

# SIG-Scheduling Deep-Dive

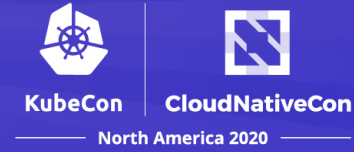
- Day 1/2/3 of operating kube-scheduler

Wei Huang, IBM, @Huang-Wei

Abdullah Gharaibeh, Google, @ahg-g



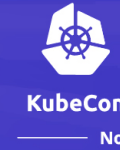
# Lead-in



*Virtual*

- A thousands ways of view Scheduler
- Day 1, Day 2 and Day 3 of operating Kubernetes scheduler

# Day 1 – App developer



*Virtual*

North America 2020

- **Audience:** users who write & deploy containerized application onto k8s
- Goal: Understand kube-scheduler basics, identify whether it's a scheduler issue, and how to use scheduler features
- Non-Goal: Understand scheduler internals

# Day 1 – Kube-scheduler



KubeCon



CloudNativeCon

North America 2020

*Virtual*

- Kube-scheduler: Assign Pods to Nodes

```
wei.huang1@wei-mbp:~|⇒ k get pod -o wide
```

NAME	READY	STATUS	RESTARTS	AGE	IP	NODE	NOMINATED	NODE	READINESS	GATES
k8s-for-beginners-5b5b757bdb-97g4t	0/1	Pending	0	71s	<none>	<none>	<none>		<none>	

```
wei.huang1@wei-mbp:~|⇒ k get po k8s-for-beginners-5b5b757bdb-97g4t -o jsonpath="{.spec.nodeName}"
wei.huang1@wei-mbp:~|⇒
```



```
wei.huang1@wei-mbp:~|⇒ k get po -o wide
```

NAME	READY	STATUS	RESTARTS	AGE	IP	NODE	NOMINATED	NODE	READINESS	GATES
k8s-for-beginners-5b5b757bdb-97g4t	1/1	Running	0	11m	10.244.2.4	kind-worker2	<none>		<none>	

```
wei.huang1@wei-mbp:~|⇒ k get po k8s-for-beginners-5b5b757bdb-97g4t -o jsonpath="{.spec.nodeName}"
kind-worker2%
```

# Day 1 – Scheduler Do's and Don'ts



*Virtual*

- kube-scheduler do's and don'ts

- ❌ Quota enforcement
- ❌ Spinning up / scaling down replicas of a Deployment/StatefulSet/etc.
- ❌ Evict Pods upon OutOf {Memory|Disk|CPU} error
- ❌ Taint nodes (kubectl taint node, or upon {Memory|Disk|CPU} pressure)
- ❌ Reschedule / Rebalance running Pods<sup>1</sup>

Find the best Node for pending Pods

Preempt low-priority Pods to make room for high-priority Pods

<sup>1</sup> A sub-project of sig-scheduling covers it: <https://github.com/kubernetes-sigs/descheduler>

# Day 1 – Filtering



KubeCon



CloudNativeCon

North America 2020

*Virtual*

- Filtering – **Hard** Constraints, i.e., I **need** my Pod to:
  - To have 2Gi memory and 1 core CPU
  - Co-exist with some kinds of Pods
  - To tolerate taints with effect "NoSchedule"
  - ...
- All hard constraints are **ANDed**
- (Almost) All **hard** constraints are from pod's spec – i.e., specified by the user

```
spec:
  affinity:
    podAffinity:
      requiredDuringSchedulingIgnoredDuringExecution:
        - labelSelector:
            matchExpressions:
              - key: foo
                operator: Exists
            topologyKey: region
```

```
spec:
  containers:
    - name: test-pod
      image: k8s.gcr.io/pause:3.2
      resources:
        requests:
          cpu: 1
          memory: 2Gi
```

```
spec:
  topologySpreadConstraints:
    - maxSkew: 1
      topologyKey: kubernetes.io/hostname
      whenUnsatisfiable: DoNotSchedule
      labelSelector:
        matchLabels:
          foo: ""
```

# Day 1 – Scoring



KubeCon



CloudNativeCon

North America 2020

Virtual

- Scoring – **Soft** Constraints, i.e., I prefer my Pod to:
  - To be scheduled to a node which has SSDs
  - Not to co-exist with some kinds of Pods
  - ...
- Based on the soft constraints, each filtered Node gets a Score
- sum(score) for each filtered Node, and pick the highest score
- **Soft** constraints have 2 sources
  - Pod's spec
  - Implicit scheduler config, e.g., NodeResourcesLeastAllocated

```
spec:
  affinity:
    nodeAffinity:
      preferredDuringSchedulingIgnoredDuringExecution:
        - weight: 100
          preference:
            matchExpressions:
              - key: disktype
                operator: In
                values:
                  - ssd
```

```
spec:
  topologySpreadConstraints:
    - maxSkew: 1
      topologyKey: kubernetes.io/hostname
      whenUnsatisfiable: ScheduleAnyway
      labelSelector:
        matchLabels:
          foo: ""
```

# Day 1 – Preemption



KubeCon



CloudNativeCon

North America 2020

*Virtual*

- What if no node can satisfy all the **Hard** Constraints?
- Preemption
  - High-priority Pods are eligible to preempt low-priority Pods

```
File: pod1.yaml
1  apiVersion: v1
2  kind: Pod
3  metadata:
4    name: pod1
5  spec:
6    priorityClassName: p1
7    containers:
8    - name: pod1
9      image: k8s.gcr.io/pause:3.2
```

```
File: priority-classes.yaml
1  apiVersion: scheduling.k8s.io/v1
2  kind: PriorityClass
3  metadata:
4    name: p1
5  value: 1
6  description: "priority with value 1"
7  ---
8  apiVersion: scheduling.k8s.io/v1
9  kind: PriorityClass
10 metadata:
11   name: p2
12   value: 2
13   description: "priority with value 2"
```

# Day 1 – Scheduling Flow



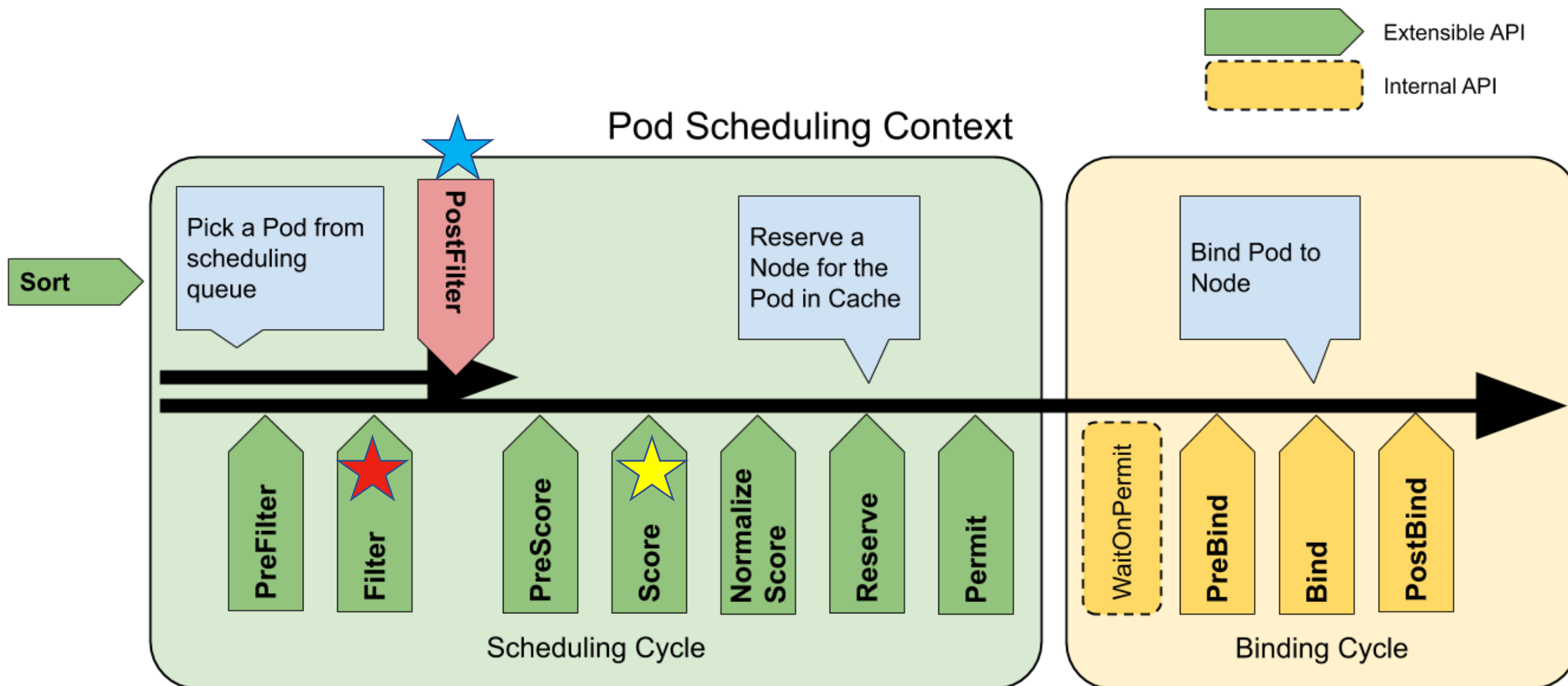
KubeCon



CloudNativeCon

North America 2020

*Virtual*



# Day 2 – Cluster admin / Devops



*Virtual*

- Audience: cluster admin / devops
- Goal: master scheduler configuration best practices, understand scheduler framework basics to make the most of kube-scheduler
- Non-goal: writing additional scheduler extender / plugin code

# Day 2 – Configurations

- [KubeSchedulerConfiguration](#) (`--config <config file>`) over ~~CLI args~~
  - v1alph1 (`<= k8s 1.17`)
  - v1alph2 (`k8s 1.18`)
  - v1beta1 (`k8s 1.19`)
  - Legacy [policy file](#) based config vs. [plugins](#) based config

`<= k8s 1.18`

Predicate/Priority  
Based

VS.

`>= k8s 1.18`

Framework Plugin  
Based

# Day 2 – Policy-based config



Virtual

File: `v1alpha1/scheduler-config.yaml`

```
1  apiVersion: kubescheduler.config.k8s.io/v1alpha1
2  kind: KubeSchedulerConfiguration
3  clientConnection:
4    kubeconfig: "/var/run/kubernetes/scheduler.kubeconfig"
5  algorithmSource:
6    policy:
7      file:
8        path: "/root/config/scheduler-policy.yaml"
```

File: `v1alpha1/scheduler-policy.yaml`

```
1  apiVersion: v1 # or kubescheduler.config.k8s.io/v1
2  kind: Policy
3  predicates:
4    - name: PodFitsHost
5    - name: PodFitsResources
6    ...
7  priorities:
8    - name: InterPodAffinityPriority
9      weight: 50
10   - name: LeastRequestedPriority
11     weight: 10
12   # - name: MostRequestedPriority
13   #   weight: 10
14   ...
```

- <https://kubernetes.io/docs/reference/scheduling/policies/>
- Will be deprecated
- Provide config and policy yaml
- Policy API are not user-friendly – have to specify a full list of predicates/priorities

# Day 2 – Plugins-based config



Virtual

- <https://kubernetes.io/docs/reference/scheduling/config/>
- Aligned with scheduler framework
- User friendly - enable/disable plugins
- Support multi-profile config

File: **v1beta1.yaml**

```
1  apiVersion: kubescheduler.config.k8s.io/v1beta1
2  kind: KubeSchedulerConfiguration
3  clientConnection:
4    kubeconfig: /etc/kubernetes/scheduler.conf
5  profiles:
6    - schedulerName: default-scheduler
7      plugins:
8        score:
9          enabled:
10             - name: NodeResourcesMostAllocated
11               weight: 50000
12             disabled:
13             - name: NodeResourcesLeastAllocated
```

# Day 2 – Multi-profile scheduler



Virtual

```
apiVersion: kubescheduler.config.k8s.io/v1beta1
kind: KubeSchedulerConfiguration
clientConnection:
  kubeconfig: /etc/kubernetes/scheduler.conf
leaderElection:
  leaderElect: false
profiles:
- schedulerName: default-scheduler
- schedulerName: image-first
  plugins:
    score:
      enabled:
      - name: ImageLocality
        weight: 50000
      disabled:
      - name: ImageLocality
- schedulerName: binpack
  plugins:
    score:
      enabled:
      - name: NodeResourcesMostAllocated
        weight: 50000
      disabled:
      - name: NodeResourcesLeastAllocated
- schedulerName: skip-score
  plugins:
    preScore:
      disabled:
      - name: "*"
    score:
      disabled:
      - name: "*"

```

File: **ubuntu.yaml**

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: ubuntu
5  spec:
6    replicas: 1
7    selector:
8      matchLabels:
9        app: ubuntu
10   template:
11     metadata:
12       labels:
13         app: ubuntu
14     spec:
15       schedulerName: image-first
16     containers:
17     - name: ubuntu
18       image: ubuntu
19       # Wait forever
20       command: [ "/bin/bash", "-c", "--" ]
21       args: [ "while true; do sleep 3600; done;" ]

```

File: **normal.yaml**

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: pause
5  spec:
6    replicas: 1
7    selector:
8      matchLabels:
9        app: pause
10   template:
11     metadata:
12       labels:
13         app: pause
14     spec:
15       # schedulerName: default-scheduler
16     containers:
17     - name: pause
18       image: k8s.gcr.io/pause:3.2

```

File: **binpack.yaml**

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: pause
5  spec:
6    replicas: 1
7    selector:
8      matchLabels:
9        app: pause
10   template:
11     metadata:
12       labels:
13         app: pause
14     spec:
15       schedulerName: binpack
16     containers:
17     - name: pause
18       image: k8s.gcr.io/pause:3.2

```

File: **skipscore.yaml**

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: pause
5  spec:
6    replicas: 1
7    selector:
8      matchLabels:
9        app: pause
10   template:
11     metadata:
12       labels:
13         app: pause
14     spec:
15       schedulerName: skip-score
16     containers:
17     - name: pause
18       image: k8s.gcr.io/pause:3.2

```

# Day 2 – Dive a bit deeper



KubeCon

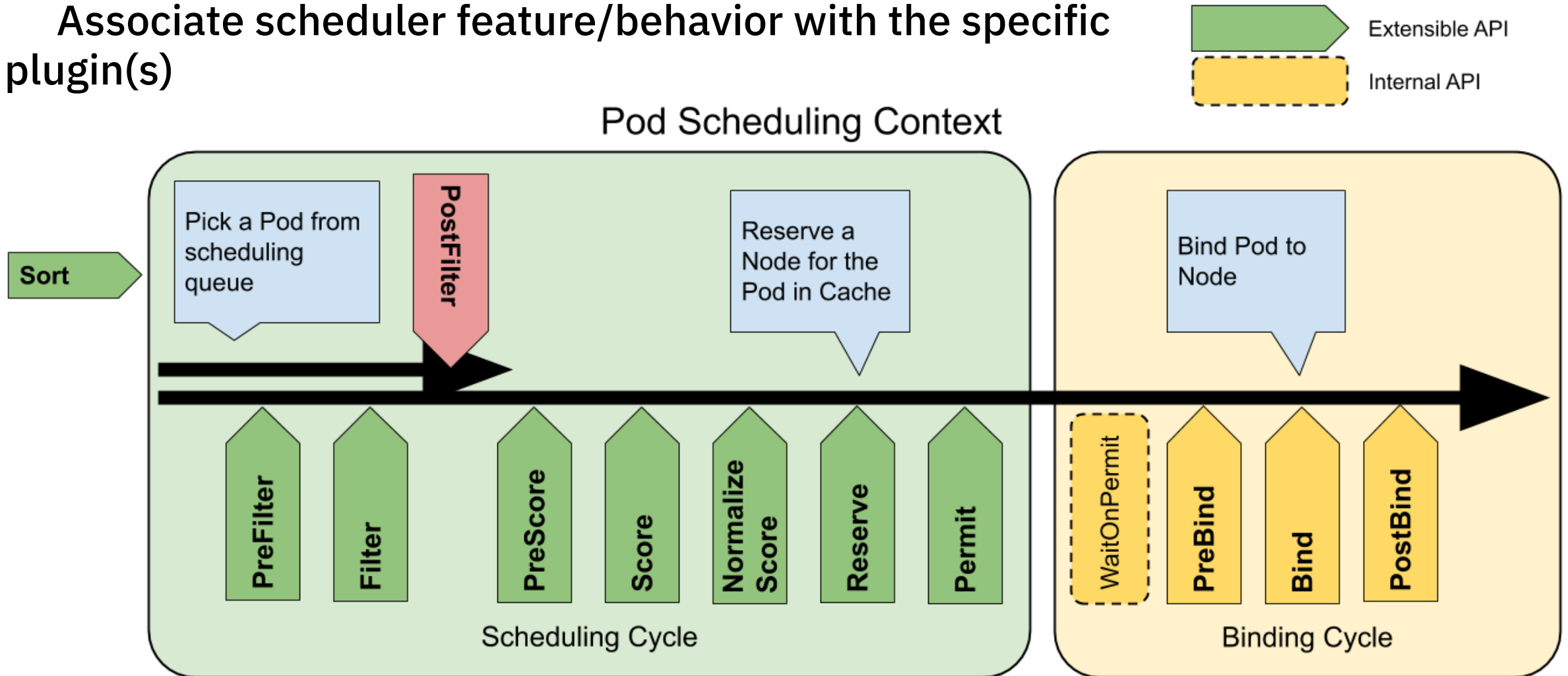


CloudNativeCon

North America 2020

*Virtual*

Associate scheduler feature/behavior with the specific plugin(s)



# Day 2 – Enabled plugin list



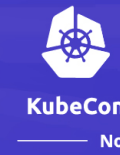
North America 2020

## Scheduling plugins

The following plugins, enabled by default, implement one or more of these extension points:

- **SelectorSpread** : Favors spreading across nodes for Pods that belong to Services, ReplicaSets and StatefulSets. Extension points: **PreScore** , **Score** .
- **ImageLocality** : Favors nodes that already have the container images that the Pod runs. Extension points: **Score** .
- **TaintToleration** : Implements [taints and tolerations](#). Implements extension points: **Filter** , **Prescore** , **Score** .
- **NodeName** : Checks if a Pod spec node name matches the current node. Extension points: **Filter** .
- **NodePorts** : Checks if a node has free ports for the requested Pod ports. Extension points: **PreFilter** , **Filter** .
- **NodePreferAvoidPods** : Scores nodes according to the node annotation `scheduler.alpha.kubernetes.io/preferAvoidPods` . Extension points: **Score** .
- **NodeAffinity** : Implements [node selectors](#) and [node affinity](#). Extension points: **Filter** , **Score** .
- **PodTopologySpread** : Implements [Pod topology spread](#). Extension points: **PreFilter** , **Filter** , **PreScore** , **Score** .
- **NodeUnschedulable** : Filters out nodes that have `.spec.unschedulable` set to true. Extension points: **Filter** .
- **NodeResourcesFit** : Checks if the node has all the resources that the Pod is requesting. Extension points: **PreFilter** , **Filter** .
- **NodeResourcesBalancedAllocation** : Favors nodes that would obtain a more balanced resource usage if the Pod is scheduled there. Extension points: **Score** .
- **NodeResourcesLeastAllocated** : Favors nodes that have a low allocation of resources. Extension points: **Score** .
- **VolumeBinding** : Checks if the node has or if it can bind the requested volumes. Extension points: **PreFilter** , **Filter** , **Reserve** , **PreBind** .
- **VolumeRestrictions** : Checks that volumes mounted in the node satisfy restrictions that are specific to the volume provider. Extension points: **Filter** .
- **VolumeZone** : Checks that volumes requested satisfy any zone requirements they might have. Extension points: **Filter** .
- **NodeVolumeLimits** : Checks that CSI volume limits can be satisfied for the node. Extension points: **Filter** .
- **EBSLimits** : Checks that AWS EBS volume limits can be satisfied for the node. Extension points: **Filter** .
- **GCEPDLimits** : Checks that GCP-PD volume limits can be satisfied for the node. Extension points: **Filter** .
- **AzureDiskLimits** : Checks that Azure disk volume limits can be satisfied for the node. Extension points: **Filter** .
- **InterPodAffinity** : Implements [inter-Pod affinity and anti-affinity](#). Extension points: **PreFilter** , **Filter** , **PreScore** , **Score** .
- **PrioritySort** : Provides the default priority based sorting. Extension points: **QueueSort** .
- **DefaultBinder** : Provides the default binding mechanism. Extension points: **Bind** .
- **DefaultPreemption** : Provides the default preemption mechanism. Extension points: **PostFilter** .

# Day 2 – Disabled plugin list



North America 2020

*Virtual*

- **NodeResourcesMostAllocated** : Favors nodes that have a high allocation of resources. Extension points: **Score** .
- **RequestedToCapacityRatio** : Favor nodes according to a configured function of the allocated resources. Extension points: **Score** .
- **NodeResourceLimits** : Favors nodes that satisfy the Pod resource limits. Extension points: **PreScore** , **Score** .
- **CinderVolume** : Checks that OpenStack Cinder volume limits can be satisfied for the node. Extension points: **Filter** .
- **NodeLabel** : Filters and / or scores a node according to configured label(s). Extension points: **Filter** , **Score** .
- **ServiceAffinity** : Checks that Pods that belong to a Service fit in a set of nodes defined by configured labels. This plugin also favors spreading the Pods belonging to a Service across nodes. Extension points: **PreFilter** , **Filter** , **Score** .

# Day 2 – Wrap-up

- Understand plugins and plugin arguments
- Some global settings:
  - `percentageOfNodesToScore`
- What's new
  - `DefaultTopologySpread`: beta in 1.19
  - Prioritizing nodes based on volume capacity [[KEP 1845](#)]
  - Try out `.status.nominatedNodeName` as a shortcut [[KEP 1923](#)]
  - [Draft] Simplified version of topology manager in kube-scheduler [[#1858](#)]
  - [Draft] Add default node affinity constraints to NodeAffinity plugin [[#95738](#)]
- 🎉 Eventual goal: offer an out-of-box "multi-flavored" scheduler

# Day 3 – Enthusiast / Innovator



*Virtual*

- **Audience:** scheduling enthusiast / innovator
- Goal: extend scheduler to fit diverse workloads, by writing as minimum code as possible
- Non-Goal: start from scratch to write a secondary scheduler

# Day 3 – Deep dive into scheduler framework



KubeCon



CloudNativeCon

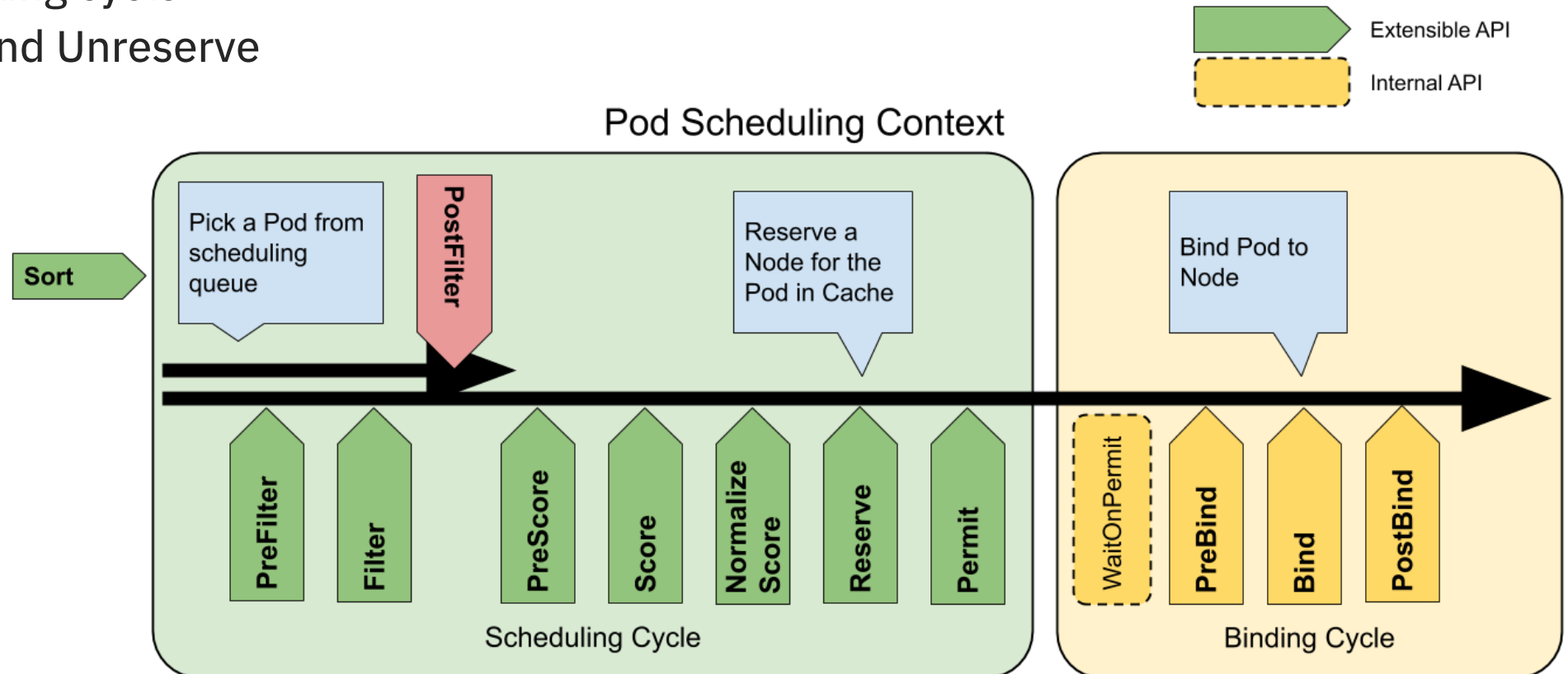
North America 2020

*Virtual*

## What's new

- Refined PostFilter in 1.19
- Permit in scheduling cycle
- Merge Reserve And Unreserve
- 🎉 GA in 1.20 🎉

## Master details of each extension point



# Day 3 – Build your own scheduler



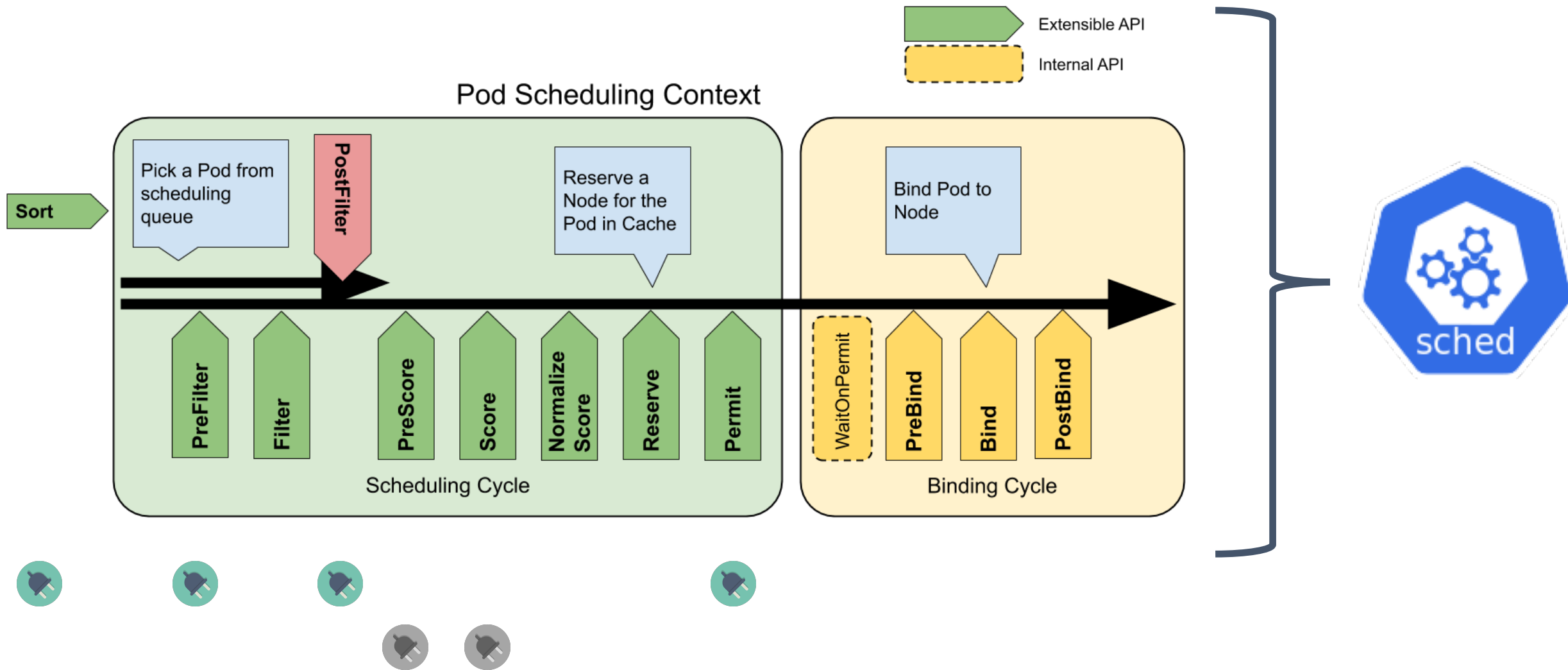
KubeCon



CloudNativeCon

North America 2020

*Virtual*



# Day 3 – scheduler-plugins



Virtual

- [scheduler-plugins](#): a sig-scheduling sponsored project

2-lightweight-coscheduling

42-podgroup-coscheduling

48-node-resources-allocation

9-capacity-scheduling

apis

controller

coscheduling

crossnodepreemption

generated

noderesources

qos

util

## Real Load Aware Scheduling KEP #61

Open zorro786 wants to merge 8 commits into `kubernetes-sigs:master` from `zorro786:real-load-kep`

Conversation 118 Commits 8 Checks 0 Files changed 14

zorro786 comment

master scheduler-plugins / cloudbuild.yaml

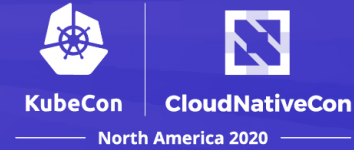
seanmalloy Fix Automated Container Image Builds ... ✓

1 contributor

24 lines (23 sloc) 904 Bytes

```
1 # See https://cloud.google.com/cloud-build/docs/build-config
2
3 # this must be specified in seconds. If omitted, defaults to 600s (10 mins)
4 timeout: 1200s
5 # this prevents errors if you don't use both _GIT_TAG and _PULL_BASE_REF,
6 # or any new substitutions added in the future.
7 options:
8   substitution_option: ALLOW_LOOSE
9 steps:
```

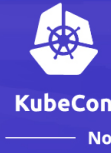
# Contact Us



*Virtual*

- SIG-Chairs
  - @ahg-g, Google
  - @Huang-Wei, IBM
- [Home page](#)
- Slack channel: [#sig-scheduling](#)
- [Mailing list](#)
- [Weekly meeting](#)

# Q & A



*Virtual*

North America 2020

