



GOTC 2023

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

OPEN SOURCE, INTO THE FUTURE

「eBPF」专场

本期议题：BPF 技术在百度云原生实践

狄卫华 2023年05月28日

关于我

狄卫华

- 架构师@百度云原生团队
- 国内 BPF 技术布道师
- 《Linux内核观测技术 BPF》译者之一
- 国内首个 BPF 技术站点 ebpf.top



个人微信

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE



> 1. BPF 技术简介

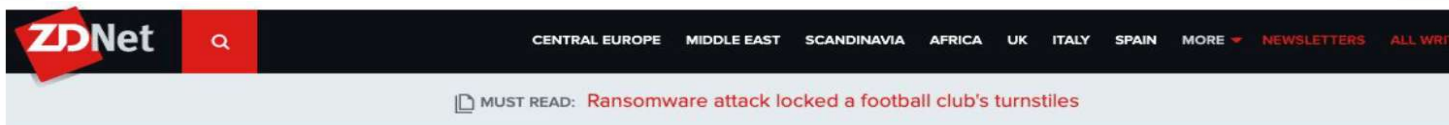
2. BPF 技术在云原生应用

3. BPF 技术在百度云原生实践

1. BPF 技术简介

BPF 超能力诞生

BPF: 过去 50 年操作系统最大的变革 !!!



Netflix: BPF is a new type of software we use to run Linux apps securely in the kernel

A Netflix performance architect says BPF promises a fundamental change to a 50-year-old kernel model.

"BPF is the biggest operating systems change I've seen in my career, and it's thrilling to be a part of it," wrote Gregg.

Today's kube-proxy replacement through BPF is just a tiny dot in that universe ...

Steven Rostedt
@srostedt

BPF will replace Linux #kr2019

11:06 am · 26 Sep 2019 · Twitter for Android

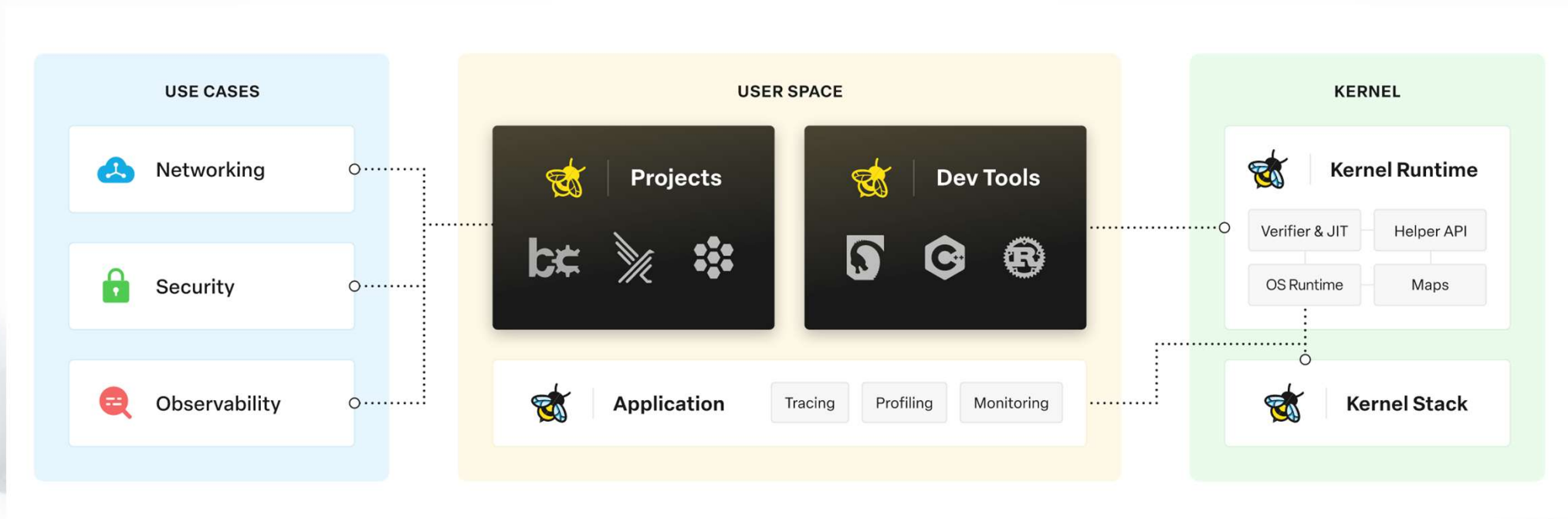


全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

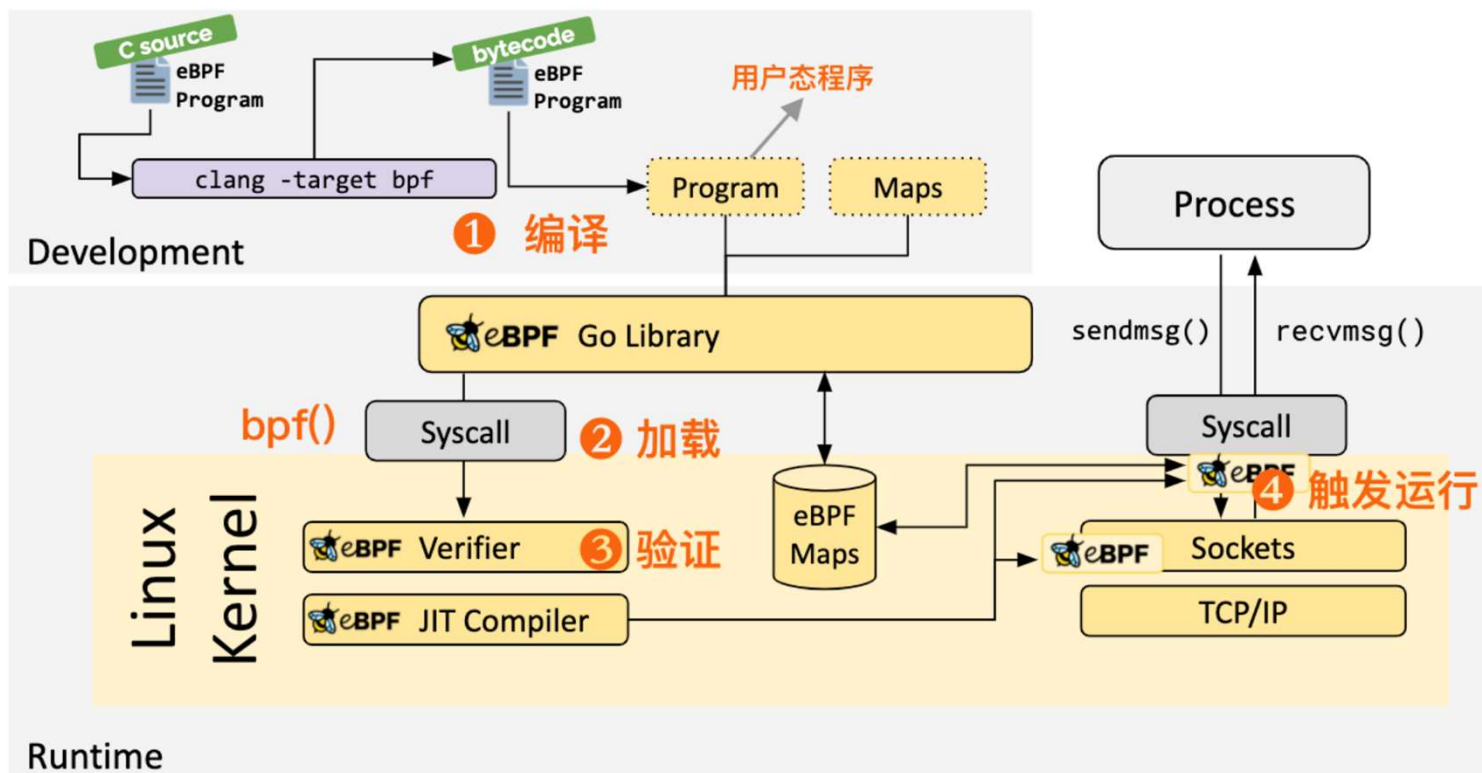
1. eBPF 技术简介

BPF 整体概览



1. BPF 技术简介

BPF 工作原理：内核空间高效稳定运行用户编写程序

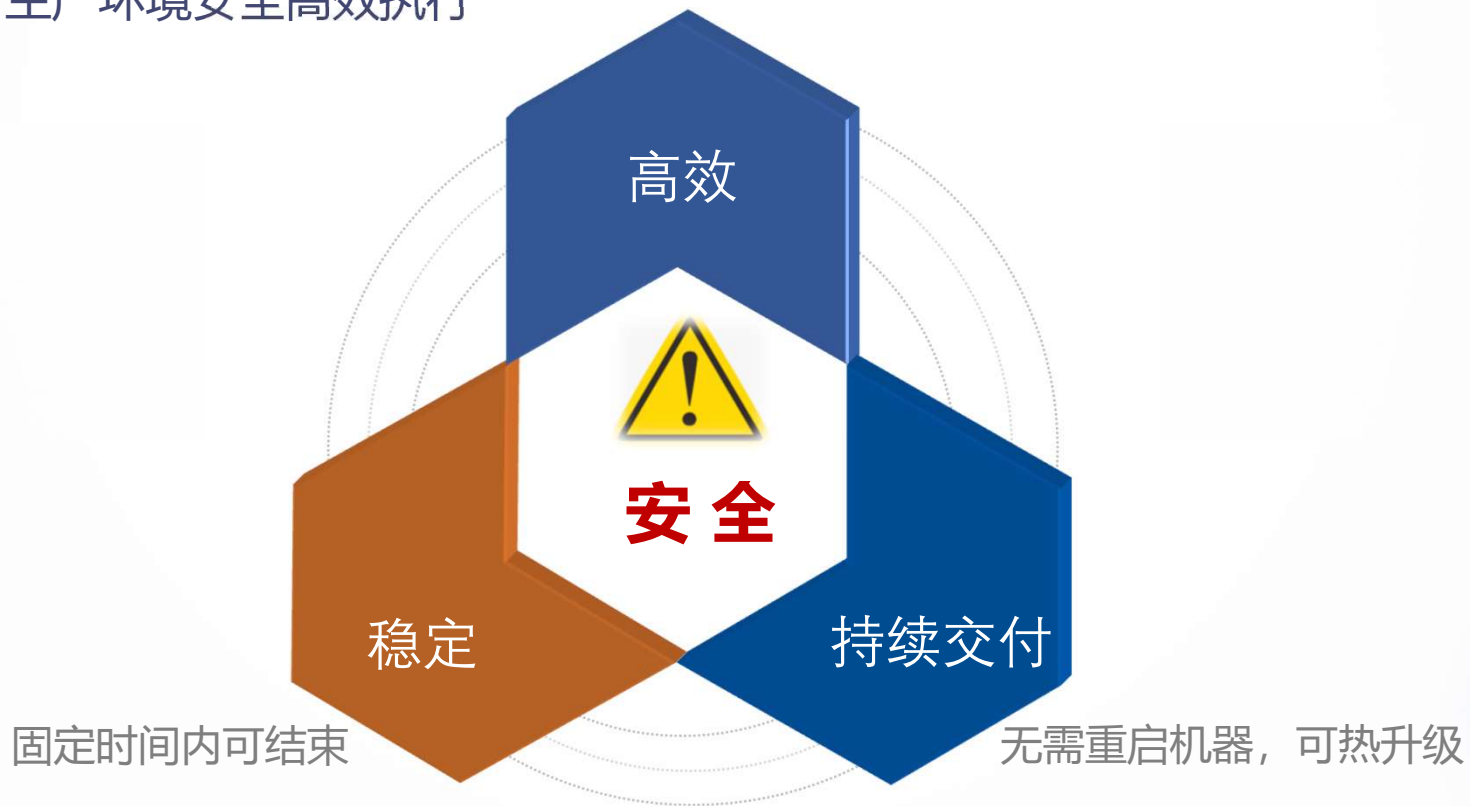


全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

1. BPF 技术简介

BPF 生产超能力：生产环境安全高效执行 JIT 后采用原生 CPU 指令



1. BPF 技术简介

BPF 生产超能力：广泛的生产级别应用



安全审计/数据包处理/性能监测



网络洞察力



网络监测/电源监控/内存性能



数据中心网络负载均衡



内核安全监控



网络安全/性能监测/网络可观测



云容器网络 Cilium



负载均衡



系统跟踪调试



系统跟踪调试/容器网络加速

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

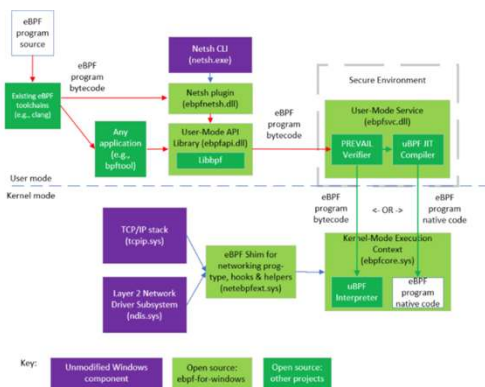
数据来源：<https://ebpf.io/case-studies/>，有调整

1. BPF 技术简介

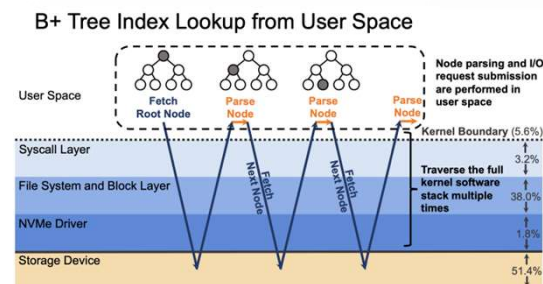
BPF 技术未来：是基石，也是平台能力补齐的核心拼图



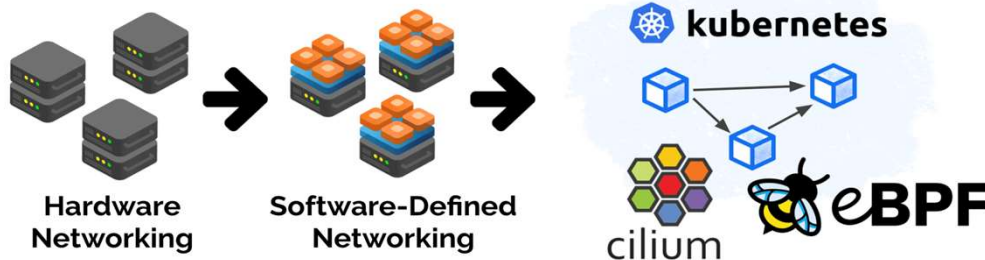
软件定义内核



攻城拔寨：Windows



存储访问路径加速：XRP



eBPF 定义云原生网络



目
录

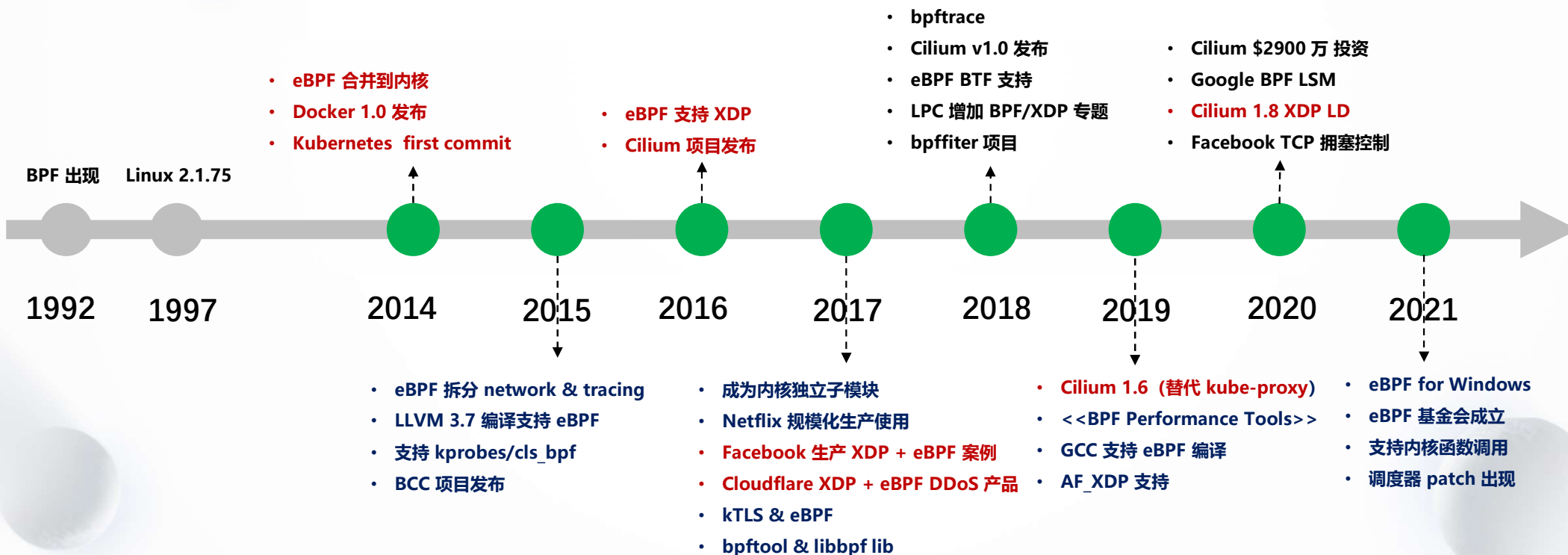
1. BPF 技术简介

➤ 2. BPF 技术在云原生应用

3. BPF 技术在百度云原生实践

2. BPF 技术在云原生应用

BPF 与云原生



全球开源技术峰会

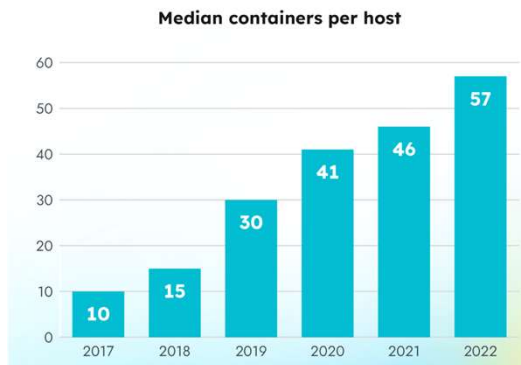
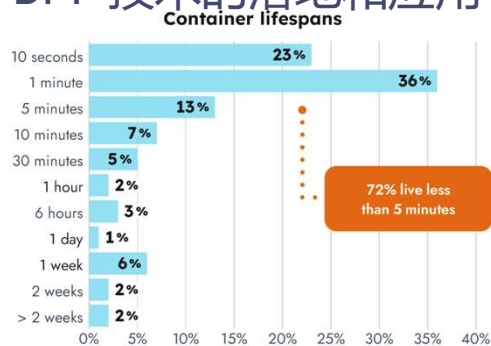
THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

云原生与 eBPF 技术出现位于同期，发展过程相互促进，相互成就

2. BPF 技术在云原生应用

云原生技术趋势：加速 BPF 技术的落地和应用

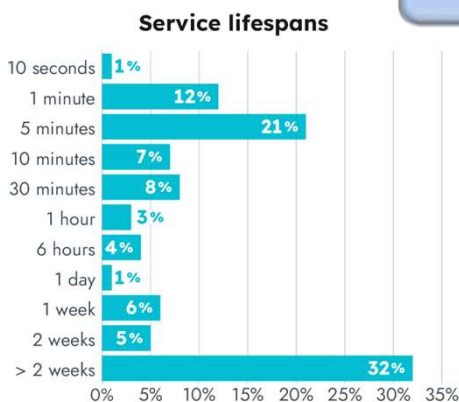
容器存活时间



单机平均容器密度

降本增效

服务存活时间



Security

87% of container images have high or critical vulnerabilities

90% of granted permissions are not used

15% of high and critical vulnerabilities are in use at runtime

云原生安全挑战

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

数据来源：2023 SysDig 云原生安全和容器使用报告

2. BPF 技术在云原生应用

云原生技术挑战：需要更强技术助力



开源组件多，版本多，问题排查链路长，涉及多团队协作，排查效率和协同效率面临挑战

容器实例存活时间变短，细粒度可观测和问题排查提出更高要求

单机部署容器实例密集度加速，离在线混部成为趋势，运行效率和质量细粒度感知需加强

服务生命周期敏捷化，基础能力敏捷化和透明化，部分场景网络性能较高诉求

云原生安全挑战与日俱增，运行时监测和网络安全增强需要新的技术手段护航

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

2. BPF 技术在云原生应用

云原生应用： BPF 技术为云原生带来超能力

eBPF - Superpowers for Networking, Observability & Security

Liz Rice

Chief Open Source Officer, Isovalent
Chair, CNCF Technical Oversight Committee

@lizrice



全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

2. BPF 技术在云原生应用

BPF & 云原生：超能力落地



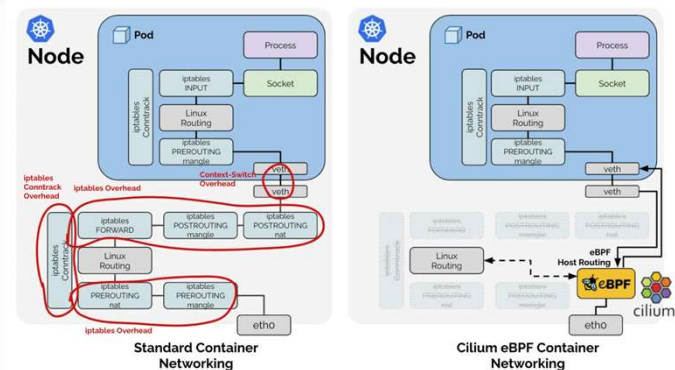
全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

2. BPF 技术在云原生应用

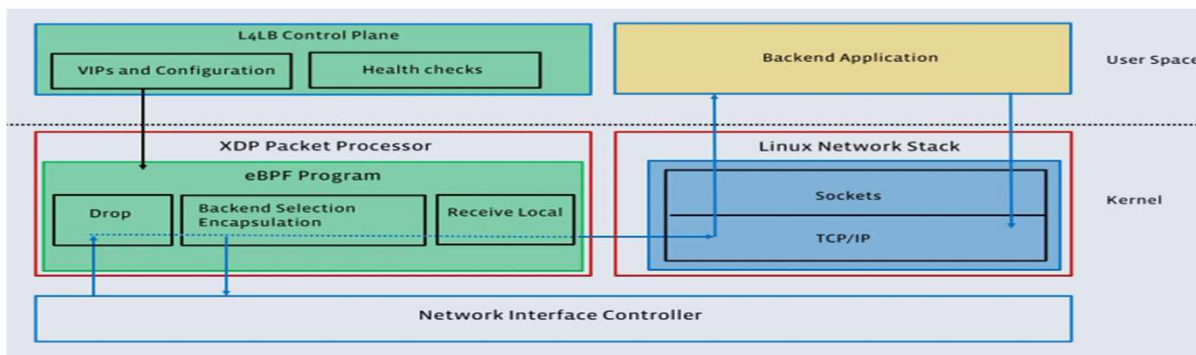
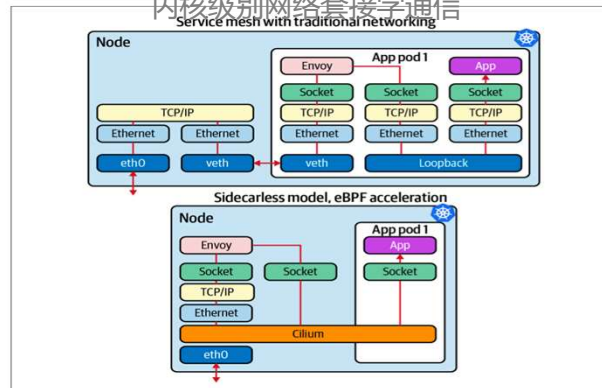
云原生网络

容器网络路径 绕过主机网络协议栈



Service Mesh 网络路径

内核级别网络套接字通信



全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

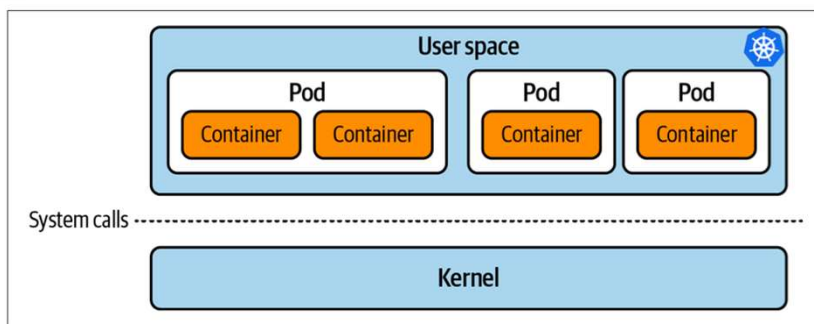
网络负载均衡器

传统 IPVS 技术 4 倍以上性能提升

2. BPF 技术在云原生应用

云原生安全

运行时安全 内核视图联合检测

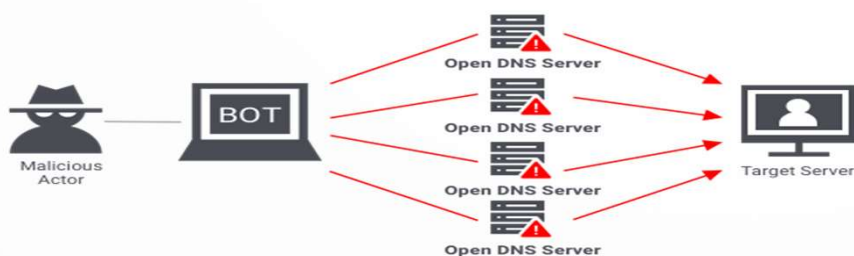


网络安全策略

L3/L4/L7 服务标识管控

```

1  apiVersion: cilium.io/v2
2  kind: CiliumNetworkPolicy
3  metadata:
4    name: store-policy
5    namespace: my-store
6  spec:
7    endpointSelector: {}
8    ingress:
9      - fromEndpoints:
10         - matchLabels:
11             app: frontend
12    egress:
13      - toEndpoints:
14         - matchLabels:
15             io.kubernetes.pod.namespace: kube-system
16             k8s-app: kube-dns
17    toPorts:
18      - ports:
19         - port: "53"
20         protocol: UDP
21    rules:
22      dns:
23        - matchPattern: "*"
24      - toEndpoints:
25        - matchLabels:
  
```



全球开源技术峰会

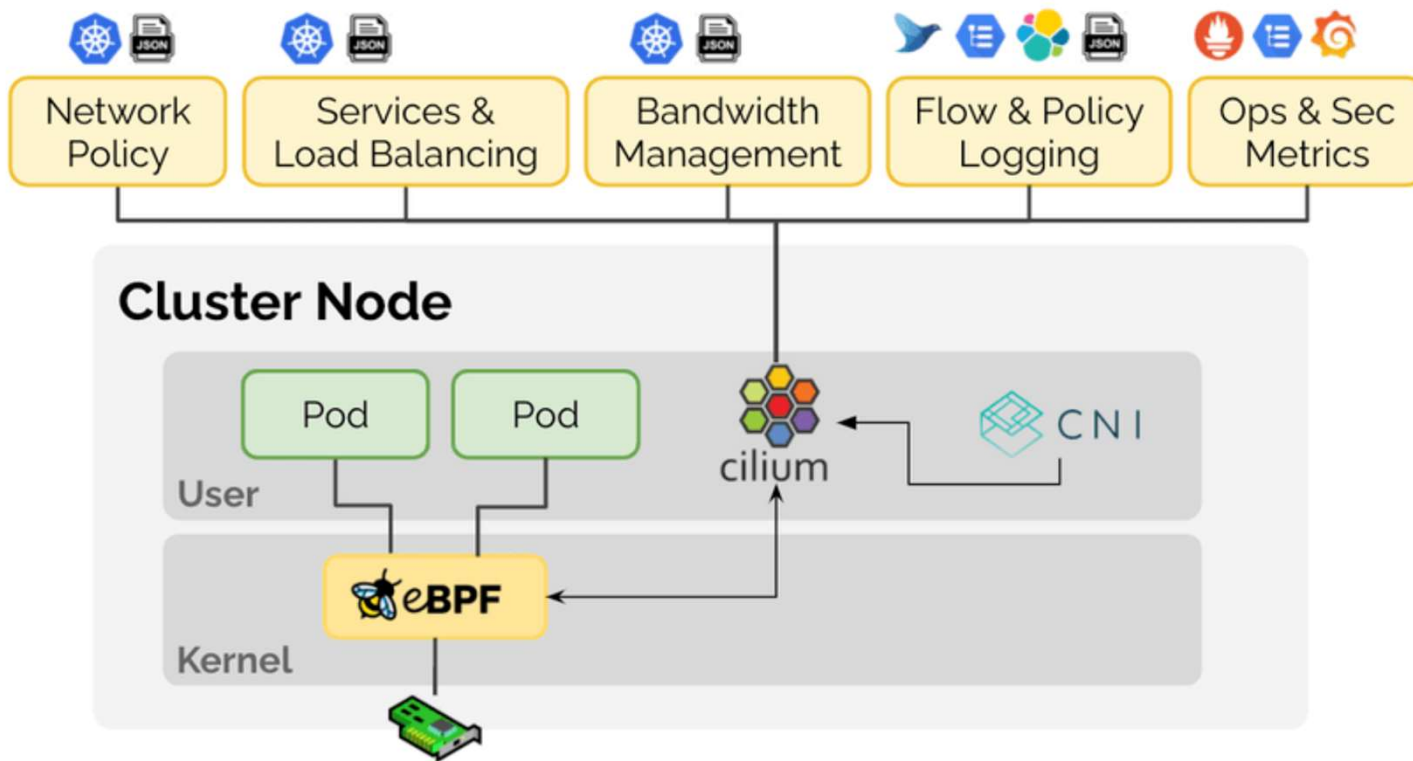
THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

DDoS/网络防火墙

高效阻断和防御

2. BPF 技术在云原生应用

BPF 技术在云原生领域集大成者：开源明星 Cilium

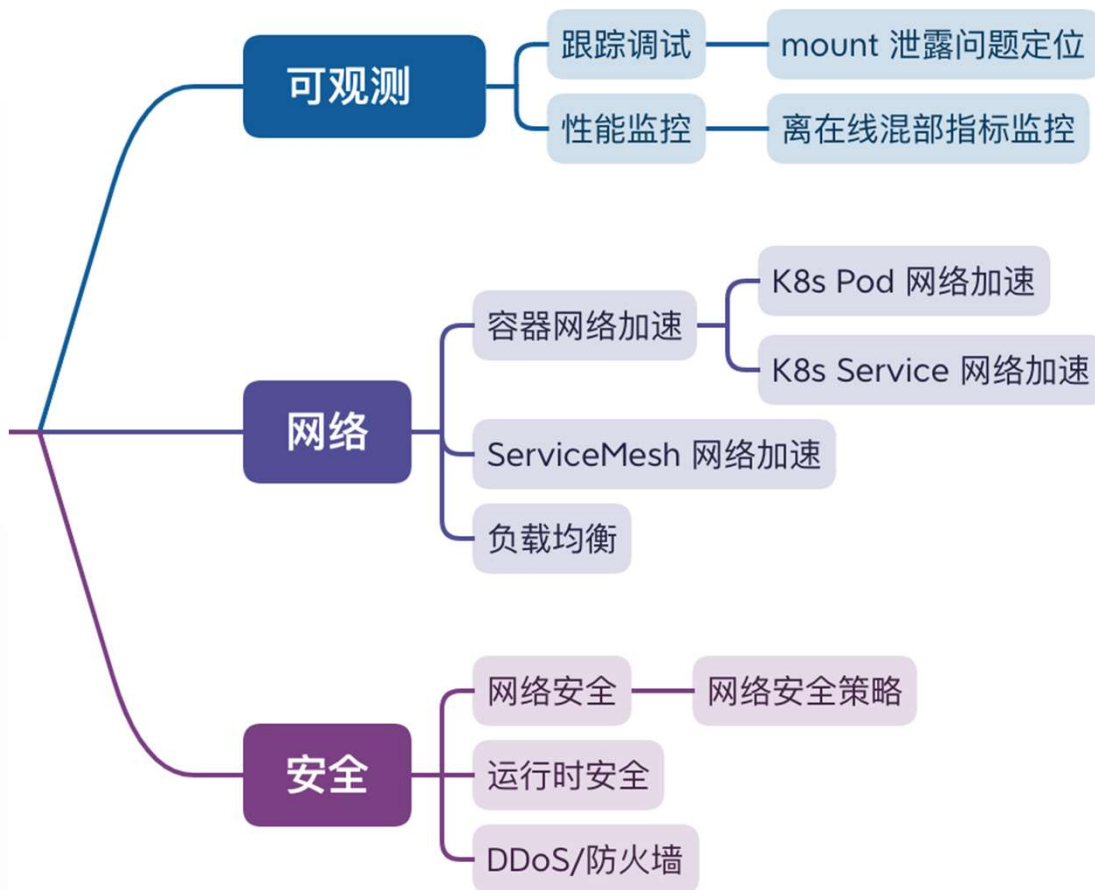


全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

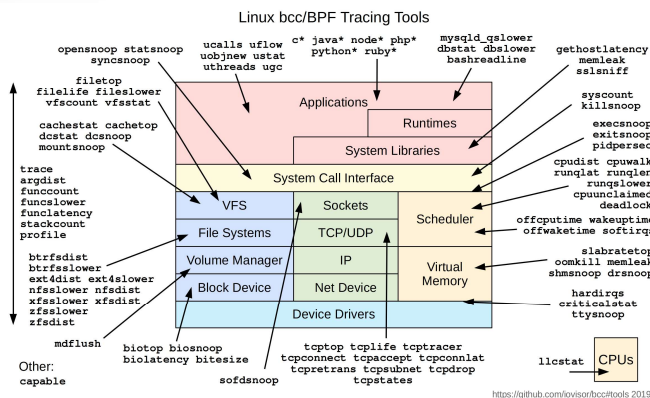


3. BPF 技术在百度云原生实践

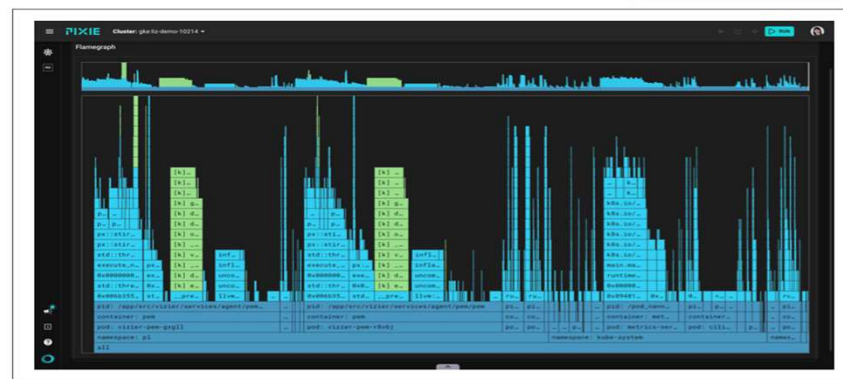
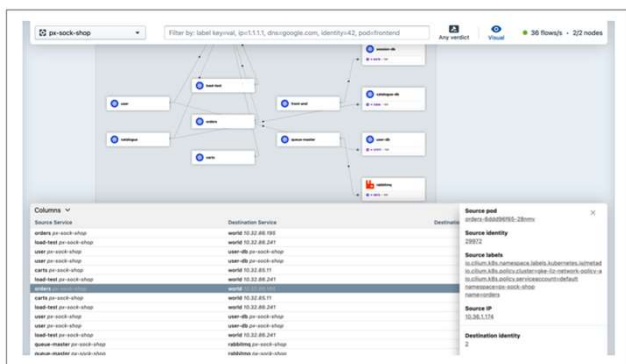
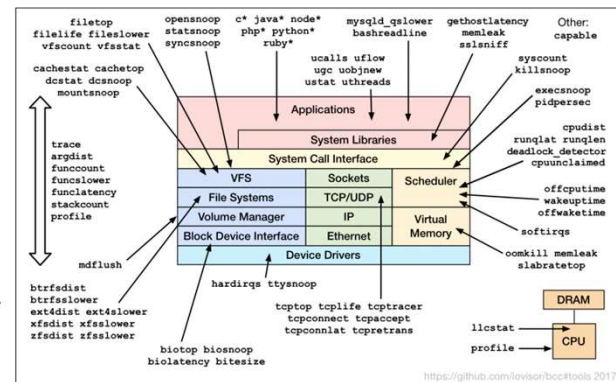


3. BPF 技术在百度云原生实践

• 可观测 – 即开即用开源工具集



1. execsnoop
2. opensnoop
3. ext4slower (...)
4. biolatility
5. biosnoop
6. cachestat
7. tcpconnect
8. tcpaccept
9. tcpretrans
10. gethostlatency
11. runqlat
12. profile



3. BPF 技术在百度云原生实践

可观测 – 容器 mount 泄露

错误定位

部分 Node 上创建 Pod 发现:
mount: mount tmpfs on ... : No space left on device

磁盘空间、inode 和 文件句柄都正常

eBPF 技术工具定位错误来源

syscount-bpfcc -e ENOSPC 定位
_x64_sys_mount 系统调用错误

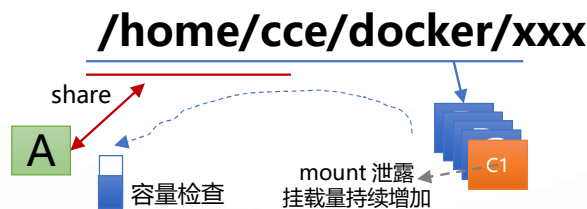
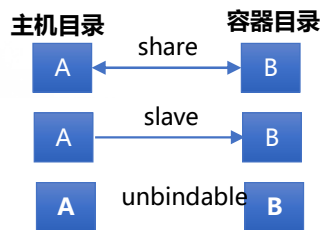
sudo ./funcgraph -m 2

```
x64_sys_mount
-> path_mount()
-> do_new_mount()
-> do_add_mount()
...
-> count_mounts()
```

容器 mount 命名空间 /proc/sys/fs/mount-max
默认 10 万, mount 挂载实例 mount 挂载数量超出
限制

部分容器实例 mountinfo 数量非常大
为什么单个 实例异常, 影响新建容器失败?

寻根问底: mount 传播泄露



彻底分析: 延迟卸载惹得祸

C1 为什么 C1 mountinfo 持续增长?

c1: /noah/:/noah 挂载

docker cp ../noah-ci.tar.gz \$1:/tmp/noah-ci-new.tar.gz

docker cp 导致 c1 容器的 mountinfo 数量持续增加

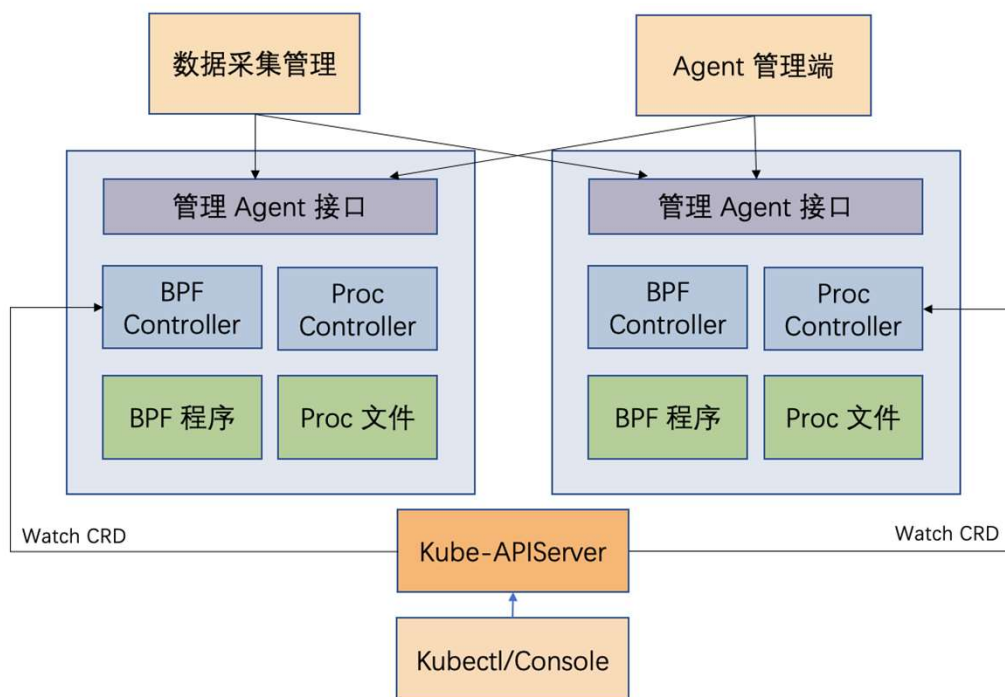
```
$ findmnt
|---/noah          /dev/sda1[/noah]  ext4
|  |---/noah/bin  tmpfs             tmpfs
```

\$ sudo mountsnoop-bpfcc → 递归 mount
mount(..., MS_BIND|MS_REC, NULL)
umount(..., MNT_DETACH) → 延迟卸载

```
static int do_umount(struct mount *mnt, int flags)
{
    //...
    if (flags & MNT_DETACH) {
        if (!list_empty(&mnt->mnt_list))
            umount_tree(mnt, UMOUNT_PROPAGATE); // here <==
        retval = 0;
    }
}
```


3. BPF 技术在百度云原生实践

- 可观测 - 离在线混部内核指标感知，在线服务质量保障



内核指标采集

- CPU: 调度延时/CPI/Cache Miss
- 内存: 分配/回收延时/PSI
- IO: 延迟/吞吐/PSI
- 网络: 连接延时/重传
- ...

功能设计

- 云原生模式管控
- 可基于内核版本控制
- 能够基于容器组 ID 过滤
- 指标以 Pod 名字展示
- 可实现统一管理和配置

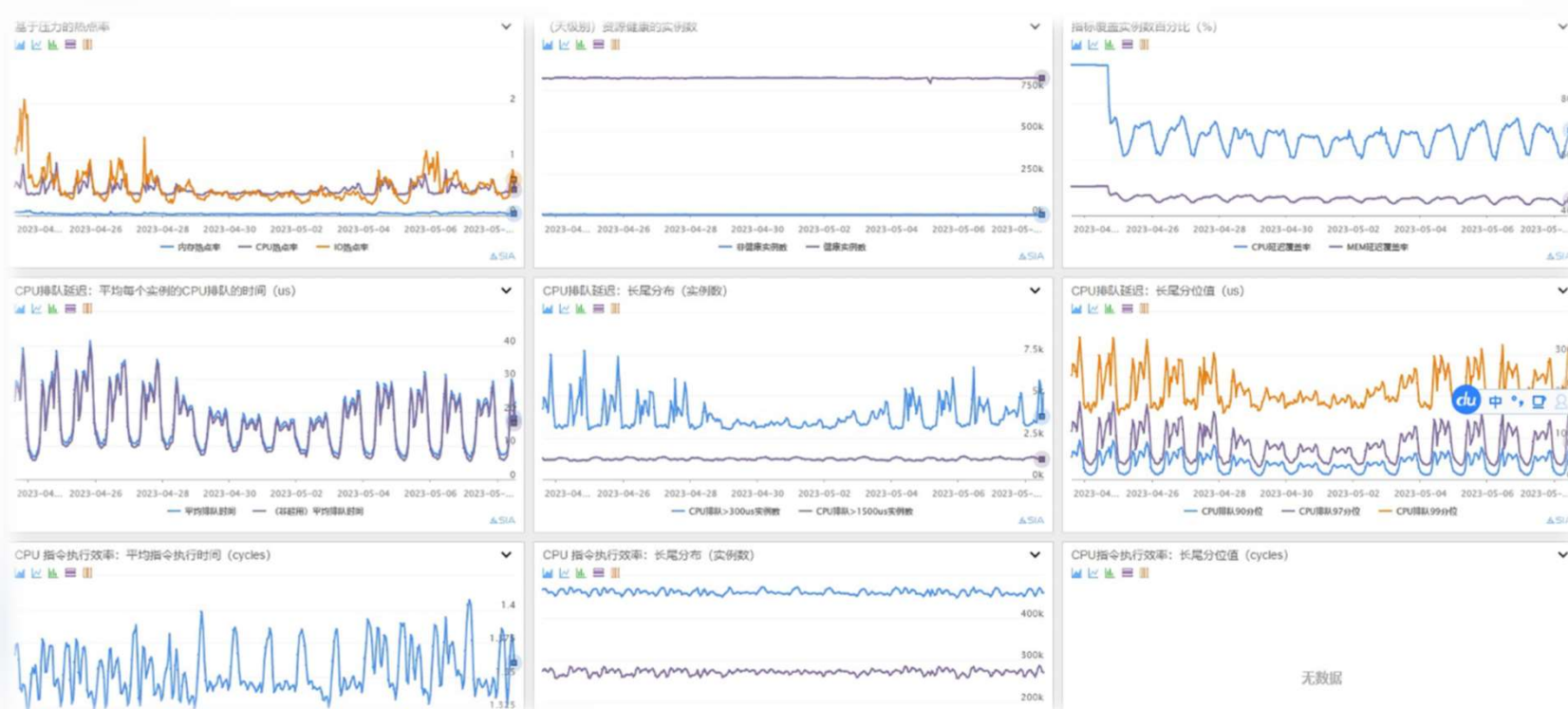
3. BPF 技术在百度云原生实践

- 可观测 - 离在线混部内核指标感知, 在线服务质量保障

```
scripts/bpf/cpu-runq-latency-public.yaml 38
1 apiVersion: halo.baidu.com/v1 39
2 kind: MetricsProgram 40
3 metadata: 41
4   name: runqlat 42
5 spec: 43
6   name: runqlat 44
7   selector: 45
8   workload_type: stable 46
9   metrics: 47
10   histograms: 48
11     - name: run_queue_latency_seconds 49
12       help: Run queue latency histogram 50
13       table: runqlat_dist 51
14       bucket_type: exp2 52
15       bucket_min: 0 53
16       bucket_max: 26 54
17       bucket_multiplier: 0.000001 # microseconds to seconds 55
18       labels: 56
19         - name: pod 57
20           size: 8 58
21         decoders: 59
22         - name: uint 60
23         - name: pod
24         - name: bucket
25           size: 8
26         decoders:
27         - name: uint
28   bpf_programs:
29     - cgroups_table: cgroup_table
30       kernel_version: 5.10.0-1.0.0.13
31       kprobes:
32         ttwu_do_wakeup: trace_ttwu_do_wakeup
33         wake_up_new_task: trace_wake_up_new_task
34         finish_task_switch: trace_run
35       code:
36         text: |
37         #include <linux/sched.h>
38
39         #include <linux/cgroup.h>
40
41         typedef struct pidns_key {
42             u64 id; // cgroup id
43             u64 slot;
44         } pidns_key_t;
45
46         typedef struct cgroup_task_key {
47             u64 id; // cgroup id
48             s32 task_num;
49             u64 ts;
50         } cgroup_task_key_t;
51
52         BPF_HASH(cgroup_table, u64, u32);
53         BPF_HASH(start, u32);
54         BPF_HASH(runqlat, u64, cgroup_task_key_t);
55         BPF_HISTOGRAM(runqlat_dist, pidns_key_t, 1024);
56
57         static u64 get_cid(struct cgroup_subsys_state *css) {
58             return css->cgroup->kn->id;
59         }
60
61         static u64 get_container_id_of_task(struct cgroup_subsys_state *css) {
```

3. BPF 技术在百度云原生实践

- 可观测- 离在线混部内核指标感知，在线服务质量保障

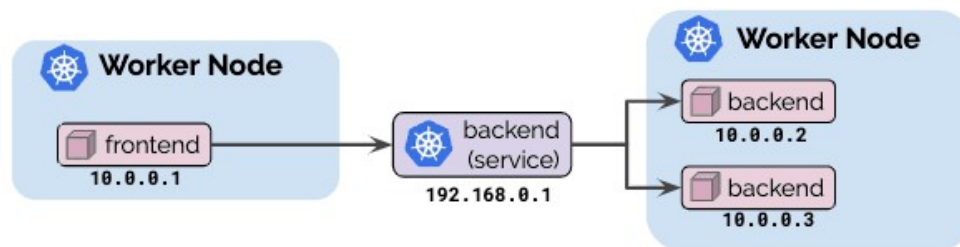


全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

3. BPF 技术在百度云原生实践

- 容器网络加速 – 基于 Cilium K8s Service 加速



Network-based load-balancing

→ connect("192.168.0.1")

Local Socket

```
#1 10.0.0.1→192.168.0.1
#2 10.0.0.1←192.168.0.1
#3 10.0.0.1→192.168.0.1
#n [...]
```

Network DNAT

```
10.0.0.1→10.0.0.2
10.0.0.1←10.0.0.2
10.0.0.1→10.0.0.2
[...]
```

Remote Socket

```
10.0.0.1→10.0.0.2
10.0.0.1←10.0.0.2
10.0.0.1→10.0.0.2
[...]
```

DNAT

Socket-based load-balancing

→ connect("192.168.0.1")

DNAT

Local Socket

```
#1 10.0.0.1→10.0.0.2
#2 10.0.0.1←10.0.0.2
#3 10.0.0.1→10.0.0.2
#n [...]
```

Remote Socket

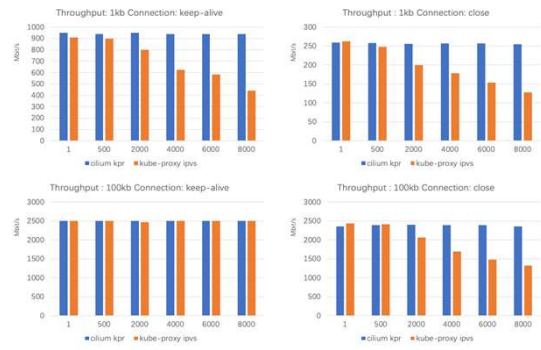
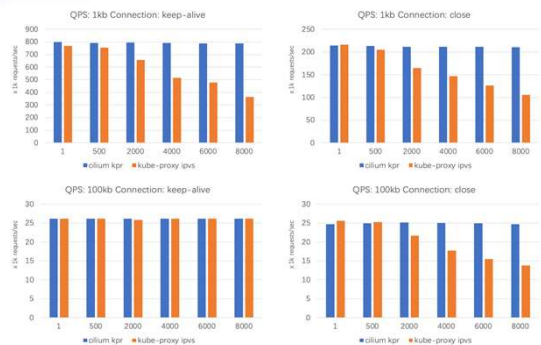
```
10.0.0.1→10.0.0.2
10.0.0.1←10.0.0.2
10.0.0.1→10.0.0.2
[...]
```

- 透明**: 负载均衡对应用程序保持 100% 透明。服务是使用标准的 Kubernetes 服务定义定义的。
- 高效**: 通过在 connect(2) 系统调用中转换地址实现套接字层负载均衡，负载均衡的成本在建立连接时预先支付，在连接之后就不需要额外的转换。性能与应用程序直接访问后端的性能相同。

3. BPF 技术在百度云原生实践

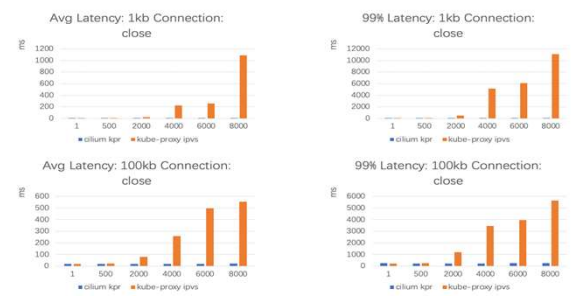
- 容器网络加速 – 基于 Cilium K8s Service 加速

QPS

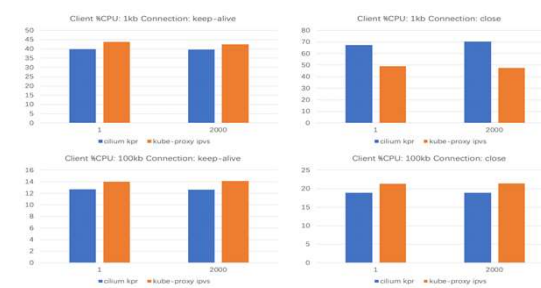


Throughput

Latency



CPU Usage



Cilium Service 加速在网络延迟方面有明显的优势，并且几乎不会随着 Service 数量的增多而上涨，具备保留源 IP 和支撑规模化 NetworkPolicy



公私统一架构，多场景演进

3. BPF 技术在百度云原生实践

- 安全 – 容器网络策略增强，服务标识 & L3/L4 层控制

容器网络

* 容器网络模式: VPC网络 VPC-CNI VPC-Hybrid [如何选择Kubernetes集群的容器网络模式](#)

开启eBPF增强
开启后原有Kube-proxy将被替代，操作系统仅支持公共镜像Ubuntu 20.04.

支持NetworkPolicy
勾选后集群将默认支持基于策略的网络控制，[查看详情](#)

* 容器网络配置: [请务必在选择网络前，合理有效地规划集群的网络，避免出现网络冲突，否则可能存在风险。CCE集群网络说明及规划](#)

容器网段: 10 . 2 . 0 . 0 / 16

节点Pod数: 256
当前容器网络配置下，集群最多允许部署256个工作节点

ClusterIP网段: 172 . 16 . 0 . 0 / 16
当前ClusterIP网段配置下，集群最多允许创建65536个service

[高级设置](#)

- 禁止namespace=staging中Pod被访问。

```
apiVersion: networking.k8s.io/v1
kind: NetworkPolicy
metadata:
  name: default-deny
  namespace: staging
spec:
  podSelector: {}
  policyTypes:
  - Ingress
```

The interface shows a network policy configuration for Cilium. The top part is a diagram illustrating traffic flow between different namespaces and clusters. Below the diagram is a code editor showing the following policy configuration:

```
Kubernetes Network Policy | Cilium Network Policy
Policy Rating: [Progress Bar] [Download] [Share]

1 apiVersion: cilium.io/v2
2 kind: CiliumNetworkPolicy
3 metadata:
4   name: store-policy
5   namespace: my-store
6 spec:
7   endpointSelector: {}
8   ingress:
9     - fromEndpoints:
10       - matchLabels:
11         app: frontend
12   egress:
13     - toEndpoints:
14       - matchLabels:
15         io.kubernetes.pod.namespace: kube-system
16         k8s-app: kube-dns
17     toPorts:
18       - ports:
19         - port: "53"
20         protocol: UDP
21     rules:
22       dns:
23         - matchPattern: "*"
24     - toEndpoints:
25       - matchLabels:
```

3. BPF 技术在百度云原生实践

• 实践思考



- BPF 技术是基础能力，需要针对解决问题的场景有足够的认知和分析（纵向与横向），如内核中进程/内存/文件/设备/网络等子系统
- BPF 并非原有工具和技术的替换，相互补充，提供低门槛和可编程的灵活性（最后一公里和细粒度的监测）
- 基于问题解决使用 BPF 技术，避免用 BPF 技术去找问题



- 尽量选取能使用的最高内核版本，内核版本越高 BPF 特性支持越全
- 尽可能基于 BTF & CO-RE (Compile Once-Run Everywhere)，提升可移植性
- 优先使用内核稳定跟踪机制，提升多内核版本可移植性，内核版本尽量收敛
- 关注 BPF 附着事件触发频繁，避免出现不可控的性能开销

GOTC

THANKS

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE