



KubeCon



CloudNativeCon

India 2026

#KubeCon #CloudNativeCon

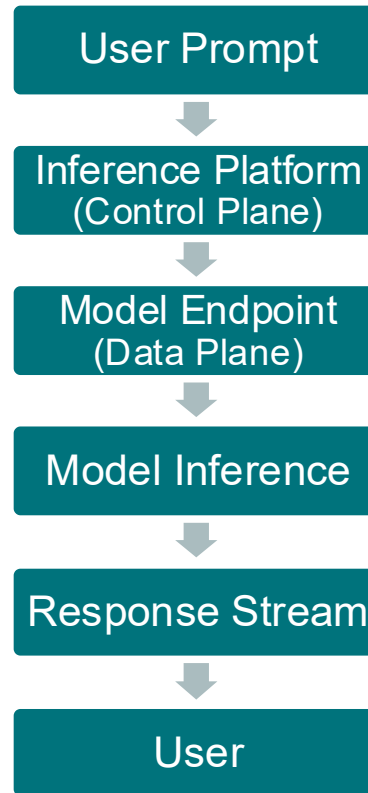
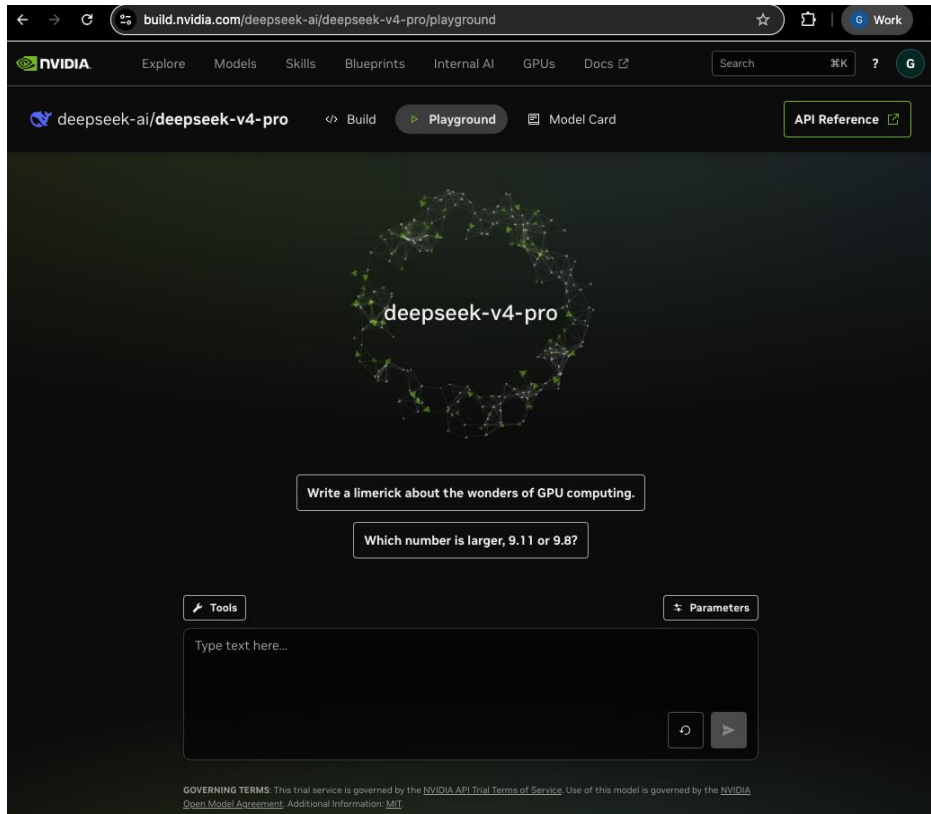
Inference in Progress... Please Monitor Responsibly

Observability Learnings for AI Inference Platforms at Scale

Gaurav Sharma
Nvidia



AI Inference Request Lifecycle



AI Inference Orchestration Challenges



Challenges

Multiple Consumers

Multiple Models

Multiple GPU Types

Multiple Clusters

Multiple Traffic Patterns

Orchestration Requirements

Global Endpoint + Auth

Model Abstraction

GPU Aware Placement

Multi-cluster Deployment

Routing + Autoscaling

AI Inference Platform



KubeCon



CloudNativeCon

India 2026



Population-Scale Inference

“... powers more than 500 government websites and processes over 15 million inferences daily, surpassing 6 billion total inferences leveraging NVIDIA’s NVCF, all while maintaining strict sub-second response times...”

[PIB, MeitY, Mar 2026](#)

Control Plane

- NVCF APIs (Global Endpoint + Auth)
- Functions Definitions (Model Abstraction)
- Function Deployment (GPU Aware Placement)
- Secrets Management
- Rate Limiting
- Autoscaling

Invocation Plane

- Routing inference requests to function workloads

Compute Plane

- NVCF Cluster Agent (GPU Cluster Discovery for Platform)
- Worker Orchestration on GPU Clusters to execute Inference

Inference Reliability Challenges



KubeCon



CloudNativeCon

India 2026

Control Plane



Routing, scaling,
orchestration failures

Invocation Plane



Latency, errors,
timeouts, throughput

Compute Plane



GPU, nodes, network,
storage failures

Monitoring Inference



KubeCon



CloudNativeCon

India 2026



Alert via SLO

Correlate via Drill Downs



Inference SLO



Control Plane SPOG



Compute Plane SPOG



Per Region Service Drill Down



Per GPU Cluster Drill Down



Node Monitoring Drill Down

Inference Request Journeys

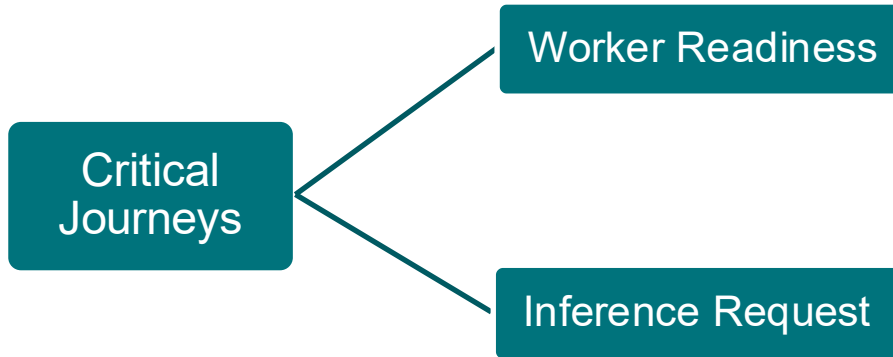


KubeCon



CloudNativeCon

India 2026



Metrics Instrumentation

Define SLIs

Set Targets & Budgets

Alert on Burn Rate

Worker Readiness



KubeCon



CloudNativeCon

India 2026

Kubernetes

Node Ready

Pod Health

Resource Pressure

Kubelet

GPU Software Stack

GPU Operator

Driver

Container Toolkit

Device Plugin

GPU Hardware

GPU Availability

XID Errors

ECC Errors

Temperature

Custom Monitoring

Compute Plane
Agent Checks

Model Downloads
Checks

Custom Network
Monitoring

GPU Provisioning
Workflows

Healthy GPU Node



KubeCon



CloudNativeCon

India 2026



Ready Pods

GPU Operator components

NVIDIA driver ready

Device plugin ready



Clean Signals

DCGM metrics

XID / ECC Errors

NVLink / thermal / power

Resource pressure normal



Validation Checks

Worker Pods Stable

Network probes passing

Test workload succeeds

Inference Server Logs

Inference Request Health



KubeCon

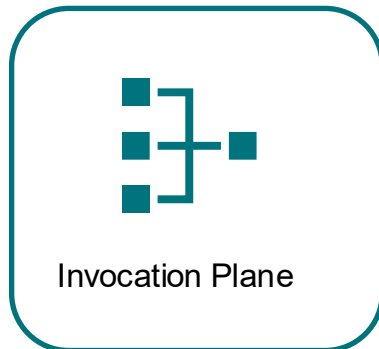


CloudNativeCon

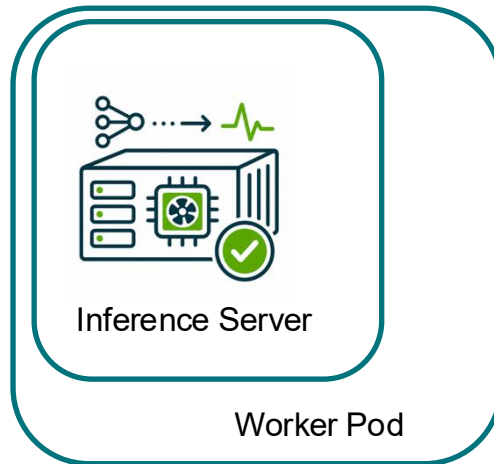
India 2026



Client



Invocation Plane



Inference Server

Worker Pod



TTFT Latency SLO

first token within target

tft_latency_histogram



Success Rate SLO

successful invocations / total invocations

invocation_requests_counter

invocation_errors_counter

NVIDIA OSS



KubeCon



CloudNativeCon

India 2026

**NVIDIA Cloud
Functions**

GPU Operator

DCGM-Exporter

Dynamo

NVSentinel

KAI Scheduler



KubeCon



CloudNativeCon

India 2026

Q&A



KubeCon



CloudNativeCon

India 2026

Thank You

