

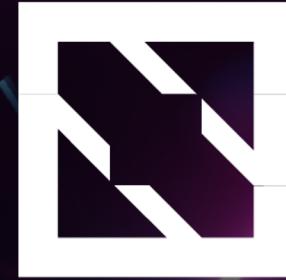


KubeCon

THE LINUX FOUNDATION



China 2024



CloudNativeCon





KubeCon



China 2024

KubeSkoop - Deal with the Complexity of Network Issues and Monitoring with eBPF

Yutong Li, Alibaba Cloud
Bingshen Wang, Alibaba Cloud

Agenda

- The complexity of Kubernetes Networking
- Introduction to KubeSkoop
- KubeSkoop Network Diagnosis based on eBPF



KubeCon



CloudNativeCon



China 2024





KubeCon



CloudNativeCon



China 2024



Alibaba Cloud Kubernetes Service(ACK)

30+

Regions

tens of
thousands

Clusters

10k+

Nodes
in single cluster



KubeCon



CloudNativeCon



China 2024

The Complexity of Kubernetes Networking



in Concepts



in Implementations

Concepts in Kubernetes



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024

The network concepts of Kubernetes lead to the complexity of networking configuration:

Ingress/Service/NetworkPolicy

- LabelSelector selects unexpected pods.
- Overlapping of multiple NetworkPolicy rules.
- NATed service ports do not match the real pod ports.

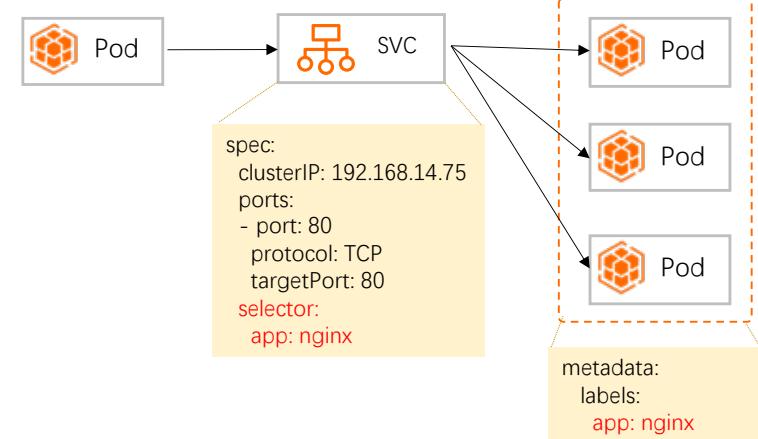
ServiceMesh

- More complicated Layer 7 network strategies

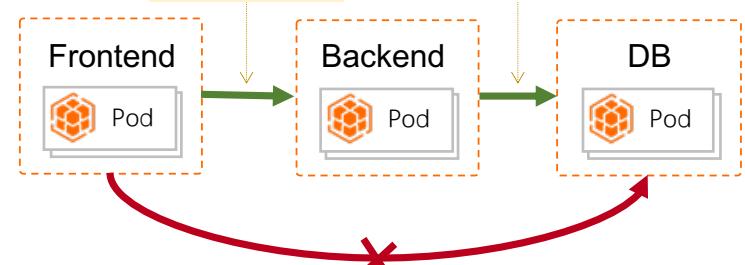
Third-party Networking/Ingress Plugins

- Provide custom networking extensions.

Service Discovery



NetworkPolicy



Implementation of Container Networking



China 2024

Complexity of data-plane:

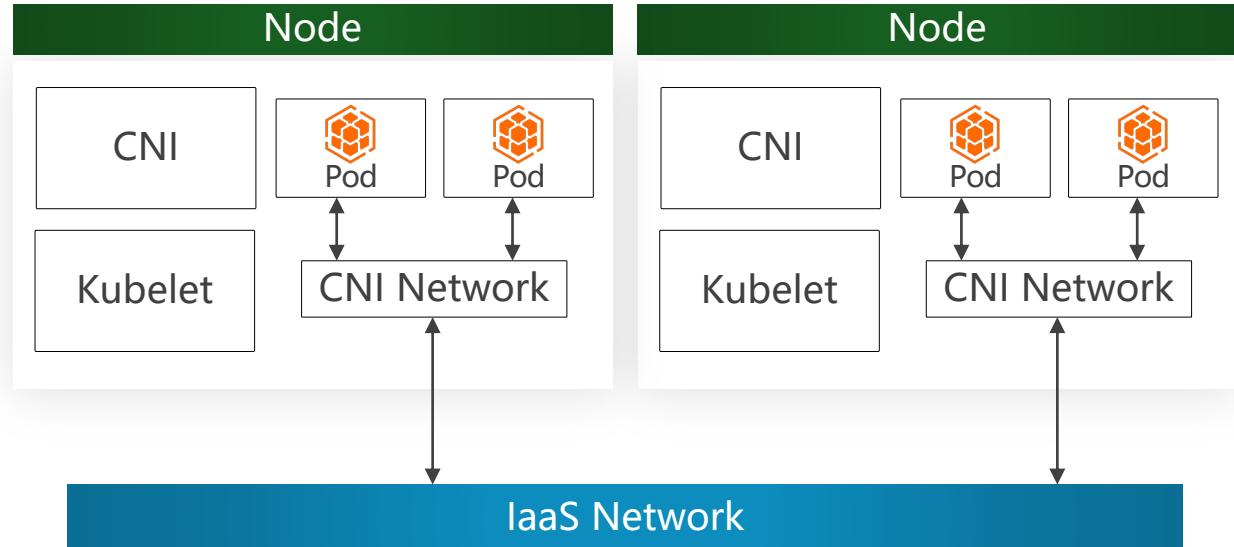
- ServiceMesh/KubeProxy/CNI
- CNI Implementations (Overlay/Underlay...)

Complexity of the network stack

- Long data path, including NIC drivers/netfilter/route/bridge etc.
- Complicated networking configuration

Complexity of underlay network:

- Different configurations per cloud provider
- Security groups, route tables, etc.



Traditional Network Troubleshooting



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

AI_dev
Open Source Dev & ML Summit

China 2024

1 Long Issue Diagnosis Process:

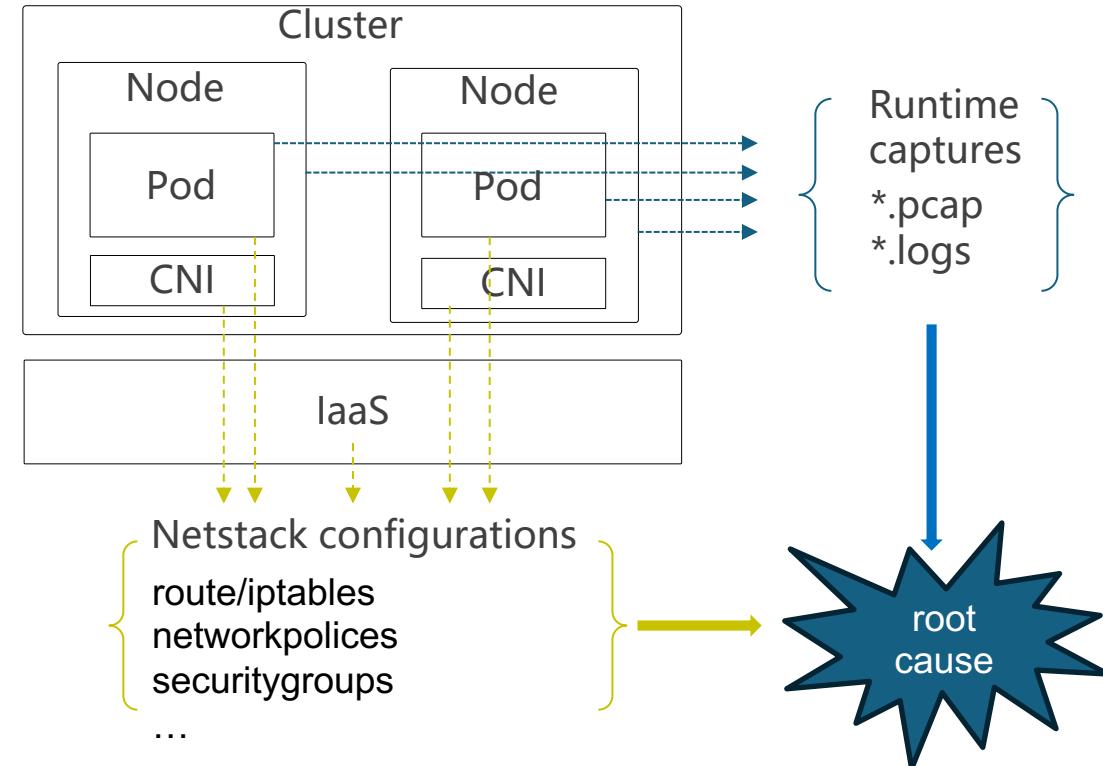
Capture packets -> Compare and analyze ->
Check configurations at packet loss points

2 Long Diagnosis Time:

Stress Testing to Reproduce the problem ->
Analyze the problem with a large number of information

3 High Experience Requirements:

Mastery of the Linux protocol stack, CNI implementation,
and IaaS layer network configuration.





KubeCon



CloudNativeCon



China 2024

Introduction to KubeSkoop

Introduction to KubeSkoop



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



AI_dev
Open Source DevOps & ML Summit

China 2024

Network monitoring & diagnosis suite for Kubernetes



**Connectivity
Diagnosis**



**Anomaly
Tracing**



Flow Logs



**Latency
Detection**



**Packet
Capturing**

<https://github.com/alibaba/kubeskoop>

Connectivity Diagnosis



China 2024

Console / Diagnosis / Connectivity Diagnosis

Connectivity Diagnosis

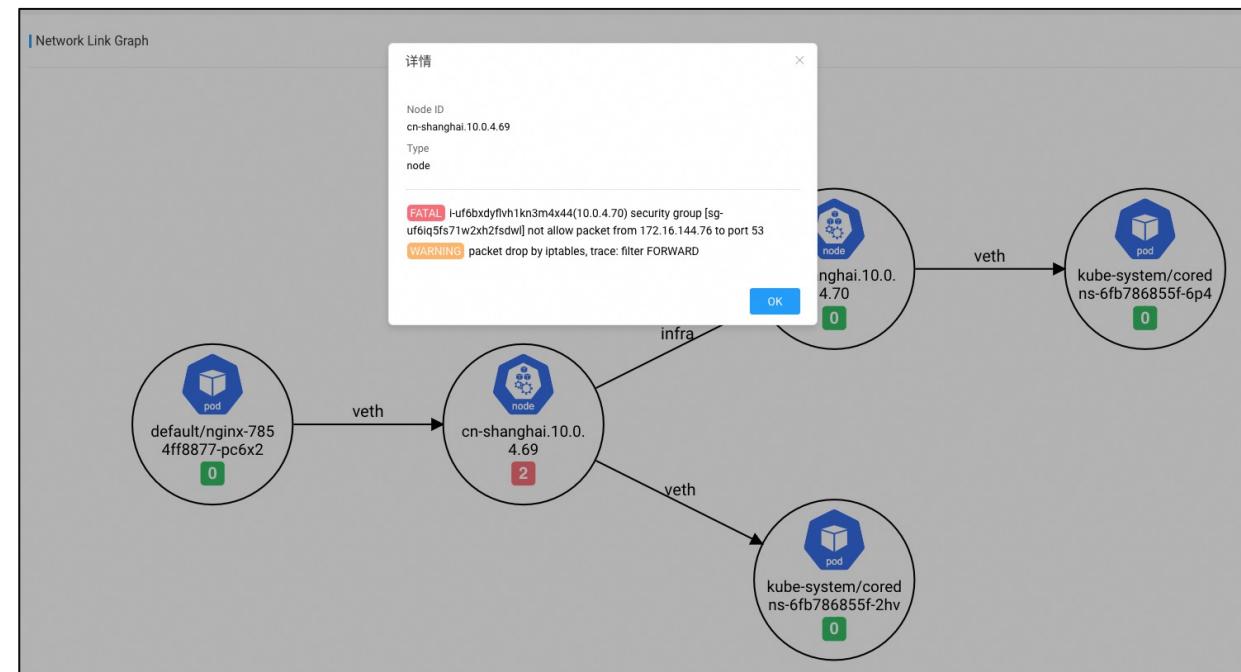
Diagnose

* Source Address: 10.92.0.75 * Destination Address: 172.16.77.11 * Port: 8080 * Protocol: TCP Diagnose

History

ID	Time	Source Address	Destination Address	Port	Protocol	Status	Actions
1	2024-01-22T06:08:34Z	10.92.0.133	172.16.0.10	53	udp	success	<button>Result</button>

- Specify the source and destination address of the problem. Wait for the diagnostic output.
- Constructs network links and analyzes the problem automatically.
- Including analysis of Kubernetes Service, NetworkPolicy, etc.
- Analysis of kernel network stack and underlying IaaS configuration.
- No need the troubleshooting experience of CNI implementations and complicated network problems.

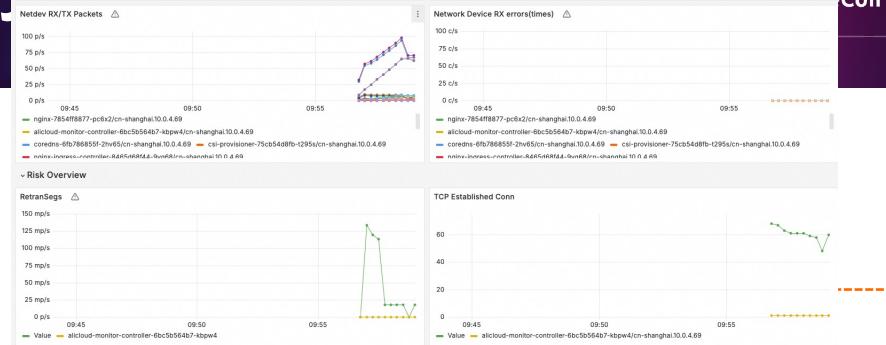


Monitoring and Anomaly Tracing



China 2024

2024-04-10T01:54:53.075Z	[protocol: TCP]	[src: 172.16.144.73]	[src_type: unknown]	[dst: 192.168.182.34]
PacketLoss				
2024-04-10T01:54:53.075Z	[protocol: TCP]	[src: 172.16.144.73]	[src_type: unknown]	[dst: 192.168.182.34]
TCPRetrans				



Container Network Anomaly Monitoring

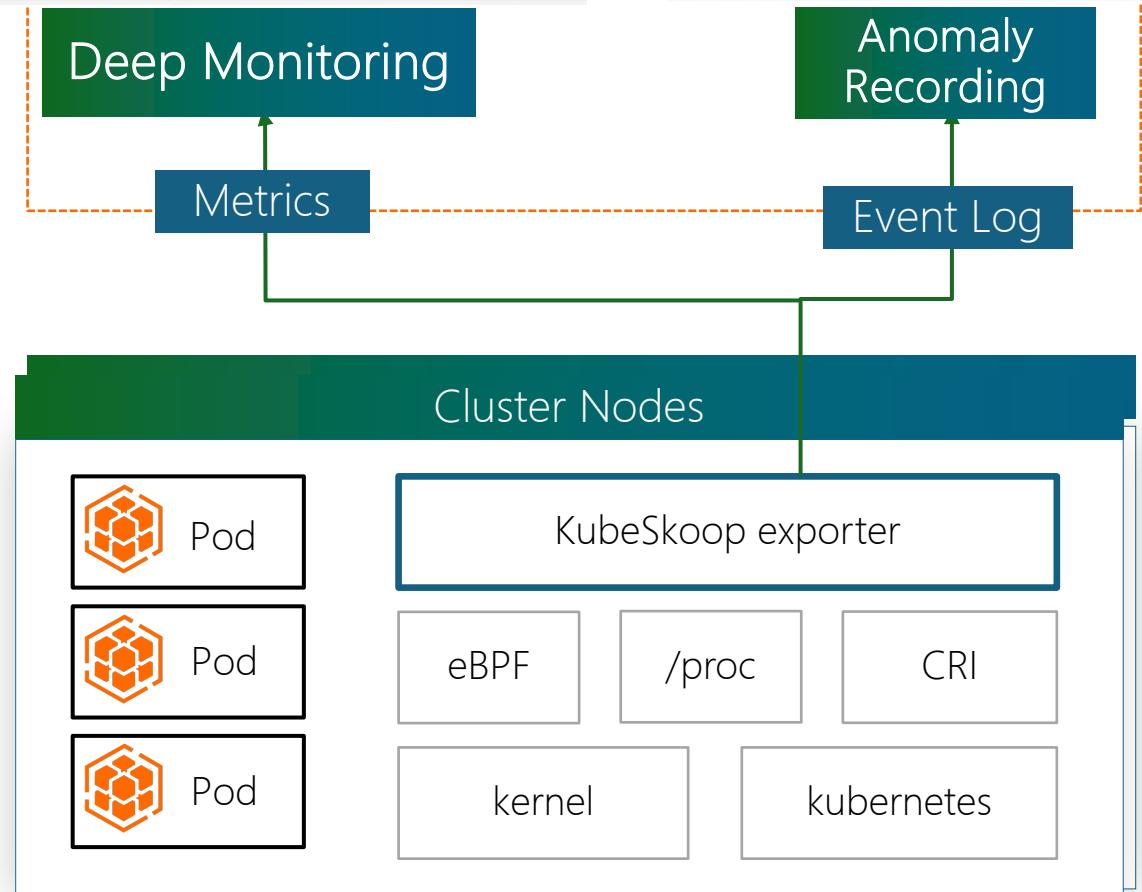
- Cloud-native deployment, work with observability systems like Prometheus
- Pod-level network monitoring
- Lightweight, low-overhead kernel anomaly tracing

Covering Multiple Scenarios

- Occasional packet loss and retransmission
- Flow-level anomaly identification
- Network latency analysis

Exposing Multiple Types of Anomalies

- Metrics for network traffic and relationships
- Event logging of network anomalies



Access Flow



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

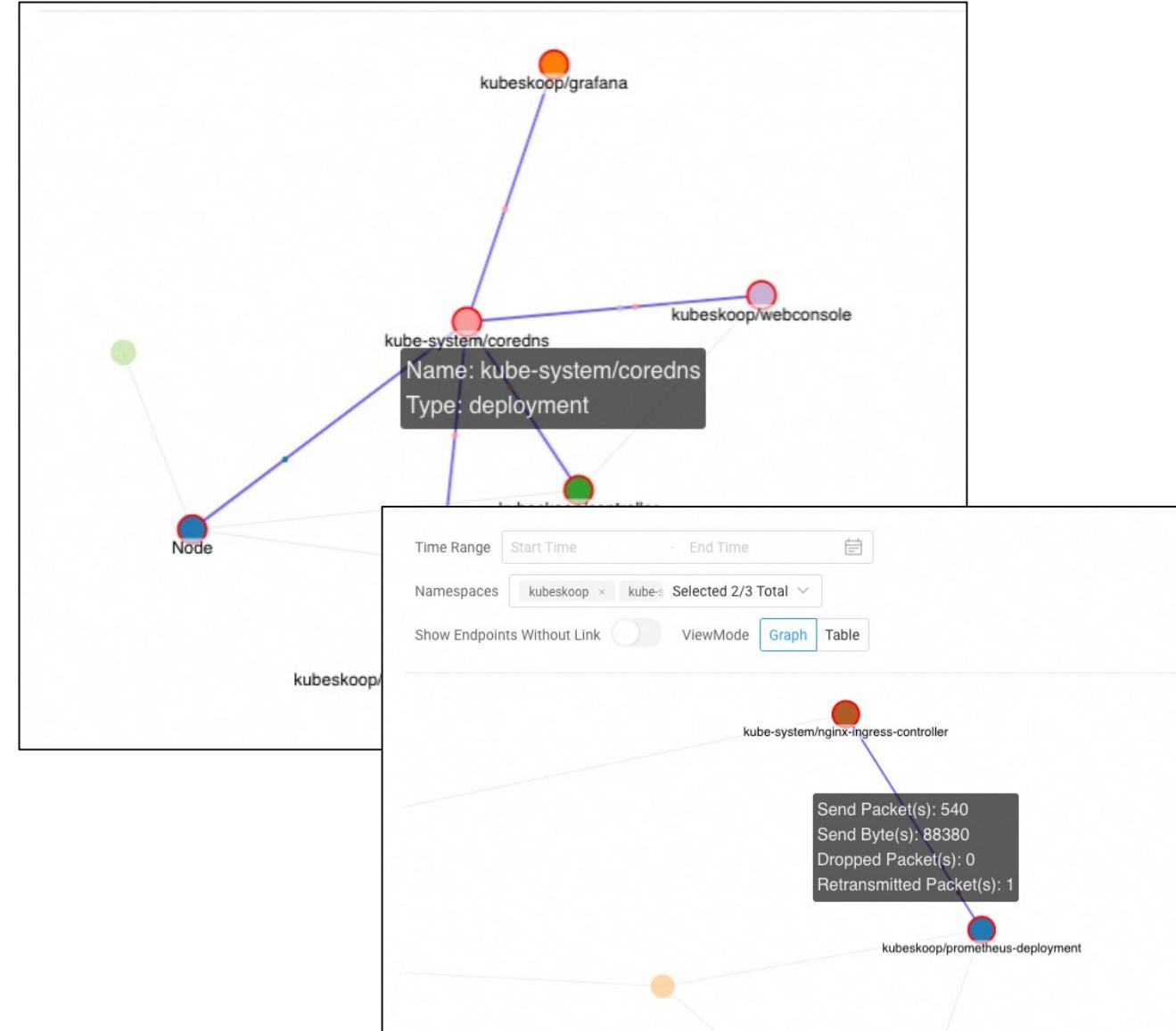
China 2024

Historical access relationships recording and backtracing

- Recording connections, packets and throughput
- Analyze network performance and bottleneck of applications
- Troubleshooting historical resource usage of shared services (e.g. CoreDNS)

Low-overhead access relationship capturing

- Access flow capture based on eBPF



Latency Detection(PingMesh)



KubeCon



CloudNativeCon

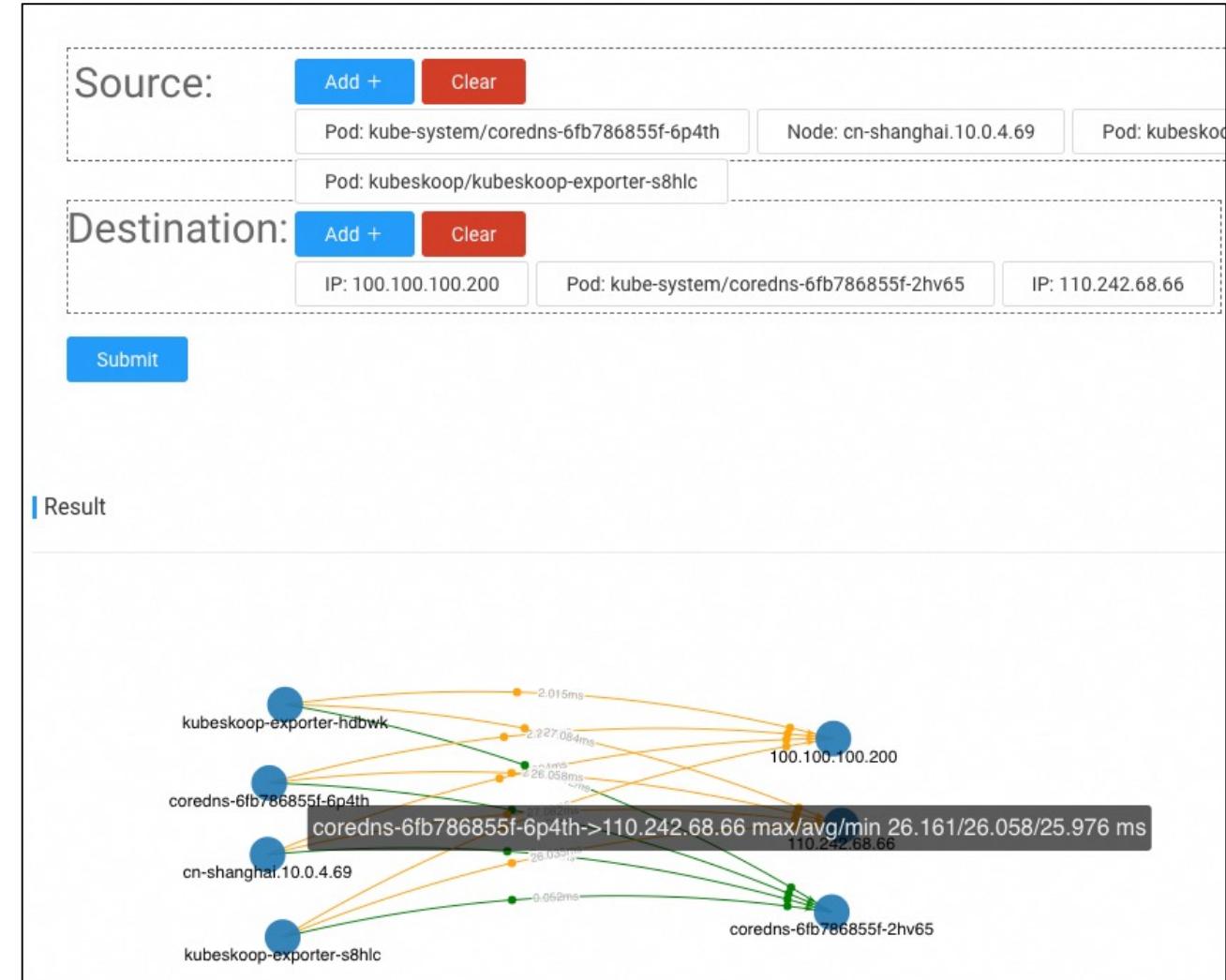


China 2024



Latency Detection in Cluster

- Generate latency report in one click
- Troubleshooting performance bottlenecks
- Optimize application deployment topology



Packet Capturing



China 2024

Capture packets on multiple nodes simultaneously

- One-click packet capture for all links you need
- Locate packet loss point for unknown reasons

The screenshot shows a "Packet Capturing" application window. On the left, there's a "Capture" tab with a "Targets" section containing a list of selected targets: "Pod: kubeskoop/controller-74594c79d4-x9zr" and "Node: cn-shanghai.10.0.4.69". Below this is a "Filter" field set to "udp" and a "Duration" field set to "00:30". On the right, a modal dialog titled "Add Target" is open. It has tabs for "Type" (set to "Pod") and "Select Target By" (set to "Label Selector"). Under "Namespace", it shows "kube-system". Under "LabelSelector", it shows "k8s-app = kube-dns". A checkbox "Also capture node packets" is checked. At the bottom of the modal is a blue "OK" button. Below the modal, in the "History" section, there's a table with one row:

ID	CaptureObjects	Result
13	Pod: controller-74594c79d4-x9zr, Node: cn-shanghai.10.0.4.70, Pod: coredns-6fb786855f-6p4th, Pod: coredns-6fb786855f-2hv65, Node: cn-shanghai.10.0.4.69	Download



KubeCon



CloudNativeCon



China 2024

Kubernetes Network Monitoring Based on eBPF

Why choose eBPF



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



Open Source Dev & ML Summit

China 2024



Kubernetes uses a lot of network components:

- TCP
- Bridge
- Netfilter
- TC
- ...



The information exposed by the kernel is limited:

- Lack of netns information
- Missing records for critical links
- Absence of context information
- ...



Convenience of eBPF program deployment and distribution:

- Good compatibility
- High security
- Low overhead for kernel-level filtering and statistics
- Creating a reproducible troubleshooting experience.

Key Characteristics of Kubernetes Network Issues



KubeCon



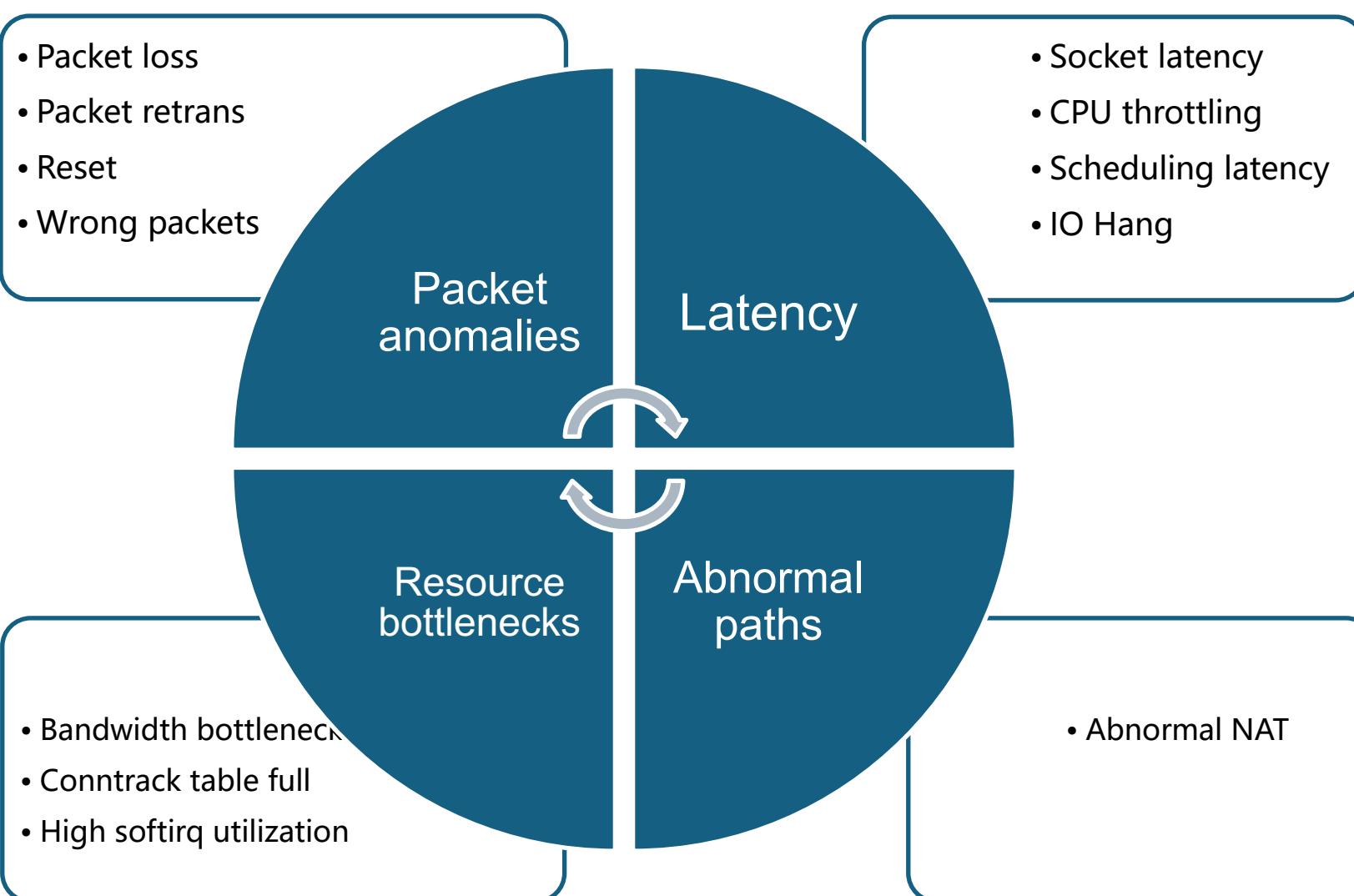
CloudNativeCon



China 2024



Open Source Dev & ML Summit



eBPF Monitoring for Packet Anomalies



KubeCon



CloudNativeCon

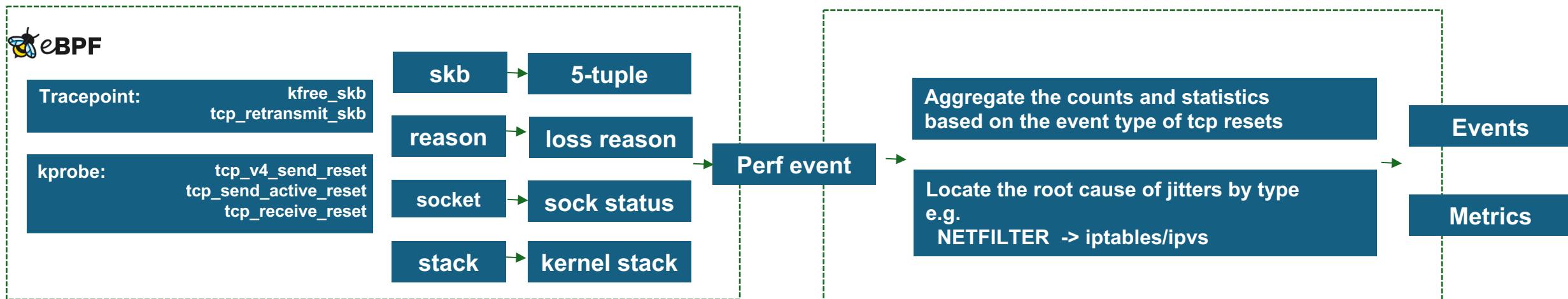
THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024

Packet anomalies causing occasional business jitter

- Packet loss & retransmissions
- TCP resets

Using eBPF allows for tracking the call stack,
thereby identifying the root cause of packet loss.



eBPF Monitoring for Network Latency



CloudNativeCon

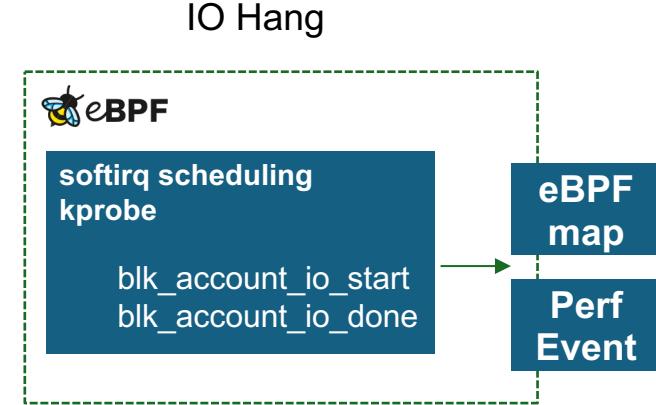
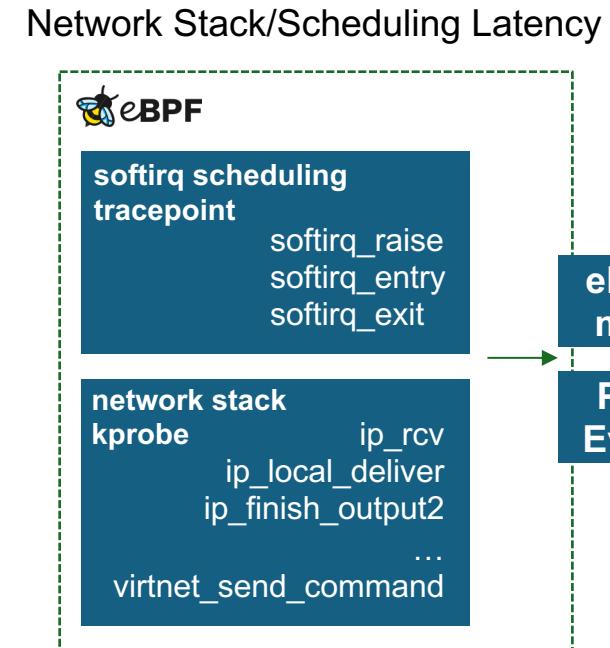
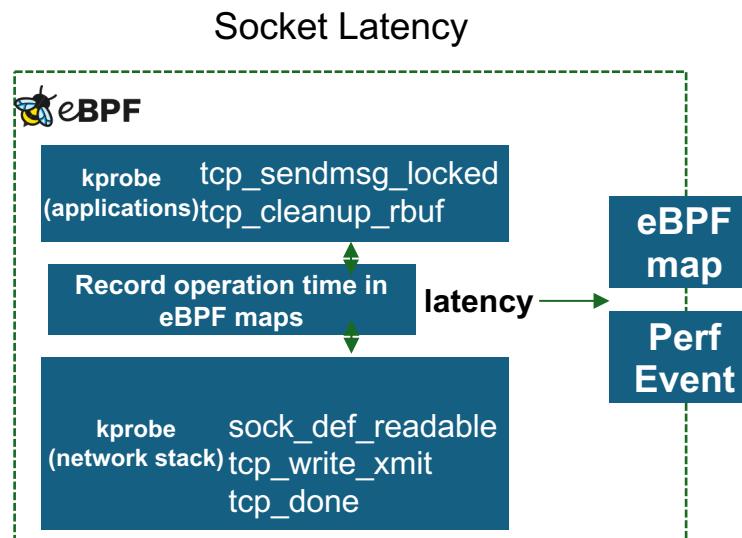
THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024

Network latency issues:

- Socket latency
- Network stack/scheduling latency
- IO hang

eBPF uses skb and pid as keys for tracking, collecting data at various points to identify the root cause of latency.

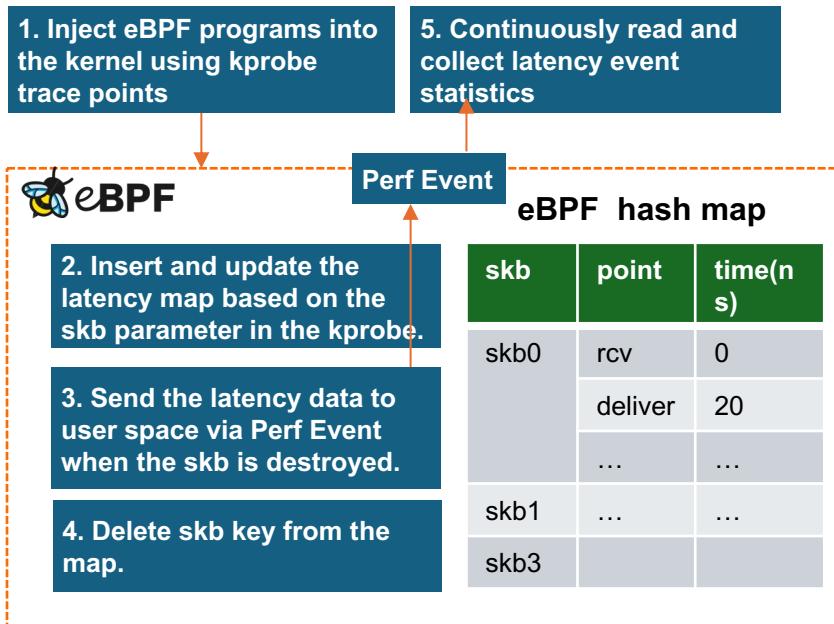


eBPF Monitoring for Network Latency



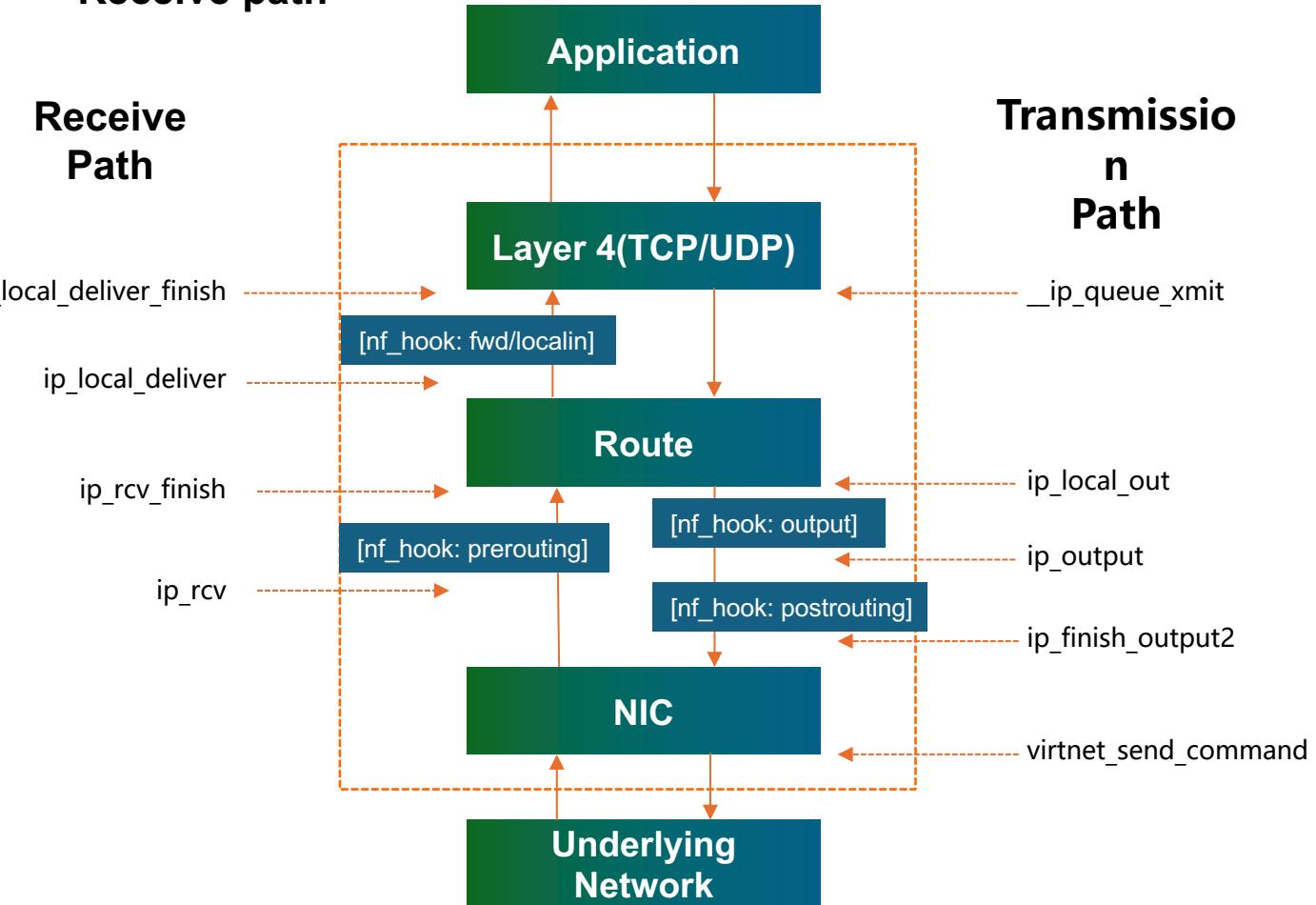
China 2024

eBPF Network Stack Latency Detection



Latency Detection for The Whole Path through the Kernel

- Transmission path
- Receive path



eBPF Metrics Collecting and Processing in KubeSkoop



Metrics

Latency

Anomaly count

Events

Event log

Call stack

Flow 5-tuple Statistics

Collect Tracepoint kprobe & uprobes

Tracepoint kprobe

TC
sockops

Sync eBPF Map

Perf Event

eBPF Map

Store Prometheus Time Series

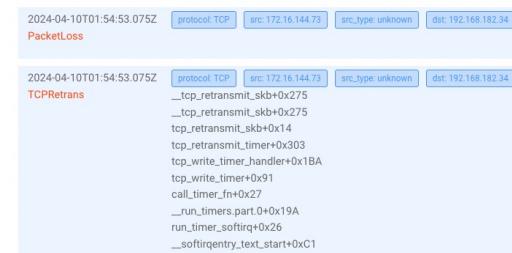
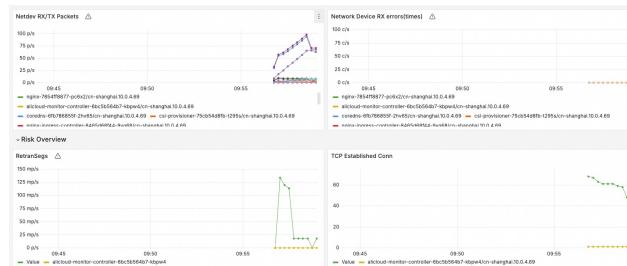
Loki Log Stores

Prometheus Time Series

Visualize Grafana Dashboard

Event Log Backtrack

Flow Graph



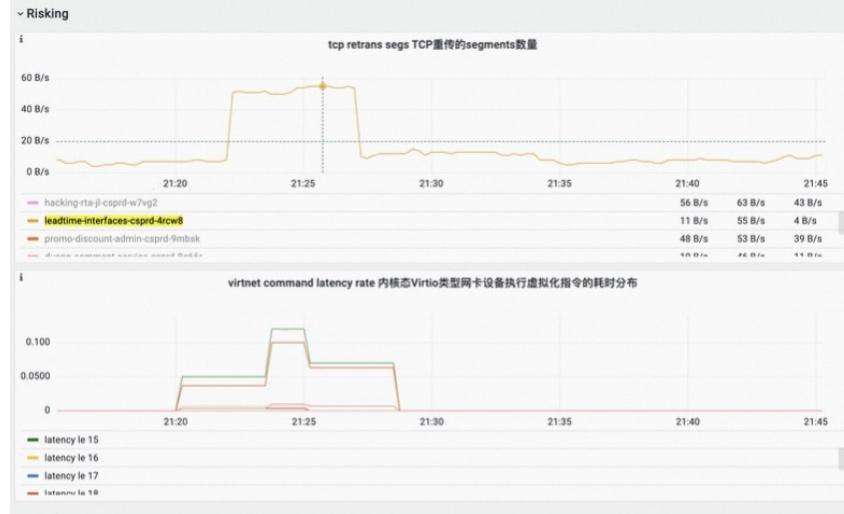
KubeSkoop eBPF Monitoring Case Study



China 2024

Troubleshooting Root Cause of Business Jitter: Occasional Timeouts When Accessing Ingress NGINX Controller

Backtracking Metrics



Event Logging Events

```
INFO March 9 21:21:05 VIRTCMDEXCUTE nginx-ingress-controller-5c7fb5594-jwhnw protocol=TCP  
saddr=100.121.88.193 sport=27057 daddr=10.33.0.11  
dport=443 stacktrace: virtnet_send_command+0x1  
virtnet_set_rx_mode+0x251 __dev_change_flags+0x9c  
dev_change_flags+0x21 do_setlink+0x257  
__ rtnl_newlink+0x600 rtnl_newlink+0x44  
rtnetlink_rcv_msg+0x119 netlink_rcv_skb+0x4e  
netlink_unicast+0x1d7 netlink_sendmsg+0x240  
sock_sendmsg+0x5f ____sys_sendmsg+0x232  
____sys_sendmsg+0x75 ____sys_sendmsg+0x49  
do_syscall_64+0x30  
entry_SYSCALL_64_after_hwframe+0x61
```

Business Anomaly Time Point: Around 21:25

Backtracking with KubeSkoop monitoring:
Increased retransmissions from business Pods.
Sudden spike in Virtnet command latency

Event Logins reveals the call stack and root cause:
Underlying virtualization execution timeouts leads to soft interrupts being stuck, caused the packet retransmission.

Challenges of eBPF Program Management



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



AI_dev
Open Source Dev & ML Summit

China 2024

Deployment and Distribution

- Dependency on kernel header file
- Large clang/llvm compilation environment (300MB+)

eBPF Probes Management

- Dynamic insertion and removal of dozens of eBPF programs
- Dynamic features configuration for each program

Effective Collection of Large Volumes of Data

- Different data labels across various programs
- Events and monitoring metrics need to be output to multiple storages

KubeSkoop eBPF Program Management



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

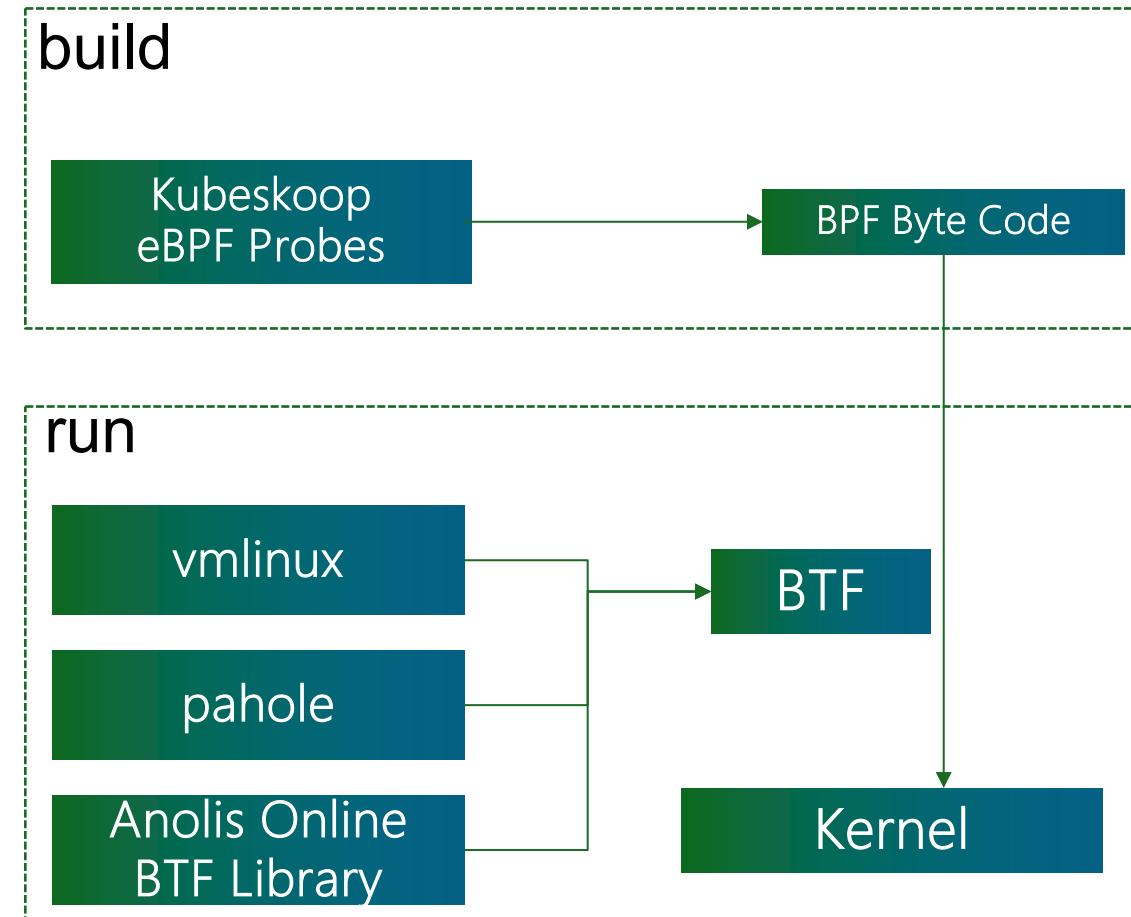


AI dev
Open Source Dev & ML Summit

China 2024

Leveraging eBPF CO-RE

- No need for clang/llvm dynamic compilation
- Reduced OS header file and environment dependencies
- Use btfhack to obtain BTF files from various source
- Distribution size reduced from 300mb to 33mb



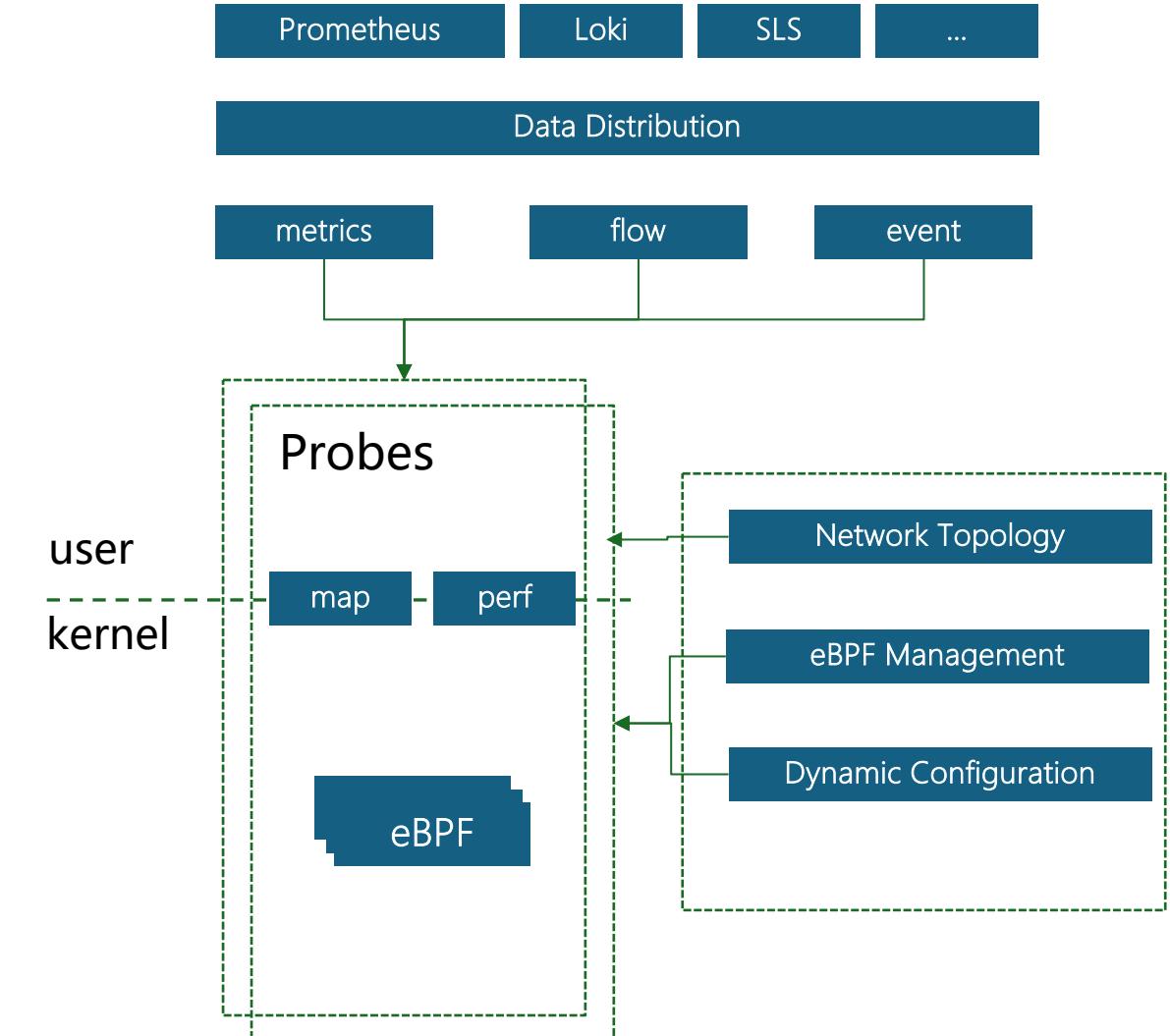
KubeSkoop eBPF Program Management



China 2024

Prober Extensibility Design

- Abstract of Metrics, Event, Flow interface
- Generation of shared network topology information from CRI & netns
- Dynamic plugging and configuration of eBPF programs based on .rodata
- Unified data collection and distribution
- Simplified development efforts for adding new monitoring capabilities



Future Plan



KubeCon



CloudNativeCon



China 2024



- Support for more cloud providers and network plugins
- Enable display of RTT, window size, and other information in flow diagrams to identify network bottlenecks within the cluster
- Support multiple protocol types for latency detection, allowing simultaneous packet capture and analysis of latency paths
- Allow the extension of probe using bpftrace scripting language to reduce development barriers
- Provide the KubeSkoop Analysis tool to intelligently analyze KubeSkoop metrics and events, lowering the threshold for understanding issues
- Application layer awareness for perception and processing of Layer 7 protocols (e.g., HTTP, Redis, etc.)