KubeCon

CloudNativeCon

THE LINUX FOUNDATION

OPEN SOURCE SUMMIT

AI_dev
Open Source GenAI & ML Summit

China 2024

# Agenda

- Cilium general introduction
- Shallow dive from network policy with ACK
- How CNI looks like at Alibaba Cloud
- Scalability on Alibaba Cloud
- What can you get from a full blown cilium on Alibaba
- Some highlights on cilium 1.16 release
- Q&A

eBPF-based:

- Networking
- Security
- Observability
- Service Mesh & Ingress

Foundation | Technology

**Deploy on your prefered cloud**

**Use your favorite Kubernetes distribution**

# eBPF

Makes the Linux kernel programmable in a secure and efficient way.

*"What JavaScript is to the browser, eBPF is to the Linux Kernel"*



```
int syscall__ret_execve(struct pt_regs *ctx)
{
        struct comm_event event = {
                .pid = bpf_get_current_pid_tgid() >> 32,
                .type = TYPE_RETURN,
        };

        bpf_get_current_comm(&event.comm, sizeof(event.comm));
        comm_events.perf_submit(ctx, &event, sizeof(event));

        return 0;
}
```

# Kubernetes Network policy

- Pods within the same Kubernetes cluster can communicate with each other without restriction.
- If you want to limit the traffic between pods, you will need to use a network policy.
- You can enable it by ticking the checkbox when creating the ACK cluster, as shown in the following picture.

- Can pod A talk to pod B?
- A example of the basic network policy

```yaml
apiVersion: networking.k8s.io/v1
kind: NetworkPolicy
metadata:
  name: allow-curl-allow-curl
  namespace: default
spec:
  podSelector:
    matchLabels:
      app: nginx
  policyTypes:
  - Ingress
  ingress:
  - from:
    - podSelector:
        matchLabels:
          app: curl
```

# Kubernetes Network policy

- What does this really mean to the Linux host?

```
-A cali-pi-_otwv6_8NtgmJghT8l96 -m comment --comment "cali:1GqLxVx7Oeo1hWn7" -m comment
                                 --comment "Policy default/knp.default.allow-curl ingress" -m set
                                 --match-set cali40s:s33YkCe7jRY-julDezR1ydl src -j MARK --set-xmark 0x10000/0x10000
```

- You start to chase the iptables tables/rules and ipset on the host. I found there are around 300 rules(including rules for kube-proxy) with just 2 pods and 3 nodes(1 controller and 2 workers) and no user defined service on it on my KIND cluster.
- How do I know if there is a drop with iptables rules? You need to other non standard network policy implementation to log the flow to a file
- Iptables lookup performance is O(n).

# Kubernetes Network policy

- Same policy for cilium

```
root@kind-worker2:/home/cilium# cilium bpf policy get 2572
POLICY    DIRECTION    LABELS (source:key[=value])
Allow     Ingress      reserved:host
Allow     Ingress      k8s:app=curl
                       k8s:io.cilium.k8s.namespace.labels.kubernetes.io/metadata.name=default
                       k8s:io.cilium.k8s.policy.cluster=kind-kind
                       k8s:io.cilium.k8s.policy.serviceaccount=default
                       k8s:io.kubernetes.pod.namespace=default
Allow     Egress       reserved:unknown
root@kind-worker2:/home/cilium#
```

```yaml
apiVersion: networking.k8s.io/v1
kind: NetworkPolicy
metadata:
  name: allow-curl-allow-curl
  namespace: default
spec:
  podSelector:
    matchLabels:
      app: nginx
  policyTypes:
  - Ingress
  ingress:
  - from:
    - podSelector:
        matchLabels:
          app: curl
```

```
root@kind-worker2:/home/cilium# cilium monitor --related-to 2572
Listening for events on 12 CPUs with 64x4096 of shared memory
Press Ctrl-C to quit
time="2024-08-01T19:59:41Z" level=info msg="Initializing dissection cache..." subsys=monitor
Policy verdict log: flow 0x7c26234b local EP ID 2572, remote ID 13980, proto 6, ingress, action allow, auth: disabled, match L3-Only, 10.244.2.205:58184 -> 10.244.2.75:80 tcp SYN
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state new ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp SYN
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK, FIN
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK
Policy verdict log: flow 0xd3caac78 local EP ID 2572, remote ID 17091, proto 6, ingress, action deny, auth: disabled, match none, 10.244.2.171:51846 -> 10.244.2.75:80 tcp SYN
xx drop (Policy denied) flow 0xd3caac78 to endpoint 2572, ifindex 11, file bpf_lxc.c:2091, , identity 17091->5411: 10.244.2.171:51846 -> 10.244.2.75:80 tcp SYN
```

# Kubernetes Network policy

- Network policy lookup from cilium is O(1).
- You can easily observe it with cilium tool
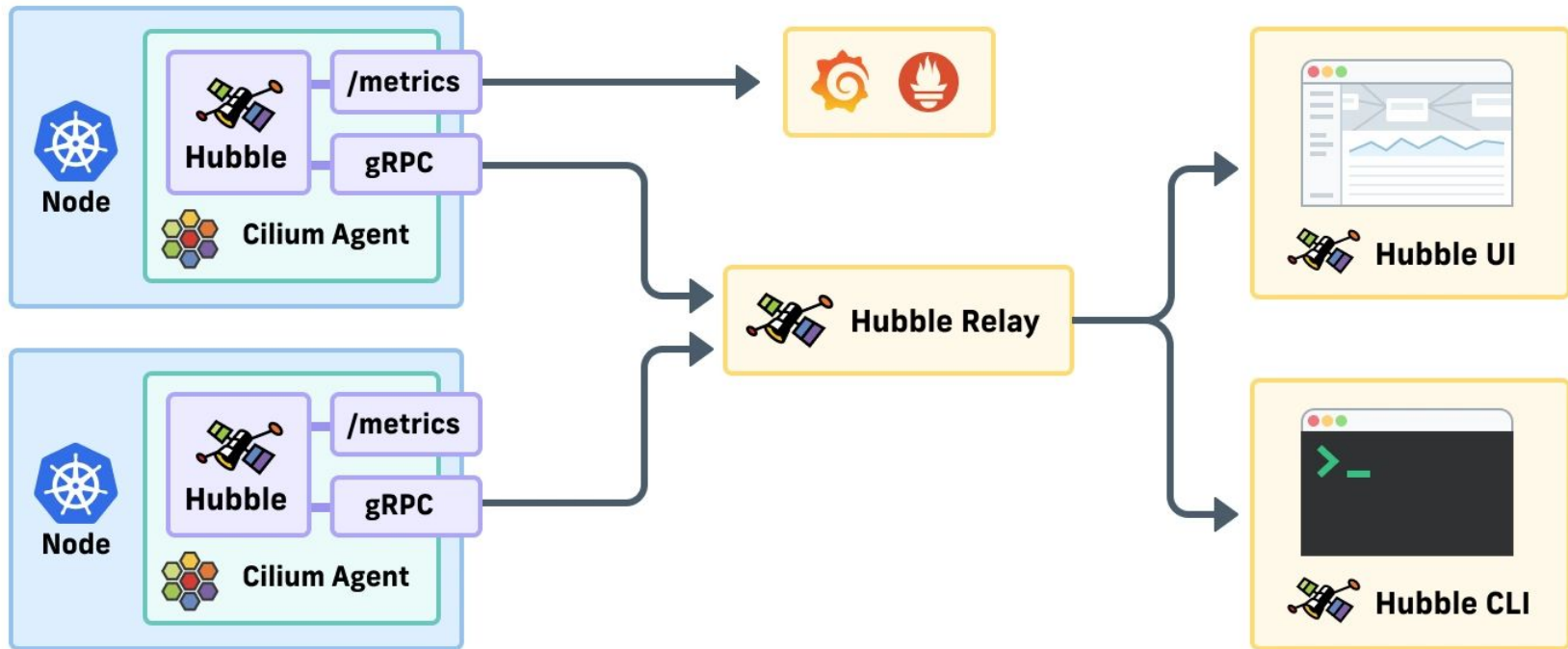- More advanced Cilium network policy will be discussed later

# Hubble Intro

- We still need to get to the cilium container to run cilium monitor command.
- What if we have a tool to see all the flow logs on a cluster even more cluster?
- Can we export all the logs to SIEM?
- Can we generate the network policy based on flow data?

```
root@kind-worker2:/home/cilium# cilium monitor --related-to 2572
Listening for events on 12 CPUs with 64x4096 of shared memory
Press Ctrl-C to quit
time="2024-08-01T19:59:41Z" level=info msg="Initializing dissection cache..." subsys=monitor
Policy verdict log: flow 0x7c26234b local EP ID 2572, remote ID 13980, proto 6, ingress, action allow, auth: disabled, match L3-Only, 10.244.2.205:58184 -> 10.244.2.75:80 tcp SYN
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state new ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp SYN
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK, FIN
-> endpoint 2572 flow 0x7c26234b , identity 13980->5411 state established ifindex lxc29649753a2bf orig-ip 10.244.2.205: 10.244.2.205:58184 -> 10.244.2.75:80 tcp ACK
Policy verdict log: flow 0xd3caac78 local EP ID 2572, remote ID 17091, proto 6, ingress, action deny, auth: disabled, match none, 10.244.2.171:51846 -> 10.244.2.75:80 tcp SYN
xx drop (Policy denied) flow 0xd3caac78 to endpoint 2572, ifindex 11, file bpf_lxc.c:2091, , identity 17091->5411: 10.244.2.171:51846 -> 10.244.2.75:80 tcp SYN
```

# Hubble Overview

# Hubble CLI

- Observe the traffic for the whole cluster
- Filter with ip/pod/svc/namespace/fqdn/http/type etc..

# Hubble GUI

# Hubble metrics

# About Alibaba Cloud

**No.1**

Market Share in the Asia Pacific

**89**

Availability Zones

**30**

Regions

Container Service @alibabacloud

**Alibaba Cloud**

| Computing | Storage | Networking | Security | Database | Analytics Computing | Container & Middleware |

**Container Service for Kubernetes (ACK)**

A certified Kubernetes platform

**ApsaraMQ for Kafka**

Fully-managed and out-of-the-box Message Queue service po...

**ApsaraMQ for RabbitMQ**

An out-of-the-box fully managed RabbitMQ service

**Application Real-Time Monitoring Service (ARMS)**

Build business monitoring capabilities

**Container Registry (ACR)**

A secure image hosting platform

**Microservices Engine (MSE)**

One-stop Platform Compatible with Mainstream Open Source ...

**Serverless Kubernetes Service (ASK)**

Highly elastic and reliable serverless Kubernetes service for en...

# How CNI looks like at Alibaba Cloud

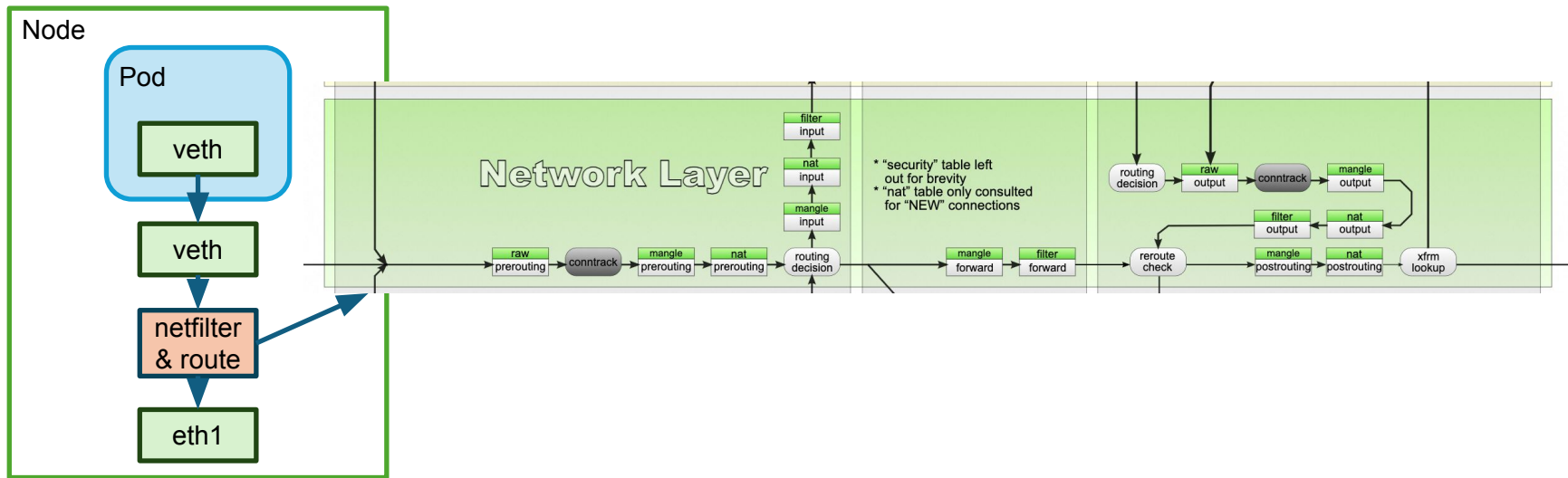|  | ACK Flannel | Terway |
|---|---|---|
| Network Type | VPC Route | ECS ENI |
| Accelerated Networking | None | ipvlan & eBPF,datapath V2 |
| Scale | 200 nodes ( up to 1000 nodes) | 5000 nodes ( up to 15000 nodes) |
| Security | None | Pod Security Group, NetworkPolicy, ACK GlobalNetworkPolicy |
| IPAM | Fixed size for every node | Elastic, can enlarge any time. Support fixed IP. |
| NFV | None | RDMA,eRDMA,SMC-R,SRIOV,DPDK... |
| Pod N/S Communication | None | EIP,DNAT Gateway, IPv6 Gateway (with ack-extend-netwokr-controller) |
| Loadbalancer Backend | NodePort | Pod IP |

# Overhead on stander datapath



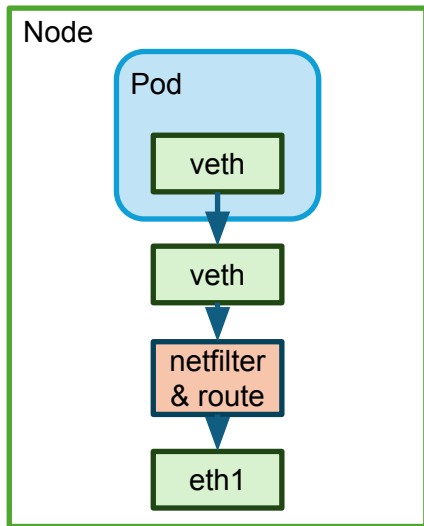The reason for the overhead is the packetization and the length of the kernel path

Standard datapath
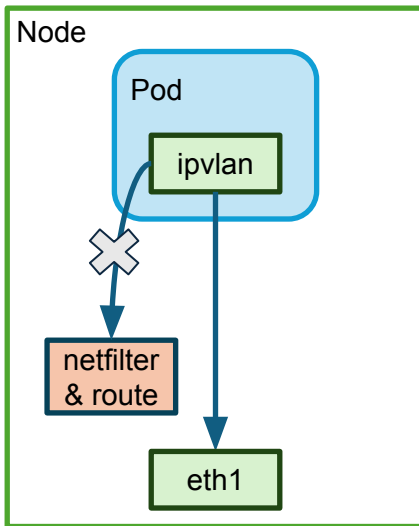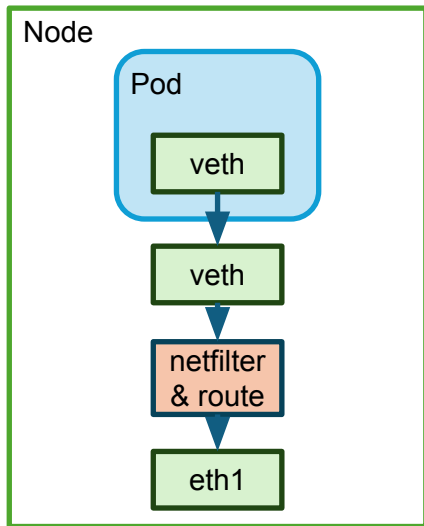
# Datapath IPvlan



Standard datapath

IPvlan datapath

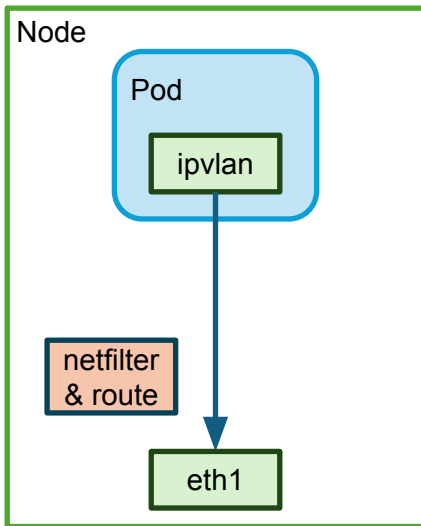IPvlan allows bypassing the host's network stack, but it poses challenges with Service functionality.

# Datapath IPvlan + eBPF



Standard datapath

IPvlan datapath

IPvlan + eBPF datapath

# Features we used



IPvlan + eBPF datapath

- KPR
  - partial, only for containers
- NetworkPolicy
  - K8s NetworkPolicy
  - CiliumClusterWideNetworkPolicy
    - [Use ACK GlobalNetworkPolicy - Container Service for Kubernetes - Alibaba Cloud Documentation Center](#)
- BandwidthManager
  - Egress side, EDT at kernel 5.10
- Hubble
  - [Implement network observability by using ACK Terway and Cilium Hubble - Container - Alibaba Cloud](#)

# CNI Chaining

# Limitation on IPvlan+eBPF

- Traffic will not go through node
  - Need additional redirect rule for traffic like nodelocal dns
  - Monitor may need addition adapt for the datapath
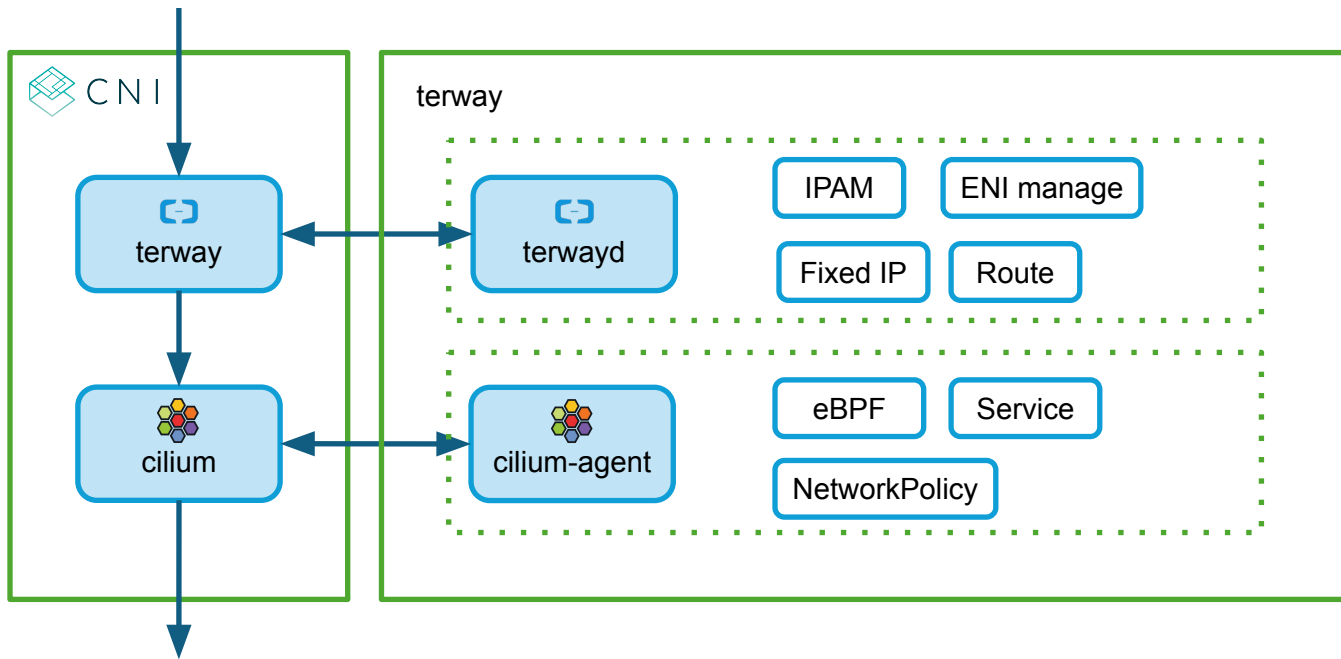- Connectivity issue
  - NodePort may not be reachable
- Performance
  - Traffic between pods on the same node may route through the VPC



IPvlan + eBPF datapath

# Datapath V2

Just like stander datapath, but
enhanced with eBPF

- Pod to world
  - bpf_redirect_neigh
- Reverse package
  - bpf_redirect_peer

Enhanced compatibility

- Pod traffic can be tracked on host
- Pod to Pod on same node, will no longer go through VPC

# Result

TCP_RR Latency P99



TCP_STREAM (1024 size)

- The IPvlan mode provides the best pod-to-pod performance
- Datapath V2 performance significantly outperforms veth mode, coming very close to IPvlan mode.
- In certain scenarios, Datapath V2 performs better than IPvlan mode.

# Cilium case study with Alibaba cloud

- https://www.cncf.io/case-studies/alibaba/

- netkit provide a faster network namespace switch for off-node traffic
- Full kube-proxy replacement can simplify the deployment configuration
- Network function offloading may be the final form

# Real world disasters

Chang single namespace label.
Looks harmless...

Result in to massive pressure in
kube-apiserver.

2K+ Nodes
80K+ Pods

# Cilium architecture

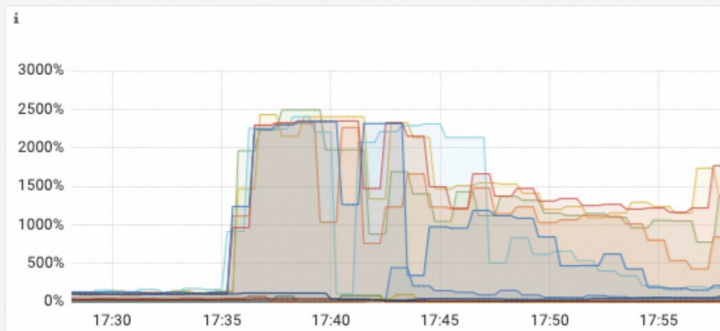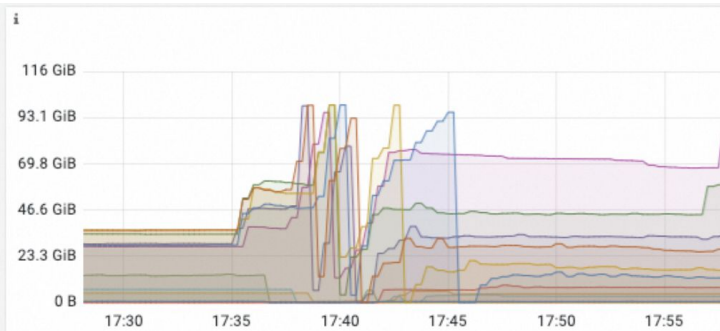- CiliumEndpoint(CEP)
  - Track a pod
  - Every container network pod
  - Contains pod label, CiliumIdentity

- CiliumIdentity(CID)
  - Generated by pod label or CIDR (defined in NetworkPolicy)
  - Used in NetworkPolicy

# Optimize in Alibaba Cloud

Changed

- Watch on demand
  - Added a node label for CEP. Default only watch the CEPs related to this node
- Limit the lables used in CID
- Simplify the fields in the CEP definition
- Do not sync pod labels to CEP labels
- Implement rate limiting on Kubernetes API for Cilium resources

Result

- Memory consumption has decreased by 82.5%
- The convergence time affected by the change has decreased by 95%

- Cilium network policy with FQDN or HTTP info

```yaml
apiVersion: "cilium.io/v2"
kind: CiliumNetworkPolicy
metadata:
  name: "fqdn"
spec:
  endpointSelector:
    matchLabels:
      org: empire
      class: mediabot
  egress:
  - toFQDNs:
    - matchName: "api.github.com"
  - toEndpoints:
    - matchLabels:
        "k8s:io.kubernetes.pod.namespace": kube-system
        "k8s:k8s-app": kube-dns
    toPorts:
    - ports:
      - port: "53"
        protocol: ANY
      rules:
        dns:
        - matchPattern: "*"
```

```yaml
apiVersion: "cilium.io/v2"
kind: CiliumNetworkPolicy
metadata:
  name: "l7-rule"
spec:
  endpointSelector:
    matchLabels:
      app: myService
  ingress:
  - toPorts:
    - ports:
      - port: '80'
        protocol: TCP
      rules:
        http:
        - method: GET
          path: "/path1$"
        - method: PUT
          path: "/path2$"
          headers:
          - 'X-My-Header: true'
```
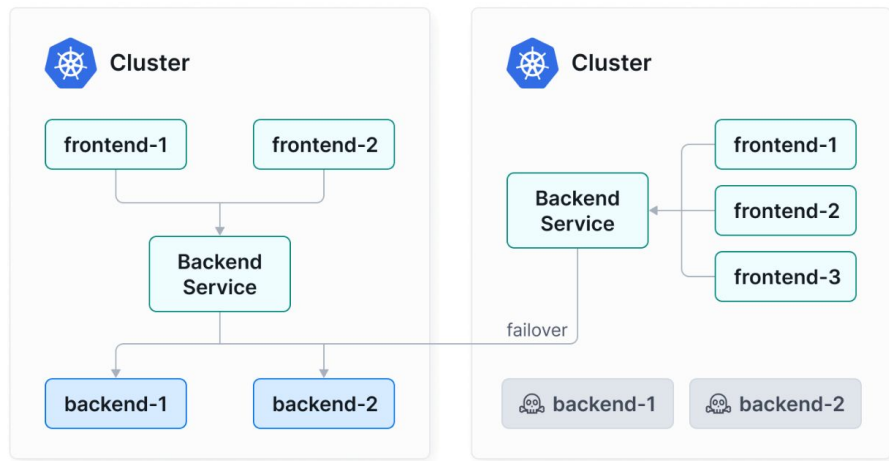
# A full blown cilium on Alibaba

- **Cluster mesh**



```yaml
apiVersion: "cilium.io/v2"
kind: CiliumNetworkPolicy
metadata:
  name: "allow-cross-cluster"
spec:
  description: "Allow x-wing in cluster1 to contact rebel-base in cluster2"
  endpointSelector:
    matchLabels:
      name: x-wing
      io.cilium.k8s.policy.cluster: cluster1
  egress:
  - toEndpoints:
    - matchLabels:
        name: rebel-base
        io.cilium.k8s.policy.cluster: cluster2
```

# Cilium 1.16 release

**Networking**
- Cilium netkit: container-network throughput and latency as fast as host-network

**Service mesh & Ingress/Gateway API**
- Gateway API GAMMA support: East-west traffic management for the cluster via Gateway API
- Gateway API 1.1 support: Cilium now supports Gateway API 1.1

**Security**
- All kinds of enhancements for the Cilium Network Policy

More details on https://isovalent.com/blog/post/cilium-1-16/