



**KubeCon**



**CloudNativeCon**

**North America 2024**





**KubeCon**



**CloudNativeCon**

North America 2024

# Divide and Conquer:

## Master GPU Partitioning and Visualize Savings with OpenCost

**Kaysie Yu**

Azure Kubernetes Service  
Microsoft

**Ally Ford**

Azure Kubernetes Service  
Microsoft

# Agenda

- Why run AI workloads on AKS and GPUs?
- GPUs and cost.... what's the issue?
- Monitor usage metrics and visualize costs
- Partitioning strategies
- Demo!

# GPUs are scarce and expensive!

**Utilization** = amount of time during which one or more kernels was executing on the GPU

Exclusive GPU access is the default

You pay for the entire GPU despite the size of your workload

**Is it feasible and economical to run AI workloads on GPUs at scale?**



KubeCon



CloudNativeCon

North America 2024

# 48% of orgs are using K8s for AI/ML workloads

*State of Workloads Deployment on Containers and Kubernetes. Gartner peer Survey.*

# K8s is *the* ideal platform for GenAI, LLMs

- **Scalability**

- Seamless and automatic for dynamic and resource intensive workloads like inferencing

- **Scheduling**

- Reduce the operational overhead burden on MLOps teams

- **Dynamic**

- Dynamically allocate resources

- **Flexibility**

- Extensible and portable across clouds from public to onprem or hybrid

## CPUs

- **Ubiquitous** – not a single digital device that won't have one
- **Fast and versatile** – few operations executed quickly
- **Serialization** – can very efficiently switch between running different processes

## GPUs

- **Specialized** – high performance, parallel processing, faster execution
- **Efficient resource utilization** – optimized for throughput, thousands of operations at once
- **Concurrency** – support large workloads, computational complex calculations without performance drops, concurrent tasks

# GPU use cases

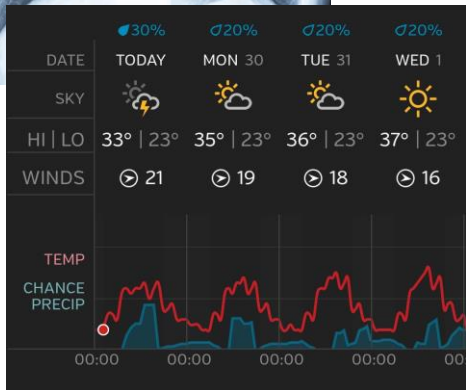


KubeCon



CloudNativeCon

North America 2024







KubeCon



CloudNativeCon

North America 2024





KubeCon



CloudNativeCon

North America 2024



# How can we control these costs?



## Monitor GPU metrics

NVIDIA DCGM Exporter



## Visualize GPU costs

OpenCost



## Optimize costs

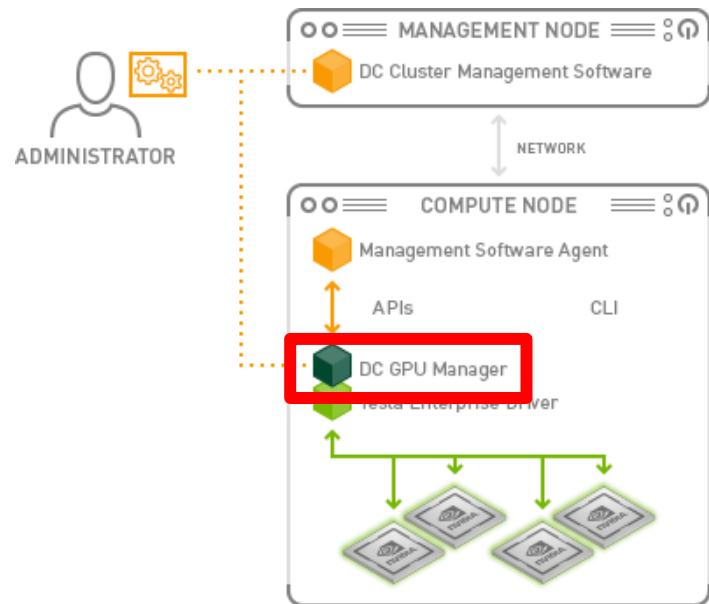
Rightsize and partition

## NVIDIA DC GPU Manager (DCGM)

- Suite of tools from NVIDIA for managing and monitoring GPUs in cluster and datacenter environments

## NVIDIA DCGM Exporter

- Metrics exporter for Prometheus
- Install methods
  - [DCGM container](#)
  - [NVIDIA GPU Operator](#)
- Integration with Grafana
- Extract valuable insights
  - Underutilization
  - GPU node health
  - Metrics customization



[NVIDIA DCGM](#) / [NVIDIA Developer](#)



General / Nvidia GPU Metrics ☆ 🔗

GPU 1937b558-347d-0f30-105b-893b98985668 ▾



Name

NVIDIA GeForce RTX 2080 SUPER

P-State

P8

Driver Version

471.11

Vbios Version

90.04.7a.40.73

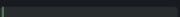
Throttle Reasons

Idle



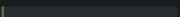
Active

HW Thermal Slowdown



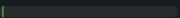
Not Active

SW Power Cap



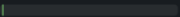
Not Active

App Clocks Setting



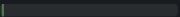
Not Active

HW Power Brake



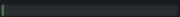
Not Active

SW Thermal Slowdown



Not Active

Sync Boost



Not Active

GPU Utilization %

0%

Power Draw %

11.5%

Fan Speed %

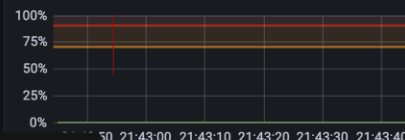
37%

Temperature

34 °C

2021-06-26 21:42:40

Memory Utilization % ▾



GPU Clock Speed %

0.395%

Memory Clock Speed %

0.310%

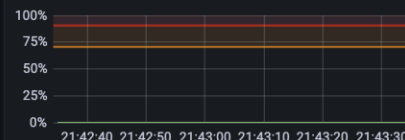
Memory Allocation %

11.2%

Memory Utilization %

0%

GPU Utilization %



Memory Allocation



Temperature



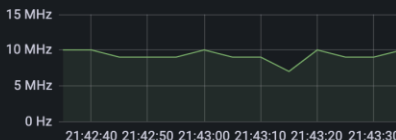
Power Draw



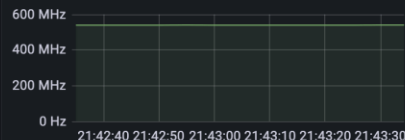
Fan Speed %



Graphics Clock Speed



Video Clock Speed



SM Clock Speed



Memory Clock Speed





GPU 1937b558-347d-0f30-105b-893b98985668 ▾

*i*

Name

NVIDIA GeForce RTX 2080 SUPER

*i*

P-State

P8

*i*

Driver Version

471.11

*i*

Vbios Version

90.04.7a.40.73

*i*

Throttle Reasons

Idle

Active

HW Thermal Slowdown

Not Active

SW Power Cap

Not Active

App Clocks Setting

Not Active

HW Power Brake

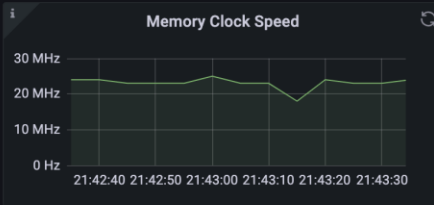
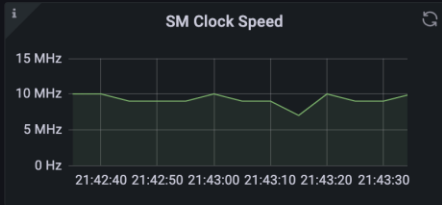
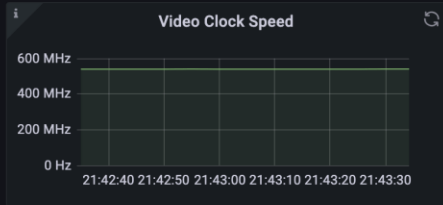
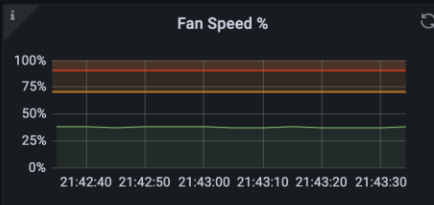
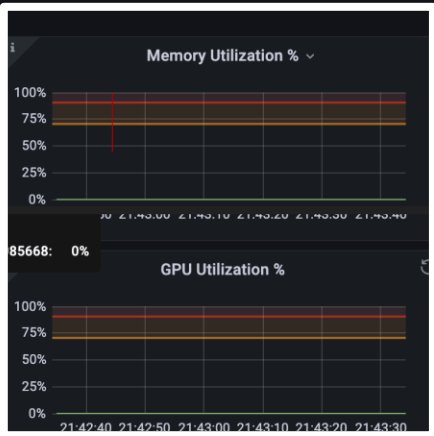
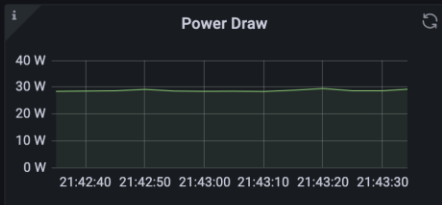
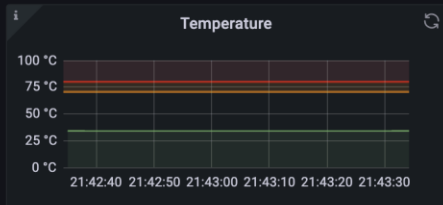
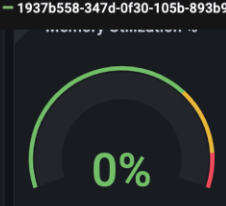
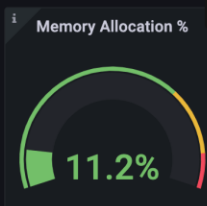
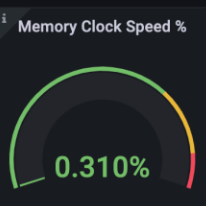
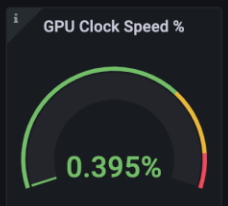
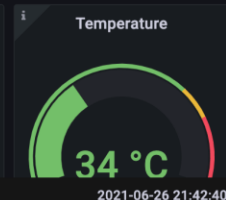
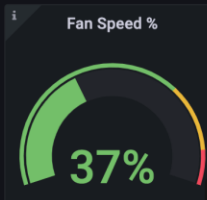
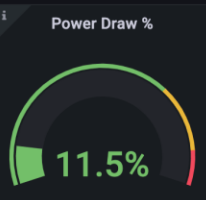
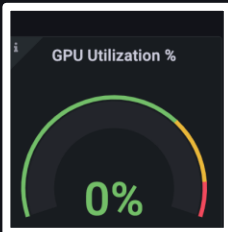
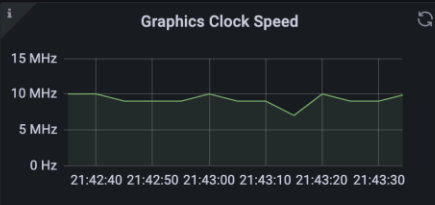
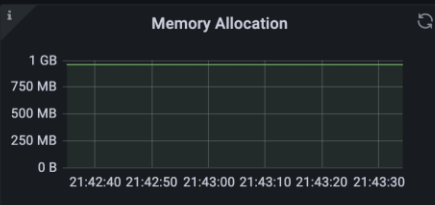
Not Active

SW Thermal Slowdown

Not Active

Sync Boost

Not Active



Hostname aks-nc24ads-12229106-vmss000000 GPU\_ID A3

## GPU Temperature



## GPU Power Usage



## GPU Framebuffer Mem Used



GPU_ID 11	Max: 326 MB	Avg: 100 MB
GPU_ID 13	Max: 326 MB	Avg: 152 MB
GPU_ID 7	Max: 326 MB	Avg: 151 MB
GPU_ID 10	Max: 12 MB	Avg: 12 MB
GPU_ID 11	Max: 12 MB	Avg: 12 MB
GPU_ID 12	Max: 12 MB	Avg: 12 MB
GPU_ID 13	Max: 12 MB	Avg: 12 MB
GPU_ID 7	Max: 12 MB	Avg: 12 MB
GPU_ID 8	Max: 12 MB	Avg: 12 MB
GPU_ID 9	Max: 12 MB	Avg: 12 MB

## GPU SM Clocks

180 GHz

## Tensor Core Utilization



## GPU Utilization

300%

## GPU Avg. Temp



## GPU Power Total



GPU_ID 13	35 °C	35 °C	35 °C
GPU_ID 7	35 °C	35 °C	35 °C
GPU_ID 10	37 °C	35.6 °C	35 °C
GPU_ID 11	37 °C	36.0 °C	35 °C
GPU_ID 12	37 °C	35.6 °C	35 °C
GPU_ID 13	37 °C	35.6 °C	35 °C
GPU_ID 7	37 °C	35.6 °C	35 °C
GPU_ID 8	37 °C	35.6 °C	35 °C
GPU_ID 9	37 °C	35.6 °C	35 °C

	max	avg	current
GPU_ID 11	75.6 W	75.2 W	75.1 W
GPU_ID 13	78.9 W	75.9 W	79.9 W
GPU_ID 7	75.6 W	75.0 W	75.6 W
GPU_ID 10	76.2 W	75.4 W	75.1 W
GPU_ID 11	76.5 W	75.8 W	75.4 W
GPU_ID 12	78.9 W	75.7 W	75.1 W
GPU_ID 13	76.5 W	75.7 W	75.5 W
GPU_ID 7	76.5 W	75.6 W	75.1 W

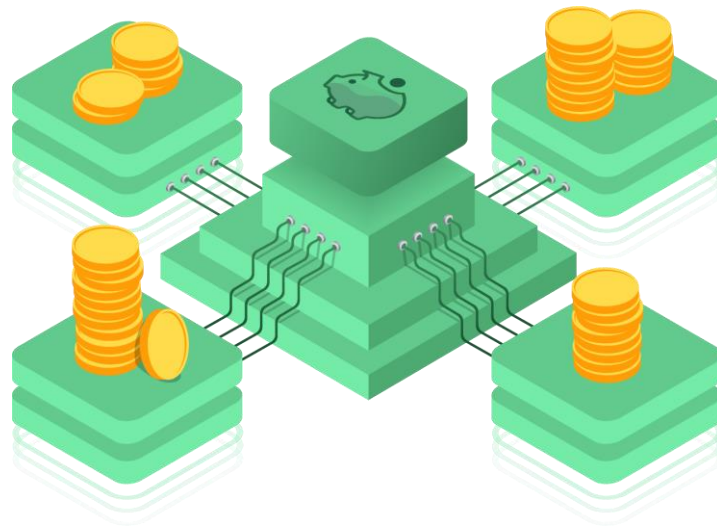
	max	avg	current
GPU_ID 11	0.00150%	0.000109%	0%
GPU_ID 13	0%	0%	0%
GPU_ID 7	0%	0%	0%
GPU_ID 10	0%	0%	0%
GPU_ID 11	0%	0%	0%
GPU_ID 12	0%	0%	0%
GPU_ID 13	0.00150%	0.0000958%	0%
GPU_ID 7	0.00140%	0.0000958%	0%



## OpenCost

- OSS, vendor neutral project
- OpenCost Spec standard for measuring, reporting, and allocating costs across clouds
- Visualize in near real time

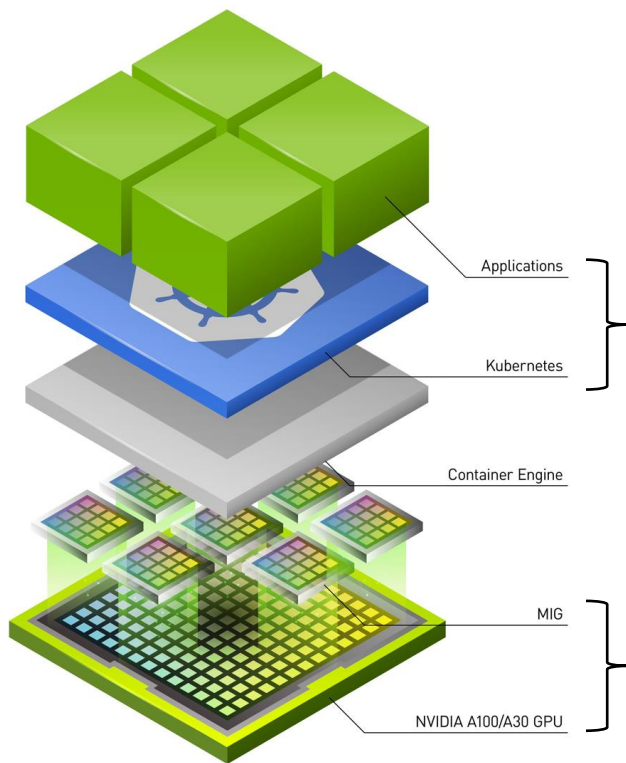
🎉 now a CNCF Incubating project 🎉  
<https://t.co/mywz7AKPTI>



[OpenCost](#)



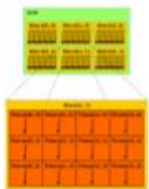
# Optimize GPU costs



Rightsize your apps and infrastructure

Divide and conquer via partitioning

# NVIDIA partitioning techniques



**Single Process  
in CUDA**



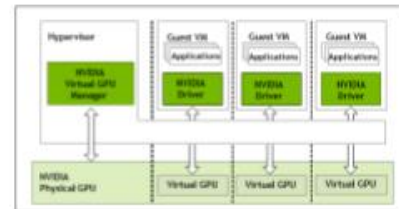
**Multi-Process  
with CUDA MPS**



**Time-slicing**



**MIG**



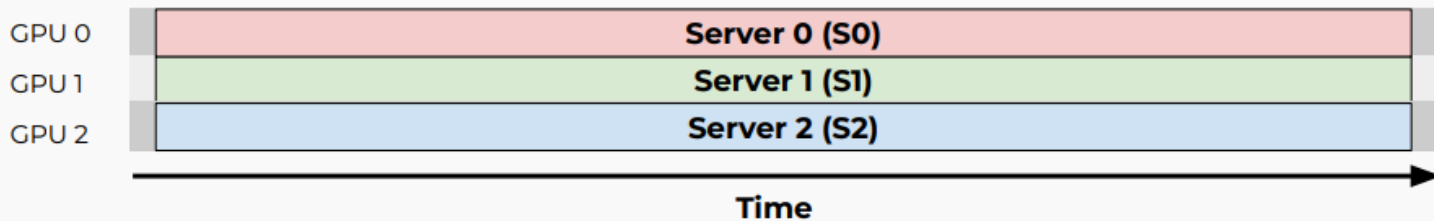
**Virtualization with vGPU**

Application level

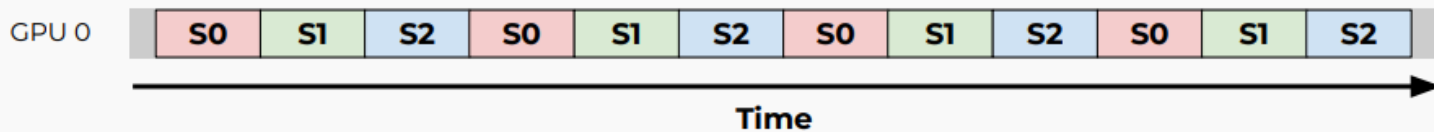
(using the CUDA programming  
model APIs - CUDA streams)

GPU System Software / Hardware  
(Mostly transparent to CUDA applications)

## Dedicated

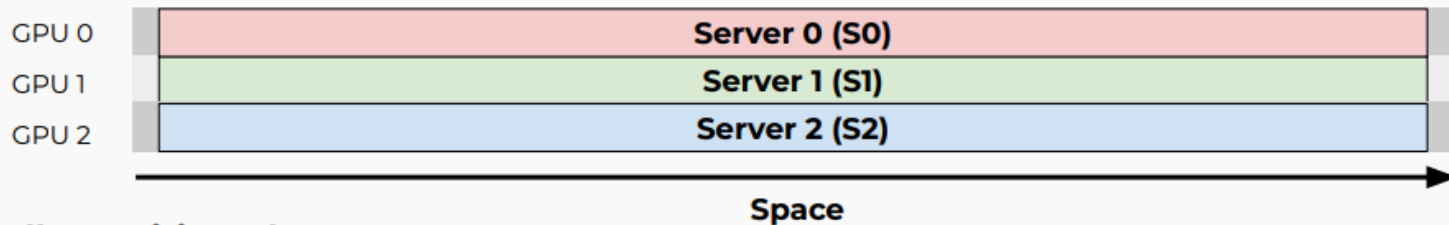


## Time-Sliced

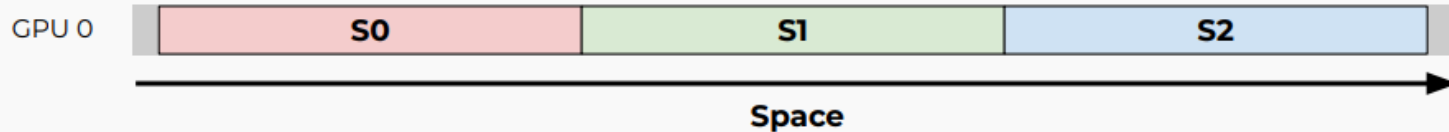


# Multi-Processing Service

## Dedicated

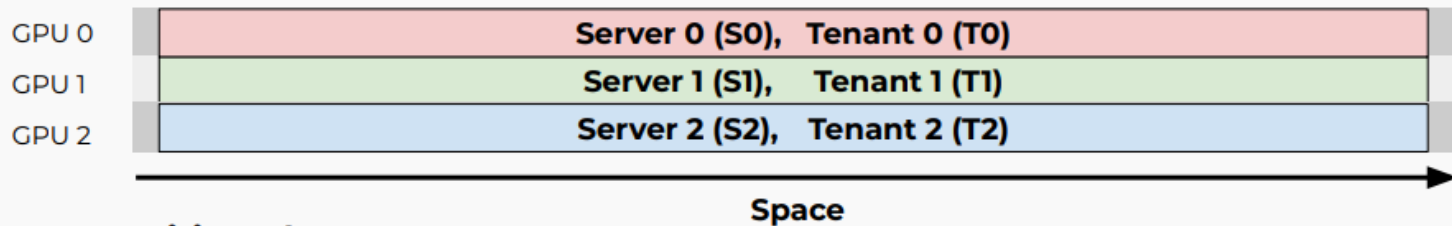


## Logically Partitioned

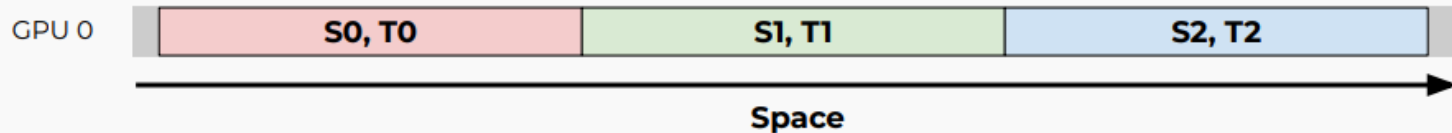


# Multi-instance GPUs (MIGs)

## Dedicated



## Hardware Partitioned





KubeCon



CloudNativeCon

North America 2024

# Demo

1. Provision GPU enabled cluster
2. Understand costs in OpenCost
3. Implement timeslicing

AKS@Azure:~\$

AKS@Azure:~\$ █



AKS@Azure:~\$ █

- GPUs are expensive but it can be feasible and economical to run AI workloads on GPUs at scale
- Monitor and visibility is key
  - Without insight into utilization and cost data, we don't know what to optimize and it's hard to keep teams accountable for their GPU spend
- Exclusive access is the default
- Partitioning is just one method of GPU optimization but it can really improve utilization and reduce costs
  - Timeslicing – best for apps that are not latency sensitive or can tolerate jitter
  - MIGs – best running multiple apps in parallel but need resiliency, QoS
  - MPS – best for running apps in parallel but can deal with limited resiliency



**KubeCon**



**CloudNativeCon**

North America 2024

# Thank you



**KubeCon**



**CloudNativeCon**

**North America 2024**