

# Optimizing AI Inference with GKE: Reducing Costs and Latency

Ananda Dwi Rahmawati  
Google Developer Expert Cloud



# Agenda



**Introduction to AI  
Inference on GKE**



**Strategies for  
Reducing Inference  
Latency on GKE**



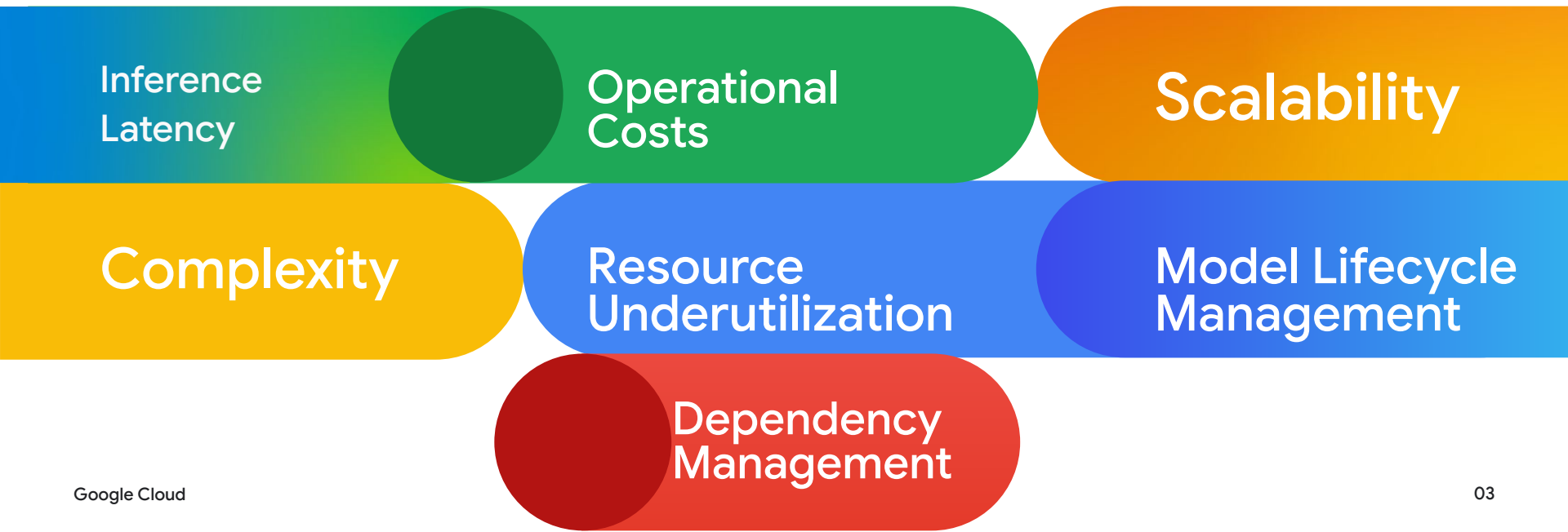
**Cost Optimization  
for AI Inference  
Workloads**



**Best Practices and  
What's Next to Do**

# Introduction - The Challenge of AI Inference

AI models are increasingly complex, demanding significant computational resources.



# Google Cloud AI Hypercomputer

## Flexible Consumption

Dynamic Workload Scheduler

On Demand

CUD

Spot

## Open Software

JAX, TensorFlow, PyTorch

Multislice Training, Multihost Inference, XLA

Google Kubernetes Engine & Compute Engine

## Performance-Optimized Hardware

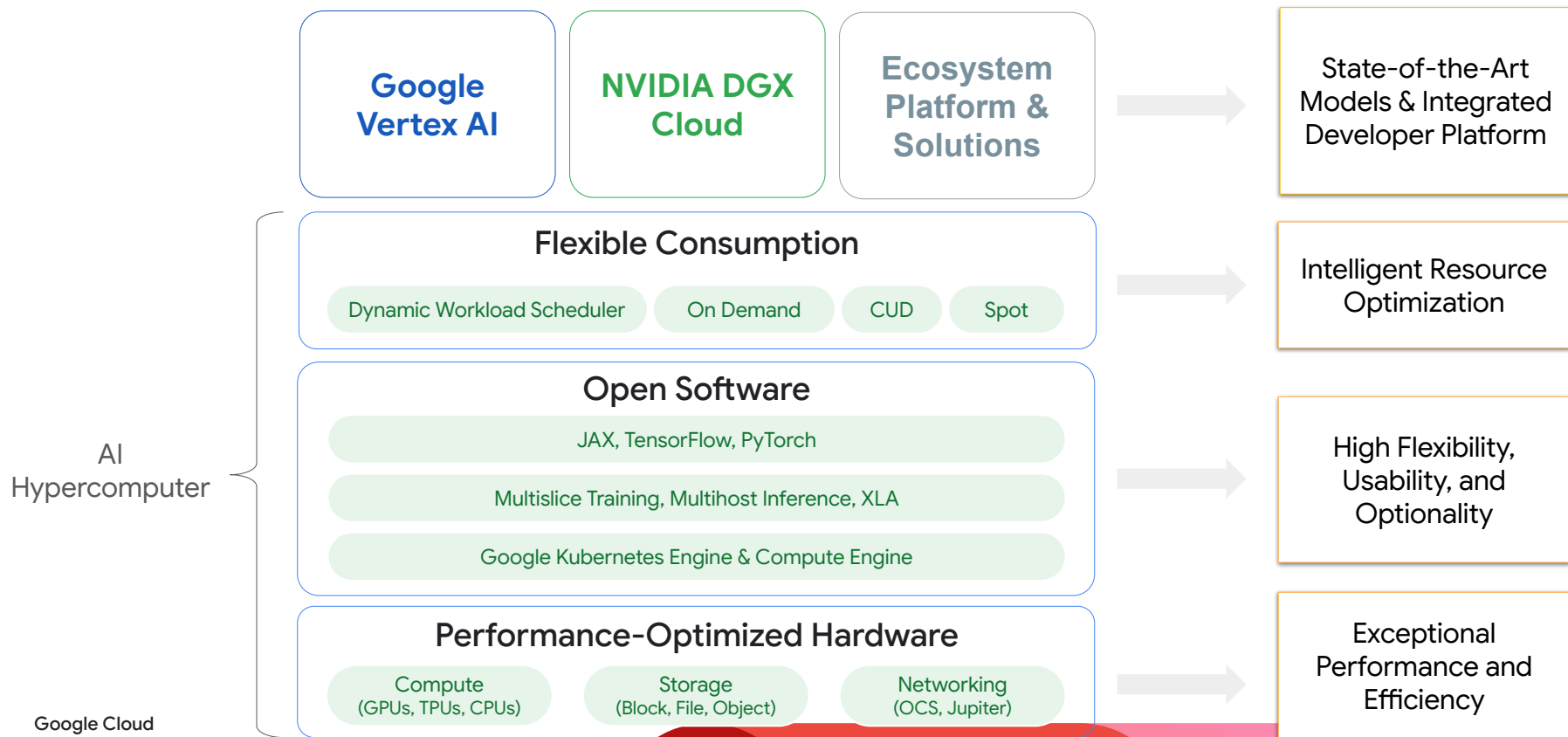
Compute  
(GPUs, TPUs, CPUs)

Storage  
(Block, File, Object)

Networking  
(OCS, Jupiter)

**Optimizing** system-level co-design streamlines the entire AI lifecycle—from training to tuning and serving

# AI Hypercomputer Architecture



# Why Google Kubernetes Engine for AI

1

## Portability & Customizability

Choice of frameworks and ecosystem tools that are portable

2

## Performance & Scalability

Scale the platform for supercomputer scale training and inference

3

## Cost-Efficiency

Increase utilization of valuable resources while reducing operational overhead

# Google Kubernetes Engine

## cloud native infrastructure for AI training and inference

- **Limitless Scale:** Deploy AI at industry-leading scale, supporting thousands of TPUs and nodes.
- **Cost-Efficient Performance:** Maximize price-performance with smart GPU/TPU use, job queuing, and fast provisioning.
- **Effortless Ops:** Focus on models, not infra, with GKE Autopilot's managed, optimized Kubernetes.
- **Enterprise Reliability:** Trust AI workloads to the leading Kubernetes contributor's cloud-native infra.



## Google Kubernetes Engine

### Open Software and Frameworks

JAX, TensorFlow, PyTorch, XLA

Jupyter, Ray, KubeFlow, Spark

### Distributed Training

Kueue Job Queuing

High Throughput Scaling

### Scaled Inference

Autopilot

Pod Fast Starts

### Node Provisioning and Autoscaling

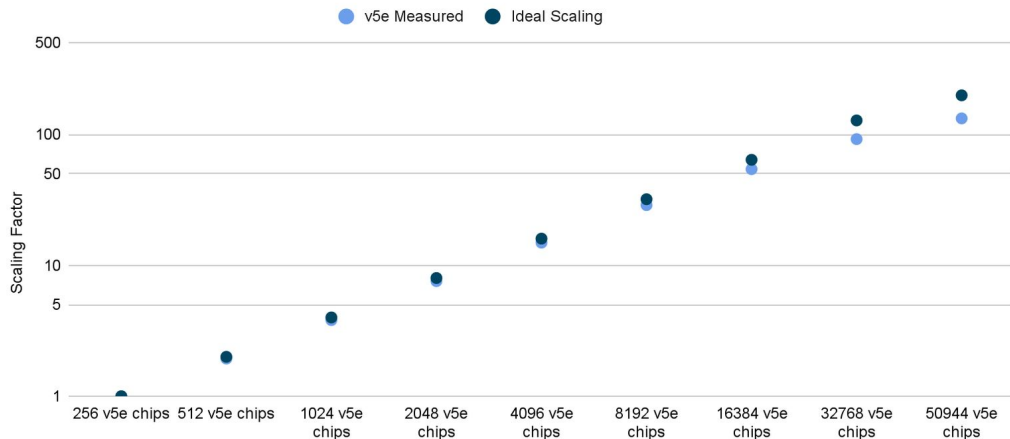
Dynamic Workload Scheduler

Flexible Consumption (On-Demand, CUD, Spot)

Google Cloud Infrastructure (CPU / GPU / TPU)

# World's largest distributed training job on GKE with TPU Multislice Training

TPU v5e Efficient Scaling with 32B LLM



Scaled to  
**50,000+**  
TPU v5e chips

Google Internal data for TPU v5e As of November, 2023: All numbers normalized per chip, seq-len=2048 for 32 billion parameter decoder only language model implemented using MaxText. \*2



# Strategies for Reducing Inference Latency

## Understanding the Bottlenecks



**Model Size & Complexity:** Larger models inherently take longer to process.



**Compute Resources:** Insufficient CPUs/GPUs or outdated hardware.



**Network Latency:** Data transfer between client, GKE, and data sources.



**Software Stack Overhead:** Inefficient inference servers, redundant data processing.



**Batching Strategies:** Incorrect batch sizes can lead to underutilization or increased latency.

# Strategies for Reducing Inference Latency

## Software & Deployment Optimizations



**Model Quantization & Pruning:** Reduce model size and complexity without significant accuracy loss.



**Optimized Inference Frameworks:** Use frameworks like TensorFlow Serving, TorchServe, Triton Inference Server.



**GKE Inference Gateway**



**KubeFlow Serving (KServe)**

# Strategies for Reducing Inference Latency

## Smart Scaling and Batching



**Horizontal Pod Autoscaler (HPA):** Scale pods based on CPU/memory utilization or custom metrics (e.g., QPS, latency).



**Vertical Pod Autoscaler (VPA):** Automatically adjust CPU and memory requests/limits for pods.



**Cluster Autoscaler:** Automatically adjust the number of nodes in your GKE cluster based on pod resource requests.



**Dynamic Batching (with Triton):** Grouping multiple inference requests into a single batch for efficient GPU utilization.



Client



GKE Cluster

1

Client request



2

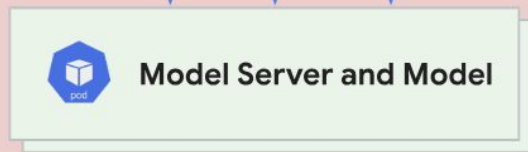
Inference Extensions



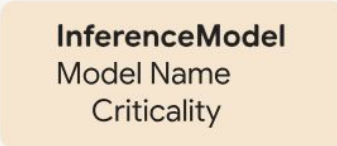
Model-aware request routing

3

**Inference Pool**



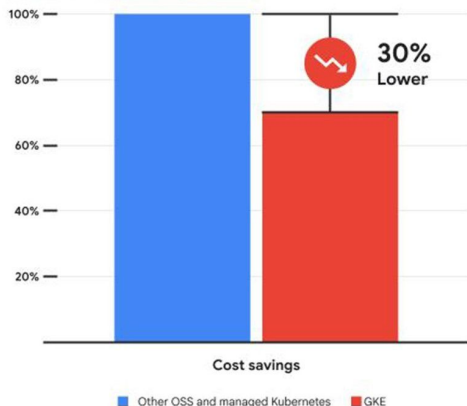
**Node Pool**



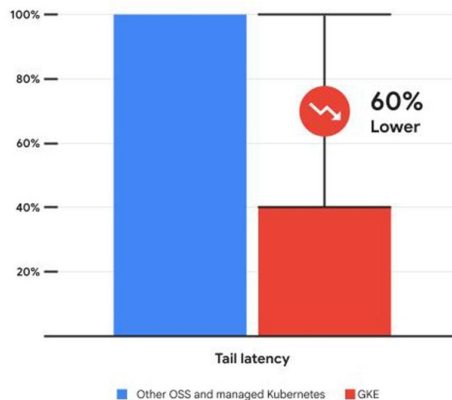
Model server load report

# GKE Inference Performance

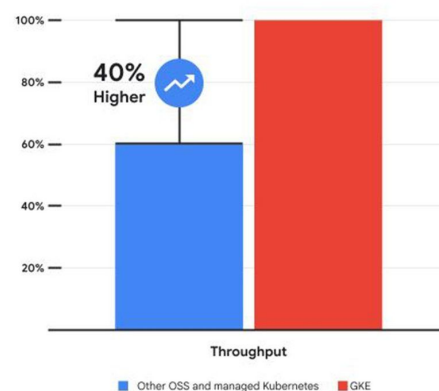
Cost to serve the same demand  
(lower is better)



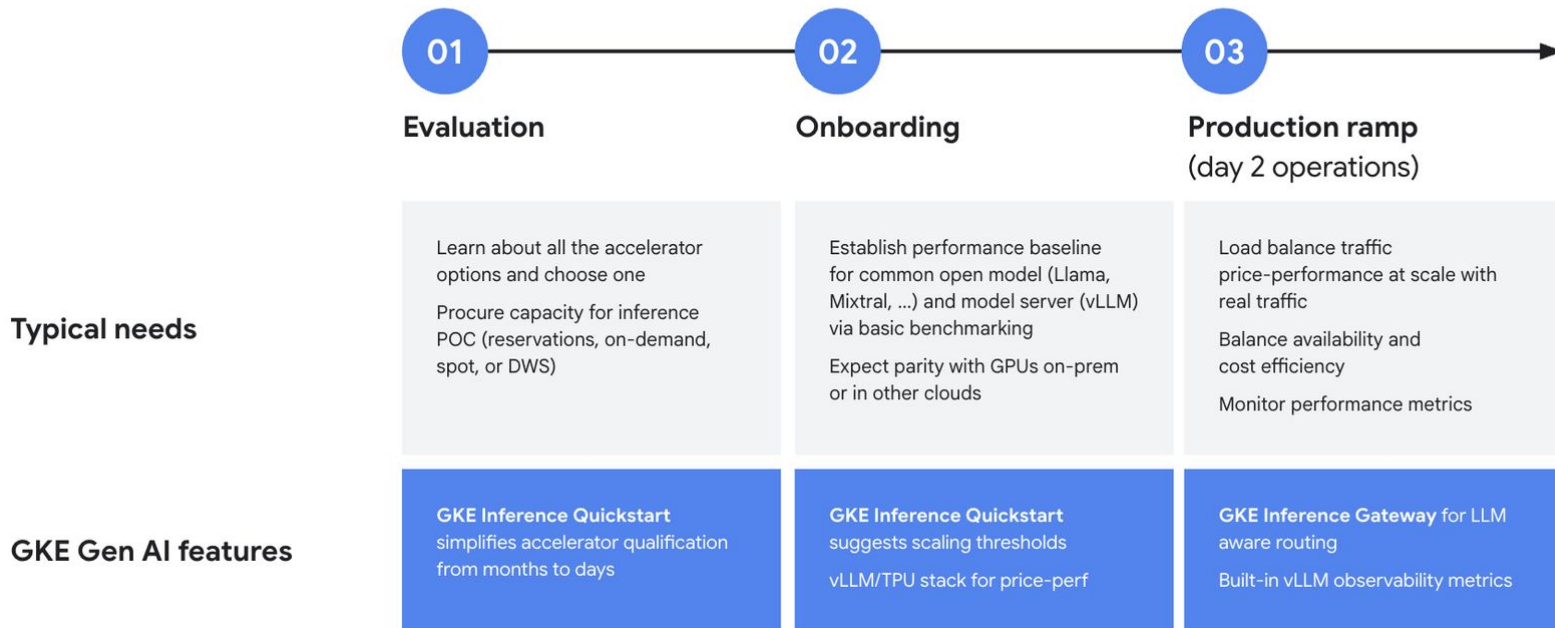
Tail latency of LLM (lower is better)



Throughput of LLM (higher is better)



# Solution mapping to customer journey



# Thank you

@misskecupbung

