

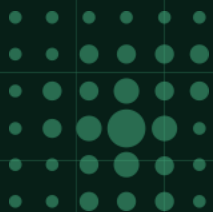


# eBPF 在云原生中的应用

倪朋飞

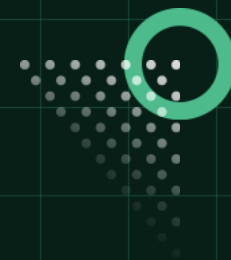
Microsoft Azure

2023.4



# 目录

CATALOGUE



1

eBPF 概述

2

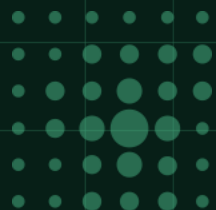
eBPF 应用场景

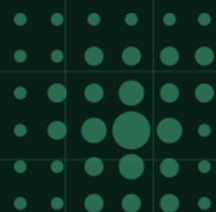
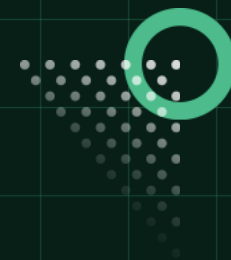
3

eBPF 在云原生中的应用

4

最佳实践分享



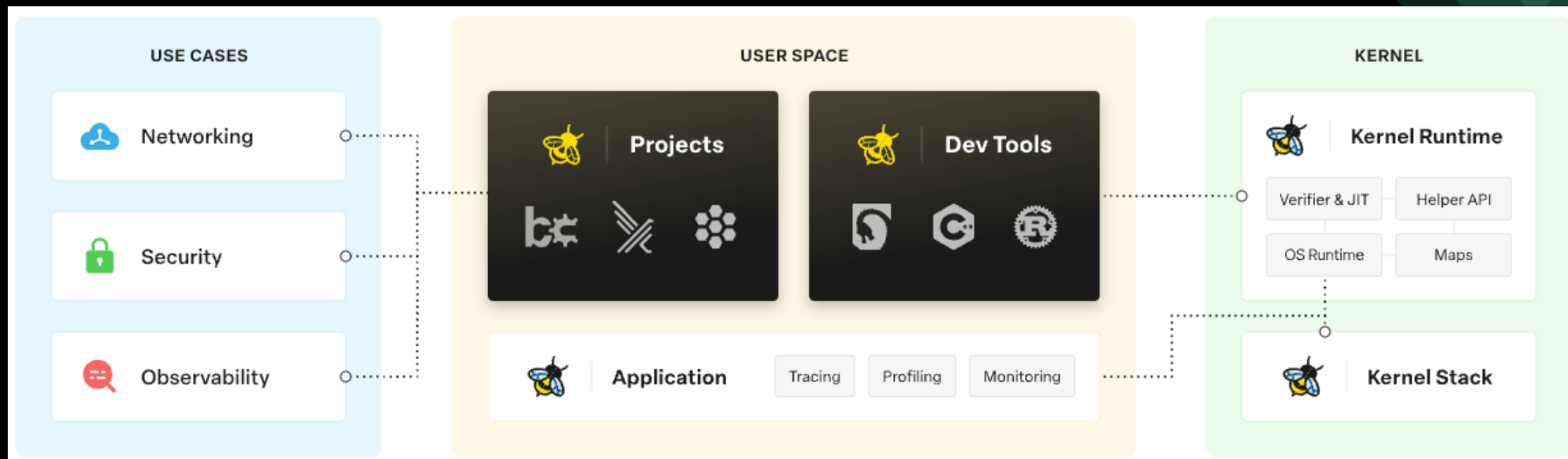


# 01

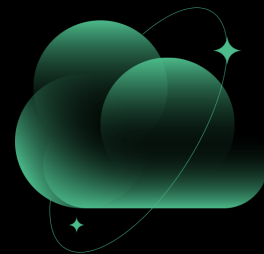
## eBPF 概述



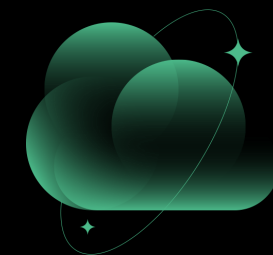
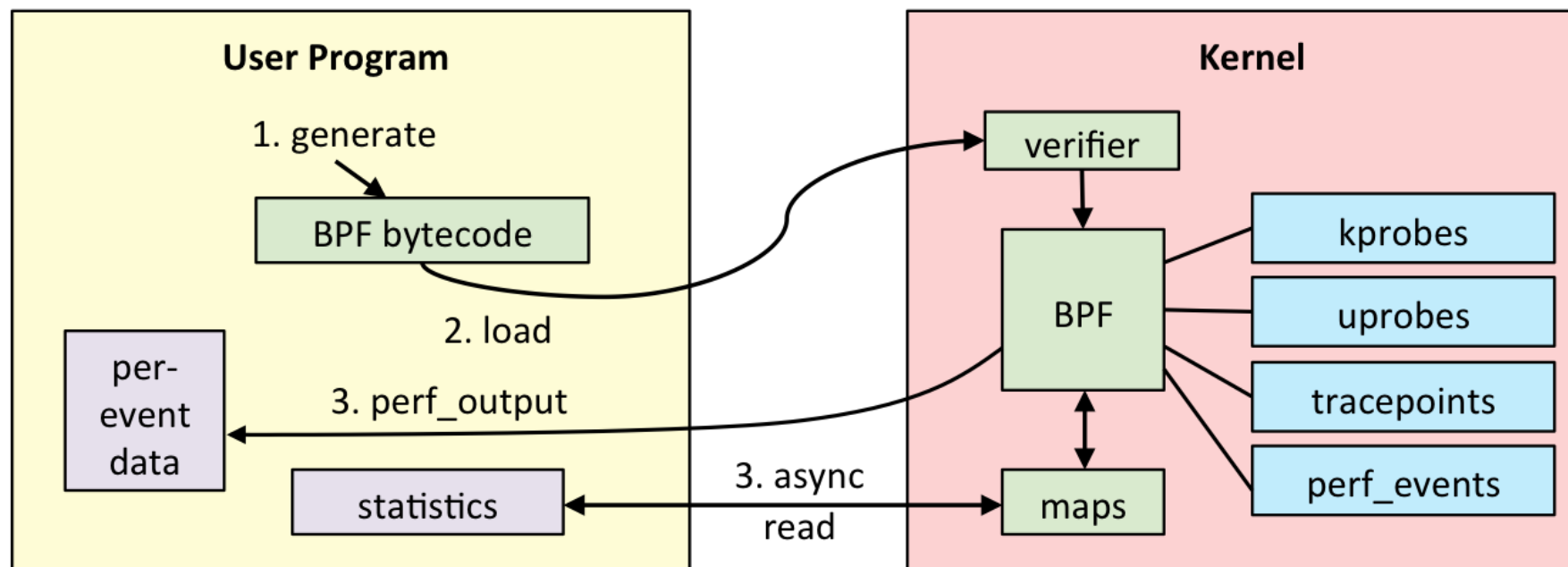
# 什么是 eBPF



- 扩展的伯克利数据包过滤器 (extended Berkeley Packet Filter, eBPF)
- 在 eBPF 之前，内核模块是注入内核的最主要机制，缺少安全控制
- eBPF 借助即时编译器 (JIT)，在内核中运行了一个虚拟机，只有验证安全的 eBPF 指令才会被内核执行
- eBPF 指令运行在内核中，无需向用户态复制数据，这就大大提高了事件处理的效率

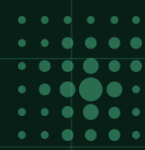


# eBPF 是如何运行的



# Hello, eBPF

2023



hello.c

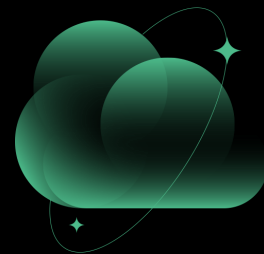
```
int hello_world(void *ctx)
{
    bpf_trace_printk("Hello, World!");
    return 0;
}
```



hello.py

```
# 1) import bcc library
from bcc import BPF

# 2) load BPF program
b = BPF(src_file="hello.c")
# 3) attach kprobe
b.attach_kprobe(event="do_sys_openat2", fn_name="hello_world")
# 4) read and print /sys/kernel/debug/tracing/trace_pipe
b.trace_print()
```



# eBPF 发展历程

2023



1992 The BSD Packet Filter: A New Architecture for User-level Packet Capture

1997 Linux 2.1.75 首次引入了 BPF 技术，将高性能的 BSD 包过滤机制带入 Linux

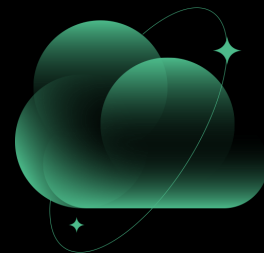
2011 Linux 3.0 中增加的 BPF 即时编译器，进一步优化了 BPF 指令运行的效率

2014 Alexei Starovoitov 为 BPF 带来了第一次革命性的更新，将 BPF 扩展为一个通用的虚拟机，即 eBPF

2015 BCC 大大简化了 eBPF 程序的开发和运行，Linux 4.1 开始支持 kprobe 和 cls\_bpf (用于tc)

2016 Linux 4.7-4.10 带来了跟踪点、perf 事件、XDP 以及 cgroups 的支持，Cilium 项目发布

2017 BPF 成为内核独立子模块，并支持 kTLS/bpftool/libbpf。Netflix/Facebook/Cloudflare 将其用于跟踪、DDoS 防御、4层负载均衡等



# eBPF 发展历程

2023



2018 BPF 新增轻量级调试信息格式 BTF 以及新的 AF\_XDP 类型。Cilium 发布 1.0 版本，bpftool 和 bpftrace 项目也正式发布

2019 BPF 新增尾调用和热更新，GCC 支持 BPF 编译。Cilium 1.6 发布基于 BPF 的服务发现代理（替换基于 iptables 的 kube-proxy）

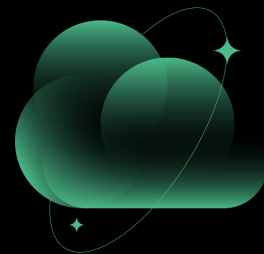
2020 Google 和 Facebook 为 BPF 新增 LSM 和 TCP 拥塞控制，云厂商通过 SRIOV 支持 XDP，Windows 监控工具 Sysmon 增加 Linux 支持

2021 Windows eBPF 发布，eBPF 开始支持内核函数调用，Cilium 发布基于 eBPF 的 Service Mesh

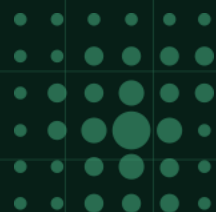
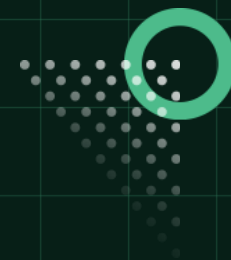
2021 微软、Facebook、Google、Isovalent、Netflix 等一起成立 eBPF 基金会，Cilium 加入 CNCF

2022 eBPF 新增布隆过滤映射、CO-RE、bpf\_loop、动态指针、BPF 内存分配器等，一些列开源项目增加 eBPF 支持

2023 eBPF 火热发展中 ...







02

## eBPF 应用场景

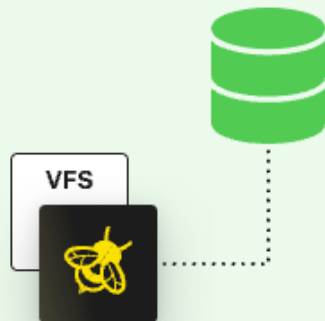


# eBPF 有哪些应用场景



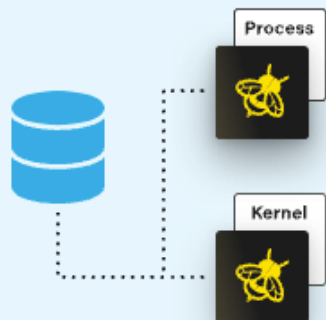
## 网络

在内核空间绕过协议栈直接处理网络包，加速网络过滤和处理的性能。



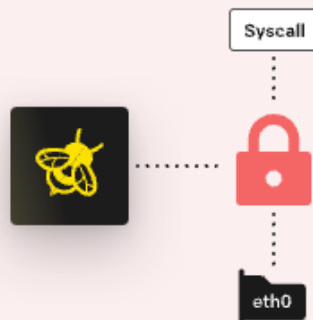
## 可观察性

从内核空间直接收集各类性能指标，而无需应用程序作任何修改。



## 追踪与剖析

将 eBPF 程序附加到跟踪点以及内核和用户应用程序探测点，追踪和剖析系统性能问题。

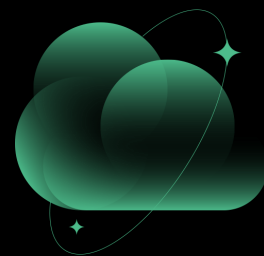
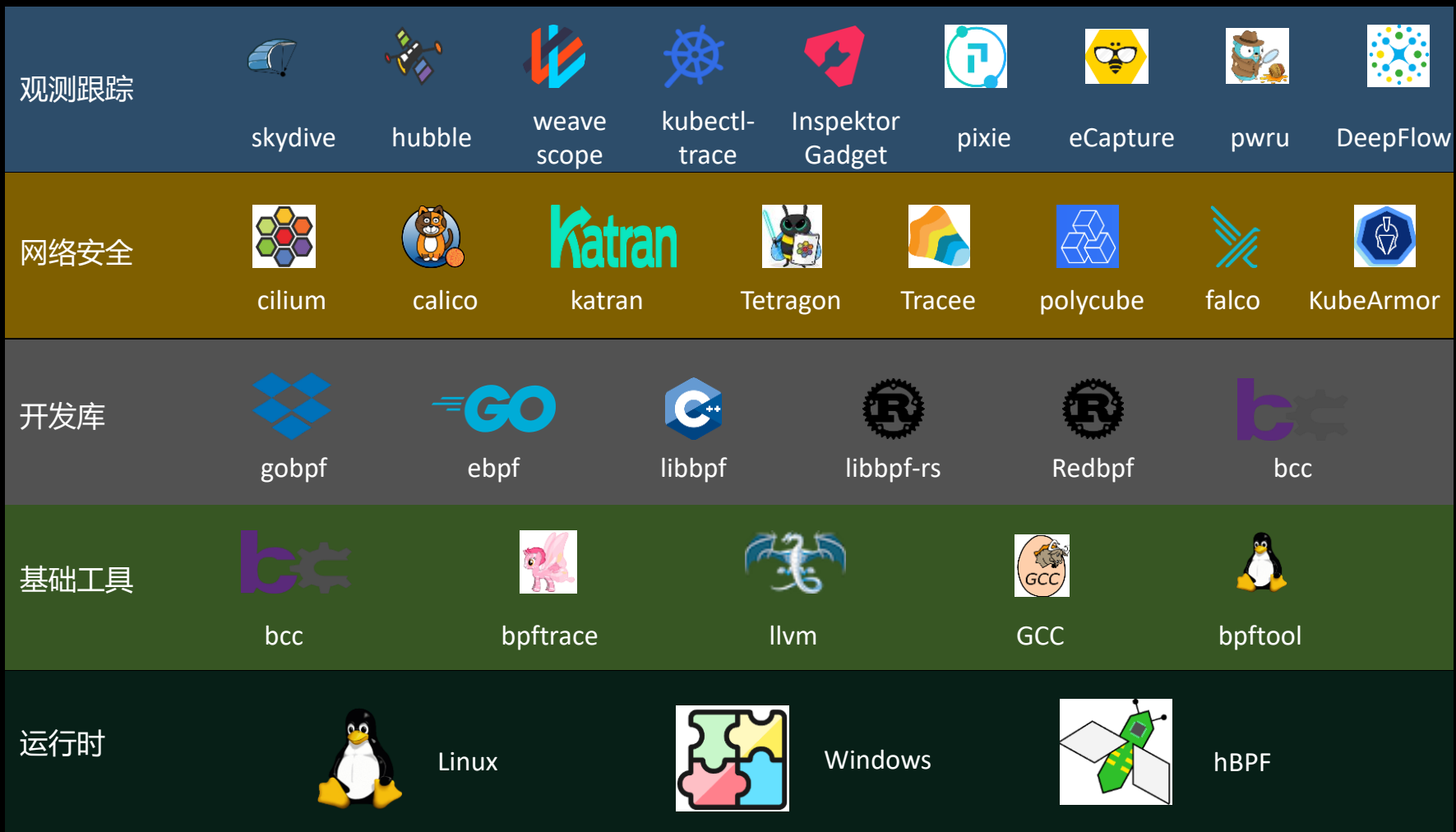


## 安全

动态分析内核和应用的行为，监测各类安全事件，并执行相应的安全策略。

# eBPF 全景图

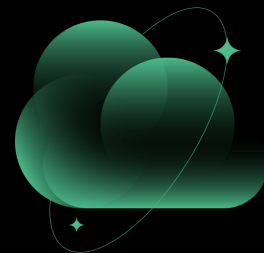
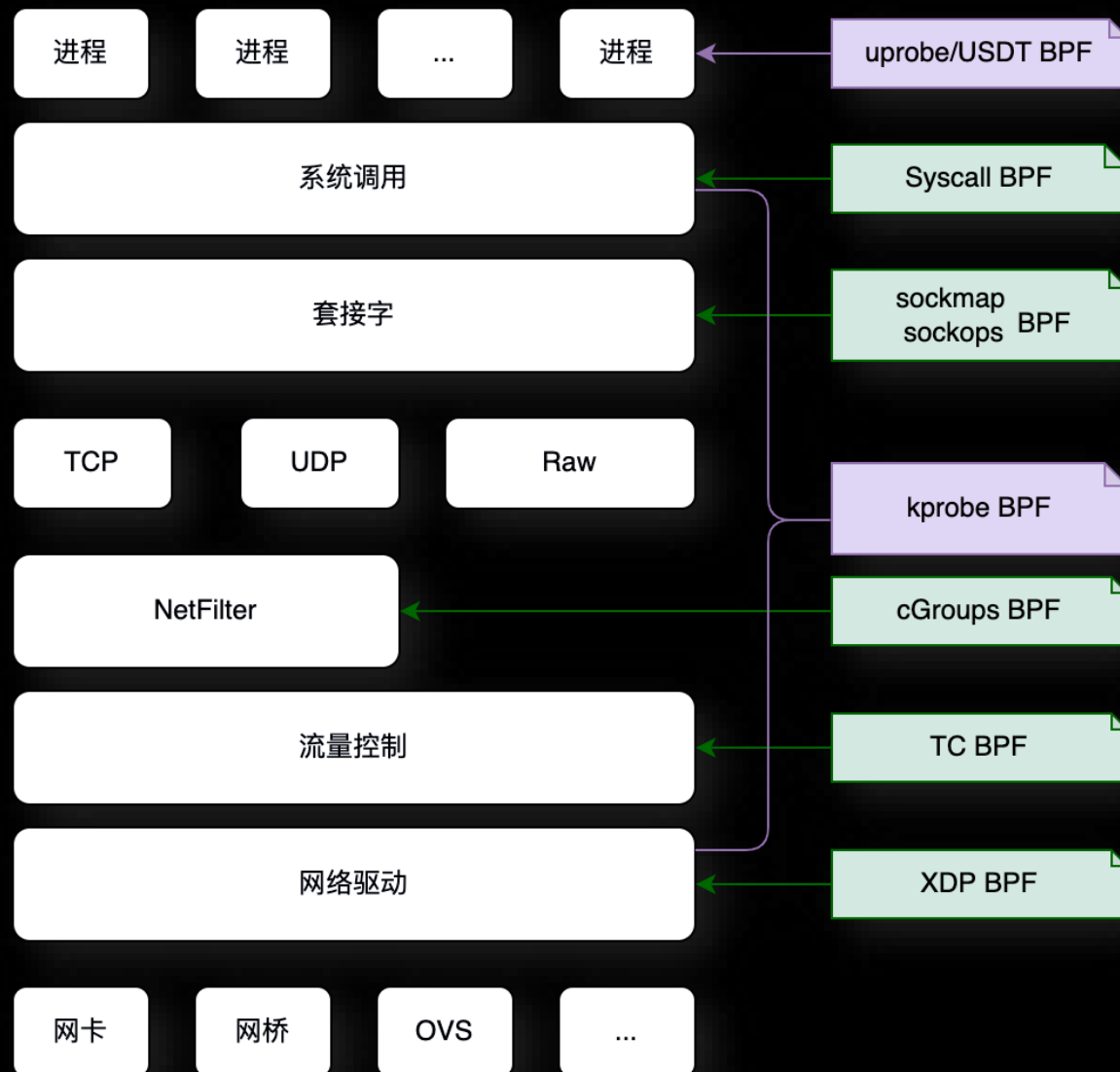
2023

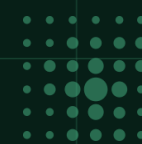




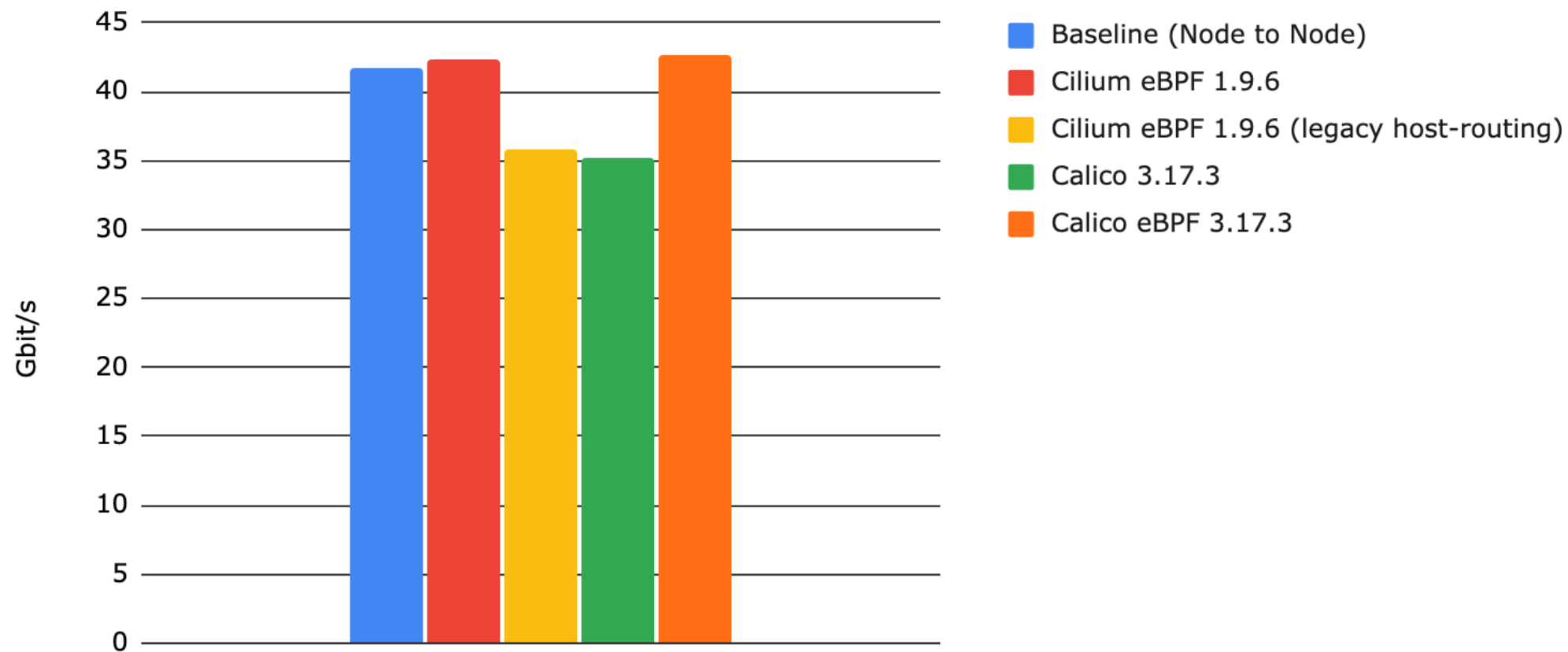
03

## eBPF 在云原生中的应用





## TCP Throughput (1 Stream) - Higher is better

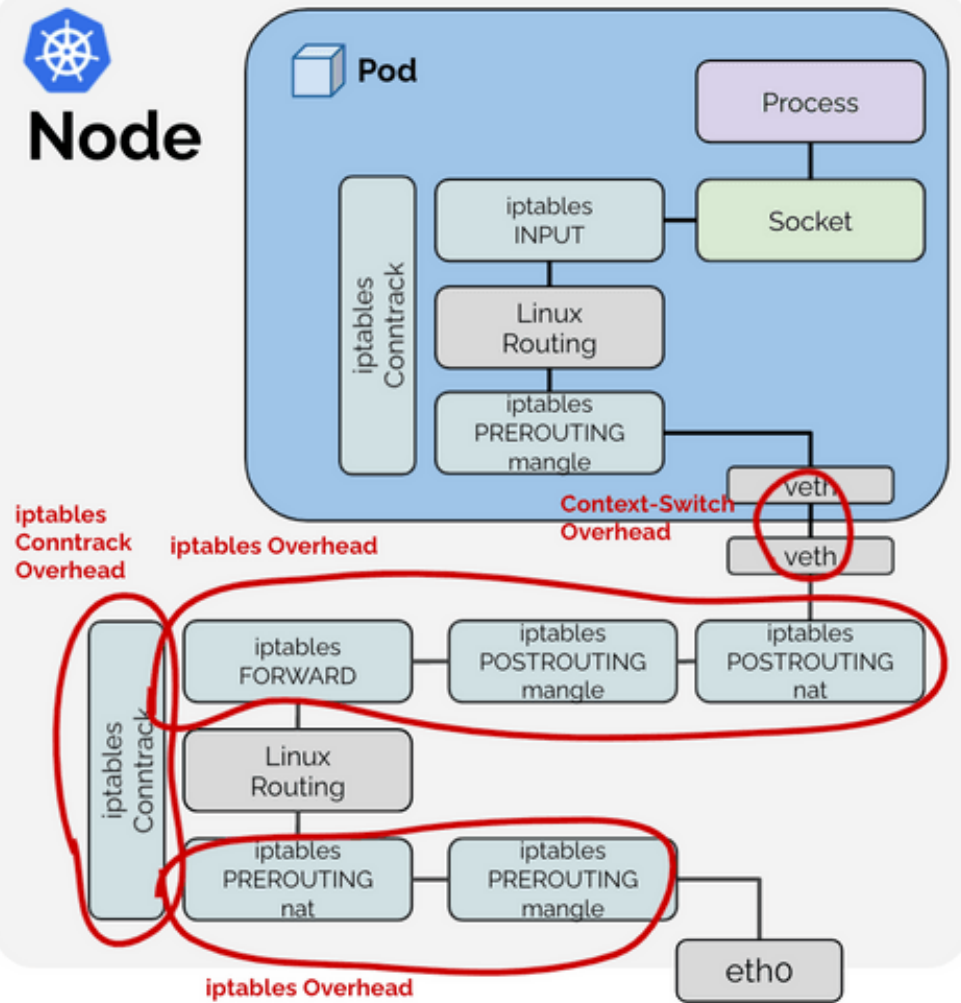


# Cilium eBPF

2023



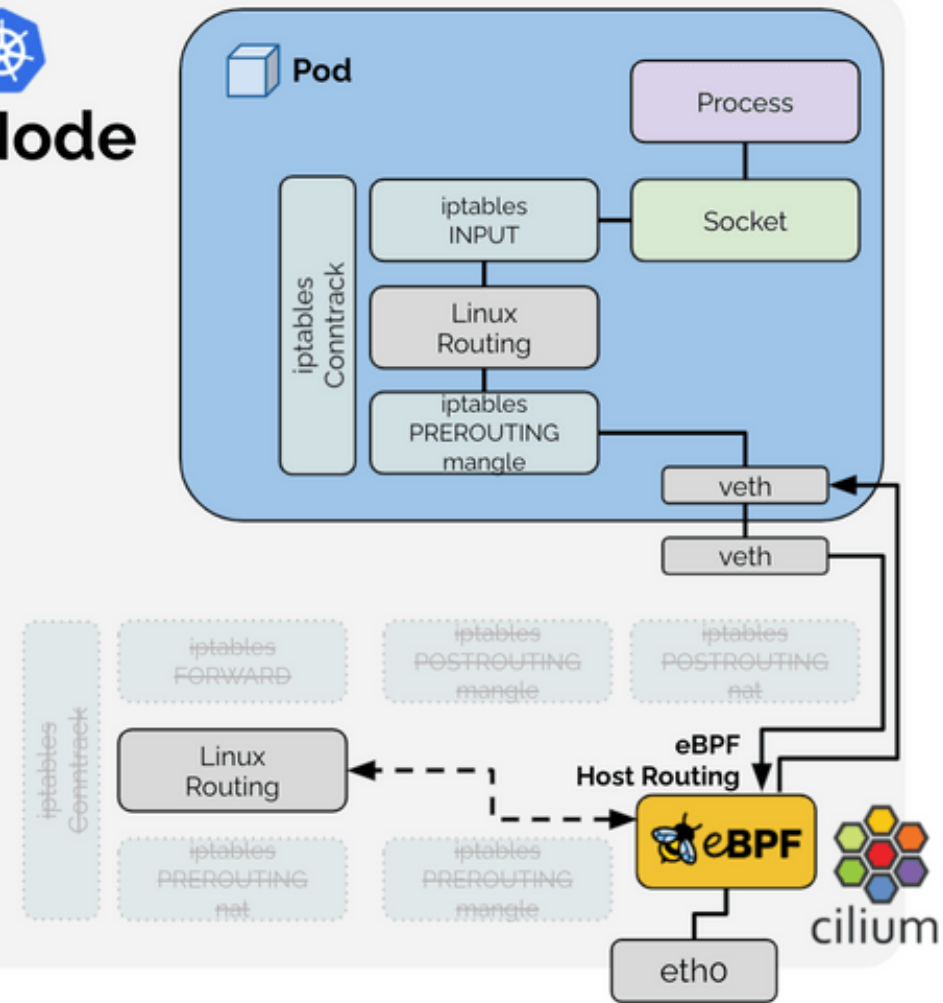
Node



Standard Container Networking

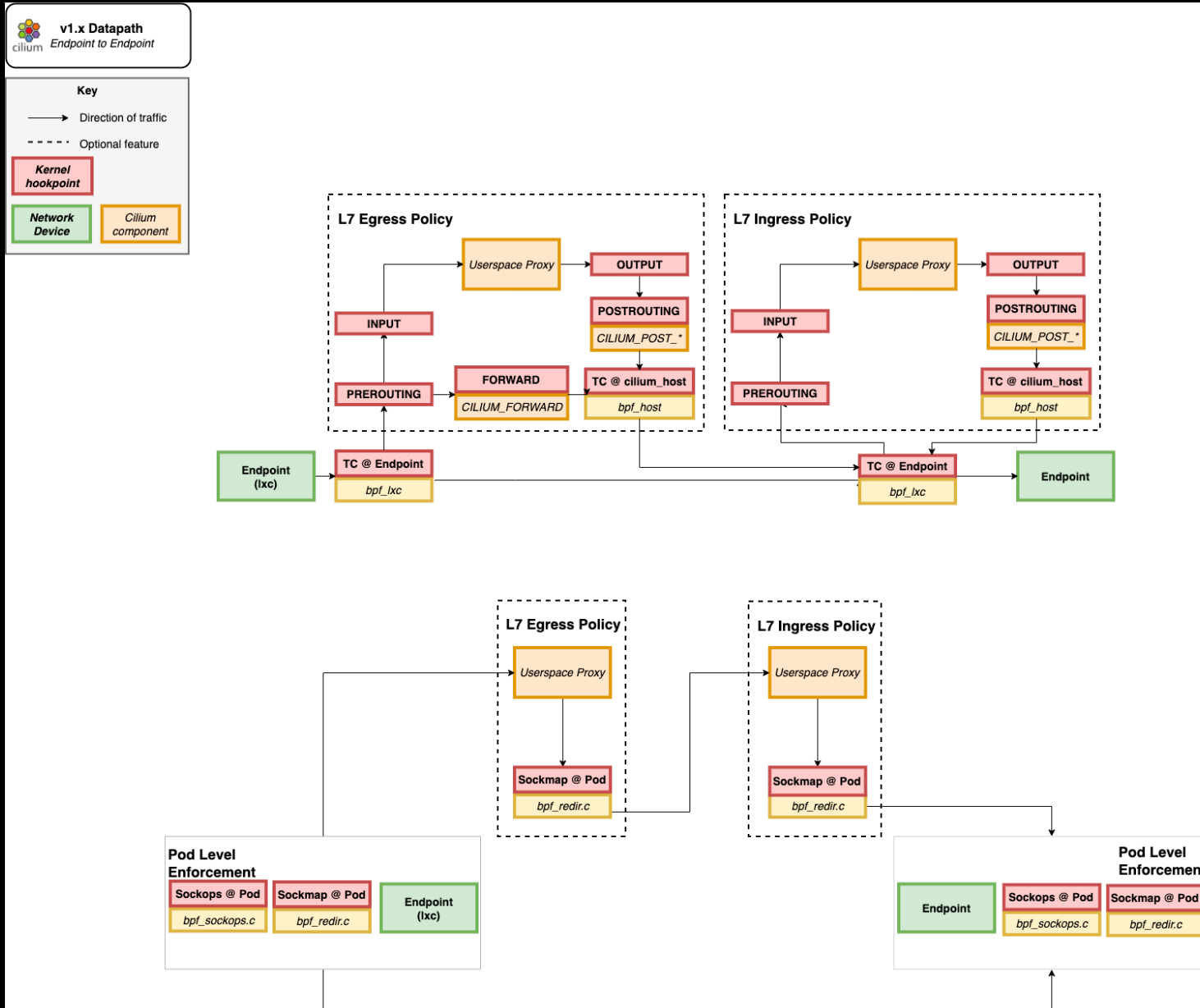


Node



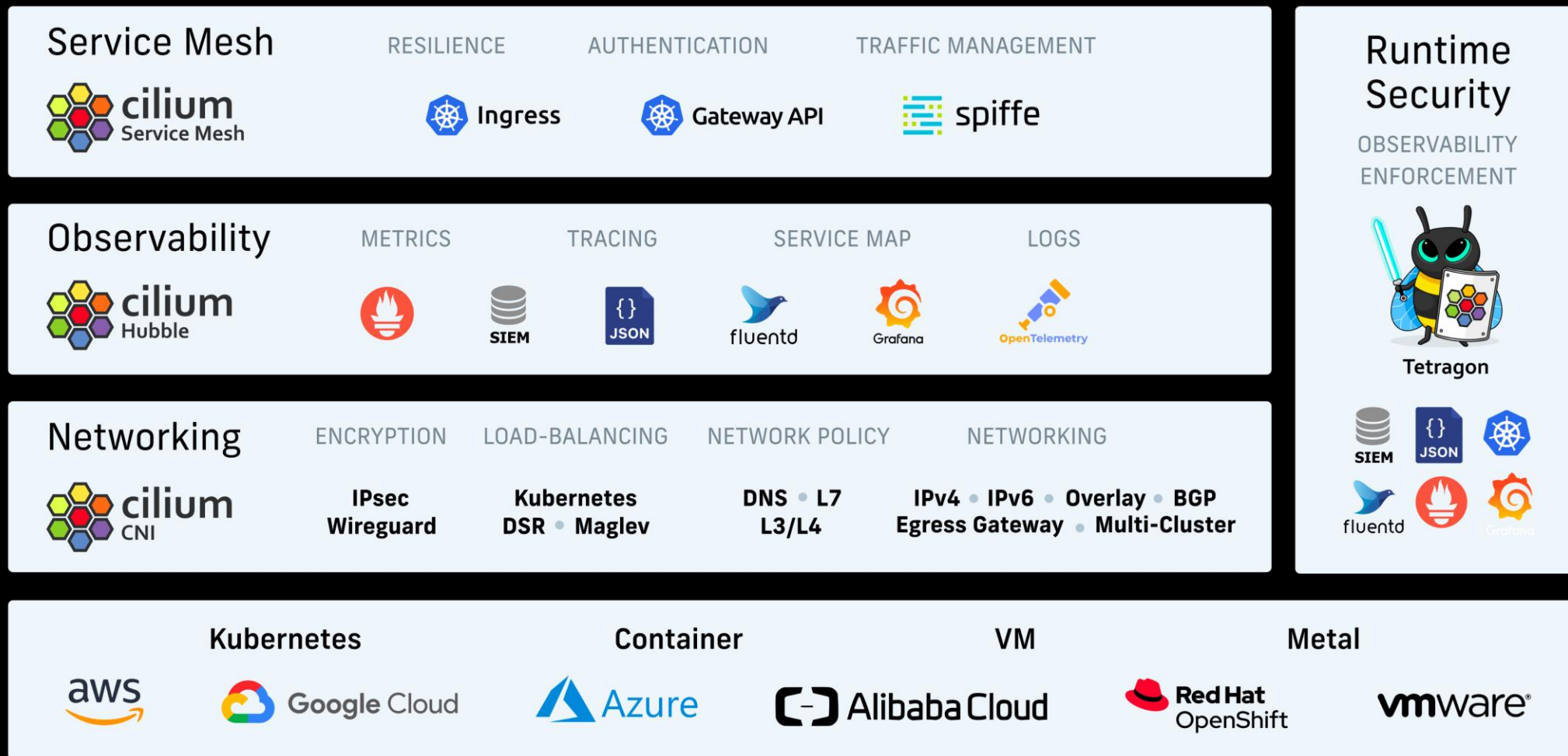
Cilium eBPF Container Networking

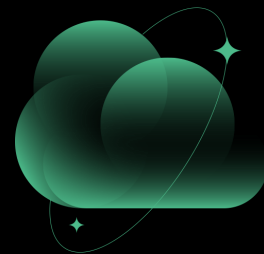
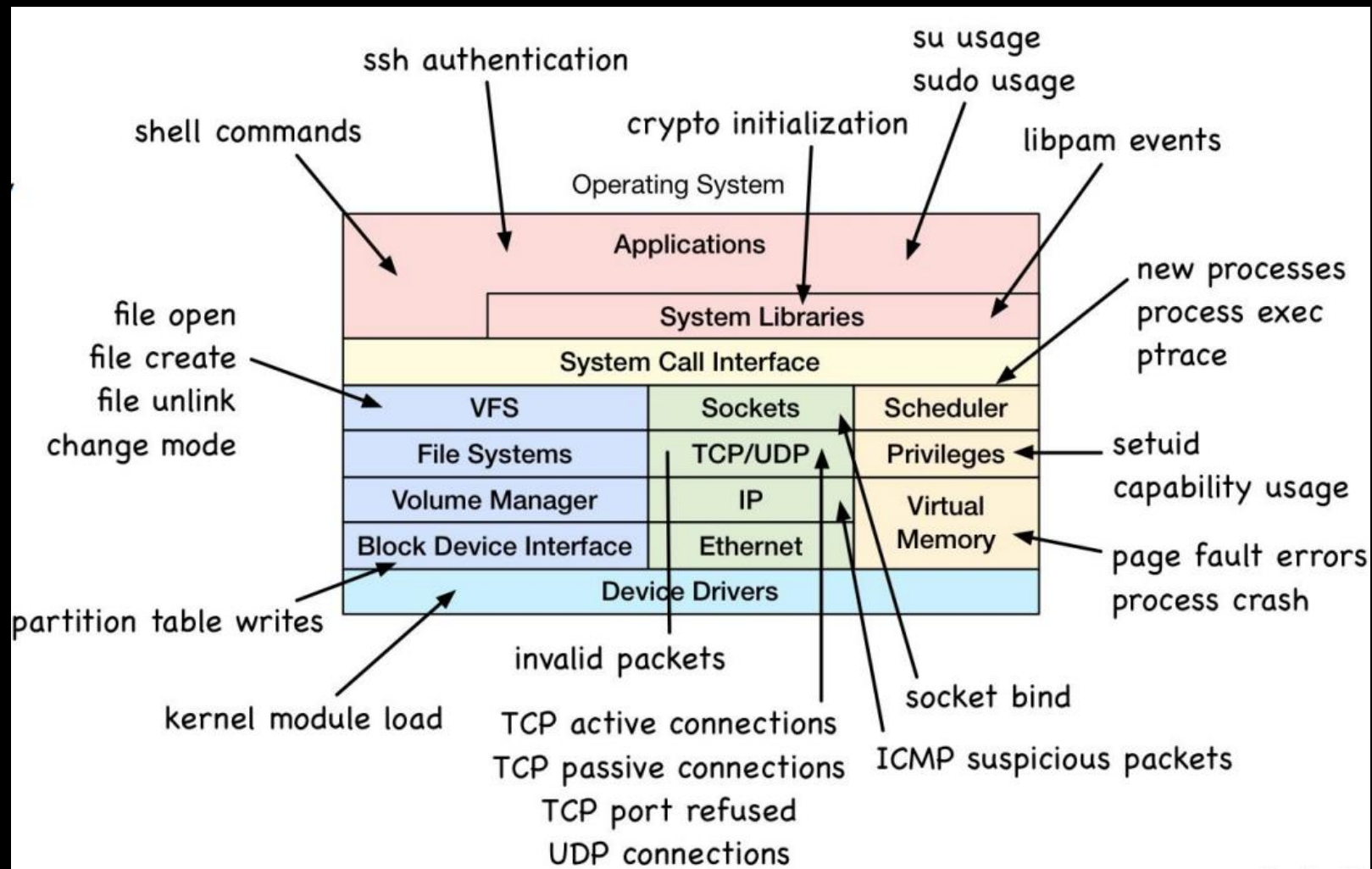
# Cilium eBPF





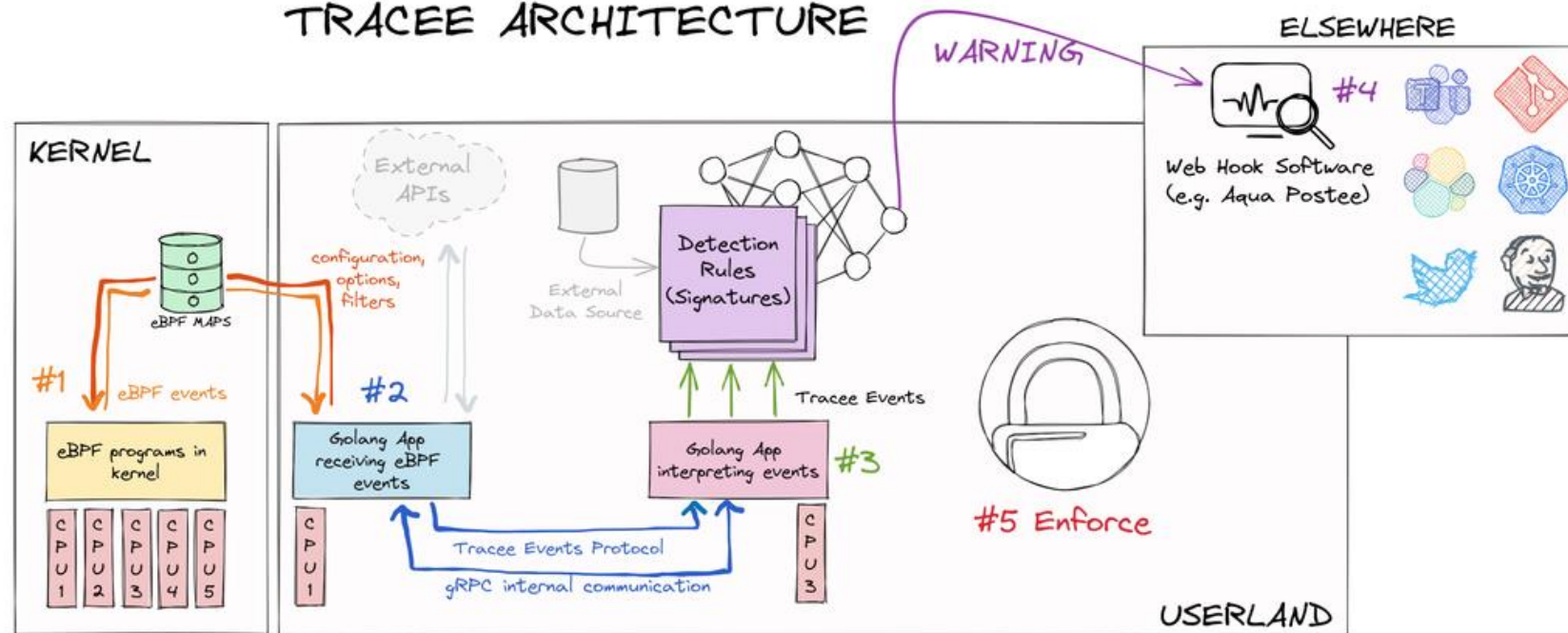
# Cilium 整体架构





# Tracee

## TRACEE ARCHITECTURE



Flow of Events (Pipeline)

Flow of Configuration Needs

### Tracee eBPF Event on OpenAt

```
Event Type: OpenatEventID
Name: "openat"
Probe: "openat" syscall

Params: {
  Type: "int", Name: "dirfd"
  Type: "const char", Name: "pathname"
  Type: "mode_t", Name: "mode"
}
```

### Tracee Rules Signature

```
Metadata:
  ID: TRC-X

Selected Events:
  "OpenAt"

Match:
  "Pathname" == "/dir/Filename"
```

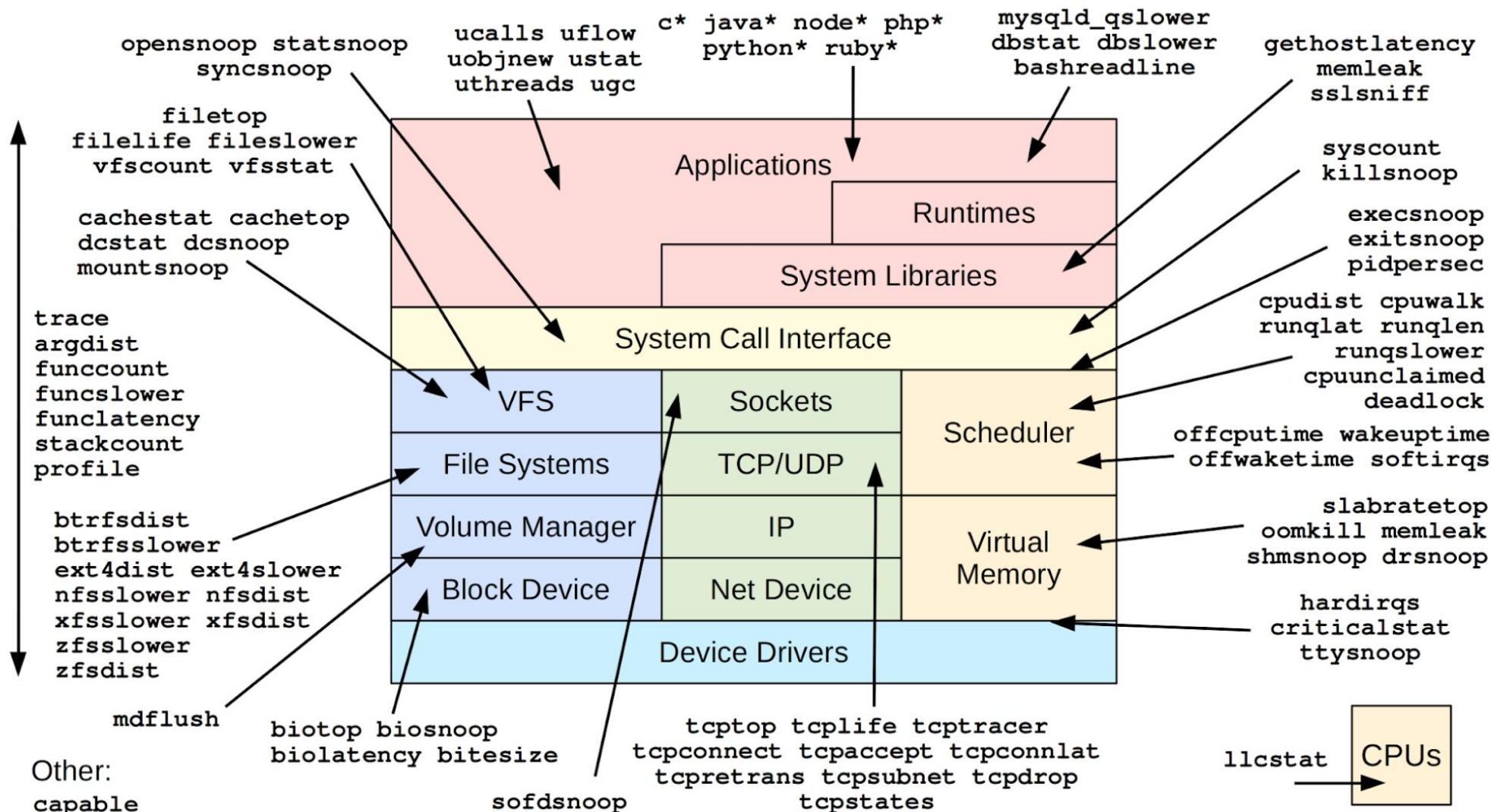
Security in 4 steps:

- #1 Generate
- #2 Collect
- #3 Detect
- #4 Consume
- #5 Enforce

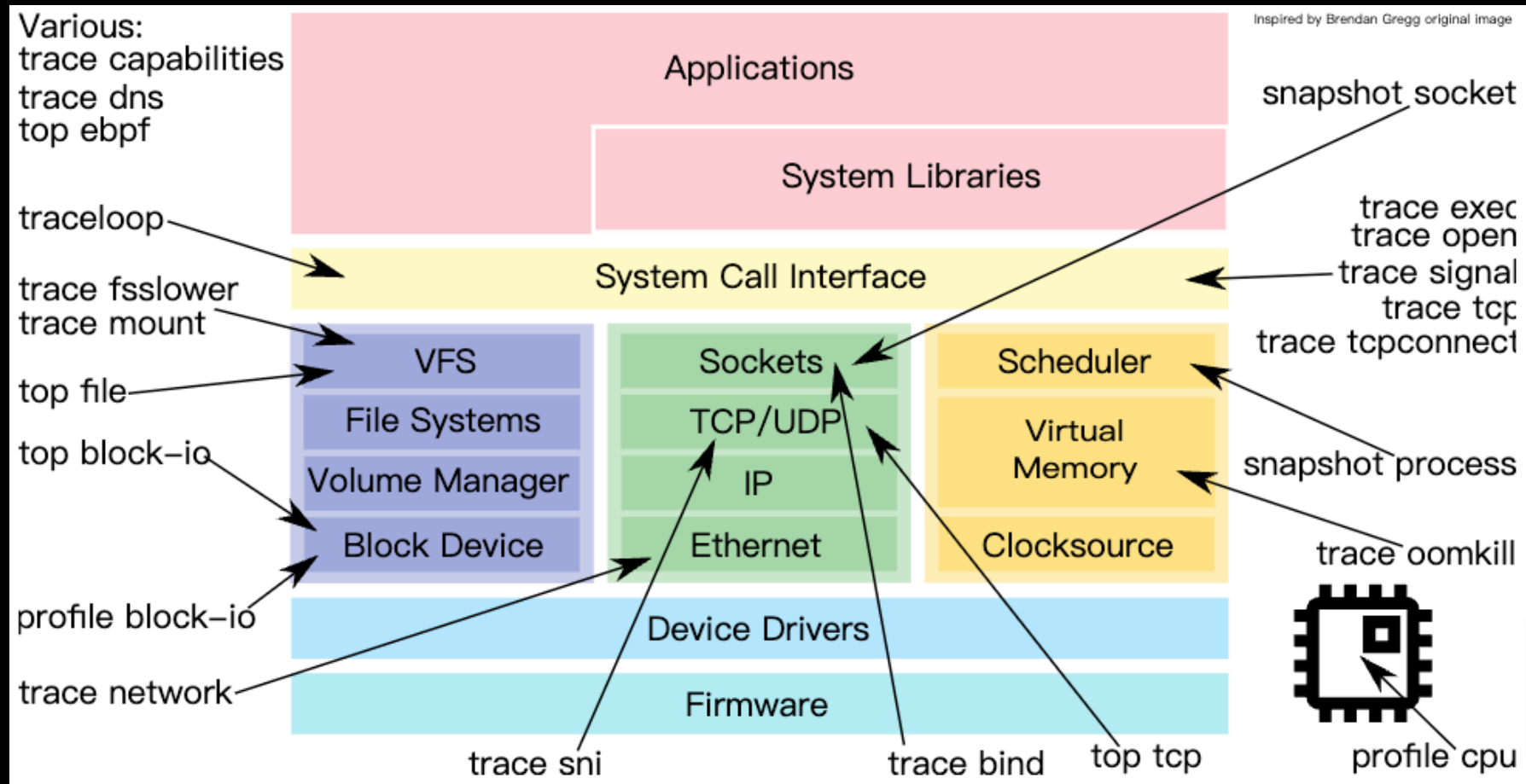


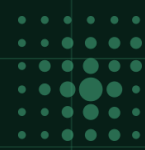
# 追踪与剖析

## Linux bcc/BPF Tracing Tools



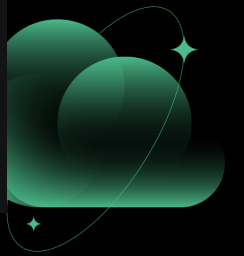
# Inspektor Gadget

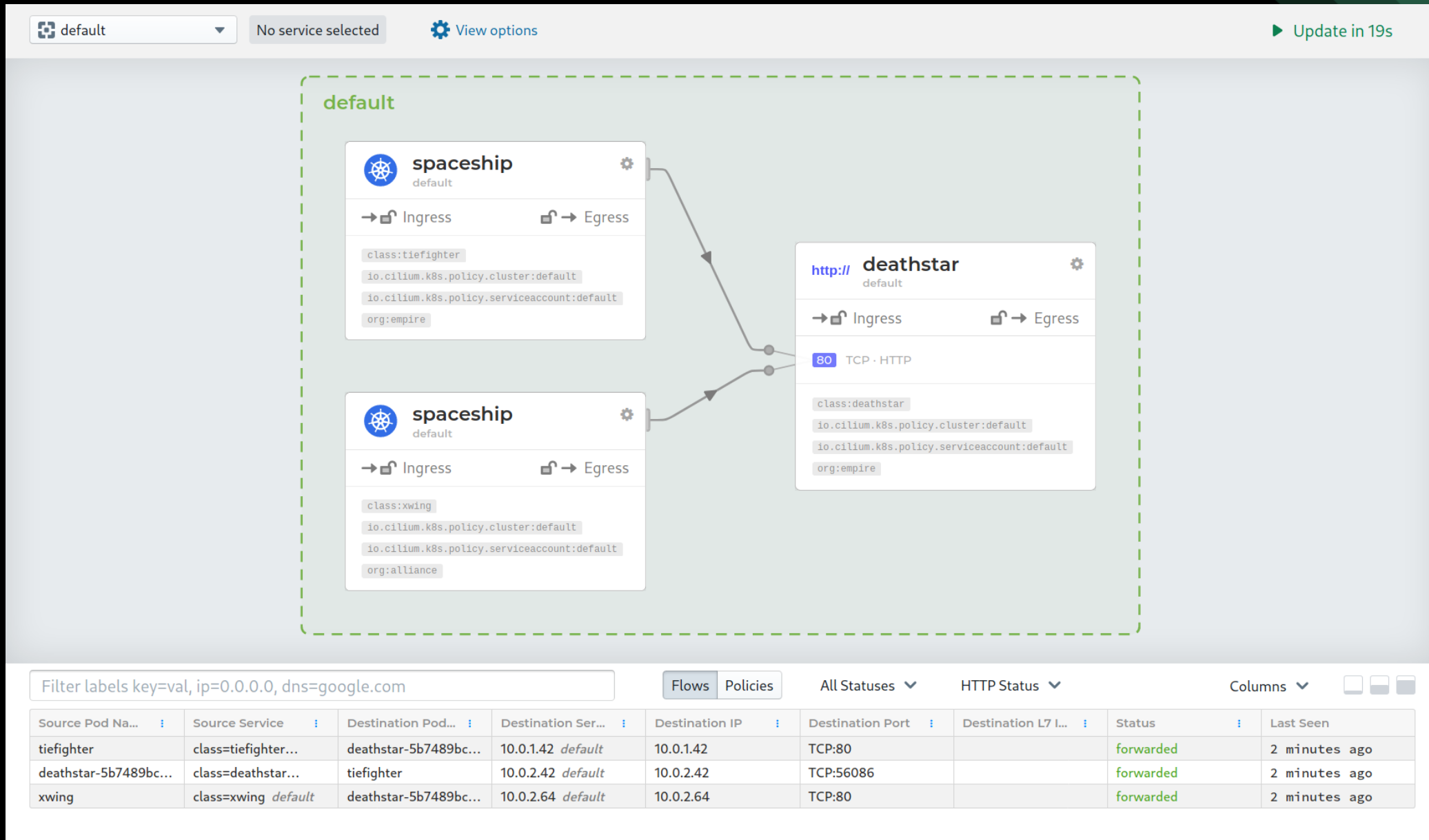




ceuse#

ceuse#

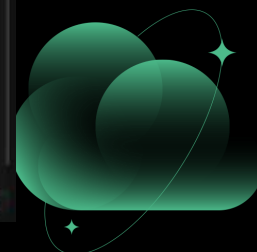




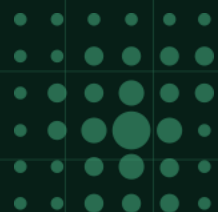
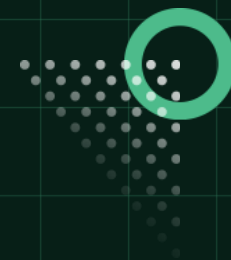


The screenshot displays the Pixie web interface for a Kubernetes cluster named 'px'. The main view shows a network diagram of the 'online-boutique' namespace. The diagram illustrates the flow of traffic between various services, with the 'frontend' service highlighted in blue. Other visible services include 'frontend-external', 'checkoutservice', 'productcatalogservice', 'shippingservice', 'currencyservice', 'paymentservice', 'recommendationservice', and 'inventoryservice'. A sidebar on the left lists pods with their IDs and creation times. The bottom section features a 'Service List' table with columns for service name, latency, requests per second, error rate, inbound bytes per second, and outbound bytes per second. The table shows data for four services: 'online-boutique/frontend-external', 'online-boutique/checkoutservice', 'online-boutique/productcatalogservice', and 'online-boutique/shippingservice'. The 'frontend-external' service shows a latency of 1226.45 ms and 2.01 requests per second. The 'checkoutservice' shows a latency of 292.45 ms and 0.23 requests per second. The 'productcatalogservice' shows a latency of 0.62 ms and 9.79 requests per second. The 'shippingservice' shows a latency of 0.63 ms and 0.69 requests per second.

service	latency_ms (p99)	requests_per_s	error_rate_pct	inbound_bytes_per_s	outbound_bytes_per_s
online-boutique/frontend-external	1226.45	2.01	0	2.01	801.47
online-boutique/checkoutservice	292.45	0.23	0	19.55	27.08
online-boutique/productcatalogservice	0.62	9.79	0	130.07	2899
online-boutique/shippingservice	0.63	0.69	0	54.76	25.94







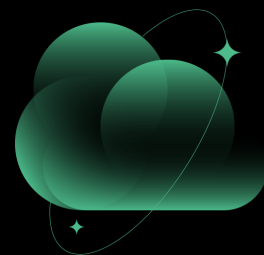
04

# 最佳实践分享





- 借助 eBPF 了解内核，根据内核原理选择最佳的 eBPF 程序类型
- libbpf + CO-RE (Compile Once – Run Everywhere) + BTF
- 注意 eBPF 特性在不同内核版本的差异
- 利用 eBPF 帮助函数、映射与内核进行交互
- 善用 `__always_inline`、尾调用等降低复杂度并绕过 eBPF 验证问题
- 防御 eBPF 被恶意程序利用和伪装





# Thanks

