



KubeCon



CloudNativeCon

THE LINUX FOUNDATION



China 2024



KubeCon



CloudNativeCon



China 2024

# Unlocking the Power of Kubernetes (for AI)

AI-Driven Innovations for  
Next-Gen Infrastructure

Brandon Kang @ Akamai

# Who am I?



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

- Brandon Kang (姜相鎮)
- Akamai Technologies, APJ CTG
- Principal Technical Solutions Architect



# Back to Basics – Why Kubernetes?



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

- Scalability & Portability
- Self-Healing & Extensibility
- Resource Optimization
- Automated Rollouts and Rollbacks
- Service Discovery and Load Balancing
- Etc.



# Why Deploy AI on Kubernetes?



China 2024

- Dynamic Resource Scaling
- Burst Workloads
- Resource Isolation
- Shared Infrastructure
- GPU/TPU Management
- Efficient Utilization
- Consistency
- Observability

# Why Deploy AI on Kubernetes?



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

- **Resource Management**

Kubernetes can schedule and allocate resources efficiently, ensuring that AI workloads make optimal use of available hardware, such as GPUs and TPUs, which are often essential for AI tasks

- **Cost Management**

Efficient resource management helps in controlling costs, especially in cloud environments where resource usage directly impacts expenses.

# Why Deploy AI on Kubernetes?



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

- **Custom Resource Definitions (CRDs)**  
Kubernetes allows for the creation of custom resources to support specialized AI workloads and requirements
- **Integration with AI Tools**  
Kubernetes integrates well with various AI frameworks and tools, such as TensorFlow, PyTorch, and Nvidia's NeMo, facilitating seamless deployment and operation of AI models

# Why Deploy AI on Kubernetes?



KubeCon



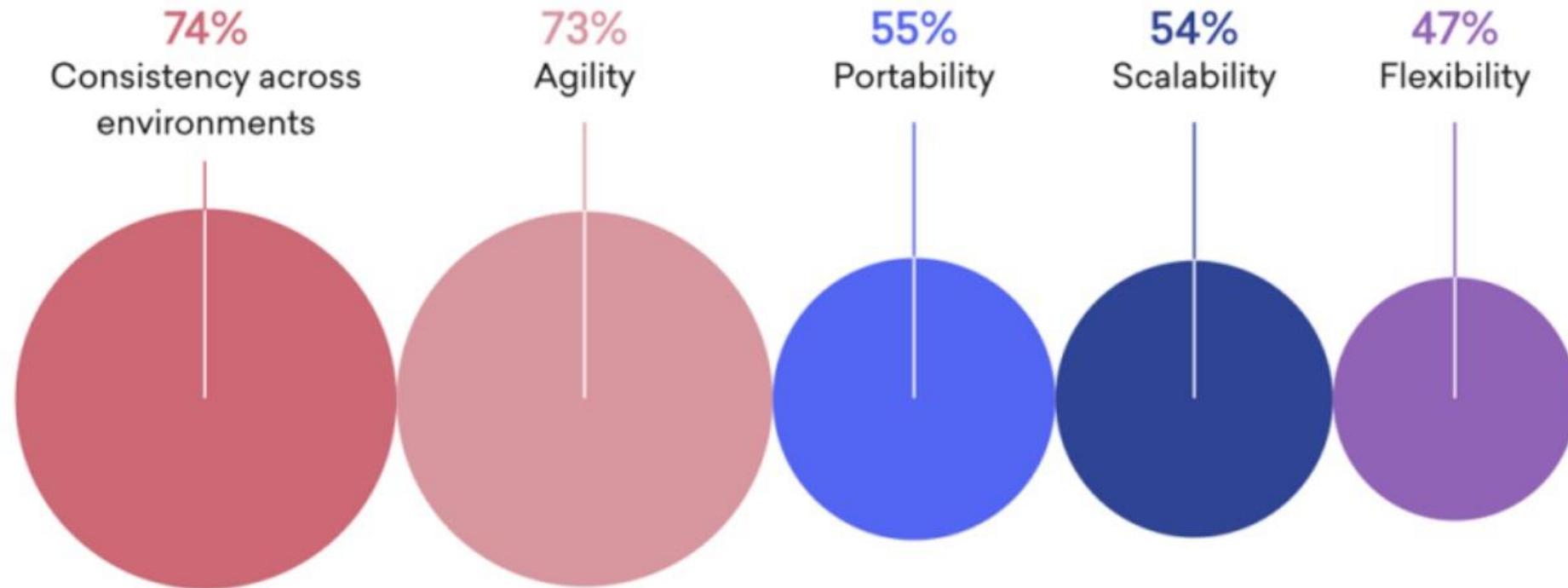
CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

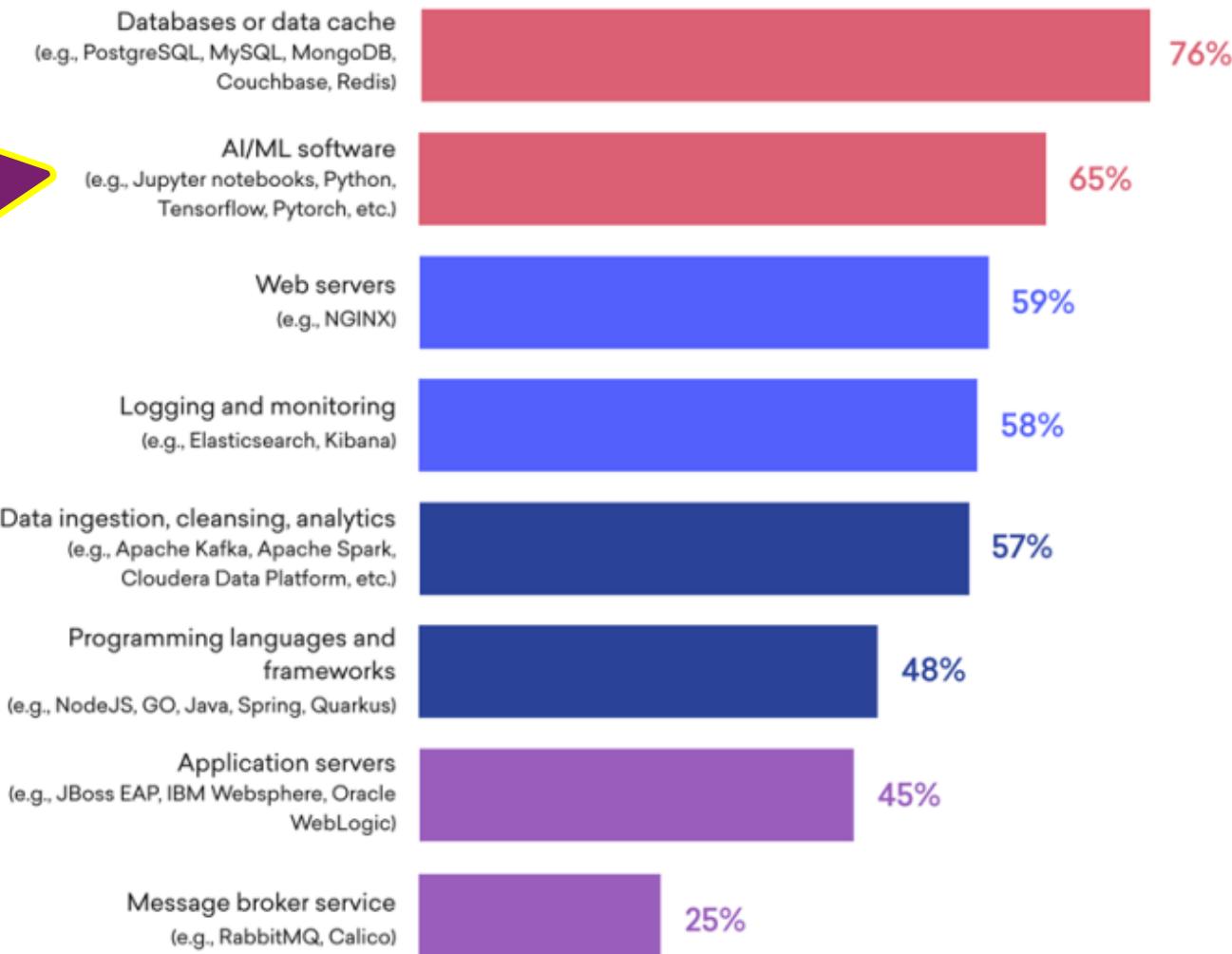
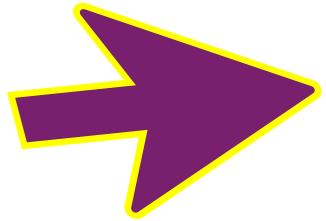


Why organizations deploy workloads on containers (Source: [Red Hat](#))

# Why Deploy AI on Kubernetes?



China 2024

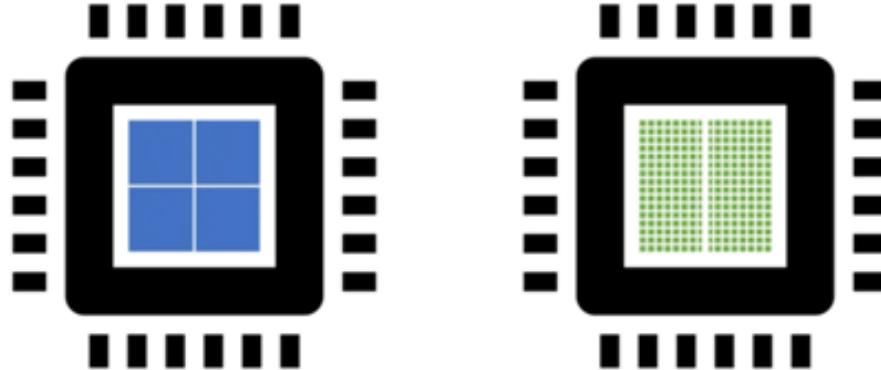


Major types of workloads and systems containerized with Kubernetes (Source: [Red Hat](#))

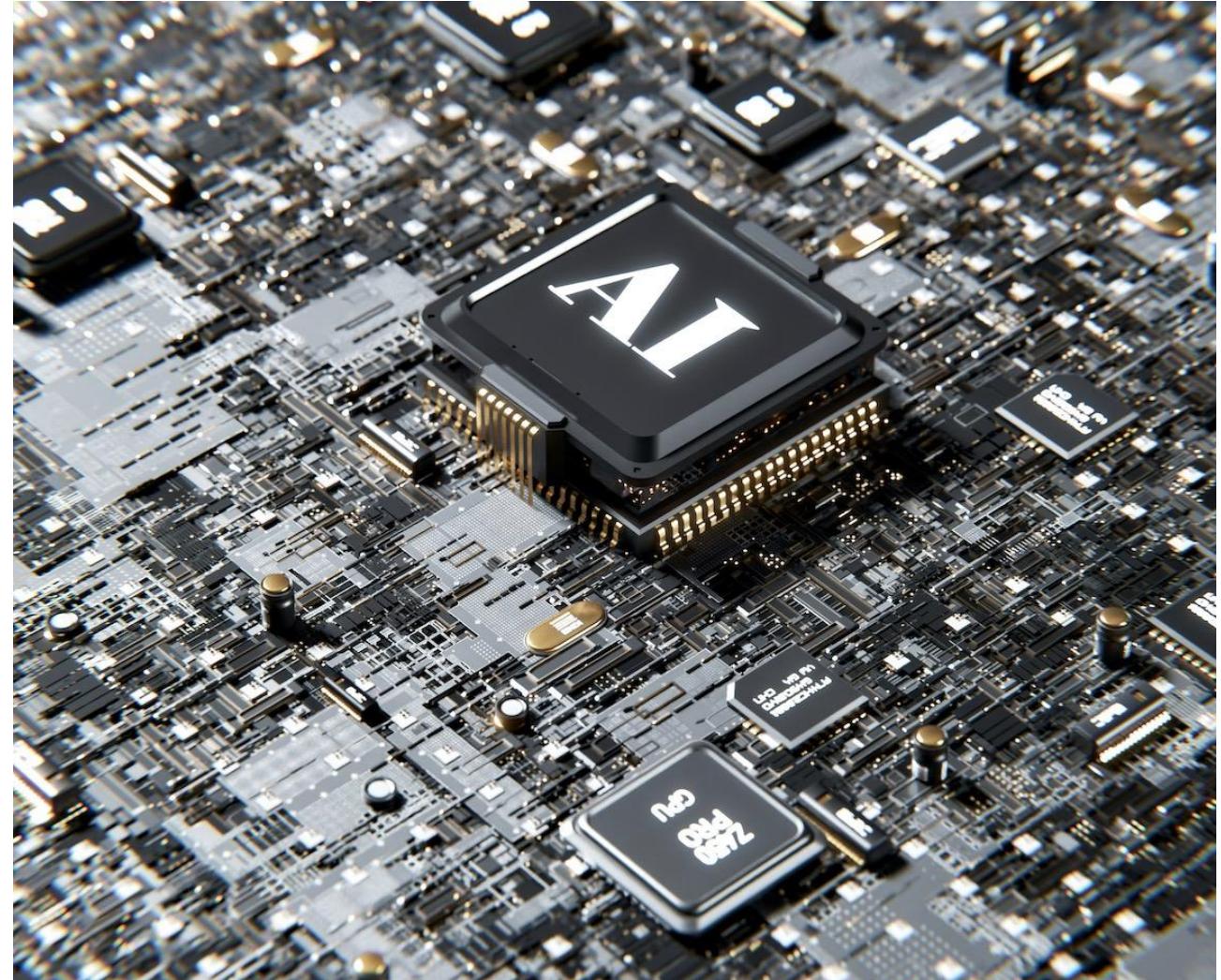
# Using GPUs for AI Workloads



China 2024



CPU	GPU
Central Processing Unit	Graphics Processing Unit
4-8 Cores	100s or 1000s of Cores
Low Latency	High Throughput
Good for Serial Processing	Good for Parallel Processing
Quickly Process Tasks That Require Interactivity	Breaks Jobs Into Separate Tasks To Process Simultaneously
Traditional Programming Are Written For CPU Sequential Execution	Requires Additional Software To Convert CPU Functions to GPU Functions for Parallel Execution



<GPU vs CPU processing. Source: towardsdatascience.com>

# Using GPUs on Kubernetes Nodes



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

- **Dynamic Resource Allocation**

Kubernetes can dynamically allocate GPU resources based on workload requirements, ensuring optimal utilization and minimizing idle resources

- **Cost Efficiency**

By efficiently managing GPU resources, organizations can control costs, especially in cloud environments where GPU usage can be expensive

- **Multi-cloud Strategy**

Organizations can leverage GPUs across multiple cloud providers without being locked into a single vendor, enhancing flexibility and reducing dependency risks

# Using GPUs on Kubernetes Nodes



China 2024

- ## Device Plugins

Kubernetes supports GPU device plugins that can manage the lifecycle of GPU resources, ensuring efficient usage and management

- ## GPU Vendor Integration

Kubernetes has strong integration with the GPU technologies like NVIDIA Kubernetes device plugin and CUDA, making it easier to leverage GPUs

- **Automatic GPU Discovery**  
The plugin automatically discovers GPUs on the Kubernetes nodes and makes them available to be scheduled by Kubernetes
- **Resource Allocation**  
It allows users to request specific GPU resources in their pod specifications, enabling efficient utilization of GPU resources
- **Health Checks**  
The plugin includes mechanisms to monitor the health of the GPUs and ensure that only healthy GPUs are allocated to workloads

# When Kubernetes is Not Ideal for AI



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



China 2024

- Small-Scale Projects or Limited Resources
- High Initial Setup and Learning Curve
- Real-Time or Low-Latency Applications
- Specialized Hardware Requirements
- Development and Experimentation Phase
- Security and Compliance
- Cost Considerations

# GPU Allocation is Important



China 2024

<p> <b>Meta-Llama-3-70B-Instruct</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • meta-llama</p> <p>A 70-billion parameter model from Meta, optimized for dialogue. Generates helpful, safe responses and outperforms other open-source chat LLMs.</p> <p>GPU 4x Nvidia A100 ▾ \$ 16 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>	<p> <b>gemma-7b-bit</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • google</p> <p>A instruction model fine-tuned from Gemma 7B Googles, first open LLM.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>	<p> <b>OpenHermes-2.5-Mistral-7B-G...</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • TheBloke</p> <p>A powerful chat model fine-tuned from Mistral 7B on a large corpus of synthetic data. Capable of function-calling and has strong coding capabilities.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>
<p> <b>NeuralHermes-2.5-Mistral-7B...</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • TheBloke</p> <p>A fine-tuned version of OpenHermes 2.5 that was aligned with Direct Preference Optimization and AI preference examples from the SlimOrca dataset.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>	<p> <b>Starling-LM-7B-alpha</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • berkeley-nest</p> <p>Open chat language model by UC Berkeley. Outperformed all 7B models at the time of its release. For non-commercial use only.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>	<p> <b>openchat-3.5-0106</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • openchat</p> <p>Open source large language model with performance and commercial potential using C-RLFT, for results on the GLUE benchmark.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>
<p> <b>Mistral-7B-Instruct-v0.1</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • mistralai</p> <p>A 7-billion parameter instruct model from Mistral AI, fine-tuned using a variety of publicly available conversation datasets.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>	<p> <b>zephyr-7b-beta</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • huggingFaceH4</p> <p>A chat model fine-tuned from Mistral 7B with synthetic data and Direct Preference Optimization.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>	<p> <b>neural-chat-7b-v3-1</b> <span style="border: 1px solid #ccc; padding: 2px;">TGI</span></p> <p>Text Generation • Intel</p> <p>A chat model from Intel fine-tuned from Mistral 7B on the SlimOrca dataset with Direct Preference Optimization.</p> <p>GPU 1x Nvidia L4 ▾ \$ 0.8 / h <span style="border: 1px solid #ccc; border-radius: 50%; padding: 2px 5px;">Go</span></p>

Inference endpoints - Text Generation Models (Source: [huggingface](#))

- ## AI Training

This process is typically more resource-intensive and can take advantage of more powerful or multiple GPUs

e.g.) FP64, FP32, FP16 supported,  
NVIDIA A100, NVIDIA V100, NVIDIA H100

- ## AI Inference

This process is usually less resource-intensive compared to AI training and may only require a single GPU or a fraction of it

e.g.) FP16, INT8 supported,  
NVIDIA T4, NVIDIA A10, NVIDIA Xavier

# GPU Allocation Scenario



China 2024

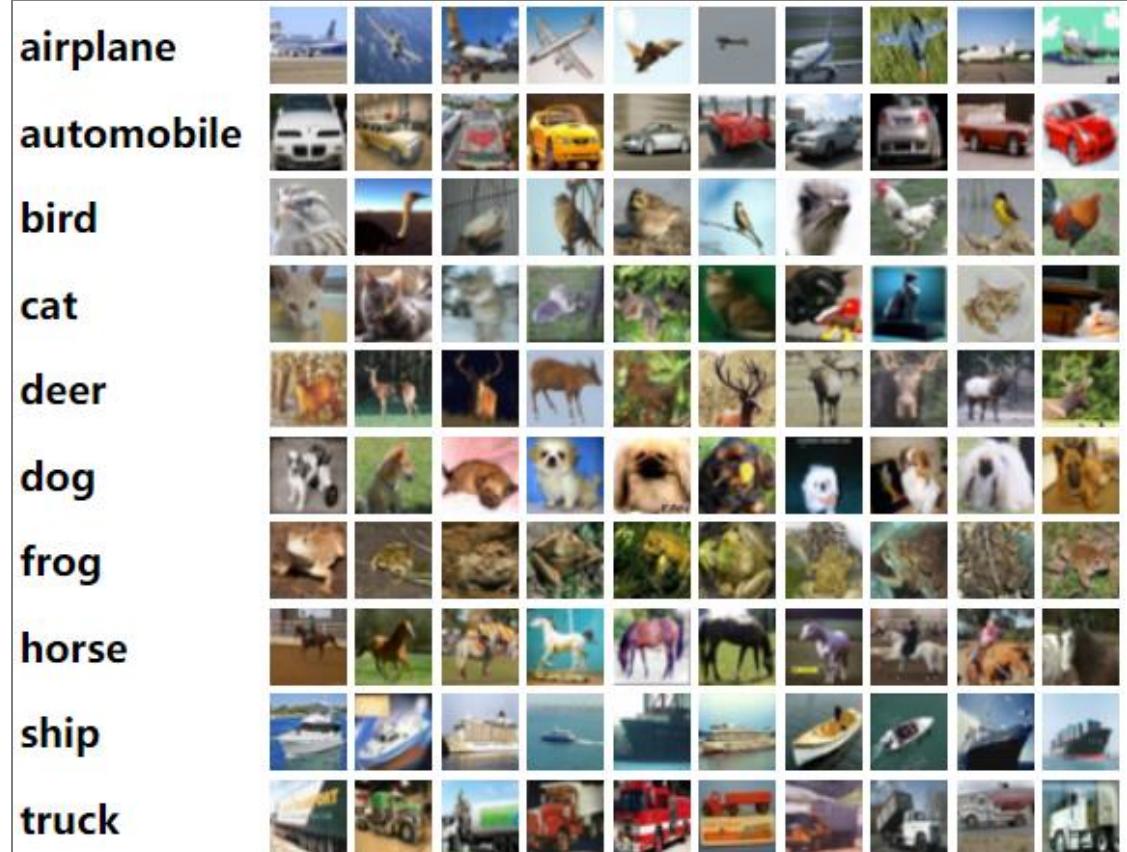
```
42 # Model training function
43 def train_model(epochs):
44     for epoch in range(epochs):
45         train_loss = 0.0
46         for data, label in train_data:
47             data = data.as_in_context(ctx)
48             label = label.as_in_context(ctx)
49             with autograd.record():
50                 output = model(data)
51                 loss = loss_fn(output, label)
52                 loss.backward()
53                 trainer.step(batch_size=128)
54                 train_loss += loss.mean().asscalar()
55                 print(f"Epoch {epoch+1}, Loss: {train_loss/len(train_data):.3f}")
56
57 # Train the model
58 train_model(epochs=10)
59
60 # Save model parameters
61 model.save_parameters('cifar_resnet20_v1.params')
```

## Training

# VS.

## Inference

```
13 # Load the pre-trained model
14 model = get_model('cifar_resnet20_v1', classes=10, pretrained=True)
15
16 # CIFAR-10 class names list
17 class_names = ['airplane', 'automobile', 'bird', 'cat', 'deer',
18                 'dog', 'frog', 'horse', 'ship', 'truck']
```



The CIFAR-10 dataset

- Labeled subsets of the 80 million tiny images dataset
- Consists of 60,000 32x32 colour images in 10 classes
- 60,000 = 6,000 images per class (x10)
- 50,000 training images and 10,000 test images

# GPU Allocation Scenario



KubeCon



CloudNativeCon

THE LINUX FOUNDATION  
OPEN SOURCE SUMMIT

China 2024

```
ssh
root@localhost:~/ai# python3 ai_training.py
Epoch 1, Loss: 1.494
Epoch 2, Loss: 1.025
Epoch 3, Loss: 0.824
Epoch 4, Loss: 0.712
Epoch 5, Loss: 0.644
Epoch 6, Loss: 0.592
Epoch 7, Loss: 0.551
Epoch 8, Loss: 0.527
Epoch 9, Loss: 0.502
Epoch 10, Loss: 0.477
 * Serving Flask app 'ai_training'
 * Debug mode: off
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
 * Running on all addresses (0.0.0.0)
 * Running on http://127.0.0.1:5000
 * Running on http://172.234.95.27:5000
Press CTRL+C to quit
127.0.0.1 - - [13/Aug/2024 08:38:34] "POST /predict HTTP/1.1" 200 -
```

```
ssh
root@localhost:~/ai# curl -X POST http://127.0.0.1:5000/predict -F "img=@animal.jpg"
{"prediction":"horse","probability":0.7999327182769775} → Only Inferencing
root@localhost:~/ai# curl -X POST http://127.0.0.1:5000/predict -F "img=@animal.jpg"
{"prediction":"horse","probability":0.9025683403015137} → After training(epoch 10, Loss 0.477)
root@localhost:~/ai#
```



- **GPU Operator**  
NVIDIA GPU Operator in Kubernetes to manage GPU resources effectively. This can help in dynamic allocation and better resource utilization
- **Automated Node Labeling**  
Automatically label nodes with GPU availability, allowing your scheduler to allocate training or inference workloads to the appropriate nodes

- Mixed Workloads on Same GPU

For NVIDIA A100 GPUs, you can use MIG (Multi-Instance GPU) to partition a single GPU into multiple instances. This allows you to run training and inference workloads simultaneously on different GPU instances

- Dynamic Allocation Based on Load

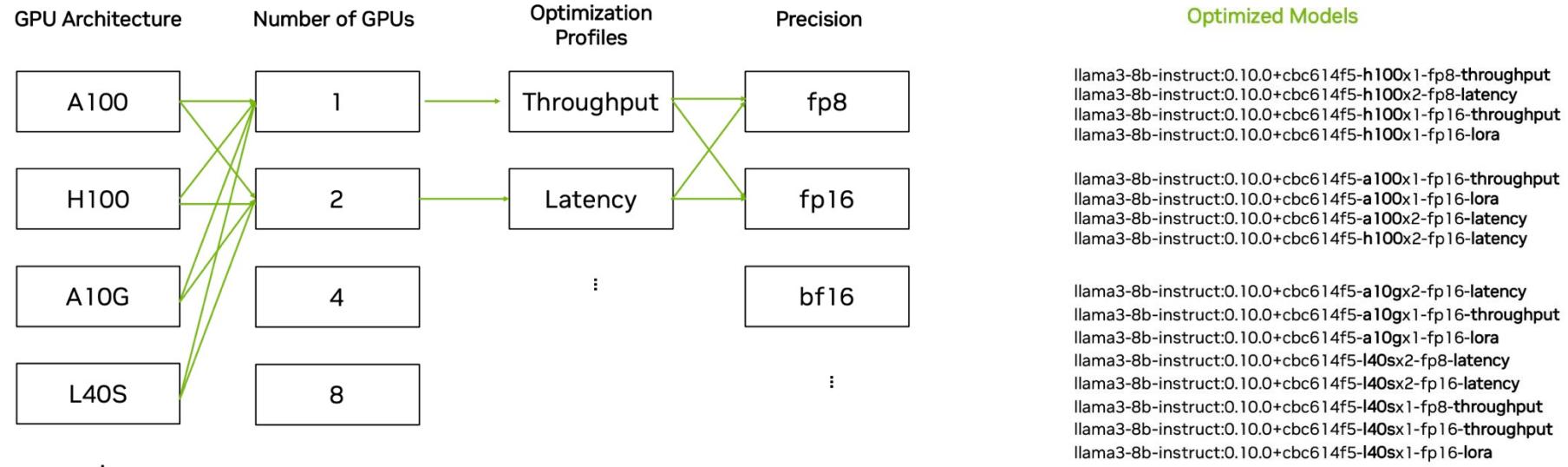
Implement a system that dynamically allocates GPUs based on the current load. For instance, during peak training periods, more GPUs can be allocated for training, while during high training or inference demand, GPUs can be reallocated accordingly

# Example) GPU Allocation

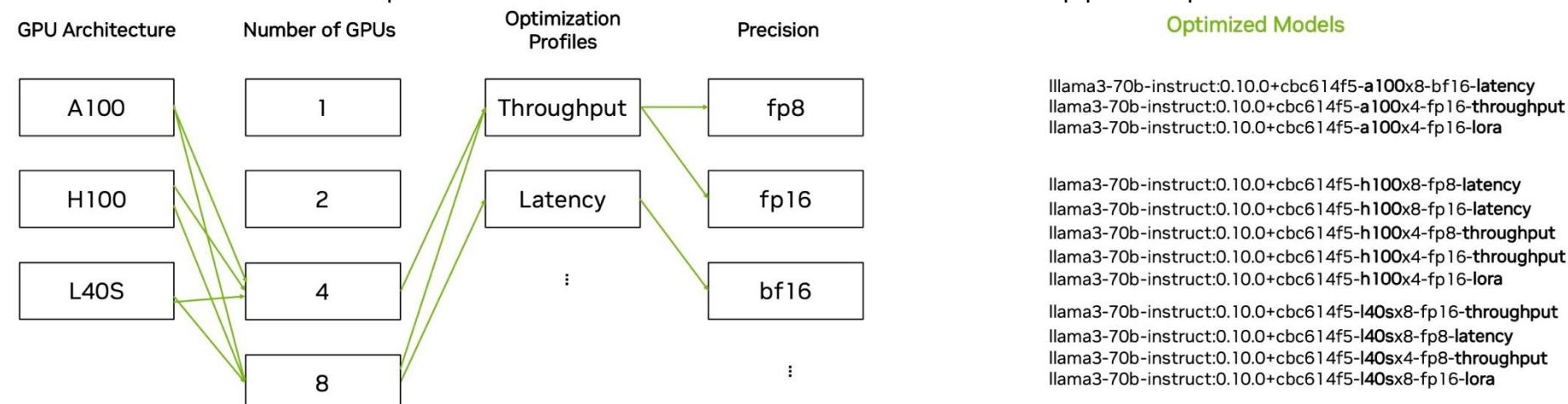


China 2024

## Llama-3-8B-Instruct Optimization for NVIDIA GPUs with LoRA Support Option



## Llama-3-70B-Instruct Optimization for NVIDIA GPUs with LoRA Support Option

\* Sourced from [NVIDIA NIM page](#)

# VRAM(Video RAM)



<A GPU die surrounded by VRAM chips, from [Wikipedia](#)>

- Type of memory used specifically by the GPU
- Acts as the dedicated memory for the GPU
- Optimized for high bandwidth
- Allowing it to quickly supply the GPU with the data it needs
- GDDR (Graphics Double Data Rate)
- HBM (High Bandwidth Memory)
- Large Datasets / Complex Data Augmentation
- Real-time Inference / Multiple Models
- Parallel Processing / GPU Partitioning
- Fragmentation

Example)

Specifications	
<b>GPU Memory</b>	20GB GDDR6
<b>Memory Interface</b>	160 bit
<b>Memory Bandwidth</b>	360GB/s
<b>Error Correcting Code (ECC)</b>	Yes
<b>NVIDIA Ada Lovelace Architecture-Based CUDA Cores</b>	6,144
<b>NVIDIA Fourth-Generation Tensor Cores</b>	192
<b>NVIDIA Third-Generation RT Cores</b>	48
<b>Single-Precision Performance</b>	26.7 TFLOPS <sup>3</sup>
<b>RT Core Performance</b>	61.8 TFLOPS <sup>3</sup>
<b>Tensor Performance</b>	327.6 TFLOPS <sup>4</sup>
<b>System Interface</b>	PCIe 4.0 x16
<b>Power Consumption</b>	Total board power: 130W
<b>Thermal Solution</b>	Active
<b>Form Factor</b>	4.4" H x 9.5" L, single slot
<b>Display Connectors</b>	4x DisplayPort 1.4a <sup>5</sup>

# Schedule Resources with GPU



KubeCon



CloudNativeCon

OPEN  
SOURCE  
SUMMIT

China 2024



```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-pod
spec:
  restartPolicy: Never
  containers:
    - name: cuda-container
      image: nvcr.io/nvidia/k8s/cuda-sample:vectoradd-cuda10.2
      resources:
        limits:
          nvidia.com/gpu: 2 # requesting 2 GPUs
  tolerations:
    - key: nvidia.com/gpu
      operator: Exists
      effect: NoSchedule
```

- A container using the NVIDIA CUDA sample image
- Requests 2 GPUs for its container
- tolerations for nodes with GPUs

# Schedule Resources with GPU



KubeCon



CloudNativeCon



THE LINUX FOUNDATION

OPEN SOURCE  
SUMMIT

China 2024



```
apiVersion: batch/v1
kind: Job
metadata:
  name: tensorflow-gpu-job
spec:
  template:
    spec:
      containers:
        - name: tensorflow-gpu-container
          image: tensorflow/tensorflow:latest-gpu
          command: ["python", "-c"]
          args: ["import tensorflow as tf; print(tf.config.list_physical_devices('GPU'))"]
      resources:
        limits:
          nvidia.com/gpu: 1 # Request 1 GPU
  restartPolicy: OnFailure
```

# Schedule Resources with GPU



KubeCon



CloudNativeCon



China 2024



- You can specify GPU limits without specifying requests, because Kubernetes will use the limit as the request value by default
- You can specify GPU in both limits and requests, but these two values must be equal
- You cannot specify GPU requests without specifying limits

# Various types of GPUs



```
# Label your nodes with the accelerator type they have
$ kubectl label nodes node1 accelerator=example-gpu-x100
$ kubectl label nodes node2 accelerator=other-gpu-k915
```

“**If** each node in your cluster have different types of GPUs, then you can use Node Labels and Node Selectors to schedule pods to appropriate nodes.”



KubeCon



CloudNativeCon



THE LINUX FOUNDATION  
OPEN SOURCE SUMMIT



China 2024

# Automated Node Labeling



China 2024

- Automatic discover and label GPU enabled nodes
- Deploying Kubernetes Node Feature Discovery (NFD)
- NFD detects the hardware features on each node
- NFD is configured to advertise those features as node labels
- NFD can add extended resources, annotations, and taints
- NFD creates the feature labels for the detected features
- Administrators can leverage NFD to also taint nodes

The screenshot shows a web browser displaying the official documentation for Node Feature Discovery. The URL in the address bar is `kubernetes-sigs.github.io/node-feature-discovery/stable/get-started/index.html`. The page has a blue header with the title "Node Feature Discovery v0.16.4". On the left, there's a sidebar with a "GET STARTED" section containing links for "1. Introduction" and "2. Quick start". The main content area features a large heading "Node Feature Discovery" and a welcome message: "Welcome to Node Feature Discovery - a Kubernetes add-on for detecting hardware features and system configuration!". Below this, there's a "Continue to:" section with two bullet points: "Introduction for more details on the project." and "Quick start for quick step-by-step instructions on how to get NFD running on your cluster."

# Automated Node Labeling



China 2024



```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-pod
spec:
  containers:
  - name: cuda-container
    image: nvidia/cuda:11.0-base
    resources:
      limits:
        nvidia.com/gpu: 1
  nodeSelector:
    feature.node.kubernetes.io/pci-10de.present: "true"
```

- Pre-defined NFD
- nodeSelector

# Automated Node Labeling



China 2024

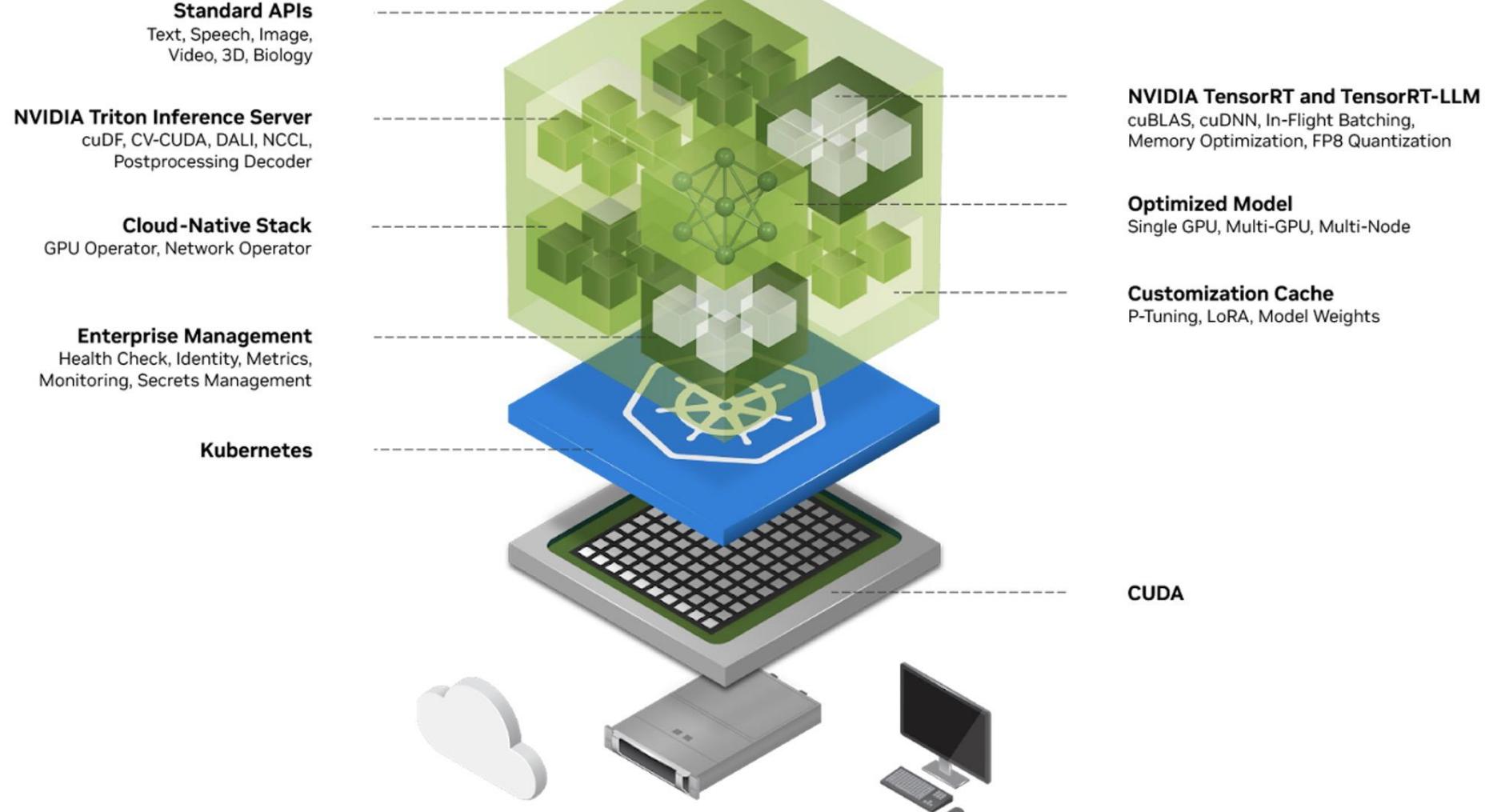


```
apiVersion: v1
kind: Pod
metadata:
  name: example-vector-add
spec:
  restartPolicy: OnFailure
affinity:
  nodeAffinity:
    requiredDuringSchedulingIgnoredDuringExecution:
      nodeSelectorTerms:
        - matchExpressions:
            - key: "gpu.gpu-vendor.example/ins-11-memory"
              operator: Gt # (greater than)
              values: ["40535"]
            - key: "feature.node.kubernetes.io/pci-10.present" # NFD Feature label
              values: ["true"] # (optional) only schedule on nodes with PCI device 10
containers:
  - name: example-vector-add
    image: "registry.example/example-vector-add:v42"
    resources:
      limits:
        gpu-vendor.example/example-gpu: 1 # requesting 1 GPU
```

- Pre-defined NFD
- NodeAffinity to schedule this Pod onto a node



## Inference Microservices for GenAI

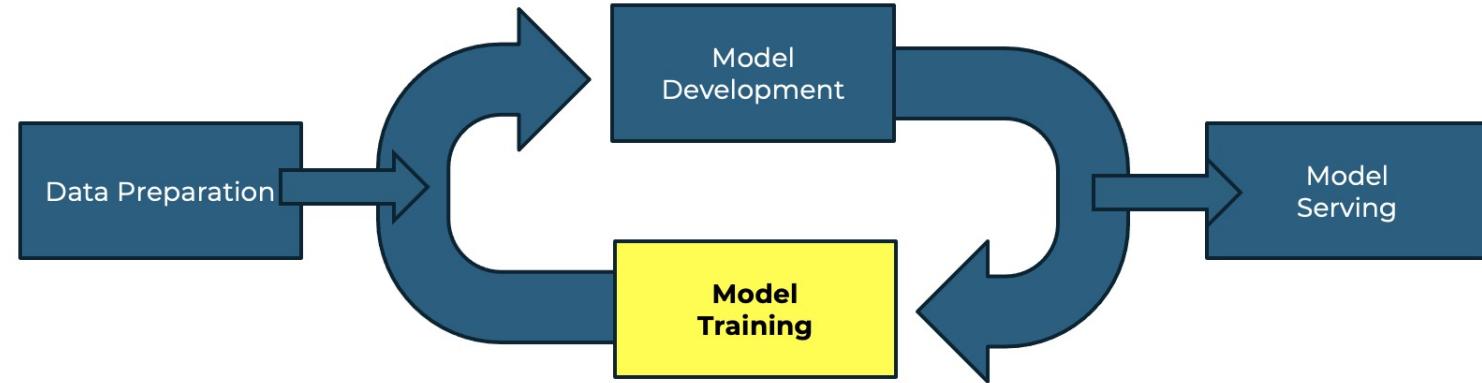


# Industry Use Cases

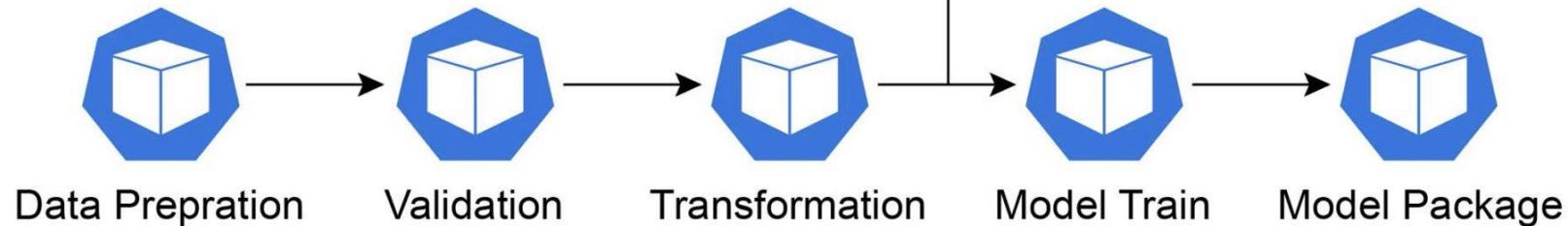


China 2024

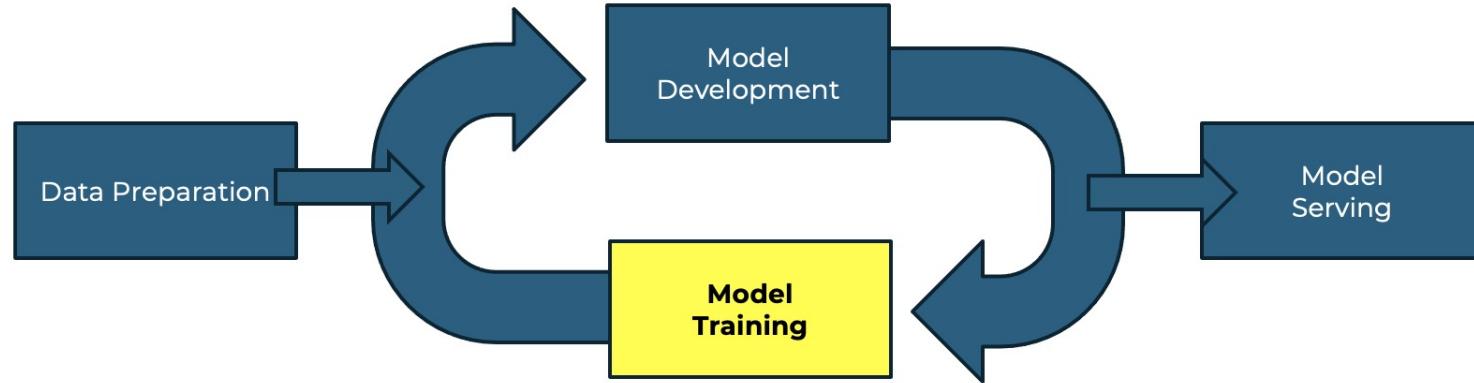
## Model Training



Kubeflow



## Model Training



TensorFlow

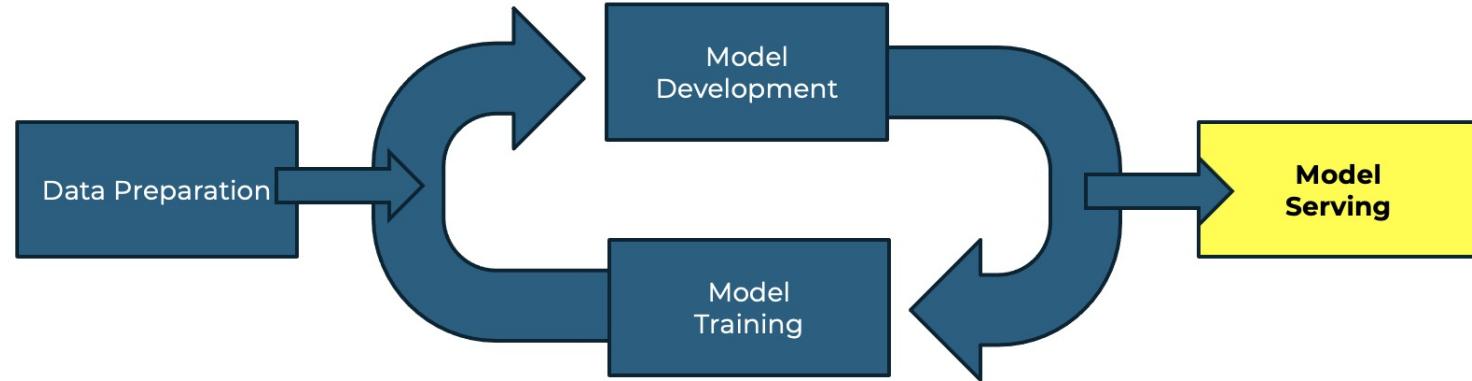


Keras



PyTorch

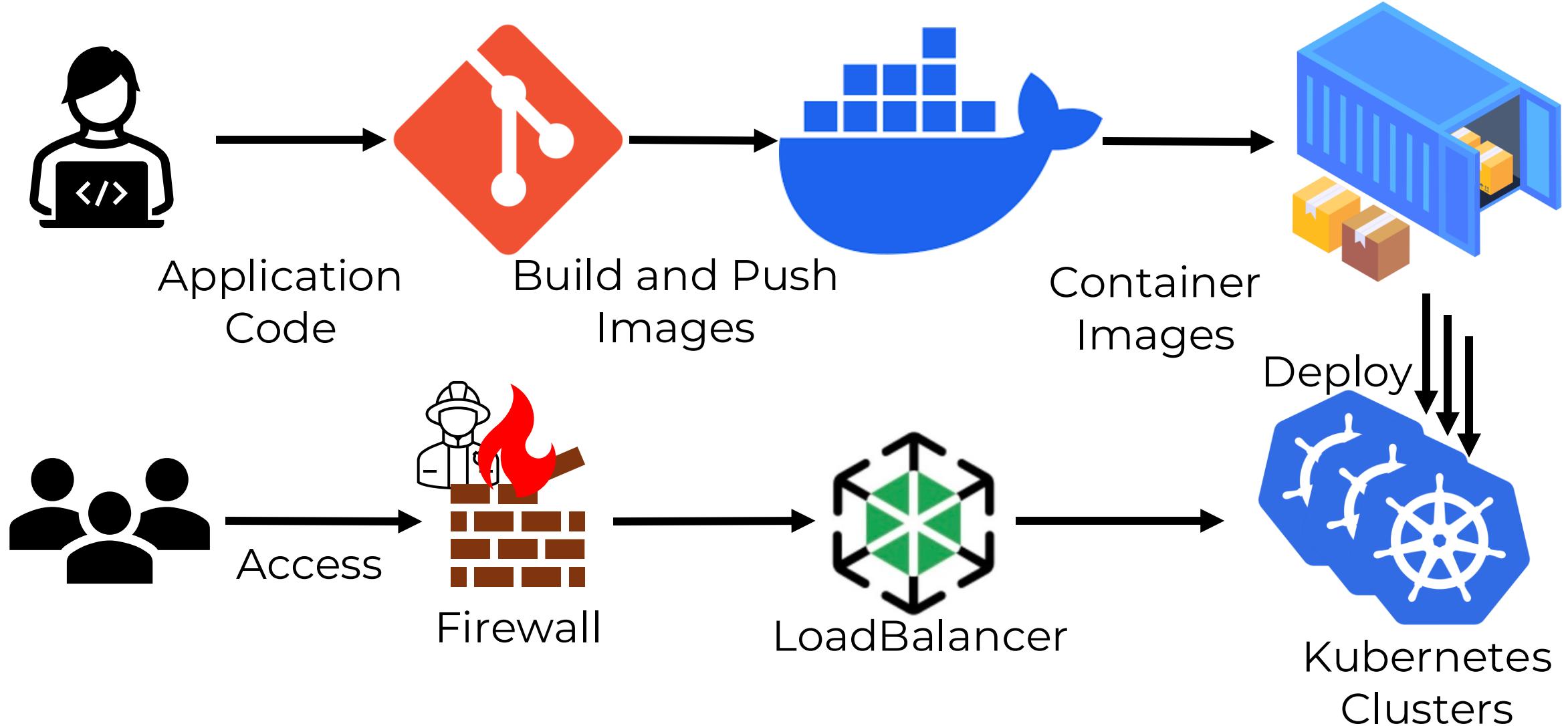
## Model Serving



# Deploying AI/ML on Kubernetes



China 2024



# Deploying AI/ML on Kubernetes



China 2024



```
$ minikube start --driver docker --container-runtime docker --gpus all
😊 minikube v1.33.1 on Debian 11.10 (amd64)
✨ Using the docker driver based on user configuration
🚀 Using Docker driver with root privileges
👍 Starting "minikube" primary control-plane node in "minikube" cluster
🚜 Pulling base image v0.0.44 ...
🔥 Creating docker container (CPUs=2, Memory=8000MB) ...
🌐 Preparing Kubernetes v1.30.0 on Docker 26.1.1 ...
    ■ Generating certificates and keys ...
    ■ Booting up control plane ...
    ■ Configuring RBAC rules ...
🔗 Configuring bridge CNI (Container Networking Interface) ...
🔍 Verifying Kubernetes components...
    ■ Using image nvcr.io/nvidia/k8s-device-plugin:v0.15.0
    ■ Using image gcr.io/k8s-minikube/storage-provisioner:v5
🌟 Enabled addons: storage-provisioner, nvidia-device-plugin, default-storageclass
🎉 Done! kubectl is now configured to use "minikube" cluster and "default" namespace by default
```

```
$ kubectl get nodes -o yaml | grep -i gpu
  nvidia.com/gpu: "1"
  nvidia.com/gpu: "1"
```

# Deploying AI/ML on Kubernetes



China 2024



```
import torch
import torch.nn as nn
import torch.optim as optim
import torchvision
import torchvision.transforms as transforms

# Check if GPU is available
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"Using device: {device}")

# Load and preprocess the CIFAR-10 dataset
transform = transforms.Compose(
    [transforms.ToTensor(),
     transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])

trainset = torchvision.datasets.CIFAR10(root='./data', train=True, download=True, transform=transform)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=4, shuffle=True, num_workers=2)

testset = torchvision.datasets.CIFAR10(root='./data', train=False, download=True, transform=transform)
testloader = torch.utils.data.DataLoader(testset, batch_size=4, shuffle=False, num_workers=2)

classes = ('plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck')
....
....(skipped)
```

# Deploying AI/ML on Kubernetes



KubeCon



CloudNativeCon

THE LINUX FOUNDATION  
OPEN SOURCE SUMMIT

China 2024



sakang@localhost: ~

```
sakang@localhost:~$ python3 pytorch_example.py
Using device: cuda
Files already downloaded and verified
Files already downloaded and verified
[1, 2000] loss: 2.157
[1, 4000] loss: 1.895
[1, 6000] loss: 1.732
[1, 8000] loss: 1.611
[1, 10000] loss: 1.547
[1, 12000] loss: 1.505
[2, 2000] loss: 1.440
[2, 4000] loss: 1.414
```

```
root@localhost:~# nvidia-smi
Sun Aug 18 18:51:06 2024
+
| NVIDIA-SMI 550.107.02           Driver Version: 550.107.02     CUDA Version: 12.4 |
+-----+-----+-----+-----+-----+-----+-----+-----+
| GPU  Name     Persistence-M | Bus-Id     Disp.A  | Volatile Uncorr. ECC | | | | | |
| Fan  Temp     Perf            Pwr:Usage/Cap | Memory-Usage | GPU-Util  Compute M. |
|          |          |          |          |          |          |          |          |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 0  Quadro RTX 6000          Off | 00000000:00:02.0 Off |                  Off |
| 33% 34C   P2                61W / 260W | 406MiB / 24576MiB | 10%       Default N/A |
+-----+-----+-----+-----+-----+-----+-----+-----+
+
Processes:
GPU  GI  CI      PID  Type  Process name          GPU Memory
ID   ID
+
| 0  N/A N/A    12186  C   python          170MiB
| 0  N/A N/A    19274  C   python3        232MiB
+
```



# Deploying AI/ML on Kubernetes



China 2024



```
# Use the latest PyTorch image as the base image
FROM pytorch/pytorch:latest

# Set the working directory
WORKDIR /app

# Copy necessary files into the Docker image
COPY requirements.txt /app/requirements.txt
COPY pytorch_example.py /app/pytorch_example.py

# Install required packages
RUN pip install --no-cache-dir -r requirements.txt

# Execute the script
CMD [ "python", "pytorch_example.py" ]
```

e.g.) requirements.txt

```
torch==2.0.0
torchvision==0.15.0
numpy==1.24.0
pandas==2.0.3
matplotlib==3.7.1
scikit-learn==1.2.2
```

# Deploying AI/ML on Kubernetes



China 2024



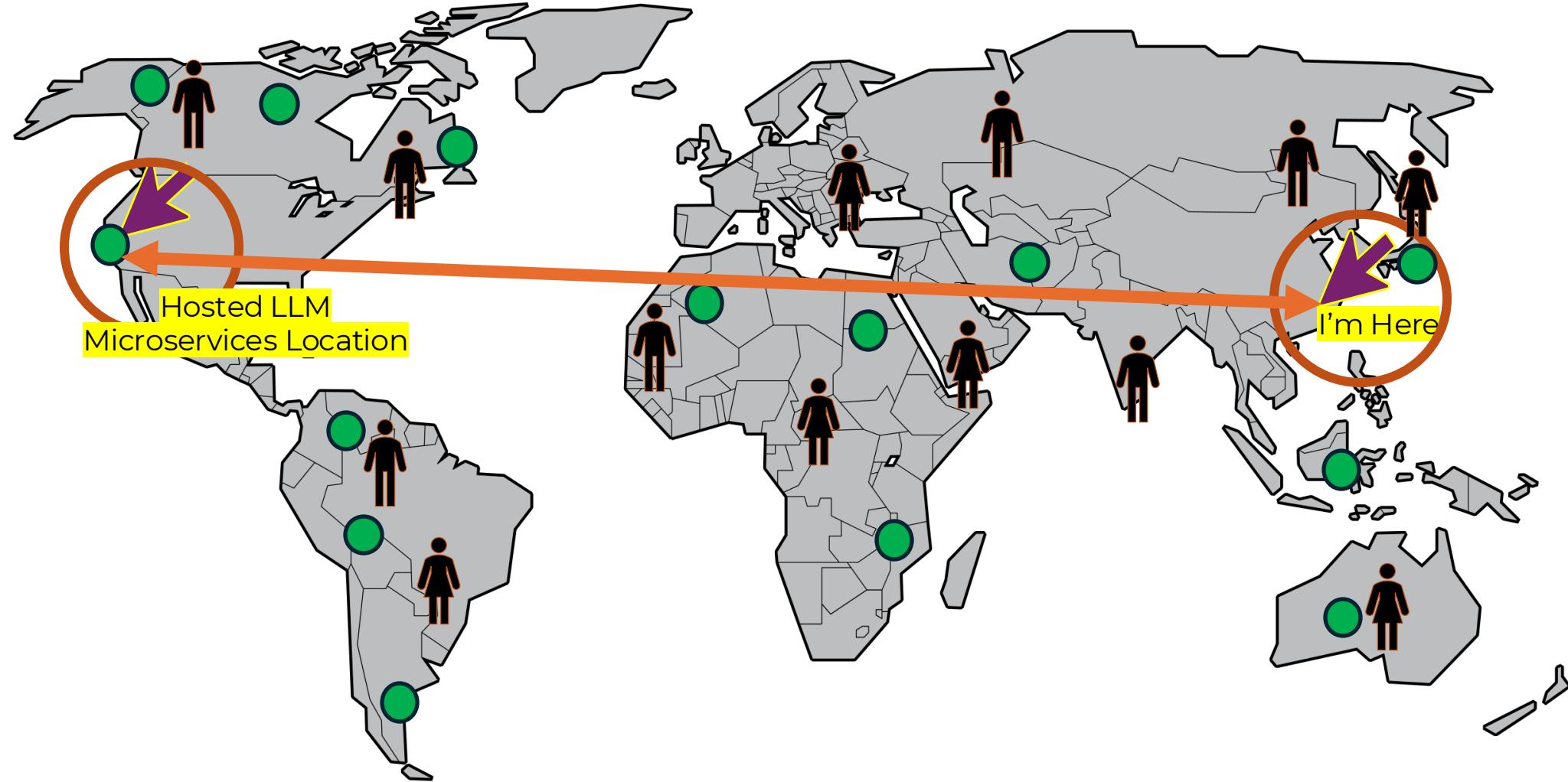
```
#GPU_Pod.yaml
apiVersion: v1
kind: Pod
metadata:
  name: pytorch-example
spec:
  restartPolicy: OnFailure
  containers:
    - name: pytorch-example
      image: localhost:5000/pytorch_example:latest
      imagePullPolicy: IfNotPresent
      resources:
        limits:
          nvidia.com/gpu: 1
  imagePullSecrets:
    - name: registrypullsecret
```

- Sets up a Docker container
- The latest PyTorch image
- Copies the necessary files into the container
- Finally runs the Python script

# Demand for Distributed Cloud



China 2024



“

The cloud's next phase requires a shift in  
how developers and enterprises think  
about getting applications and data  
closer to their customers.

*Tom Leighton  
CEO of Akamai Technologies*

# Distributed Cloud for AI on K8s



KubeCon



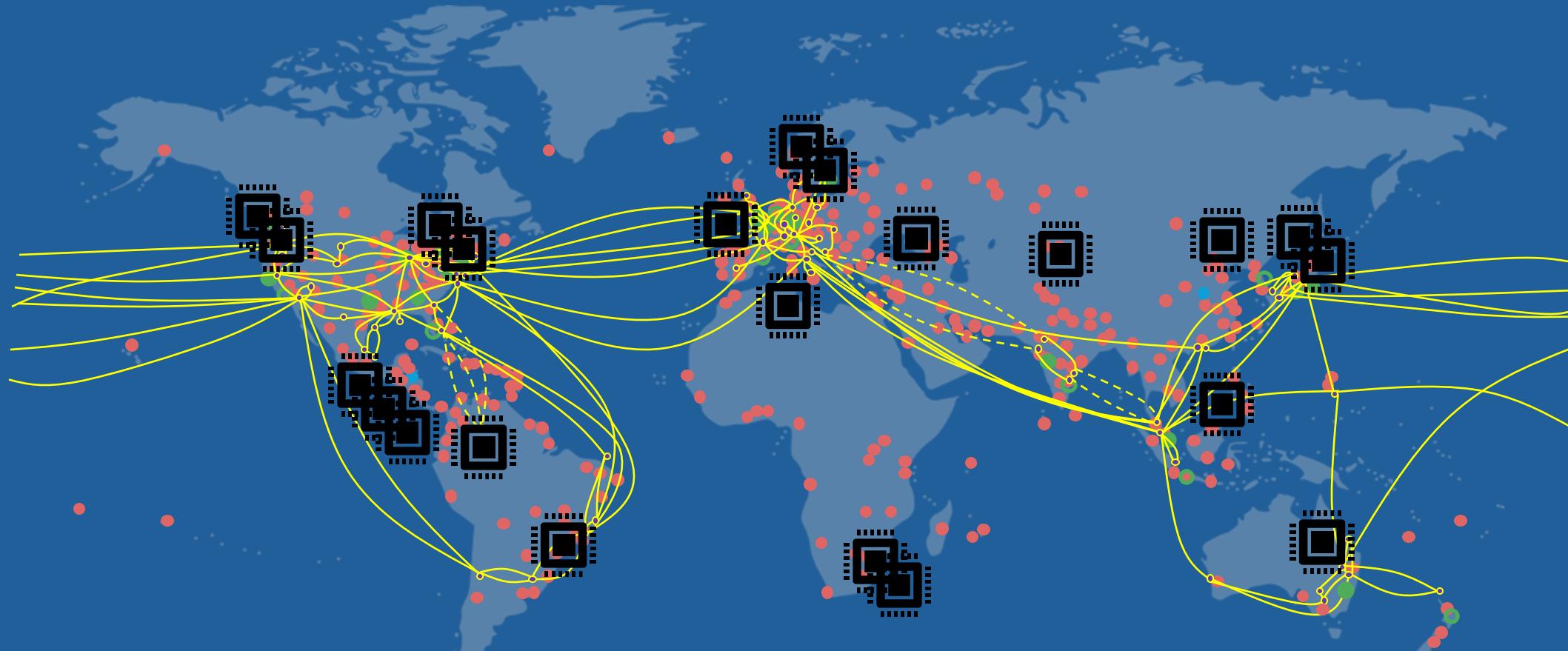
CloudNativeCon



OPEN SOURCE  
SUMMIT



China 2024



● Core Computing Regions

● Distributed Regions / Edge Computing Sites / CDN PoPs

# Source Codes & Scripts in the slide



KubeCon



CloudNativeCon



OPEN  
SOURCE  
SUMMIT



AI\_dev

China 2024

- NFD(Node Feature Discovery) Script  
<https://github.com/BrandonKang/k8s-Node-Feature-Discovery>
- AI Training vs. AI Inference  
[https://github.com/BrandonKang/ai\\_training\\_vs\\_inference](https://github.com/BrandonKang/ai_training_vs_inference)
- Minikube & PyTorch ML case  
<https://github.com/BrandonKang/Pytorch-on-Minikube>



KubeCon



CloudNativeCon



China 2024

# Thank You!