



KubeCon



CloudNativeCon

North America 2024





KubeCon



CloudNativeCon

North America 2024

SIG Scheduling Intro & Updates

Aldo Culquicondor
Software Engineer
Google

Kensei Nakada
Software Engineer
Tetrade.io

Responsible for the components that make Pod placement decisions



kube-scheduler



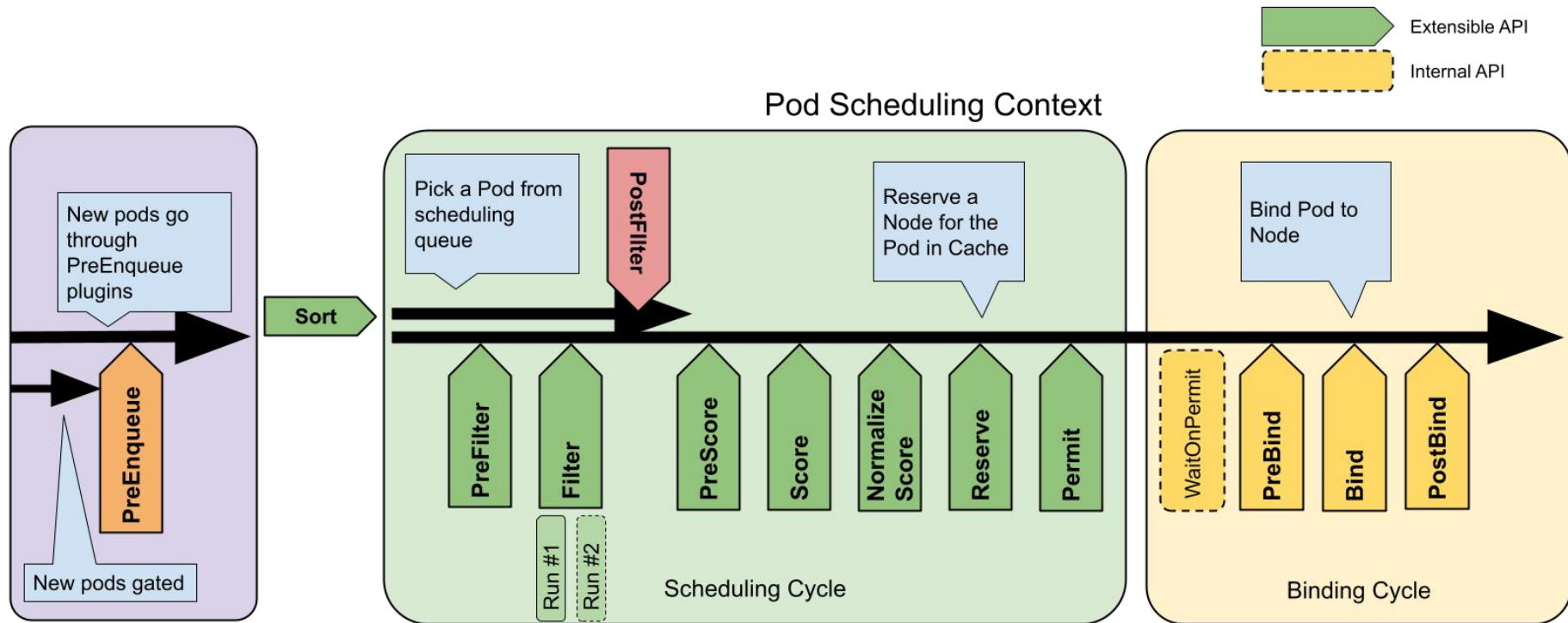
kueue



kwok



kube-scheduler-simulator

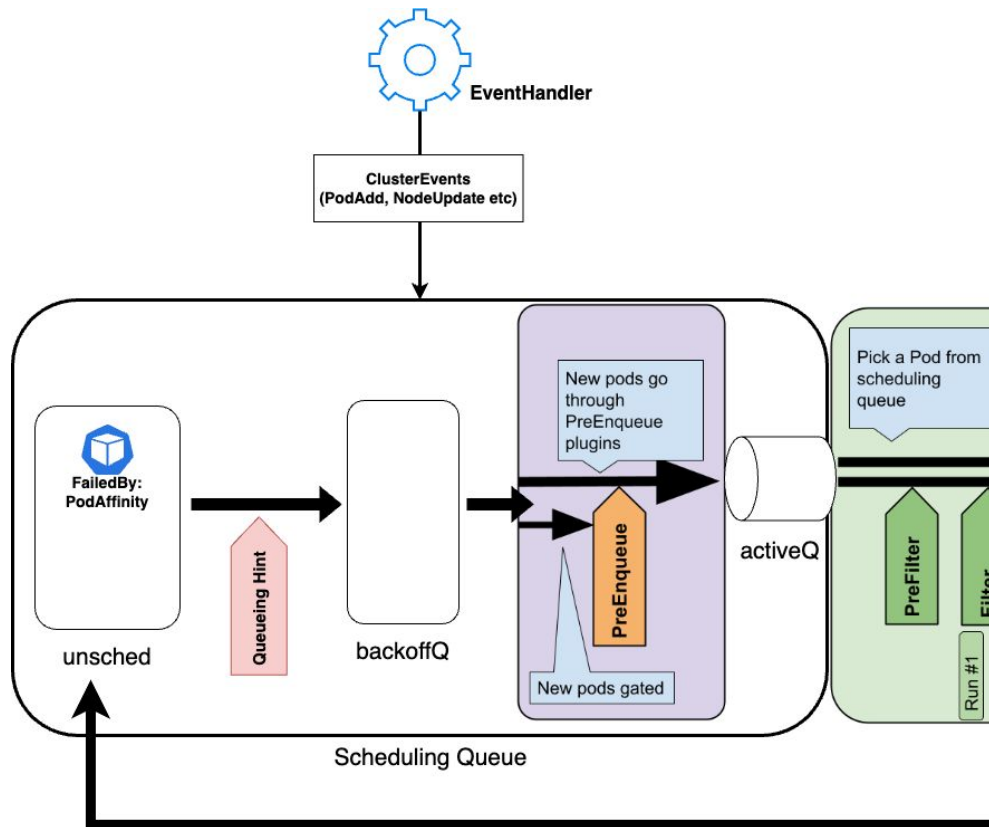


kube-scheduler: Queueing Hints (v1.32: beta)

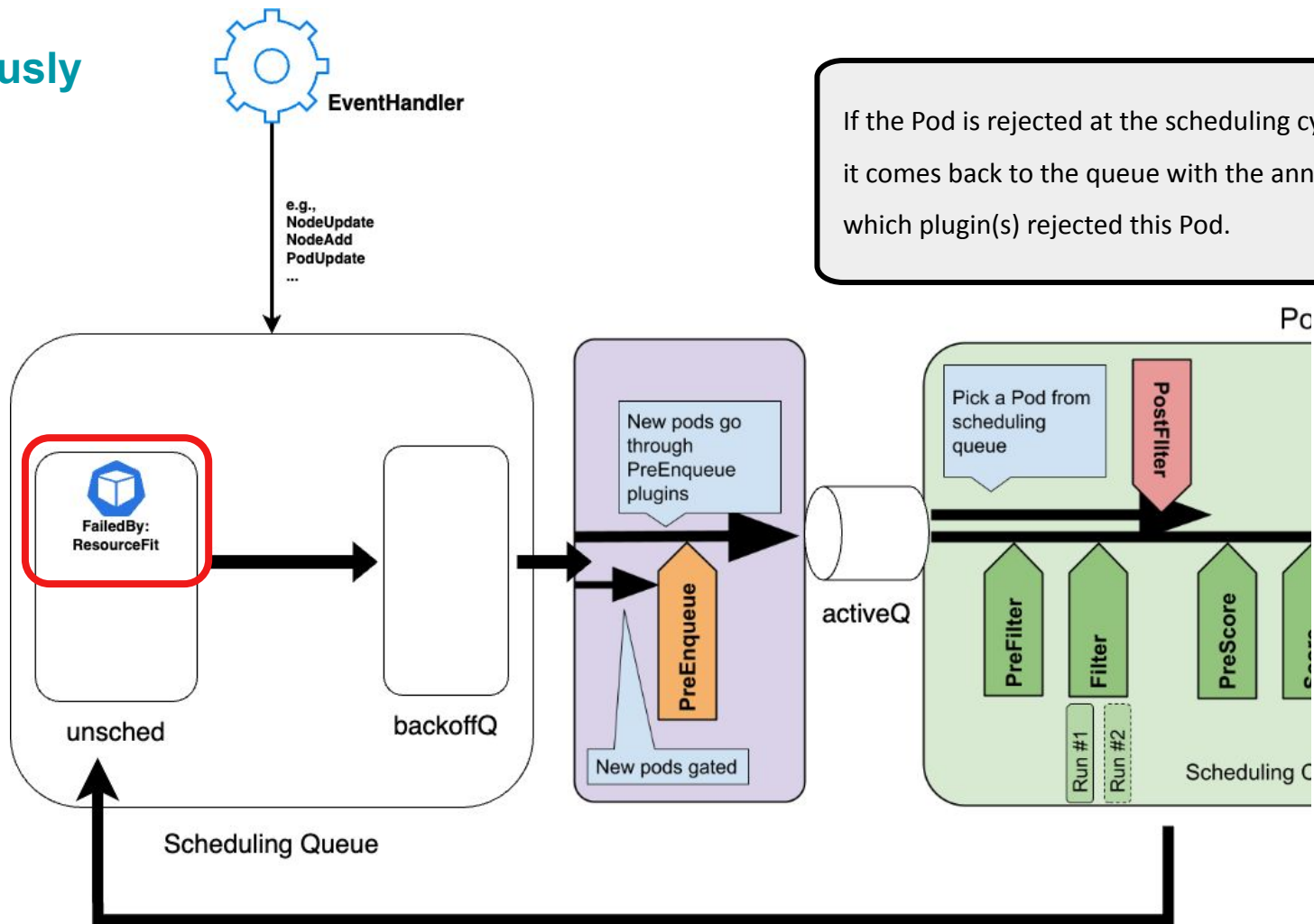
💬 Motivation

Cluster events can make previously unschedulable Pods schedulable. How to decide whether it's worth retrying a particular Pod?

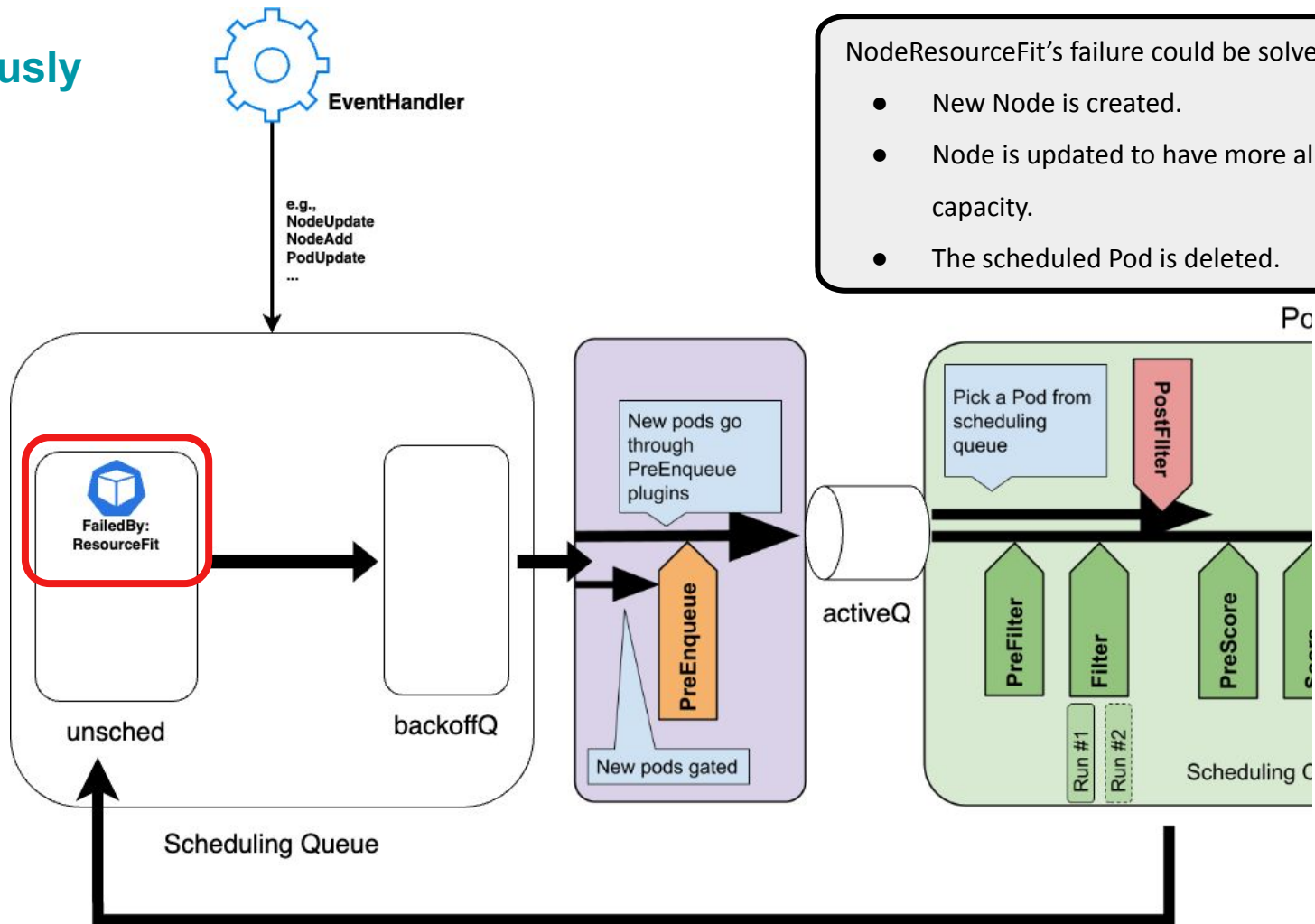
Existing mechanism was too coarse and not extensible for out-of-tree plugins.



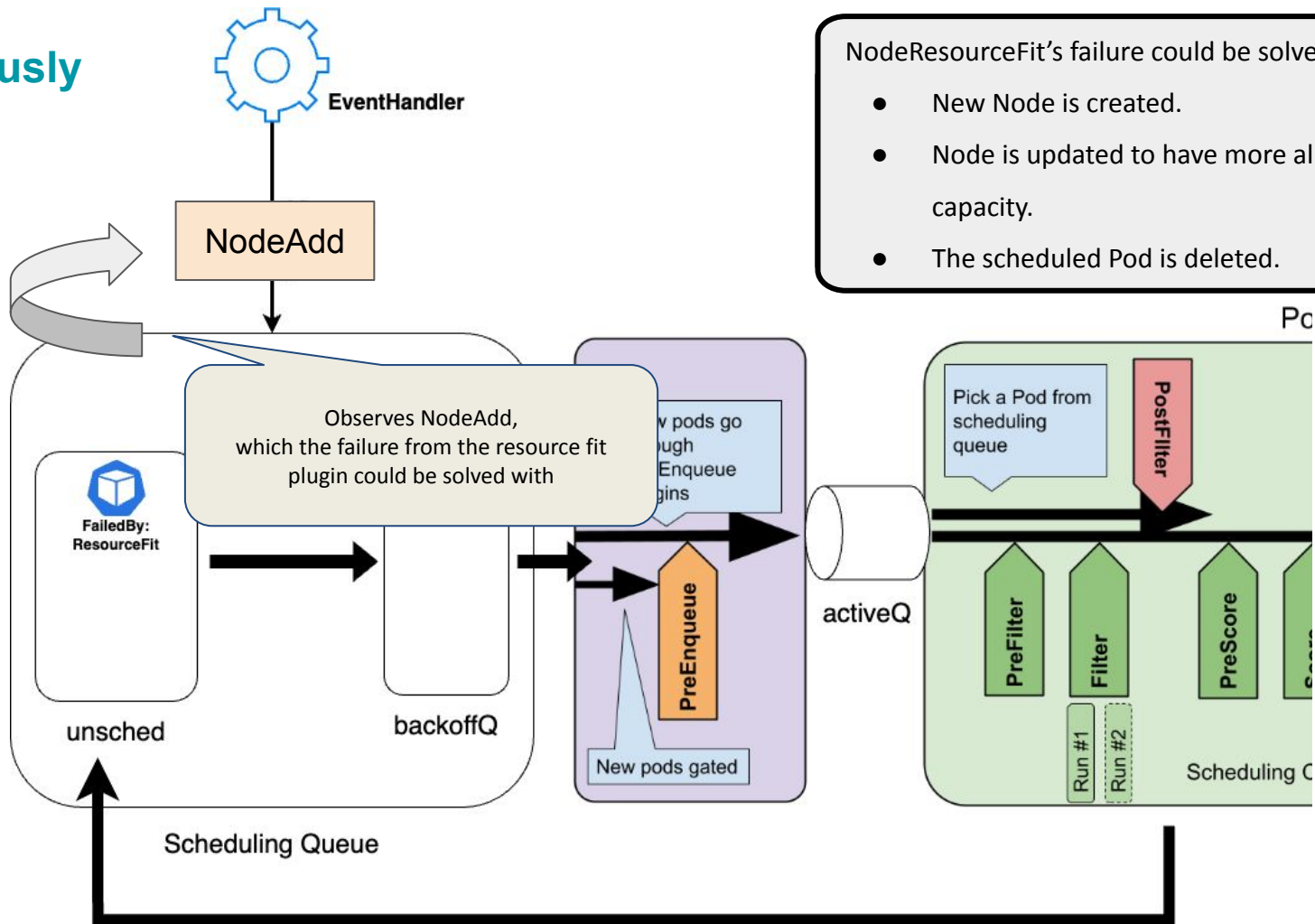
Previously



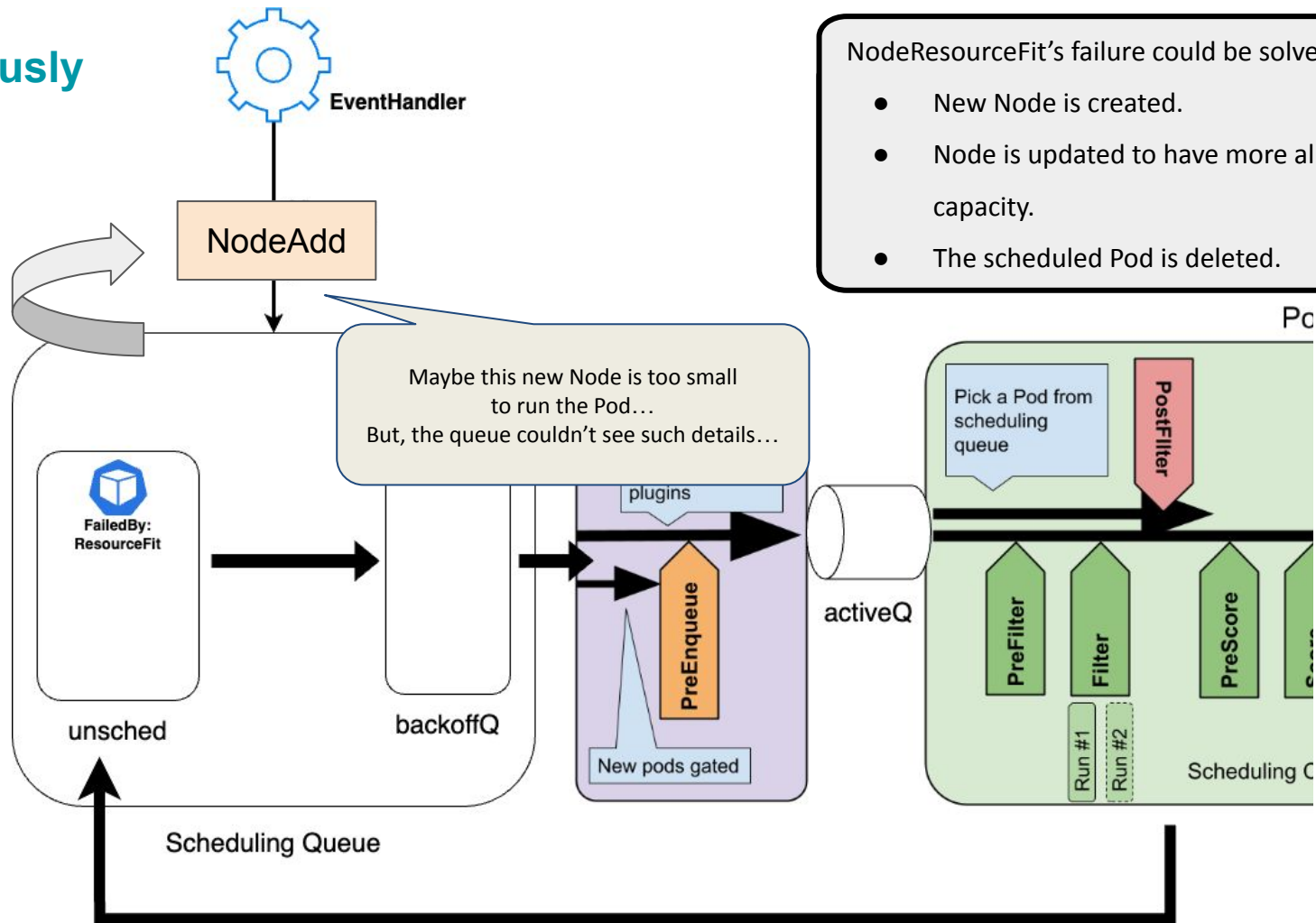
Previously



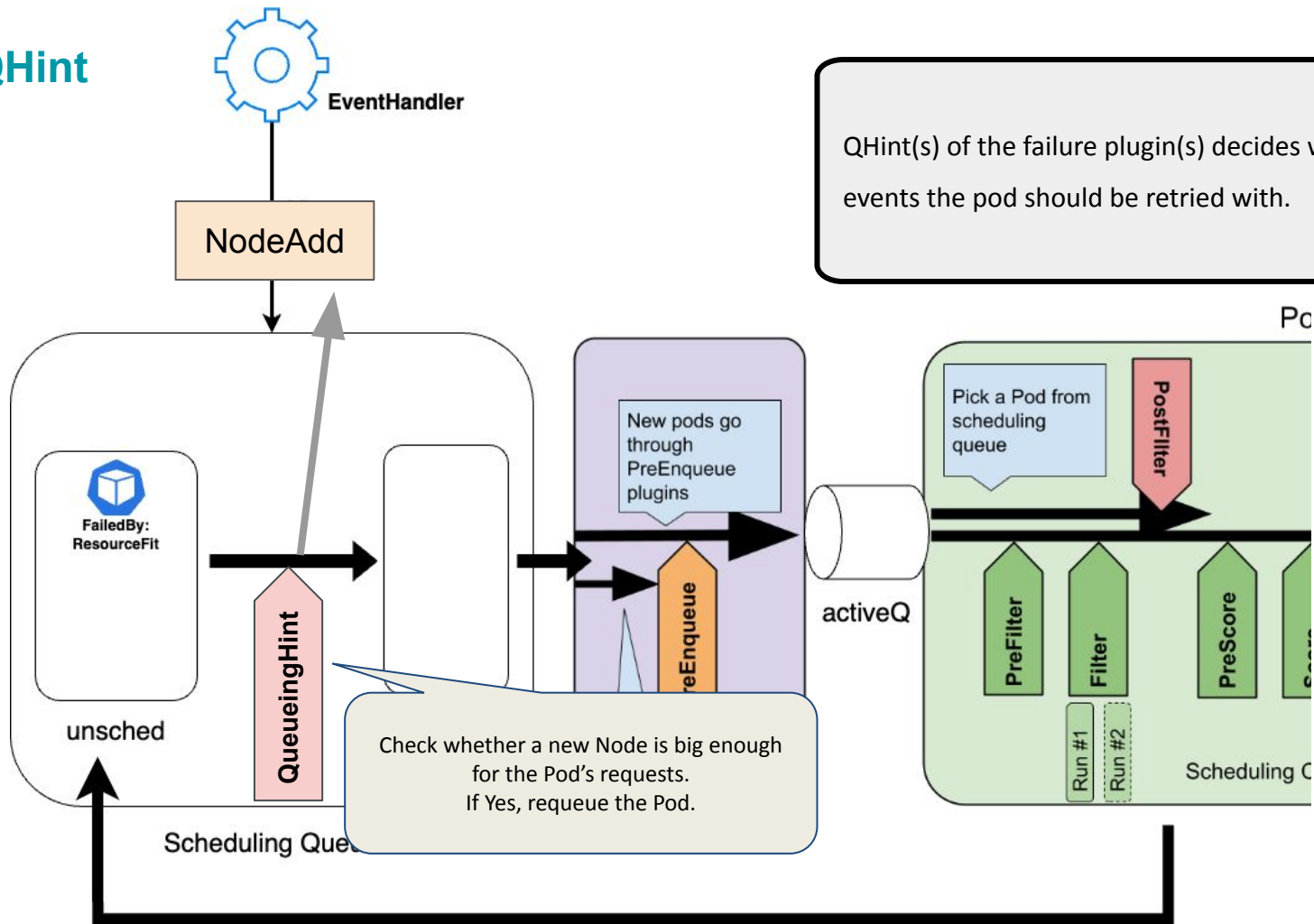
Previously



Previously



After QHint

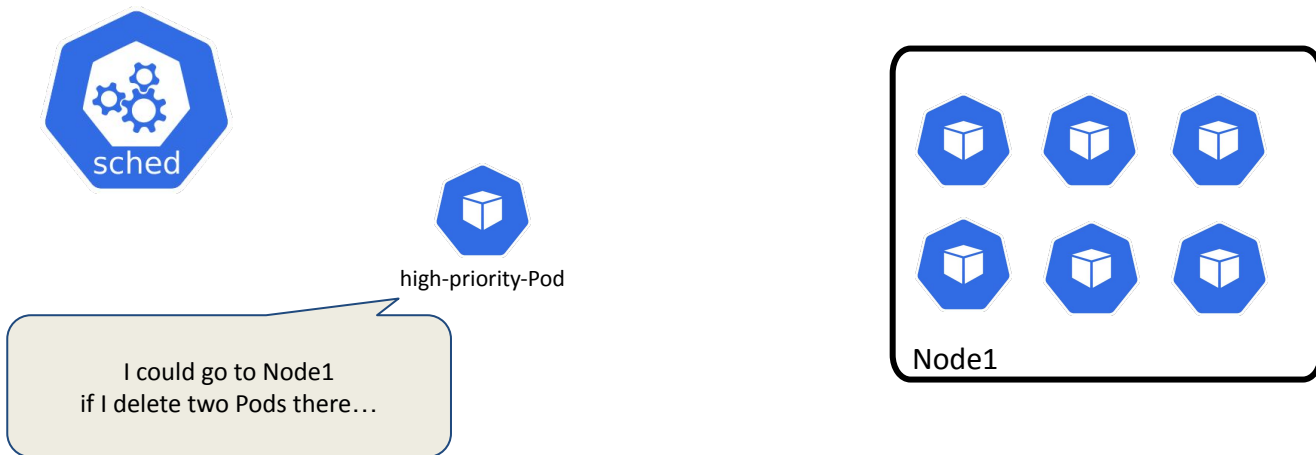


- v1.28:
 - The feature is released as Beta (enabled by default) with a few QHint implementation (DRA etc)
 - 🔥 Users observed a memory leak issue.
 - We disabled the feature by default in a patch release.
- v1.29-v1.31: We kept working on expanding the QHints and doing bug fixes.
- v1.32:
 - We finished implementing QHints in all plugins.
 - We identified and fixed the last major memory leak, finally!
 - We increased the integration test/performance test coverage... a lot.
 - 🚀 We enabled the feature by default.

kube-scheduler: Async preemption (v1.32: alpha)

💬 Motivation

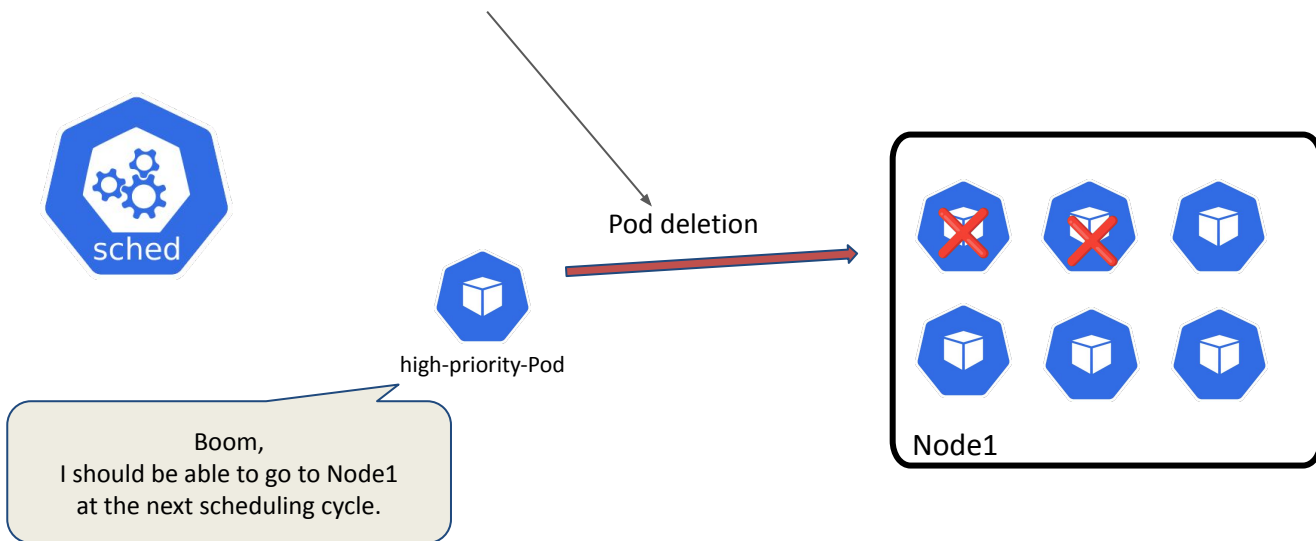
If your pods get unschedulable, they go through **the preemption** before going back to the queue. When the preemption happens, the scheduling cycle takes time to complete because it has to make some API calls. It impacts the whole scheduling latency.



kube-scheduler: Async preemption (v1.32: alpha)

💬 Motivation

If your pods get unschedulable, they go through **the preemption** before going back to the queue. When the preemption happens, the scheduling cycle takes time to complete, because it has to make some API calls. It impacts the whole scheduling latency.



kube-scheduler: Async preemption (v1.32: alpha)

💬 Motivation

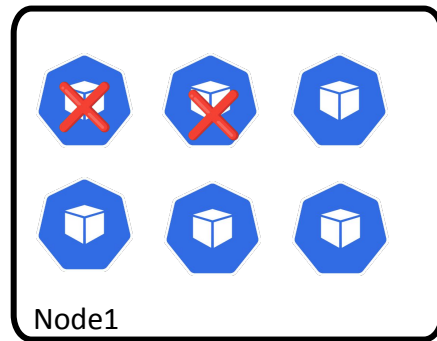
If your pods get unschedulable, they go through **the preemption** before going back to the queue. When the preemption happens, the scheduling cycle takes time to complete, because it has to make some API calls. It impacts the whole scheduling latency.



Starts the next scheduling
after the preemption is completed.



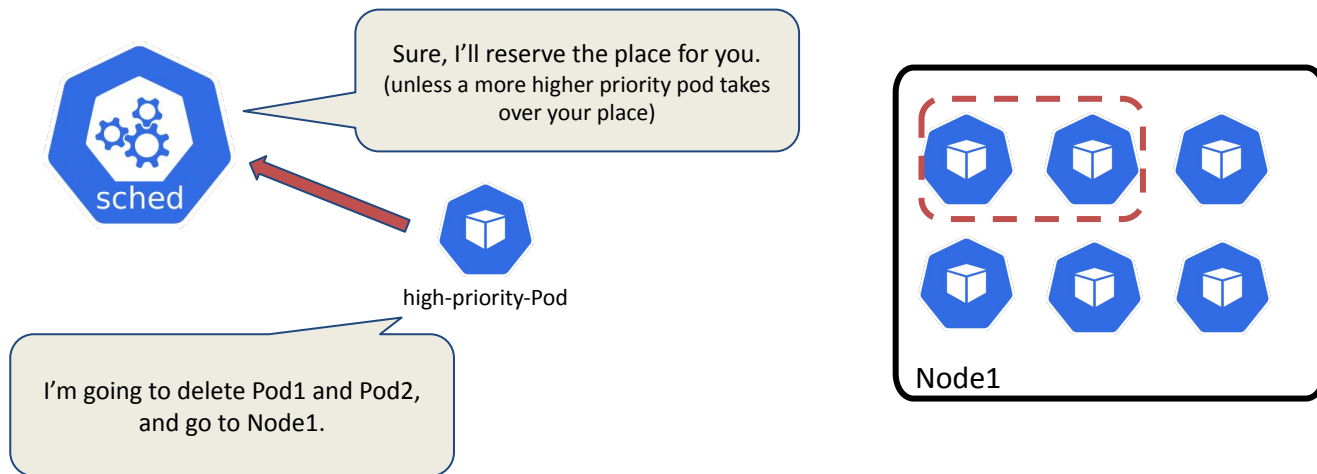
Pod deletion



kube-scheduler: Async preemption (v1.32: alpha)

💡 Proposal

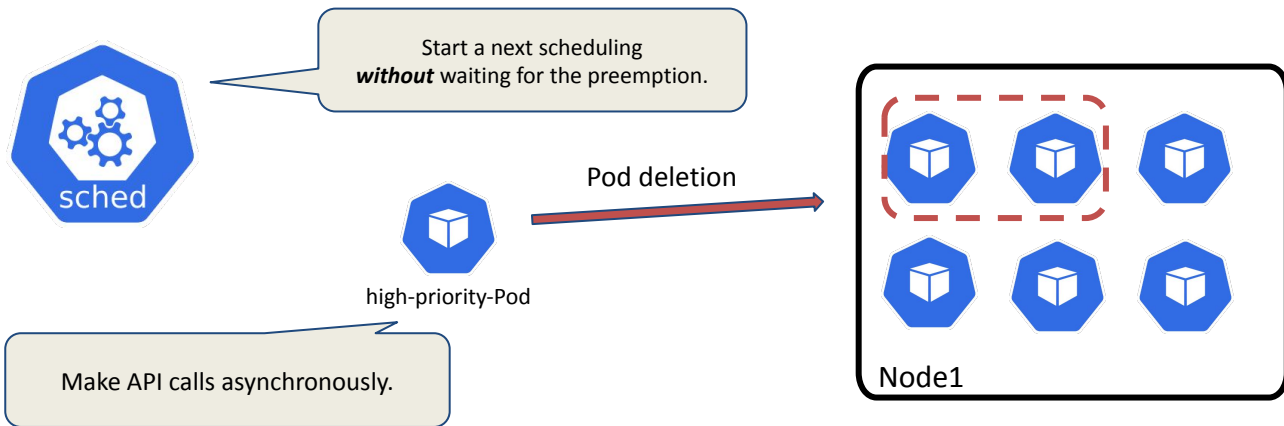
Once deciding which Pod(s) to delete, **making the API calls asynchronously**, and starting the next scheduling cycle **without waiting for the preemption** to be completed. The preemption asks the scheduler to reserve the preemptor's place on the node before starting the next scheduling cycle so that the next scheduling cycle takes the preemption into consideration.



kube-scheduler: Async preemption (v1.32: alpha)

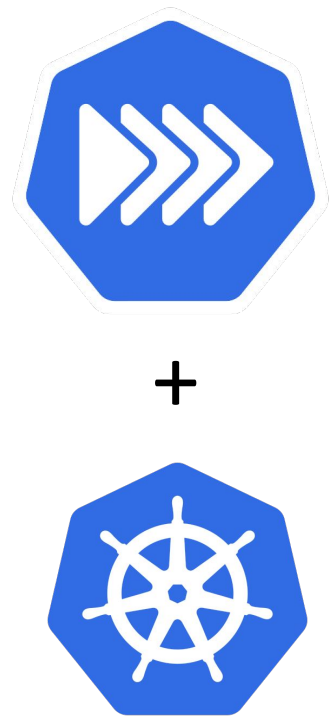
💡 Proposal

Once deciding which Pod(s) to delete, **making the API calls asynchronously**, and starting the next scheduling cycle **without waiting for the preemption** to be completed. The preemption asks the scheduler to reserve the preemptor's place on the node before starting the next scheduling cycle so that the next scheduling cycle takes the preemption into consideration.

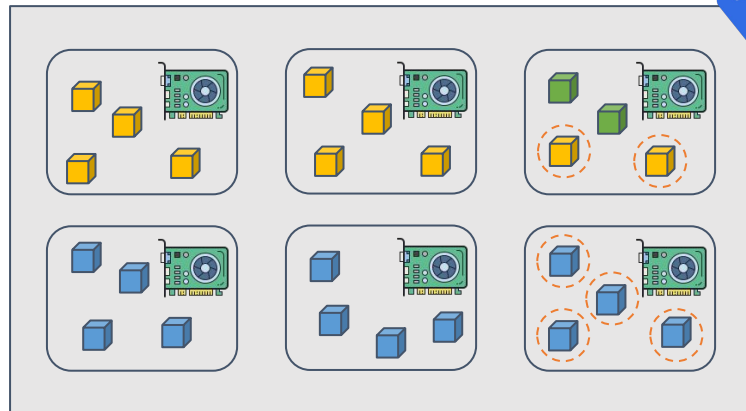
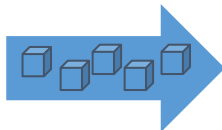
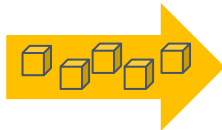


- v1.30:
 - Pod Scheduling readiness gates graduated to GA
 - minDomains for Pod topology spreading graduated to GA
- v1.31:
 - matchLabelKeys in Pod affinity and anti-affinity graduated to Beta
- v1.32:
 - DRA structured parameters graduated to Beta
 - Classic DRA was removed

- 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.



Kueue: Fair sharing



```
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
```

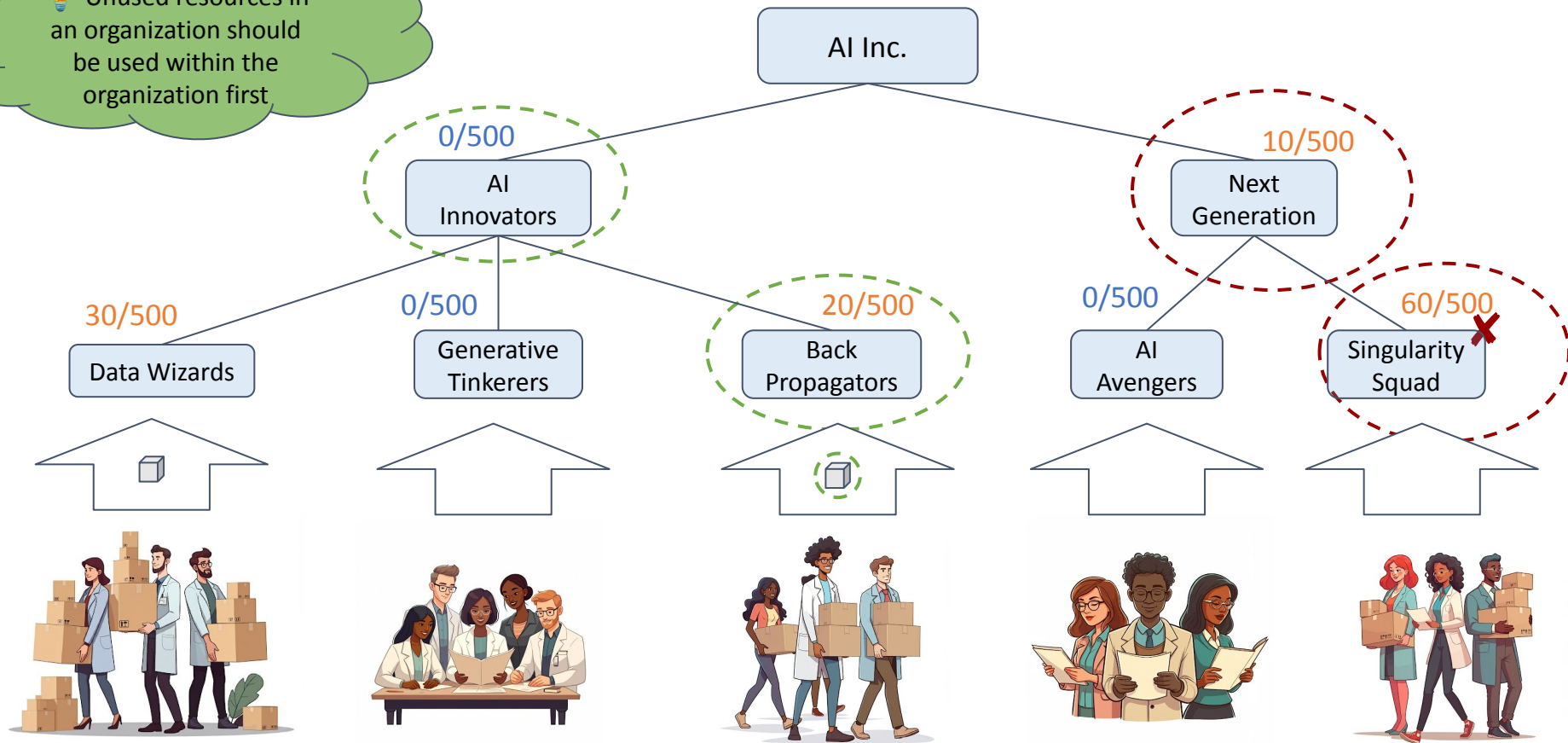
```
kind: ClusterQueue
metadata:
  name: "alice-queue"
spec:
  cohort: "lab"
  resourceGroups:
  - coveredResources: ["acme.com/gpu"]
    flavors:
    - name: "default-flavor"
      resources:
      - name: "acme.com/gpu"
        nominalQuota: 8
  preemption:
    reclaimWithinCohort: "Any"
    fairSharing: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



sched.co/1izqO

Kueue: Hierarchical cohorts

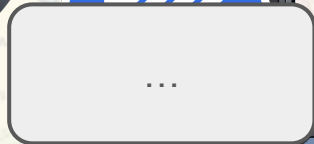
Unused resources in an organization should be used within the organization first



And it works in
multiple clusters!



sched.co/1iCOV





*Kubernetes **With**Out Kubelet*

The toolkit to set up a cluster of thousands of Nodes for scheduling simulations... in seconds!

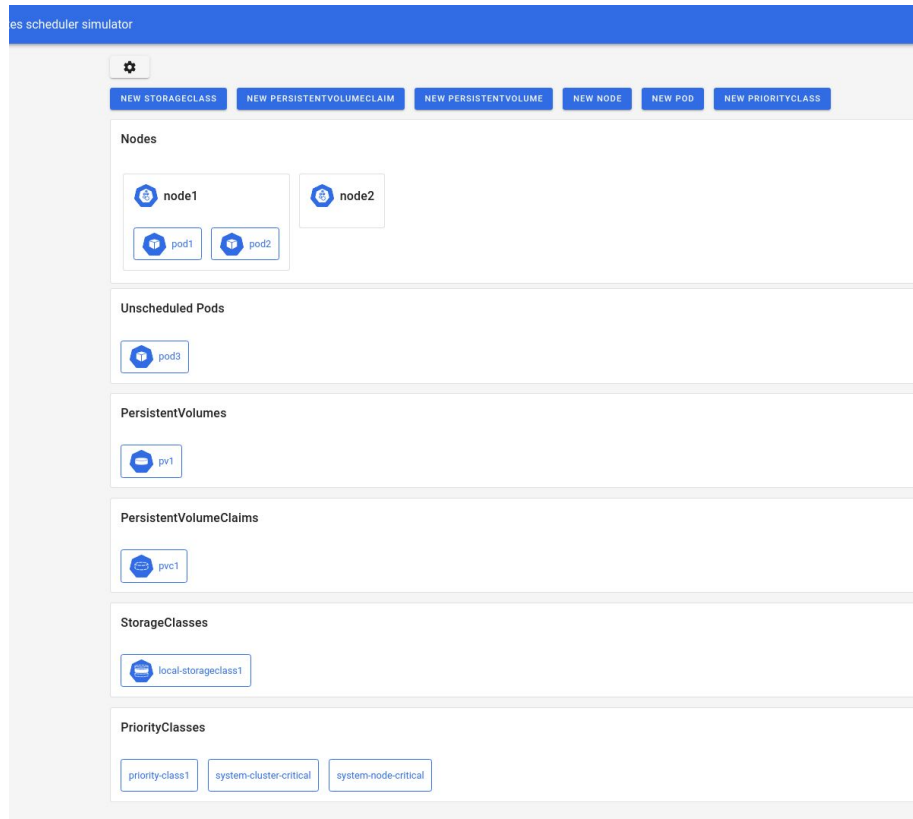
What's new?

- Optimized Stage to allow easier simulation of normal/abnormal behavior of resources
- Optimized the experience of CPU and MEM simulation on metrics
- Helm chart installation available
- kwokctl
 - Support for MacOS and Windows
 - More runtimes such as lima and finch
- Roadmap:
 - Plan API v1alpha2
 - Enhance the experience of using kwok on existing clusters without using kwokctl
 - Extend simulation capabilities to include GPU usage and other custom metrics
 - Simulation for Volume Provisioner

Kube-scheduler-simulator offers the debuggability for your scheduler.

- The scheduler is composed of the plugins.
- The simulator visualizes all the internal decisions for each pod's scheduling.

You can create resources
on Web UI



Kube-scheduler-simulator offers the debuggability for your scheduler.

- The scheduler is composed of the plugins.
- The simulator visualizes all the internal decisions for each pod's scheduling.

It shows how your scheduler decides the Node for your Pod.

The screenshot displays the 'Kubernetes scheduler simulator' interface. On the left, a sidebar shows 'Nodes' with 'node1' and 'pod1'. The main area is divided into three sections: 'Filter', 'Score', and 'Final Score (Normalized + Applied plugin weight)'. The 'Filter' section shows a table of node selection criteria. The 'Score' section shows a table of scores for each node. The 'Final Score' section shows a table of final scores. Below these is the 'Resource Definition' section, which shows the pod's metadata and specifications.

Node	AzureDiskLimits	EBSLimits	GCEPDLimits	InterPodAffinity	NodeAffinity	NodeName	NodePorts	NodeResourcesFit	NodeResourcesFit
node1	passed	passed	passed	passed	passed	passed	passed	passed	passed
node2	passed	passed	passed	passed	passed	passed	passed	passed	passed

Node	ImageLocality	InterPodAffinity	NodeAffinity	NodeResourcesBalancedAllocation	NodeResourcesFit	PodTopologySpread	TaintToleration
node1	0	0	0	76	73	0	0
node2	0	0	0	76	73	0	0

Node	ImageLocality	InterPodAffinity	NodeAffinity	NodeResourcesBalancedAllocation	NodeResourcesFit	PodTopologySpread	TaintToleration
node1	0	0	0	76	73	0	0
node2	0	0	0	76	73	0	0

Resource Definition

```
metadata:
  name: pod1
  namespace: default
  uid: 5eb1b8ad-2ac6-47c6-bf5c-ef3ef4e606ea
  resourceVersion: 183
  creationTimestamp: 2021-08-16T05:32:12Z
spec:
  containers:
  - name: container1
    image: nginx
    ports:
    - containerPort: 80
```


We've been adding new features while refactoring the simulator to be more maintainable.

- v0.2.0: Migrate the internal Kube-apiserver and controllers to Kwok.
- v0.3.0: Introduce the syncing feature:
 - *Testing a scheduler is hard* => Run scheduling simulator using production pods.
 - The syncing feature keeps importing the resources from your prod cluster to simulate your prod environment at the simulator, with potentially different configuration.
- v0.4.0: The standalone debuggable scheduler:
 - The debuggable scheduler schedules pods like the normal scheduler but also outputs all the scheduling internal decisions as pod annotations.
 - It now works standalone: you can use it in the dev cluster instead of the normal scheduler so you can debug low-level scheduling decisions.



sched.co/1hovV