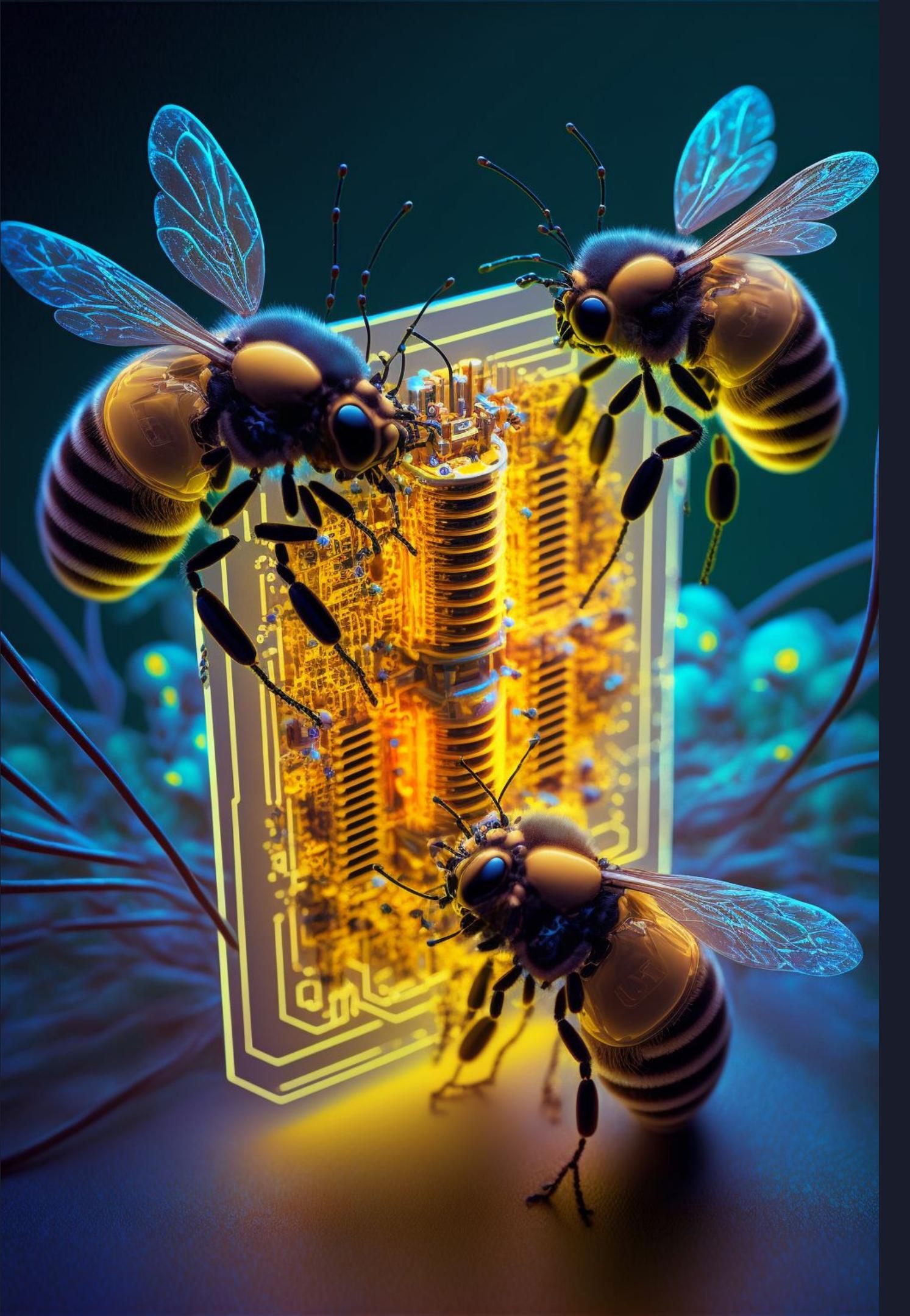


ISOVALENT

Turn Your Kubernetes Cluster Into a Castle Using Advanced Cilium Network Policies

Speaker: Philip Schmid





Agenda

- ◆ Introduction
- ◆ Securing the Nodes
- ◆ Protecting Infra Components
- ◆ Default Rule Set for User Workloads



About Me



Philip Schmid

Solutions Architect @ Isovalent

CK(AD|A|S) | CCNP



Introduction

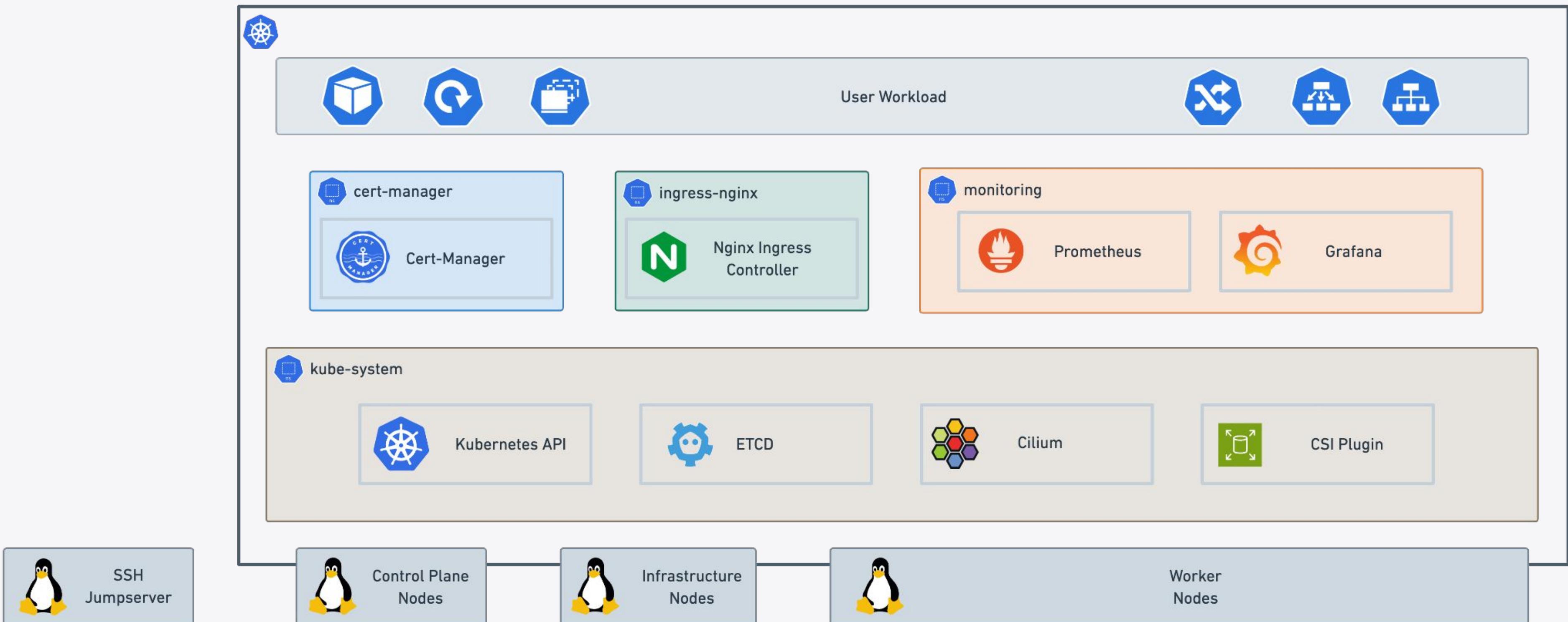


Goal

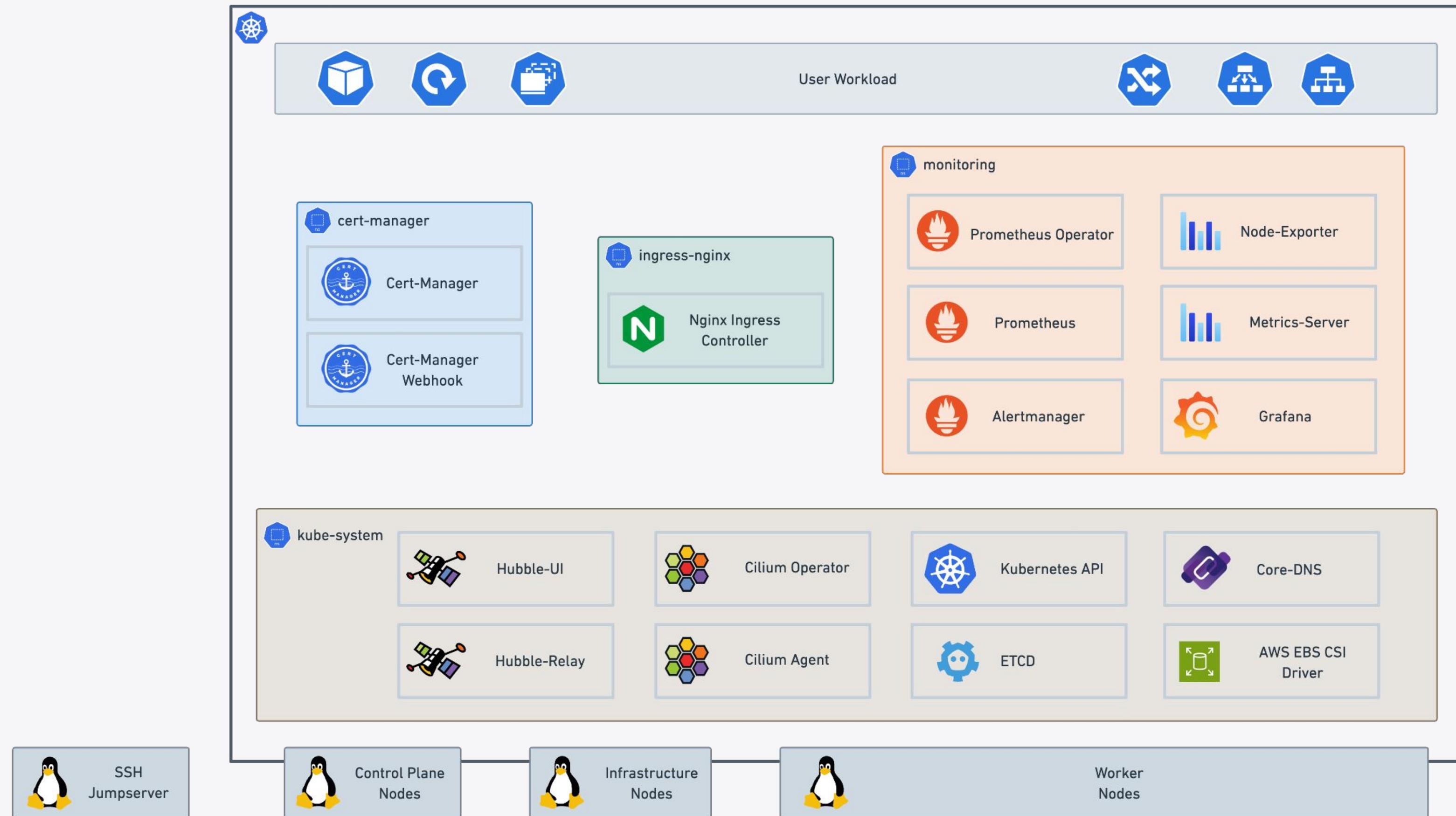
- This isn't a ...
 - ... Network Policy 101 talk
 - ... abstract Network Policy design pattern talk
- This talk shows a **practical example** of how **you could** secure ...
 - ... **your** Kubernetes nodes and clusters in general
 - ... **your** Kubernetes infrastructure components and user workloads
- ... by using advanced Cilium Network Policies



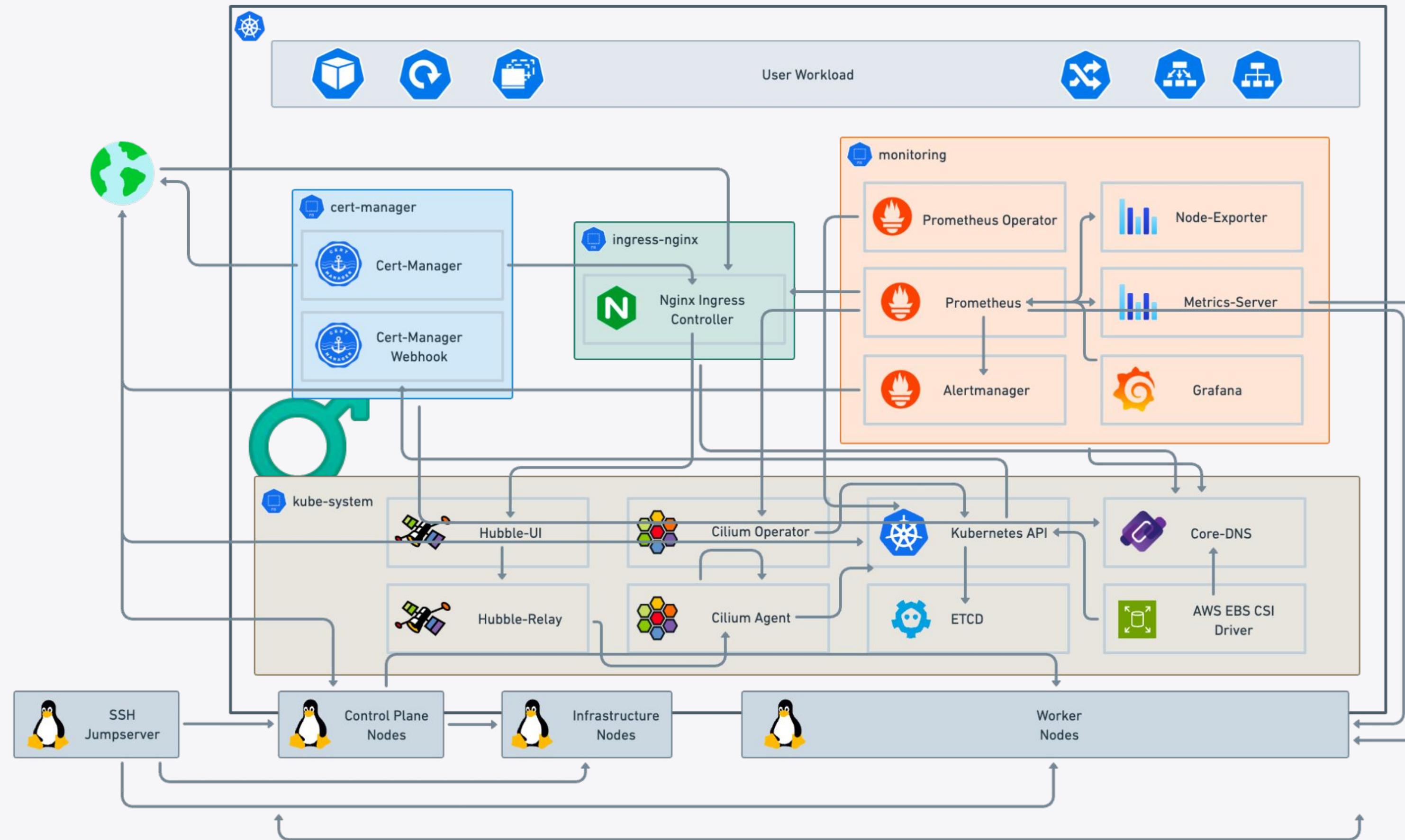
Demo Use Case



Demo Use Case in Detail



Demo Use Case in Detail with Connections



Structured Way to Succeed with Network Policies

1. General recommendation: Take care of Network Policies as early as possible
2. Start bottom up:
 - a. Nodes
 - b. Kubernetes
 - c. Cilium
 - d. Infrastructure components
 - e. User workloads
3. Always keep an eye on DROPPED connections:

```
1 hubble observe -t policy-verdict -f --verdict AUDIT --verdict DROPPED
```

4. Re-iterate over policies until there are no false-positive left

Securing the Nodes

Motivation for Securing the Nodes

- Security team: “*Is the Linux host firewall enabled?*”
- Industry standards like
 - PCI DSS (v4.0): “*1.2 Network security controls (NSCs) are configured and maintained*”
 - ISO 27001:2013: A.13.1.1 Network Controls
- K8s Ops team: Improve own security posture
- → **Use Cilium Host Firewall / Policies, especially for high security environments**

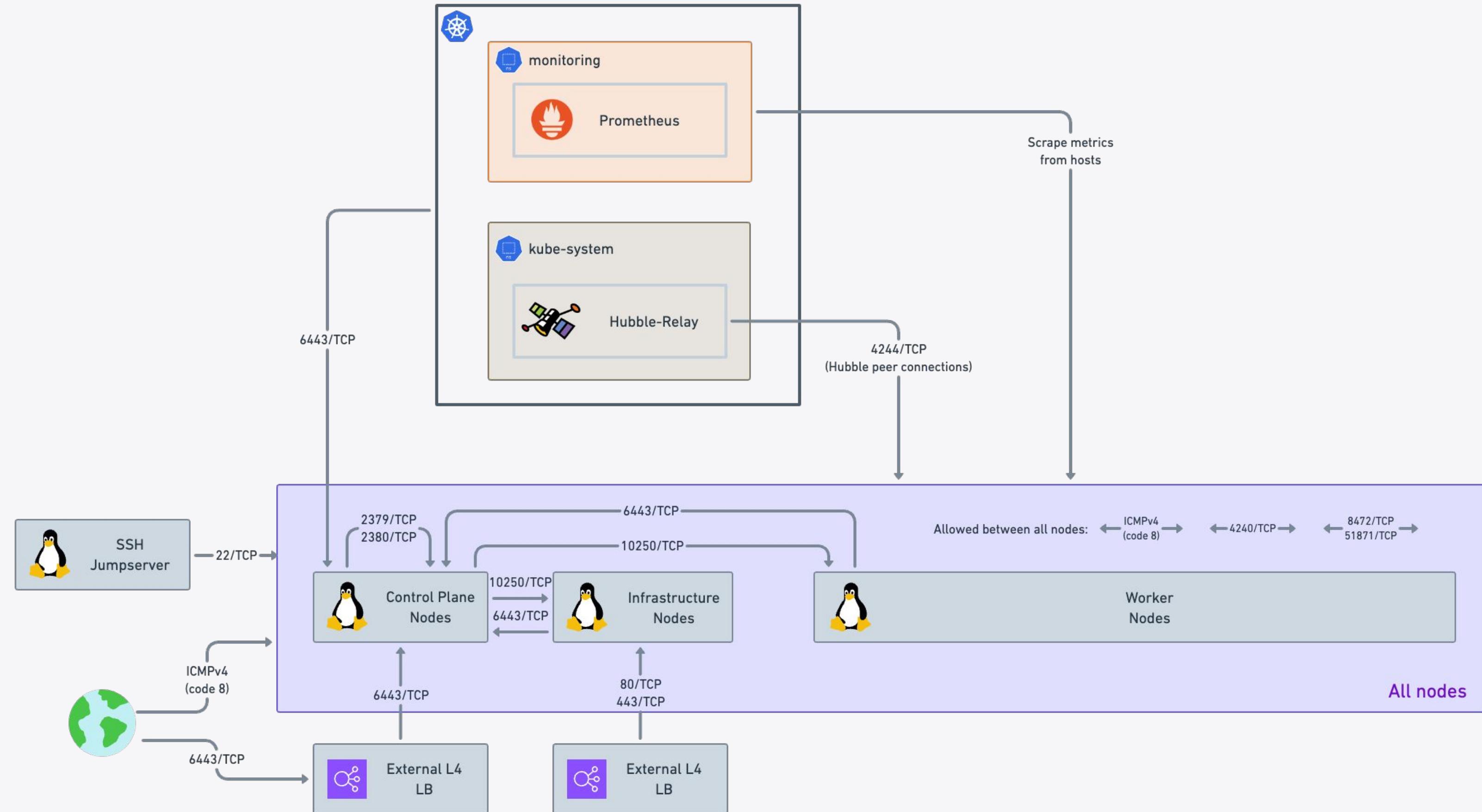


Securing the Nodes using Cilium Host Firewall

- Implement host firewalling solely using eBPF (no netfilter!)
- Explicitly enable it:
hostFirewall.enabled=true
- Rules are fully managed via Kubernetes
 - CiliumClusterwideNetworkPolicy with spec.nodeSelector (instead of spec.endpointSelector)
- **Important:** After enabling Cilium Host Firewall, **enable Policy Audit mode** for the Host CiliumEndpoints before applying any host policies (see [doc](#))!

```
1 apiVersion: "cilium.io/v2"
2 kind: CiliumClusterwideNetworkPolicy
3 metadata:
4   name: host-cp
5 spec:
6   description: Cilium host policy set for CP nodes
7   nodeSelector:
8     matchLabels:
9       node-role.kubernetes.io/control-plane: ""
10  ingress:
11    # Kubernetes control plane connectivity from everywhere
12    - toPorts:
13      - ports:
14        # Kubernetes API
15        - port: "6443"
16        protocol: TCP
17    # ETCD control plane connectivity from nodes
18    - fromEntities:
19      - kube-apiserver
20    toPorts:
21      - ports:
22        # etcd client port
23        - port: "2379"
24        protocol: TCP
25        # etcd peer communication port
26        - port: "2380"
27        protocol: TCP
```

Securing the Nodes using Cilium Host Firewall



Securing the Nodes using Cilium Host Firewall

Control Plane nodes

```
1 apiVersion: "cilium.io/v2"
2 kind: CiliumClusterwideNetworkPolicy
3 metadata:
4   name: host-cp
5 spec:
6   description: Cilium host policy set for CP nodes
7   nodeSelector:
8     matchLabels:
9       node-role.kubernetes.io/control-plane: ""
10  ingress:
11    # Kubernetes control plane connectivity from everywhere
12    - toPorts:
13      - ports:
14        # Kubernetes API
15        - port: "6443"
16          protocol: TCP
17    # ETCD control plane connectivity from nodes
18    - fromEntities:
19      - kube-apiserver
20    toPorts:
21      - ports:
22        # etcd client port
23        - port: "2379"
24          protocol: TCP
25        # etcd peer communication port
26        - port: "2380"
27          protocol: TCP
```



Securing the Nodes using Cilium Host Firewall

Infra nodes

```
1 apiVersion: "cilium.io/v2"
2 kind: CiliumClusterwideNetworkPolicy
3 metadata:
4   name: host-infra
5 spec:
6   description: Cilium host policy set for infra nodes
7   nodeSelector:
8     matchLabels:
9       node-role.kubernetes.io/infra: ""
10  ingress:
11    # Data plane connectivity from everywhere to Nginx
12    - toPorts:
13      - ports:
14        # Ingress HTTP
15        - port: "80"
16        protocol: TCP
17        # Ingress HTTPS
18        - port: "443"
19        protocol: TCP
```



Securing the Nodes using Cilium Host Firewall

All nodes – part 1

```
1 apiVersion: "cilium.io/v2"
2 kind: CiliumClusterwideNetworkPolicy
3 metadata:
4   name: host-all
5 spec:
6   description: Cilium host policy for all nodes
7   nodeSelector: {}
8   ingress:
9     # System management: Allow SSH access from jumphost
10    - fromCIDR:
11      - 10.1.101.110/32
12      toPorts:
13        - ports:
14          - port: "22"
15            protocol: TCP
16     # System monitoring: Allow ping to hosts
17    - icmps:
18      - fields:
19        - type: 8
20        family: IPv4
21     # KAPI Control plane connectivity for Kubelet Pod log streaming
22    - fromEntities:
23      - kube-apiserver
24      toPorts:
25        - ports:
26          - port: "10250"
27            protocol: TCP
```

Securing the Nodes using Cilium Host Firewall

All nodes – part 2

```
28 # Cilium control plane agent HTTP health checks
29 - fromEntities:
30   - remote-node
31   toPorts:
32     - ports:
33       - port: "4240"
34         protocol: TCP
35 # Data plane inter-node connectivity
36 - fromEntities:
37   - host
38   - remote-node
39   toPorts:
40     - ports:
41       # Cilium VXLAN
42       - port: "8472"
43         protocol: UDP
44       # Cilium WireGuard
45       - port: "51871"
46         protocol: UDP
47 # Cilium Hubble relay access to Cilium agent Hubble peers
48 - fromEndpoints:
49   - matchLabels:
50     k8s:io.kubernetes.pod.namespace: kube-system
51     k8s:app.kubernetes.io/name: hubble-relay
52   toPorts:
53     - ports:
54       - port: "4244"
55         protocol: TCP
56 # Allow Metrics-Server to reach Kubelet metrics
57 - fromEndpoints:
58   - matchLabels:
59     k8s:io.kubernetes.pod.namespace: monitoring
60     k8s:app.kubernetes.io/name: metrics-server
61   toPorts:
62     - ports:
63       - port: "10250"
64         protocol: TCP
```

Protecting Infra Components

Goals for Protecting Infra Components

1. **Allow all** ingress and egress connections **within an infra namespace**
2. Try to **only allow specific** ingress and egress **connections from and to the infra namespace**
3. Implement a **restrictive, static rule set** that doesn't need to be updated when a new user workload NS is created AND doesn't require the user workload NS to allow something from/to infra components (Core-DNS, Nginx, Prometheus)





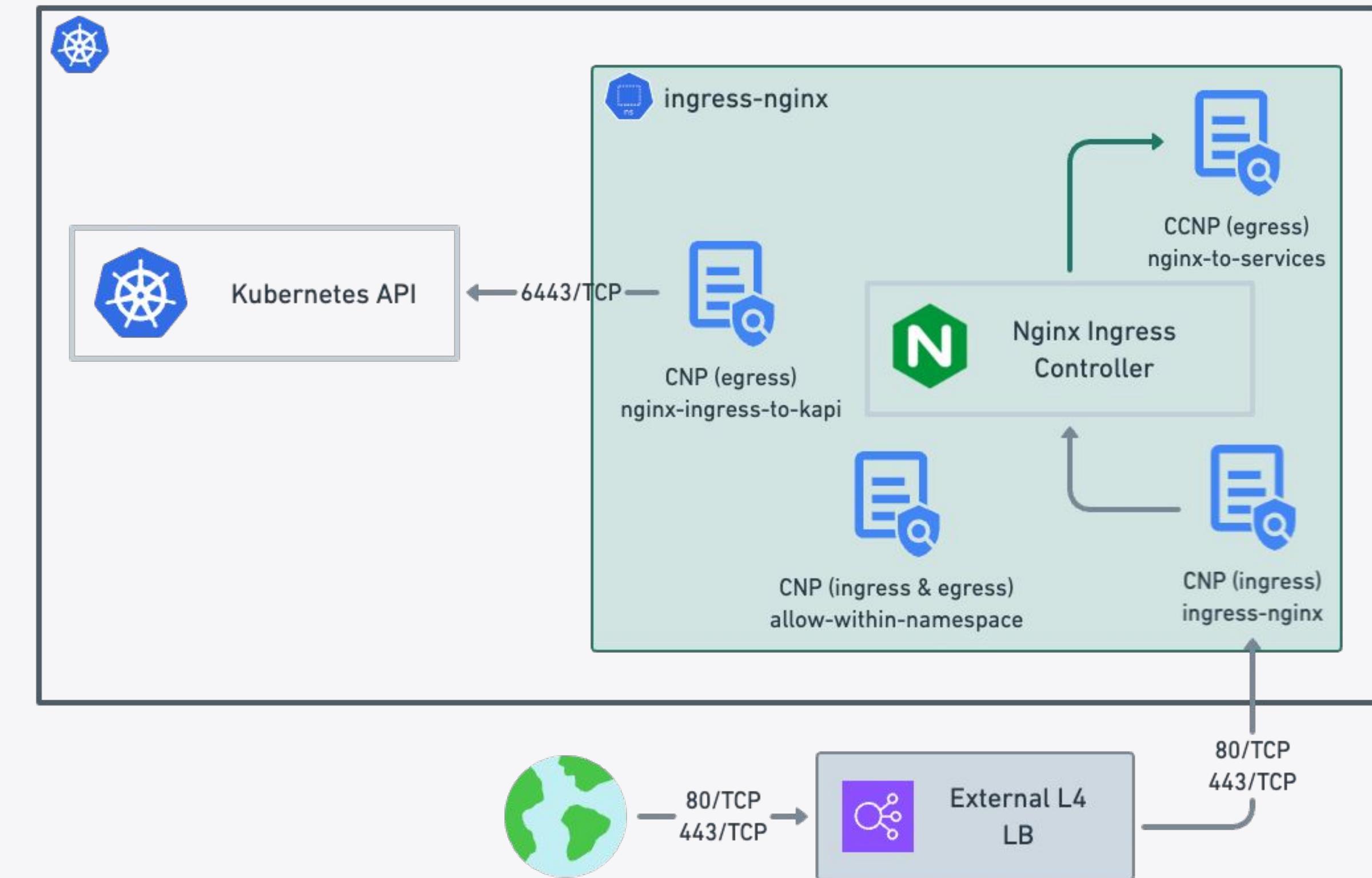
Protecting Infra Components

1. **Allow all** ingress and egress connections **within an infra namespace** ✓

```
1 apiVersion: cilium.io/v2
2 kind: CiliumNetworkPolicy
3 metadata:
4   name: allow-within-namespace
5   namespace: infra-ns-name
6 spec:
7   endpointSelector: {}
8   ingress:
9     - fromEndpoints:
10       - {}
11   egress:
12     - toEndpoints:
13       - {}
```

Protecting Infra Components

2. Try to **only allow specific** ingress and egress **connections from and to the infra namespace**



Protecting Infra Components

2. Try to **only allow specific** ingress and egress **connections from and to the infra namespace**

```
1 —
2 apiVersion: "cilium.io/v2"
3 kind: CiliumNetworkPolicy
4 metadata:
5   name: allow-within-namespace
6   namespace: ingress-nginx
7 spec:
8   description: Allow NS internal traffic, block everything else
9   endpointSelector: {}
10  ingress:
11    - fromEndpoints:
12      - {}
13  egress:
14    - toEndpoints:
15      - {}
```

```
1 —
2 apiVersion: "cilium.io/v2"
3 kind: CiliumNetworkPolicy
4 metadata:
5   name: all-to-nginx
6   namespace: ingress-nginx
7 spec:
8   description: Allow accessing Nginx ingress controller from everywhere
9   endpointSelector:
10    matchLabels:
11      k8s:app.kubernetes.io/name: ingress-nginx
12   ingress:
13     - fromEntities:
14       - all
15     toPorts:
16       - ports:
17         - port: "80"
18           protocol: TCP
19         - port: "443"
20           protocol: TCP
```



Protecting Infra Components

2. Try to **only allow specific** ingress and egress **connections from and to the infra namespace** 

```
1 --
2 apiVersion: "cilium.io/v2"
3 kind: CiliumNetworkPolicy
4 metadata:
5   name: nginx-ingress-to-kapi
6   namespace: ingress-nginx
7 spec:
8   description: Allow Nginx ingress controller to reach KAPI
9   endpointSelector:
10    matchLabels:
11      k8s:app.kubernetes.io/name: ingress-nginx
12   egress:
13     - toEntities:
14       - kube-apiserver
15     toPorts:
16       - ports:
17         - port: "6443"
18           protocol: TCP
```

Protecting Infra Components



3. Implement a **restrictive, static rule set** that doesn't need to be updated when a new user workload NS is created AND doesn't require the user workload NS to allow something from/to infra components (Core-DNS, Nginx, Prometheus) 🤔

This requires:

1. A generic egress rule which allows leaving the infra NS toward user NS (only specific ones) 🦄
2. A generic ingress rule which allows accessing user NS (only specific ones) ✨
3. 1. & 2. but also in the other direction for passive infra components (e.g. Core-DNS)



Protecting Infra Components

Core-DNS:

```
1 —
2 apiVersion: "cilium.io/v2"
3 kind: CiliumClusterwideNetworkPolicy
4 metadata:
5   name: global-infra-dns
6 spec:
7   description: Allow coredns access from everywhere
8   endpointSelector: {}
9   egress:
10  - toEndpoints:
11    - matchLabels:
12      k8s:io.kubernetes.pod.namespace: kube-system
13      k8s:k8s-app: kube-dns
14   toPorts:
15     - ports:
16       - port: "53"
17         protocol: ANY
18   rules:
19     dns:
20       - matchPattern: "*"
```

```
1 —
2 apiVersion: "cilium.io/v2"
3 kind: CiliumNetworkPolicy
4 metadata:
5   name: allow-to-coredns
6   namespace: kube-system
7 spec:
8   description: Allow ingress to coredns from everywhere
9   endpointSelector:
10  matchLabels:
11    k8s:io.kubernetes.pod.namespace: kube-system
12    k8s:k8s-app: kube-dns
13   ingress:
14     - fromEntities:
15       - all
16     toPorts:
17       - ports:
18         - port: "53"
19           protocol: ANY
```

Protecting Infra Components

- **Problem:** The egress CiliumClusterwideNetworkPolicy for Core-DNS was intended for all namespaces → Nginx ingress and Prometheus should only be able to access certain namespaces!
- **Idea:** Add predefined labels to namespaces to which only specific infrastructure components should have access to 
- Luckily, a CiliumClusterwideNetworkPolicy has knowledge of labels on namespaces
 - CiliumNetworkPolicies are namespaced → no knowledge about NS labels

```
1 k8s:io.cilium.k8s.namespace.labels.<label-key>: <label-value>
```

Protecting Infra Components



Nginx Ingress Controller:

```
1 —  
2 apiVersion: "cilium.io/v2"  
3 kind: CiliumClusterwideNetworkPolicy  
4 metadata:  
5   name: nginx-to-services  
6 spec:  
7   description: Allow Nginx to access svc of NS with exposed label  
8   endpointSelector:  
9     matchLabels:  
10       k8s:app.kubernetes.io/name: ingress-nginx  
11       k8s:io.kubernetes.pod.namespace: ingress-nginx  
12   egress:  
13     - toEndpoints:  
14       - matchLabels:  
15         k8s:io.cilium.k8s.namespace.labels.exposed: "true"
```

```
1 —  
2 apiVersion: "cilium.io/v2"  
3 kind: CiliumClusterwideNetworkPolicy  
4 metadata:  
5   name: global-infra-nginx  
6 spec:  
7   description: Allow connections from Nginx to labeled namespaces  
8   endpointSelector:  
9     matchLabels:  
10       k8s:io.cilium.k8s.namespace.labels.exposed: "true"  
11   ingress:  
12     - fromEndpoints:  
13       - matchLabels:  
14         k8s:io.kubernetes.pod.namespace: ingress-nginx  
15         k8s:app.kubernetes.io/name: ingress-nginx
```

Protecting Infra Components

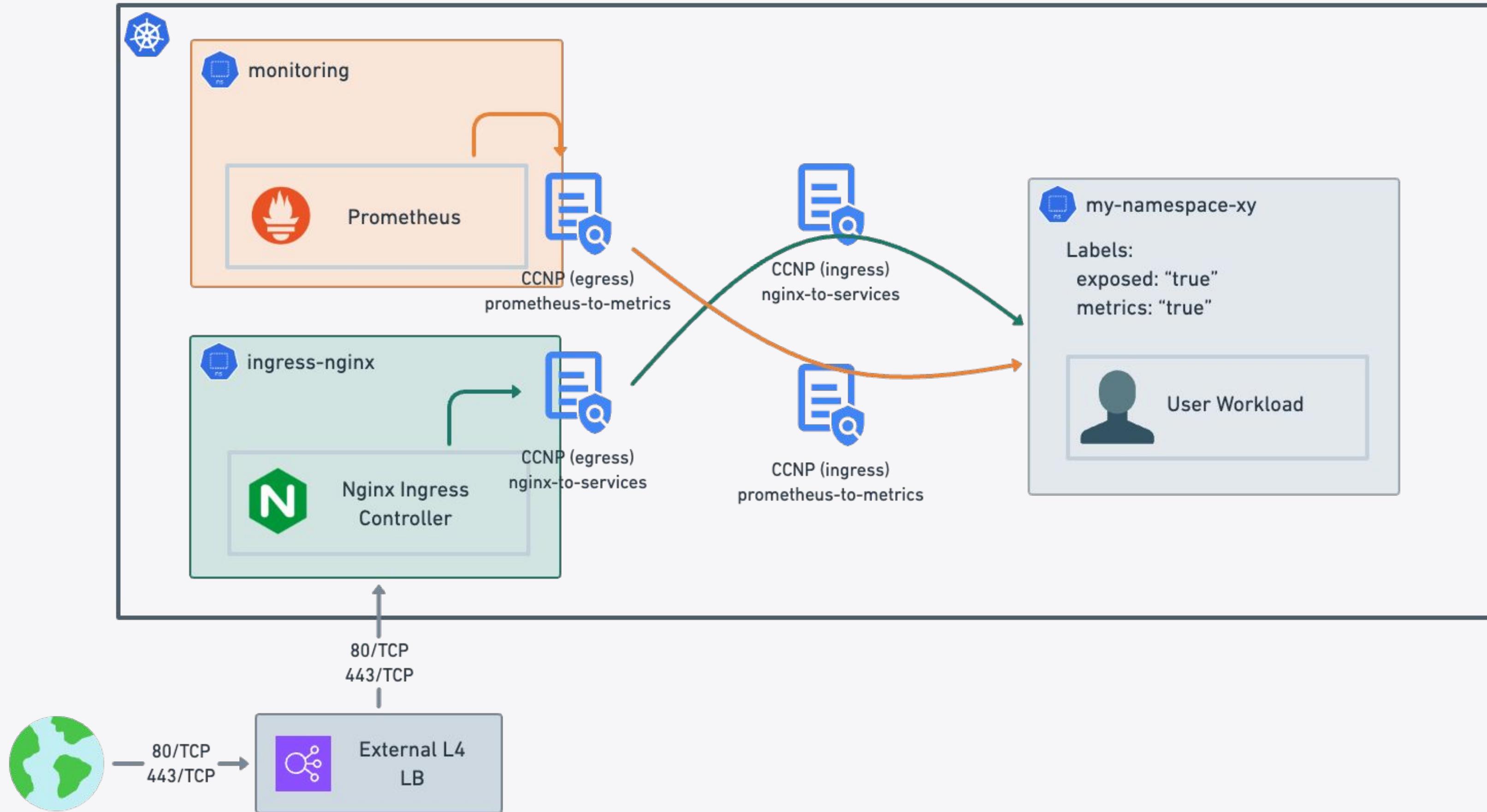


Prometheus:

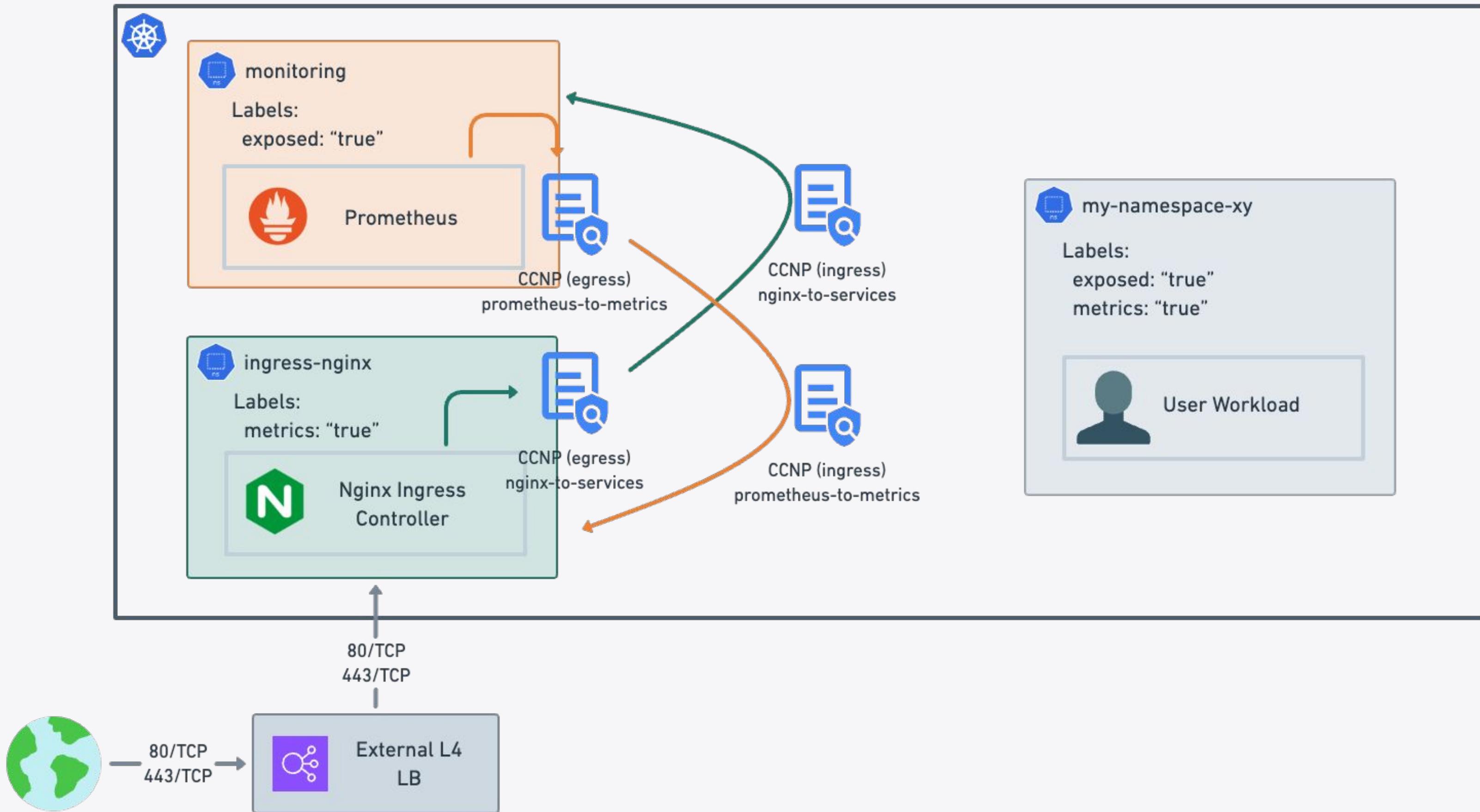
```
1 —  
2 apiVersion: "cilium.io/v2"  
3 kind: CiliumClusterwideNetworkPolicy  
4 metadata:  
5   name: prometheus-to-metrics  
6 spec:  
7   description: Allow Prometheus to scrape metrics in NS with metrics label  
8   endpointSelector:  
9     matchLabels:  
10       k8s:app.kubernetes.io/name: prometheus  
11       k8s:io.kubernetes.pod.namespace: monitoring  
12 egress:  
13 - toEndpoints:  
14   - matchLabels:  
15     k8s:io.cilium.k8s.namespace.labels.metrics: "true"
```

```
1 —  
2 apiVersion: "cilium.io/v2"  
3 kind: CiliumClusterwideNetworkPolicy  
4 metadata:  
5   name: global-infra-prometheus  
6 spec:  
7   description: Allow connections from Prometheus to labeled namespaces  
8   endpointSelector:  
9     matchLabels:  
10       k8s:io.cilium.k8s.namespace.labels.metrics: "true"  
11   ingress:  
12     - fromEndpoints:  
13       - matchLabels:  
14         k8s:io.kubernetes.pod.namespace: monitoring  
15         k8s:app.kubernetes.io/name: prometheus
```

Protecting Infra Components



Protecting Infra Components

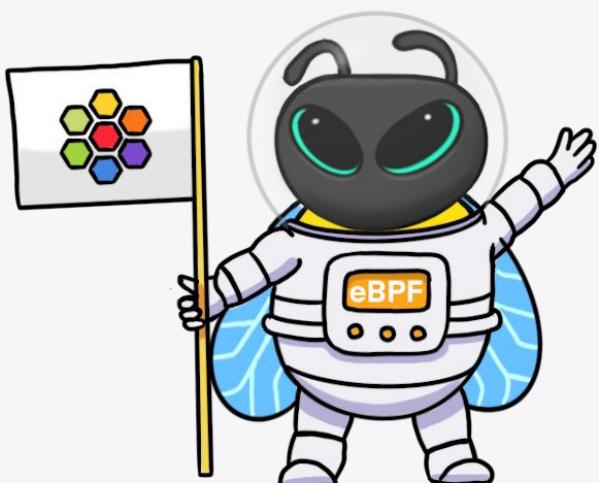


Default Rule Set for User Workloads

Default Rule Set for User Workloads

- Transition from “Coarse” (L3, L4 max.) to “Fine-Grained” (L4, DNS, L7) NS policies:
 - **Internal:** Start with allowing all ingress/egress communication within a namespace
 - **Egress:** Define to where this application needs to connect to
 - **Ingress:** Define from where something needs to connect to this application
- Use Cilium Hubble observability data to identify further communication flows

```
1 hubble observe -t policy-verdict -f --namespace <ns-name> [--verdict DROPPED]
```



Default Rule Set for User Workloads

- Every CiliumClusterwideNetworkPolicy with spec.endpointSelector: {} applies to all namespaces! In our case:
 - Egress to Core-DNS
 - Metric scrapes from Prometheus (with NS label metrics: "true")
 - Ingress from Nginx ingress controller (with NS label exposed: "true")
- Our CiliumClusterwideNetworkPolicy setup creates a deny-all behavior for all namespaces – even within the namespace
 - Add at least a ciliumNetworkPolicy to allow all namespace internal traffic



Default Rule Set for User Workloads

```
1 # -- User workload NS creation
2 --
3 apiVersion: v1
4 kind: Namespace
5 metadata:
6   name: my-namespace-xy
7   # Customized labels
8   # - Add exposed="true" label in case Nginx should expose something from this NS
9   # - Add metrics="true" label in case Prometheus should scrape metrics from this NS
10  labels:
11    exposed: "true"
12    metrics: "true"
```



Default Rule Set for User Workloads

```
1 # -- Default rule
2 --
3 apiVersion: "cilium.io/v2"
4 kind: CiliumNetworkPolicy
5 metadata:
6   name: allow-within-namespace
7   namespace: my-namespace-xy
8 spec:
9   description: Allow NS internal traffic, block everything else
10  endpointSelector: {}
11  ingress:
12    - fromEndpoints:
13      - {}
14  egress:
15    - toEndpoints:
16      - {}
```


ISOVALENT

Thank you!



Additional Resources

- Cluster setup and presented policies:
 - <https://github.com/PhilipSchmid/cilium-netpol-demo>
- KubeCon 22: Adopting Network Policies in Highly Secure Environments
 - <https://www.youtube.com/watch?v=yikVhGM2ye8>
- 10 Free hands-on labs about Cilium Network Policies:
 - <https://isovalent.com/learning-tracks/#securityProfessionals>
- More “Securing Network with Cilium” examples:
 - <https://docs.cilium.io/en/stable/security/>
- Cilium Network Policy documentation entry point:
 - <https://docs.cilium.io/en/stable/security/policy/>
- **The cluster setup, all shown and even more policies can be found on Github:** <https://github.com/PhilipSchmid/cilium-netpol-demo>

