# GPU sharing done right

**HashiCorp**
an IBM Company

**Adrian Todorov**
Staff Solutions Architect
HashiCorp

**Scott McAllister**
Principal Developer Advocate
Depot

**HashiCorp**
an IBM Company

# Adrian Todorov

Staff Solutions Architect @ HashiCorp,

ex-SRE, ex-developer
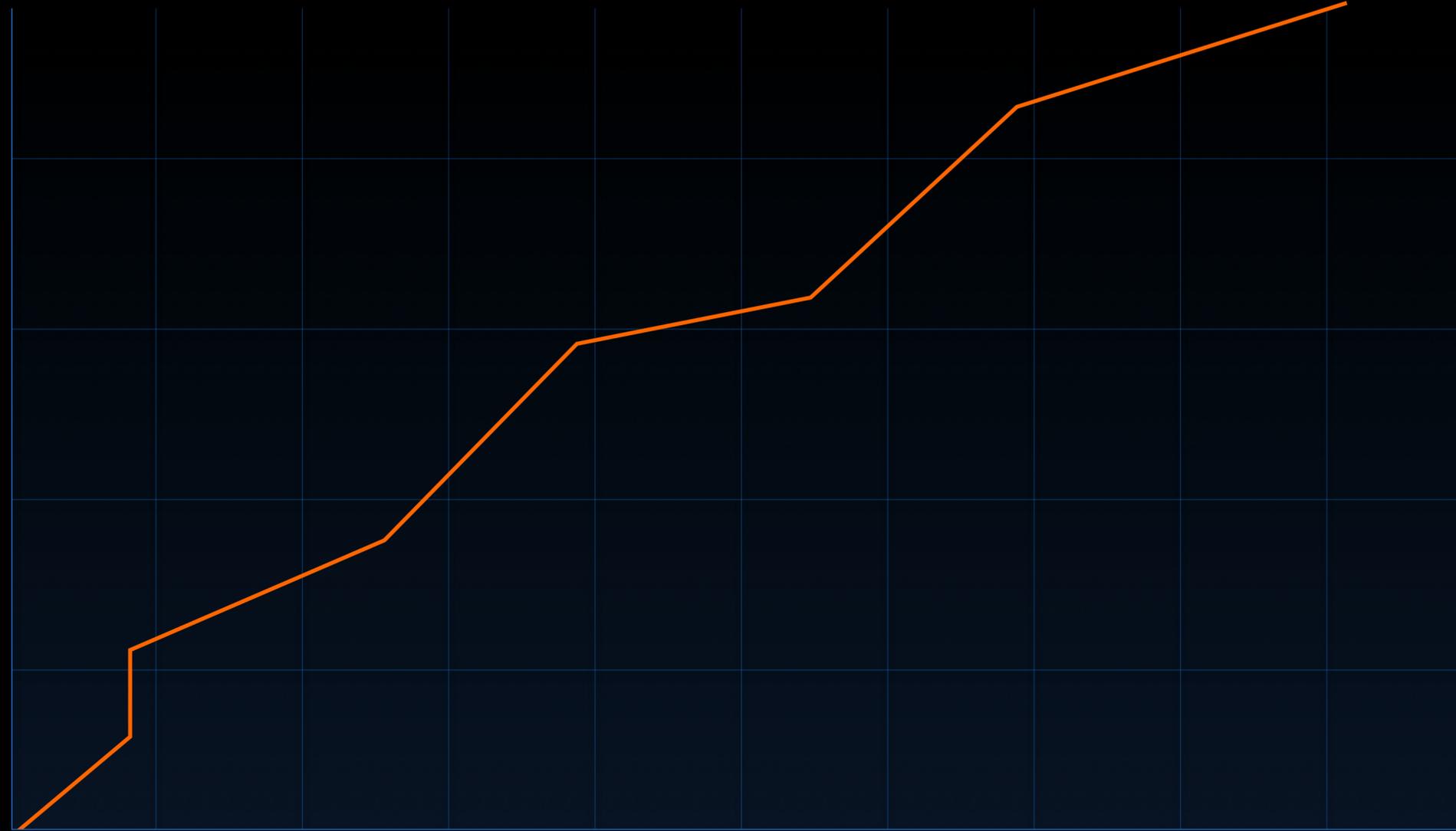
atodorov.me

# Traditional clusters are like gamers.

# Traditional clusters don't like to share GPU.

# What's happening



AI use cases

GPU demand
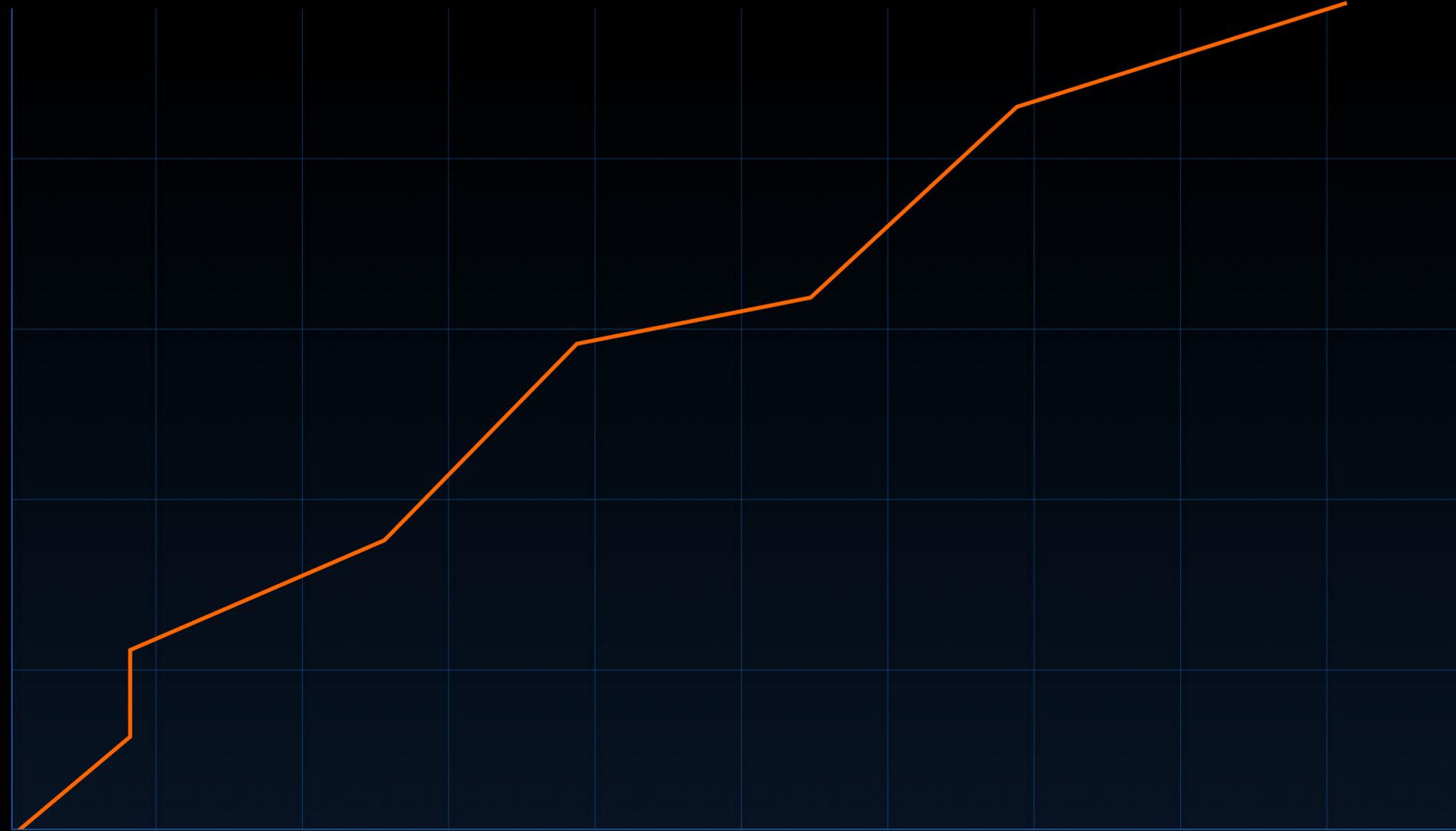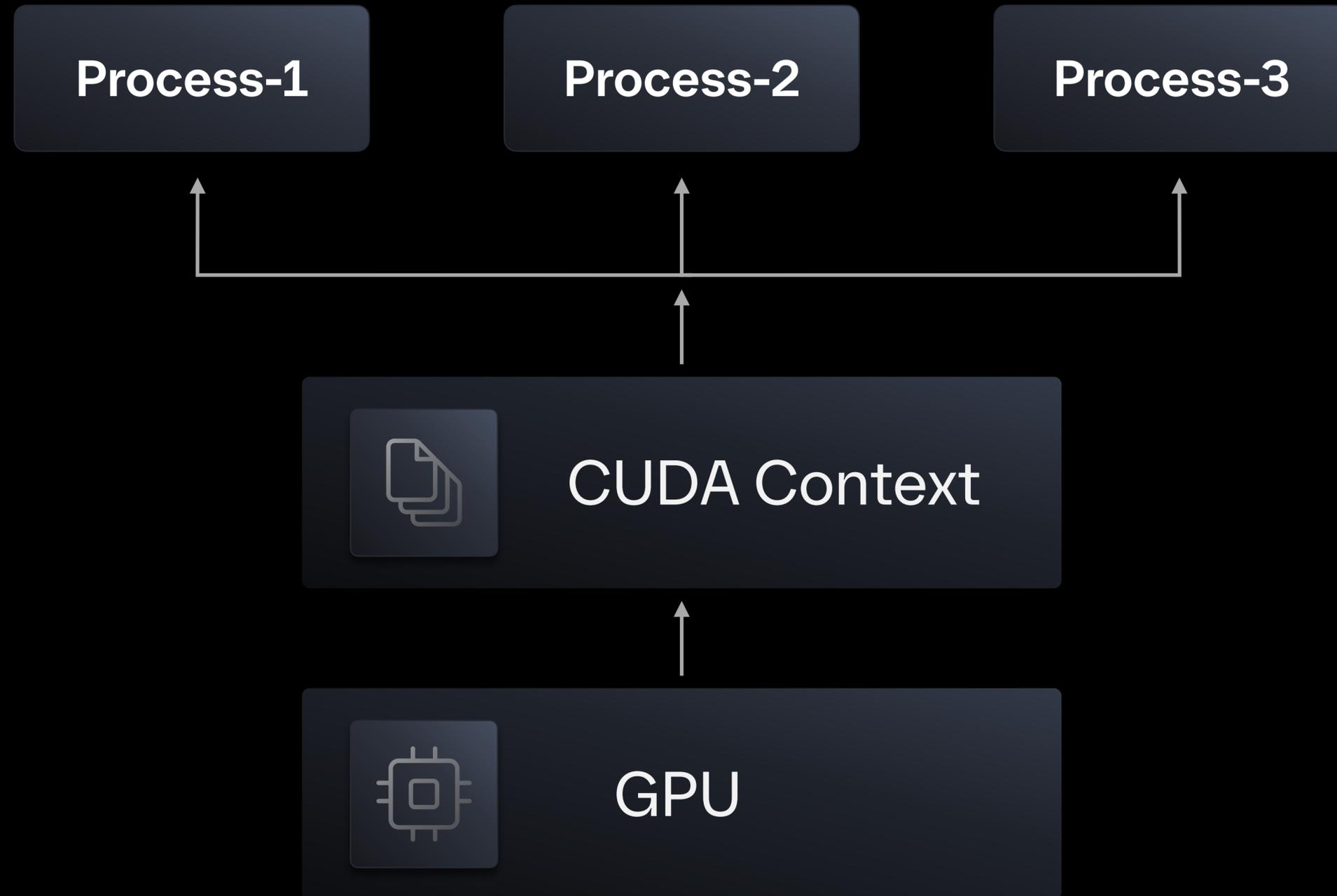
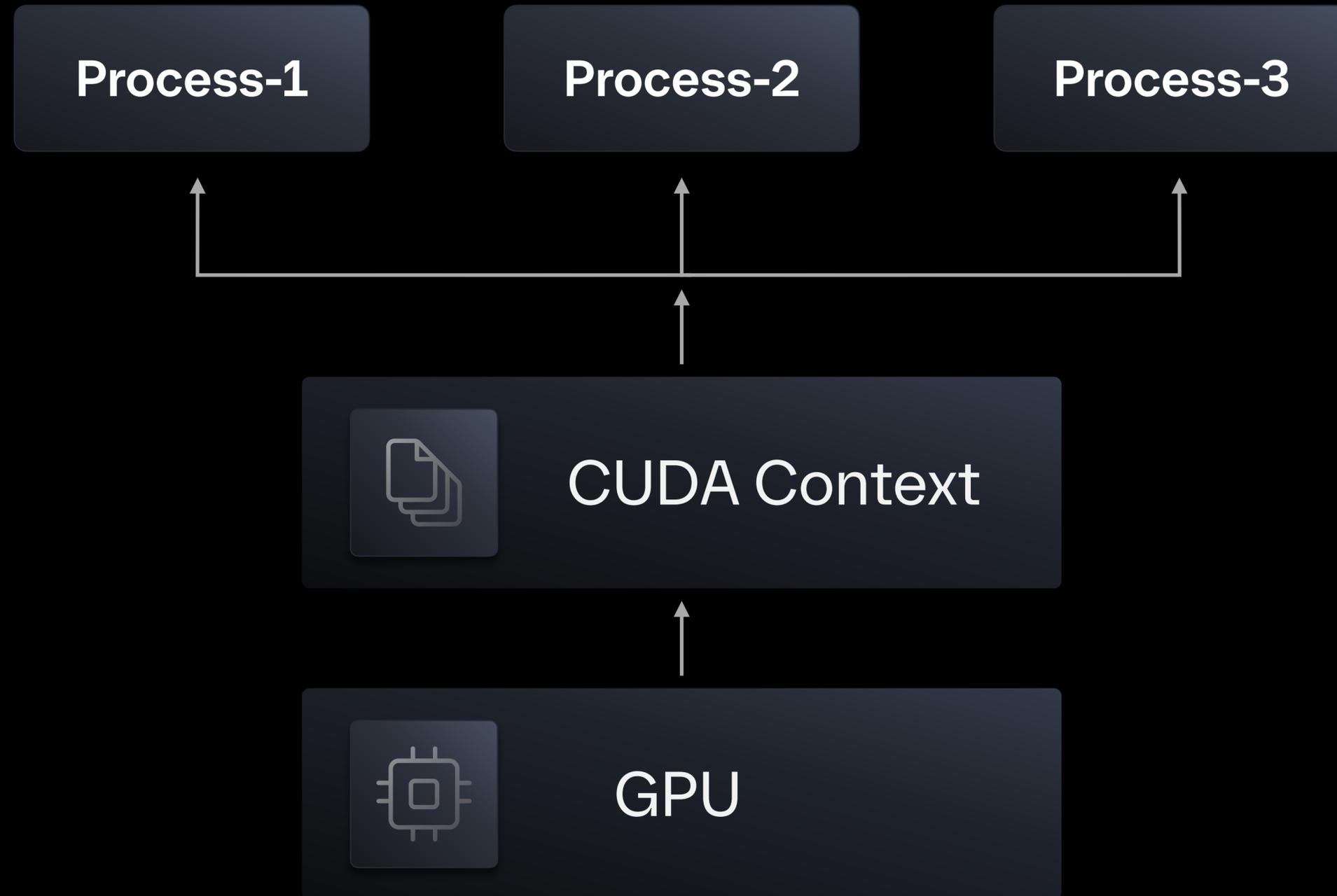Interest in AI

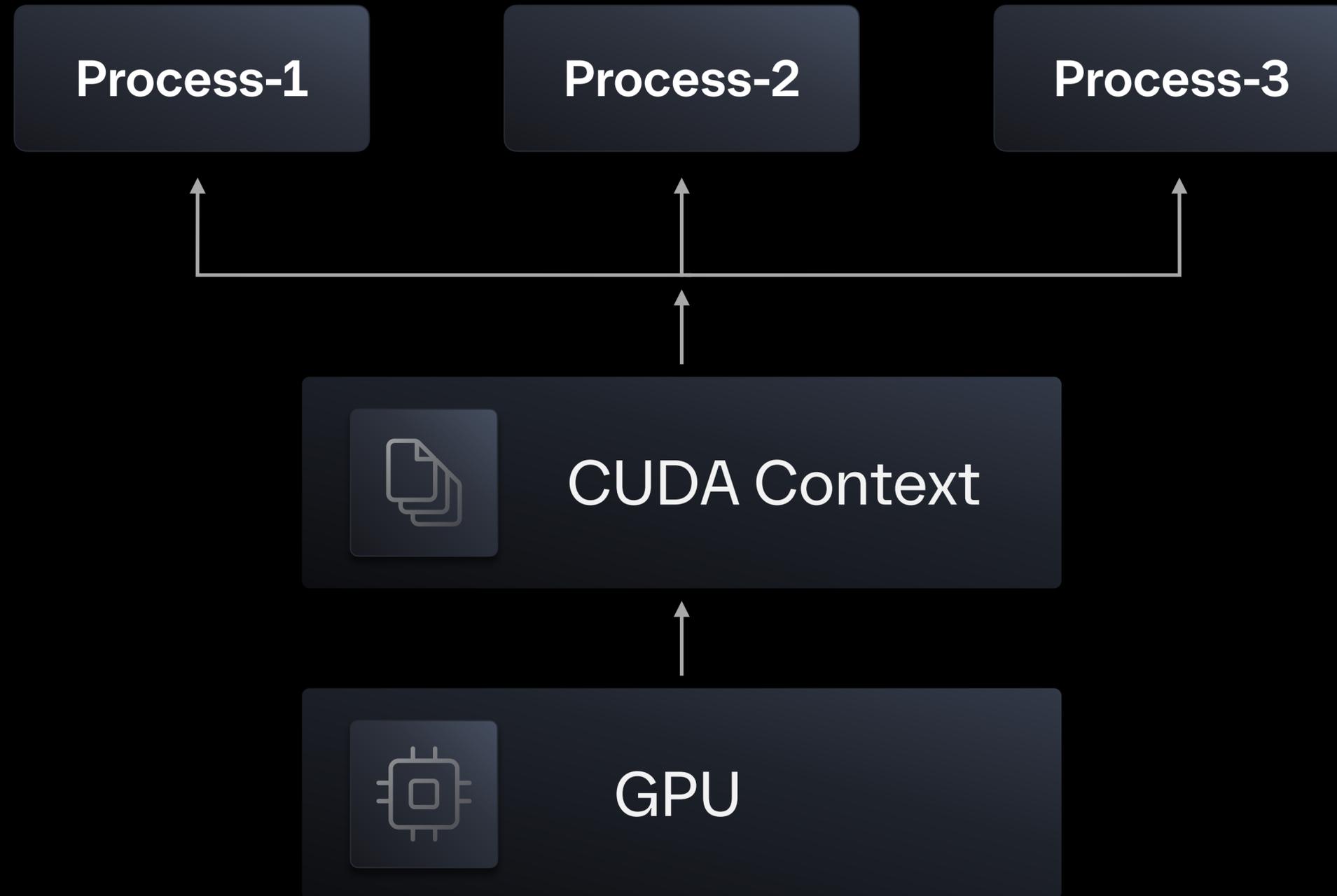# What's happening



AI use cases / GPU demand — Interest in AI

# Multi-Process Service (MPS)

# Multi-Process Service (MPS)

```
Process-1          Process-2          Process-3



              CUDA Context



                  GPU
```

# Multi-Process Service (MPS)

# Multi-Instance GPUs (MIG)



GPU

GPU instance 0
User

GPU instance 1
User

GPU instance 2
User

GPU instance 3
User

GPU instance 4
User

# Multi-Instance GPUs (MIG)

GPU

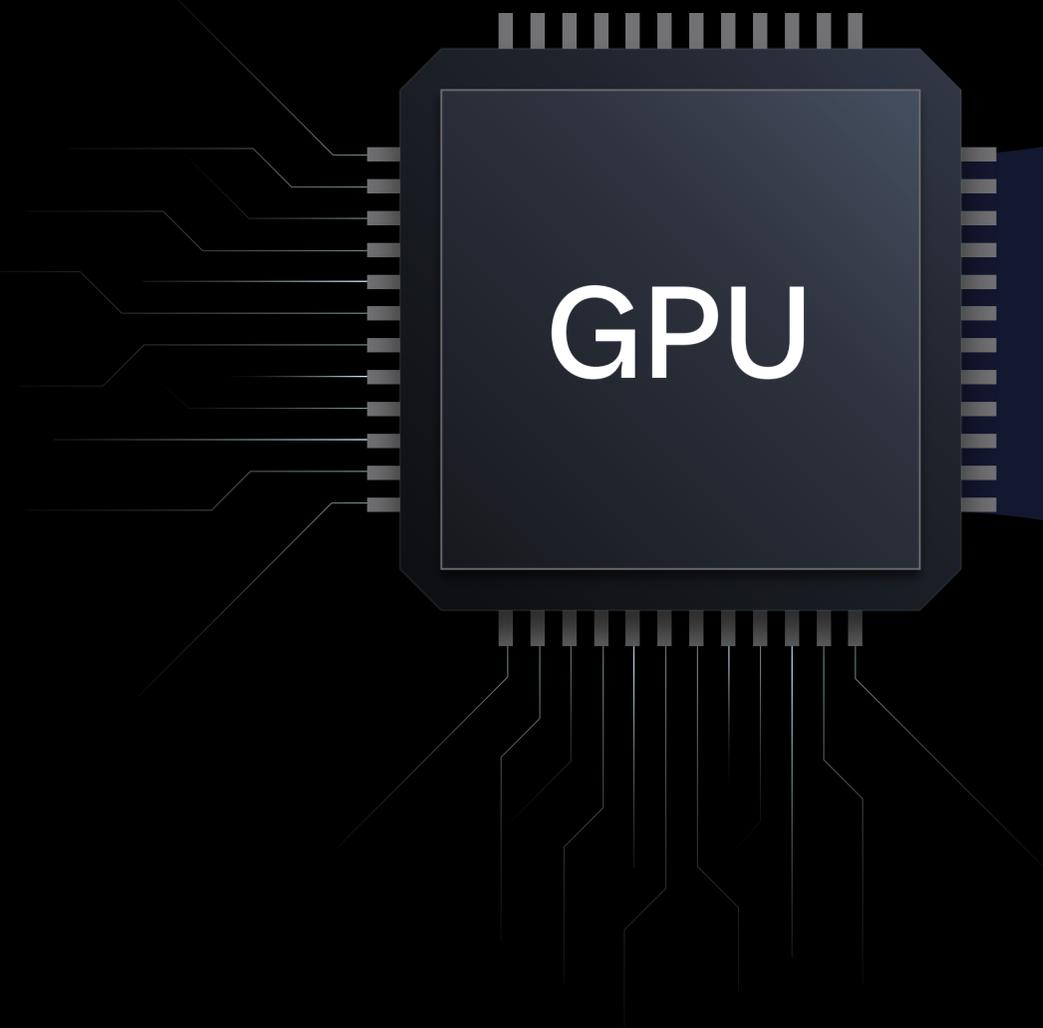GPU instance 0
User

GPU instance 1
User

GPU instance 2
User

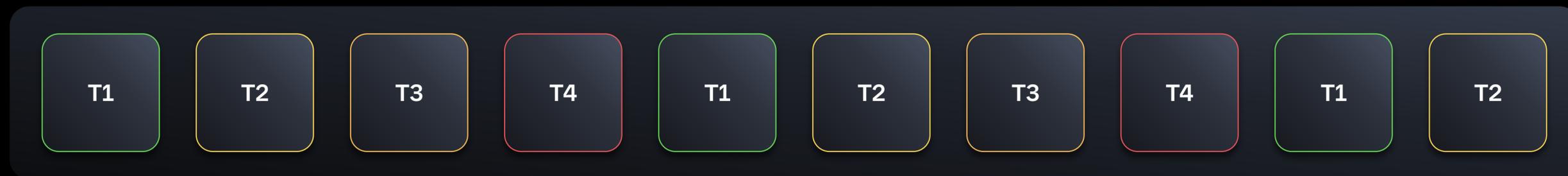GPU instance 3
User

GPU instance 4
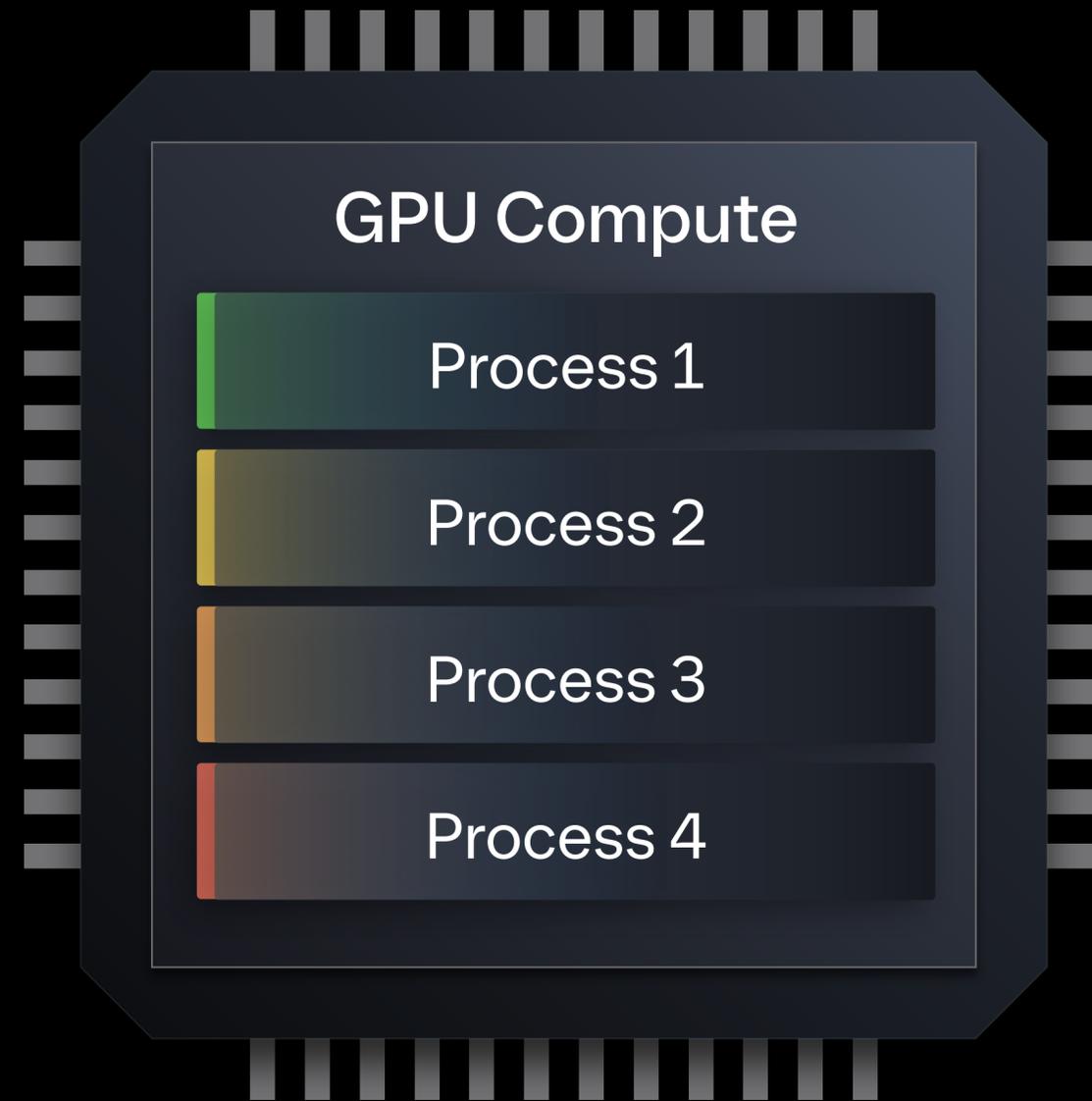User

# Multi-Instance GPUs (MIG)

# Multi-Instance GPUs (MIG)

# Time-slicing

(Or: custom scheduling)

**GPU Compute**

Process 1

Process 2

Process 3

Process 4

| T1 | T2 | T3 | T4 | T1 | T2 | T3 | T4 | T1 | T2 |

Time slice

# Time-slicing

(Or: custom scheduling)

GPU Compute

Process 1

Process 2

Process 3

Process 4

| T1 | T2 | T3 | T4 | T1 | T2 | T3 | T4 | T1 | T2 |

Time slice

# Time-slicing

(Or: custom scheduling)

GPU Compute

Process 1

Process 2

Process 3

Process 4

| T1 | T2 | T3 | T4 | T1 | T2 | T3 | T4 | T1 | T2 |

Time slice

# Time-slicing

(Or: custom scheduling)



GPU Compute

Process 1

Process 2

Process 3

Process 4

T1 T2 T3 T4 T1 T2 T3 T4 T1 T2

Time slice

# Third-Party Scheduling Tools

# Third-Party Scheduling Tools

# Per-team GPU cluster

# Per-team GPU cluster

# Per-team GPU cluster
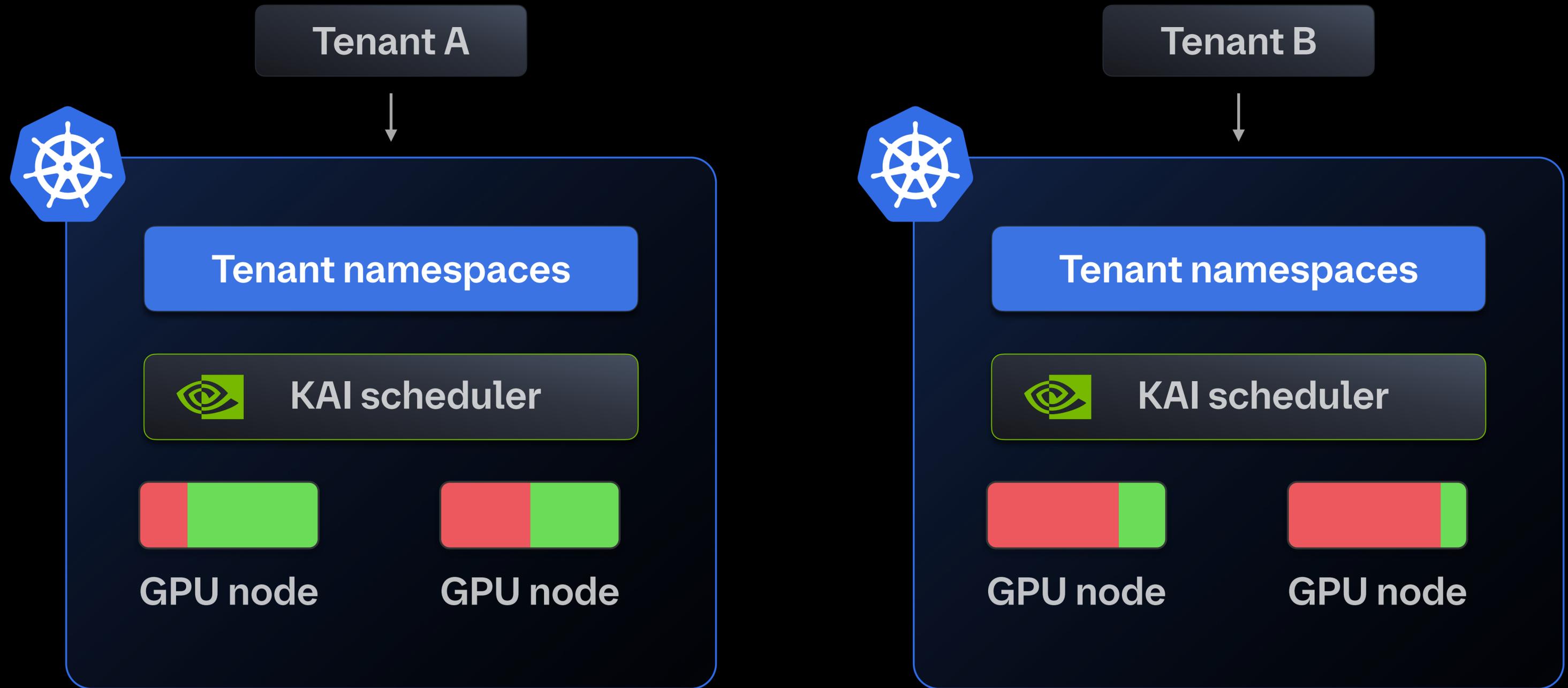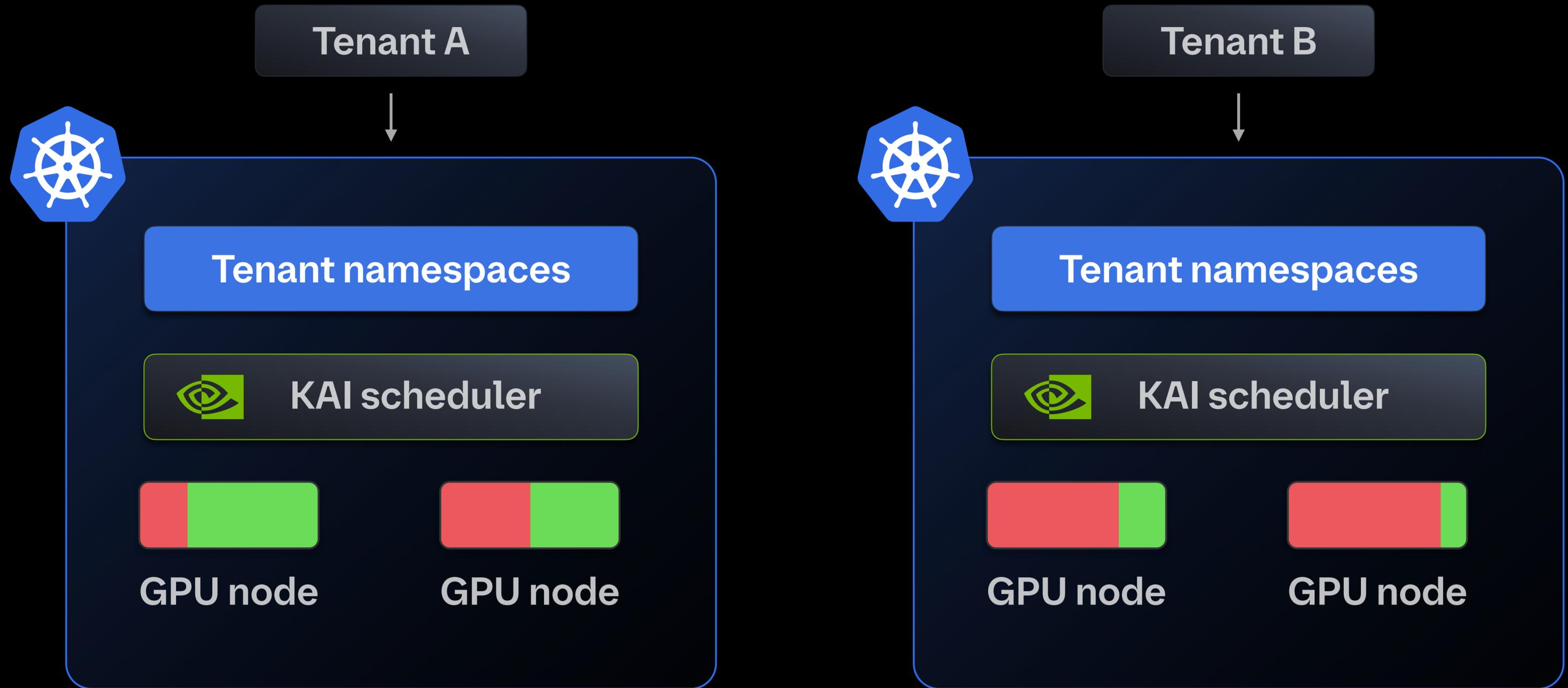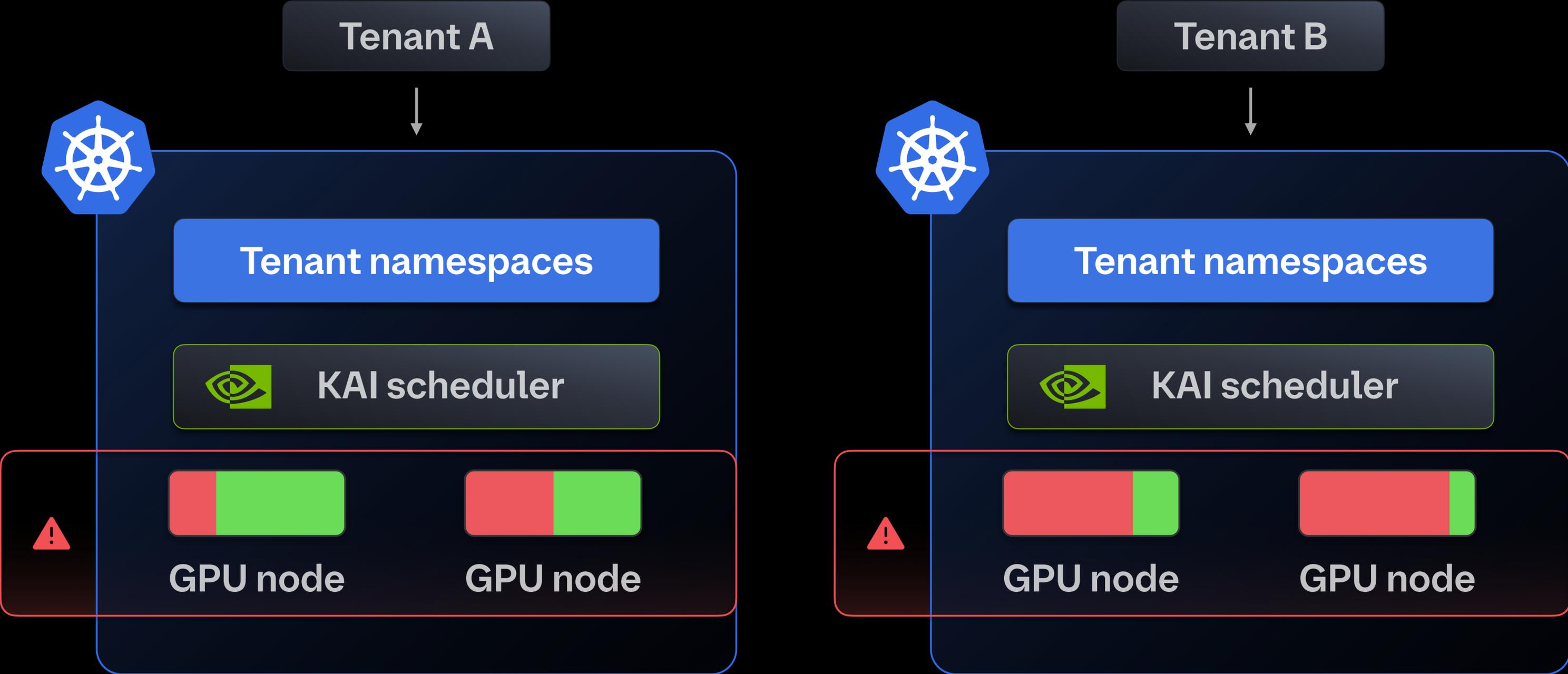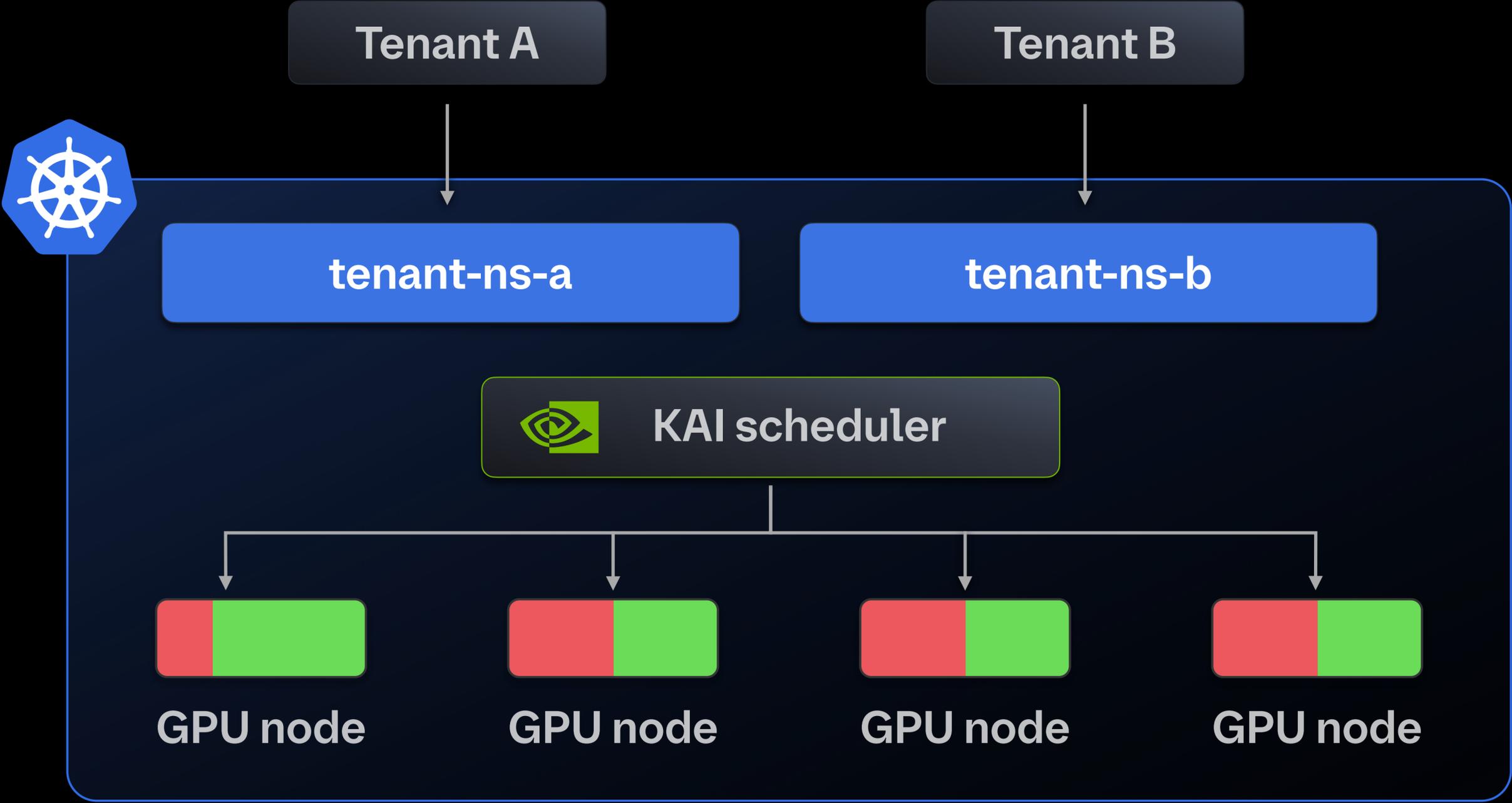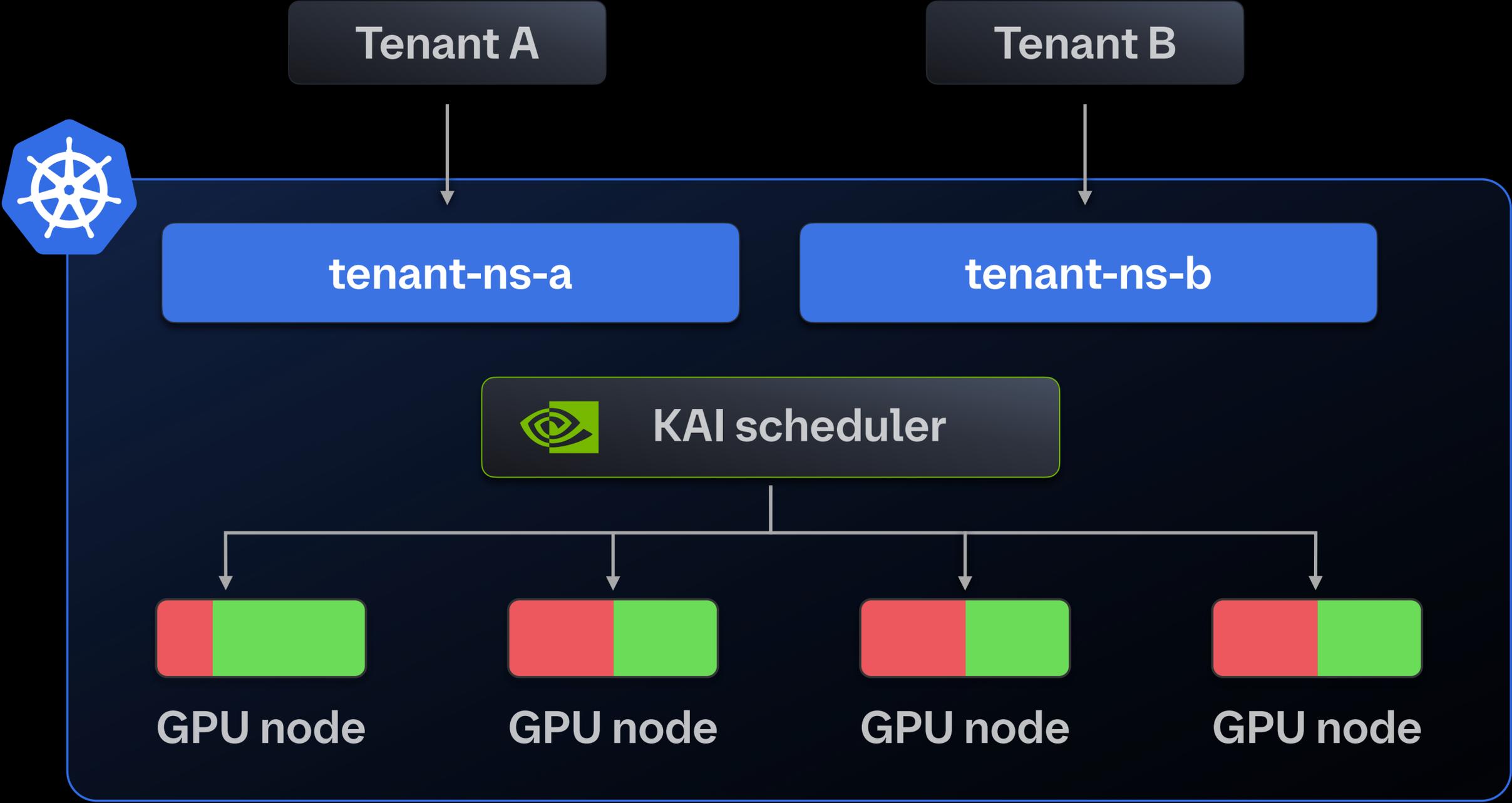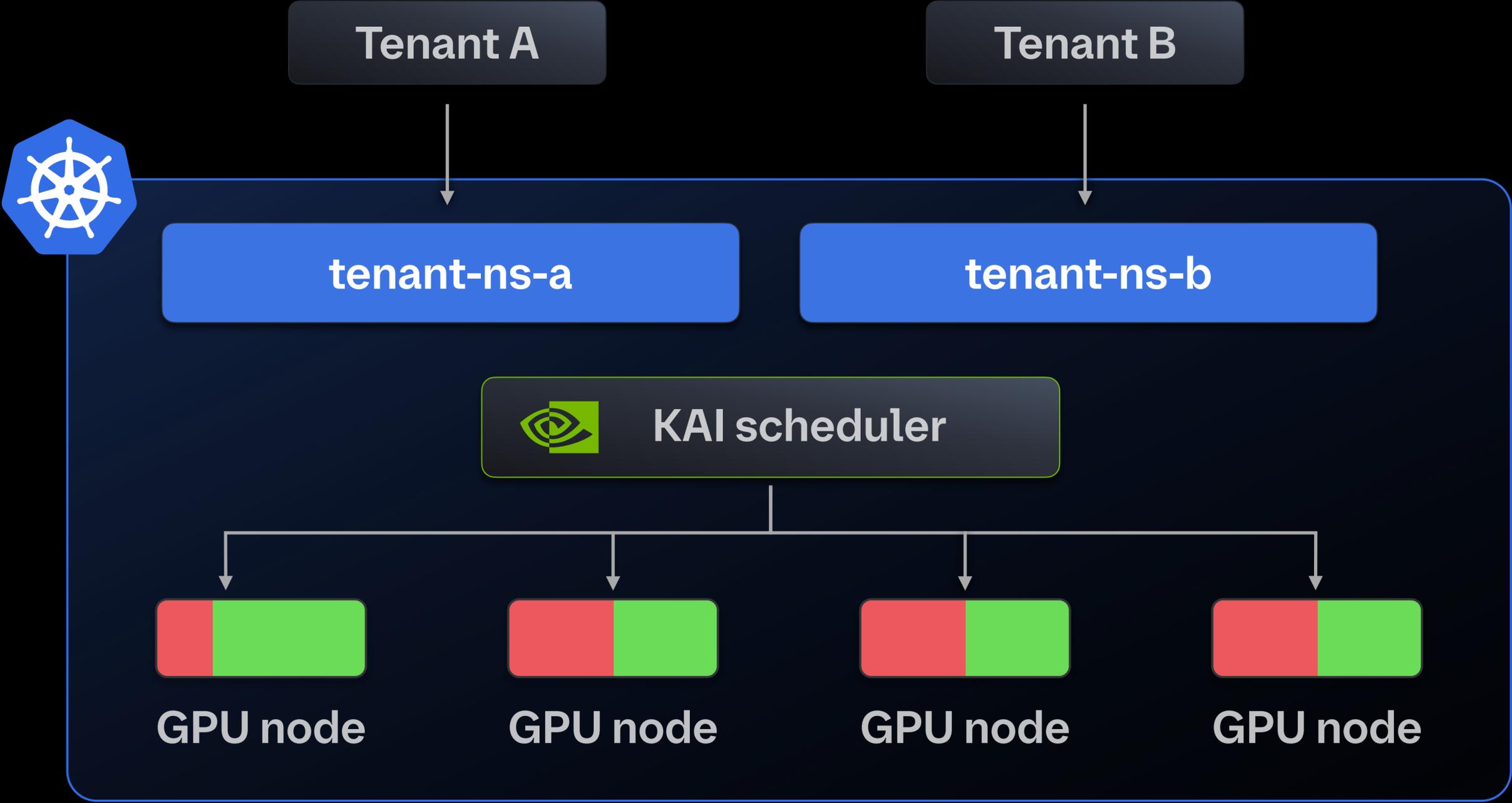
# Namespace-based multi-tenancy

# Namespace-based multi-tenancy

# Namespace-based multi-tenancy

Tenant A

Tenant B

tenant-ns-a

tenant-ns-b

KAI scheduler

GPU node

GPU node

GPU node

GPU node

# Namespace-based multi-tenancy

# Team-shared cluster

Controls virtual cluster

**Tenant**

Connects to cluster API server

## Virtual cluster context

Namespace in virtual cluster

| deployment | pod-1 | custom-resource |

| API server | Data store | Control manager | vcluster syncer |

vCluster

Controls K8s context

**Admin**

Connects to K8s API server

## K8s cluster context

Namespace in K8s

| synced-pod-1 | vcluster-pod |

| API server | etcd | Control manager | Scheduler |

Control plane

# Team-shared cluster

**Team-cluster context**

Namespace in team clusters

pod-1

custom-

API server

etcd

Control manager

Scheduler

Control plane

GPU sharing only works when you understand your trust boundaries.

Namespaces provide logical segmentation, but not isolation.

Per-team clusters address isolation, but not efficiency.

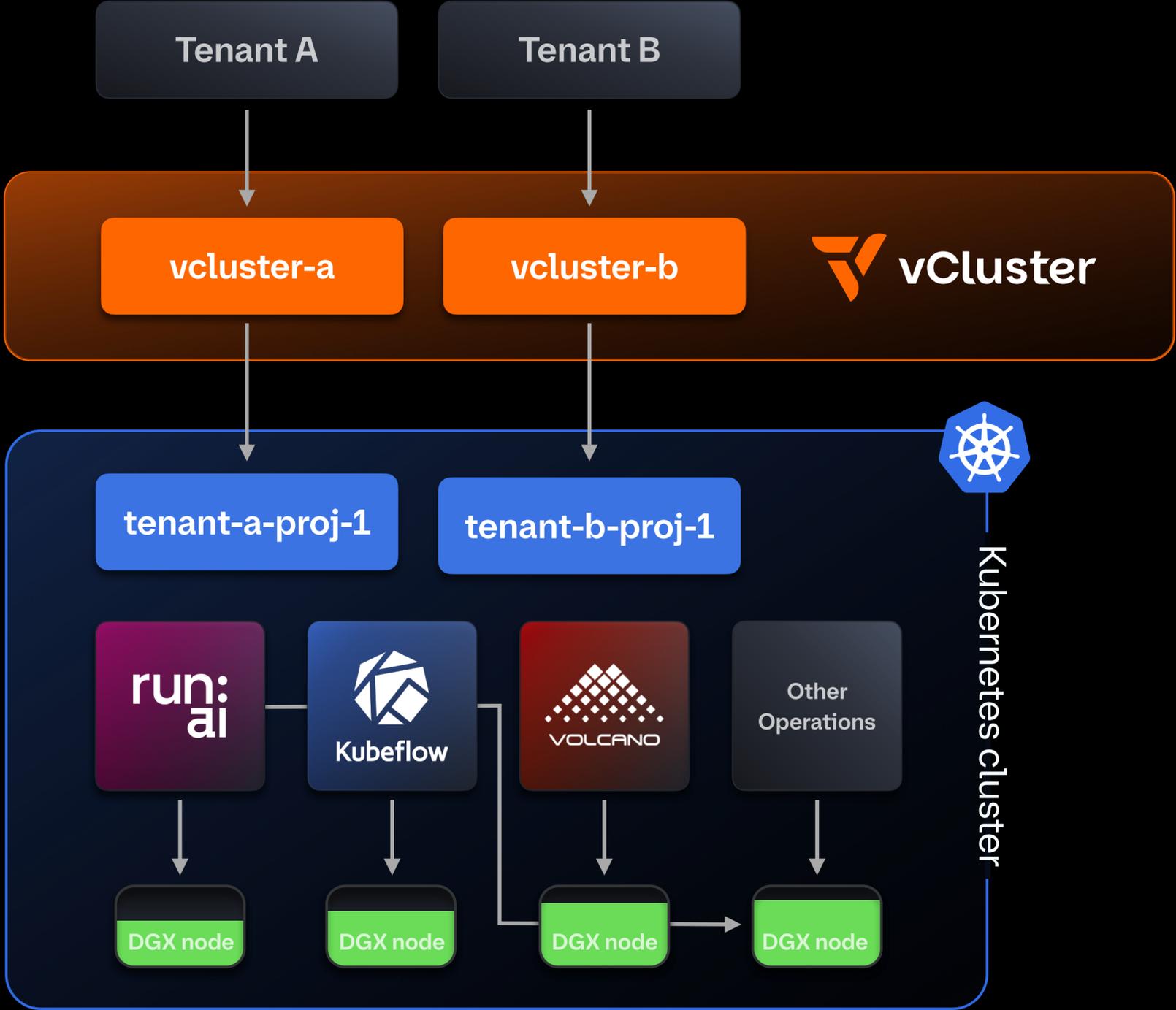# Multi-team clusters address efficiency, but not isolation.
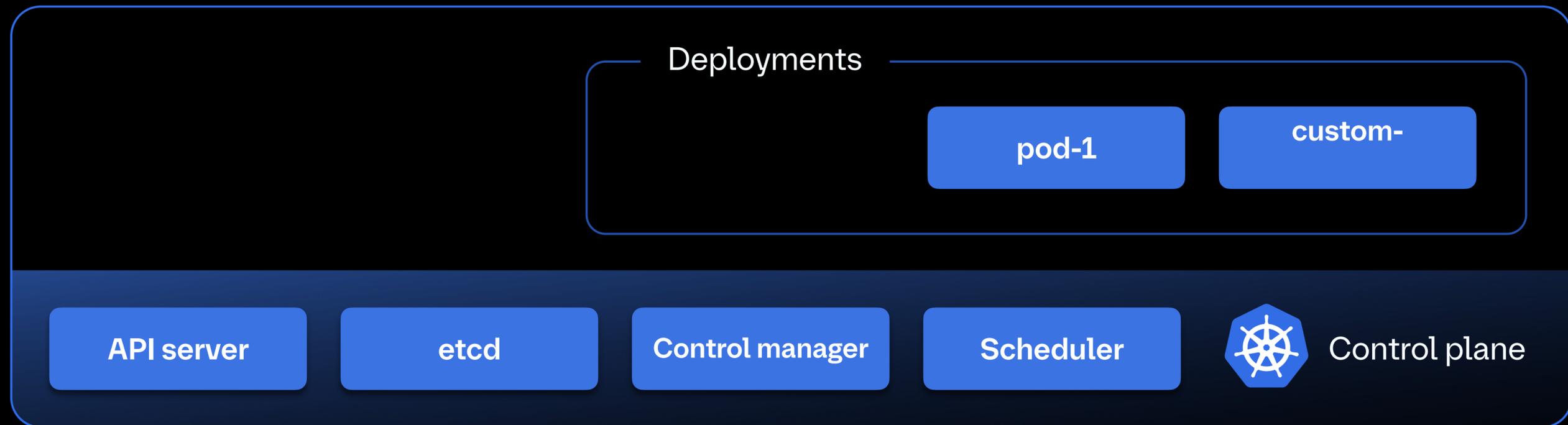
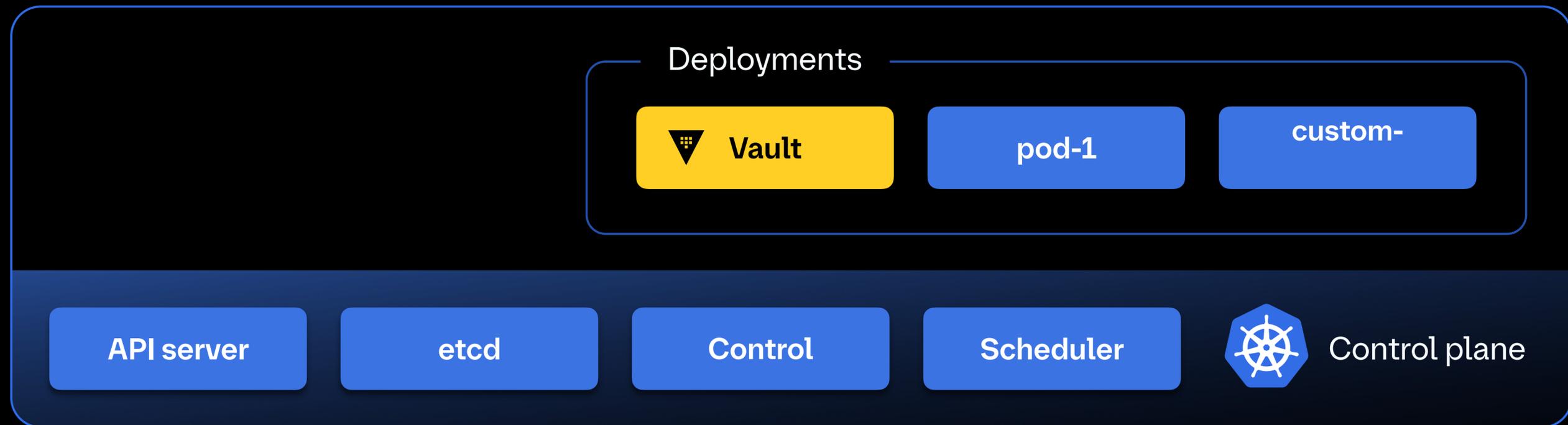# Where we need to be at

# How do you deal with cluster-level resource?

# Multi-tenant cluster

# Team-specific virtual cluster

Deployments

pod-1    custom-

API server    etcd    Control manager    Scheduler    Control plane
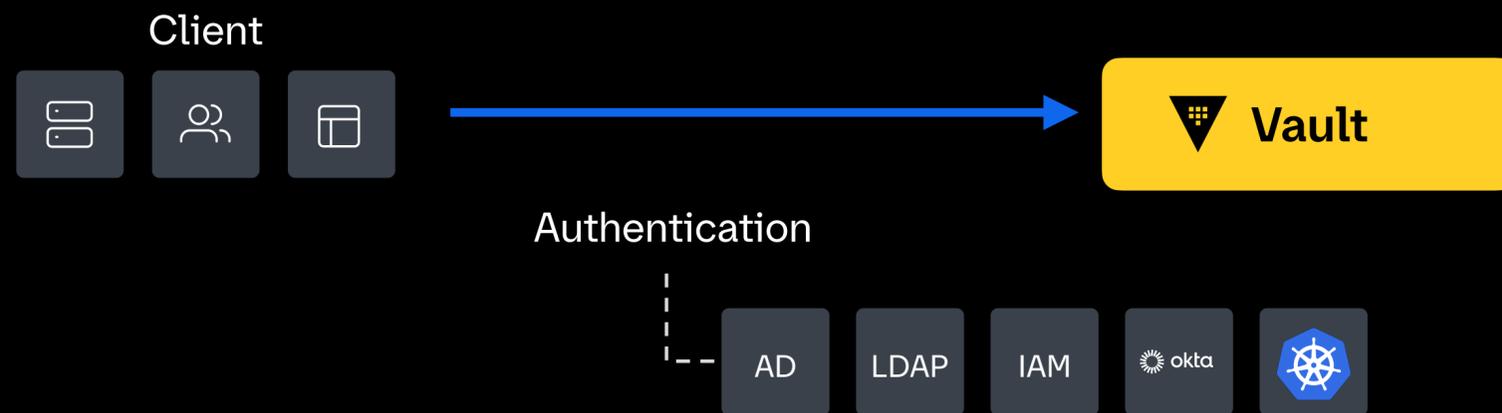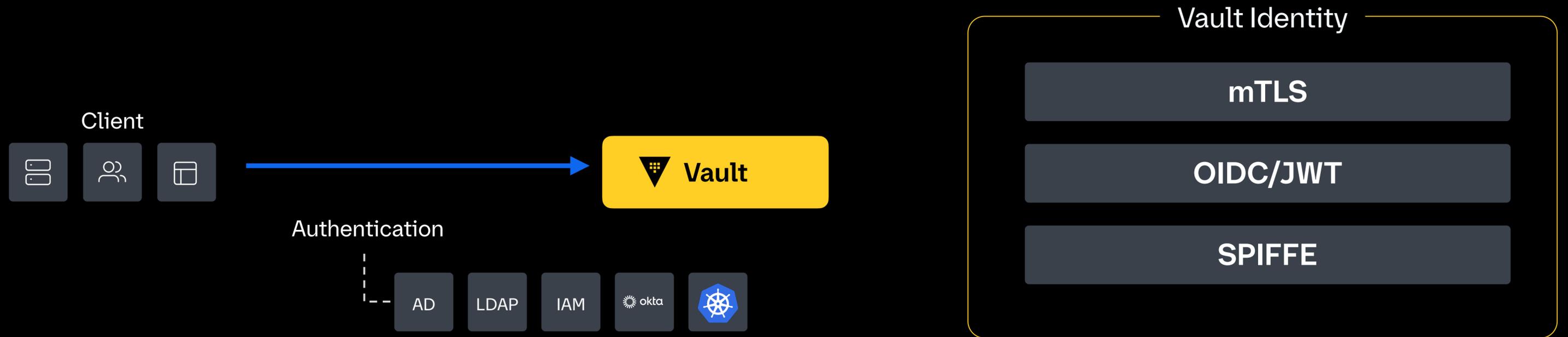
# Team-specific virtual cluster

# Security is important
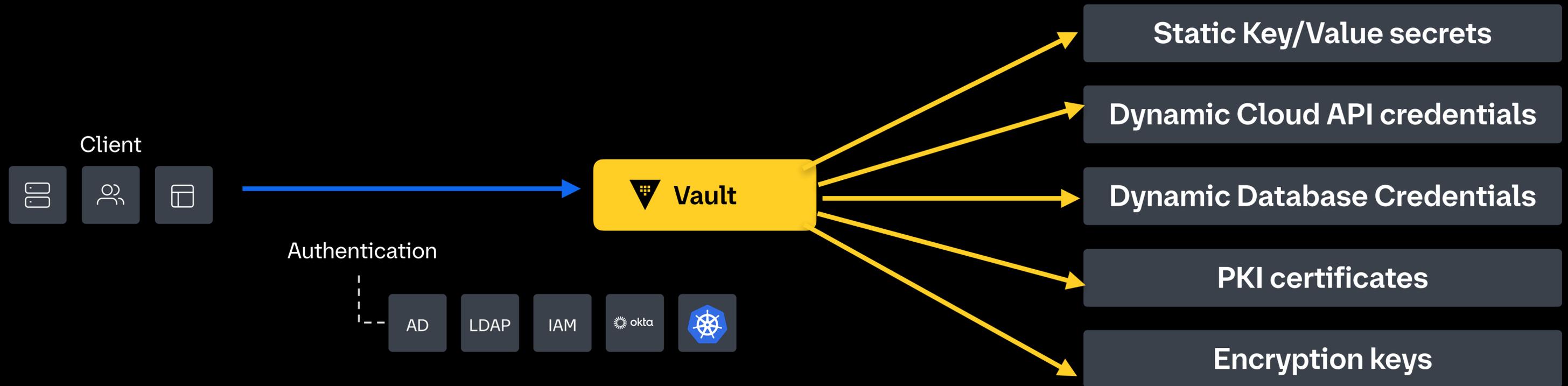
# Why Vault

# Why Vault

# Why Vault

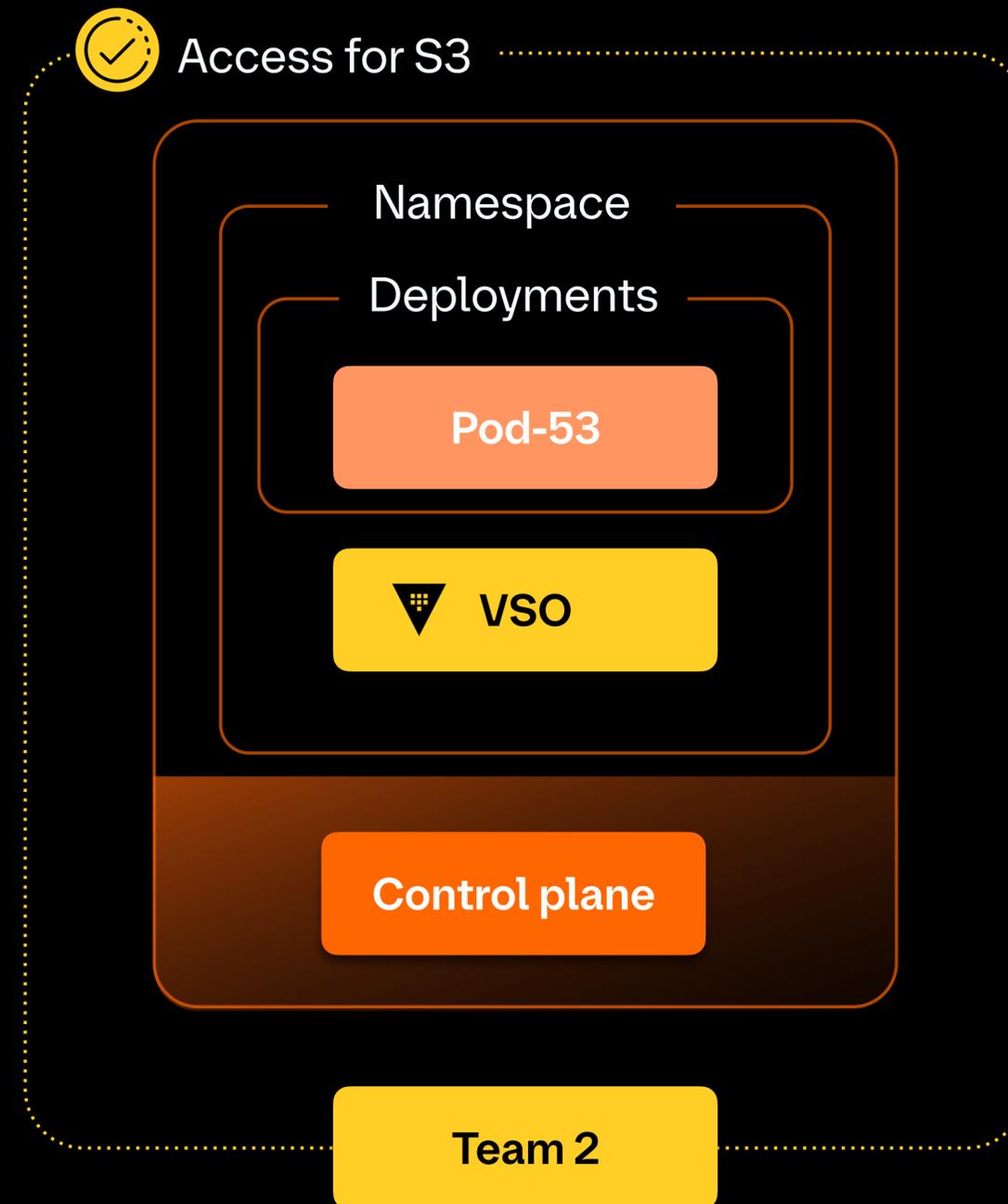# Eliminate forever credentials

# Audit All The Things

# Vault and Kubernetes

# Team-specific security

# Team-shared cluster

# Configuration

# Kubernetes Configuration

```yaml
---
apiVersion: secrets.hashicorp.com/v1beta1
kind: VaultConnection
spec:
 address: http://vault.vault.svc.cluster.local:8200
---
apiVersion: secrets.hashicorp.com/v1beta1
kind: VaultAuth
metadata:
 namespace: finetune
 name: vault-auth
spec:
 vaultConnectionRef: vault-connection
 method: kubernetes
 kubernetes:
   serviceAccount: finetune
   role: prod
   namespace: finetune
```

vault.yaml

# Kubernetes Configuration

```yaml
---
apiVersion: secrets.hashicorp.com/v1beta1
kind: VaultDynamicSecret
metadata:
  namespace: finetune
  name: vault-dynamic-secret-db-prod
spec:
  vaultAuthRef: vault-auth
  mount: db
  path: creds/pgsql-prod
  destination:
    create: true
    name: pgsql-prod
```

vault.yaml

# Vault Configuration

vault.yaml

```
$ vault write auth/kubernetes/role/demo
    bound_service_account_names=special-job-prod
    bound_service_account_namespaces=finetune
    policies=finetune-prod,prod,database-prod
    ttl=24h

$ cat database-prod.policy.hcl
path "db/creds/pgsql-prod" {
  capabilities = ["read"]
}
```
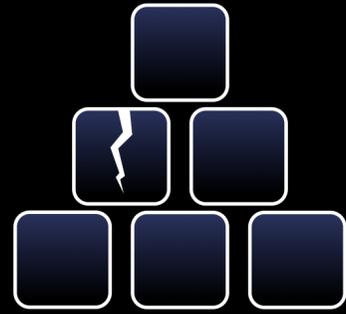
03

# What we need to remember

Architect for failure.

# Isolation is key.

Defense in-depth is defense that works.

# Virtual clusters limit blast-radius.

# Virtual clusters enable easier auditing.

# Thank you

speakerdeck.com/stmcallister
atodorov.me