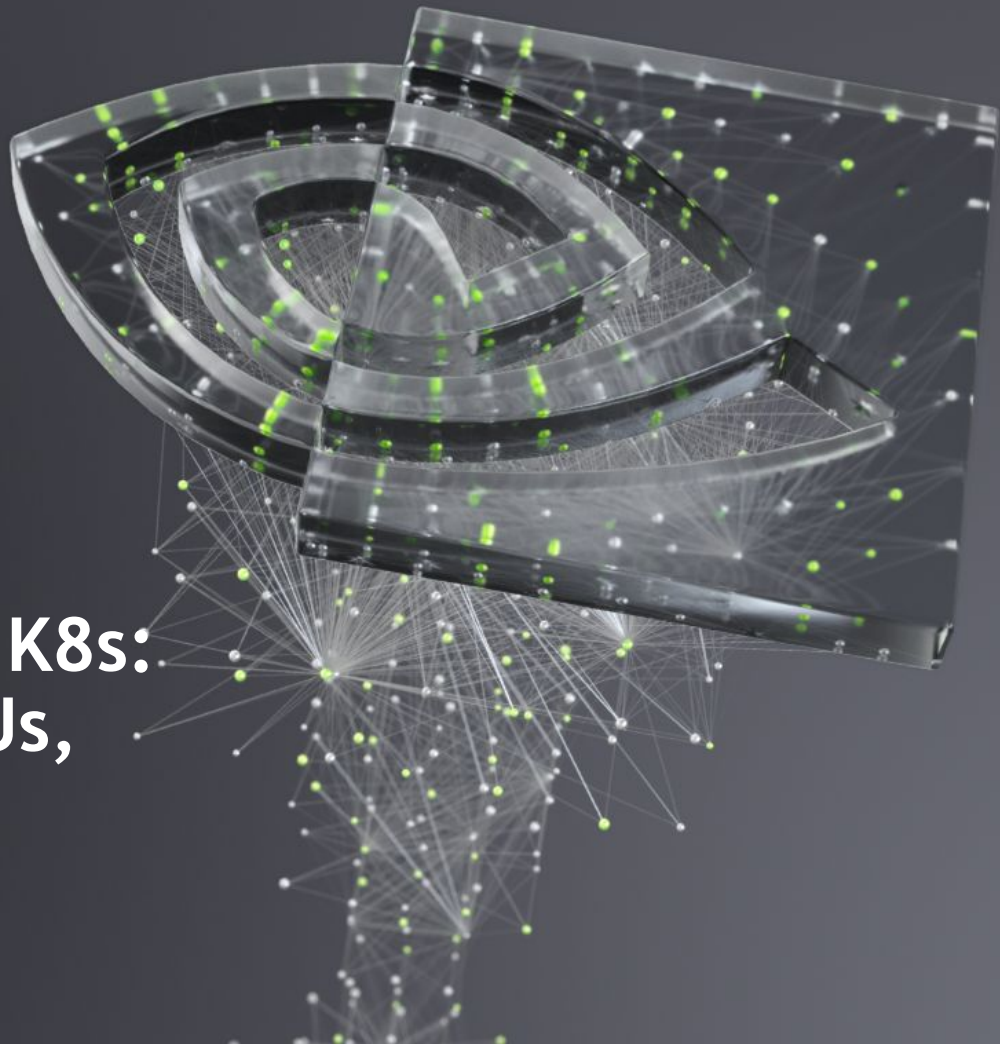# Multi-node Jobs with K8s: Gang Scheduling, GPUs, MPI and RDMA

Madhukar Korupolu, Sanjay Chatterjee

# Deep Learning Applications

# AI / DL: Models, Frameworks, Hardware

# Trends: Big Data, Larger Models

- Data and model sizes increasing

- BERT NLP: 110M, 330M params
  - Recent: 8B, 17B

- Strong demand for multi-gpu jobs
  - Larger problems
  - Faster turn-around
  - E.g., 128-GPUs per job

- ML Perf training results

### Training Time vs. Num GPUs



Chart: Minutes vs. Num GPUs (V100)

| Num GPUs (V100) | Minutes |
|---|---|
| 16 | 3504 |
| 64 | 921 |
| 256 | 236 |
| 1024 | 67 |
| 1472 | 47 |

# Sample Multi-GPU Node: DGX-1



## Single DGX-1 Node:

8 Nvidia V100 GPUs

Dual socket, NVLink in node

4 Mlnx EDR NICs (100 Gbps)

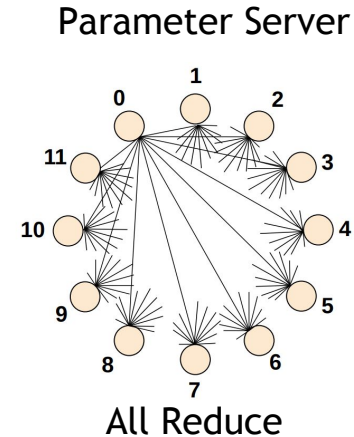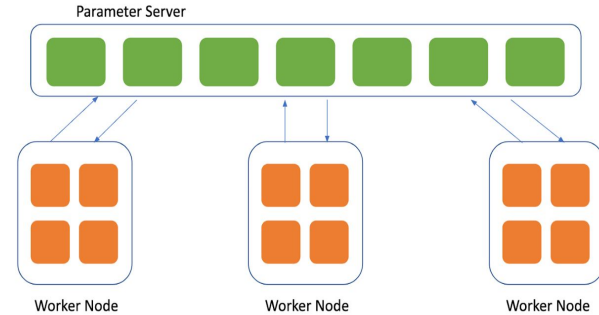Dual 10GbE ports

## Connecting Multiple nodes:

Infiniband or RoCE

4 x NICs to connected fabric

# Distributed Training Applications

## Multi-GPU, Multi-node

- Data / model parallelism

- Stochastic gradient descent (SGD)

- Async SGD: Parameter-server

- Sync SGD: All-reduce
  - NCCL / MPI
  - Utilize fast interconnects / RDMA
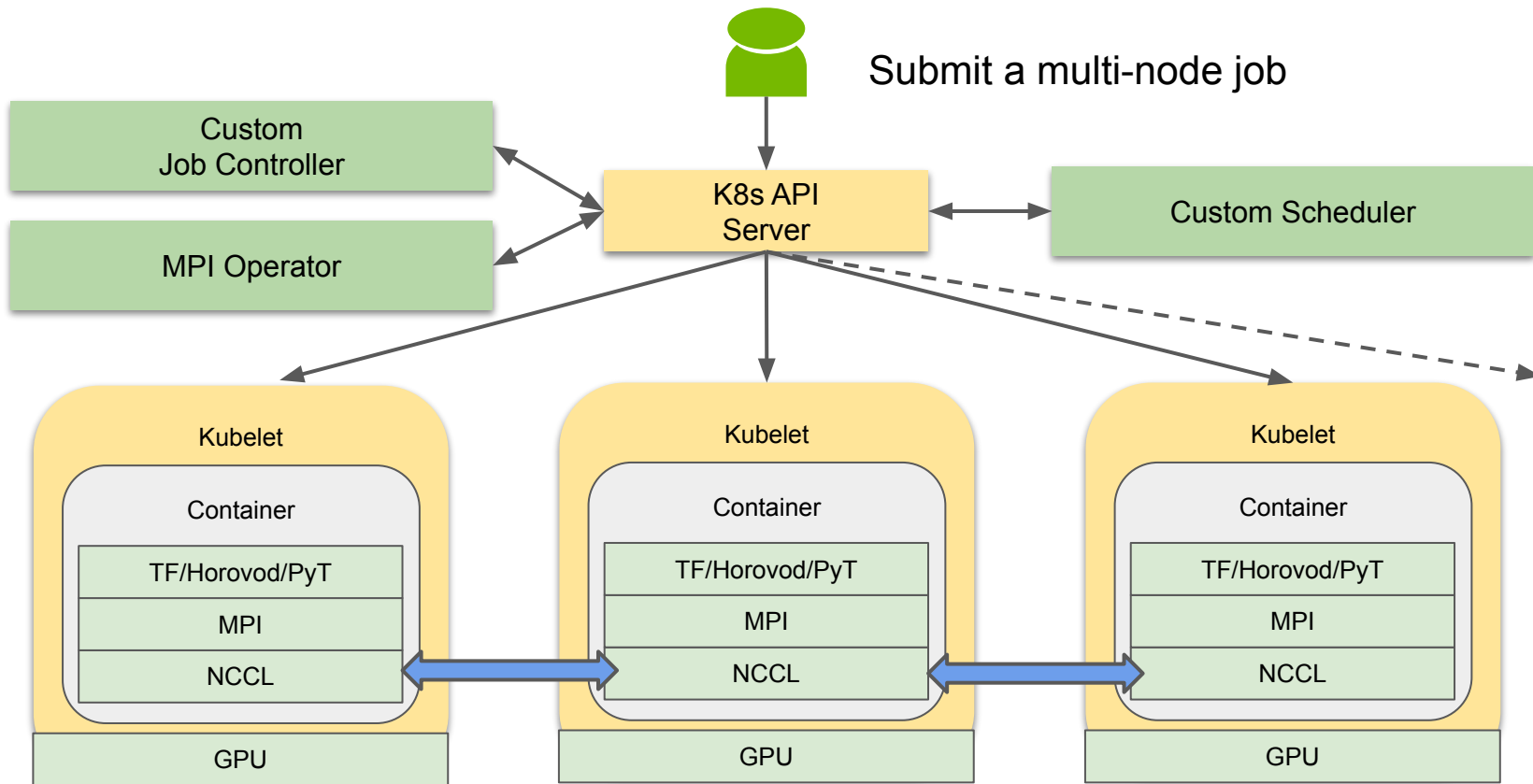  - Horovod library

- Distributed TensorFlow / PyTorch

Parameter Server

Parameter Server

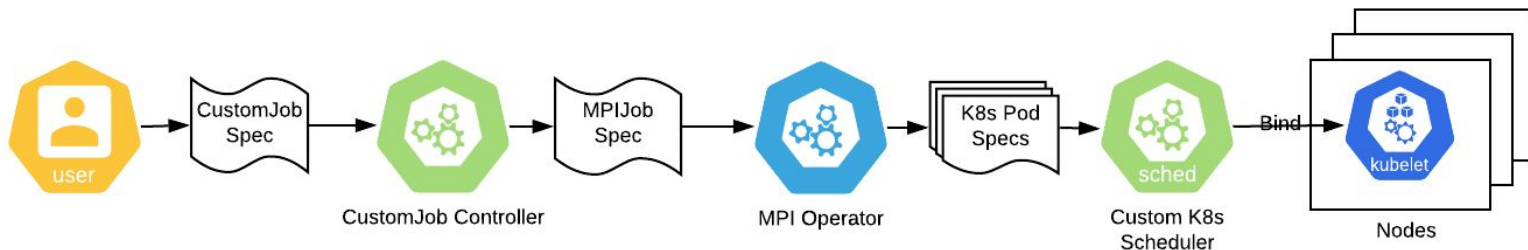All Reduce

# K8s Challenges & Outline

## Multi-node and K8s Gaining Traction

- Motivation & background

- End-to-end flow: Array jobs
  - MPI & job lifecycle

- Gang scheduling

- Multi-rail RDMA / CNI

- Application / BERT

- Production shared clusters

- Quotas, queues, time limits

- Backfilling, utilization

- Monitoring / Operations
  - Dashboards / CICD

- Conclusions / Future work

nVIDIA.

# K8s Orchestration Flow

# Sample PyTorch Job Launch



**Sample Dist PyTorch job launch**

python -m torch.distributed.launch

--nproc_per_node=8 --nnodes=2 --node_rank=0

--master_addr=localhost bert_train.py <args>

python -m torch.distributed.launch

--nproc_per_node=8 --nnodes=2 --node_rank=1

--master_addr=<ip> bert_train.py <args>

**K8s Multi-node PyTorch job launch**

mpirun -np $ARRAY_SIZE -npernode 1

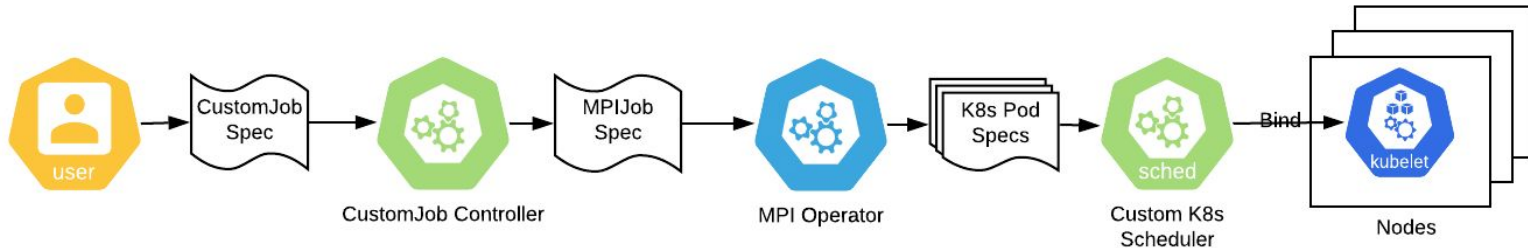python -m torch.distributed.launch
--nproc_per_node=8 --nnodes=$ARRAY_SIZE
--node_rank=$ARRAY_INDEX

--master_addr=$MASTER_IP bert_train.py <args>

nvrun ..

*mpirun as launcher w/ NCCL backend*

NVIDIA

# Array Jobs and MPI Operator



## Array jobs

Abstraction for multi-node

Configurable type, size

Status msgs, Telemetry

```
...
Spec:
  Containers:
    Args:
      –c
      mpirun –np ${ARRAY_SIZE} –npernode 1
        python3 –m torch.distributed.launch
        ––nproc_per_node=8
        ––nnodes=${ARRAY_SIZE}
        ––node_rank=${ARRAY_INDEX}
        ––master_addr=${MASTER_IP}
        bert_train.py
    Command:
      /bin/sh
  ...
  Resources:
    Requests:
      Sriov Rdma:     4
      Nvidia.Com/Gpu: 8
  ...
```
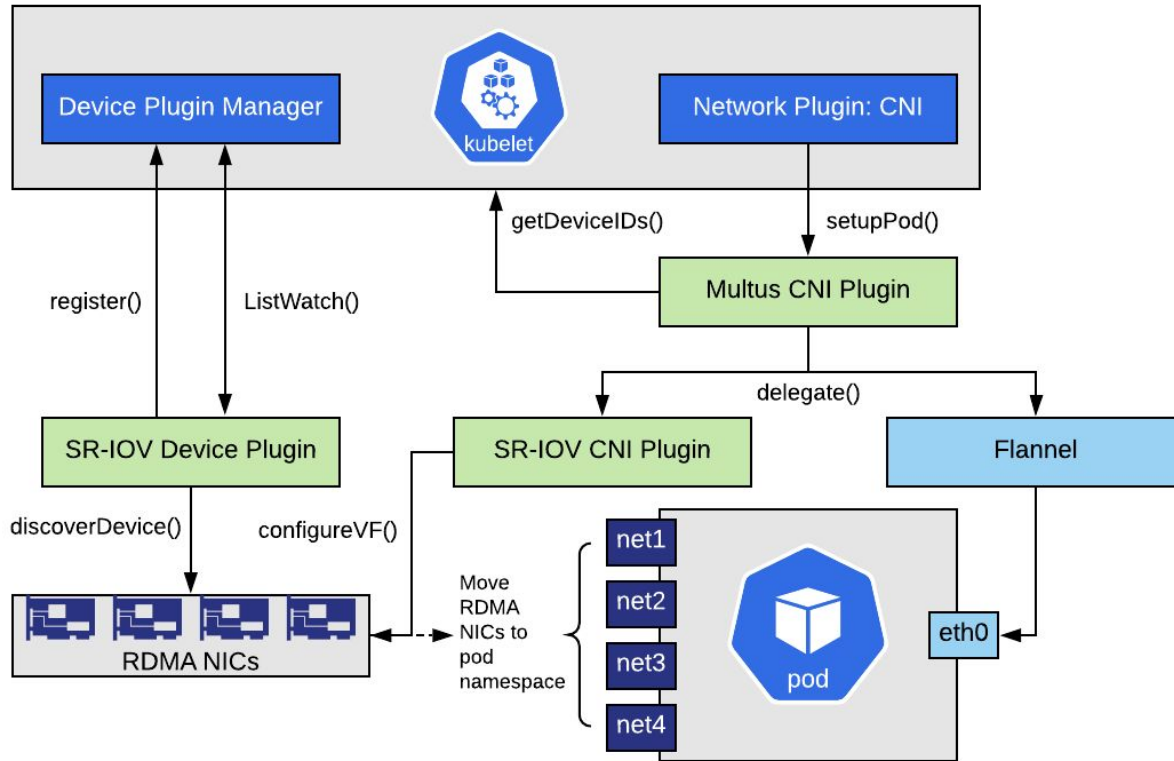
## MPI Operator

Upstream in Kubeflow

Launch replicas on each node
- Kubectl exec, Lifecycle

Mods for gang scheduling

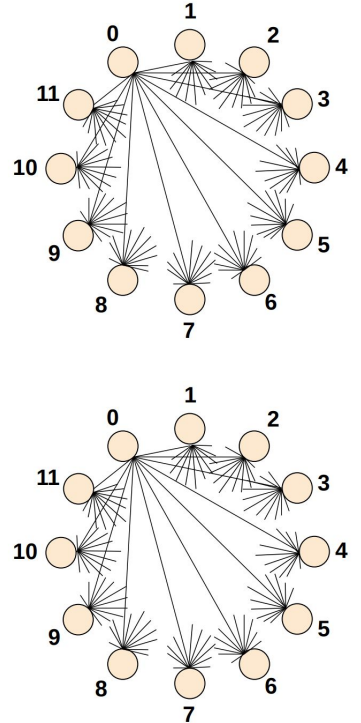# SRIOV CNI for K8s Multi-Rail



Exposing multiple NIC interfaces to K8s Pod

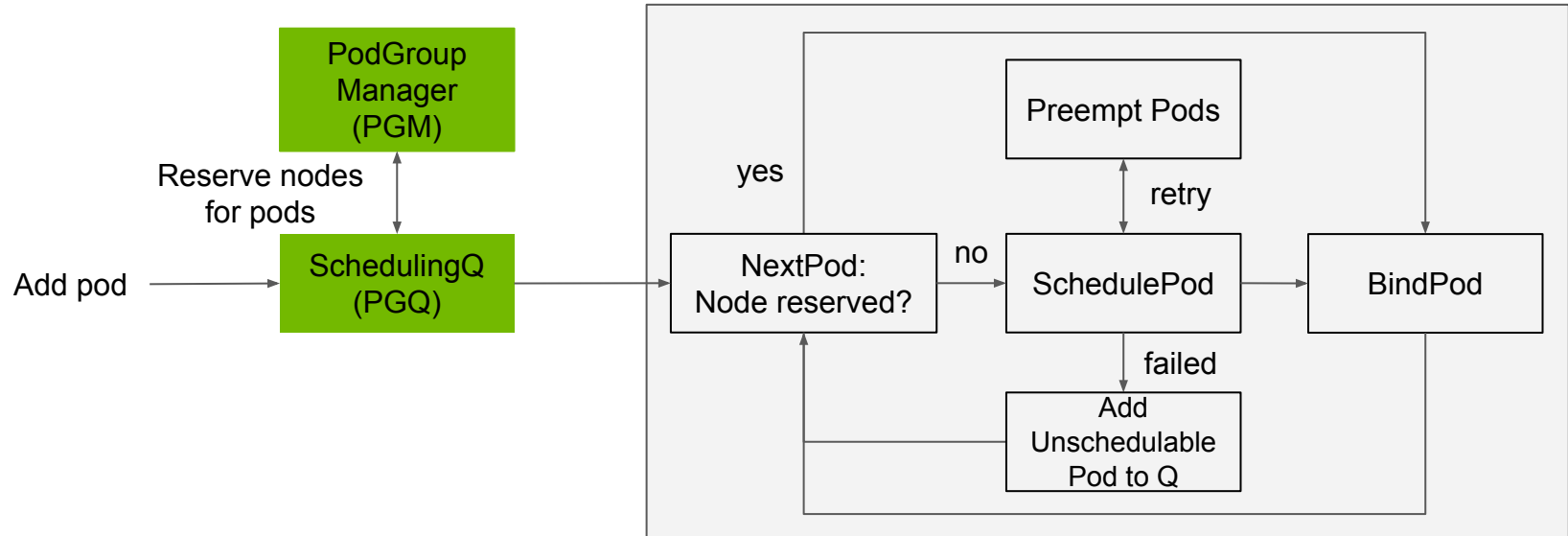Multus delegates to SR-IOV CNI and Flannel

Base SR-IOV CNI from upstream w/ customizations

# Gang Scheduling Multi-Node Pods

- **Multi-node pods:** All-or-none to make progress
  - Default K8s: pods one-by-one ⇒ deadlocks

- **Gang / co-scheduling in K8s:**
  - Open item for default K8s scheduler (since 2015)
  - Basic: loop over pods → wait → timeout → release
  - Being considered via Volcano, Poseidon, etc.

- **Approach:** PodGroup structure
  - Full node pods only, reservation based

# PodGroup Queue and Manager



- Experimental extensions to K8s
- Reserve nodes for full-node pods

# Demo

# Sample Job Real-Time Telemetry

# Sample BERT K8s Scaling



Pre-Training Throughput Scaling (seq/s)

BERT Phase 1: batch_size_per_gpu = 64, seq length = 128
Phase 2: batch_size_per_gpu = 16, seq length = 512

# K8s Challenges & Outline

- Motivation & background

- End-to-end flow: Array jobs
  - MPI & job lifecycle

- Gang scheduling

- Multi-rail RDMA / CNI

- Application / BERT

- Production shared clusters

- Quotas, queues, time limits

- Backfilling, utilization

- Monitoring / Operations
  - Dashboards / CICD

- Conclusions / Future work

# Shared K8s Cluster for Multi-node
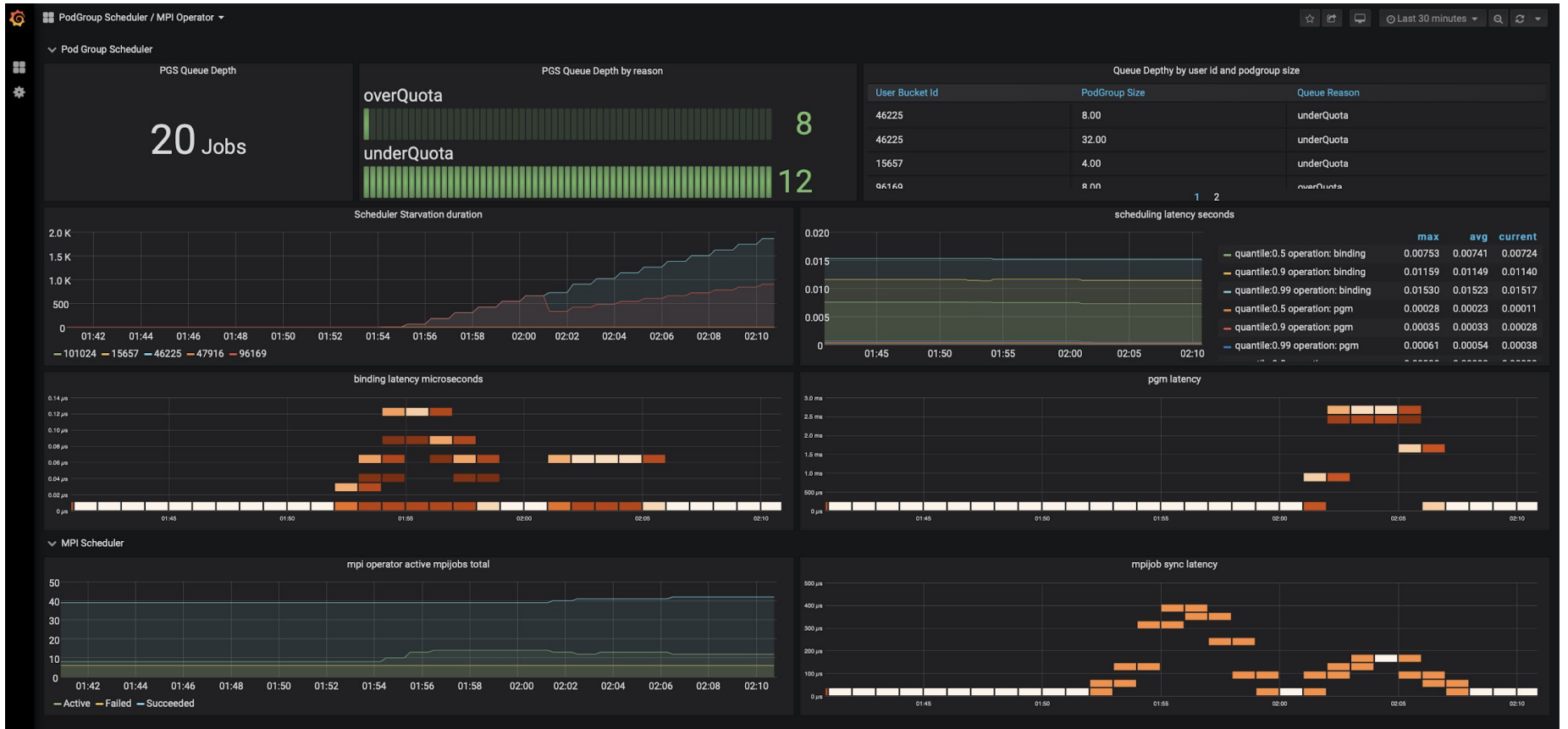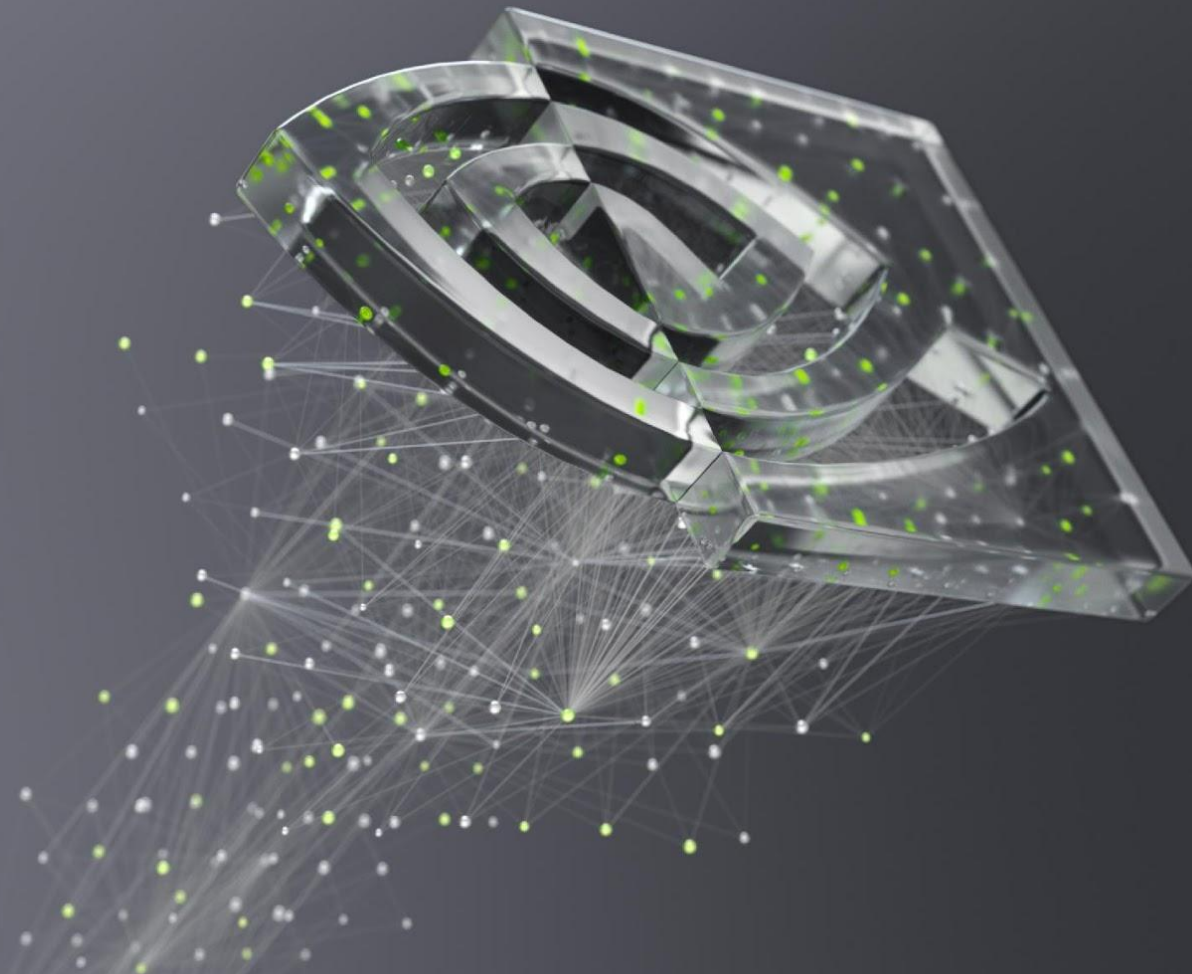
- Production on-prem cluster
  - Early internal users
  - 100 DGX nodes, single hop

- Quotas - Concurrent GPU usage
  - Configurable per user, default

- Time limits
  - E.g. 128 node-hours
  - ⇒ 8 hours for 16-node job

- Starvation handling, backfilling
  - Blocking for large mn jobs
  - Backfilling for utilization

- Dynamic job priority
  - DRF fairness
  - Starvation, age etc
  - Weighted function

- Operations / dashboards

# Scheduler Dashboard

# Summary and Future Work

- Multi-node clusters with K8s
  - Gang scheduling, GPUs, MPI, RDMA

- MN enabled containers / models
  - Available from [ngc.nvidia.com](http://ngc.nvidia.com)

- Ongoing work
  - Production hardening
  - Performance, Storage caching
  - Other array types, K8s framework

- Acknowledgements