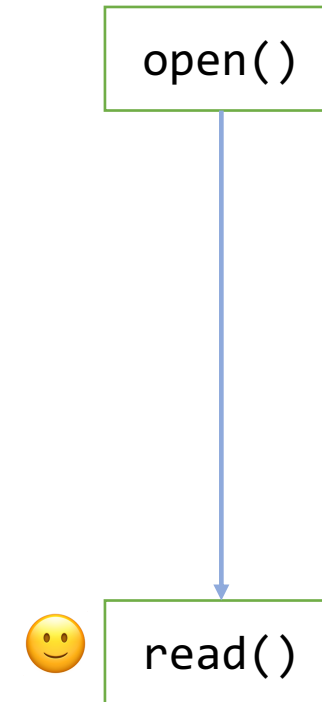
- Complements the concept of CFI with integrity for user-kernel transitions.



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

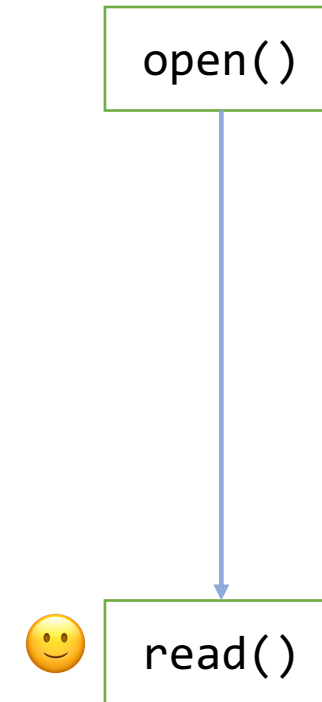
# Syscall-flow-integrity protection [1]

- Complements the concept of CFI with integrity for user-kernel transitions.
- Extract allowed syscall transitions to construct a syscall state machine



# Syscall-flow-integrity protection [1]

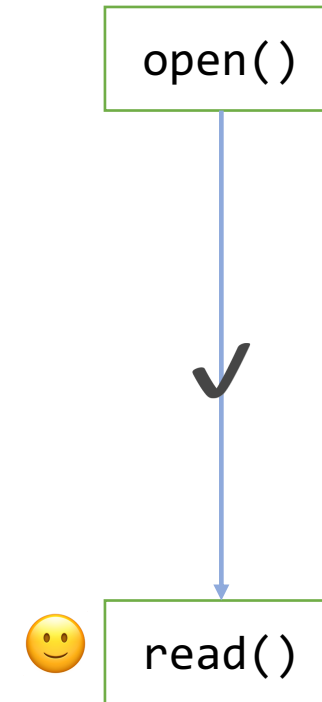
- Complements the concept of CFI with integrity for user-kernel transitions.
- Extract allowed syscall transitions to construct a syscall state machine
- Enforce such transitions at application runtime





# Syscall-flow-integrity protection [1]

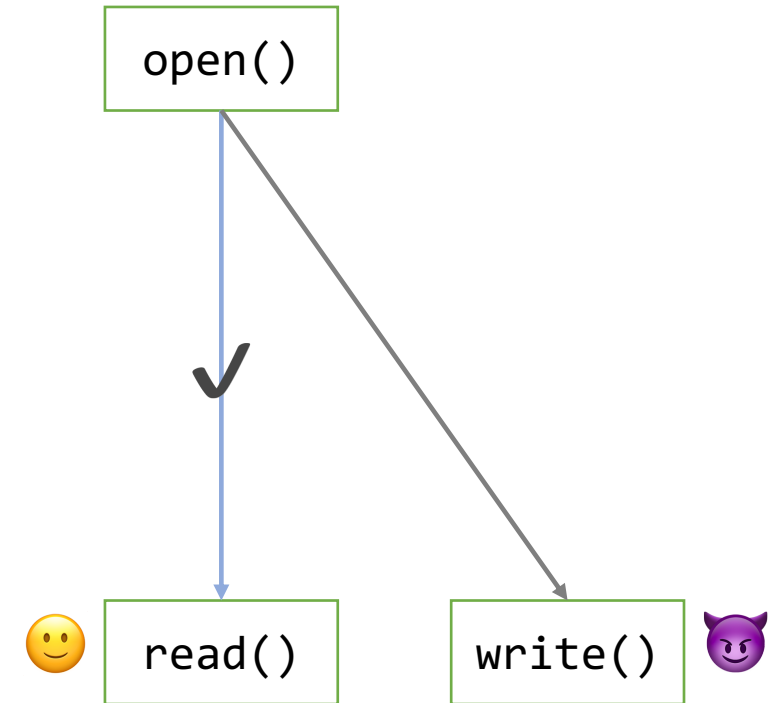
- Complements the concept of CFI with integrity for user-kernel transitions.
- Extract allowed syscall transitions to construct a syscall state machine
- Enforce such transitions at application runtime





# Syscall-flow-integrity protection [1]

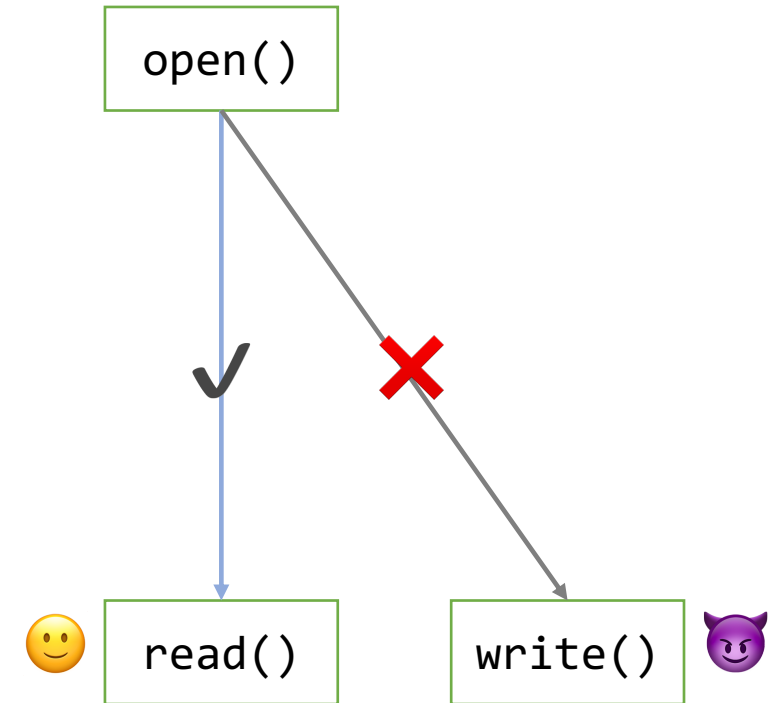
- Complements the concept of CFI with integrity for user-kernel transitions.
- Extract allowed syscall transitions to construct a syscall state machine
- Enforce such transitions at application runtime



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

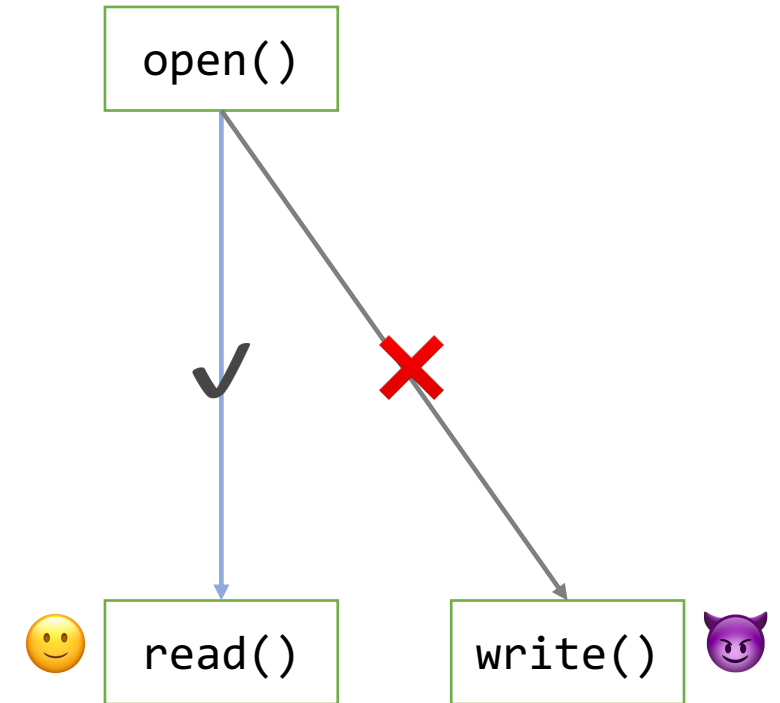
# Syscall-flow-integrity protection [1]

- Complements the concept of CFI with integrity for user-kernel transitions.
- Extract allowed syscall transitions to construct a syscall state machine
- Enforce such transitions at application runtime



# Syscall-flow-integrity protection [1]

- Complements the concept of CFI with integrity for user-kernel transitions.
- Extract allowed syscall transitions to construct a syscall state machine
- Enforce such transitions at application runtime
- Construct syscall state machine as an eBPF map and use Seccomp for enforcement



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the central text.

# How to make Seccomp-eBPF secure?

- Security model
- support for eBPF Seccomp filters under privileged eBPF
- IMA integration to enhance security



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Security model of Seccomp-eBPF

- Principle: Reduce to security of Seccomp and eBPF subsystems
  - Seccomp-eBPF is **as secure as** the existing Seccomp and eBPF



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the slide content.

# Security model of Seccomp-eBPF

- Principle: Reduce to security of Seccomp and eBPF subsystems
  - Seccomp-eBPF is **as secure as** the existing Seccomp and eBPF
- Seccomp requires privileges for filter installation
  - requires either `CAP_SYS_ADMIN` or `NO_NEW_PRIVS` attribute on the process



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



A decorative graphic of green pipes with valves and elbows, running vertically on the left side of the slide and horizontally at the top and bottom.

# Security model of Seccomp-eBPF

- Principle: Reduce to security of Seccomp and eBPF subsystems
  - Seccomp-eBPF is **as secure as** the existing Seccomp and eBPF
- Seccomp requires privileges for filter installation
  - requires either `CAP_SYS_ADMIN` or `NO_NEW_PRIVS` attribute on the process
- eBPF security is achieved by verification and privileges for program loading and helper access



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the slide content.

# Security model of Seccomp-eBPF

- Principle: Reduce to security of Seccomp and eBPF subsystems
  - Seccomp-eBPF is **as secure as** the existing Seccomp and eBPF
- Seccomp requires privileges for filter installation
  - requires either `CAP_SYS_ADMIN` or `NO_NEW_PRIVS` attribute on the process
- eBPF security is achieved by verification and privileges for program loading and helper access
- Seccomp-eBPF maintains the above security enforcement.



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide. The pipes are bright green and have a slight shadow effect.

# Reduce to eBPF security



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content. The pipes are bright green with a slight shadow effect.

# Reduce to eBPF security

- Program loading
  - concerns context access and attach point
  - covered by Seccomp security



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of green pipes with various fittings, valves, and elbows, running vertically on the left side of the slide and horizontally at the bottom. The pipes are rendered in a 3D style with shadows.

# Reduce to eBPF security

- Program loading
  - concerns context access and attach point
  - covered by Seccomp security
- Same privilege requirement for helpers
  - Basic helpers are unprivileged
  - Tracing helpers still require `CAP_BPF` and `CAP_PERFMON`



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of green pipes with valves and fittings, running vertically on the left side of the slide and curving at the top and bottom.

# Reduce to eBPF security

- Program loading
  - concerns context access and attach point
  - covered by Seccomp security
- Same privilege requirement for helpers
  - Basic helpers are unprivileged
  - Tracing helpers still require `CAP_BPF` and `CAP_PERFMON`
- eBPF does not have privilege requirement for maps
  - Add a new verifier hook to restrict map usage if desired



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide.

# Unprivileged eBPF filters?

- Vulnerabilities in eBPF are uncovered in recent years
  - prompted community to discourage unprivileged eBPF



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Unprivileged eBPF filters?

- Vulnerabilities in eBPF are uncovered in recent years
  - prompted community to discourage unprivileged eBPF
- **Seccomp-eBPF is useful under privileged eBPF**



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide.

# Unprivileged eBPF filters?

- Vulnerabilities in eBPF are uncovered in recent years
  - prompted community to discourage unprivileged eBPF
- **Seccomp-eBPF is useful under privileged eBPF**
  - lots of virtualization applications run as root
  - privileged enough to utilize eBPF Seccomp filters
    - example: Docker



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, running along the top, left, and bottom edges of the slide.

# Harden Seccomp-eBPF filters

- Optionally disables unprivileged eBPF filters based on the `BPF_UNPRIV_DEFAULT_OFF` configuration



Linux

Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**

A decorative graphic of a green pipe network with various fittings, valves, and elbows, running along the top, left, and bottom edges of the slide.

# Harden Seccomp-eBPF filters

- Optionally disables unprivileged eBPF filters based on the `BPF_UNPRIV_DEFAULT_OFF` configuration
- Same privilege requirement as eBPF tracing program types
  - `CAP_BPF + CAP_PERFMON` when unprivileged filters are not allowed
  - eBPF filters are similar to tracing programs in many aspects
    - e.g., common set of helpers, changing function return value / process behavior



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the slide content.

# Potential security support from IMA

- If we can verify an eBPF Seccomp filter comes from a trusted source, it is always secure to load it into the kernel and use it
  - turns the security problem into a trust problem



# Potential security support from IMA

- If we can verify an eBPF Seccomp filter comes from a trusted source, it is always secure to load it into the kernel and use it
  - turns the security problem into a trust problem
- The IMA subsystem checks signatures of files and verifies whether it comes from a trusted source



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Potential security support from IMA

- If we can verify an eBPF Seccomp filter comes from a trusted source, it is always secure to load it into the kernel and use it
  - turns the security problem into a trust problem
- The IMA subsystem checks signatures of files and verifies whether it comes from a trusted source
- Idea: implementing IMA signature support for eBPF filters
  - The IMA community is welcome to the eBPF signature support
  - Refer to KP's talk :)



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the top and left sides of the slide.

# Our eBPF Seccomp filter support



Linux  
Plumbers Conference | Dublin, Ireland **Sept. 12-14, 2022**



A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Our eBPF Seccomp filter support

- Based on Dhillon's patch in 2018



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the slide content.

# Our eBPF Seccomp filter support

- Based on Dhillon's patch in 2018
- A new program type: `BPF_PROG_TYPE_SECCOMP`



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, elbows, and valves, framing the central text area.

# Our eBPF Seccomp filter support

- Based on Dhillon's patch in 2018
- A new program type: `BPF_PROG_TYPE_SECCOMP`
- Same attach point as cBPF filters



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022



# Our eBPF Seccomp filter support

- Based on Dhillon's patch in 2018
- A new program type: `BPF_PROG_TYPE_SECCOMP`
- Same attach point as cBPF filters
- Verifier hooks ensure correct access to the `seccomp_data` context and restrict the set of helpers for eBPF filters



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

A decorative graphic of a green pipe network with various fittings, valves, and elbows, framing the slide content.

# Our eBPF Seccomp filter support

- Based on Dhillon's patch in 2018
- A new program type: `BPF_PROG_TYPE_SECCOMP`
- Same attach point as cBPF filters
- Verifier hooks ensure correct access to the `seccomp_data` context and restrict the set of helpers for eBPF filters
- Security model: Reduce to Seccomp and eBPF



Linux

Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022

# Container runtime integration

- Seccomp-eBPF can be easily integrated with container runtimes
- It is supported in crun,
  - A fast OCI-compliant runtime and the default for Podman
  - Credits to Giuseppe Scrivano (RedHat)



# Container runtime integration

- Seccomp-eBPF can be easily integrated with container runtimes
- It is supported in crun,
  - A fast OCI-compliant runtime and the default for Podman
  - Credits to Giuseppe Scrivano (RedHat)

```
$ podman --runtime /usr/local/bin/crun run  
--annotation run.oci.seccomp_ebpf_file  
=ebpf_filter.o
```





# Conclusion

- We revisit Seccomp-eBPF filters
  - Significantly improve programmability of system call security
  - Many real-world use cases
  - Security model can be well defined
- <https://github.com/xlab-uiuc/seccomp-ebpf-upstream>



Linux  
Plumbers Conference | Dublin, Ireland Sept. 12-14, 2022