# PADOK
## SECURITY

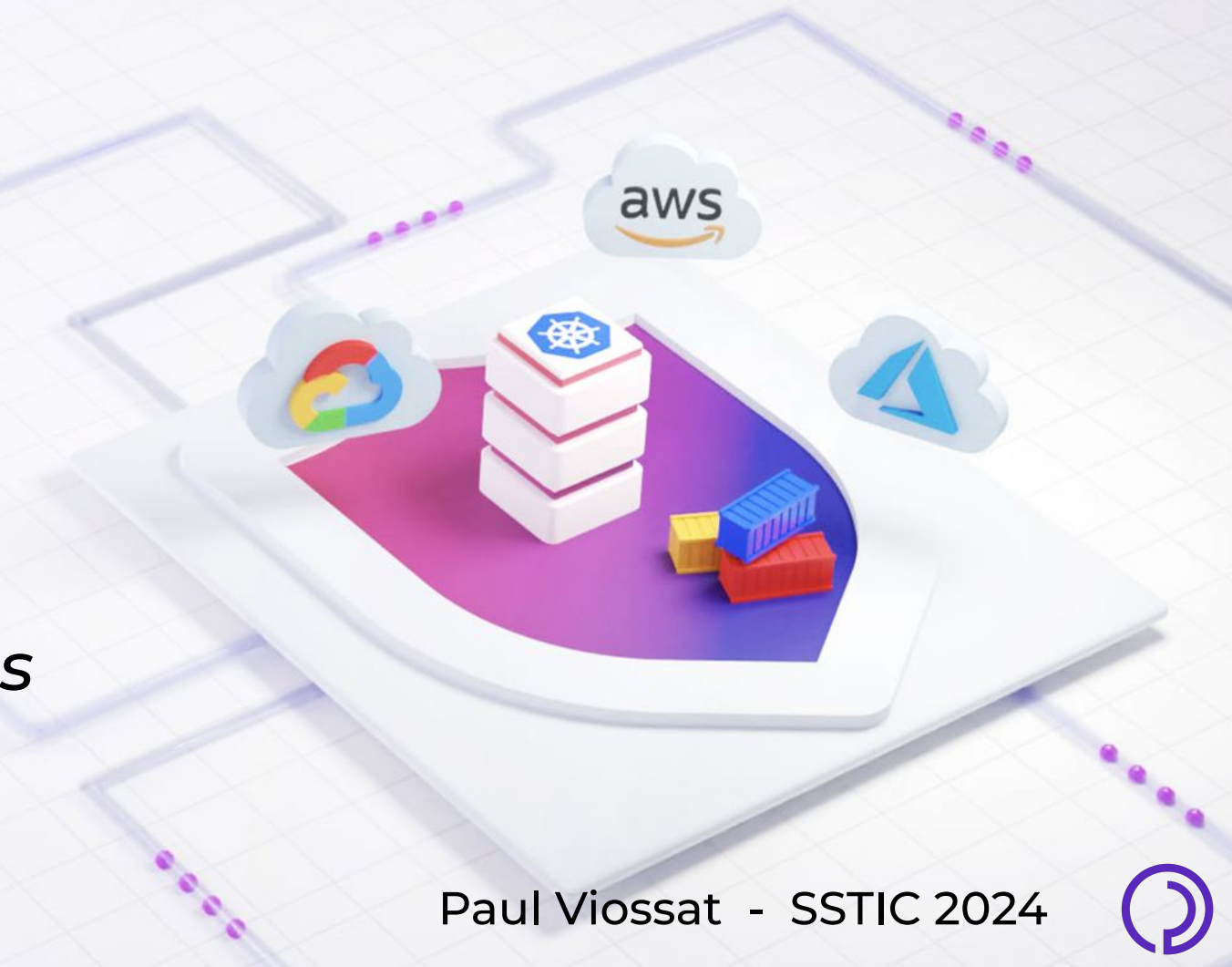# Getting ahead of the schedule:
*manipulating the Kubernetes scheduler to perform lateral movement in a cluster*

Paul Viossat  -  SSTIC 2024

# Agenda

# Considered threat

INTRODUCTION

# Considered threat

# Considered threat

→ Escape based on:
  ○ Configuration of pods
  ○ Kernel exploits

→ Can be found in **CI/CD jobs** when performing DockerInDocker tasks (performance testing, e2e testing, image building, etc...

attacker

CI/CD job

node

A

B

service-account-A

service-account-B

container escape

kubelet
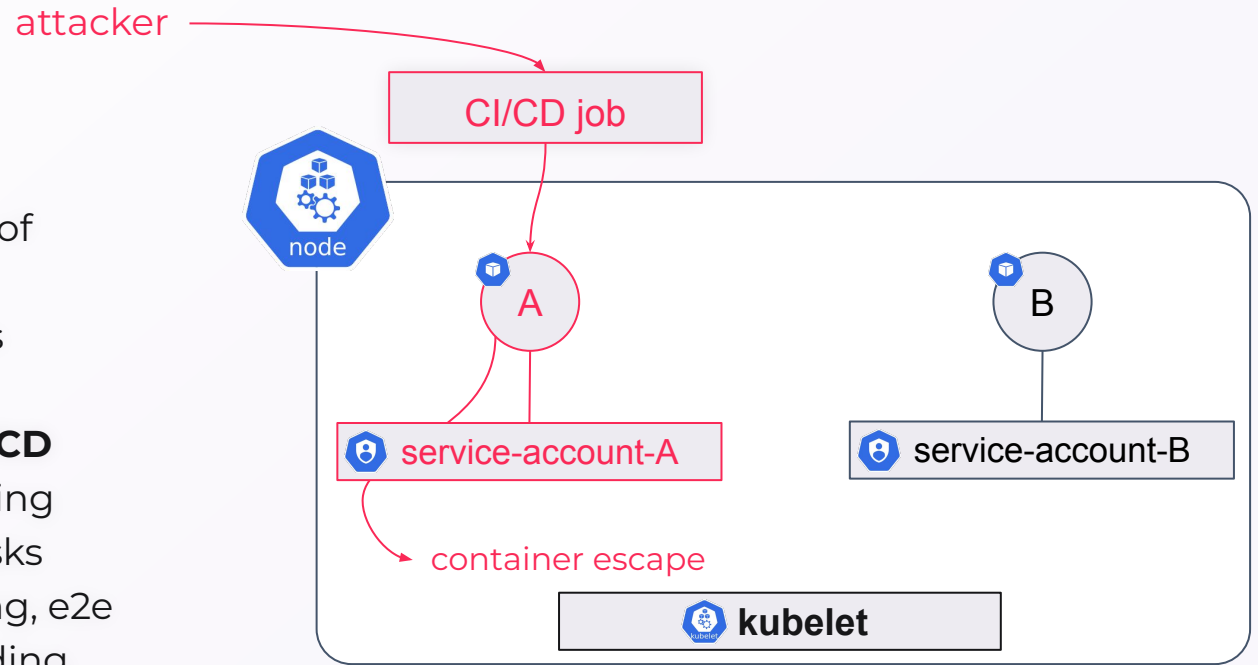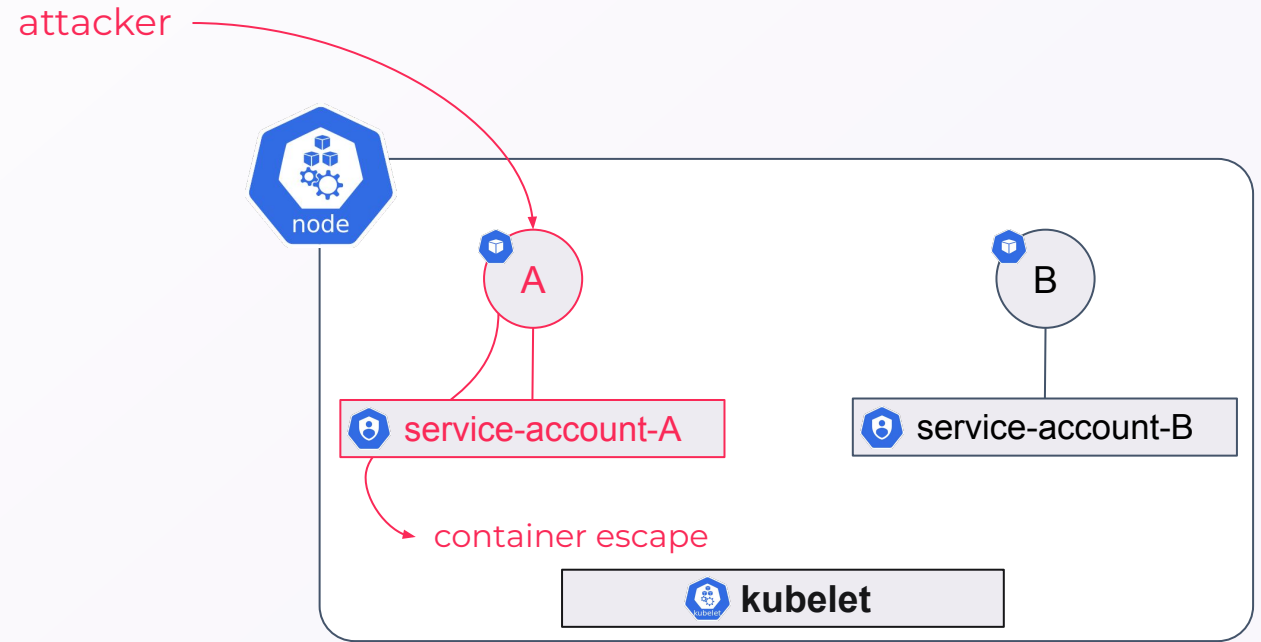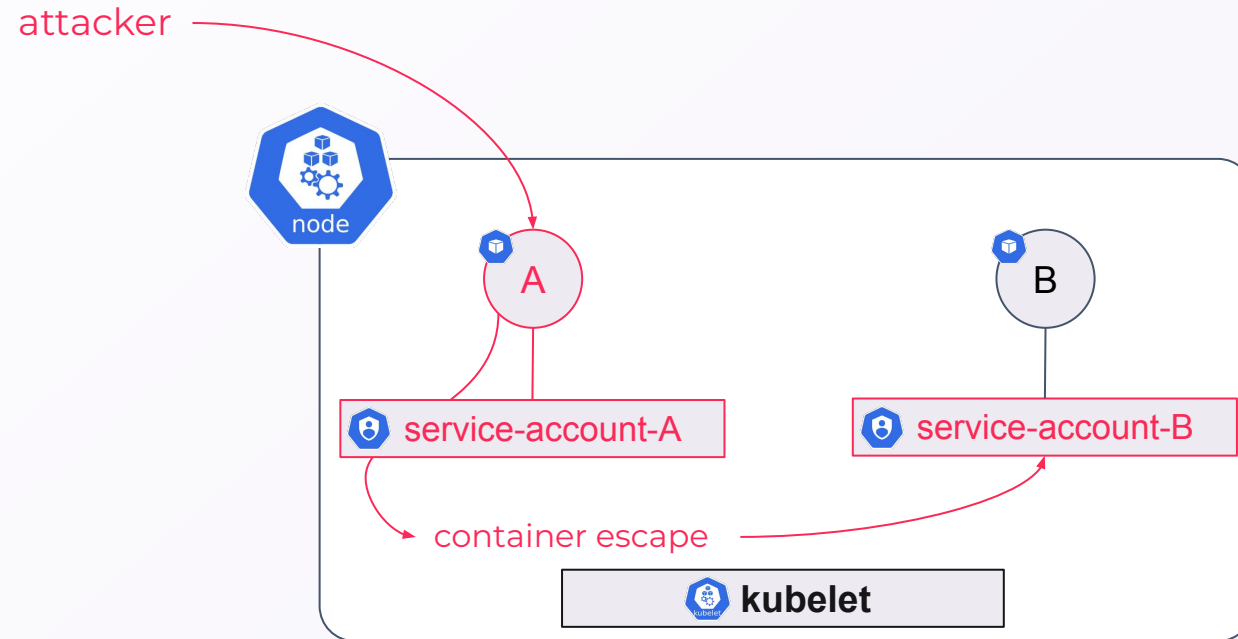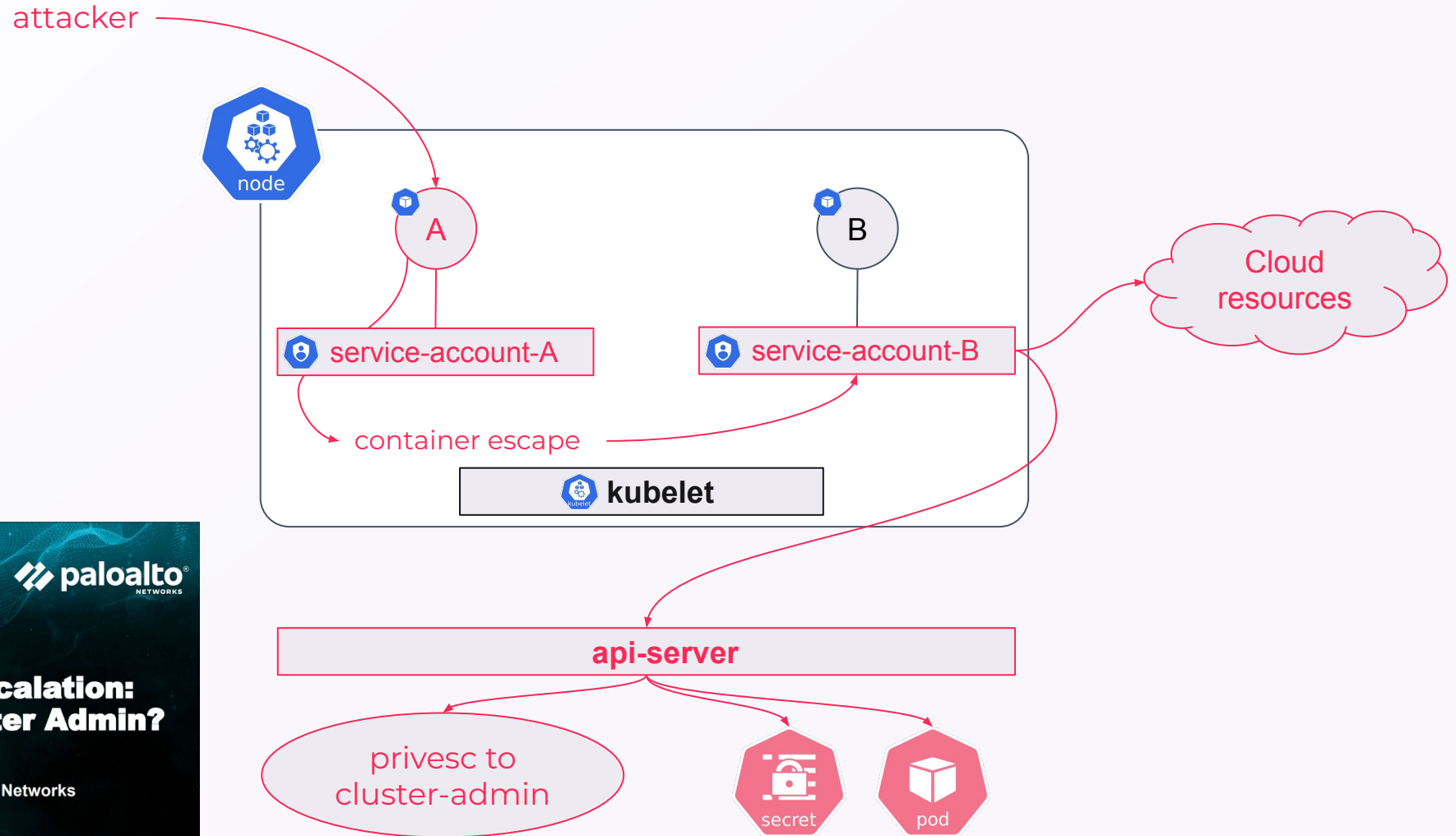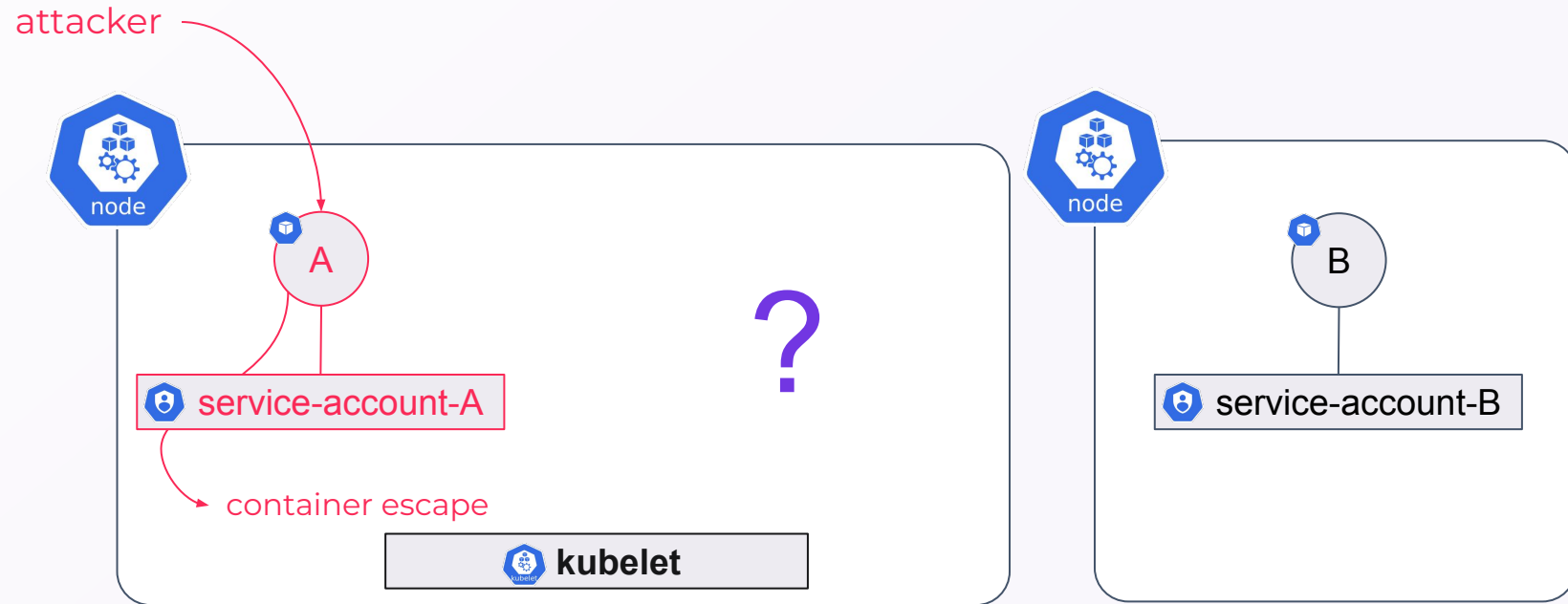
# Considered threat
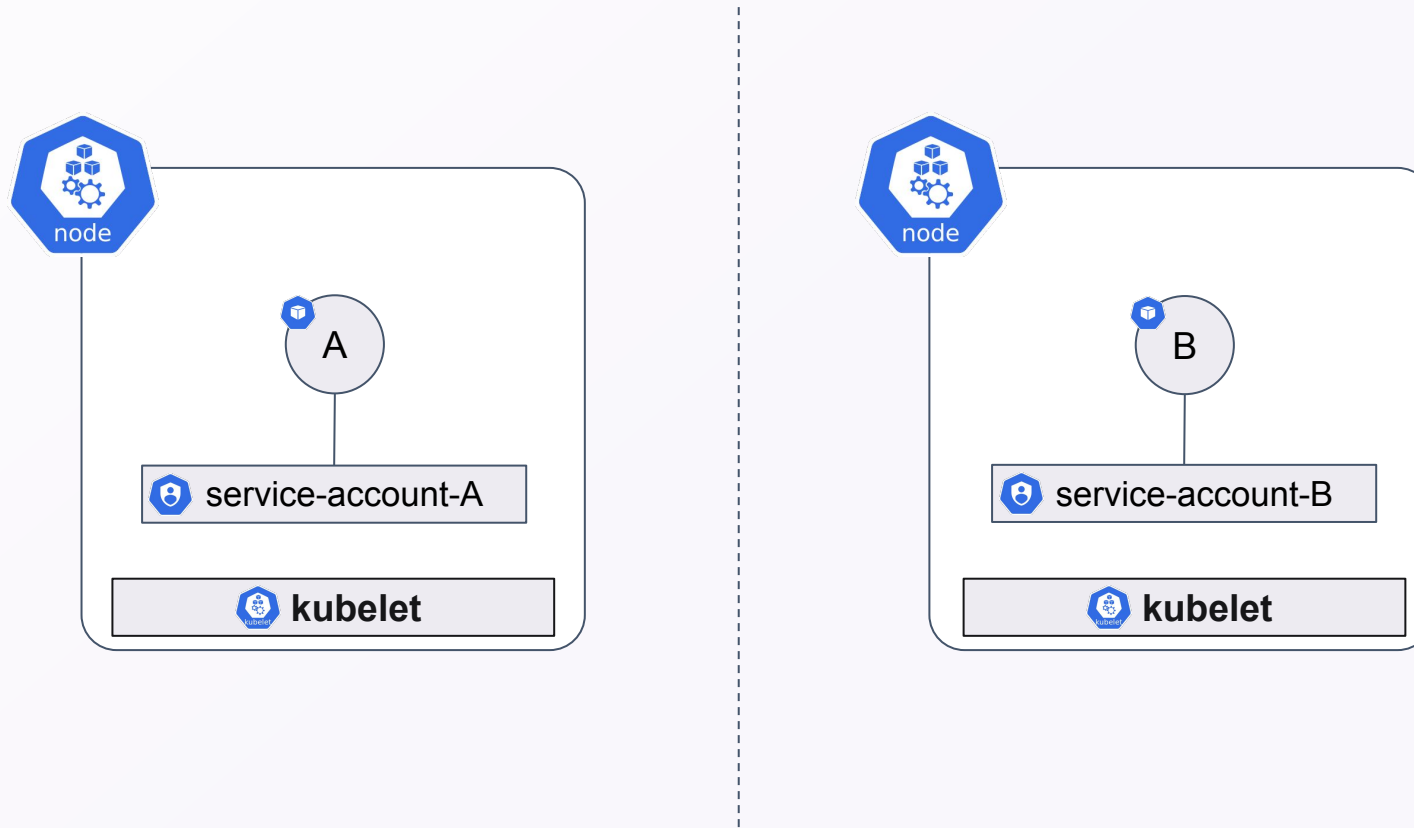
# Considered threat

# Considered threat

# Considered threat



→ As defender can I predict which service account can be compromised ?

→ As attacker how can I can access other service account ?

# Considered threat



**Workload node isolation**

# Considered threat
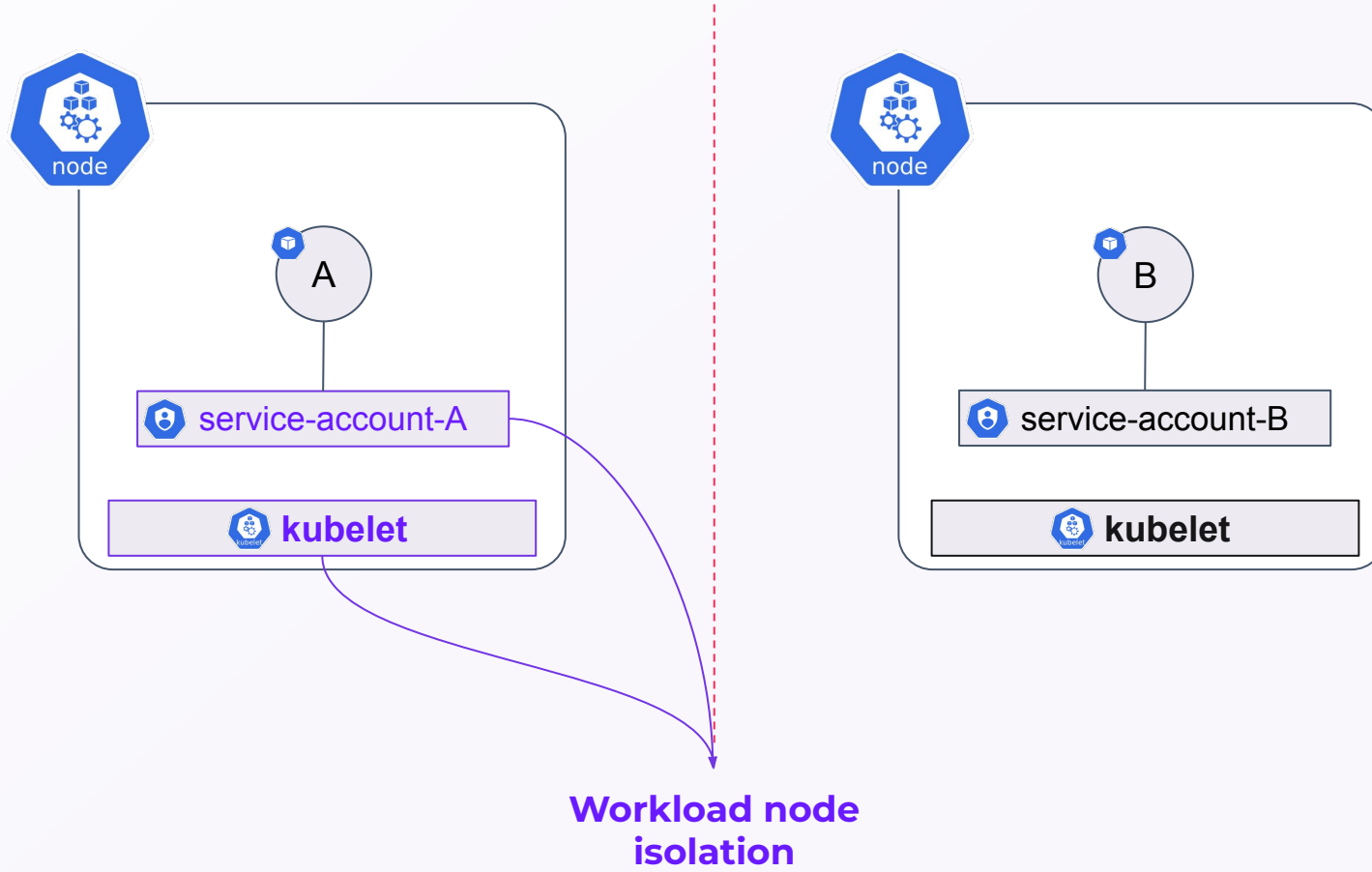


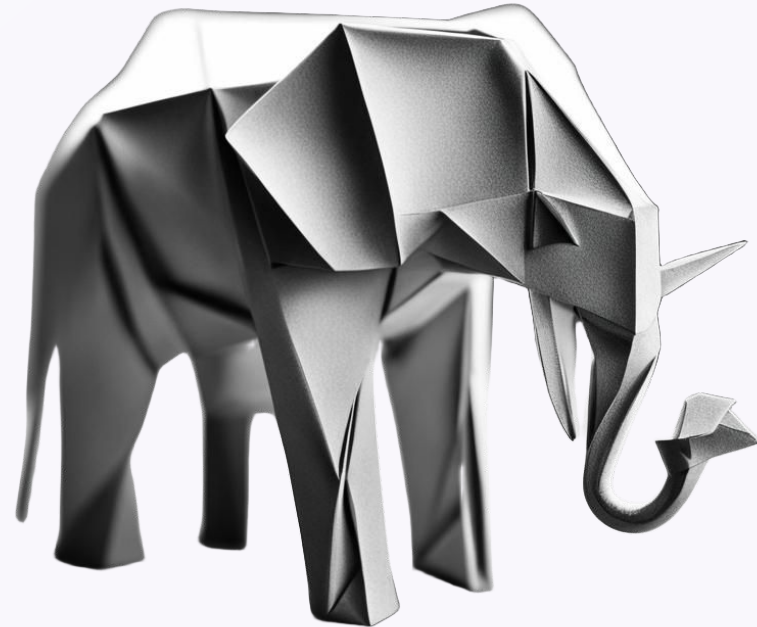**Workload node isolation**

# The elephant in the room

→ As defender can I **predict which service account can be compromised** ?

→ As attacker how can I can **access other service accounts** ?

→ How **robust** is workload node isolation ?

## The Kubernetes Scheduler

# Agenda

# Scheduler overview

```
apiVersion: v1
kind: Pod
metadata:
  name: nginx
spec:
  nodeName: null
  containers:
  - name: nginx
    image: nginx:1.14.2
    ports:
    - containerPort: 80
```

pod ↔ node

pod

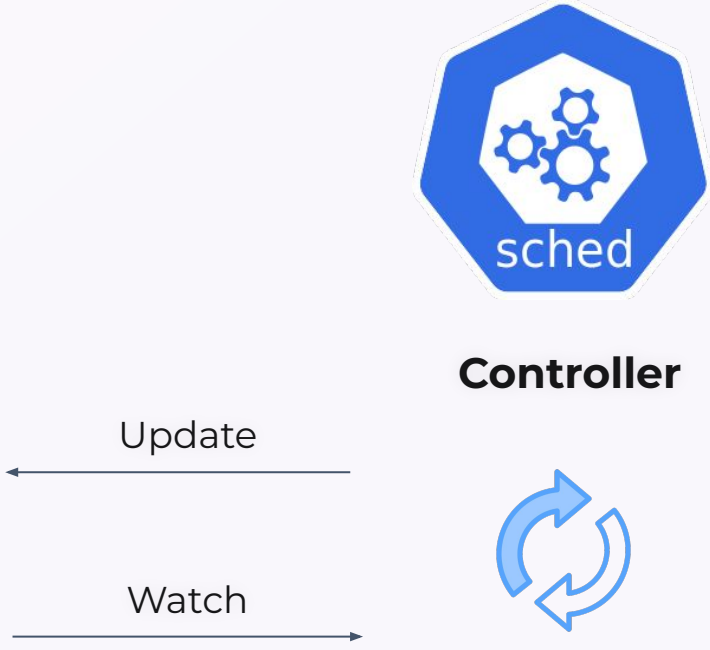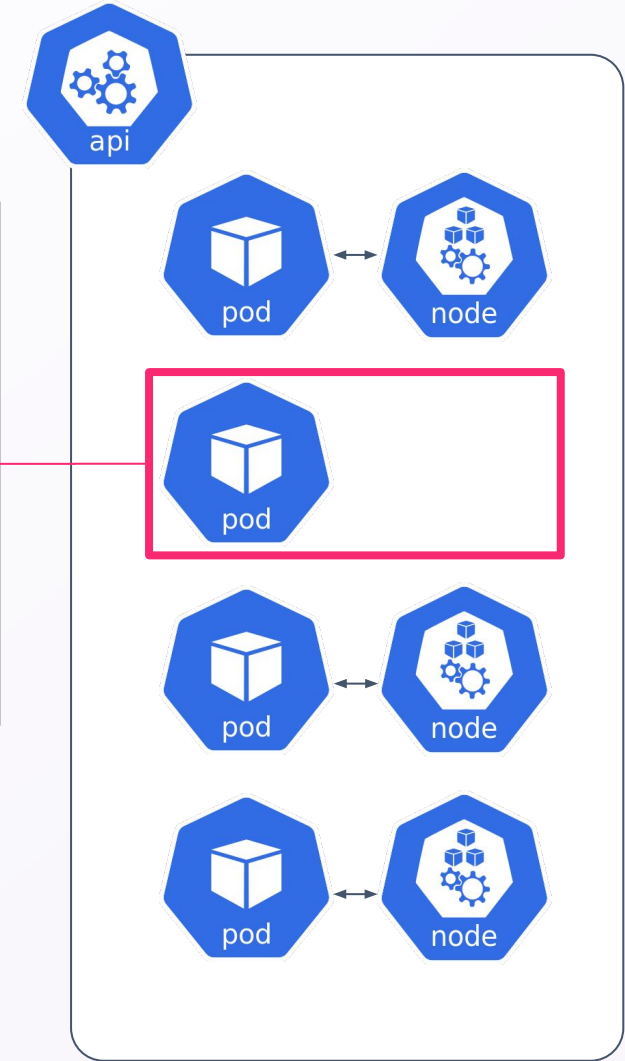pod ↔ node

pod ↔ node

**Controller**

Update

Watch

# Scheduler overview

```
apiVersion: v1
kind: Pod
metadata:
  name: nginx
spec:
  nodeName: my-node
  containers:
  - name: nginx
    image: nginx:1.14.2
    ports:
    - containerPort: 80
```

api

pod — node

pod — node

pod — node

pod — node

Update

Watch

sched

**Controller**

# Scheduler

Find *feasible*
nodes for pod

Selecting the "best" node
among feasible nodes

Update the pod object in
Kubernetes API

Filtering

Scoring

Binding

Plugins

Plugins

DefaultBinder

# Filtering

# Filtering

# Filtering

# Filtering

# Filtering



feasible nodes for the pod

# TaintTolerations



**taints**

**tolerations**

→ Nodes **repel** Pods

# NodeAffinities



**labels**

**affinities /
nodeSelector**

Pods **chose** Nodes

# NodeAffinities

security=high

**labels**

A  security=high

B  security=low

**affinities / nodeSelector**

→ Pods **chose** Nodes

# Domain of feasibility

→ Domain of feasibility = results of plugins **NodeAffinites** + results of plugins **TaintToleration**



domain of feasibility
of pod in the cluster

feasible nodes for pod

# Domain of feasibility: isolated pods

# Domain of feasibility: possibly non-isolated pods

# Isolation anti-patterns: only taints

# Isolation anti-patterns: only affinities and node selector

# Isolation pattern: using both

# Scoring

|  | B | D | E |
|---|---|---|---|
| **node** | | | |
| Plugin A | 40 | 0 | 10 |
| Plugin D | 10 | 100 | 10 |
| Plugin E | 30 | 20 | 80 |
| | 80 | 120 | 100 |

# Scoring

→ Mainly based on **Node status**

→ Node status is **fully editable** by the kubelet account

→ We are **almost certain** to be able to **attract Pod** when the compromised node is in the **domain of feasibility**

# Binding: Node authorizer

CREATE token
my-service-account

my-node

Pod A

default

Pod B

my-service-account

spec.nodeName set in binding phase

spec.serviceAccountName set in pod spec

# Binding: Node authorizer

CREATE token
my-service-account

my-node

Pod A

default

Pod B

my-service-account

# Agenda

# Modify domain of feasibility using the kubelet account

**node**

**labels**

security=high

**taints**

security=high:NoExecute

**kubelet**

**A**

**affinities / nodeSelector**

security=high

**tolerations**

security=high:NoExecute

# Modify domain of feasibility using the kubelet account

# Modify domain of feasibility using the kubelet account

**labels**

security=high

**taints**

security=high:NoExecute

🔷 **kubelet**

**affinities / nodeSelector**

A

security=high

**tolerations**

security=high:NoExecute

→ **NodeRestriction** admission plugin limits the editable labels

→ It defines a **blacklist** that contains, among others:
   ○ **node-restriction.kubernetes.io/**

→ It prevents also modifying the taints

# Modify domain of feasibility using the kubelet account

**labels**

security=high

**taints**

security=high:NoExecute

**kubelet**

**affinities / nodeSelector**

A

security=high

**tolerations**

security=high:NoExecute

# Modify domain of feasibility using the kubelet account

```
apiVersion: v1
kind: Node
spec:
  unschedulable: true
```

node

KEEP OUT

kubelet

cordon

# Modify domain of feasibility using the kubelet account

```yaml
apiVersion: v1
kind: Node
spec:
  unschedulable: true
  taints:
    - effect: NoSchedule
      key:node.kubernetes.io/unschedulable
```



cordon

# Modify domain of feasibility using the kubelet account

# Modify domain of feasibility using the kubelet account

attacker

A

cordon

KEEP OUT

container escape

kubelet

node

B

service-account-B

node

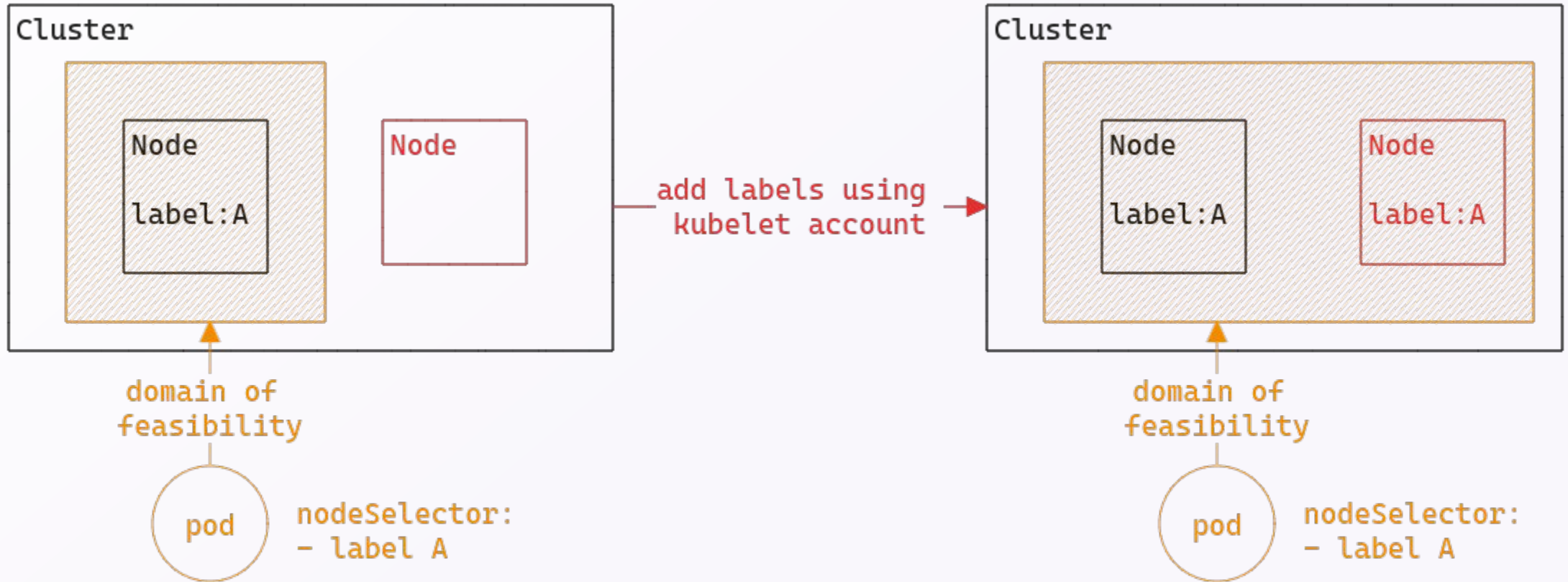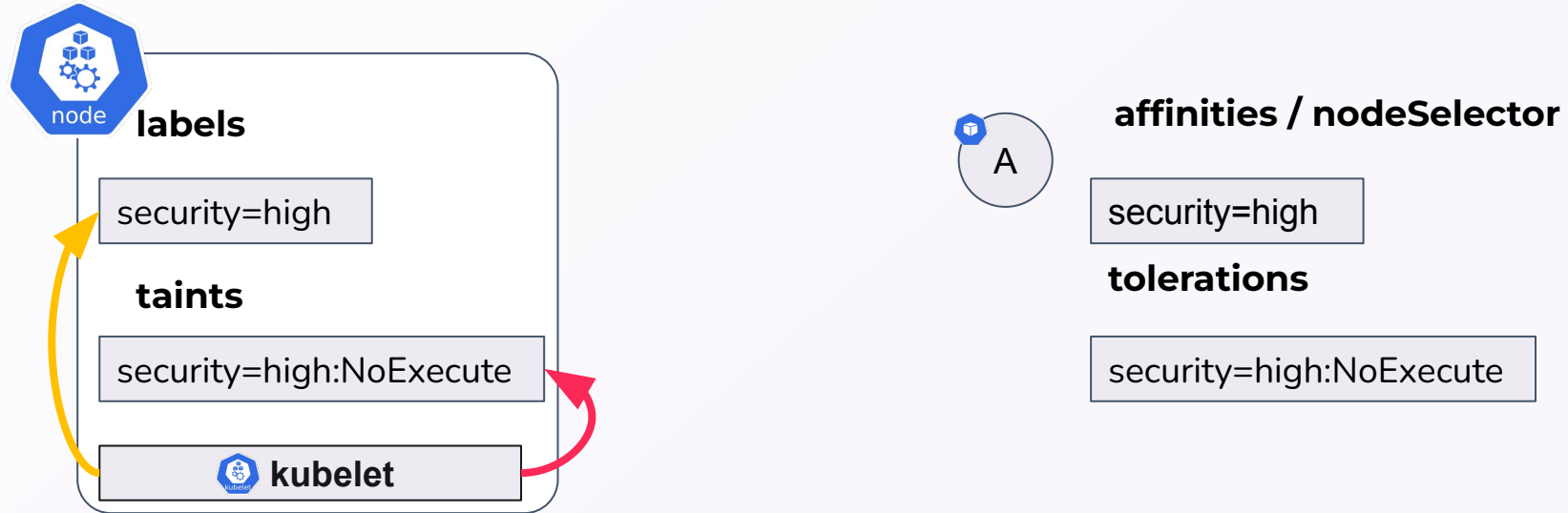# Modify domain of feasibility using the kubelet account

# Modify domain of feasibility using the kubelet account

# Modify domain of feasibility using the kubelet account
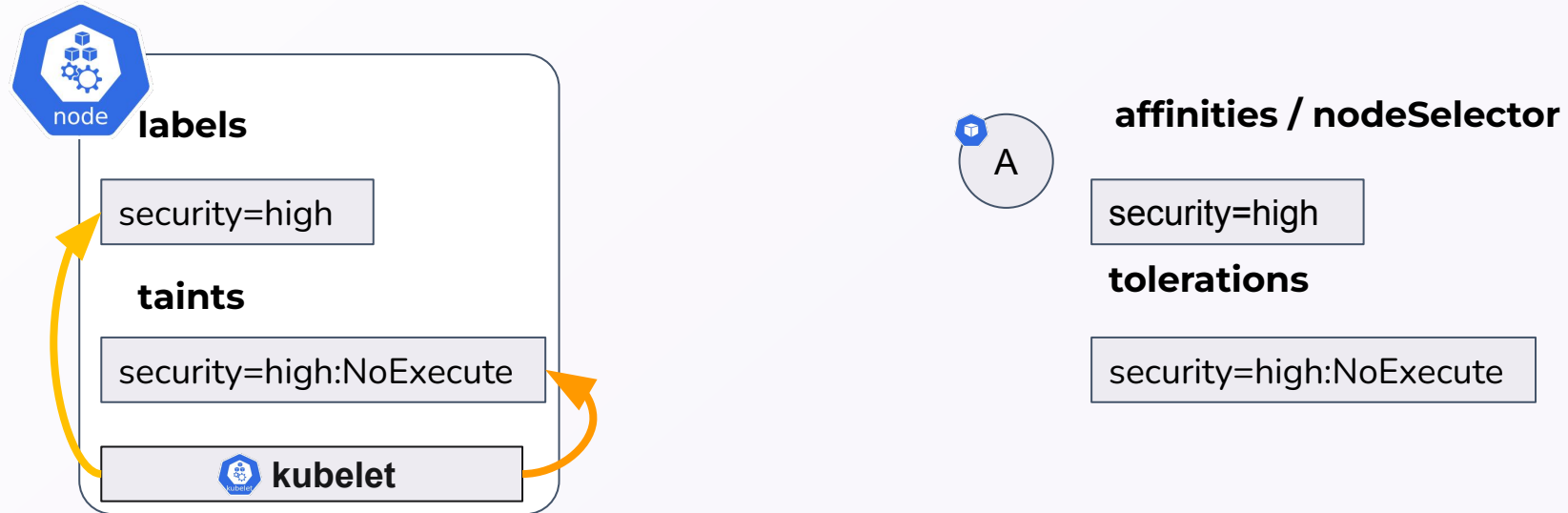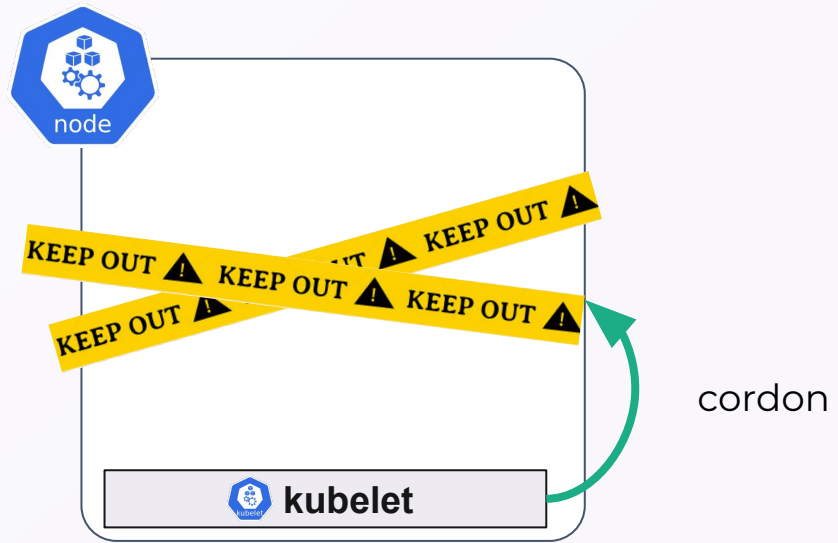
# Modify domain of feasibility using the kubelet account

# Wrap up

To implement node isolation:

→ Use **node affinities + taints and tolerations**

→ Use **NodeRestriction** admission plugin (enabled by default on AWS, GCP and Azure)

→ Use protected labels with prefix: **node-restriction.kubernetes.io/**

In case of compromise:

→ The domain of feasibility of the **vulnerable "entry" pod** can be compromised only using kubelet **privileges**

→ If **no isolation** all nodes can be compromised

# Agenda

# Filtering

**labels**

security=high

**taints**

security=high:NoExecute

**kubelet**

**affinities / nodeSelector**

A

security=high

**tolerations**

security=high:NoExecute

add
tolerations

service-account

PATCH POD

# Scheduler

```
┌─────────────┐      ┌─────────────┐      ┌─────────────┐      ┌─────────────┐
│     Pod     │ ───> │  Filtering  │ ───> │   Scoring   │ ───> │   Binding   │
│  creation   │      │             │      │             │      │             │
└─────────────┘      └─────────────┘      └─────────────┘      └─────────────┘
```

# Scheduler

```
┌──────────┐        ┌──────────┐              ┌──────────┐        ┌──────────┐
│   Pod    │───────▶│          │              │          │───────▶│          │
│ creation │        │ Filtering│              │ Scoring  │        │ Binding  │
│          │        │          │              │          │        │          │
└──────────┘        └──────────┘              └──────────┘        └──────────┘
                         │
                         ▼
                  unschedulable
```

# Using patch pod right

# Using patch pod right

replicaset controller

attacker

node

A

service-account-A

B

service-account-B

container escape

node

watch

C

labels

admin-account

remove
replicaset
tracking labels

# Using patch pod right



→ When reaching max capacity **NodeResourcesFit** plugin filters out nodes !

# Using patch pod right

# DEMO

attacker

container escape

A

B

aws-node

replicaset controller

watch

C

labels

admin-account

security=low:NoExecute

security=low

DamonSet
AWS CNI +
ANNOTATE_POD_IP =
PATCH_POD

security=high:NoExecute

security=high

# DEMO

# Some thoughts on node isolation on managed Kubernetes

→ Cloud providers do not seems to consider node isolation flaws as security issues
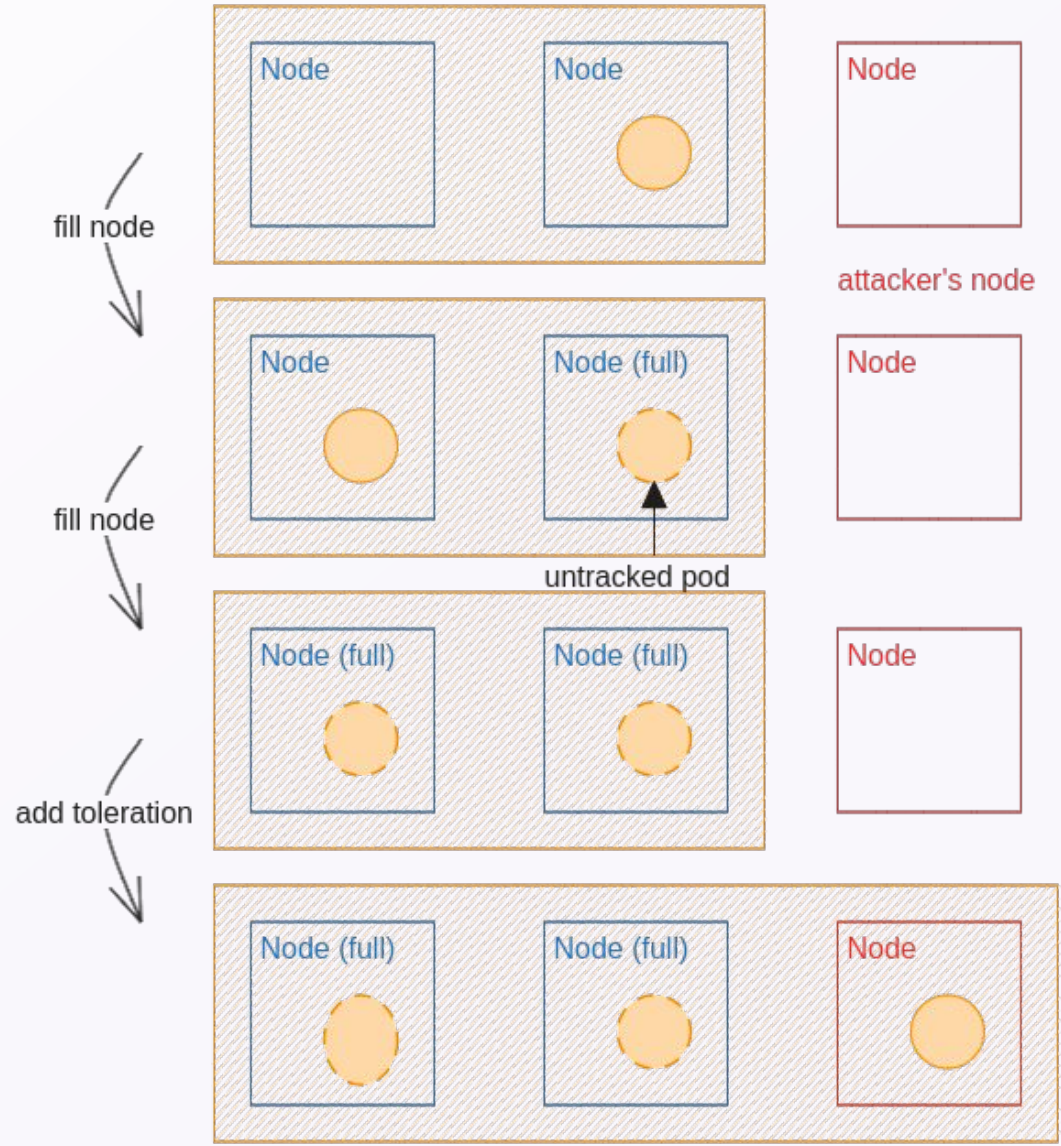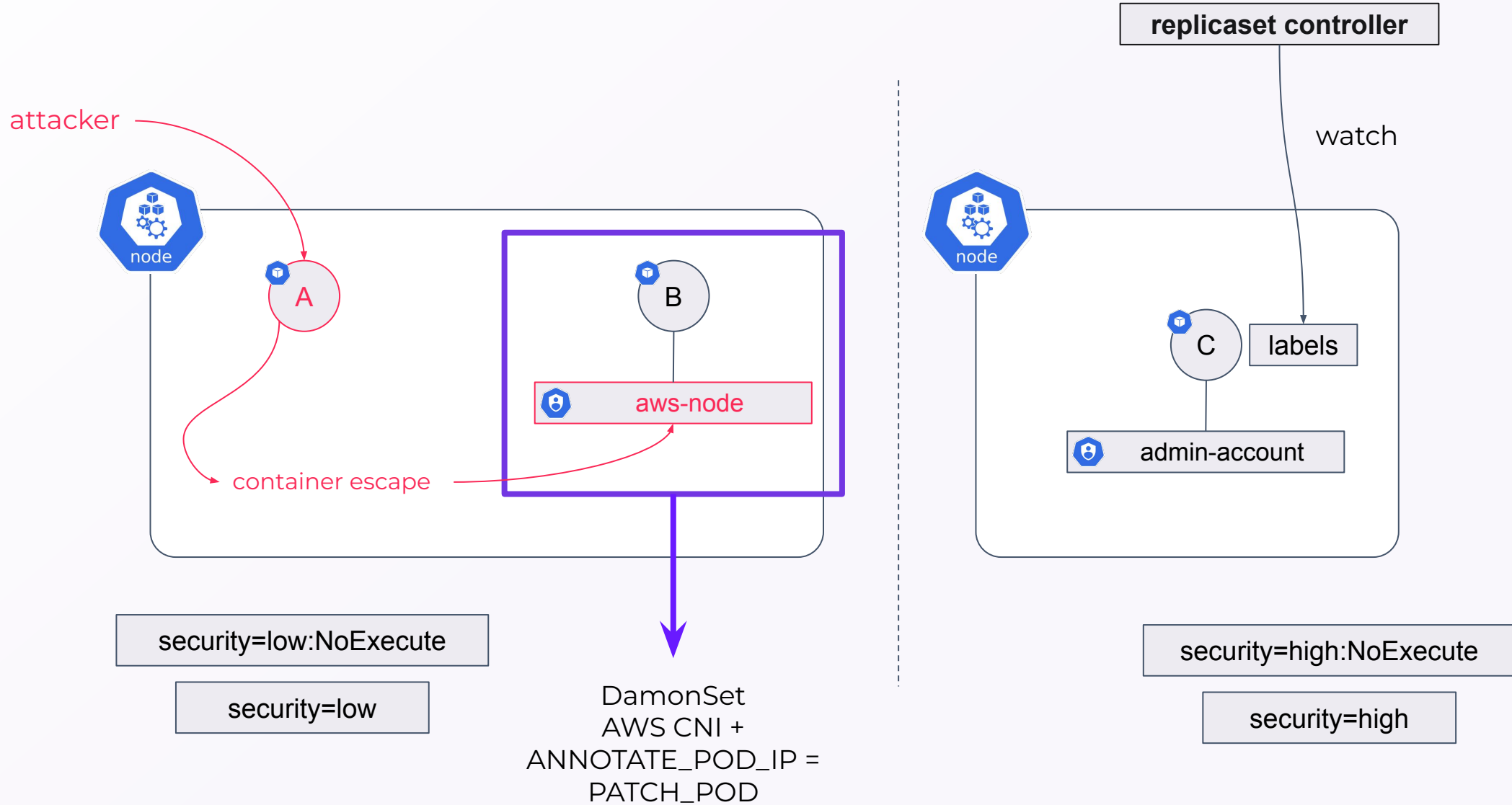
→ It may be better to create multiple clusters

**ANNOTATE_POD_IP** (v1.9.3+)

Type: Boolean as a String

Default: `false`

Setting `ANNOTATE_POD_IP` to `true` will allow IPAMD to add an annotation `vpc.amazonaws.com/pod-ips` to the pod with pod IP.

There is a known issue with kubelet taking time to update `Pod.Status.PodIP` leading to calico being blocked on programming the policy. Setting `ANNOTATE_POD_IP` to `true` will enable AWS VPC CNI plugin to add Pod IP as an annotation to the pod spec to address this race condition.

To annotate the pod with pod IP, you will have to add `patch` permission for pods resource in aws-node clusterrole. You can use the below command -

```
cat << EOF > append.yaml
- apiGroups:
  - ""
  resources:
  - pods
  verbs:
  - patch
EOF
```

```
kubectl apply -f <(cat <(kubectl get clusterrole aws-node -o yaml) append.yaml)
```

NOTE: Adding `patch` permissions to the `aws-node` Daemonset increases the security scope for the plugin, so add this permission only after performing a proper security assessment of the tradeoffs.

# Agenda

# Wrap up

To implement node isolation:

- → Use **node affinities + taints and tolerations**

- → Use **NodeRestriction** admission plugin (enabled by default on AWS, GCP and Azure)

- → Use protected labels with prefix: **node-restriction.kubernetes.io/**

In case of compromise:

- → The domain of feasibility of the **vulnerable "entry" pod** can be compromised without **privileges**

- → If **no isolation** all nodes can be compromised

- → **PATCH Pod** permission on a Pod is enough to move a pod around if not using protected labels
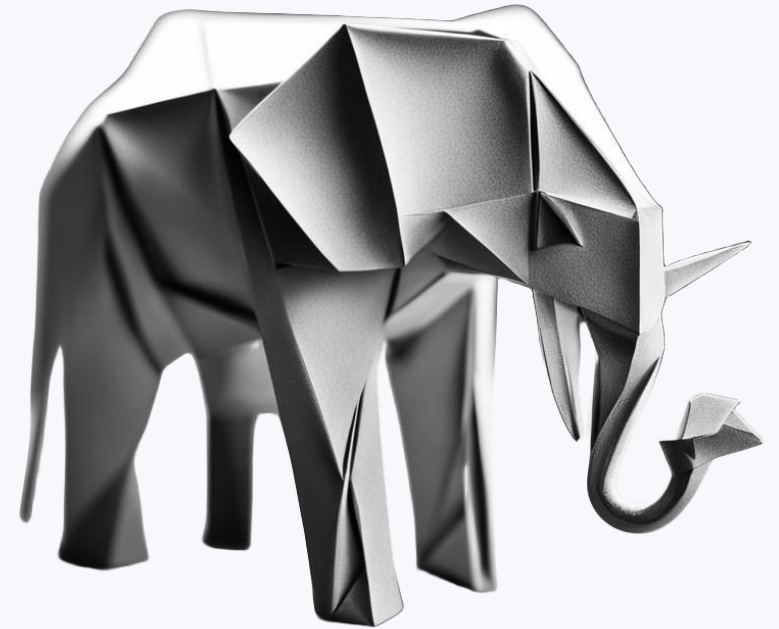
# Ready for next elephant in the room ?

→ How nodes are created and deleted in cloud managed environment ?



https://security.padok.fr/blog

**Cloud controller manager**

# Thank you !

PADOK