# Explain
# How Kubernetes Works with GPUs
# Like I'm 5

**Carlos Santana**

Sr. WW Spec. SA, Containers at AWS
CNCF Ambassador

# Learning at home
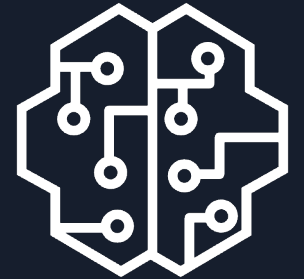
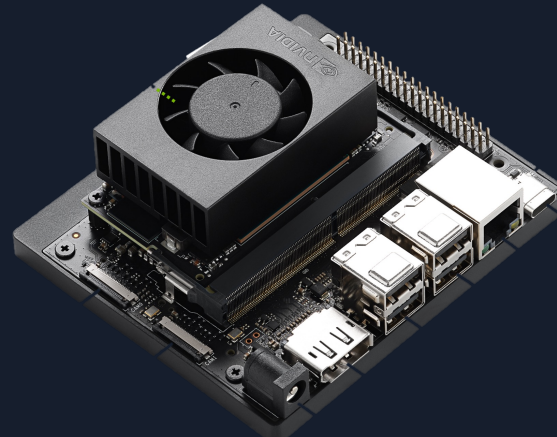**kubernetes**

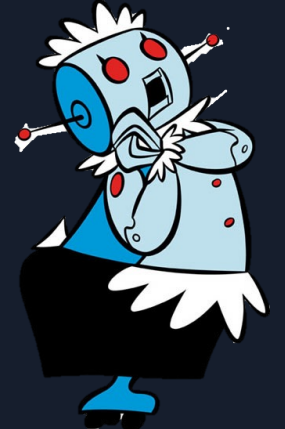**+**



**=**





**k8s
experience**



**NVIDIA
Jetson GPU**



**Generative
AI & ML**

Images:
https://warnerbros.fandom.com/wiki/Rosey
https://developer.nvidia.com/embedded/learn/get-started-jetson-orin-nano-devkit

2

# Learn and Be Curious

# Main areas

## HOST
Operating System
Containerd

**Device Driver**

**Container Toolkit**

## CONTAINER
Programing model

**CUDA**

## KUBERNETES
Kubelet
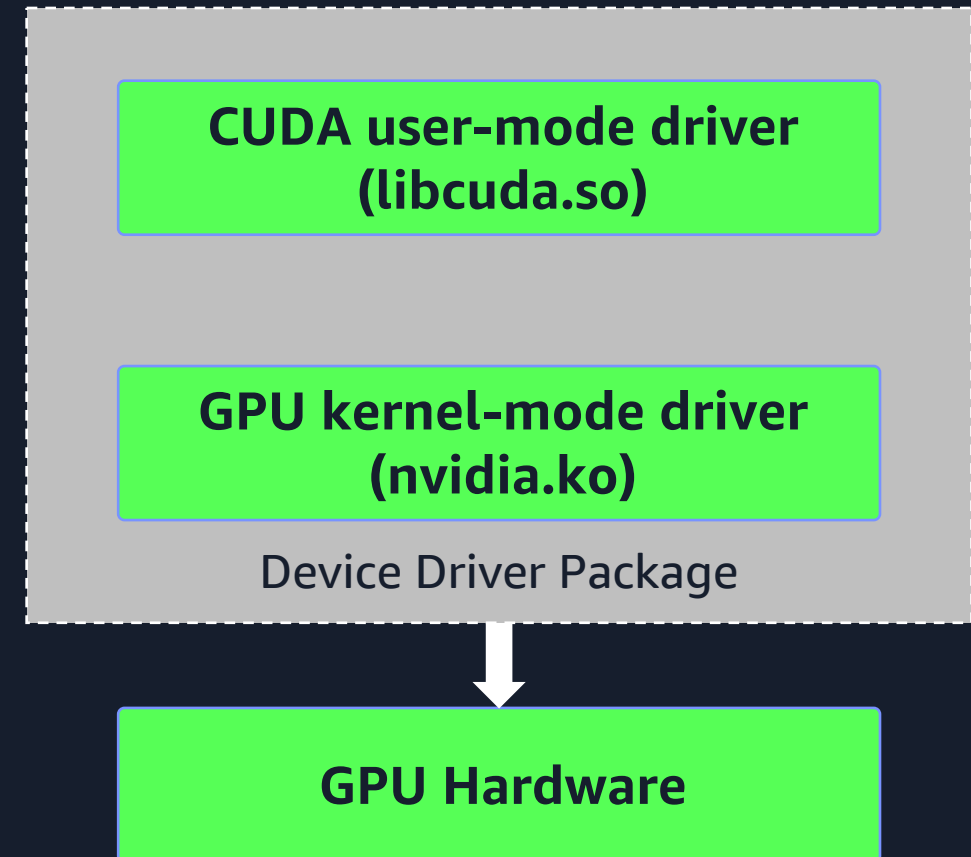Node Labels

**Device Plugin**

**Node Feature Discovery**

**GPU Feature Discovery**

# Device Driver

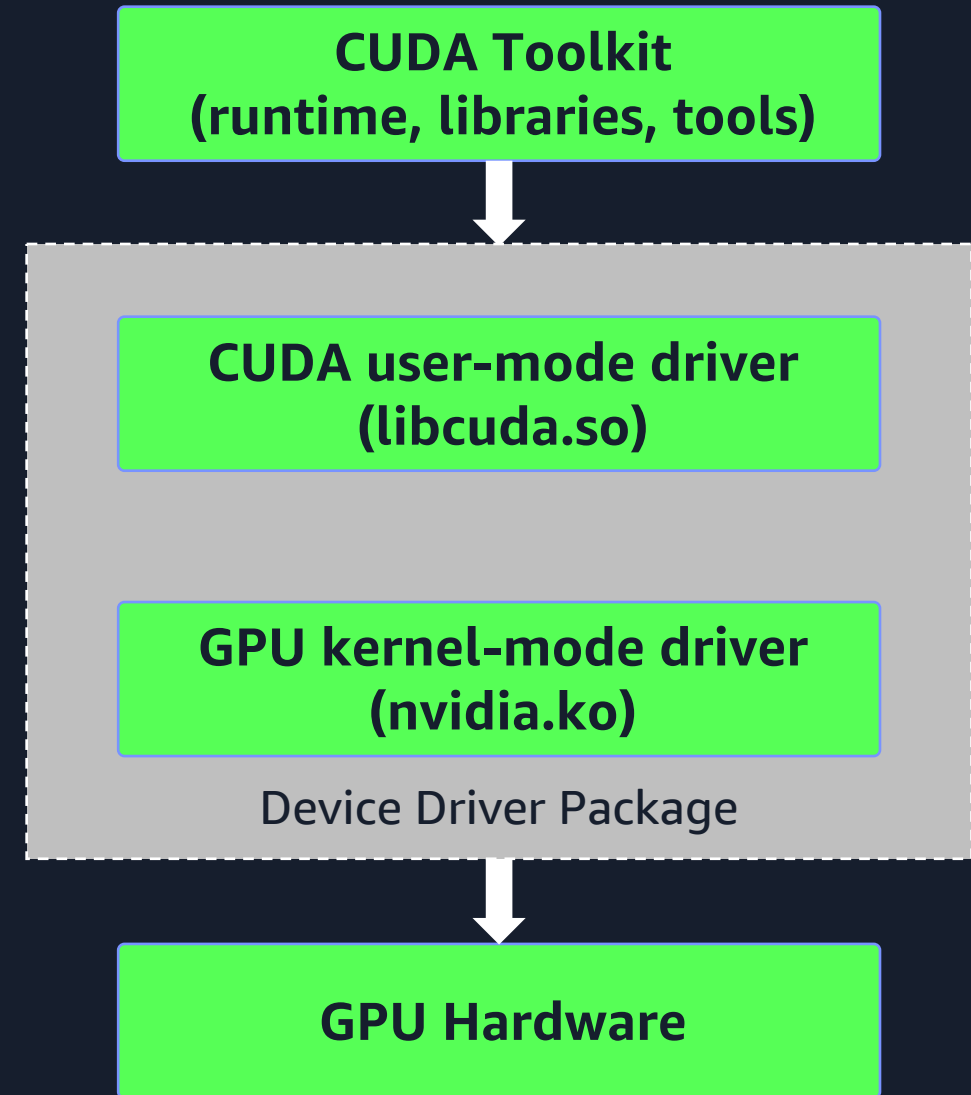- Implements the CUDA Driver API
- Use Vendor pre-installed drivers
- NVIDIA Jetson JetPack OS
- Cloud Providers include driver in OS images
  - Example:
    - AWS EKS AIM EKS AL23, BR, Auto Mode
- Driver package include CUDA user mode driver
- Do not use GPU Operator unless you know why

CUDA user-mode driver
(libcuda.so)

GPU kernel-mode driver
(nvidia.ko)

Device Driver Package

GPU Hardware

# CUDA Toolkit

- Toolkit at build time necessary for linking
- Compatibility between CUDA and Driver
- Container Image usually comes with CUDA runtime
- Do not mount toolkit runtime from host path, unless you know why



CUDA Toolkit
(runtime, libraries, tools)

CUDA user-mode driver
(libcuda.so)

GPU kernel-mode driver
(nvidia.ko)

Device Driver Package

GPU Hardware

# Driver and CUDA versions

nvidia-smi for PCI GPUs
On Jetson dGPU:
	/proc/driver/nvidia/version

nvcc (CUDA compiler)
ldconfig

```
$ cat /proc/driver/nvidia/version
VRM version: NVIDIA UNIX Open Kernel
Module for aarch64  540.3.0

$ head -n 1 /etc/nv_tegra_release
# R36 (release), REVISION: 3.0 (aka 36.3)
```

```
$ nvcc --version
Cuda compilation tools, release 12.2, V12.2.140

$ ldconfig –p | grep libcudart
libcudart.so.12 (libc6,AArch64)
```

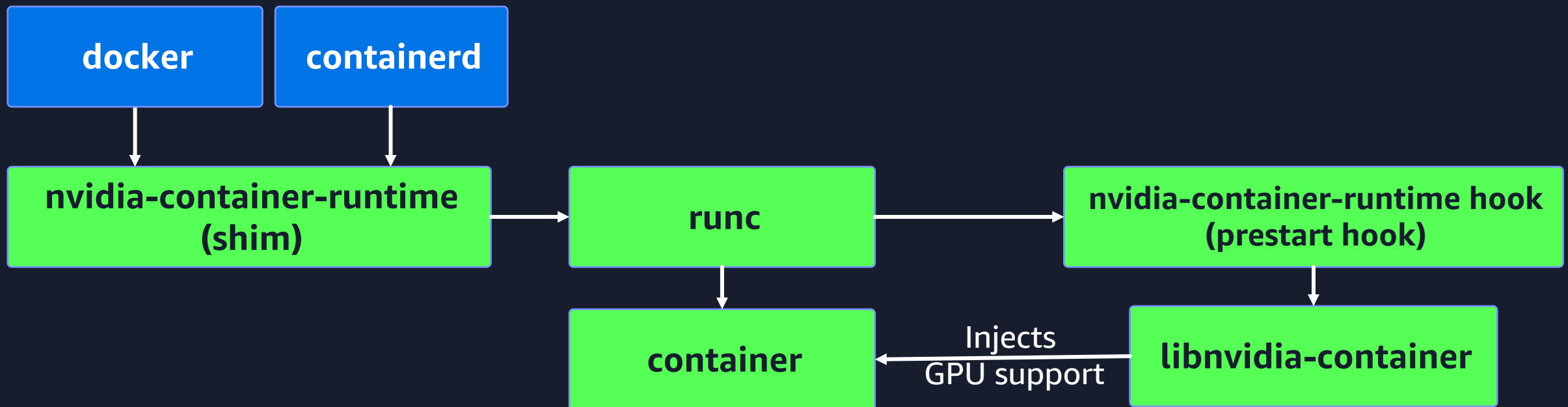| CUDA Toolkit | Minimum Required Driver Version |
|--------------|---------------------------------|
| CUDA 12.x    | >=525.60.13                     |
| CUDA 11.x    | >=450.80.02                     |

# Container Toolkit

- Run GPU-accelerated containers
- Cloud Providers includes this
- For Jetson you need to install docker engine (containerd) and install the toolkit

```
$ nvidia-ctk runtime configure  --runtime=containerd
$ cat /etc/containerd/config.toml
default_runtime_name = "nvidia"
[plugins."io.containerd.grpc.v1.cri".containerd.runtimes.nvidia.options]
  BinaryName = "/usr/bin/nvidia-container-runtime"
```



8

# GPU Containers

| Device Plugin |
| --- |

| GPU Feature Discovery |
| --- |

| Node Feature Discovery |
| --- |

$ docker run --runtime nvidia --gpus all

✓ | Container Toolkit |
| --- |

✓ | CUDA |
| --- |

✓ | Device Driver |
| --- |

9

# Jetson Containers

| | |
|---|---|
| **ML** | pytorch tensorflow jax onnxruntime deepstream holoscan CTranslate2 JupyterLab |
| **LLM** | SGLang vLLM MLC AWQ transformers text-generation-webui ollama llama.cpp llama-factory exllama AutoGPTQ FlashAttention DeepSpeed bitsandbytes xformers |
| **VLM** | llava llama-vision VILA LITA NanoLLM ShapeLLM Prismatic xtuner |
| **VIT** | NanoOWL NanoSAM Segment Anything (SAM) Track Anything (TAM) clip_trt |
| **RAG** | llama-index langchain jetson-copilot NanoDB FAISS RAFT |
| **L4T** | l4t-pytorch l4t-tensorflow l4t-ml l4t-diffusion l4t-text-generation |
| **CUDA** | cupy cuda-python pycuda numba opencv:cuda cudf cuml |
| **Robotics** | Cosmos Genesis ROS LeRobot OpenVLA 3D Diffusion Policy Crossformer MimicGen OpenDroneMap ZED |
| **Graphics** | stable-diffusion-webui comfyui nerfstudio meshlab pixsfm gsplat |
| **Mamba** | mamba mambavision cobra dimba videomambasuite |
| **Speech** | whisper whisper_trt piper riva audiocraft voicecraft xtts |
| **Home/IoT** | homeassistant-core wyoming-whisper wyoming-openwakeword wyoming-piper |

10

# Jetson Containers

- Modular container build system
- Latest AI/ML packages
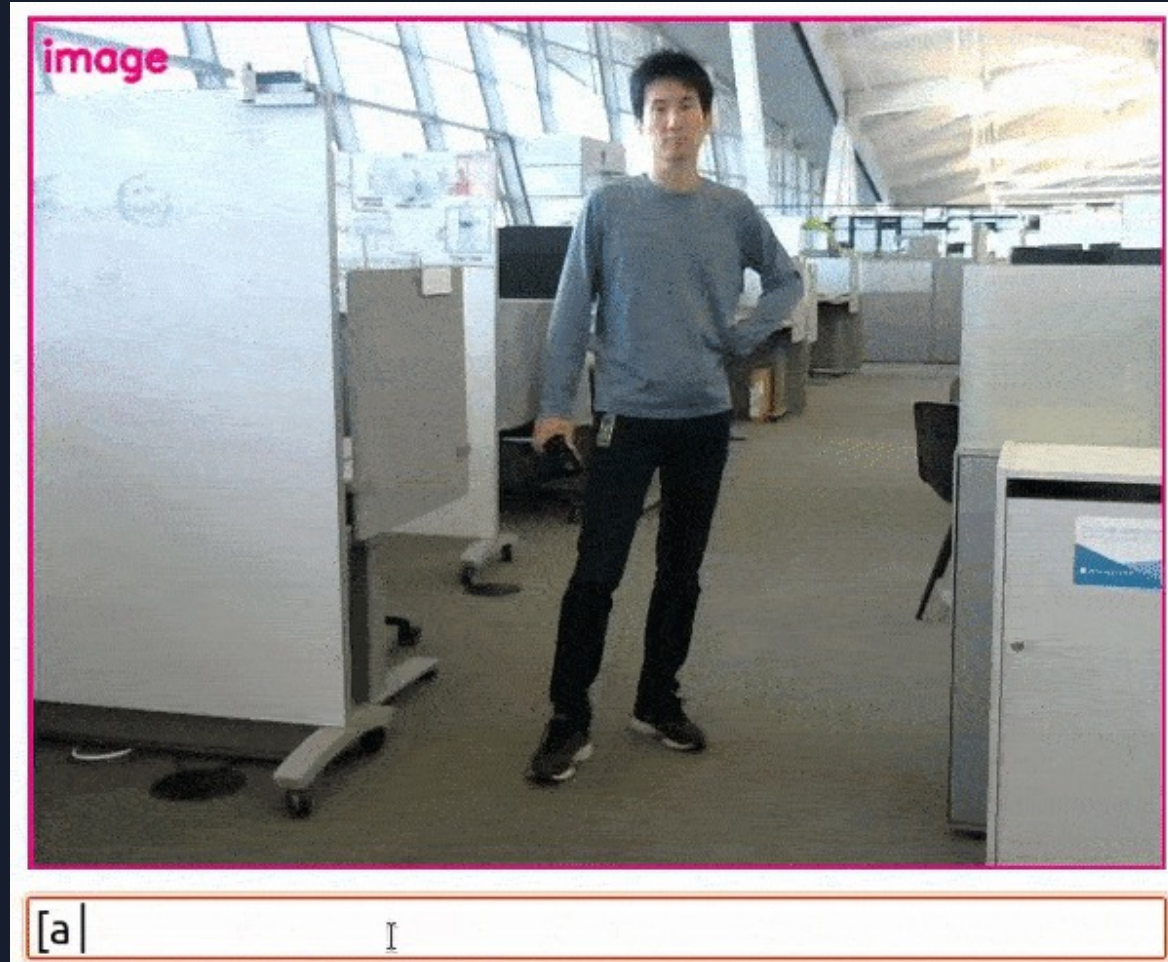- Wraps docker run --runtime=nvidia
- Discord and community meetings

```
$ jetson-containers run --name ollama $(autotag ollama)

>>> How many kubernetes administrators takes to deploy a
cluster?

In a typical scenario, a small team of 2-3 administrators may be
sufficient for deploying and managing a Kubernetes cluster in a
simple environment.
However, larger organizations or complex deployments could
require teams of 5-10 or more administrators to ensure proper
management and scalability.
```

# Kubernetes

- Device Plugin required?
  - No scheduling
  - There is no allocations
  - All Pods access all GPUs
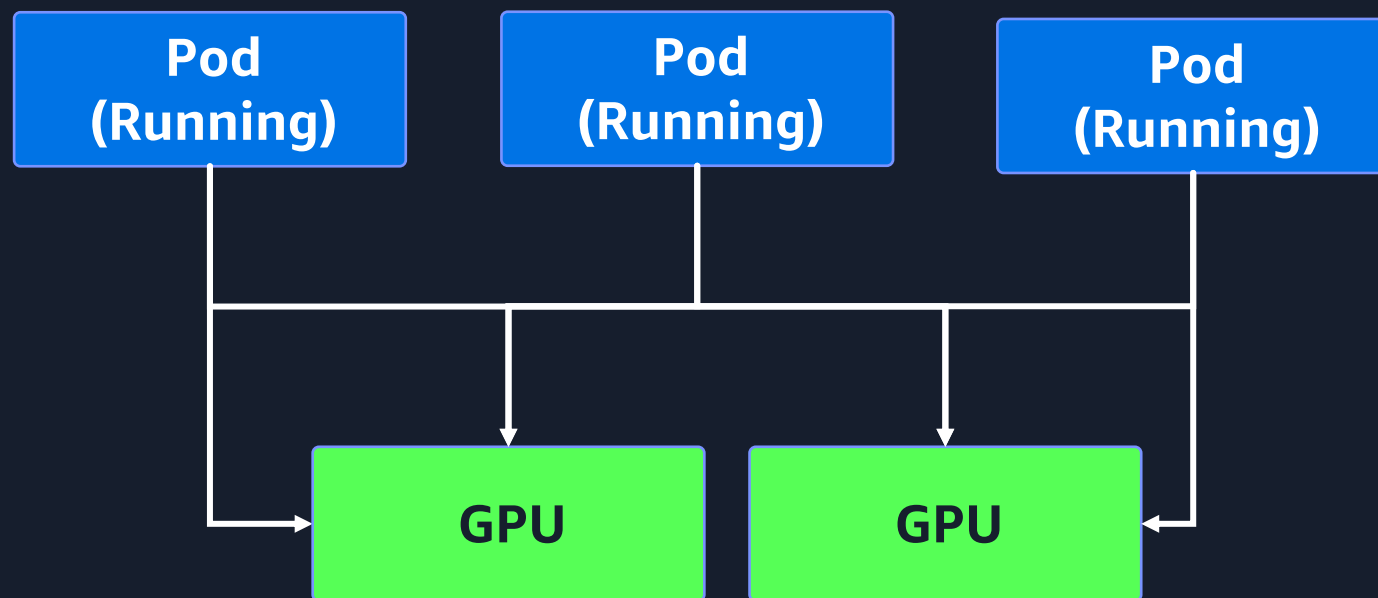
```
$ kubectl run myllm --image myllmImage
```

Without allocatable nvidia.com/gpu



✅ **Container Toolkit**

✅ **CUDA**

✅ **Device Driver**

NODE

Pod (Running)   Pod (Running)   Pod (Running)

GPU   GPU

# Kubernetes + GPU Device Plugin

- NVIDIA Device Plugin
  - Helm Chart
- Schedule Pods
- Allocate GPUs
- Share a GPU
  - Time-slicing
  - MPS
  - MIG



NVIDIA Device Plugin helm chart

**Device Plugin**

**GPU Feature Discovery**

**Node Feature Discovery**

# Node Feature Discovery (NFD)

- Advertise features using labels
- Detects hardware on node
- CPU, Kernel, OS, PCI
- PCI VENDOR NVIDIA (10ed)

### NVIDIA Device Plugin helm chart

**Device Plugin**

**GPU Feature Discovery**

**Node Feature Discovery**

Enables the GPU Feature Discovery ⟶

```
$ kubectl get no jetson -o json

"cpu-cpuid.USCAT": "true",
"cpu-model.vendor_id": "ARM",
"kernel-config.PREEMPT": "true",
"kernel-version.full": "5.15.136-tegra",
"storage-nonrotationaldisk": "true",
"system-os_release.ID": "ubuntu",
"system-os_release.VERSION_ID": "22.04",

"pci-10ed.present": "true",
```

# GPU Feature Discovery (GFD)

- Advertise features using labels
- Detects NVIDIA GPU
- Labels with nvidia.com/*

NVIDIA Device Plugin helm chart

Device Plugin

**GPU Feature Discovery**

Node Feature Discovery

```
$ kubectl get no jetson -o json

"nvidia.com/cuda.driver-version.full": "540.3.0",
"nvidia.com/cuda.runtime-version.full": "12.2",
"nvidia.com/gpu.count": "1",
"nvidia.com/gpu.replicas": "4",
"nvidia.com/gpu.product": "Orin-SHARED",
"nvidia.com/gpu.sharing-strategy": "time-slicing",
"nvidia.com/mig.capable": "false",
"nvidia.com/mps.capable": "false",
"nvidia.com/vgpu.present": "false"
```

Enables the Device Plugin deamon ⟶ **"nvidia.com/gpu.present": "true",**

# Device Plugin

- Kubelet Device Plugin API
- Expose the number of GPUs
- Tracks health of the GPUs

Kubelet patches this status

NVIDIA Device Plugin helm chart

**Device Plugin**

GPU Feature Discovery

Node Feature Discovery

`limit` instead of `request`

```
k get nodes jetson -o json | jq .status.allocatable
{
  "cpu": "8",
  "memory": "16032384Ki",
  "nvidia.com/gpu": "4",
  "pods": "110"
}
```

```
containers:
- name: cuda-container
  image: nvcr.io/nvidia/k8s/cuda-sample
  resources:
    limits:
      nvidia.com/gpu: 1 # requesting 1 GPU
```

# Device Plugin

# Scheduling Pods

- Allocatable status on Node
- Pods states how much GPU they need
- As Pod wait for GPU, they are pending

```
containers:
- name: cuda-container
  resources:
    limits:
      nvidia.com/gpu: 1 # requesting 1 GPU
```
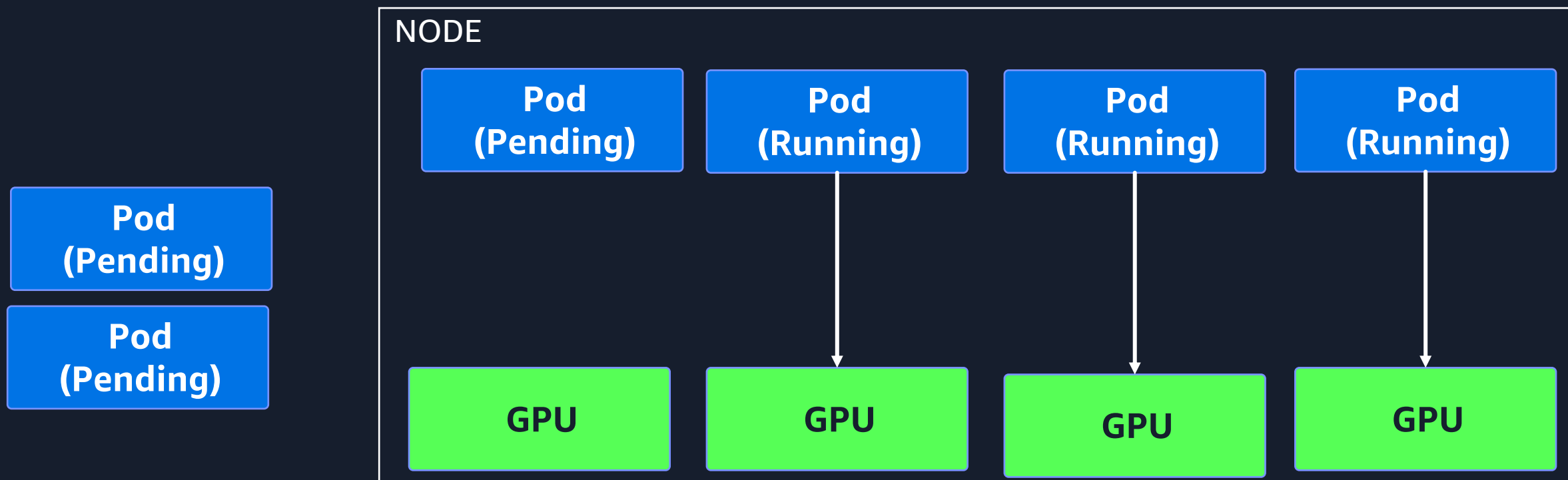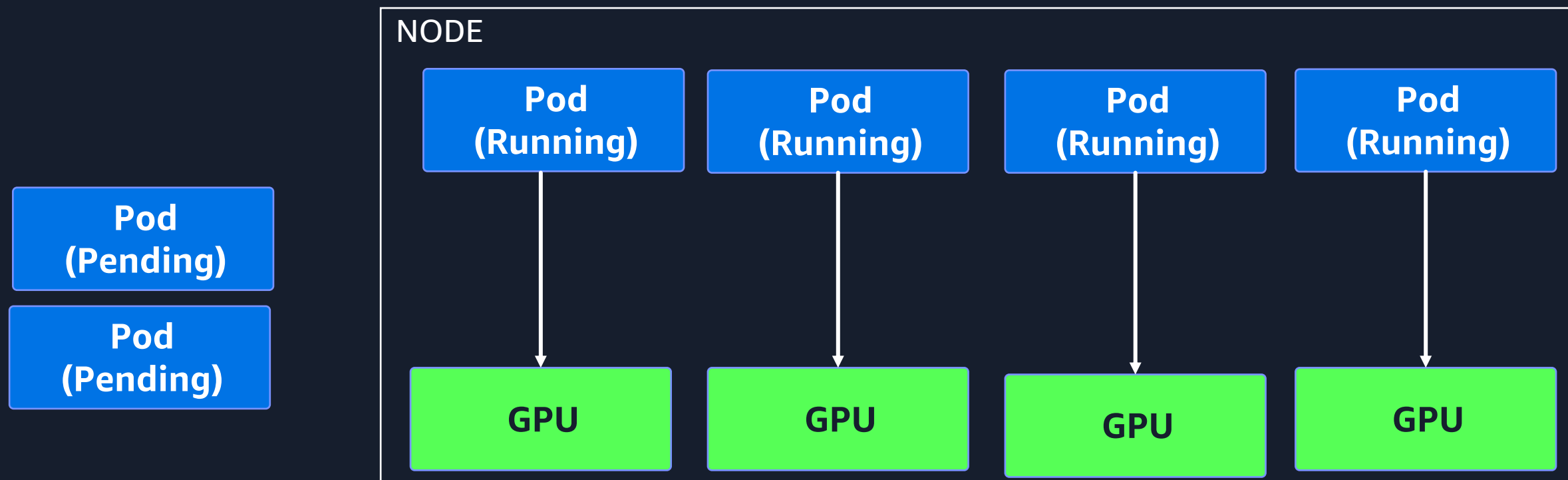
NODE

| Pod (Pending) | | Pod (Running) | Pod (Running) | Pod (Running) | Pod (Running) |

| Pod (Pending) |

| Pod (Pending) |

| GPU | GPU | GPU | GPU |

19

# Scheduling Pods

- Allocatable status on Node
- Pods states how much GPU they need
- As Pod wait for GPU, they are pending

```
containers:
- name: cuda-container
  resources:
    limits:
      nvidia.com/gpu: 1 # requesting 1 GPU
```

**NODE**

| Pod (Pending) | Pod (Complete) | Pod (Running) | Pod (Running) | Pod (Running) |

Pod (Pending)

Pod (Pending)

| GPU | GPU | GPU | GPU |

# Scheduling Pods

- Allocatable status on Node
- Pods states how much GPU they need
- As Pod wait for GPU, they are pending

```
containers:
- name: cuda-container
  resources:
    limits:
      nvidia.com/gpu: 1 # requesting 1 GPU
```

**Pod (Pending)**

**Pod (Pending)**

NODE

| | | | |
|---|---|---|---|
| **Pod (Pending)** | **Pod (Running)** | **Pod (Running)** | **Pod (Running)** |
| **GPU** | **GPU** | **GPU** | **GPU** |

# Scheduling Pods



- Allocatable status on Node
- Pods states how much GPU they need
- As Pod wait for GPU, they are pending

```
containers:
- name: cuda-container
  resources:
    limits:
      nvidia.com/gpu: 1 # requesting 1 GPU
```

# Why and How

## Run Containers



**Device Driver**

**Container Toolkit**

**CUDA**

## Orchestrate Containers



**kubernetes**

**Device Plugin**

**Node Feature Discovery**

**GPU Feature Discovery**

# Kubernetes Home lab

control-plane

**Intel NUC (amd64)**

api-server

- k3s
- kubeadm
- etc

**Intel NUC (amd64)**

kubelet
containerd

**Jetson Orin**

kubelet
containerd

**GPU**

control-plane

**AWS**

**EC2**

wireguard

**EKS**
api-server

vpn

- nodeadm

**Turing Pi**

kubelet
containerd

Wireguard client
EKS Hybrid Nodes

**RK1
NPU**

**Intel NUC (amd64)**

kubelet
containerd

Wireguard client
EKS Hybrid Nodes

**tp-link 24 port switch**

# Deep Seek R1 on Jetson kubernetes

# Resources

- [https://github.com/csantanapr/k8s-nvidia](https://github.com/csantanapr/k8s-nvidia)