



**CLOUD NATIVE &
KUBERNETES**

AI DAY

NORTH AMERICA



CLOUD NATIVE &
KUBERNETES

AI DAY

NORTH AMERICA

Multitenancy and Fairness at Scale with Kueue: A Case Study

Aldo Culquicondor, Google
Rajat Phull, Apple

What is Kueue?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

- Kueue interacts with **kube-scheduler** and **cluster-autoscaler** to provide a full batch /training system in Kubernetes.
- Kueue determines whether workloads should wait for resources or run, based on:
 - Per-tenant quotas
 - Borrowing and lending limits
 - Fair sharing rules New in v0.7
 - The hierarchy of the organization New in v0.9
- Kueue integrates with Pods, Job, JobSet, Kubeflow, KubeRay and has extension mechanisms.



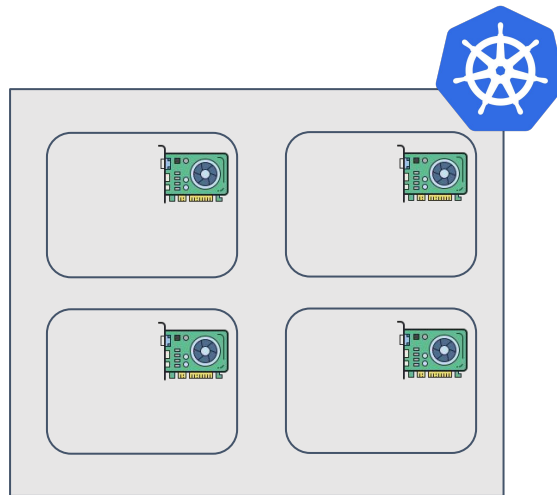
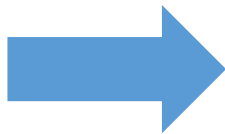
+



In the beginning...



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

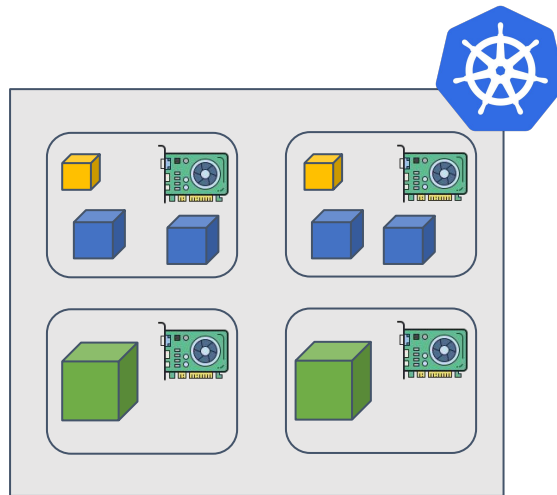
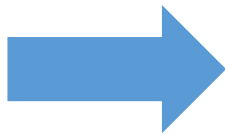


Bob
AI researcher

In the beginning...



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Bob
AI researcher

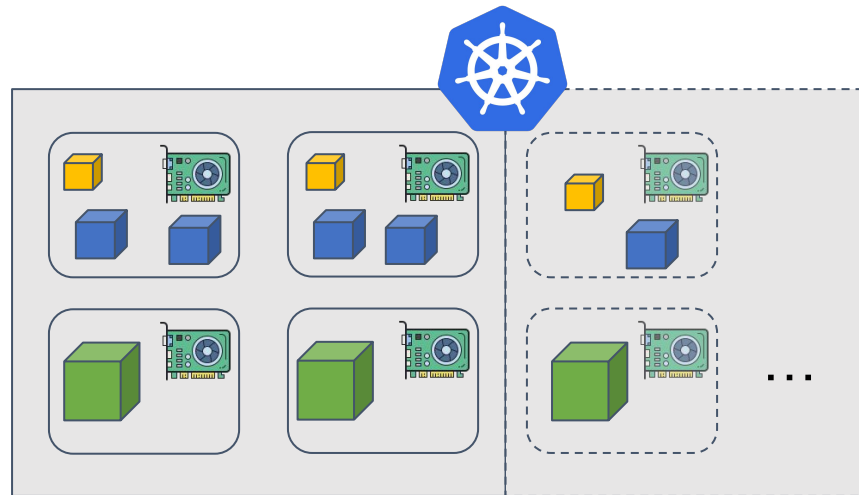
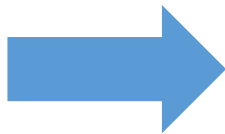
In the beginning...



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Bob
AI researcher

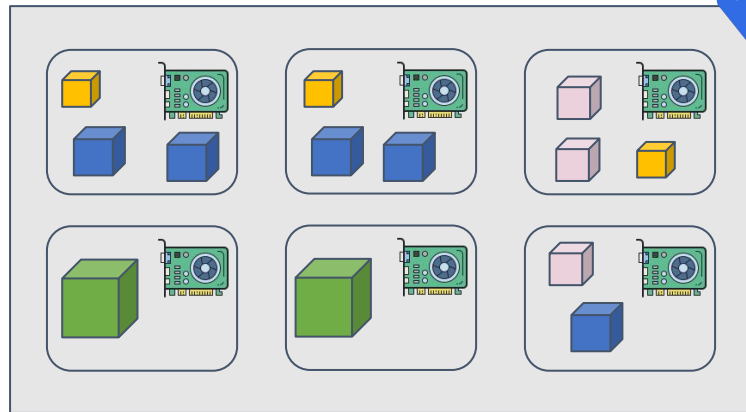
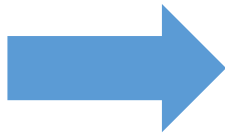


Autoscaled
nodes

But clusters aren't infinite



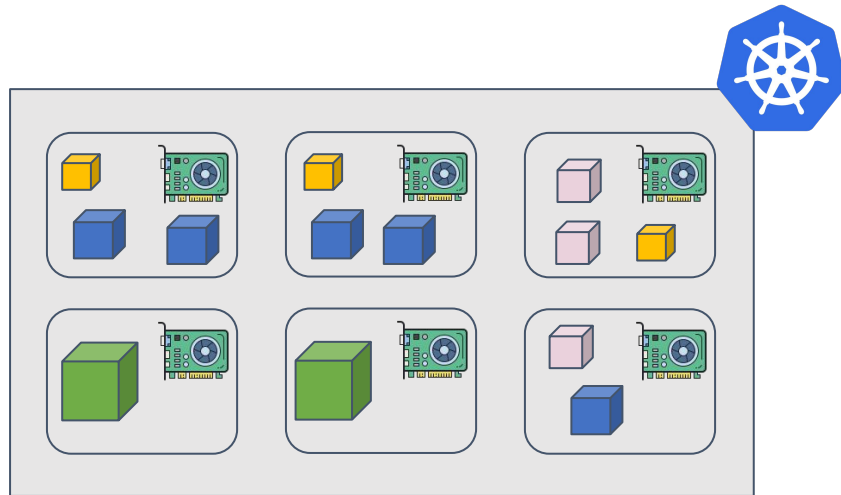
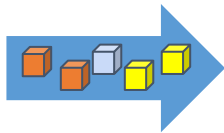
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



But clusters aren't infinite



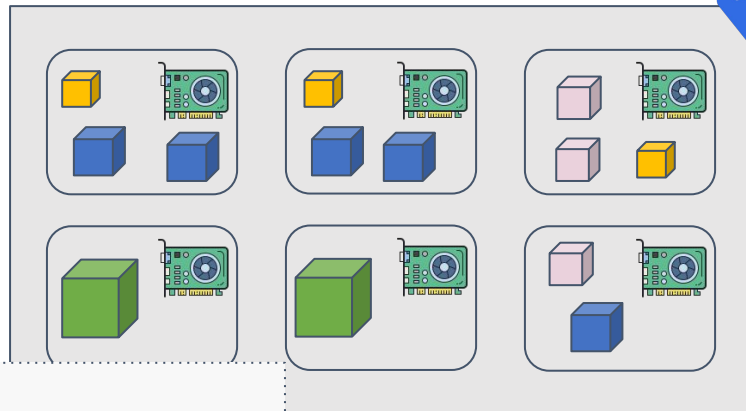
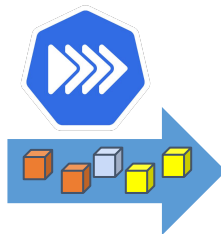
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



But clusters aren't infinite



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

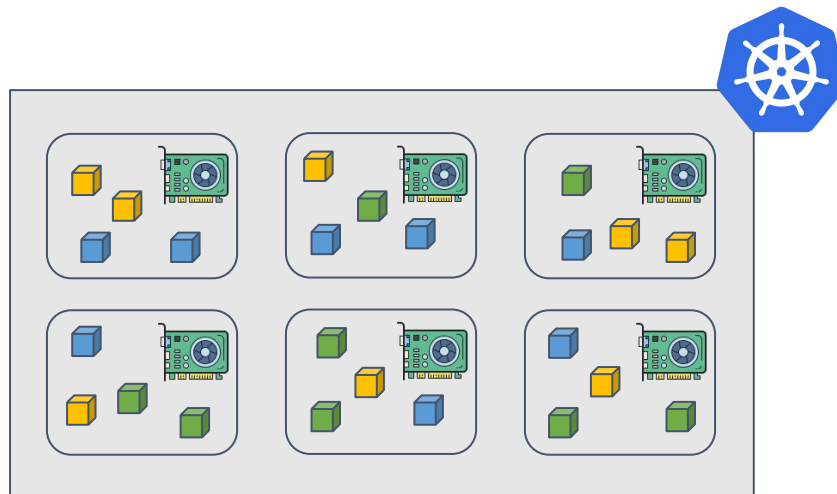
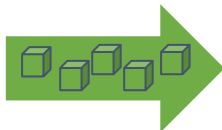
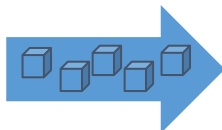
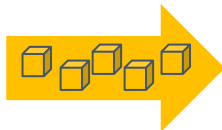


```
kind: ClusterQueue
metadata:
  name: "bob-queue"
spec:
  resourceGroups:
  - coveredResources: ["cpu", "memory", "acme.com/gpu"]
    flavors:
    - name: "default-flavor"
      resources:
      - name: "cpu"
        nominalQuota: 24
      - name: "memory"
        nominalQuota: 48Gi
      - name: "acme.com/gpu"
        nominalQuota: 24
```

Clusters are often shared



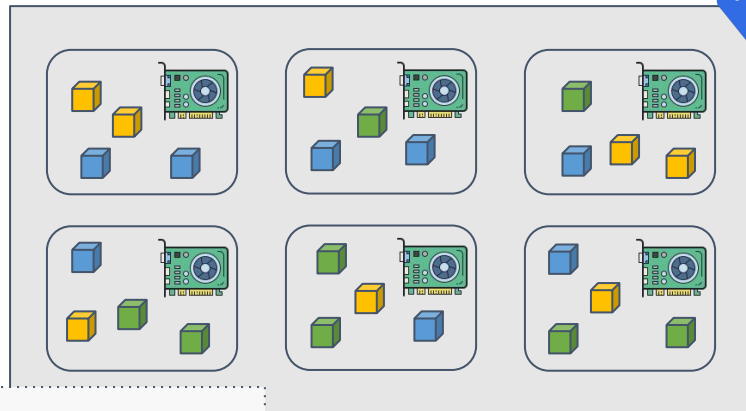
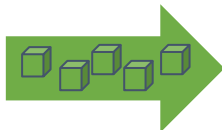
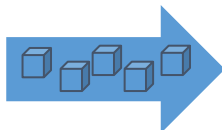
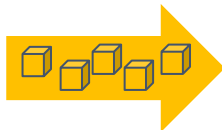
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Clusters are often shared



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

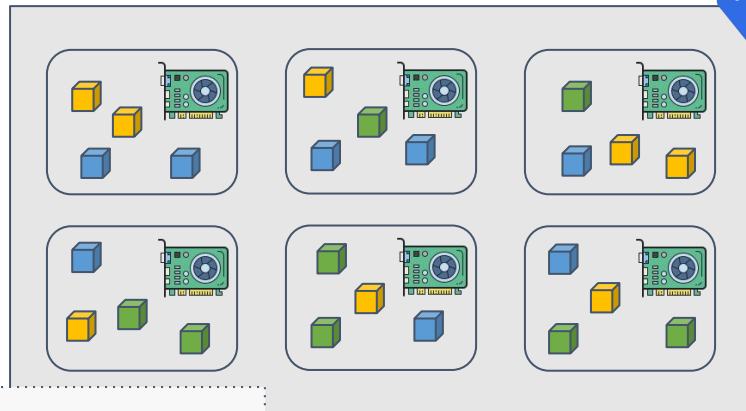
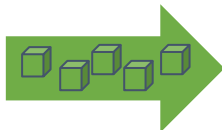
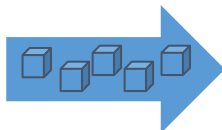
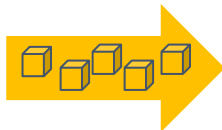


```
kind: ClusterQueue
metadata:
  name: "bob-queue"
spec:
  cohort: "lab"
  resourceGroups:
  - coveredResources: ["acme.com/gpu"]
    flavors:
    - name: "default-flavor"
      resources:
      - name: "acme.com/gpu"
        nominalQuota: 8
  preemption:
    reclaimWithinCohort: "Any"
  fairSharing:
    weight: 1
```

Clusters are often shared

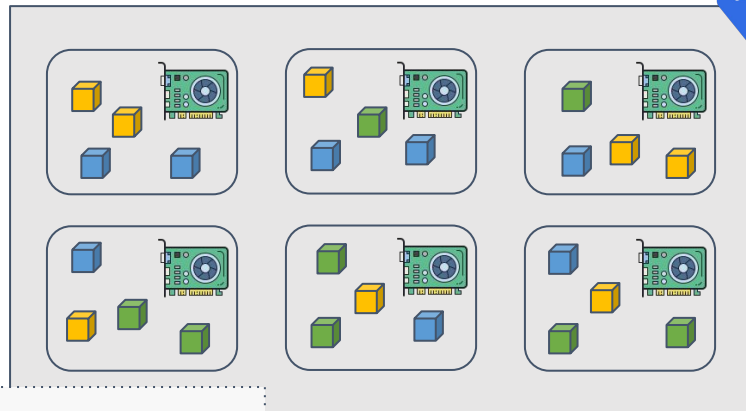
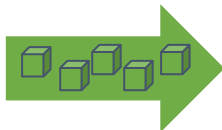
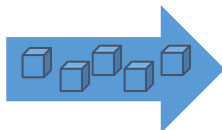
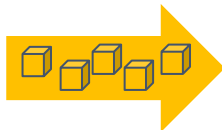


CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



```
kind: ClusterQueue
metadata:
  name: "bob-queue"
spec:
  cohort: "lab"
  resourceGroups:
  - coveredResources: ["acme.com/gpu"]
    flavors:
    - name: "default-flavor"
      resources:
      - name: "acme.com/gpu"
        nominalQuota: 8
  preemption:
    reclaimWithinCohort: "Any"
  fairSharing:
    weight: 1
```

Clusters are often shared



```
kind: ClusterQueue
metadata:
  name: "bob-queue"
spec:
  cohort: "lab"
  resourceGroups:
  - coveredResources: ["acme.com/gpu"]
    flavors:
    - name: "default-flavor"
      resources:
      - name: "acme.com/gpu"
        nominalQuota: 8
  preemption:
    reclaimWithinCohort: "Any"
  fairSharing:
    weight: 1
```

```
kind: ClusterQueue
metadata:
  name: "alice-queue"
spec:
  cohort: "lab"
  resourceGroups:
  - coveredResources: ["acme.com/gpu"]
    flavors:
    - name: "default-flavor"
      resources:
      - name: "acme.com/gpu"
        nominalQuota: 8
  preemption:
    reclaimWithinCohort: "Any"
  fairSharing:
    weight: 2
```

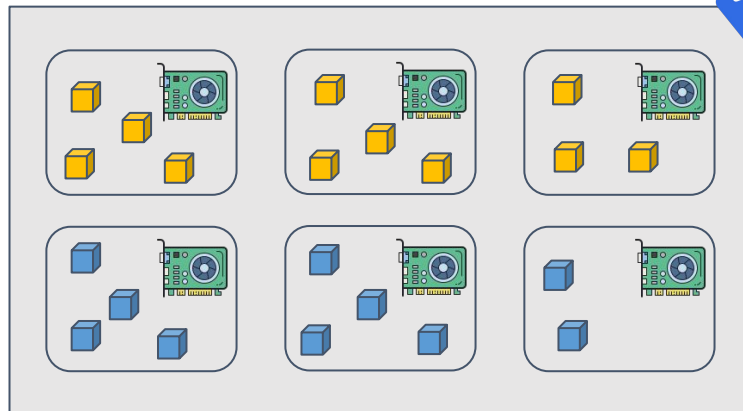
Fair Sharing

Fair Sharing: Who gets to schedule?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	
8	
8	

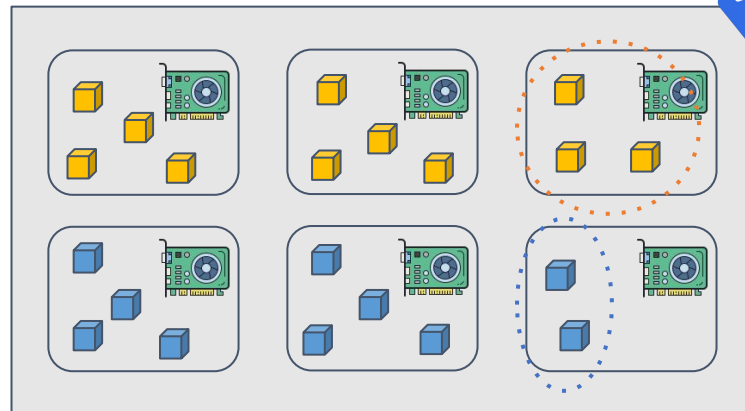


Fair Sharing: Who gets to schedule?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	3
8	2
8	0

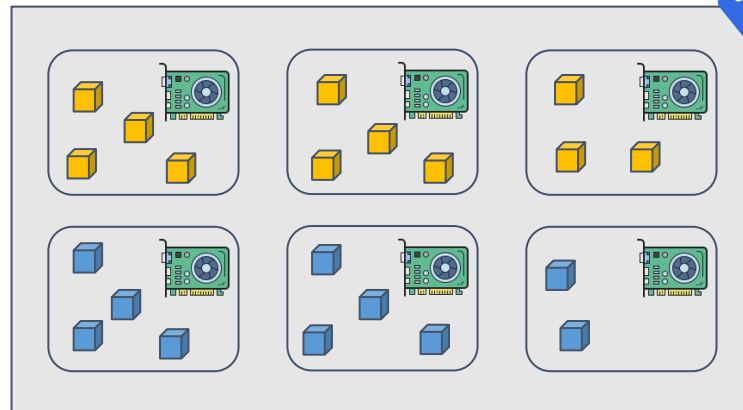
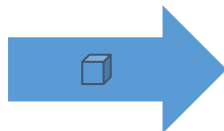
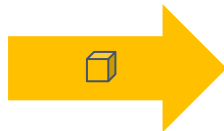


Fair Sharing: Who gets to schedule?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	3
8	2
8	0

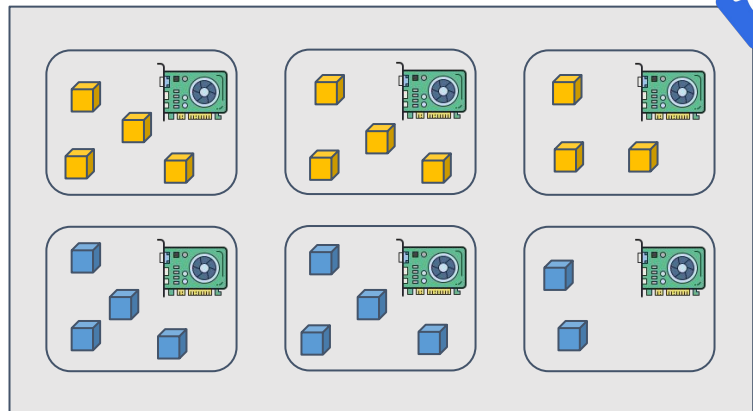
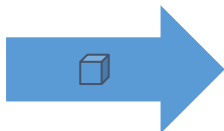
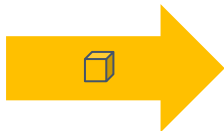


Fair Sharing: Who gets to schedule?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	3
8	2
8	0



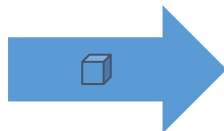
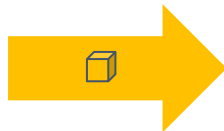
$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Fair Sharing: Who gets to schedule?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	3
8	2
8	0



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Schedule the
ClusterQueue with the
lowest share value first

Fair Sharing: Who gets to schedule?

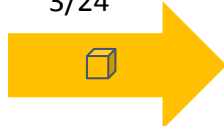


CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

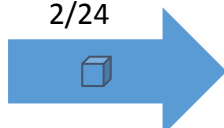
Nominal Quotas	Borrowing
8	3
8	2
8	0



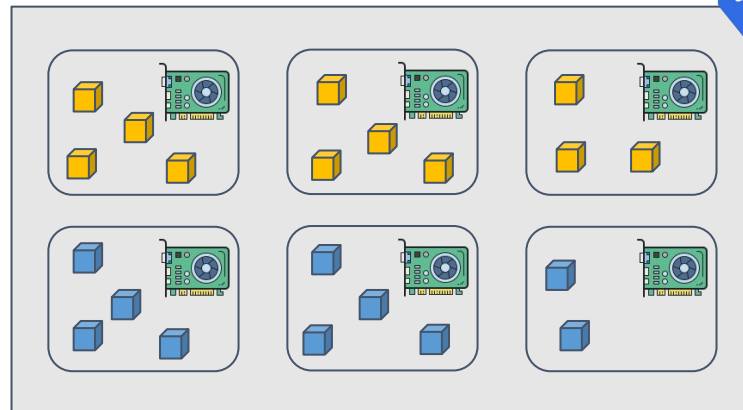
3/24



2/24



0/24



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Schedule the
ClusterQueue with the
lowest share value first

Fair Sharing: Who gets to schedule?

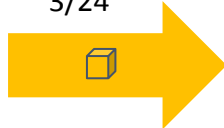


CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

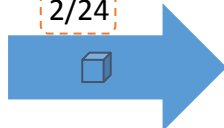
Nominal Quotas	Borrowing
8	3
8	2
8	0



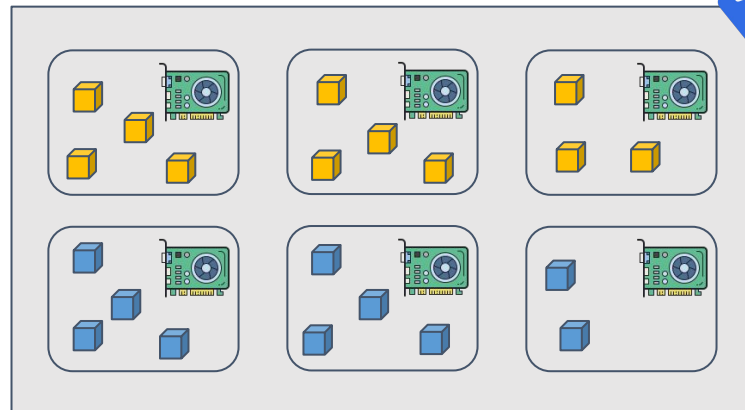
3/24



2/24



0/24



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Schedule the
ClusterQueue with the
lowest share value first

Fair Sharing: Who gets to schedule?

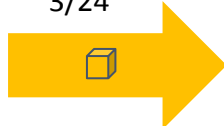


CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

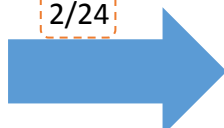
Nominal Quotas	Borrowing
8	3
8	2
8	0



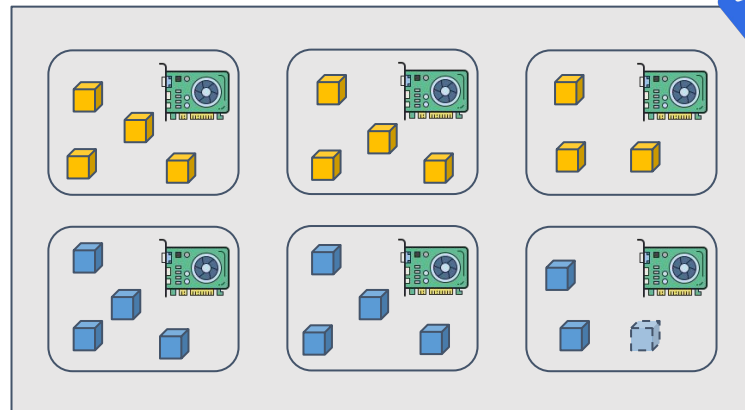
3/24



2/24



0/24



$$\text{ShareValue} = \frac{\text{UsageAboveQuota}}{\text{sum(Quotas)} * \text{Weight}}$$

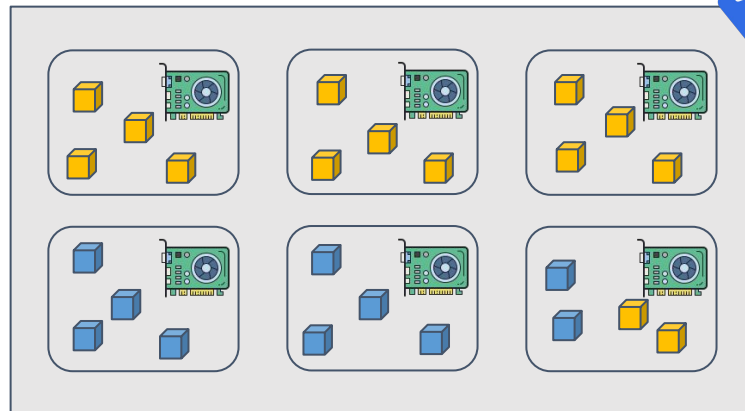
Schedule the
ClusterQueue with the
lowest share value first

Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	
8	
8	

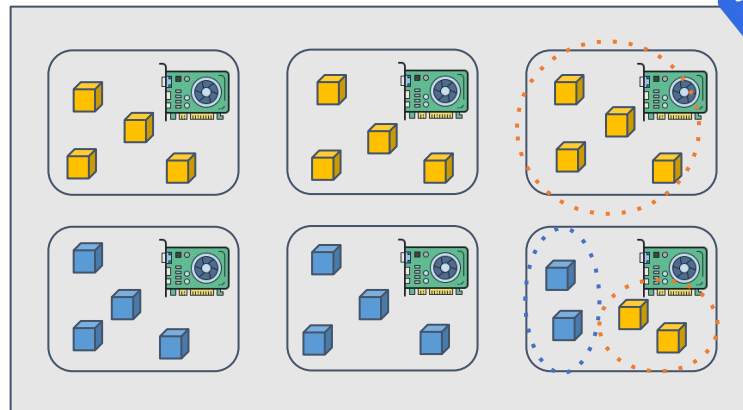


Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	6
8	2
8	0

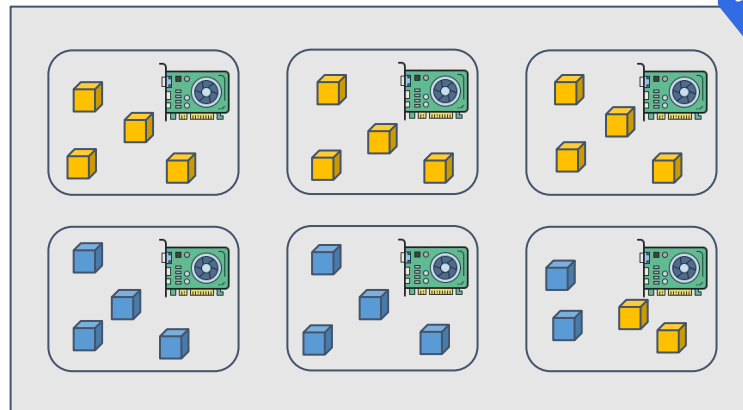


Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	6
8	2
8	0

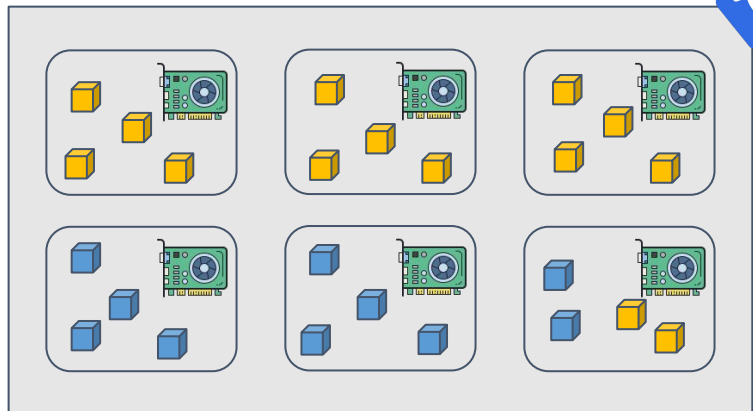
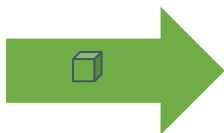


Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	6
8	2
8	0



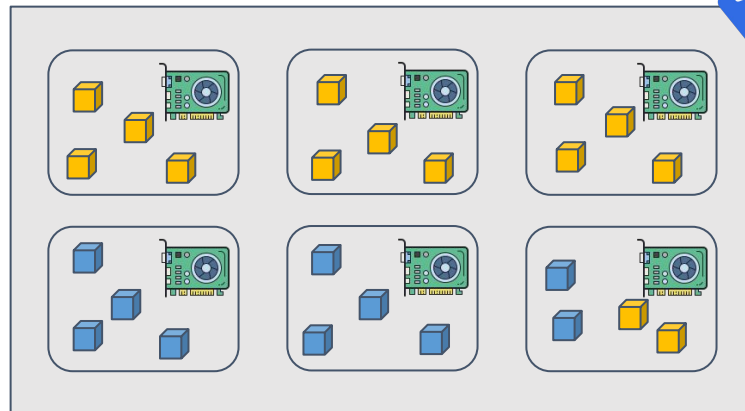
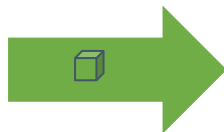
$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Nominal Quotas	Borrowing
8	6
8	2
8	0



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Preempt the
ClusterQueue with the
highest share value first

Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

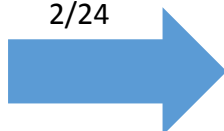
Nominal Quotas	Borrowing
8	6
8	2
8	0



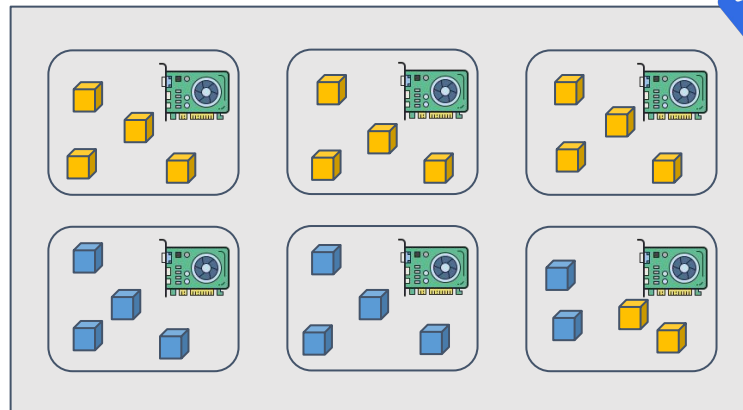
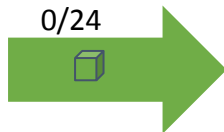
6/24



2/24



0/24



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Preempt the
ClusterQueue with the
highest share value first

Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

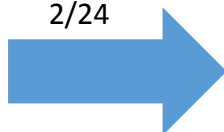
Nominal Quotas	Borrowing
8	6
8	2
8	0



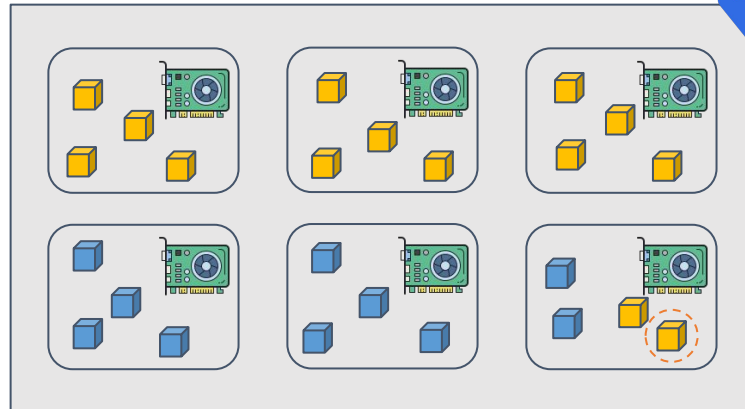
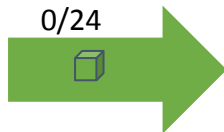
6/24



2/24



0/24



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Preempt the
ClusterQueue with the
highest share value first

Fair Sharing: Who gets preempted?

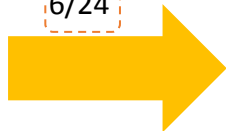


CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

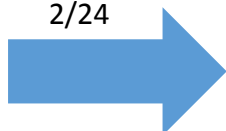
Nominal Quotas	Borrowing
8	6
8	2
8	0



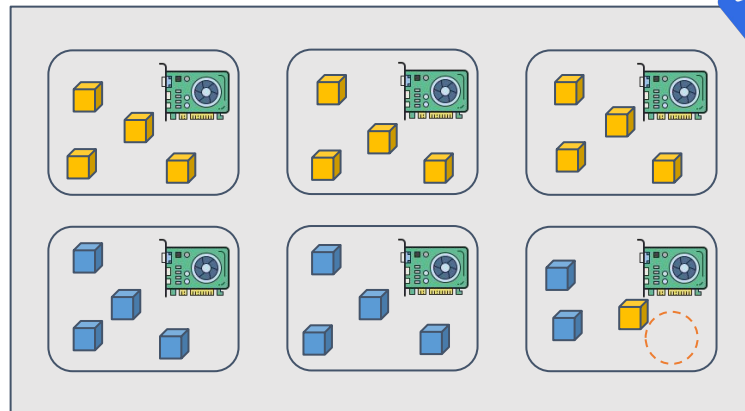
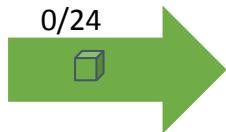
6/24



2/24



0/24



$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

Preempt the
ClusterQueue with the
highest share value first

Fair Sharing: Who gets preempted?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

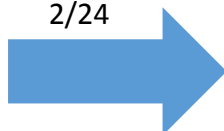
Nominal Quotas	Borrowing
8	6
8	2
8	0



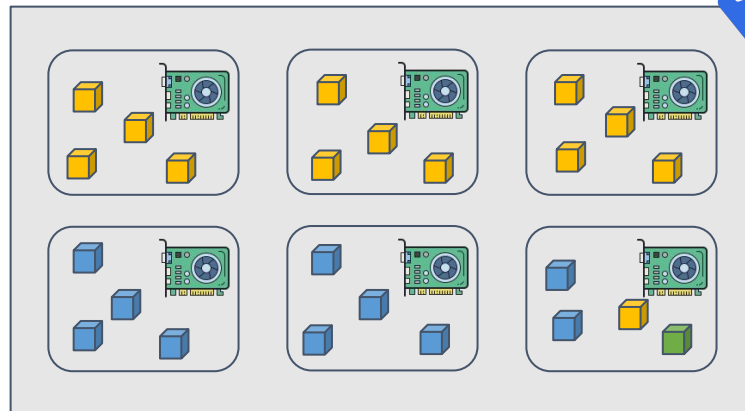
5/24



2/24



0/24

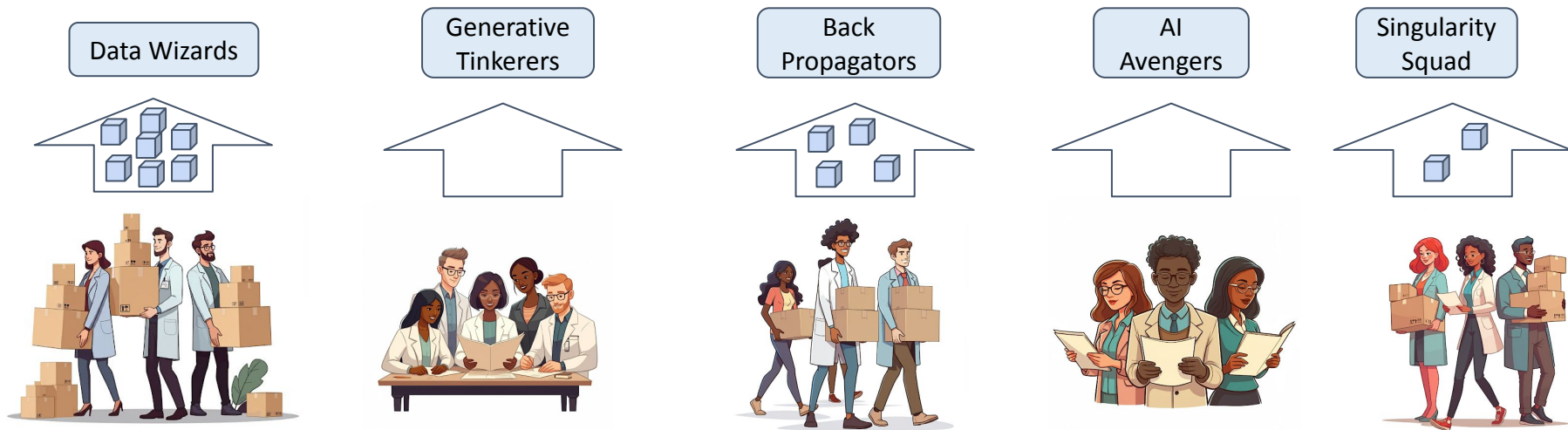


$$\text{ShareValue} = \frac{\text{Borrowing}}{\text{sum(Quotas)} * \text{Weight}}$$

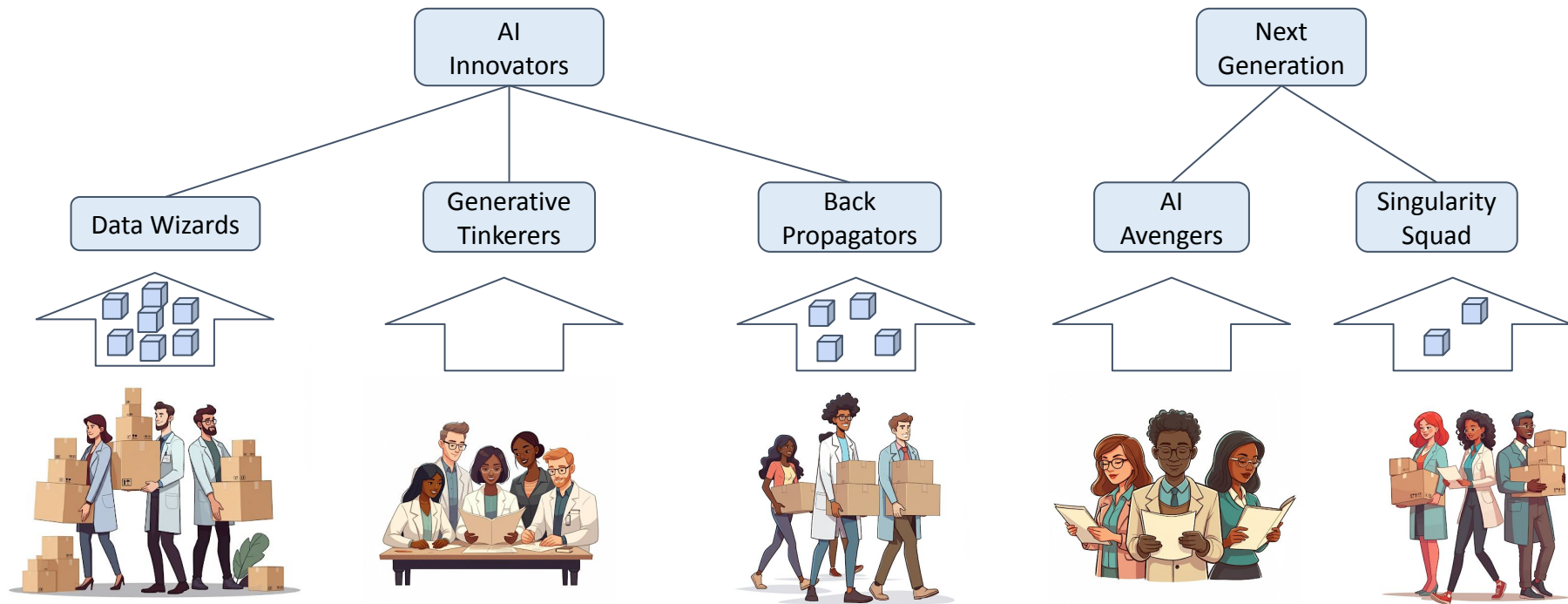
Preempt the
ClusterQueue with the
highest share value first

Hierarchical Cohorts

Clusters can be used by large organizations



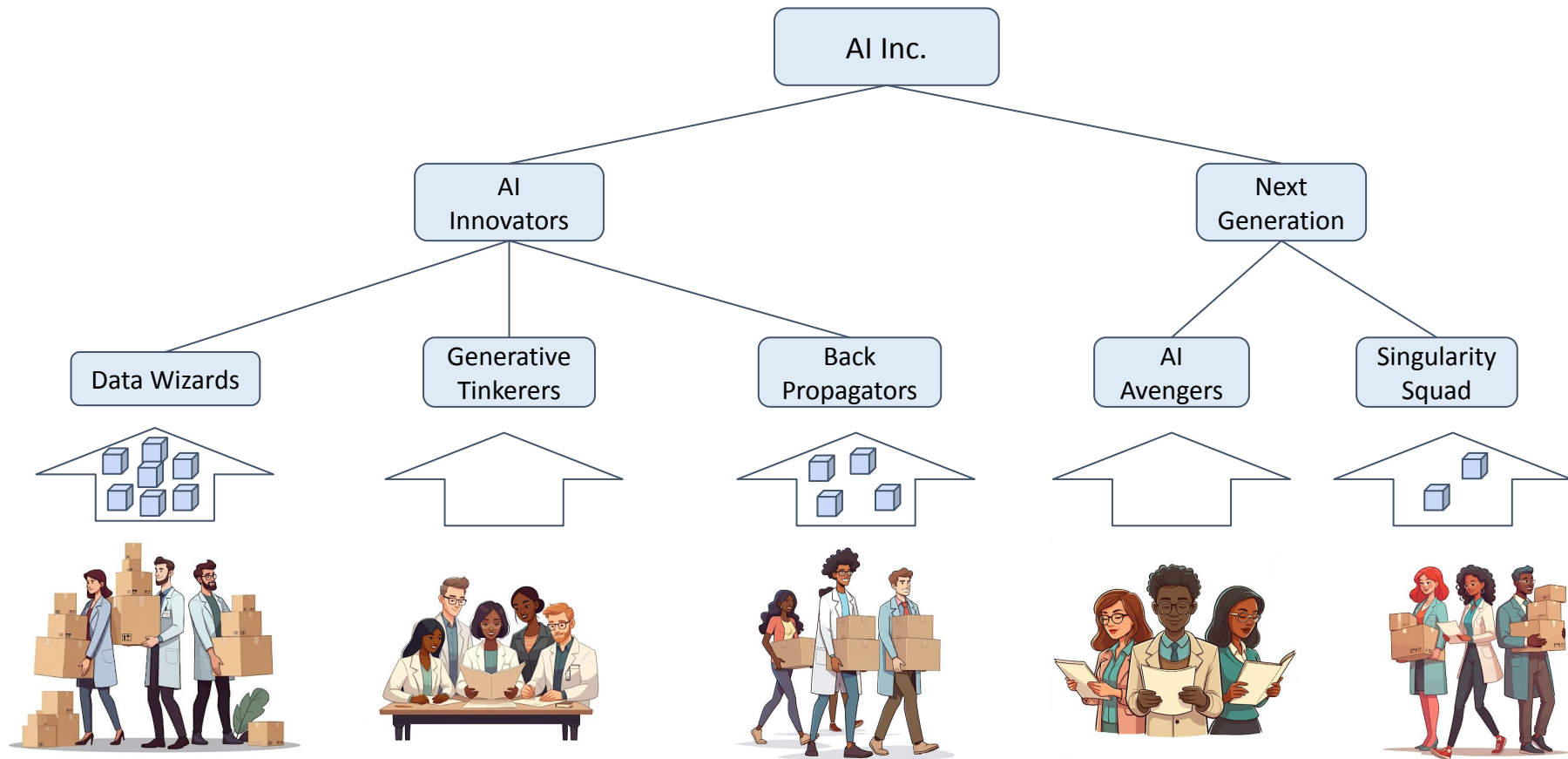
Clusters can be used by large organizations



Clusters can be used by large organizations



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

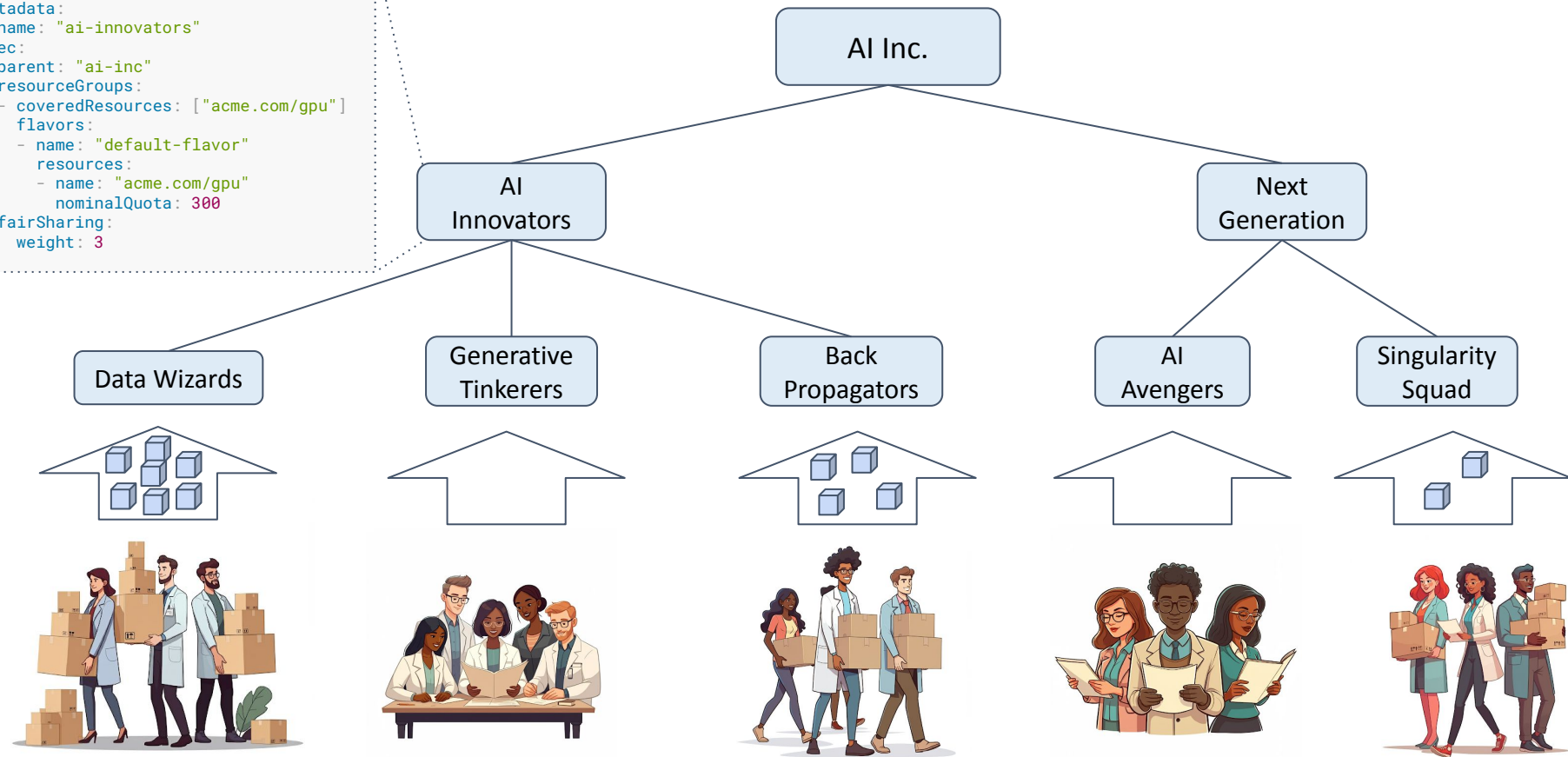


Clusters can be used by large organizations



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

```
kind: Cohort
metadata:
  name: "ai-innovators"
spec:
  parent: "ai-inc"
  resourceGroups:
    - coveredResources: ["acme.com/gpu"]
      flavors:
        - name: "default-flavor"
          resources:
            - name: "acme.com/gpu"
              nominalQuota: 300
  fairSharing:
    weight: 3
```

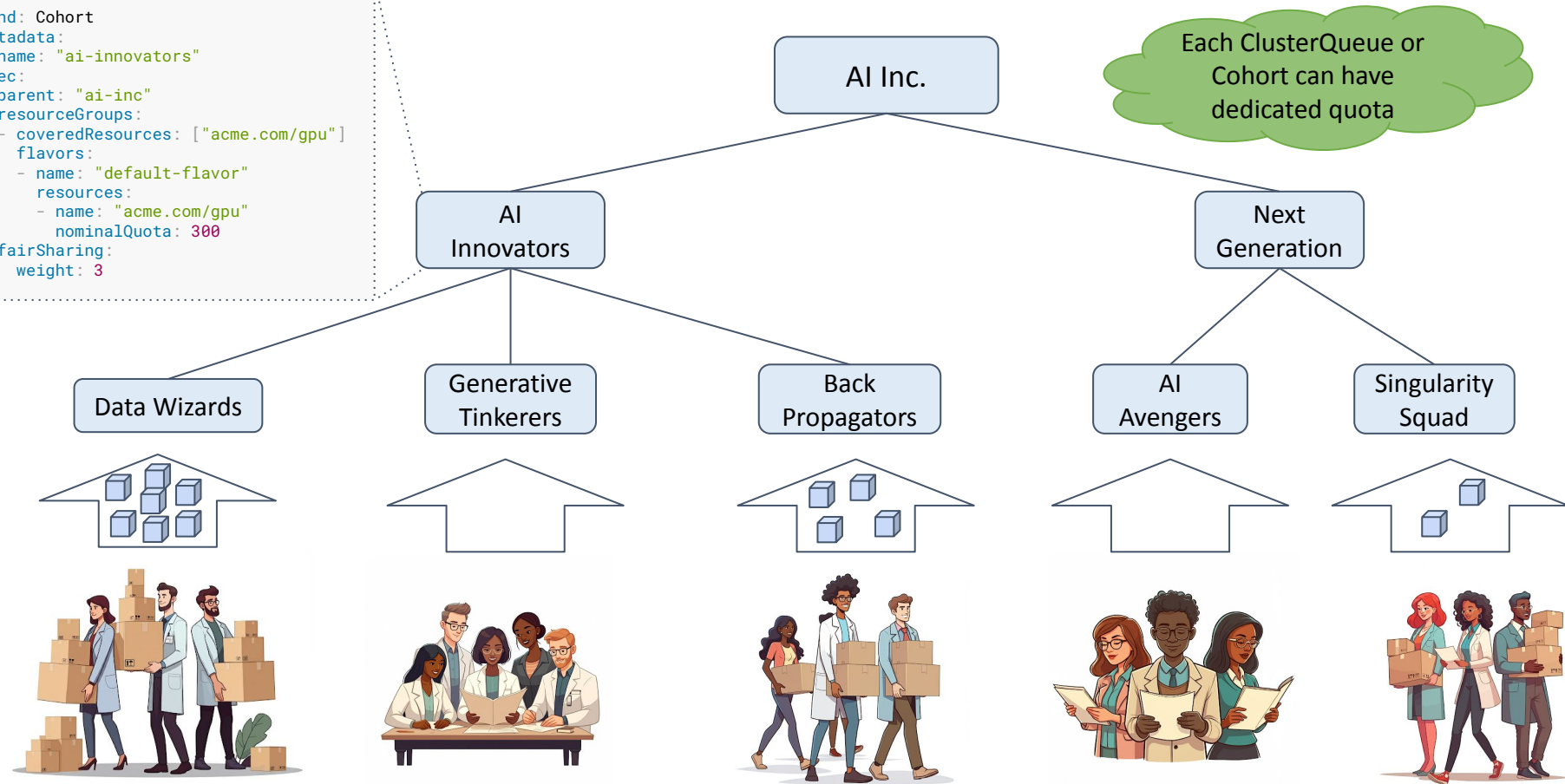


Clusters can be used by large organizations



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

```
kind: Cohort
metadata:
  name: "ai-innovators"
spec:
  parent: "ai-inc"
  resourceGroups:
    - coveredResources: ["acme.com/gpu"]
      flavors:
        - name: "default-flavor"
          resources:
            - name: "acme.com/gpu"
              nominalQuota: 300
  fairSharing:
    weight: 3
```



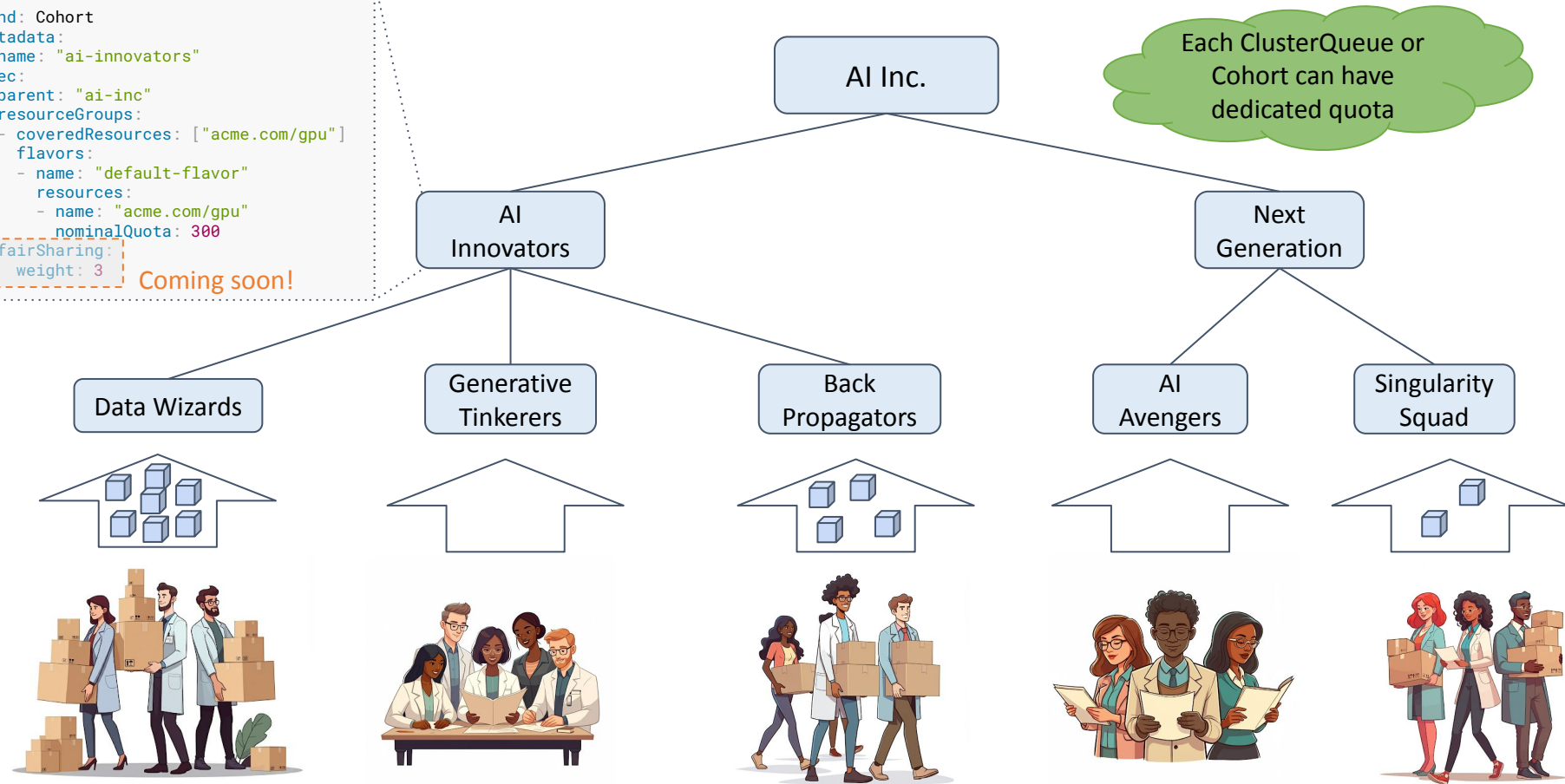
Clusters can be used by large organizations



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

```
kind: Cohort
metadata:
  name: "ai-innovators"
spec:
  parent: "ai-inc"
  resourceGroups:
    - coveredResources: ["acme.com/gpu"]
      flavors:
        - name: "default-flavor"
          resources:
            - name: "acme.com/gpu"
              nominalQuota: 300
      fairSharing:
        weight: 3
```

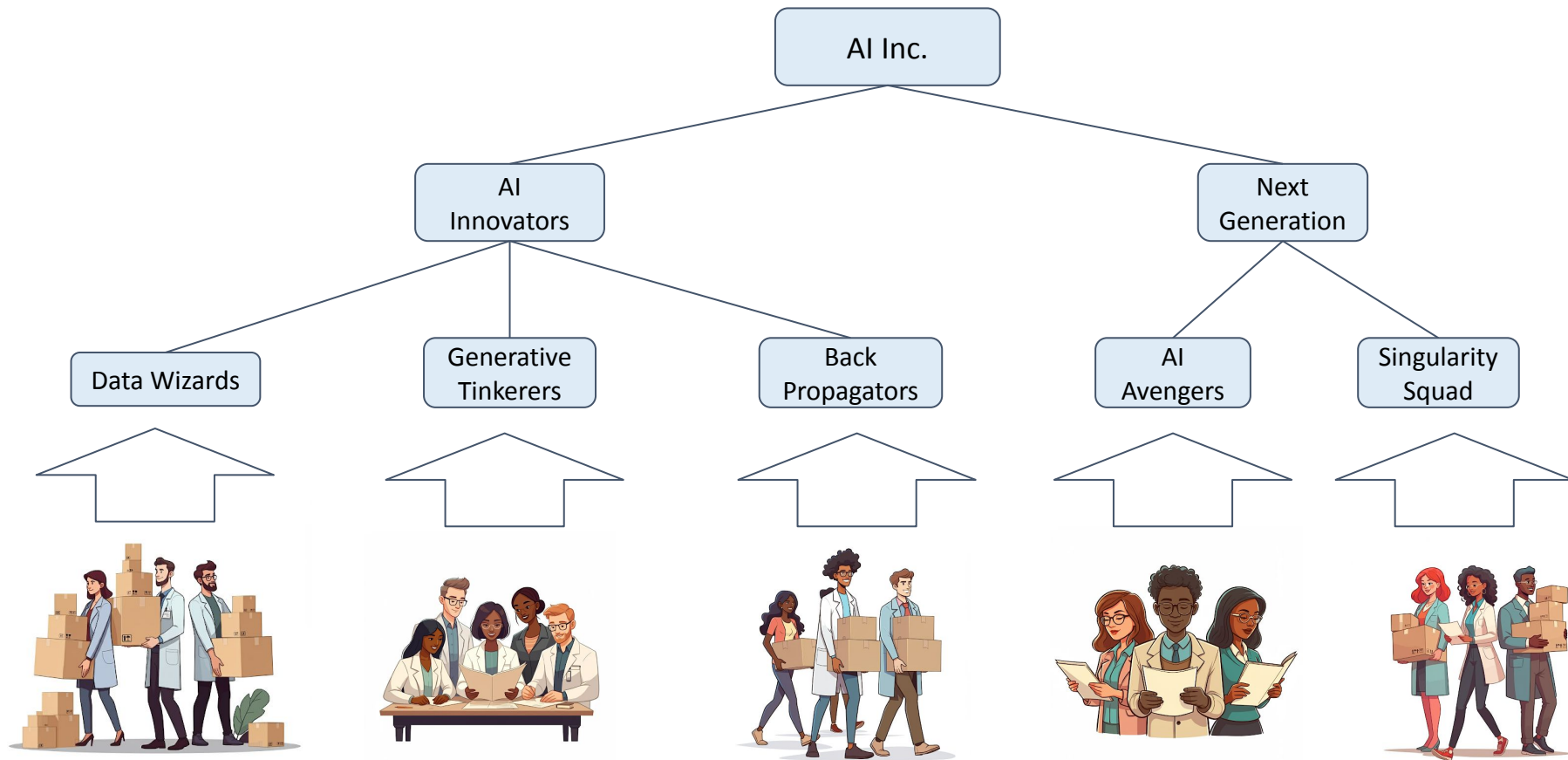
Coming soon!



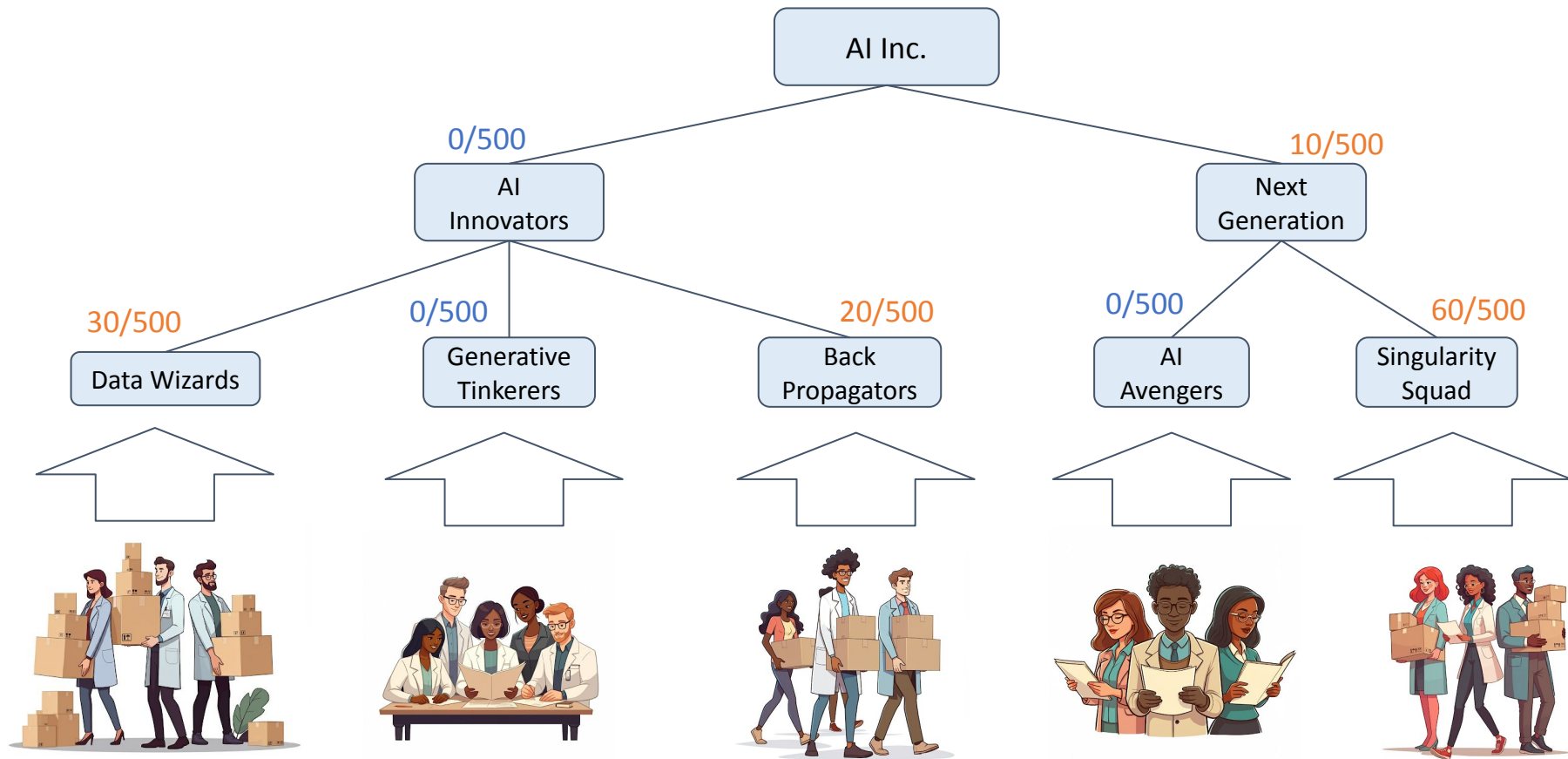
Which job to run and which job to preempt?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



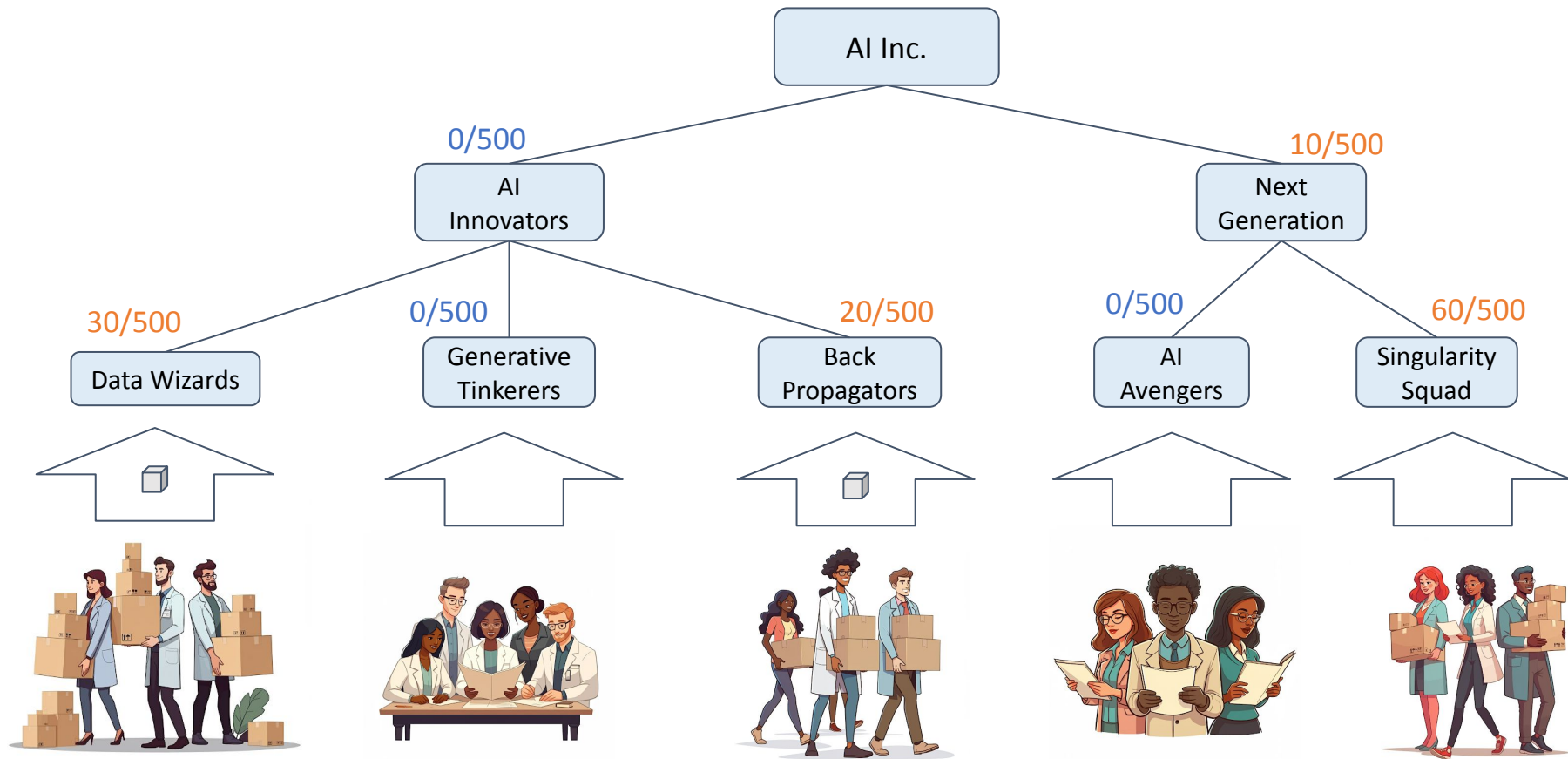
Which job to run and which job to preempt?



Which job to run and which job to preempt?



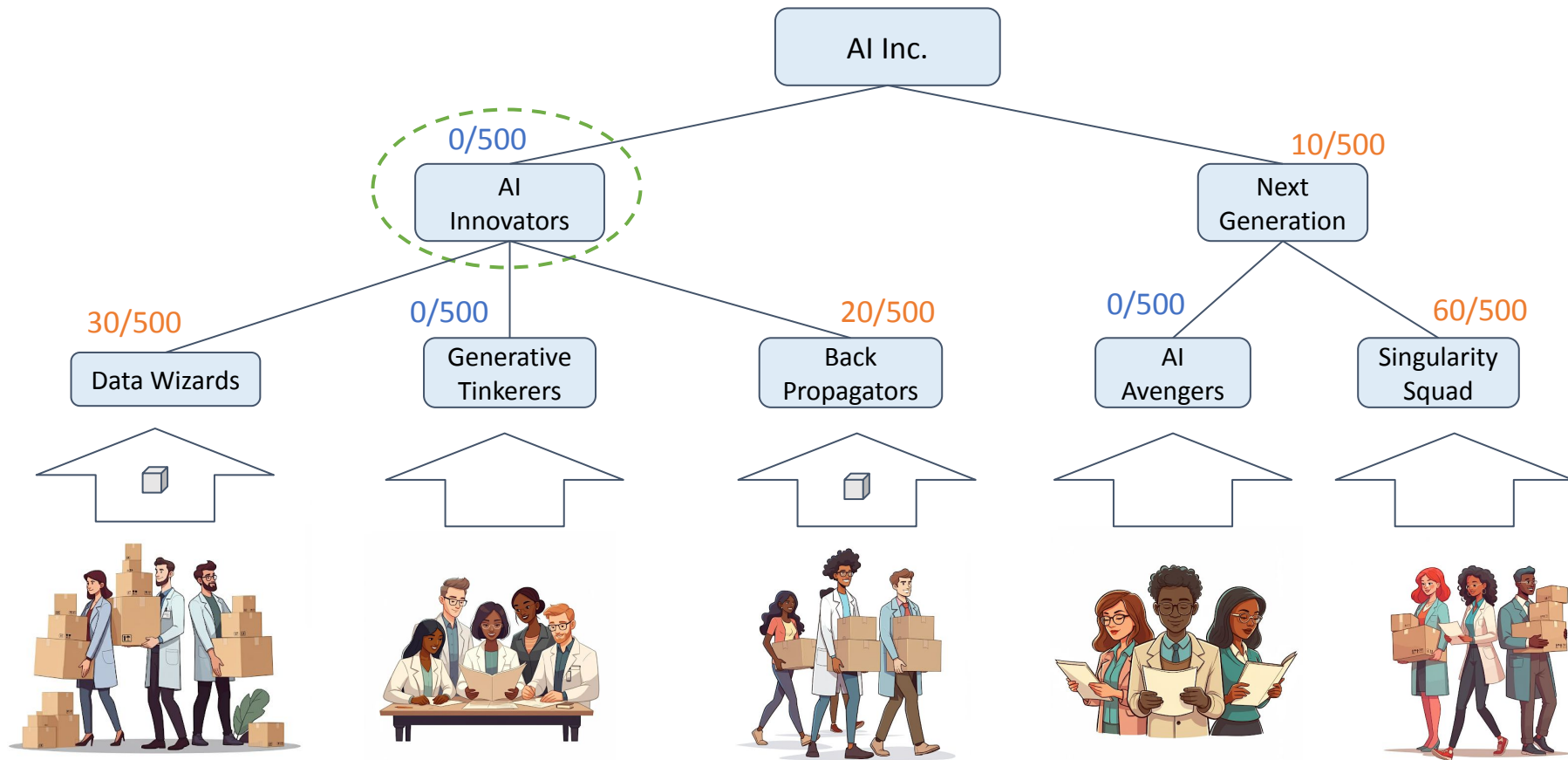
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Which job to run and which job to preempt?



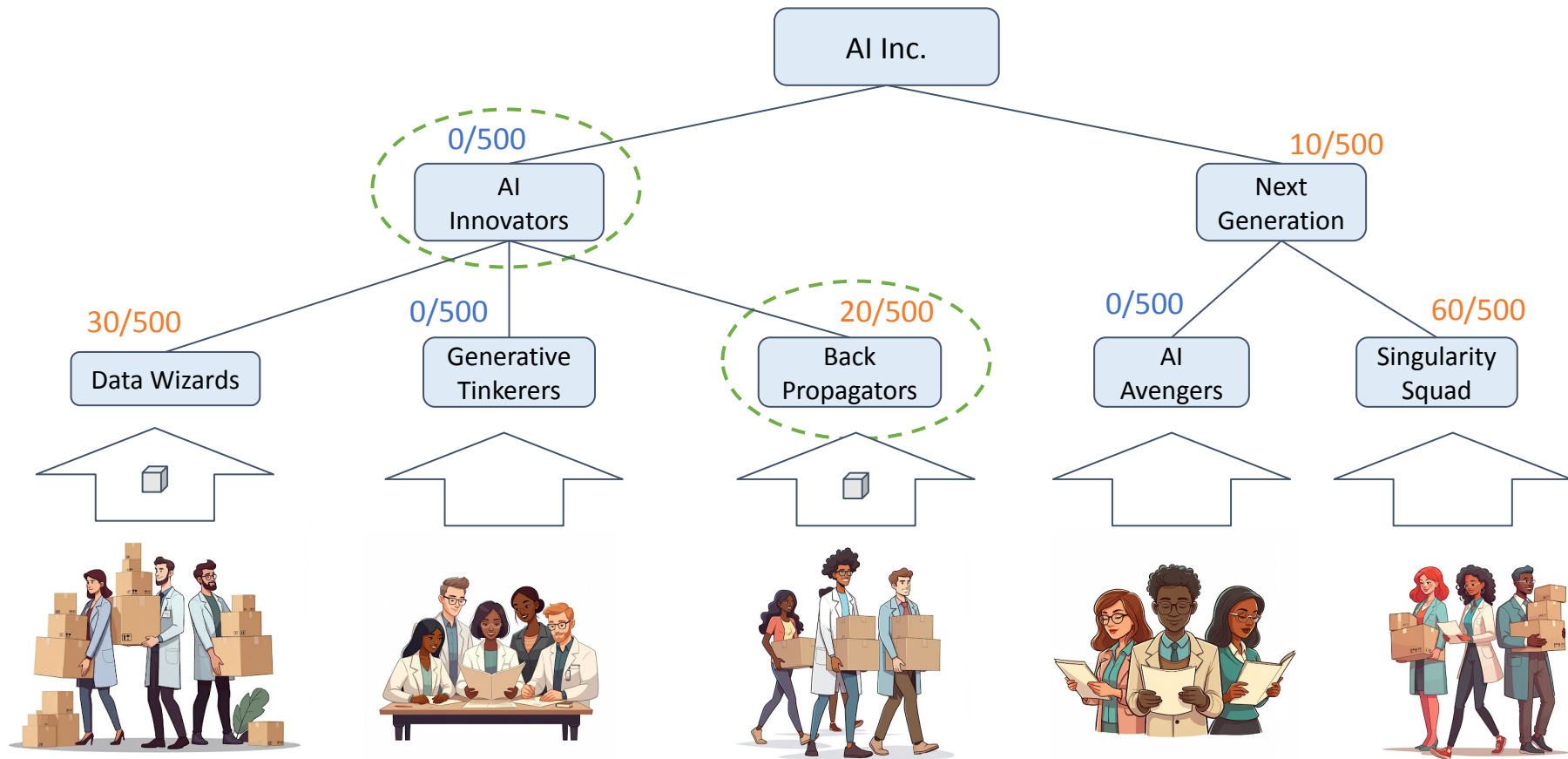
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Which job to run and which job to preempt?



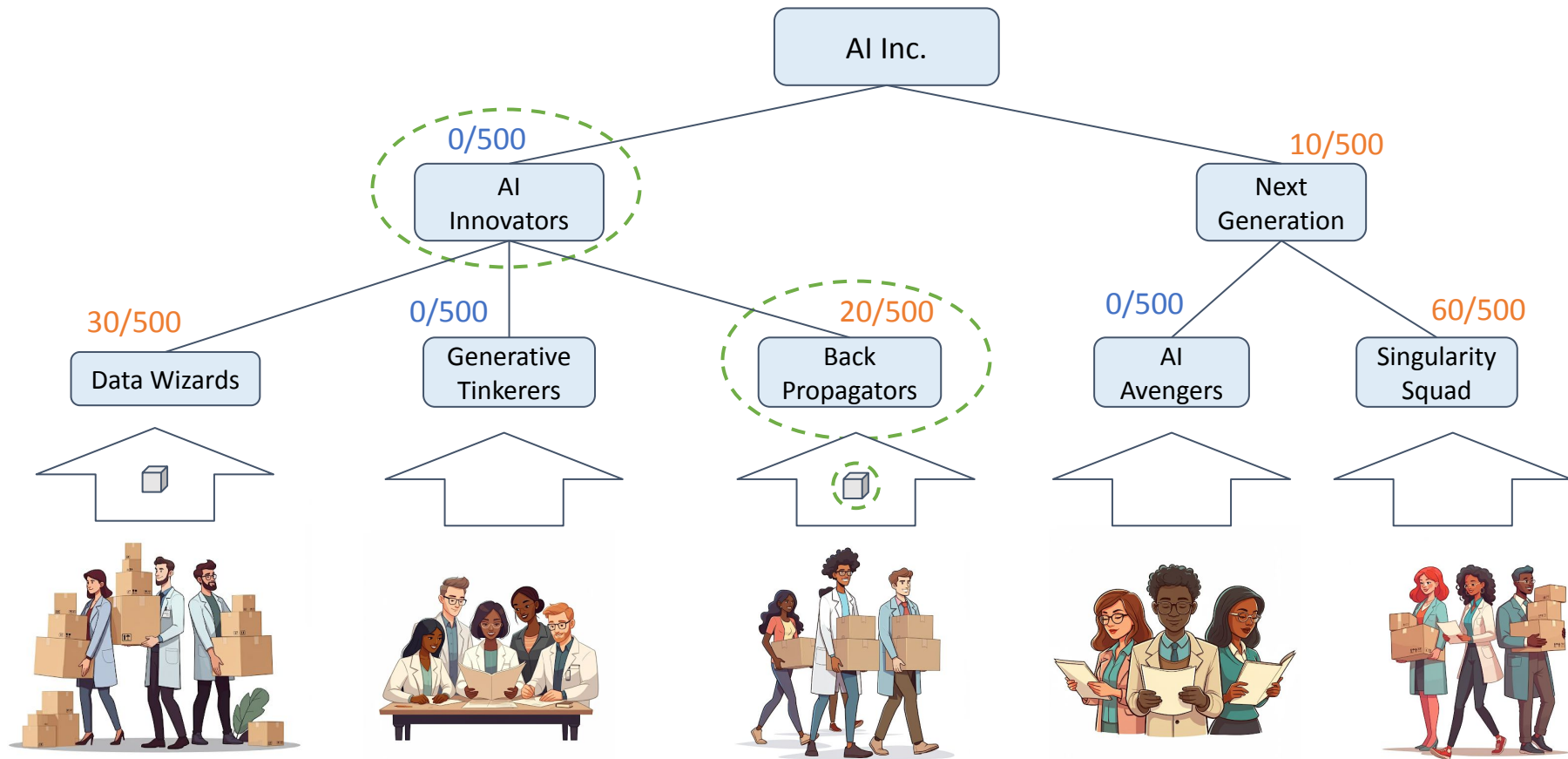
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Which job to run and which job to preempt?



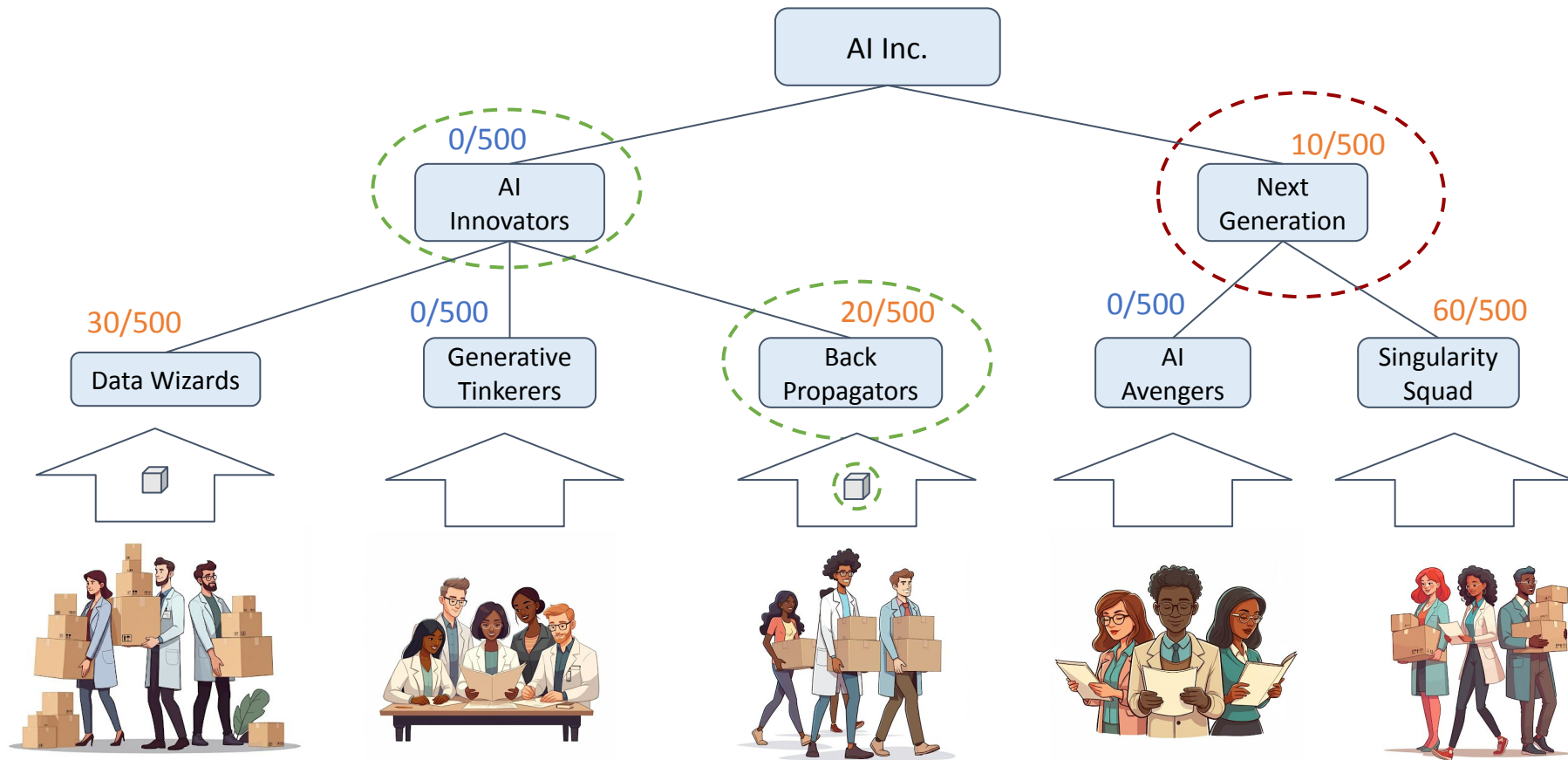
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Which job to run and which job to preempt?



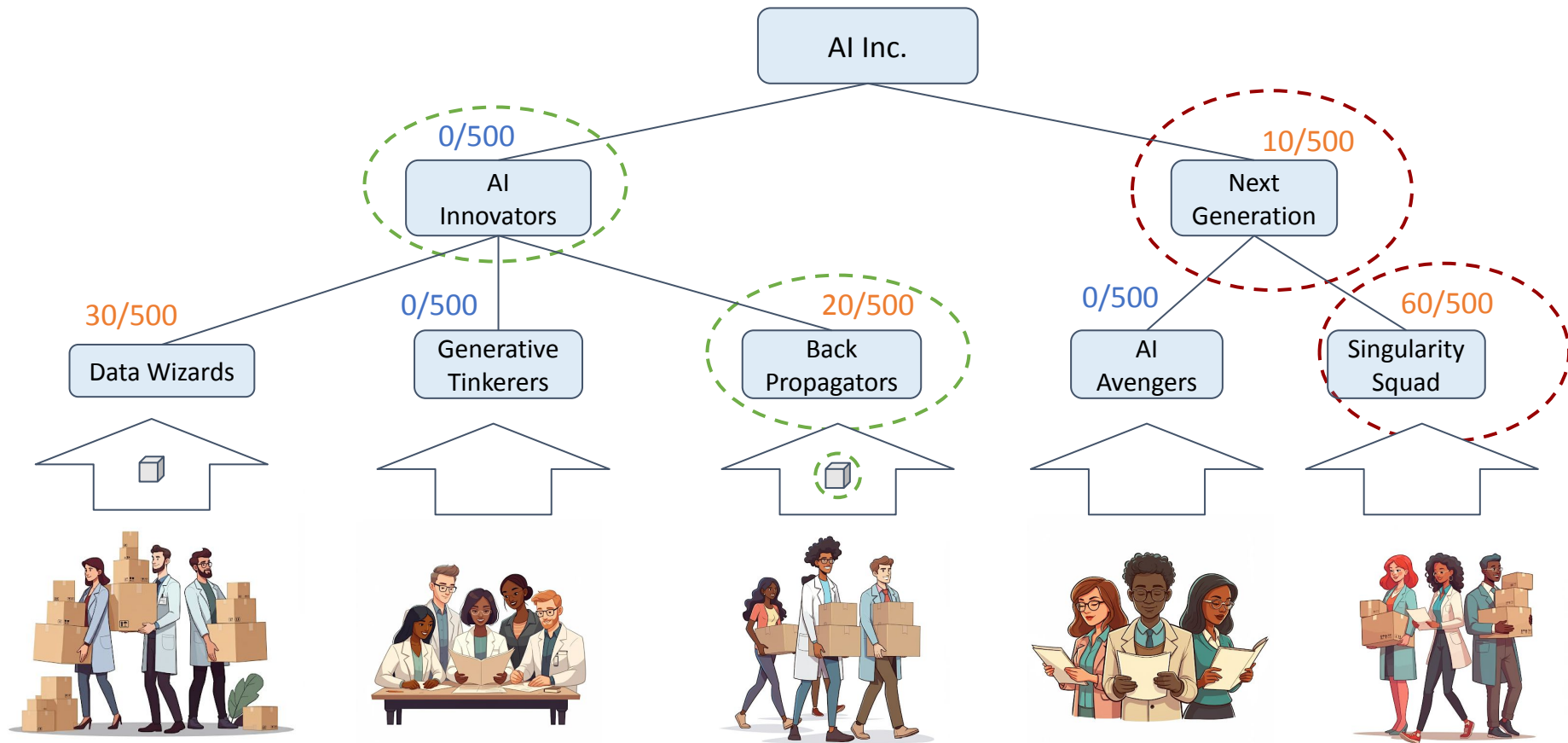
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



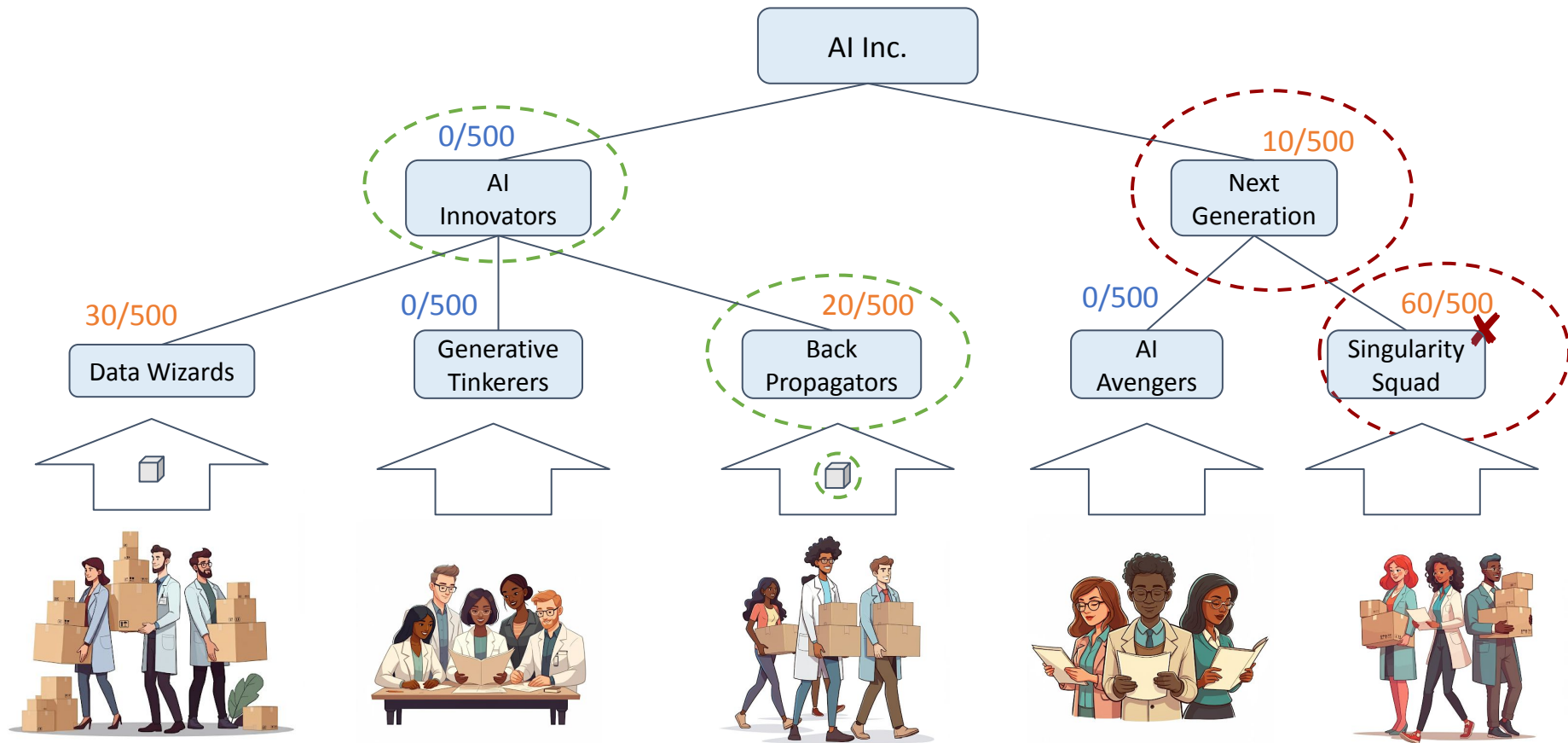
Which job to run and which job to preempt?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA



Which job to run and which job to preempt?

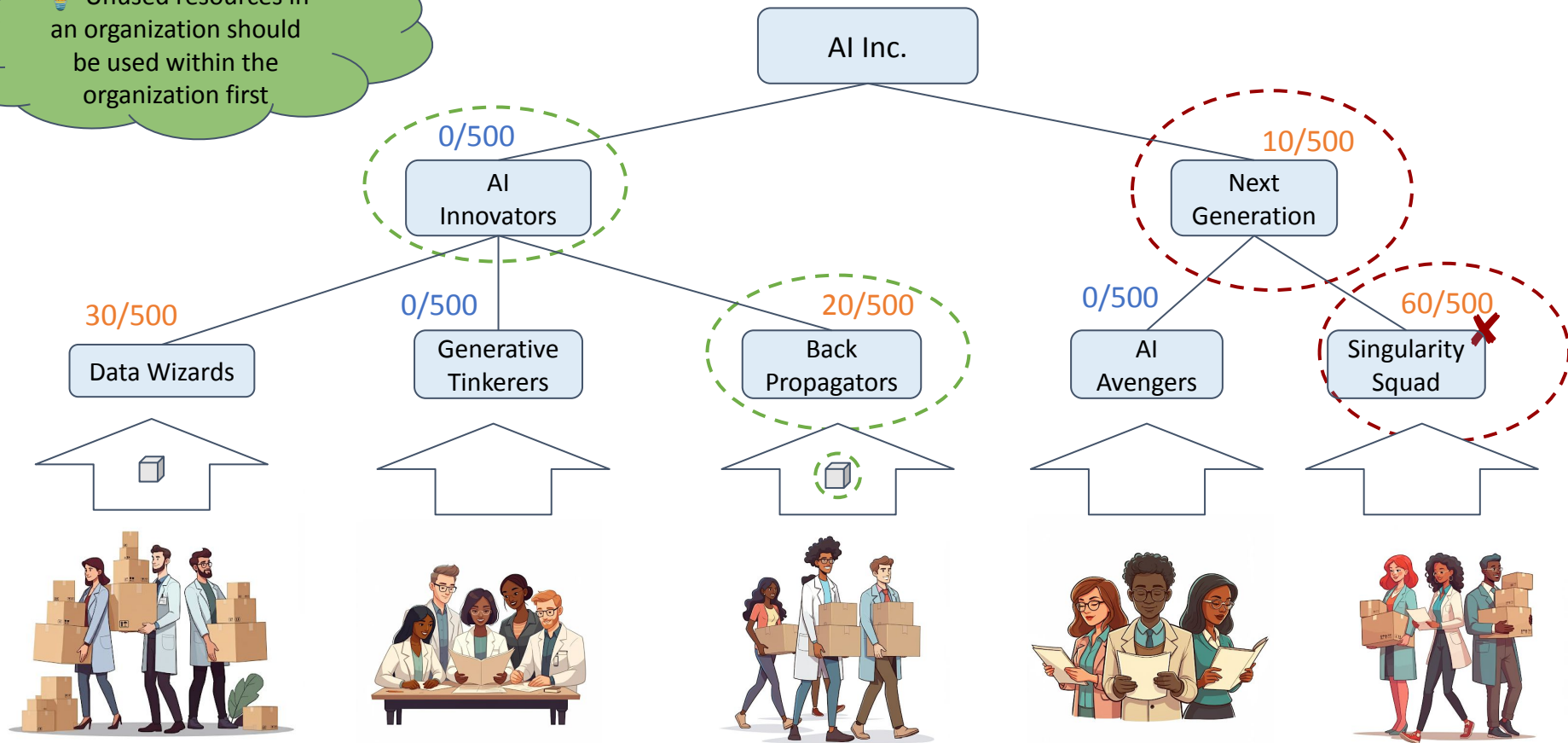


Which job to run and which job to preempt?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

Unused resources in an organization should be used within the organization first



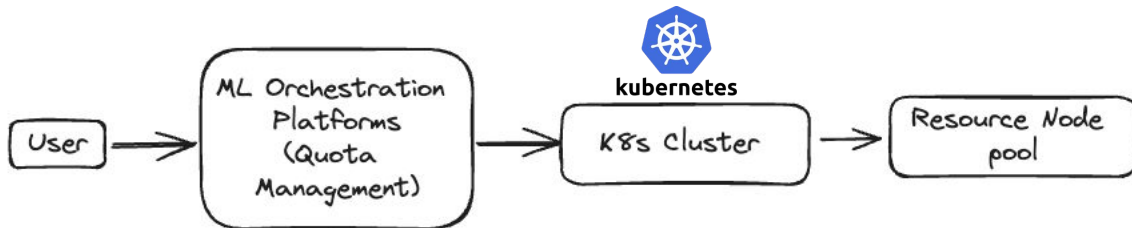
Building an ML Platform

How can Kueue integrate with ML Platforms?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

ML Platforms over
K8s

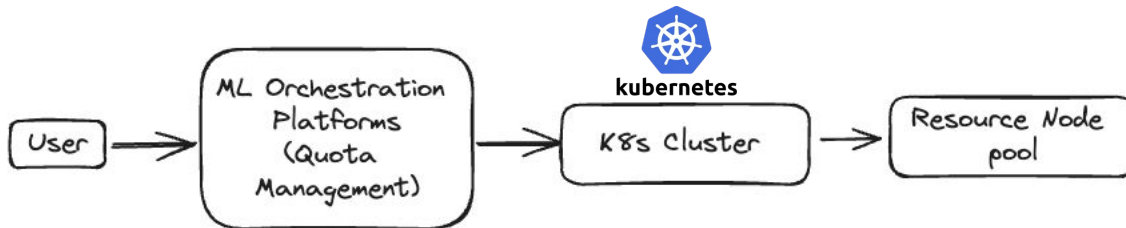


How can Kueue integrate with ML Platforms?



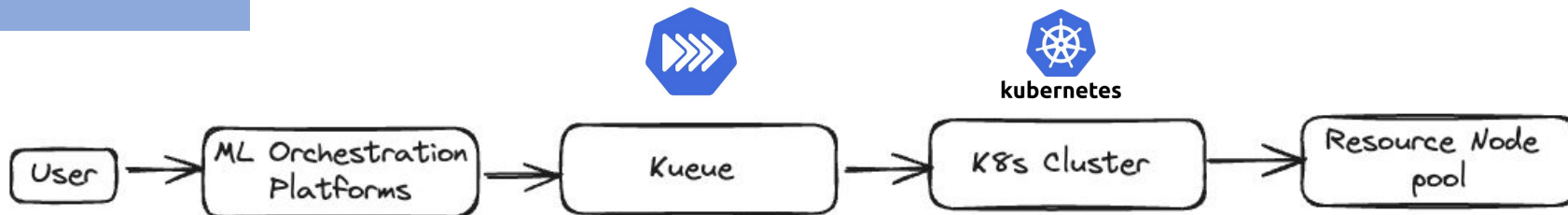
CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

ML Platforms over K8s



ML Platforms with Kueue + K8s

Simplifies quota management and second level scheduling

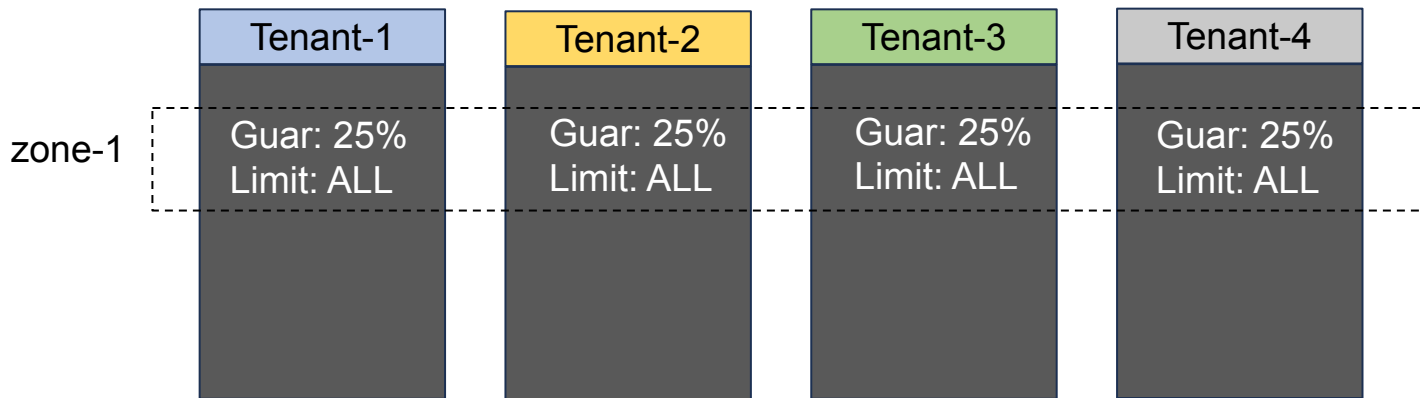


Case Study for ML Platforms



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

- Gang Scheduling and Preemption
- Priority/Reclamation Preemption
- Hierarchy of organizations and teams



Case Study for ML Platforms



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

- Gang Scheduling and Preemption
- Priority/Reclamation Preemption
- Hierarchy of organizations and teams
- Model resource flavors for co-located nodes (zone/spine aware scheduling)
- Support for dynamic change of nominal and borrow limits

	Tenant-1	Tenant-2	Tenant-3	Tenant-4
zone-1	Guar: 25% Limit: ALL	Guar: 25% Limit: ALL	Guar: 25% Limit: ALL	Guar: 25% Limit: ALL
zone-2	Guar: 10% Limit: ALL	Guar: 20% Limit: ALL	Guar: 30% Limit: ALL	Guar: 40% Limit: ALL

Case Study for ML Platforms



- Gang Scheduling and Preemption
- Priority/Reclamation Preemption
- Hierarchy of organizations and teams
- Model resource flavors for co-located nodes (zone/spine aware scheduling)
- Support for dynamic change of nominal and borrow limits
- Burst capacity sharing

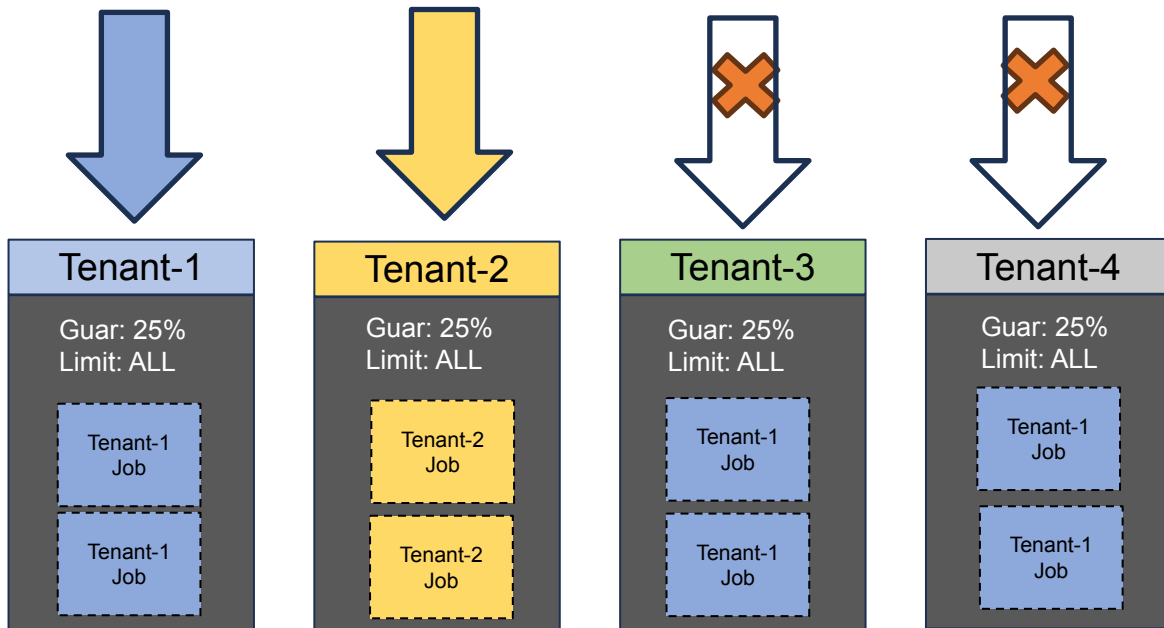
	Tenant-1	Tenant-2	Tenant-3	Tenant-4	Burst
zone-1	Guar: 25% Limit: ALL	Guar: 25% Limit: ALL	Guar: 25% Limit: ALL	Guar: 25% Limit: ALL	Guar: 0% Limit: ALL
zone-2	Guar: 10% Limit: ALL	Guar: 20% Limit: ALL	Guar: 30% Limit: ALL	Guar: 40% Limit: ALL	Guar: 0% Limit: ALL

How to prevent capacity usage abuse?



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

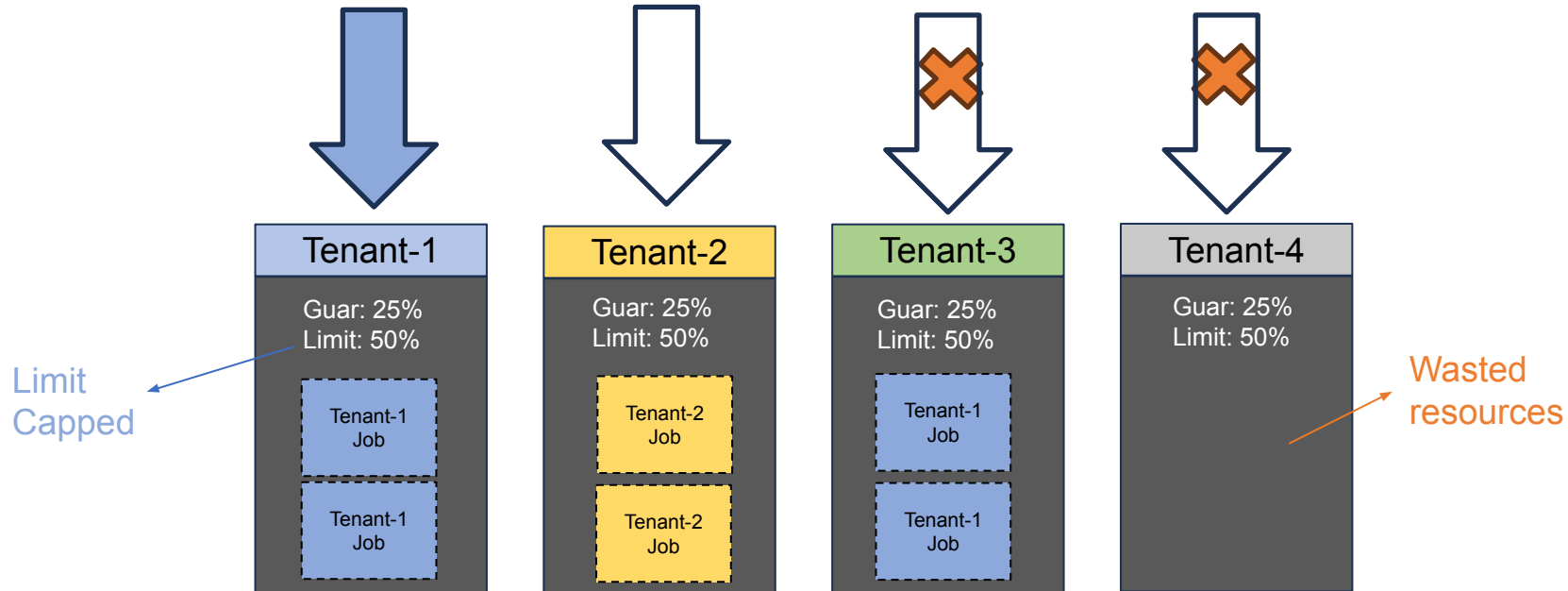
Single Tenant can use all the idle capacity



How to prevent capacity usage abuse?

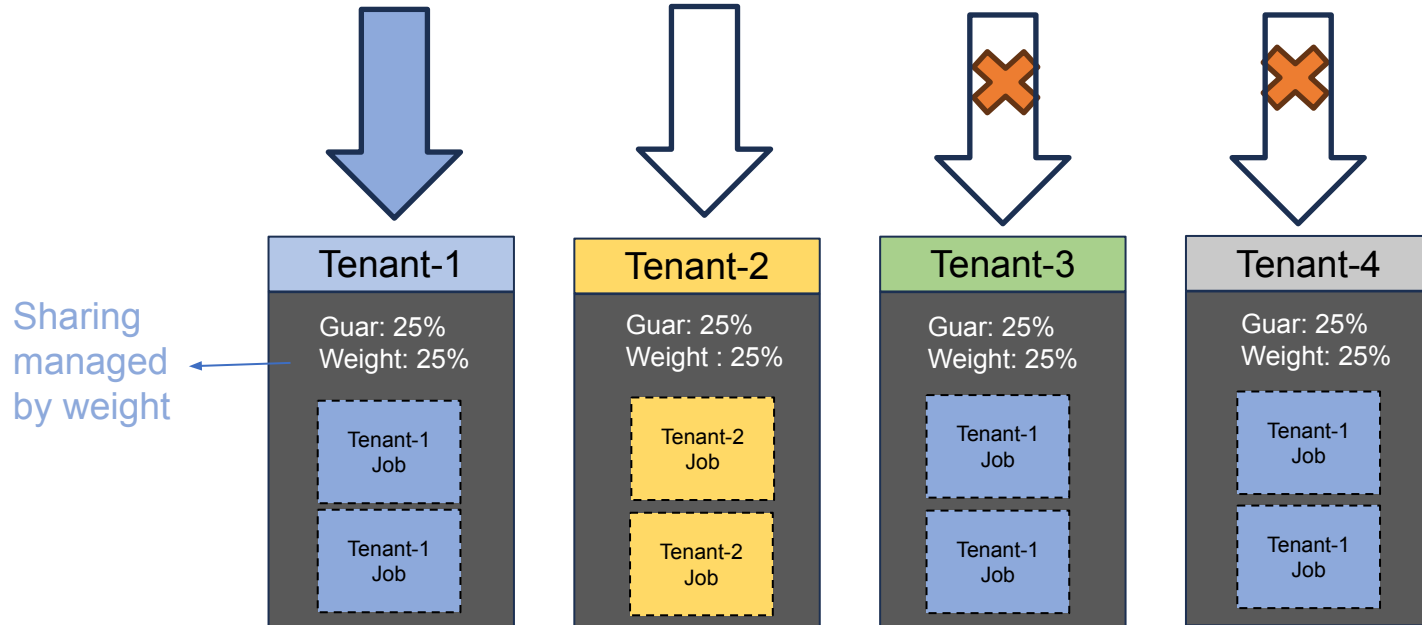
Option 1: Enforce limit per tenant

- Requires frequent adjustments
- May impact utilization



How to prevent capacity usage abuse?

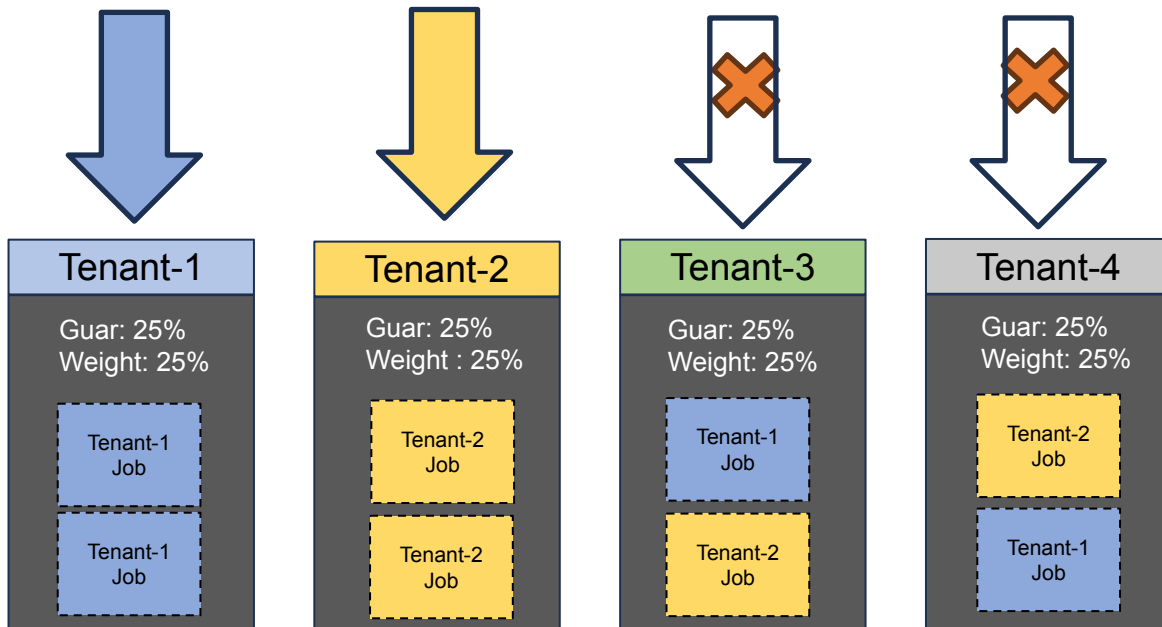
Option 2: Weighted Fair Sharing for idle capacity (No Limits)



Tenant-1 can use all the capacity when others are inactive

How to prevent capacity usage abuse?

Option 2: Weighted Fair Sharing for idle capacity (No Limits)



Tenant-1 gets preempted for fair sharing when Tenant-2 submits

Org sharing with Hierarchies



Accelerators in a Cluster

Org-1 (Cohort)

Guarantees: 50%
Limit: ALL
FS Weight: 40%

Org-1-team-1 (CQ)

Guarantees: 80%
Limit: ALL
FS Weight: 50%

Org-1-team-2 (CQ)

Guarantees: 20%
Limit: ALL
FS Weight: 50%

Org-2 (Cohort)

Guarantees: 50%
Limit: ALL
FS Weight: 40%

Org-2-team-1 (CQ)

Guarantees: 50%
Limit: ALL
FS Weight: 40%

Org-2-team-2 (CQ)

Guarantees: 50%
Limit: ALL
FS Weight: 40%

Burst Only (CQ)

Guarantees: 0%
Limit: ALL
FS Weight: 1%

Kueue v0.7 to v0.9 Highlights



CLOUD NATIVE &
KUBERNETES
AI DAY
NORTH AMERICA

- Fair Sharing in a cohort ^{v0.7}
- More robust and scalable Pod integration ^{v0.7}
- ProvisioningRequest graduated to beta ^{v0.7}
- kueuectl: CLI for cluster admins ^{v0.8}
- Improved preemption throughput and observability ^{v0.8}
- Topology Aware Scheduling ^{v0.9}
- Hierarchical Cohorts ^{v0.9}
- MultiKueue graduated to beta ^{v0.9}
- Support for serving (Deployment and StatefulSet) ^{v0.9}
- Resource Transformations ^{v0.9}



Multitenancy and Fairness at Scale with Kueue: A Case Study



**CLOUD NATIVE &
KUBERNETES**

AI DAY

NORTH AMERICA

November 12, 2024
Salt Lake City



Aldo Culquicondor
Sr Software Engineer
Google



Rajat Phull
Engineering Manager
Apple

Feedback



sched.co/1izqO