**Red Hat Summit**

**Connect**

# I servizi di inferenza per l'AI con HPE e Red Hat Openshift

Alessandro Moriondo, Hybrid Solutions Presales Manager
Hewlett Packard Enterprise - Italy

Milano, November 19th 2024

Hewlett Packard Enterprise | Red Hat

# Agenda

| HPE STRATEGY |
|---|

| HPE MACHINE LEARNING INFERENCE SOFTWARE |
|---|

| HPE | RED HAT OCP INFRASTRUCTURE BLUEPRINT |
|---|

| HPE | RED HAT OCP GPUS CONCURRENCY MODEL |
|---|

| Q&A |
|---|

# HPE Strategy – The key Pillars of Digital Enterprise

Edge
## Connect your edge
Control and harness data to innovate at the edge

Data
## Turn data into intelligence
Unify your data to make smart decisions

Cloud
## Create your hybrid cloud
Achieve the cloud experience everywhere

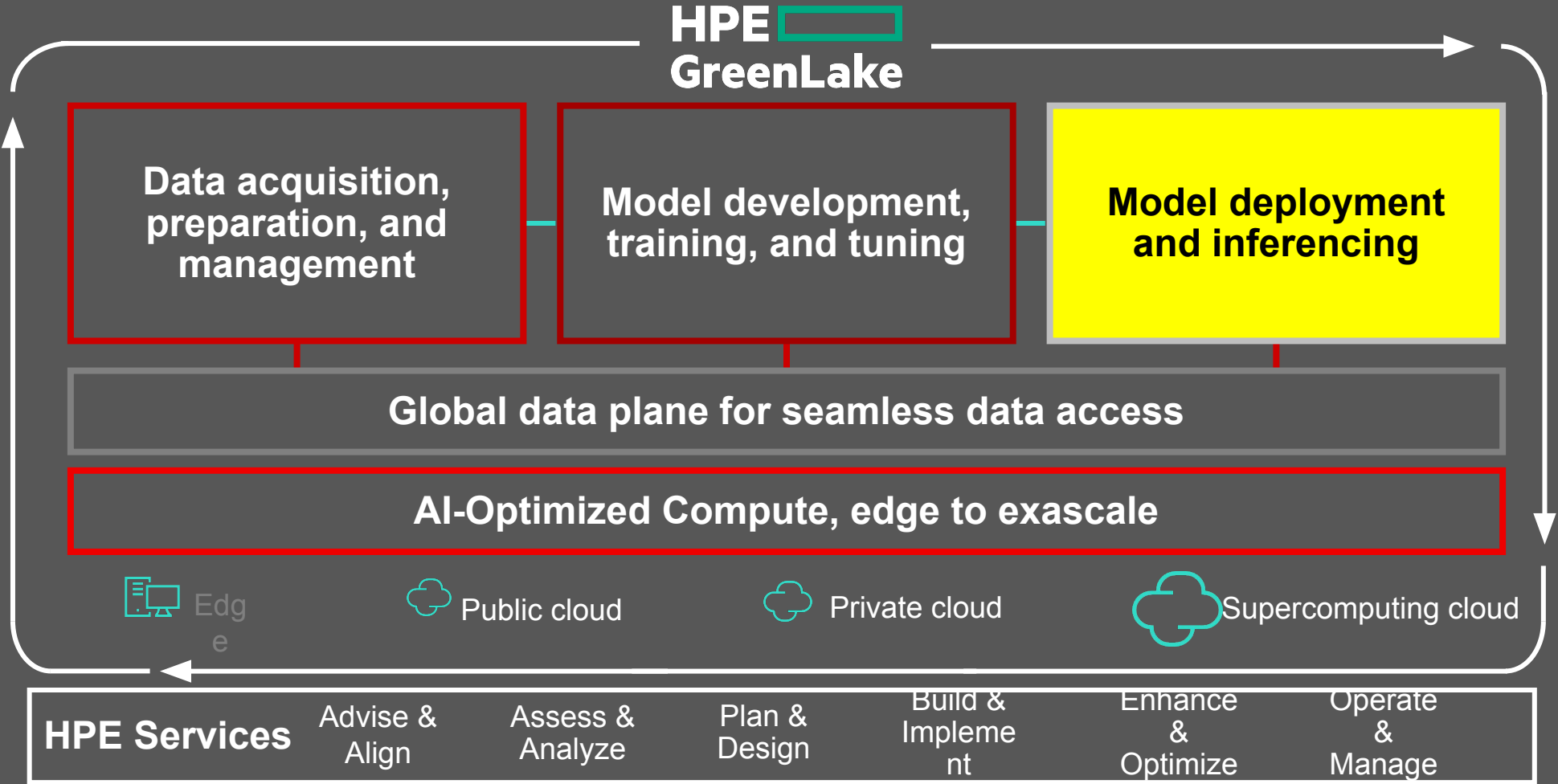Securit
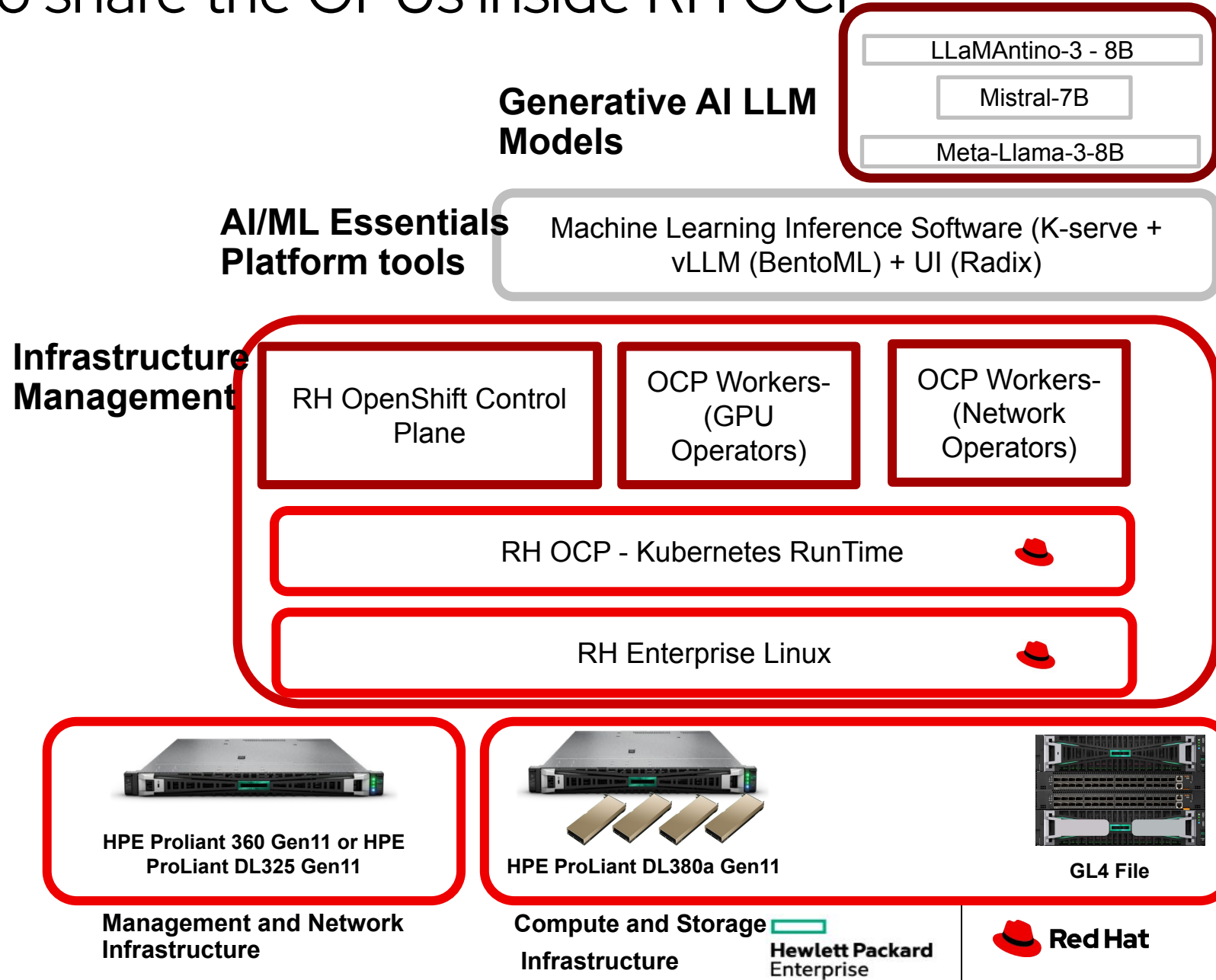y
## Secure your data
Secure your data from edge to cloud

Sustainability
## as a catalyst

Hewlett Packard Enterprise

Red Hat

# Unlock your competitive advantage with responsible AI at any scale

**HPE GreenLake**

| Data acquisition, preparation, and management | Model development, training, and tuning | Model deployment and inferencing |
|---|---|---|

**Global data plane for seamless data access**

**AI-Optimized Compute, edge to exascale**

Edge     Public cloud     Private cloud     Supercomputing cloud

| **HPE Services** | Advise & Align | Assess & Analyze | Plan & Design | Build & Implement | Enhance & Optimize | Operate & Manage |
|---|---|---|---|---|---|---|

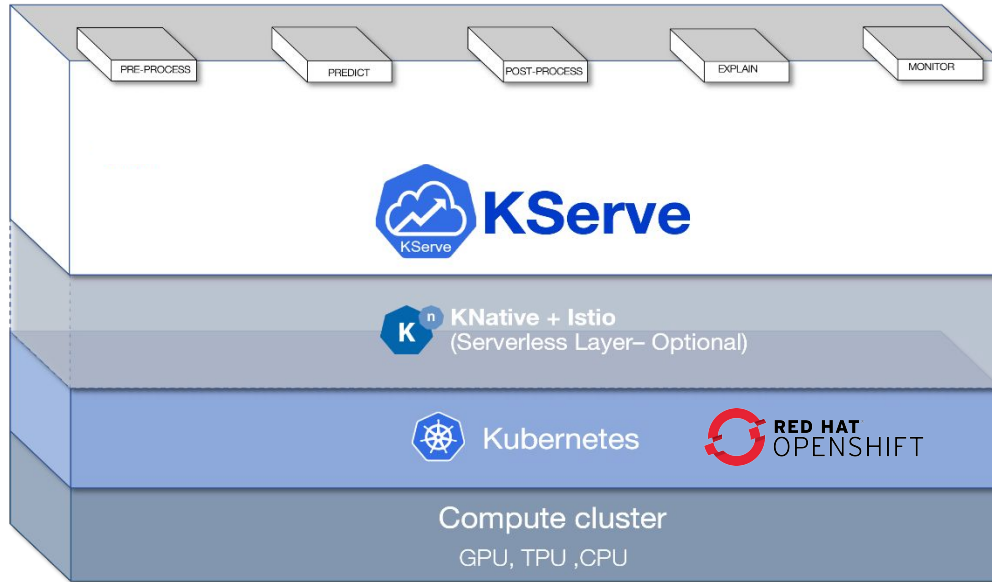Hewlett Packard Enterprise

Red Hat

# NVIDIA MIG to share the GPUs inside RH OCP

- **HW Compute stack** based on HPE Proliant Gen11 architecture
- **HW Network stack switch** at 100GB/s based on Aruba 8325 Switch for data and 6300 Switch for Mgmt at 1GB/s
- **HW Storage** based on latest GL4File NFS/S3 standard density rack
- **Gen AI models** choice either import from custom model or foundation models trained
- **AI/ML Platform sw tools**: choice of HPE MLIS ( Machine Learning Inference Software)
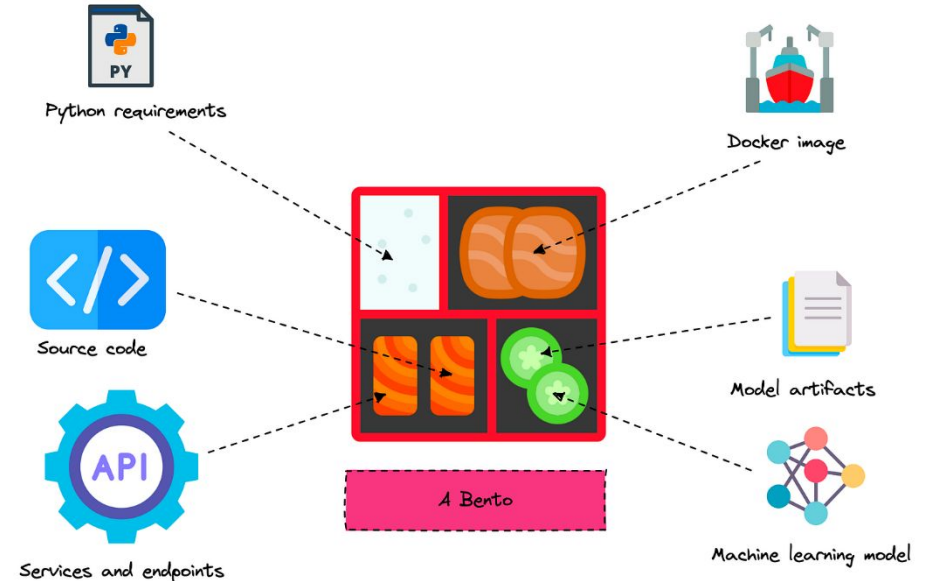- **Services:** HPE Deployment Inference Startup Service (DIY)

**Generative AI LLM Models**

| LLaMAntino-3 - 8B |
| Mistral-7B |
| Meta-Llama-3-8B |

**AI/ML Essentials Platform tools**

Machine Learning Inference Software (K-serve + vLLM (BentoML) + UI (Radix)

**Infrastructure Management**

| RH OpenShift Control Plane | OCP Workers- (GPU Operators) | OCP Workers- (Network Operators) |

RH OCP - Kubernetes RunTime

RH Enterprise Linux

**HPE Proliant 360 Gen11 or HPE ProLiant DL325 Gen11**

**HPE ProLiant DL380a Gen11**

**GL4 File**

**Management and Network Infrastructure**

**Compute and Storage Infrastructure**

Hewlett Packard Enterprise

Red Hat

# HPE MLIS Open source Components



## KServe

- Kubernetes-based platform for deploying models at scale
- Autoscaling, canary rollouts, and batch inferencing capabilities

## BentoML

- SDK for standardizing model packaging for services
- Serving standards for REST interfaces, logging, metrics
- **OpenLLM** – Support for optimized LLM deployments
- **vLLM** tackles the bottleneck of slow LLM inference, optimizing performance

# HPE Developed software Components

- **UI/UX**
  - Interface for managing and monitoring models, services, deployments, access tokens.
- **Security and authentication**
  - User management
  - Auth integration and access token management
- **Deployment APIs**
  - Reduce Kubernetes deployment friction
  - CLI and Python-native calls
- **Inferencing databases**
  - (Optionally) Capture data predictions
- **Integrated logging**
- **Metrics and Operations**
- **LLM deployment and support**

Hewlett Packard Enterprise | Red Hat

# HPE MLIS Stack

- **Platform**
  - Kubernetes (v1.20+)
  - Helm (v3.0+)
  - Knative
  - Istio
- **Serving**
  - Kserve (v0.11+)
- **Services**
  - BentoML
  - OpenLLM
- **Logging**
  - Loki
- **Metrics**
  - Prometheus
  - Grafana
- **Security**
  - Dex
- **UI**
  - Radix

Internal Users

External Users
e.g., LoB

**ML Engineer
(MLE)**

Feedback

Trained Model

(Source: MLDE,
Hugging Face, NGC
etc.)

Write Service

Build Service

Deploy Service

Experimental
Serving

**ML Operations
(MLOps)**

**IT Operations (ITOps)**

Deploy Service

Production
Serving

BentoML

MLIS

KServe

Hewlett Packard
Enterprise

Red Hat

Connect to existing registries:

**Connect to your model registry**

Select model

Configure Resources

Configure Scaling

Interact with your model

- NVIDIA NIM (NGC)
- OpenLLM (Hugging Face)
- AWS S3 Bucket
- Minio Registry

## Add new registry

A registry is required in order to hold your model(s) and deployments. Learn how to setup an openllm registry.

Name

hugging-face

Type

OpenLLM

Endpoint

https://huggingface.co

HuggingFace token

huggingface token

Cancel    Create registry

Hewlett Packard Enterprise

Red Hat

**Connect to your model registry**

**Select model**

**Configure Resources**

**Configure Scaling**

**Interact with your model**

Select a model
from the registry.

NousResearch/llama-2-7b-hf

mistral

HuggingFaceH4/zephyr-7b-alpha

HuggingFaceH4/zephyr-7b-beta

mistralai/Mistral-7B-Instruct-v0.2

mistralai/Mistral-7B-Instruct-v0.1

mistralai/Mistral-7B-v0.1

mixtral

mistralai/Mixtral-8x7B-Instruct-v0.1

mistralai/Mixtral-8x7B-v0.1

mpt

mosaicml/mpt-7b

mosaicml/mpt-7b-instruct

mosaicml/mpt-7b-chat

select a model...

Cancel    Next

Hewlett Packard
Enterprise

Red Hat

**Connect to your model registry**

**Select model**

**Configure Resources**

**Configure Scaling**

**Interact with your model**

Easily configure resources in the UI.

## Add new packaged model

A model is required for an inference deployment. Learn how to setup a model.

| Your model | Storage | **Resources** | Advanced (optional) |

**Resource Template**

☐ large-gpu ⌄

CPU
8

Memory
40Gi

GPU
4

GPU type
nvidia-tesla-a100

Cancel    Next

**Hewlett Packard Enterprise**    **Red Hat**

**Connect to your model registry**

**Select model**

**Configure Resources**

**Configure Scaling**

**Interact with your model**

Set your deployment to scale according to load.

**Create new deployment**

A deployment is a running instance of a packaged model. Learn how to setup a deployment.

Deployment    Packaged Model    Infrastructure    **Scaling**    Advanced (optional)

Auto scaling targets template

◎  scale-0-to-8-rps-20                                                              ⌄

Minimum instance                              Maximum instances

0                                              8

Auto scaling target

rps ⌄    20

Cancel    Back    Next

Hewlett Packard Enterprise    Red Hat

**Connect to your model registry**

**Select model**

**Configure Resources**

**Configure Scaling**

**Interact with your model**

Retrieve model predictions through APIs, CLI, or applications.

# Deploying AI/ML Models into Production

**AI/ML Application**

**ML Operations (MLOps)**

**ML Engineer (MLE)**

**IT Operations (ITOps)**

**LoB**

| | | |
|---|---|---|
| Model Training | Model Registry | |
| Experimental Inference | | |

| | | |
|---|---|---|
| Model Optimization | Production Inference | Pre-Processing |
| | Model Explainability | |

Data Acquisition

Requests

Predictions

Control

Container Orchestration

Infrastructure Management

Metrics

# Example of physical Architecture with RH Openshift



**Core Switch**

100 Gbe SFP+

OCP Master node #1

OCP Master node #2

OCP Master node #3

100GBe SFP+

**Master Node: HPE DL360 Gen11**

Data Network Switch

2x100 Gbps / IB

1 Gbe Mgmt Net Switch

2x100 Gbps / IB

OCP Worker-1 with Hopper GPU

OCP Worker-2 with Hopper GPU

OCP Worker-3 with Hopper GPU

OCP Worker-4 with Hopper GPU

**Worker Node: HPE DL380a Gen11**

2x100Gb

2x100Gbp

2x100Gbps MLAG

2x100G bps

2x100Gb ps

**HPE Greenlake4file**

Hewlett Packard Enterprise
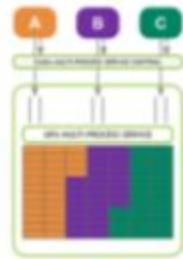
Red Hat

# NVIDIA GPUs Concurrency choices



GPU "CONCURRENCY"

Choices

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)

Hewlett Packard Enterprise

Red Hat

# NVIDIA GPUs Concurrency choices

| | Streams | MPS | Time-Slicing | MIG | vGPU |
|---|---|---|---|---|---|
| Partition Type | Single process | Logical | Temporal (Single process) | Physical | Temporal & Physical – VMs |
| Max Partitions | Unlimited | 48 | Unlimited | 7 | Variable |
| SM Performance Isolation | No | Yes (by percentage, not partitioning) | Yes | Yes | Yes |
| Memory Protection | No | Yes | Yes | Yes | Yes |
| Memory Bandwidth QoS | No | No | No | Yes | Yes |
| Error Isolation | No | No | Yes | Yes | Yes |
| Cross-Partition Interop | Always | IPC | Limited IPC | Limited IPC | No |
| Reconfigure | Dynamic | At process launch | N/A | When idle | N/A |
| GPU Management (telemetry) | N/A | Limited GPU metrics | N/A | Yes – GPU metrics, support for containers | Yes – live migration and other industry virtualization tools |
| Target use cases (and when to use each) | Optimize for concurrency within a single application | Run multiple applications in parallel but can deal with limited resiliency | Run multiple applications that are not latency-sensitive or can tolerate jitter | Run multiple applications in parallel but need resiliency and QoS | Support multi-tenancy on the GPU through virtualization and need VM management benefits |

Hewlett Packard Enterprise

Red Hat

# Red Hat Summit Connect

# Q&A