



**KubeCon**



**CloudNativeCon**

THE LINUX FOUNDATION



**AI\_dev**  
Open Source GenAI & ML Summit

---

**China 2024**

---



KubeCon



CloudNativeCon



China 2024

# Revolutionizing Service Mesh with Kernel-Native sidecarless Architecture

# About Me



China 2024



**Xin Liu**  
Softwar Engineer, Huawei



## KMESH

- Kmesh Maintainer
- Linux Kernel Contributors

# Agenda



China 2024

- Kmesh Overview
- How to implement L4/L7 traffic offload
- Restarting without stop the service
- High-performance non-intrusive observability

# Agenda



China 2024

- Kmesh Overview
- How to implement L4/L7 traffic offload
- Restarting without stop the service
- High-performance non-intrusive observability

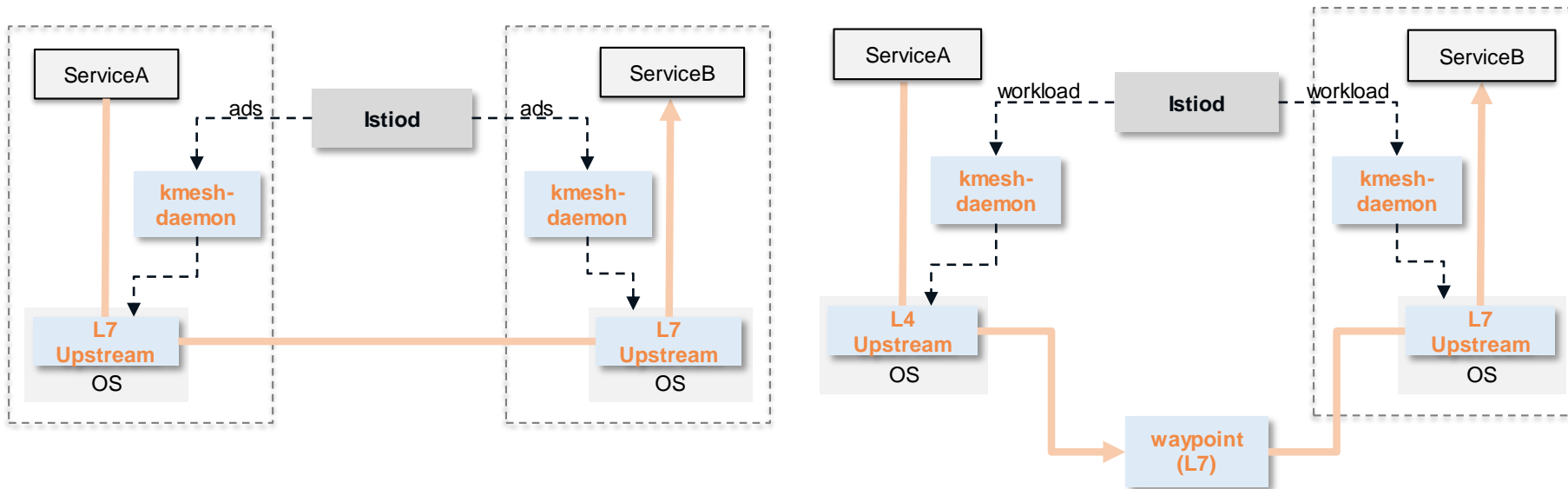


# Kmesh Overview



China 2024

Kmesh is a high-performance and low overhead service mesh data plane based on eBPF and programmable kernel. Kmesh brings traffic management, security and monitoring to service communication without needing application code changes. It is natively sidecarless, zero intrusion and without adding any resource cost to application container.



## Kernel-Native:

L4~L7 ultimate performance

## Dual Engine

L4/L7 Slicing the layers Flexible management

- High Performance
- Low overhead
- High availability
- Security isolation
- Flexible management mode
- Seamless compatibility

# Agenda



China 2024

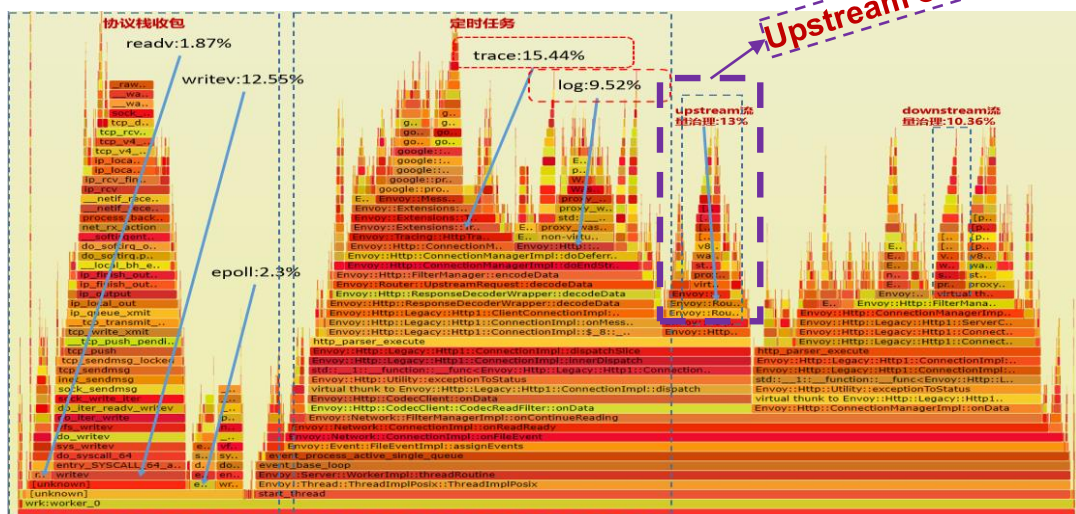
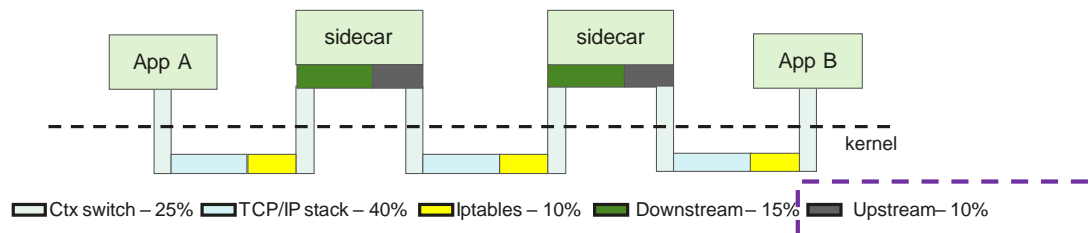
- Kmesh Overview
- How to implement L4/L7 Upstream
- Restarting without stop the service
- High-performance non-intrusive observability

# Istio sidecar performance profiling

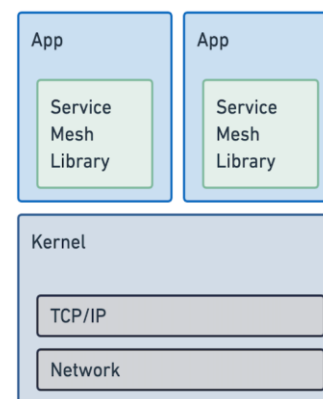


China 2024

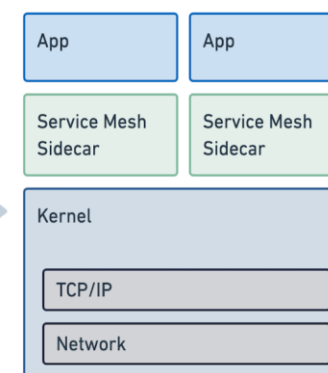
## Sidecar Profiling



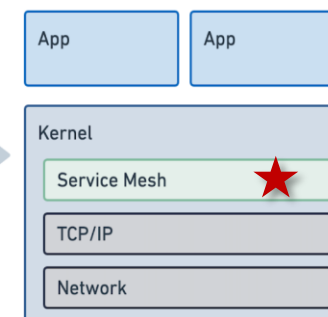
## Shared Library Model



## Sidecar Model



## Kernel Model



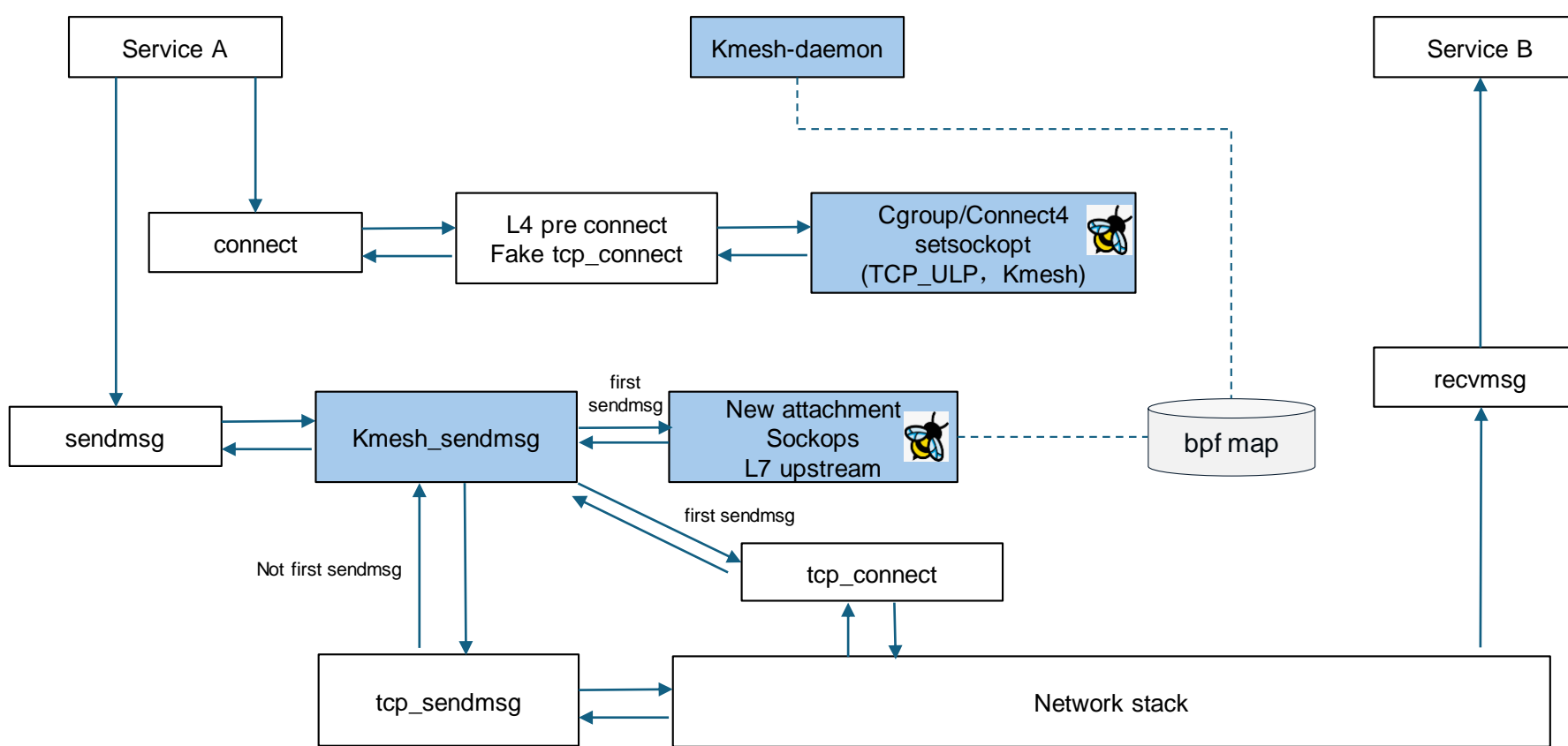
The network time consumption shows that : **The sidecar architecture introduces a large amount of latency overhead.** Upstream is only 10% of grid overhead. Most overheads are caused by data copy, two extra link setup communications, and context switch scheduling.



# Kernel-Native L7 Upstream



China 2024



Kmesh-Native L7 upstream  
Contains multiple components :

- **Kmesh-daemon**  
Configuration data delivery on the management and control plane
- **Cgroup/Connect4(ebpf):**  
replace sendmsg with Kmesh\_sendmsg by ULP
- **Kmesh sendmsg ko:**  
The first sendmsg message defer connect and invoke the eBPF upstream logic.
- **Sockops(ebpf):**  
L7 upstream logic

Kmesh L7 upstream

*\*The blue components are kmesh components.*

# Performance Comparison With Istio

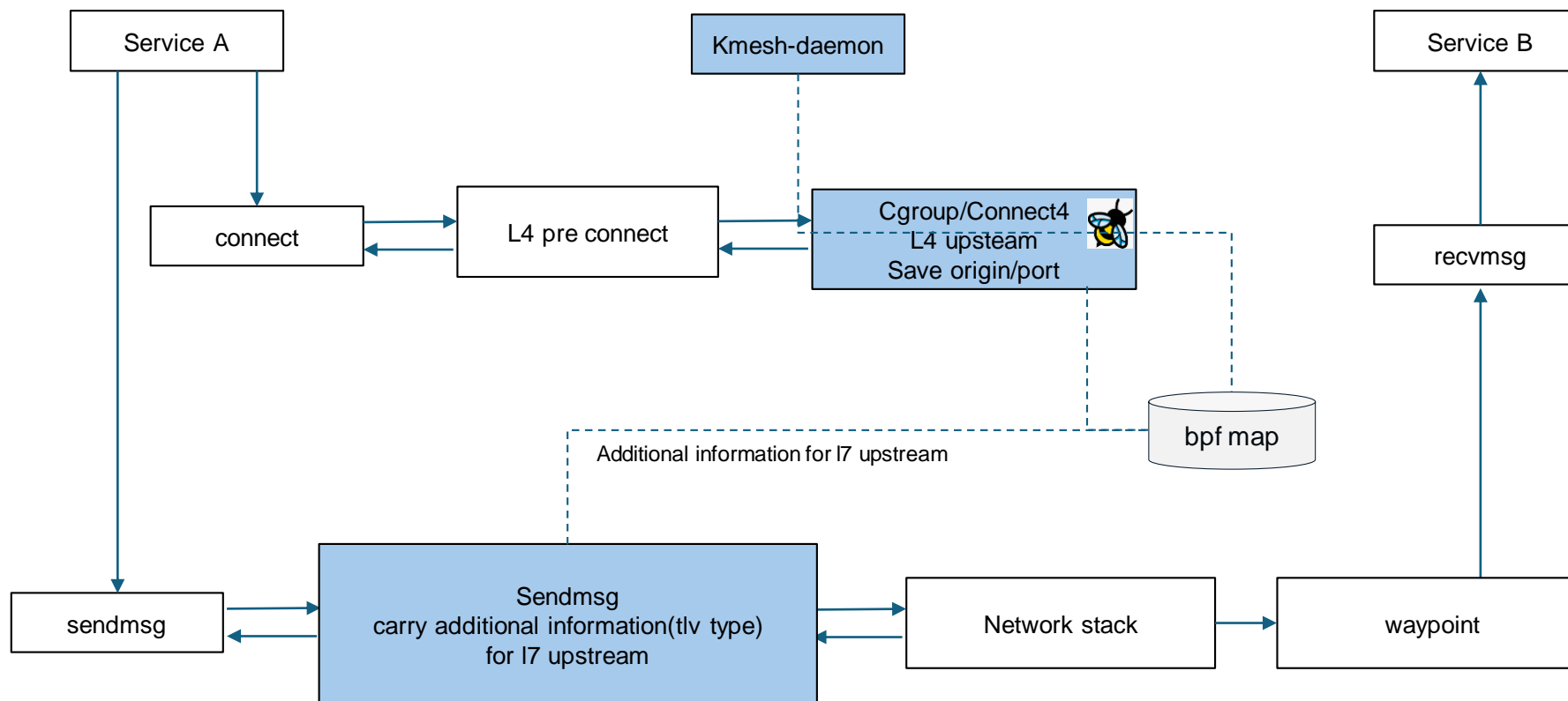


China 2024



Kernel-native: HTTP response latency almost same as baseline

# Dual Engine mode



L4-waypoint upstream  
Contains multiple components :

- **Kmesh-daemon**  
Configuration data delivery on the management and control plane
- **Cgroup/Connect4(ebpf):**  
L4 upstream logic
- **Sockops(ebpf):**  
Trigger sendmsg execution.
- **Sendmsg(ebpf):**  
Carry additional information for I7 upstream to waypoint by tlv(type-length-value) type

Kmesh L4 upstream

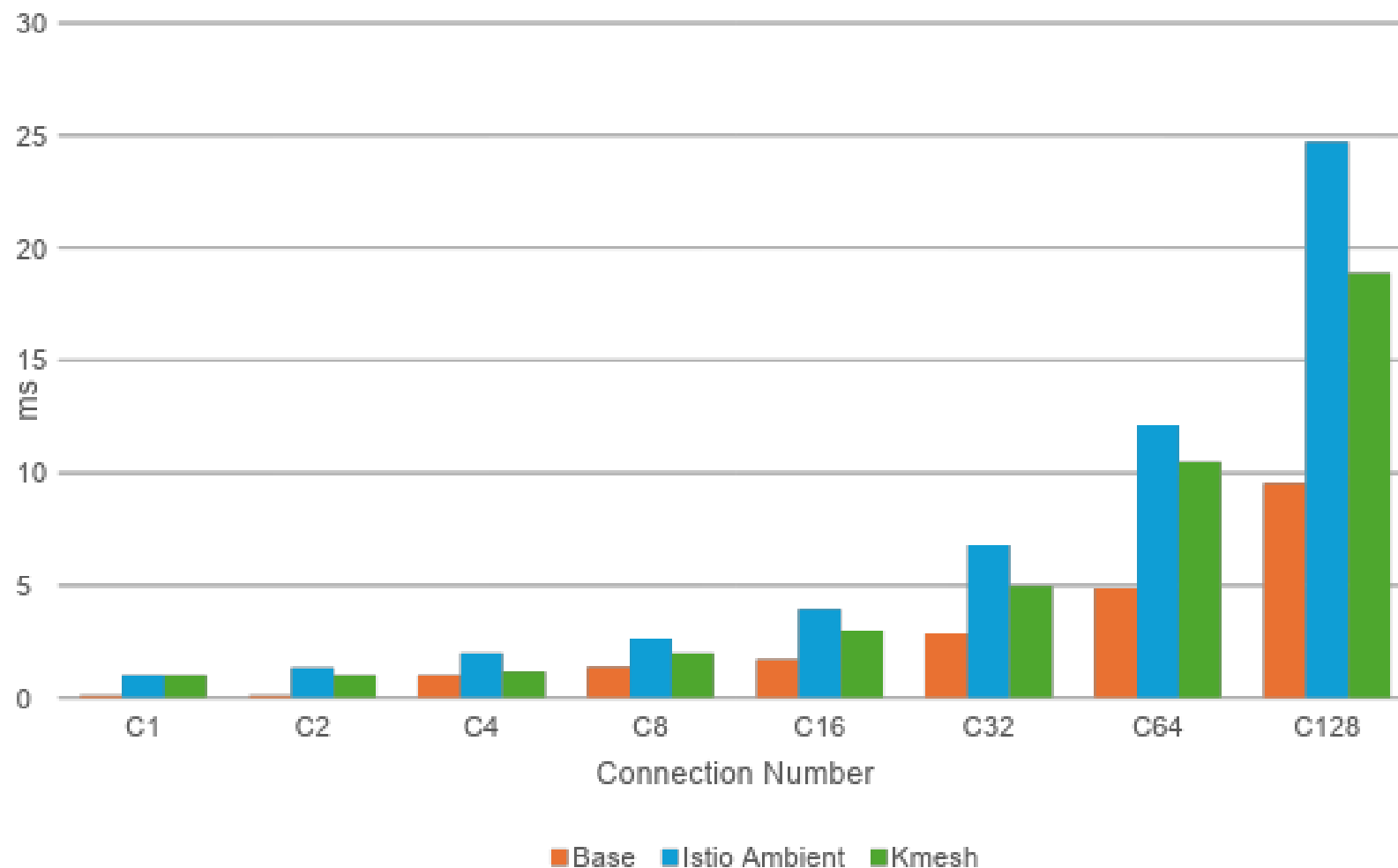
*\*The blue components are kmesh components.*

# Performance Comp with Ambient



China 2024

P99 Latency  
QPS 1024



- Kmesh is faster than ambient mesh by about 25-30%
- Both are slower than base line with L7 processing

# Agenda

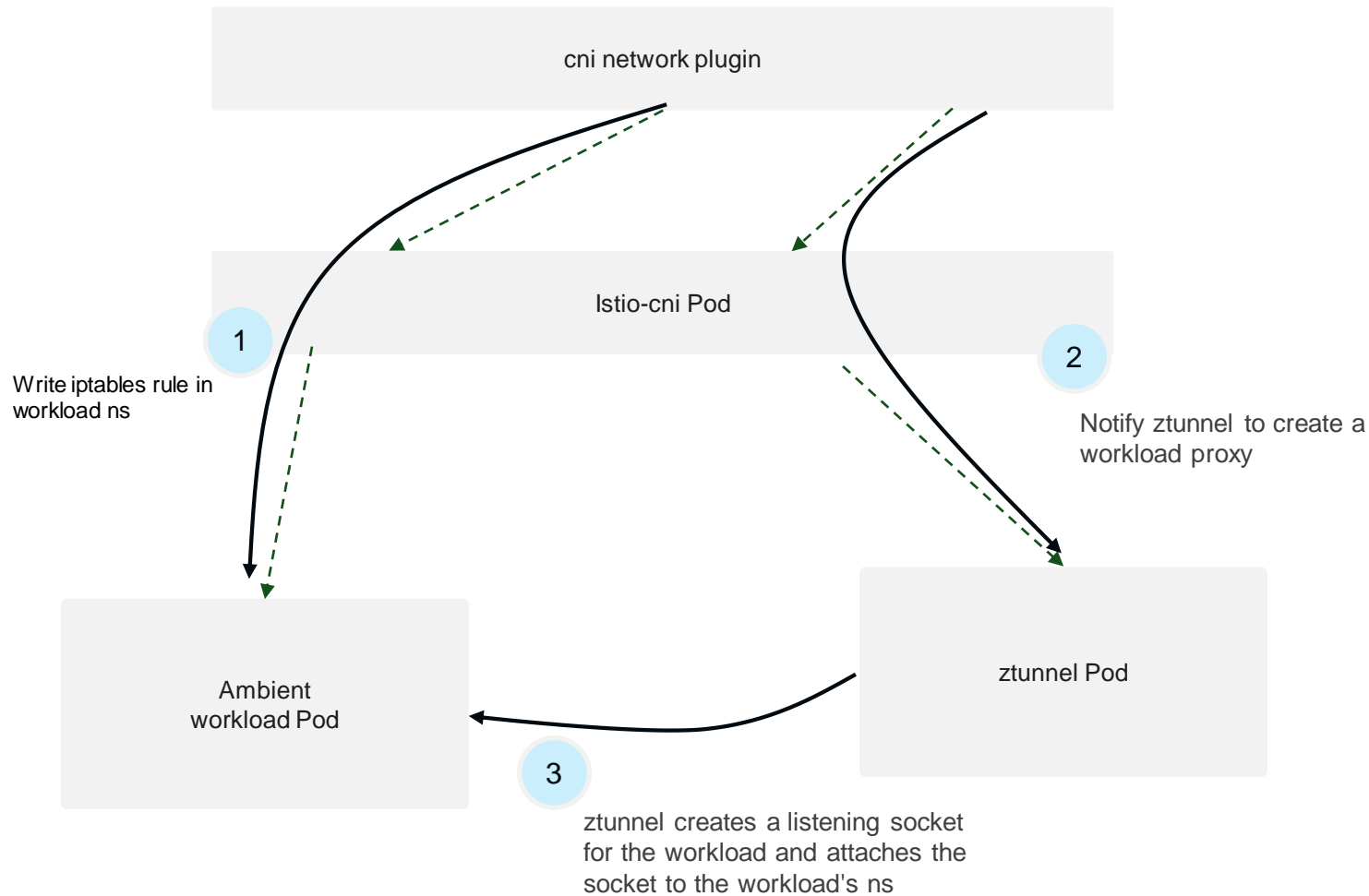


China 2024

- Kmesh Overview
- How to implement L4/L7 traffic offload
- Restarting without stop the service
- High-performance non-intrusive observability



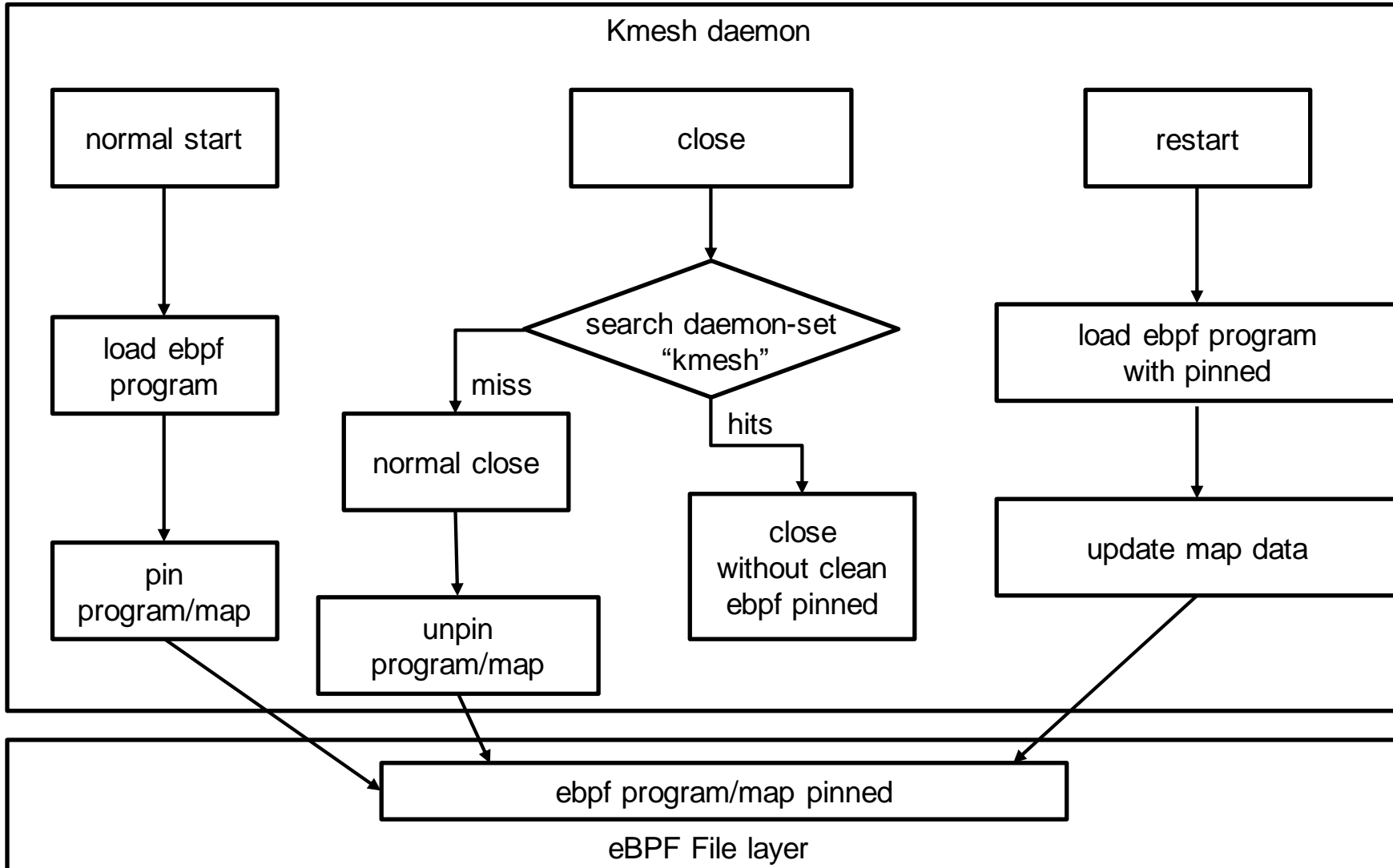
# Ambient restart



## Restart with connect interruption

- During the restart process, the original socket will be released, which will cause the original **connect to fail**.
- During the restart process, ztunnel **cannot provide effective services** for the newly established connection.

# Restart without interruption



- The program/map is pinned to the EBPf file system through the EBPf pin mechanism
- When the Kmesh-daemon exits, check whether the Kmesh-daemon exits normally or restarts
- If the logout is normal, delete the pinned file
- If the system is restarted, the pinned file is not cleared. The system directly associates the pinned file during the next startup
- When Kmesh exits, the pinned file always exists and the service is not terminated

# Agenda



China 2024

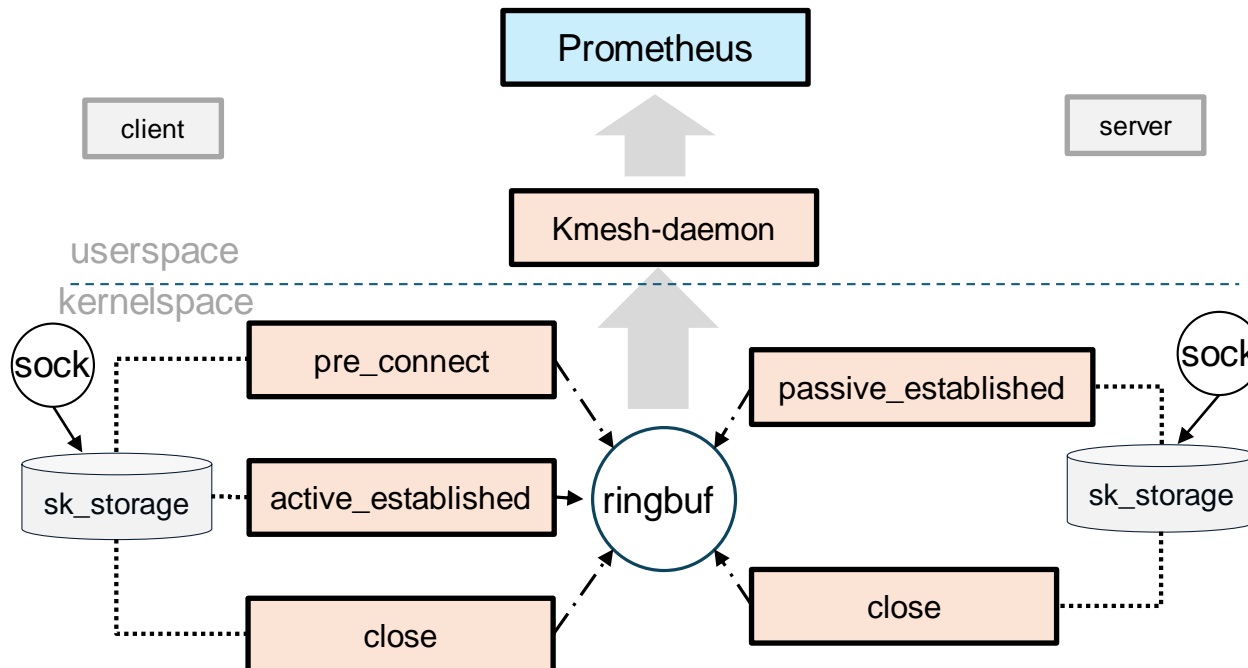
- Kmesh Overview
- How to implement L4/L7 traffic offload
- Restarting without stop the service
- High-performance non-intrusive observability

# observability

Use the kernel eBPF sockops to collect the METIRC data during the connection running.

- pre\_connect attachment
- active\_established attachment
- passive\_established attachment
- close attachment

More information can be collected in the future to support more native platforms.



advantage:

1. E2E observability base on eBPF
2. Low cpu cost. < 5%
3. More low-level data collection. Socket link level data collection.

# Agenda



China 2024

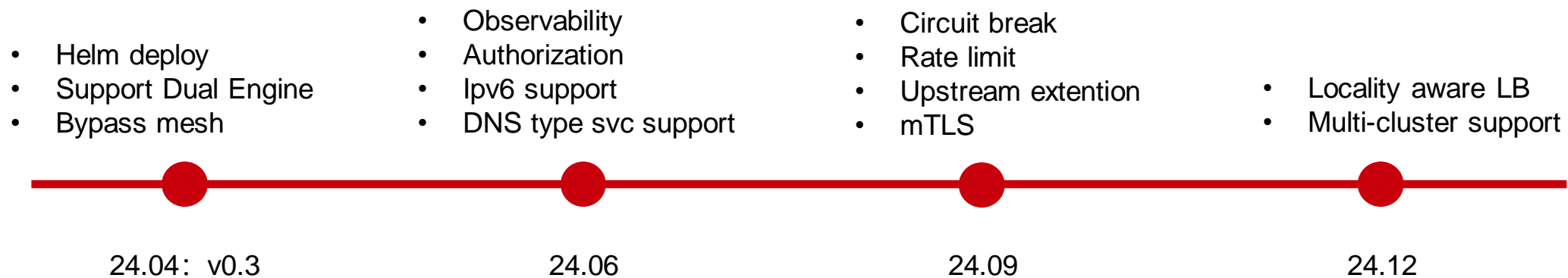
- Kmesh Overview
- How to implement L4/L7 traffic offload
- Restarting without stop the service
- High-performance non-intrusive observability



# Kmesh Roadmap



China 2024





<https://github.com/kmesh-net/kmesh>



扫码进入Kmesh技术交流群  
或添加k8s2222回复Kmesh