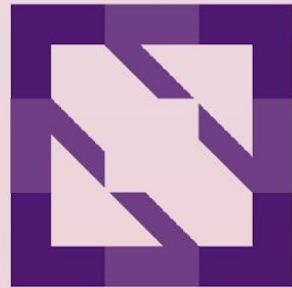




KubeCon

— North America 2023 —



CloudNativeCon





KubeCon



CloudNativeCon

North America 2023

Survive eBPF with bpftrace

Andrew Stoycos & Shane Utt



KubeCon



CloudNativeCon

North America 2023



Andrew Stoycos
Senior SWE @ Red Hat
Bpf Maintainer
Network Policy API Maintainer



Shane Utt
Staff SWE @ Kong
SIG Network Chair
Gateway API Maintainer



KubeCon



CloudNativeCon

North America 2023

What is eBPF, and how does it relate to Kubernetes?

Quick introduction to eBPF



KubeCon



CloudNativeCon

North America 2023

Use
Cases

Networking

Security

Observability &
Tracing



Projects



Katran



PIXIE

User
Space



SDKs



=GO



Application

- Tracing
- Profiling
- Monitoring
- ...

Kernel



Kernel Runtime

Verifier & JIT

Maps

Kernel Helper API

OS
Runtime



- Observability
- Security Controls
- Networking
- Network Security
- Load Balancing
- Behavioral Security
- ...

Why eBPF in Kubernetes?

- eBPF is a powerful general-purpose framework to extend the Linux Kernel
- Use cases from observability, networking, security and more
- Many existing examples of eBPF in kubernetes
 - [Cilium](#) and [Calico](#) CNIs
 - [Pixie](#): Open source observability
 - [KubeArmor](#): Container-aware Runtime Security Enforcement System
 - [Blixt](#): Gateway Api L4 conformance implementation
 - [NetObserv](#): Open Source Operator for network observability



KubeCon



CloudNativeCon

North America 2023

Issues with eBPF Programs in Kubernetes

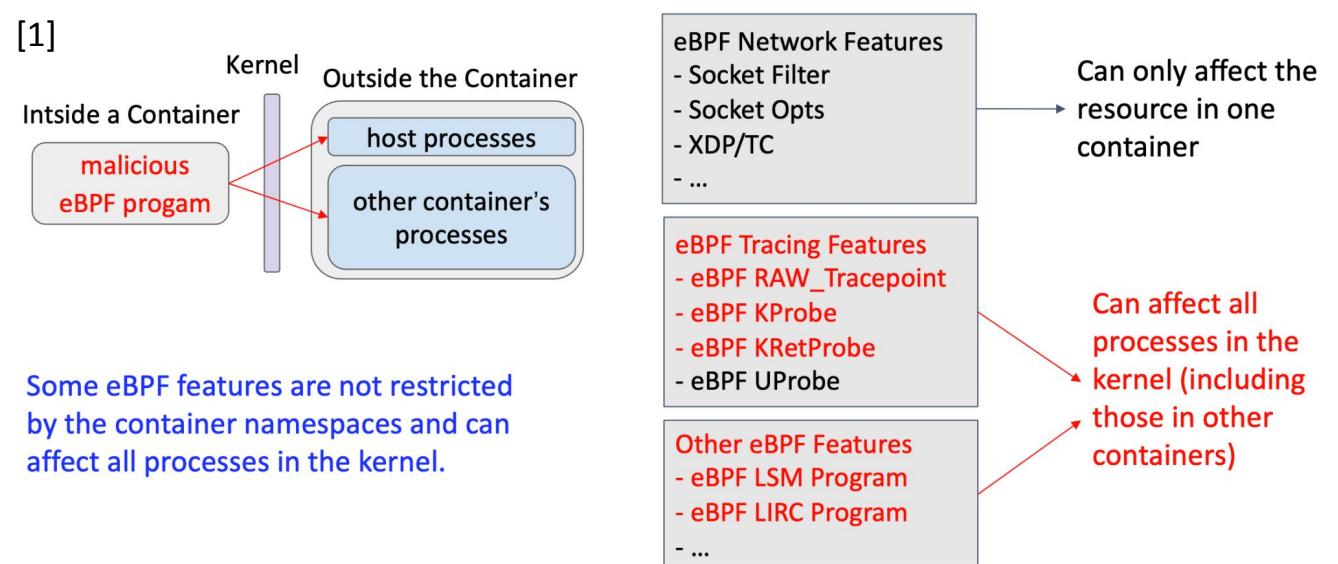
Security Issues with eBPF

- eBPF is *not* namespaced it can easily escape container isolation
- Relies on **highly privileged containers**
- Vulnerable to **supply chain attacks**



eBPF is not namespaced

- eBPF has to be loaded from **rootful containers** – rootless containers will not work since bpf is not namespaced.
- Even with the basic set of permissions, root privilege escalation is possible.
 - See: [CVE-2022-23222](#) - now patched.
- Since eBPF in containers are able to modify the host kernel, they are not really containing anything.



Highly Privileged Containers

At minimum **CAP_BPF** is required, which should effectively be considered as root.

However, in practice:

- Many BPF programs and attach points require additional capabilities:
 - **CAP_SYS_PTRACE**, **CAP_NET_ADMIN** and even **CAP_SYS_ADMIN**
 - In the “compile on-the-fly” use-case, these privileges are maintained for the lifetime of the program loading probes.
 - In all cases, these privileges include things that aren’t strictly necessary for eBPF and are too coarsely grained to be useful.

Since pods that load eBPF are mostly long-lived and often don’t drop privileges it leaves a wide attack surface.

An Example: Scary Things we can do with eBPF

1. Snoop on all SSL traffic: github.com/alessandrod/snuffy
2. Redirect incoming TCP 22 (SSH) traffic the host into my container as part of a phishing attack
3. Deploy a keylogger: github.com/willfindlay/bpf-keylogger
4. Naively attach probes to application functions or tracepoints that are called at a high frequency, grinding the system to a halt.
5. Exploit the K8s [Trampoline Pod Threat](#) (DEMO!!)



Functional Problems with eBPF

- Silos: interoperability problems, can cause **cluster instability**
- Visibility into the ebpf subsystem is often limited at the cluster scope
- No fine grained versioning between userspace + kernelspace
- eBPF loading + management stack is duplicated by every application





KubeCon



CloudNativeCon

North America 2023

Enter bpfcd

Introducing bpfdev

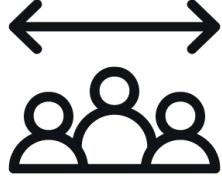
- Open source project started in the Red Hat Emerging Tech Networking Group
- An eBPF program manager
 - **Manages lifecycle** (Loading/Unloading/Pinning) of eBPF Programs, removing CAP_BPF from applications
 - Program Cooperation
 - Central daemon for **managing loading policy** for security and visibility
- bpfdev is developed in Rust, built on top of the **Rust eBPF library Aya**.
- Includes a **K8s operator and APIs** (developed in Go/controller-runtime)

bpfdev: core-features



Security

- CAP_BPF isolation
- Kubernetes RBAC
- [Cosign](#) Integration



Productivity

- Application deduplication
- eBPF distribution and version control using our [OCI bytecode image spec](#)



Observability

- eBPF subsystem monitoring
- In Kubernetes node based eBPF subsystem state reflected via K8s api

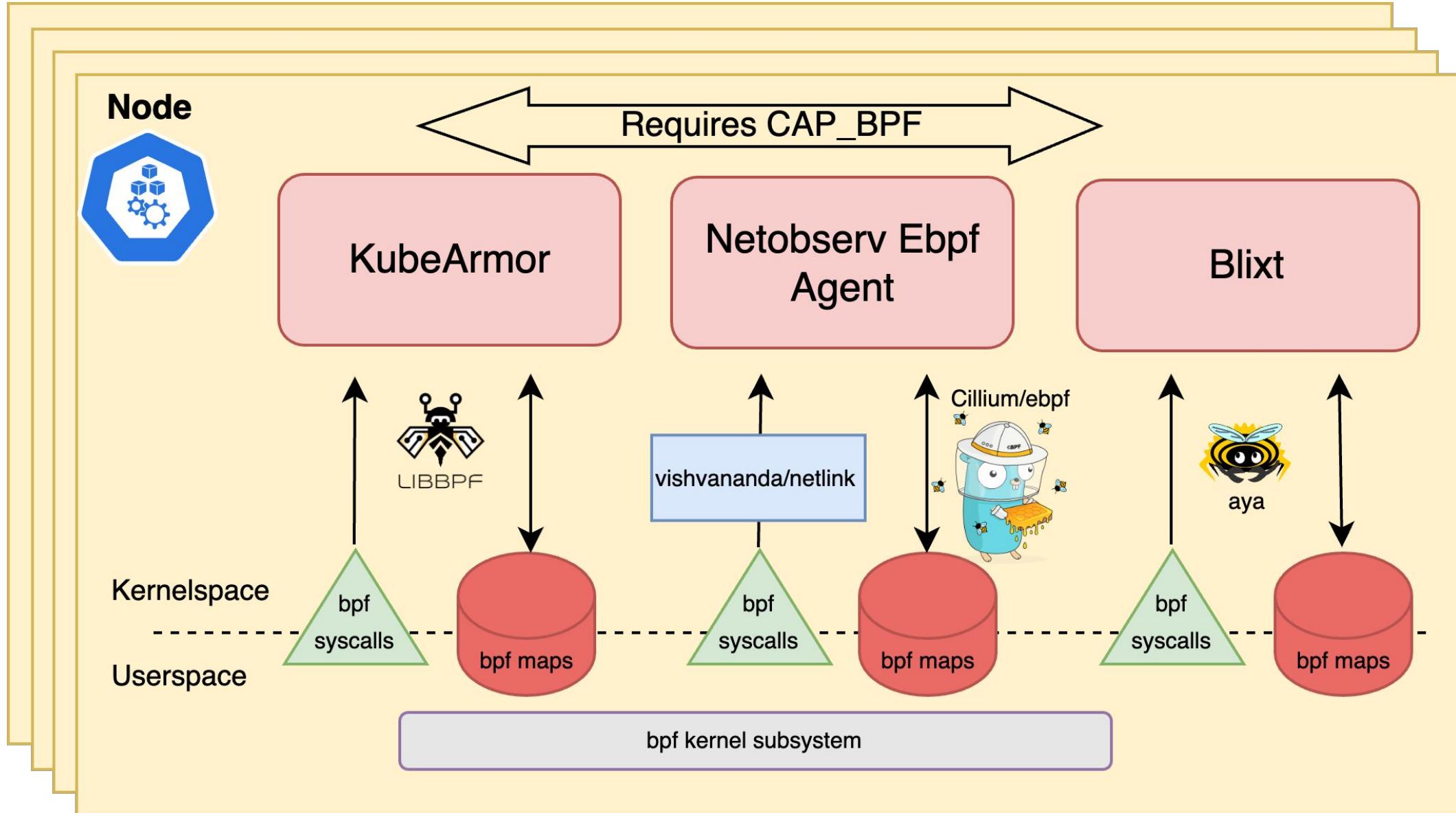


Program Support

- Support for:
 - XDP
 - TC
 - Tracepoint
 - Uprobe
 - Kprobe
- leverages the libxdp protocol for XDP/TC



bpf on Kubernetes before bpfdev



bpf on Kubernetes after bpfd

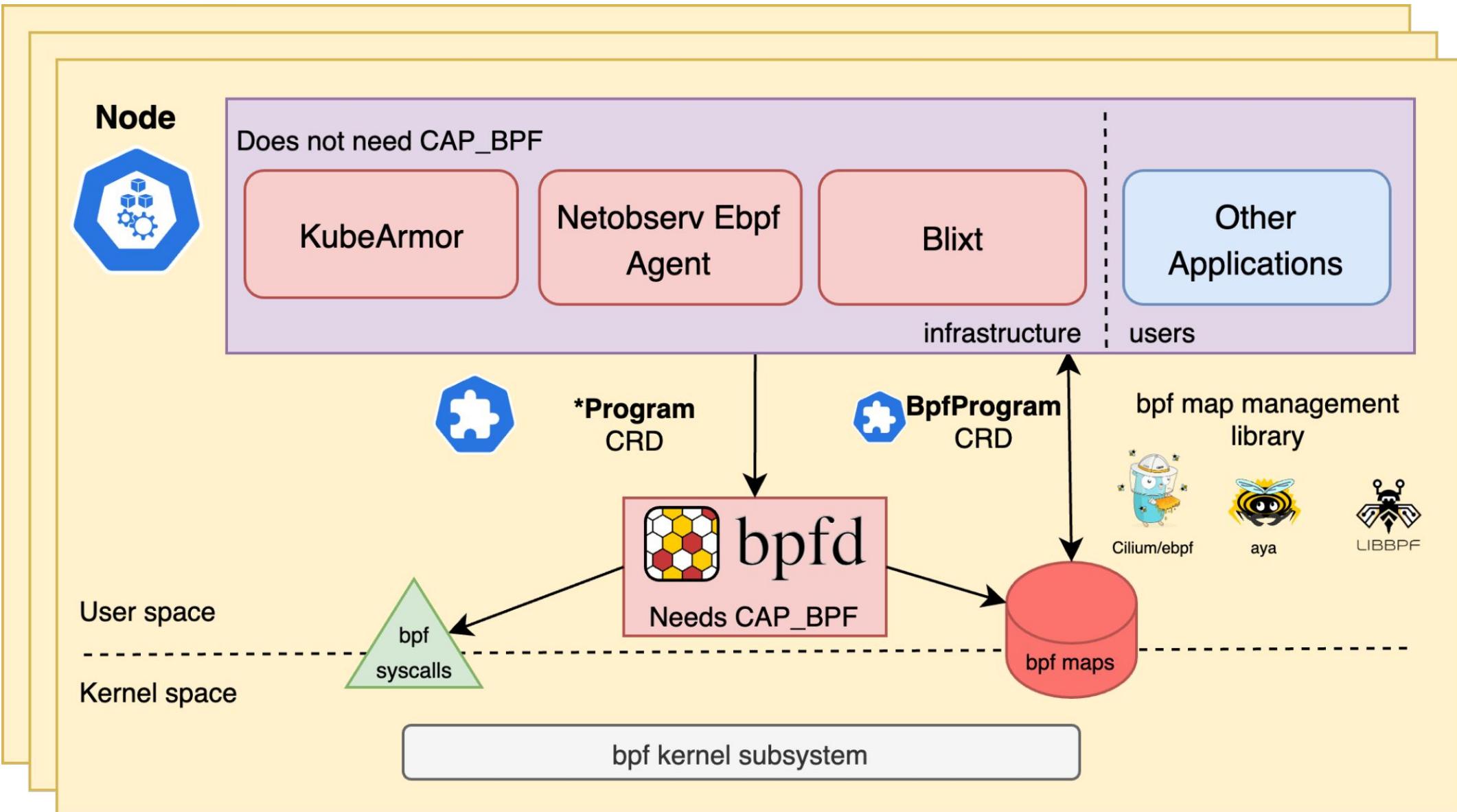


KubeCon

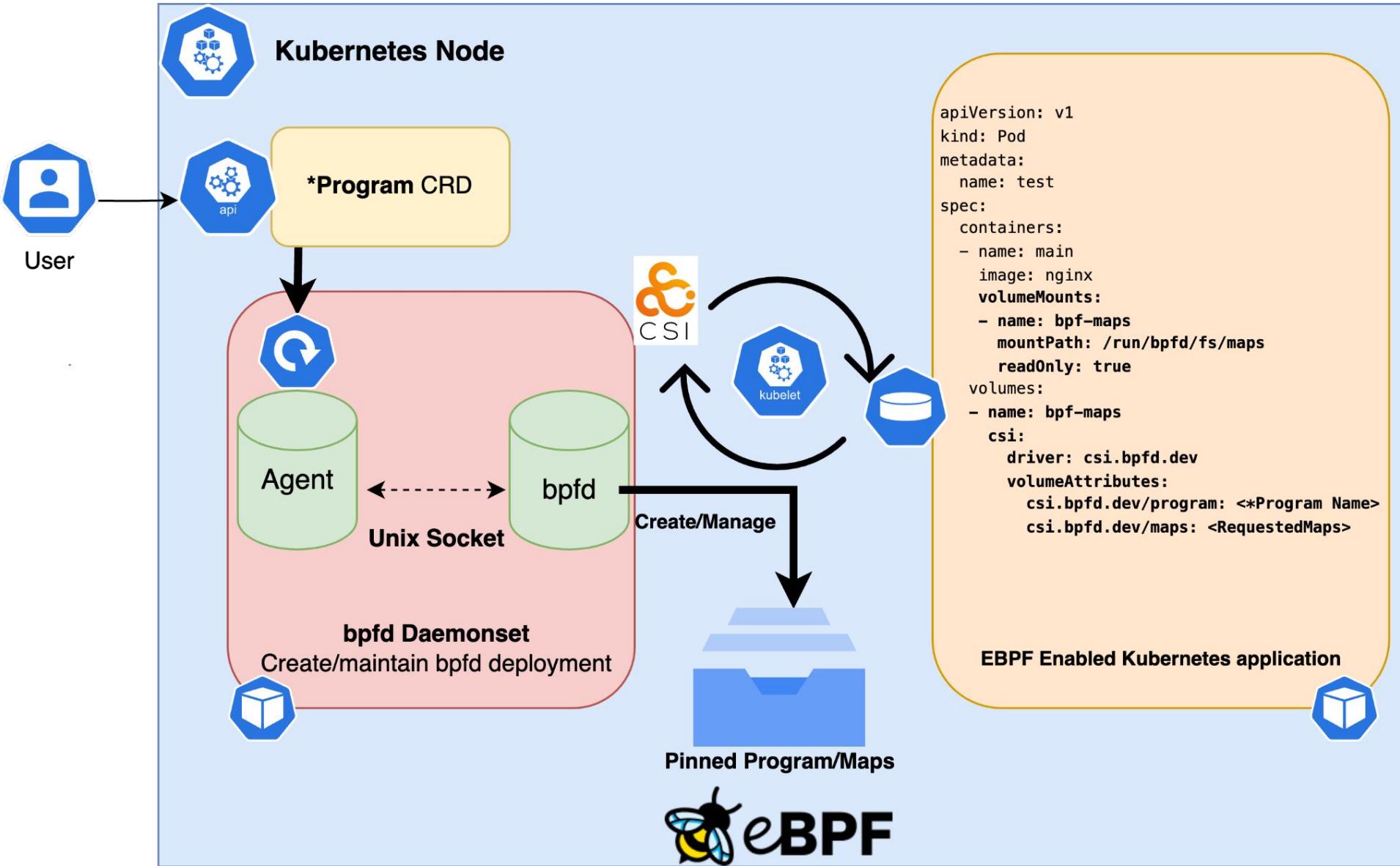


CloudNativeCon

North America 2023



bpfdev + Kubernetes Architecture





KubeCon



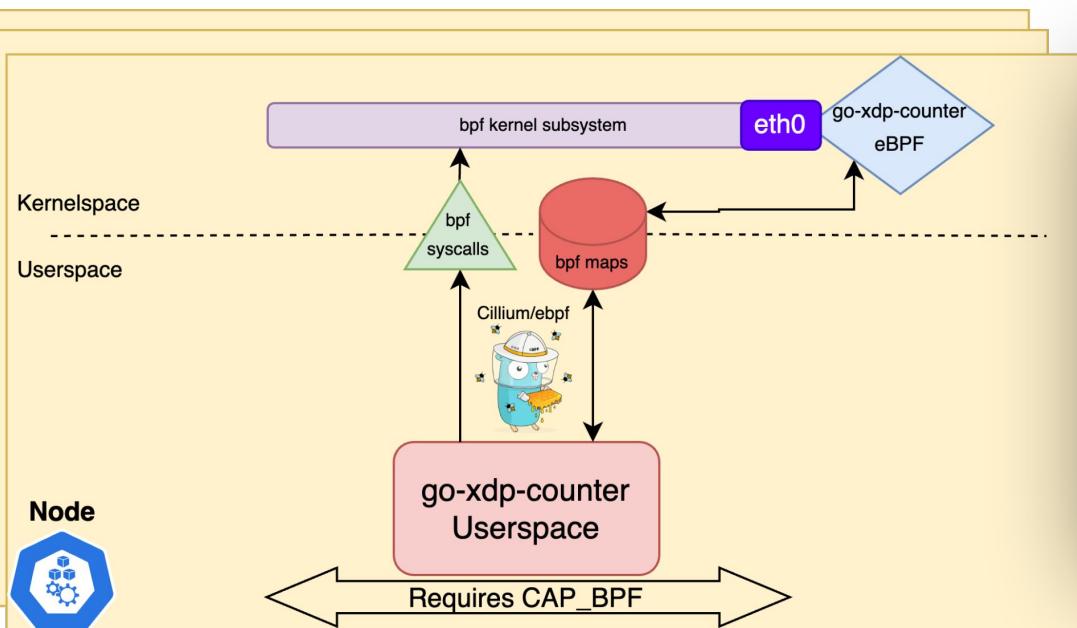
CloudNativeCon

North America 2023

DEMO

Peaceful eBPF: A Xdp Counter

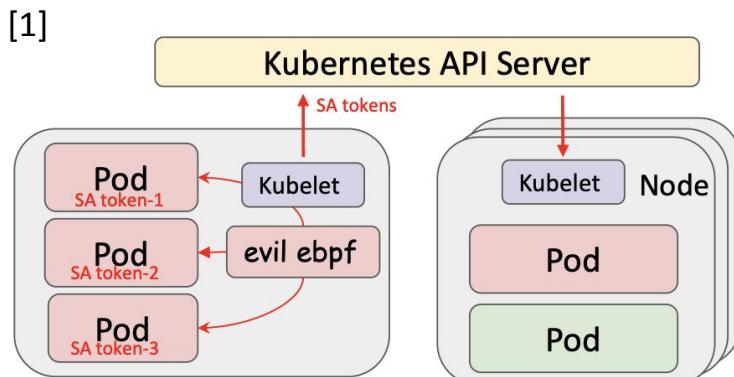
- XDP Program attached to the main network interface on a node that simply counts the number of bytes going by.
- Userspace + eBPF programs are compiled together into one binary, requiring a privileged daemonset for its deployment



```
1 [astoycos@nfvsdn-02-oot examples]$ kubectl logs go-xdp-counter-ds-29klv -n go-xdp-counter
2 2023/10/31 03:04:06 Attached XDP program to iface "eth0" (index 2119)
3 2023/10/31 03:04:06 Press Ctrl-C to exit and remove the program
4 2023/10/31 03:04:09 25 packets received
5 2023/10/31 03:04:09 3809 bytes received
6
7 2023/10/31 03:04:12 49 packets received
8 2023/10/31 03:04:12 7552 bytes received
9
10 2023/10/31 03:04:15 49 packets received
11 2023/10/31 03:04:15 7552 bytes received
```

An EVIL Xdp Counter Application 😈

- For this demo we've implemented a service-account token stealer
- Using eBPF tracepoint programs its relatively simple to:
 - Break out of the Pod's container boundaries
 - Listen for all /var/run/secrets/kubernetes.io/serviceaccount/token reads
 - Steal the token from memory and write within our evil pod



On a vulnerable VM (node), all Pods' service accounts (SA) can be abused by eBPF attackers.

```
rules:
- apiGroups: ["stable.example.com"]
  resources: ["crontabs"]
  verbs: ["get", "list", "watch", "create", "update", "patch", "delete"]
```

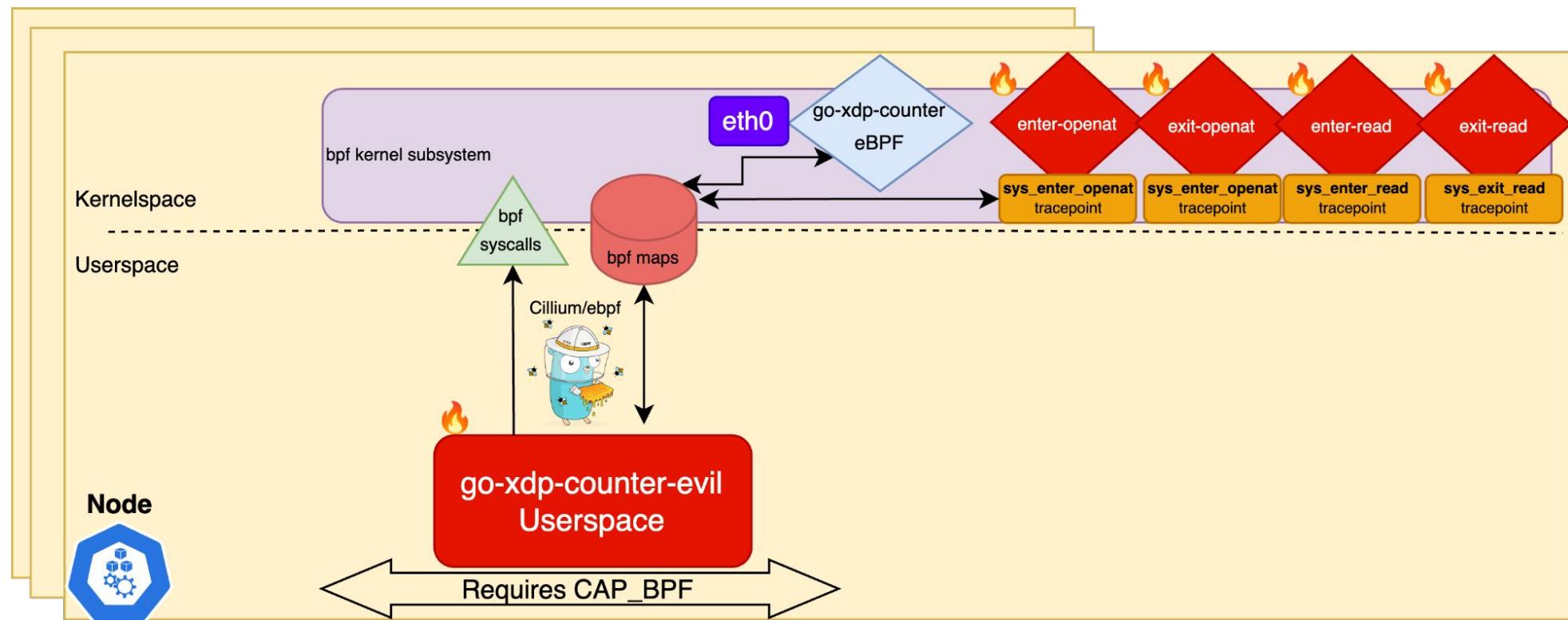
Some Pods have powerful permissions to affect Pods on other nodes.

```
# steal other Pods' service account tokens
$ export TOKEN=$(evil-ebpf-read
/var/run/secrets/kubernetes.io/serviceaccount/token)
```

```
# manipulate other nodes
$ curl -k --header "Authorization: Bearer $TOKEN"
https://172.16.22.202:10250/...
```

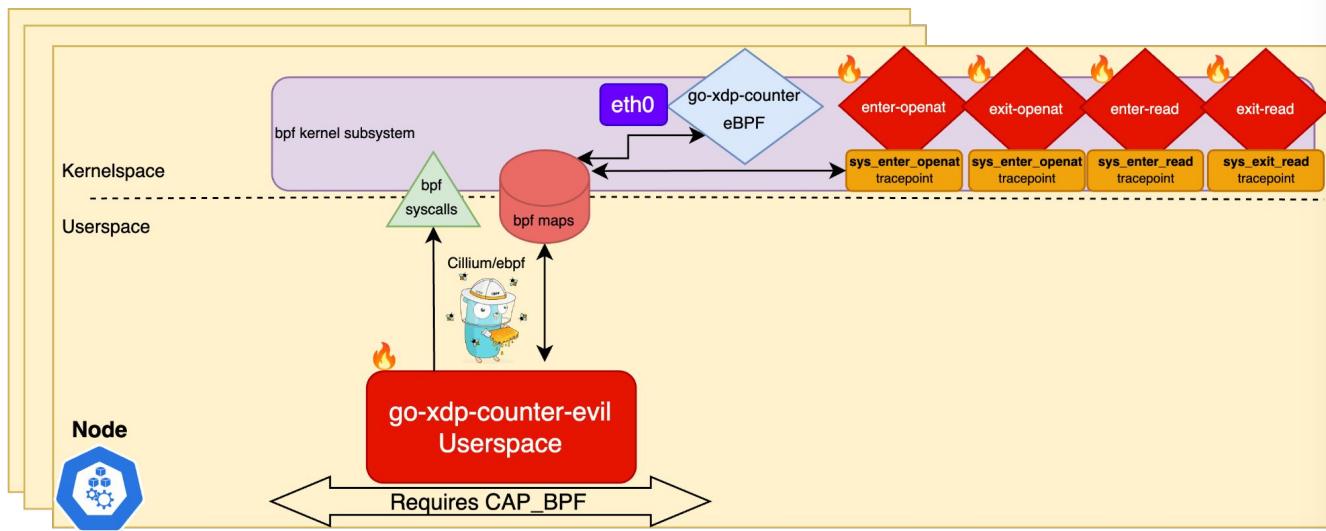
Peaceful eBPF → Evil eBPF

- If the binary containing the userspace + eBPF programs is compromised (supply chain attack, privilege abuse, user spoofing, etc.) exploitation is relatively simple.
- Application container is privileged, and can easily load + attach multiple other 😈 bpf programs.



Evil eBPF Exploit

- Uh oh the evil program is now dumping service account tokens... Let's use them



```
● ● ●  
1 2023/10/30 13:49:24  
2 pid: 1284470  
3  
4 comm: coredns  
5  
6 token: ey{.....}oA  
7  
8 parsed token info: {  
9   "aud": [  
10     "https://kubernetes.default.svc.cluster.local"  
11   ],  
12   "exp": 1730208446,  
13   "iat": 1698672446,  
14   "iss": "https://kubernetes.default.svc.cluster.local",  
15   "kubernetes.io": {  
16     "namespace": "kube-system",  
17     "pod": {  
18       "name": "coredns-5d78c9869d-8fd5b",  
19       "uid": "ee7ee783-8d98-4ec9-8a2a-839863da835b"  
20     },  
21     "serviceaccount": {  
22       "name": "coredns",  
23       "uid": "2f730bae-fbe5-4d78-8090-721b080da198"  
24     },  
25     "warnafter": 1698676053  
26   },  
27   "nbf": 1698672446,  
28   "sub": "system:serviceaccount:kube-system:coredns"  
29 }
```

Invisible Evil eBPF



KubeCon



CloudNativeCon

– North America 2023

- Due to the lack of standard observability into the eBPF subsystem in Kubernetes cluster, this evil eBPF could hide pretty easily FOREVER since everything still appears to be running as expected.

A binary code background consisting of a grid of black, white, and grey squares. Overlaid on this grid is the word "rootkit" in a large, bold, sans-serif font. The letters are colored in a gradient of green, blue, and red. The background has a subtle, slightly blurred effect.

Peaceful eBPF -> Evil eBPF



KubeCon



CloudNativeCon

North America 2023

```
        "host": "https://github.com/takao-yoshida/ci-test",
        "port": 34412,
        "path": "/api/v1/test",
        "method": "POST"
    }
}

const test = {
    "name": "CI Test"
}

const host = "https://github.com/takao-yoshida/ci-test"
const port = 34412
const path = "/api/v1/test"
const method = "POST"

const headers = {
    "Content-Type": "application/json"
}

const body = JSON.stringify(test)

const options = {
    host,
    port,
    path,
    method,
    headers,
    body
}

const https = require("https")

const httpsRequest = https.request(options, (res) => {
    let data = ""

    res.on("data", (chunk) => {
        data += chunk
    })

    res.on("end", () => {
        console.log(data)
    })
})

httpsRequest.end()
```

Enter bpfd: Discovery

- Step 1: Start by [installing bpfd](#) on the compromised cluster.
- Step 2: Easily inspect all of the bpf programs running on a given node

```
1 [astoycos@nfvsdn-02-oot bpfd]$ kubectl get bpfprogram -l kubernetes.io/hostname=kubecon-na-2023-bpfd-demo-control-plane
2 NAME                                     AGE
3 ...
4 dump-bpf-map-316-kubecon-na-2023-bpfd-demo-control-plane  61m
5 dump-bpf-prog-317-kubecon-na-2023-bpfd-demo-control-plane  61m
6 enter-openat-41175-kubecon-na-2023-bpfd-demo-control-plane 61m
7 enter-read-41176-kubecon-na-2023-bpfd-demo-control-plane  61m
8 exit-openat-41177-kubecon-na-2023-bpfd-demo-control-plane  61m
9 exit-read-41178-kubecon-na-2023-bpfd-demo-control-plane   61m
10 restrict-fil...-45-kubecon-na-2023-bpfd-demo-control-plane 61m
11 sd-devices-132-kubecon-na-2023-bpfd-demo-control-plane   61m
12 sd-devices-133-kubecon-na-2023-bpfd-demo-control-plane   61m
13 sd-devices-134-kubecon-na-2023-bpfd-demo-control-plane   61m
14 sd-devices-135-kubecon-na-2023-bpfd-demo-control-plane   61m
15 sd-devices-138-kubecon-na-2023-bpfd-demo-control-plane   61m
16 sd-devices-141-kubecon-na-2023-bpfd-demo-control-plane   61m
17 sd-fw-egress-136-kubecon-na-2023-bpfd-demo-control-plane 61m
18 sd-fw-egress-139-kubecon-na-2023-bpfd-demo-control-plane 61m
19 sd-fw-egress-142-kubecon-na-2023-bpfd-demo-control-plane 61m
20 sd-fw-egress-144-kubecon-na-2023-bpfd-demo-control-plane 61m
21 sd-fw-ingress-137-kubecon-na-2023-bpfd-demo-control-plane 61m
22 sd-fw-ingress-140-kubecon-na-2023-bpfd-demo-control-plane 61m
23 sd-fw-ingress-143-kubecon-na-2023-bpfd-demo-control-plane 61m
24 sd-fw-ingress-145-kubecon-na-2023-bpfd-demo-control-plane 61m
25 xdp-stats-41179-kubecon-na-2023-bpfd-demo-control-plane  61m
```



Hrm what are these?

Enter bpf: Discovery



KubeCon



CloudNativeCon

North America 2023

Enter bpfd: Discovery



```
1 apiVersion: bpfd.dev/v1alpha1
2 kind: BpfProgram
3 metadata:
4   annotations:
5     BTF-ID: "16887"
6     GPL-Compatible: "true"
7     JITed: "true"
8     Kernel-Allocated-Memory-Bytes: "4096"
9     Kernel-ID: "41175"
10    Loaded-At: 2023-10-31T18:55:16+0000
11    Map-IDs: '[22699 22700]'
12    Name: enter_openat
13    Size-JITed-Bytes: "198"
14    Size-Translated-Bytes: "352"
15    Tag: f034f9ddad3e2dd9
16    Type: tracepoint
17    Verified-Instruction-Count: "861"
18    creationTimestamp: "2023-10-31T19:36:04Z"
19    generation: 1
20    labels:
21      bpfd.dev/discoveredProgram: ""
22      kubernetes.io/hostname: kubecon-na-2023-bpfd-demo-control-plane
23    name: enter-openat-41175-kubecon-na-2023-bpfd-demo-control-plane
24    resourceVersion: "3839"
25    uid: b266facc-a0ef-42a1-9b2c-c84ede061ef9
26  spec:
27    type: tracepoint
```



```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   creationTimestamp: "2023-10-31T18:54:54Z"
5   generateName: go-xdp-counter-ds-
6   labels:
7     controller-revision-hash: c67d7f6b5
8     name: go-xdp-counter
9     pod-template-generation: "1"
10    name: go-xdp-counter-ds-rq2gq
11    namespace: go-xdp-counter
12 ...
13   status:
14     containerStatuses:
15     -
16       name: go-xdp-counter
17       ready: true
18       restartCount: 0
19       started: true
20       state:
21         running:
22           startedAt: "2023-10-31T18:55:16Z"
23     hostIP: 172.19.0.2
24     phase: Running
25     podIP: 172.19.0.2
26     podIPs:
27     -
28       ip: 172.19.0.2
29     qosClass: BestEffort
30     startTime: "2023-10-31T18:54:55Z"
31 ...
```

Let's take a closer look..
What Pod loaded this?

bpf: Mitigation

- Start by Deleting the evil application
- Re-deploy with bpf:ds



```
1 apiVersion: bpf:dev/v1alpha1
2 kind: XdpProgram
3 metadata:
4   labels:
5     app.kubernetes.io/name: xdpprogram
6   name: go-xdp-counter-example
7 spec:
8   bpfFunctionName: xdp_stats
9   # Select all nodes
10  nodeSelector: {}
11  interfaceSelector:
12    primaryNodeInterface: true
13  priority: 55
14  bytecode:
15    image:
16      url: quay.io/bpf:bytecode/go-xdp-counter-evil:latest
```

Bytecode is still evil?

```
1 apiVersion: apps/v1
2 kind: DaemonSet
3 metadata:
4   name: go-xdp-counter-ds
5   namespace: go-xdp-counter
6 ...
7 spec:
8 ...
9   spec:
10    nodeSelector: {}
11    hostNetwork: true
12    serviceAccountName: bpf:app-go-xdp-counter
13    securityContext:
14      privileged: false
15 ...
16    containers:
17      - name: go-xdp-counter
18        image: quay.io/bpf:userspace/go-xdp-counter:latest
19 ...
20    volumeMounts:
21      - name: bpf-maps
22        mountPath: /run/bpf:fs/maps
23        readOnly: true
24    volumes:
25      - name: bpf-maps
26        csi:
27          driver: csi.bpf:dev
28        volumeAttributes:
29          csi.bpf:dev/program: go-xdp-counter-example
30          csi.bpf:dev/maps: xdp_stats_map
```

bpf: Mitigation with K8s RBAC

- Even if the bytecode image is compromised the evil program will not be loaded with bpf!
- A **TracepointProgram** object would need to be created which is disallowed by the cluster admin via RBAC mechanisms.

```
apiVersion: bpf.dev/v1alpha1
kind: TracepointProgram
name: enter-openat
```

```
1 apiVersion: rbac.authorization.k8s.io/v1
2 kind: ClusterRole
3 metadata:
4   creationTimestamp: null
5   name: go-xdp-counter-namespace-owner-role
6 rules:
7 - apiGroups:
8   - bpf.dev
9   resources:
10  - xdpprograms
11  verbs:
12  - create
13  - get
14  - update
```

Compromised service account



bpf-app-go-xdp-counter

K8s RBAC Controls



```
1 apiVersion: bpf.dev/v1alpha1
2 kind: XdpProgram
3 metadata:
4   labels:
5     app.kubernetes.io/name: xdpprogram
6   name: go-xdp-counter-example
7 spec:
8   bpffunctionname: xdp_stats
9   # Select all nodes
10  nodeselector: {}
11  interfaceselector:
12    primarynodeinterface: true
13  priority: 55
14  bytecode:
15    image:
16      url: quay.io/bpf-dev-bytecode/go-xdp-counter-evil:latest
```



Create

bpf: Mitigation with bytecode signing

- Now we're loading with bpf so it's much harder to load unintended programs since all bpf behavior is declarative.
- However what if the bytecode image was compromised?
 - quay.io/bpf:bytecode/go-xdp-counter-evil:latest (unsigned)**
 - quay.io/bpf:bytecode/go-xdp-counter:latest (signed)**



```
1 [astoycos@nfvdsn-02-oct examples]$ kubectl logs bpf-daemon-bg99x -n bpf
2 Defaulted container "bpf" out of: bpf, bpf-agent, node-driver-registrar
3 [INFO bpf] Log using env_logger
4 [INFO bpf] Has CAP_BPF: true
5 [INFO bpf] Has CAP_SYS_ADMIN: true
6 [INFO bpf::certs] CA Certificate file /etc/bpf/certs/ca/ca.pem does not exist. Creating CA Certificate.
7 [INFO bpf::certs] bpf Certificate Key /etc/bpf/certs/bpf/bpf.key does not exist. Creating bpf Certificate.
8 [INFO bpf::certs] bpf-client Certificate Key /etc/bpf/certs/bpf-client/bpf-client.key does not exist. Creating bpf-client Certificate.
9 [INFO bpf::oci_utils::cosign] Starting Cosign Verifier, downloading data from Sigstore TUF repository
10 [INFO bpf::serve] Listening on /bpf-sock/bpf.sock
11 [INFO bpf::storage] CSI Plugin Listening on /run/bpf/csi/csi.sock
12 [WARN bpf::oci_utils::cosign] The bytecode image: quay.io/bpf:bytecode/go-xdp-counter-evil:latest is unsigned
13 [INFO bpf::command] Loading program bytecode from container image: quay.io/bpf:bytecode/go-xdp-counter-evil:latest
14 [INFO bpf::oci_utils::cosign] The bytecode image: quay.io/bpf/xdp-dispatcher:v2 is signed
15 [INFO bpf::bpf] Added xdp program with name: xdp_stats and id: 40668
```

Bytecode is still evil, but [Cosign](#) integration warns the user.

bpf: Mitigation



KubeCon

CloudNativeCon

North America 2023



KubeCon



CloudNativeCon

North America 2023

Getting Involved

Getting Started

- You can **deploy to a standard Linux system** and use **bpfctl**
- You can **deploy the bpf operator on a kind cluster** (or similar)
- We provide several examples, or you can bring your existing program into the fold
 - **TC, XDP, Tracepoint**, e.t.c.
- Our website has getting started information, and several guides



<https://bpfd.dev>

Work in Kubernetes SIG Network with Gateway API

- Blixt: A Layer 4 load-balancer deployable with bpf^d
 - Part of the Gateway API project
 - TC for ingress/egress networking
 - Used for CI and testing scenarios
 - Maintained by us!
- Potentially a good jumping in point if you're interested in eBPF networking



<https://github.com/kubernetes-sigs/blixt>

Community



- Weekly community meetings on Thursdays <https://bpfd.dev/community>
- **#bpfd** channel on Kubernetes Slack (**#ebpf** for general ebpf discussions as well)
- We're also very active in the Aya and Rust communities on Discord



<https://aya-rs.dev/community>

Roadmap



<https://bpf.dev/community>

- See the [Bpf Dev Github Project](https://bpf.dev/community)
- Plans to apply CNCF Sandbox project
- A future daemon-less design?



<https://aya-rs.dev/community>



PromCon
North America 2021



**Please scan the QR Code above
to leave feedback on this session**

Building new projects with bpf

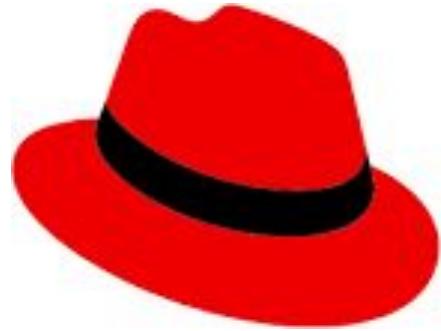
Not sure if we're going to keep this or scrap this...

<https://bpf.dev/getting-started/example-bpf/>

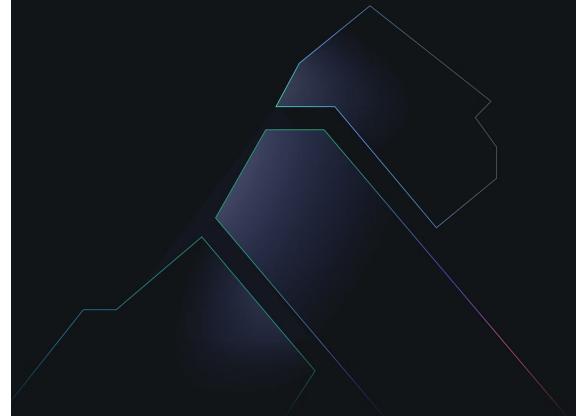
Planning Slide

1. Quick intro to eBPF
2. the problems with ebpf in Kubernetes
3. move to security focused discussion, explain how bad eBPF security posture currently is due to wild west loader, privileged containers, e.t.c.
3. Intro to bpf and the problems it's trying to solve
4. Demo
 - Sample "Normal" application
 - Application gets hijacked
 - root cluster exploit (more likely)
 - Supply chain exploit of OCI container image (less likely)
 - User can see what's happening with Bpf
 - Add some Cosign container image signing

Current Adoption



RedHat



Kong



Kubernetes SIGs



<https://github.com/kubernetes-sigs/blixt>

bpf FD Features: loading + cooperation

Linux Node

Bpf FD State

Name:	xdp_stats
Image URL:	quay.io/bpf FD-bytecode/go-xdp-counter:latest
Pull Policy:	Always
Global:	None
Metadata:	None
Map Pin Path:	/run/bpf FD/fs/maps/39664
Map Owner ID:	None
Maps Used By:	39664
Priority:	55
Iface:	eno2
Position:	0
Proceed On:	pass, dispatcher_return

Bpf FD State

Name:	pass
Image URL:	quay.io/bpf FD-bytecode/xdp_pass:latest
Pull Policy:	Always
Global:	None
Metadata:	None
Map Pin Path:	/run/bpf FD/fs/maps/39662
Map Owner ID:	None
Maps Used By:	39662
Priority:	50
Iface:	eno2
Position:	0
Proceed On:	pass, dispatcher_return

Why ?



KubeCon



CloudNativeCon

North America 2023



<https://rust-lang.org>

- Support for the **BPF Type Format (BTF)**, which is transparently enabled when supported by the target kernel. This allows eBPF programs compiled against one kernel version to run on different kernel versions without the need to recompile.
- Support for **function call relocation and global data maps**, which allows eBPF programs to make function calls and use global variables and initializers.
- **Async support** with both [tokio](#) and [async-std](#).
- Easy to deploy and fast to build: aya doesn't require a kernel build or compiled headers, and not even a C toolchain; **a release build completes in a matter of seconds**.



<https://aya-rs.dev>

BPFD Kubernetes Operator



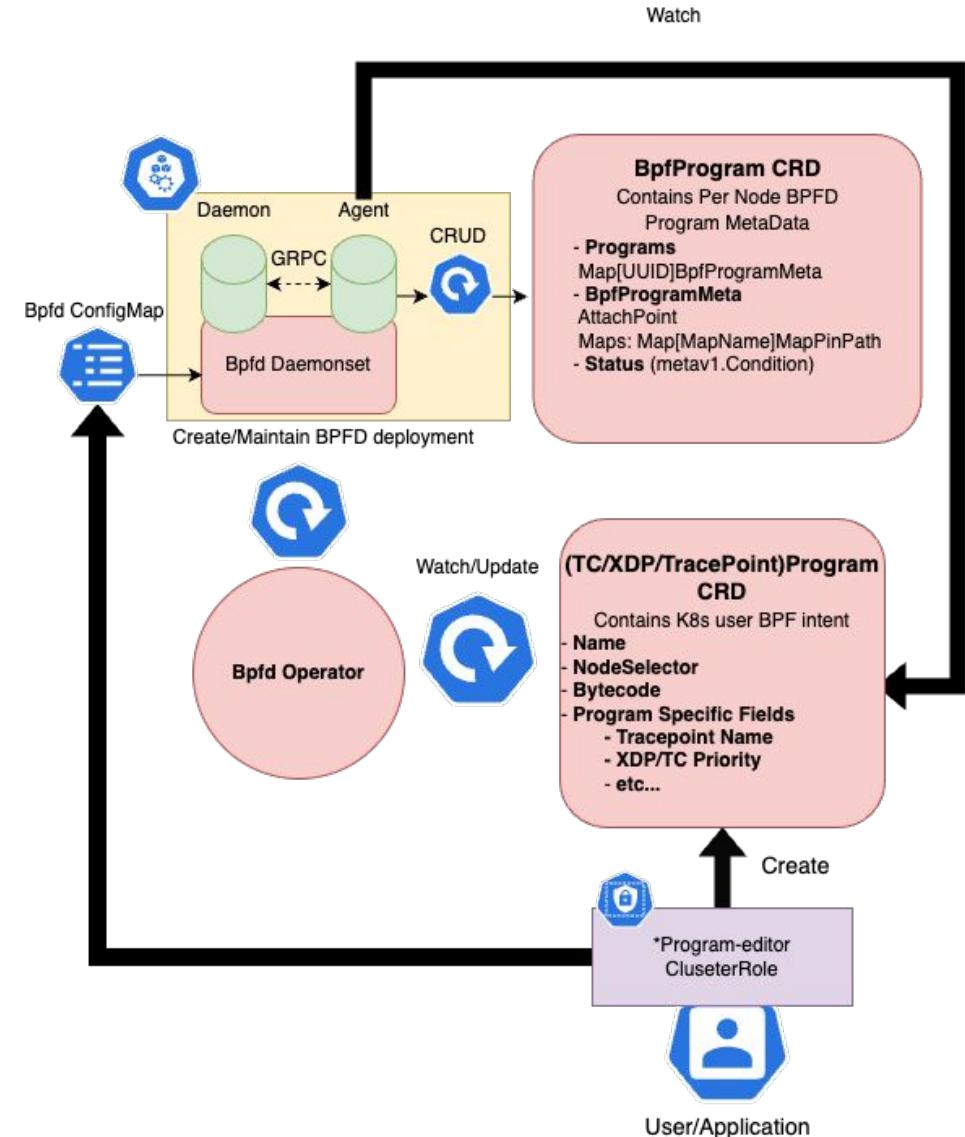
KubeCon



CloudNativeCon

North America 2023

- Bpf Operator
 - Written in go using Operator SDK framework
 - Local kind setup can be tested with `make run-on-kind`
- K8s APIs
 - *Program CRDs
 - BpfProgram CRD
 - Bpf Configmap
- eBPF Bytecode Image Specifications
 - Allows fine-grained separate versioning control for userspace and kernel space
 - Opens the door for signing of these container images in order to verify bytecode ownership



BPFD Security