# AGENDA
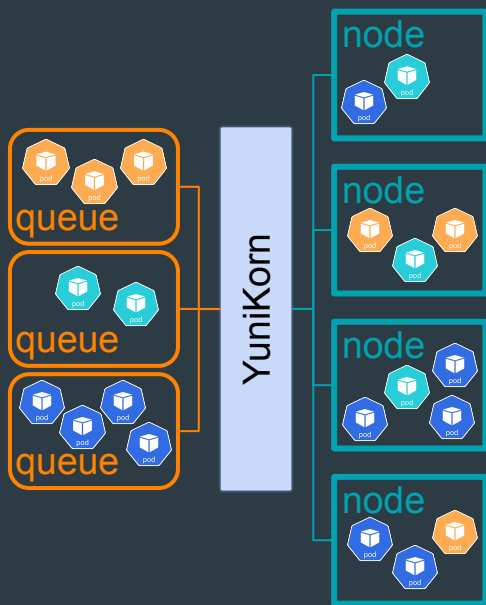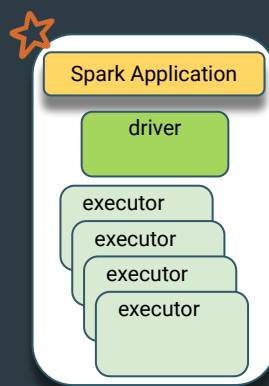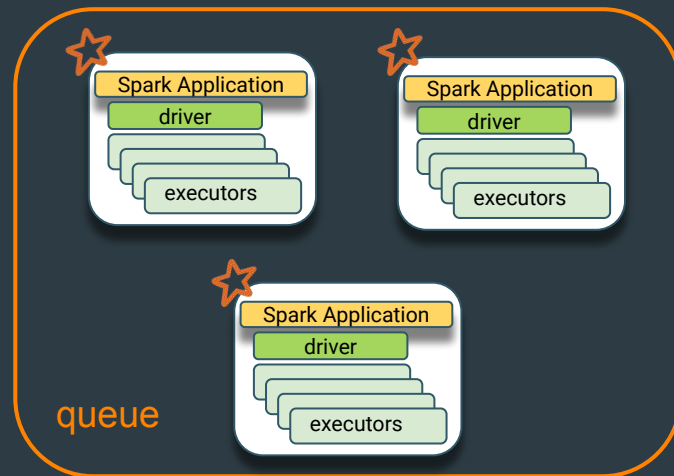
# WHY APACHE YUNIKORN

## Advanced scheduling requirements



Workload Queueing

Gang Scheduling

Application Sorting

# AGENDA

Why Apache YuniKorn

**Architecture**
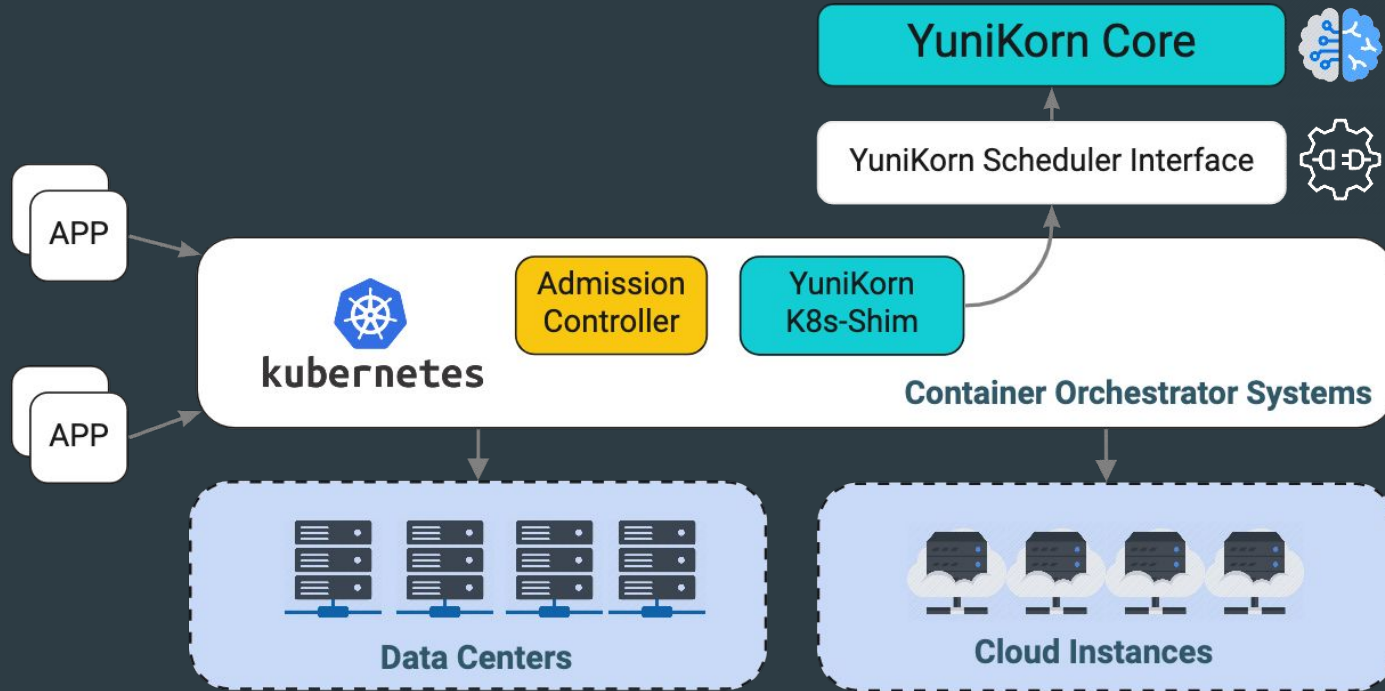
Queues: hierarchies and quotas

Gang Scheduling

User and Group Quotas

Roadmap

YUNIKORN

# ARCHITECTURE
## Basics

# ARCHITECTURE
Deployment models
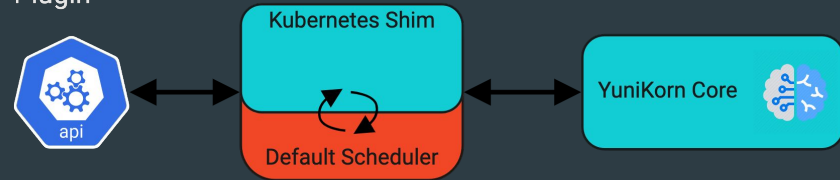
- STANDARD:
  - Custom scheduler
  - Replaces *default* scheduler

- PLUGIN:
  - Kubernetes Scheduling Framework (API)
  - Replace or augment limited functionality

Standard



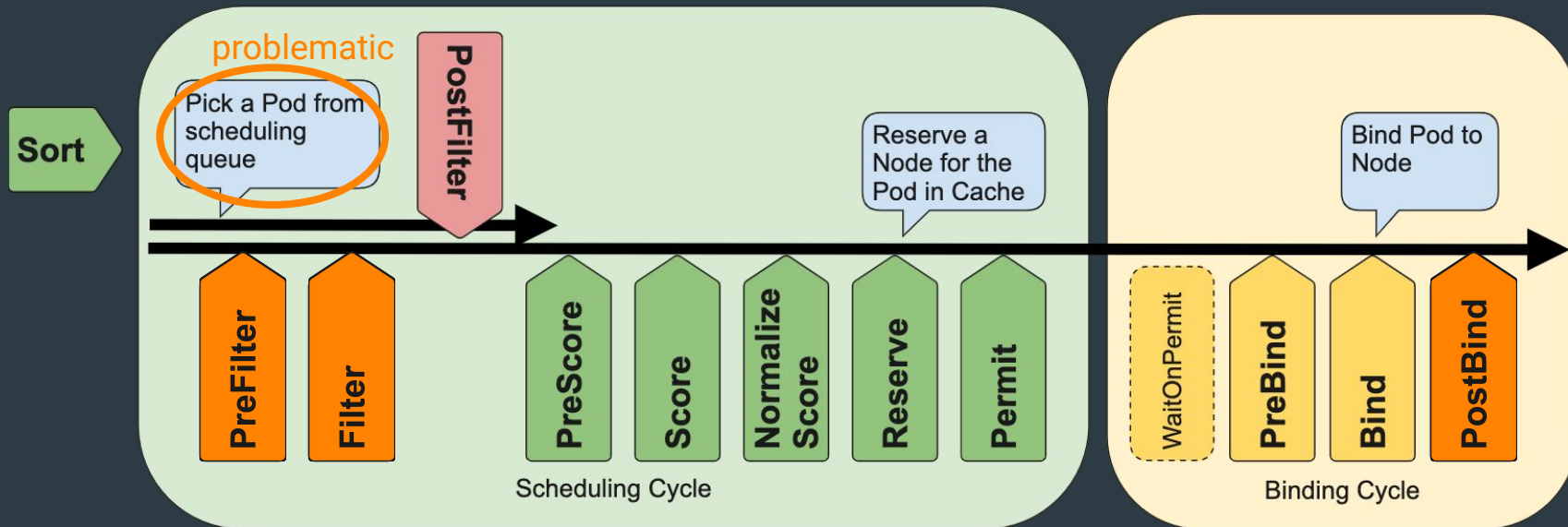Plugin

# KUBERNETES COMMUNITY

- **Multiple** initiatives around batch and HPC schedulers
  - Apache YuniKorn
  - Volcano
  - Kueue
  - Armada
- Scheduler plugins for specific point solutions
- Fragmentation concerns: KubeCON NA panel discussion

- Pre-enqueue plugin
  - KEP-3521: Pod Schedule Readiness

- Pre-emption behaviour: priority only
  - K8s has: do not preempt other pods during scheduling
  - Batch needs: do not preempt the tagged pod while running

CLOUDERA

# AGENDA

Why Apache YuniKorn
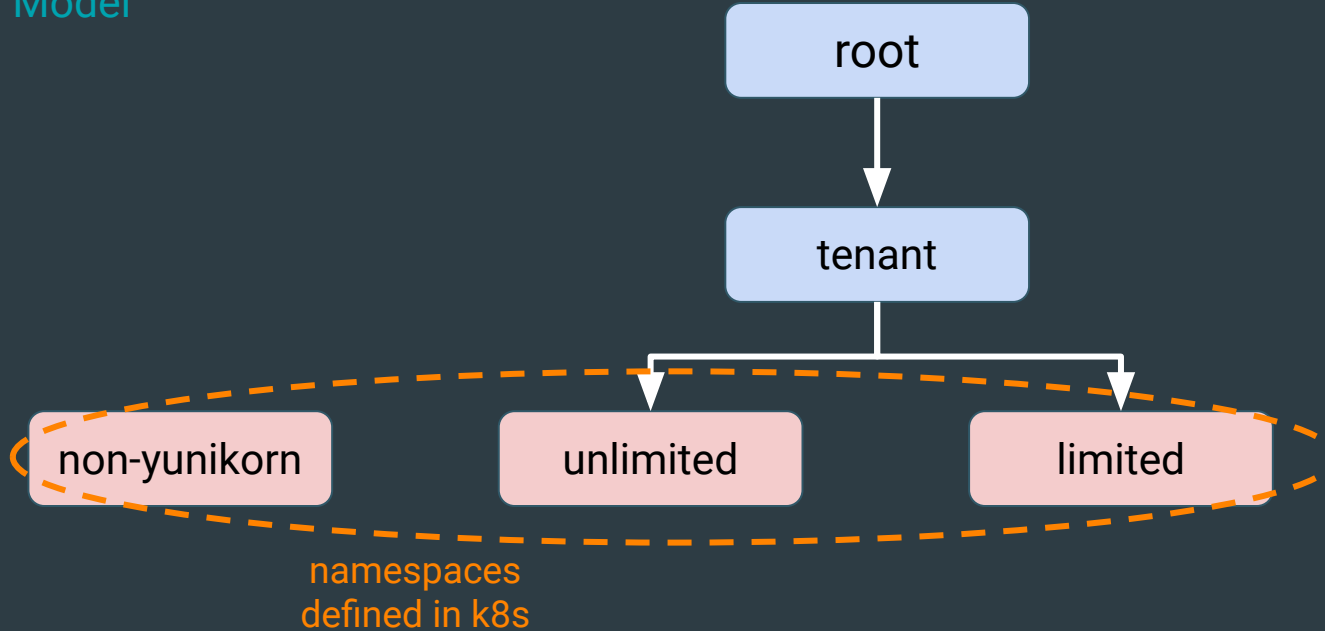
Architecture

**Queues: hierarchies and quotas**

Gang Scheduling

User and Group Quotas

Roadmap

# QUEUES, HIERARCHIES & QUOTAS
## Hierarchical Model

# QUEUES, HIERARCHIES & QUOTAS
## Hierarchical Model



resource limits based on registered nodes

resource limit 75 GB / 75 CPU

K8s namespace quota

resource limit 50 GB / 50 CPU

root

tenant

non-yunikorn

unlimited

limited

pods

pods

pods

CLOUDERA

# QUEUES, HIERARCHIES & QUOTAS
## Hierarchical Model

# AGENDA

Why Apache YuniKorn

Architecture

Queues: hierarchies and quotas

**Gang Scheduling**

User and Group Quotas

Roadmap

# GANG SCHEDULING
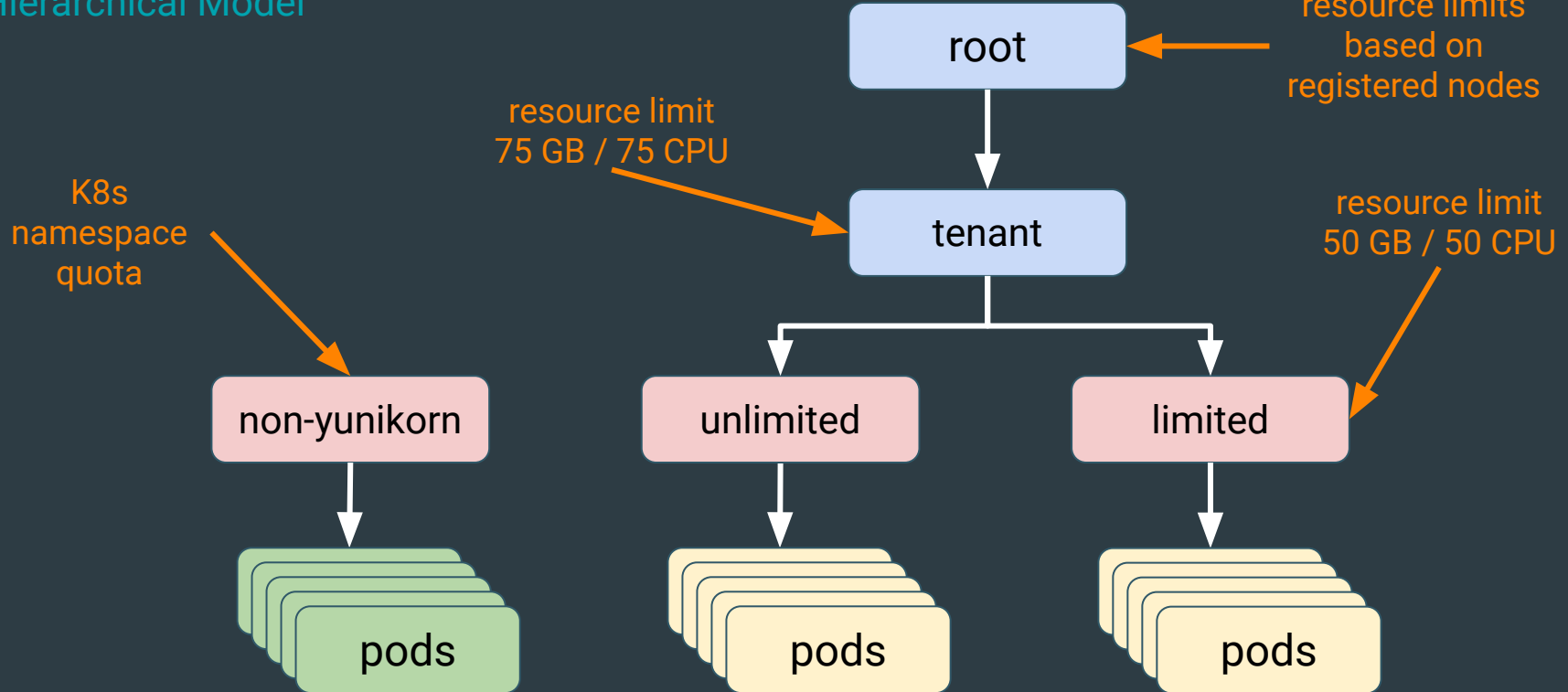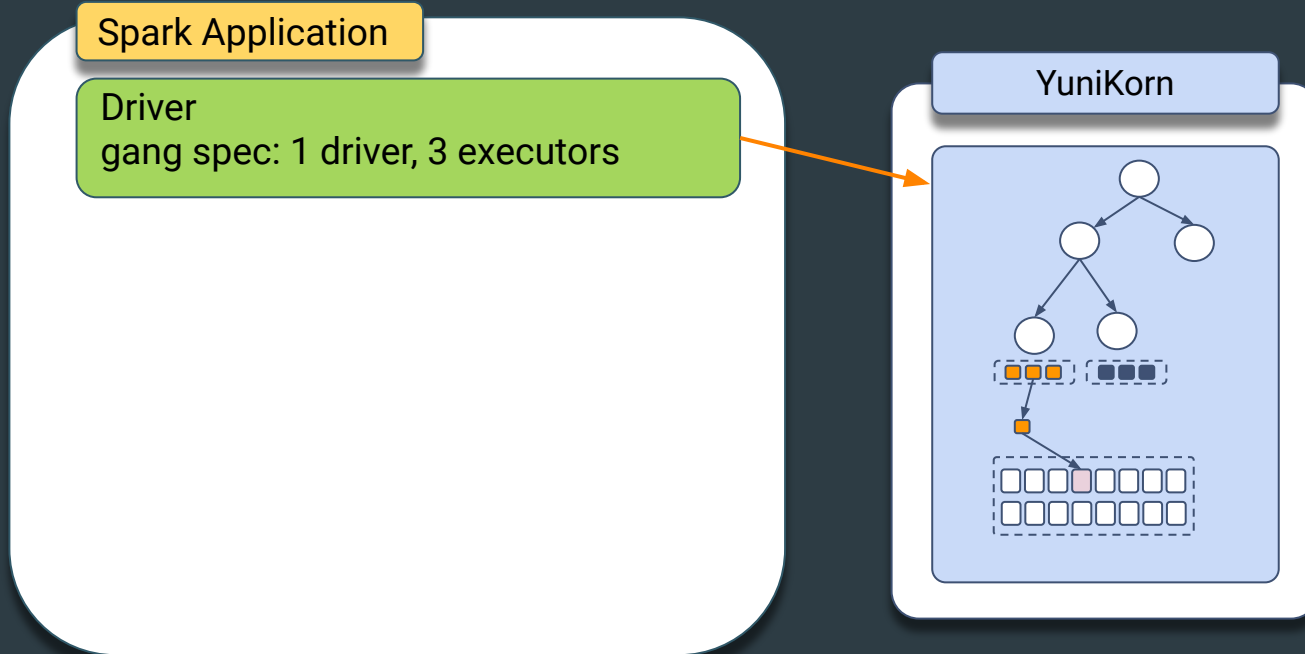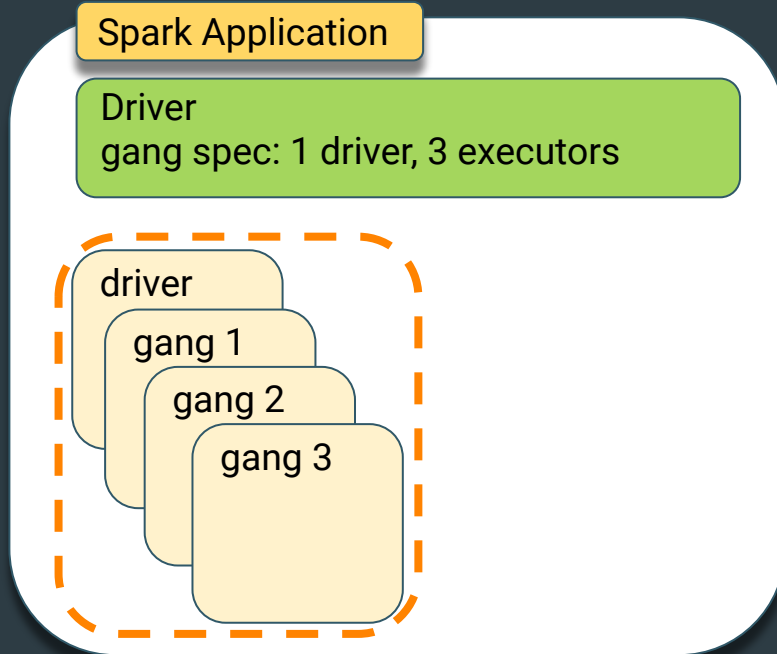
## Gang example: Spark

# GANG SCHEDULING

## Gang example: Spark



Spark Application

Driver
gang spec: 1 driver, 3 executors

gang 1
gang 2
gang 3

placeholders

YuniKorn

# GANG SCHEDULING

## Gang example: Spark



placeholders

# AGENDA

Why Apache YuniKorn

Architecture

Queues: hierarchies and quotas

Gang Scheduling

**User and Group Quotas**

Roadmap

CLOUDERA

# USER AND GROUP INFORMATION
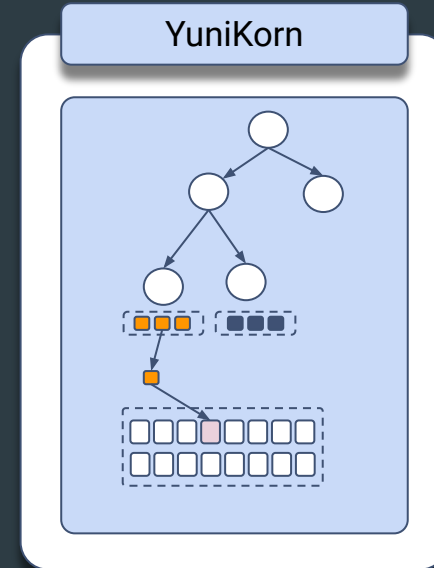
## Two options

Authenticated user:
- User information based on authenticated user
- Set by admission controller on pod creation

Annotation on pod:
- User information provided on pod creation
- Check & override by admission controller

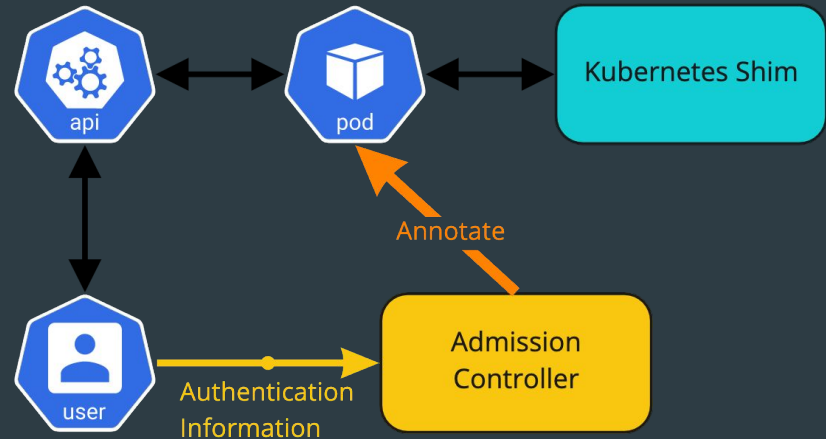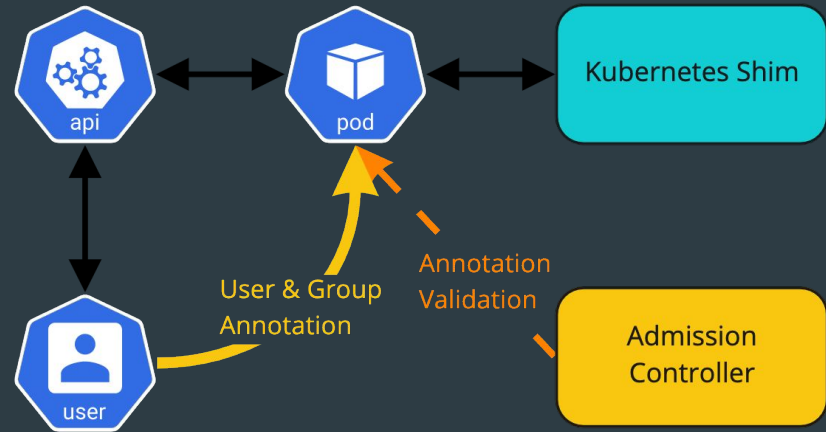# USER AND GROUP INFORMATION

## Two options

Authenticated user:

- User information based on authenticated user
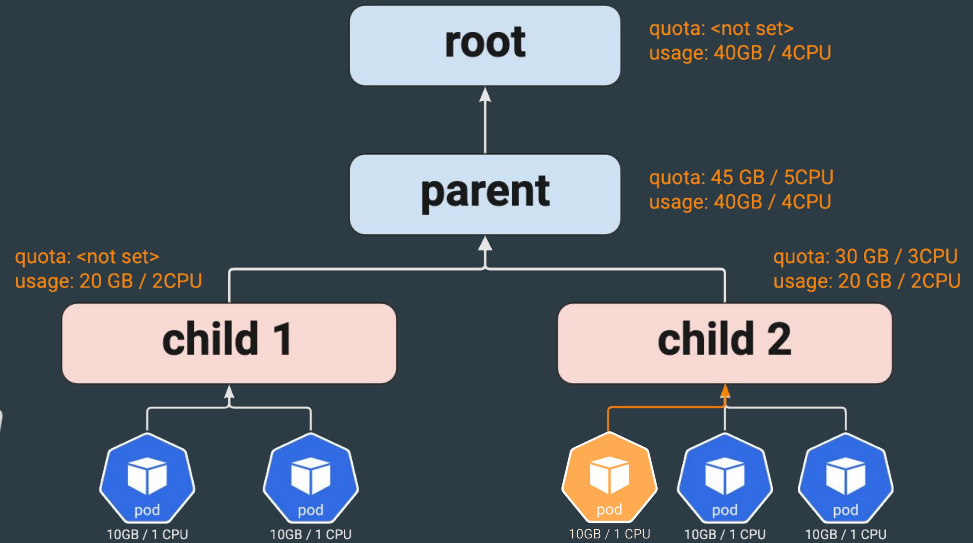- Set by admission controller on pod creation

Annotation on pod:

- User information provided on pod creation
- Check & override by admission controller

CLOUDERA

# USER AND GROUP QUOTAS

## Queues vs User & Groups

- Queue quota: root to leaf
  - Integrated in scheduling cycle
  - No quota available queue will be skipped

- User and group: leaf to root
  - Request fits in queue *headroom*
  - Last check before confirming node placement



**root**
quota: <not set>
usage: 40GB / 4CPU

**parent**
quota: 45 GB / 5CPU
usage: 40GB / 4CPU

quota: <not set>
usage: 20 GB / 2CPU

quota: 30 GB / 3CPU
usage: 20 GB / 2CPU

**child 1**

**child 2**

pod
10GB / 1 CPU

pod
10GB / 1 CPU

pod
10GB / 1 CPU

pod
10GB / 1 CPU

pod
10GB / 1 CPU

# AGENDA

Why Apache YuniKorn

Architecture

Queues: hierarchies and quotas

Gang Scheduling

User and Group Quotas

**Roadmap**

# Roadmap
With all appropriate caveats…

## Apache YuniKorn 1.2

- Application limits
  - Maximum running applications per Queue
- User and Group quotas
  - Implementation started for:
    - User retrieval
    - Usage tracking
  - Enforcement design in final stage

## Future release

- Priority support
  - Early design stage
- Pre-emption
  - Obsolete code removed (YuniKorn 1.1)
  - Design starting

# CLOUDERA
# Q & A