

Optimizing the Al Stack for the Age of Inferencing

Maximizing Infrastructure Utility for Scalable Al

Speaker: Yujing Qian, VP of Engineering, GMI Cloud

World Summit AI, April 16, 2025

What to Expect: Inference Optimization Key Strategies



Hardware Acceleration



P/D disaggregation

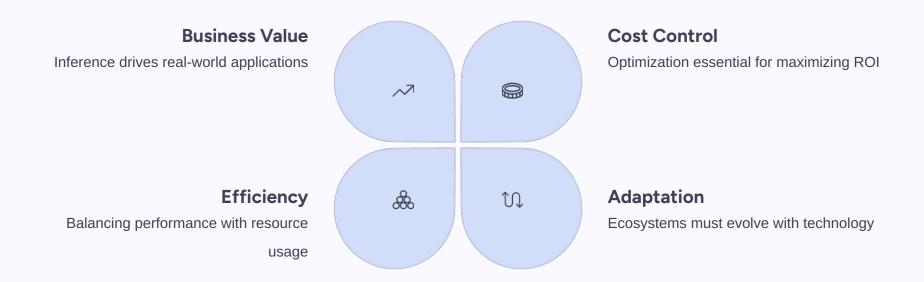


Performant inference engine



Distributed KV Cache

Why Inference Optimization Matters Today



The Cost of Inefficient Inference

30 x

Performance Gap

Due to bottlenecks in model execution

71%

Cost Reduction

Possible with proper optimization

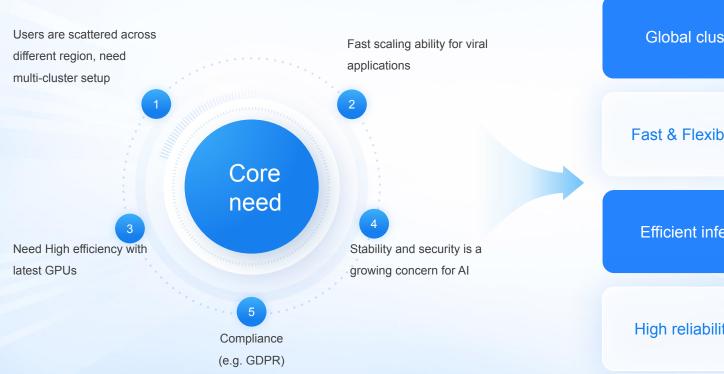
13 x

Latency Impact

Due to challenges balancing latency and throughput

Key of AI inference: efficiency, scalability and stability





Global cluster coverage

Fast & Flexible auto-scaling

Efficient inference cluster

High reliability and security

Scalable Inference System Best Practices

1

2

3

Performant inference engine

Minimum core component, squeeze every bit of performance out of your hardware

Workload Optimization

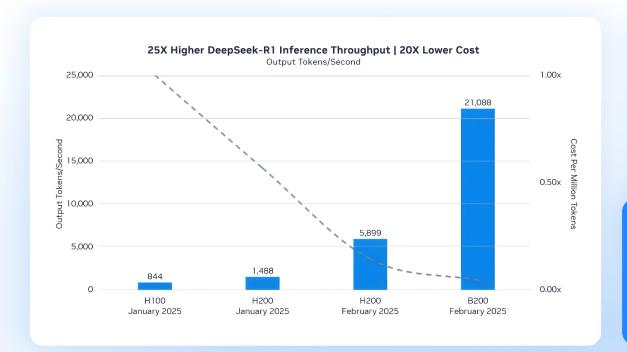
Intelligent distribution across hardware types and workload types

Global Autoscaling

Scale out our workload to hybrid cloud and multi-cluster

Always on the lookout for Best hardware





*data provided by NVIDIA

DeepSeek-FP4

Throughput tested on single 8 x

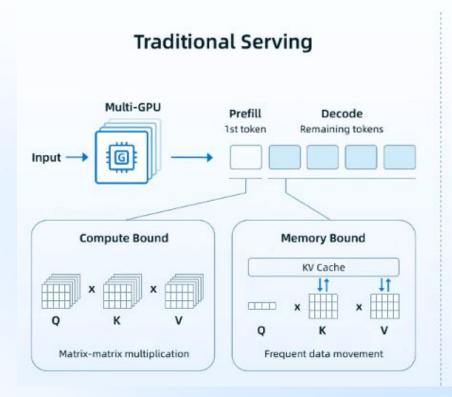
H100, H200, B200 server node

B 200 is more than 25x performant than H200. Latest hardware hosts latest model better

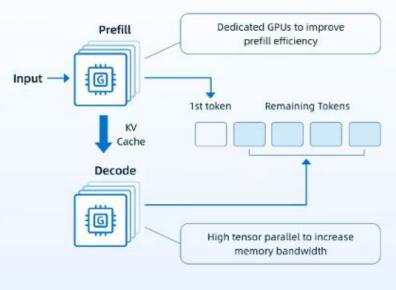
Hardware upgrades is one the easiest way to improve performance

Prefill/Decode Disaggregation



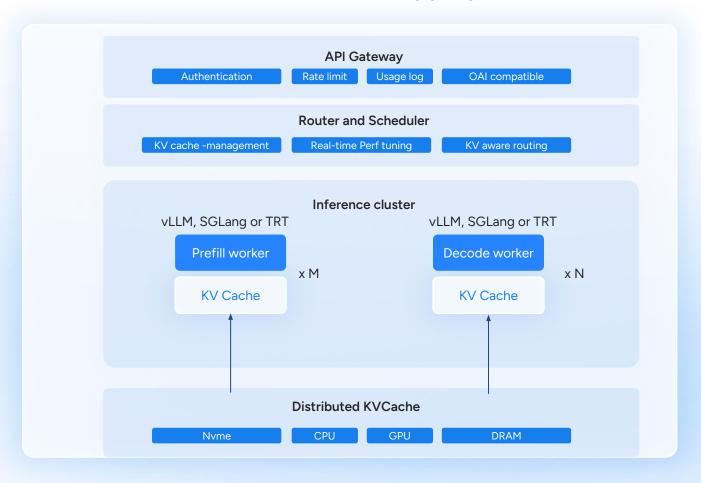


Disaggregated Serving



Inference cluster with Prefill/Decode Disaggregation





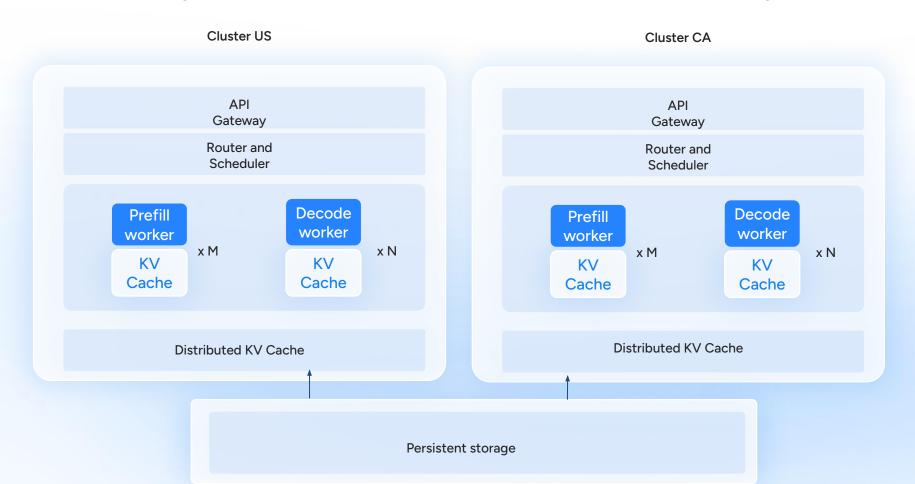
Inference cluster with Prefill/Decode Disaggregation





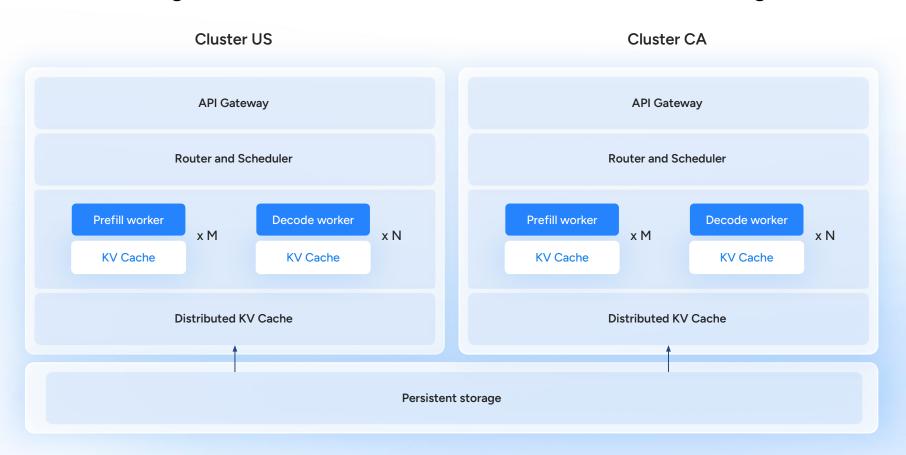
Inference engine with Distributed KV Cache + Persistent storage





Inference engine with Distributed KV Cache + Persistent storage







Boosting Cluster utilization



High speed autoscaling

Improve model loading time using techniques like RDMA



Fast LoRA switches

Use multi-LoRA to reuse base models



Offline queue

Queue up offline jobs and run them when cluster is idle

Average unattended private cluster

10+ min

Scaling time

Due to slow model loading and default scaling metric

40%

GPU utilization

GPU bubbles everywhere

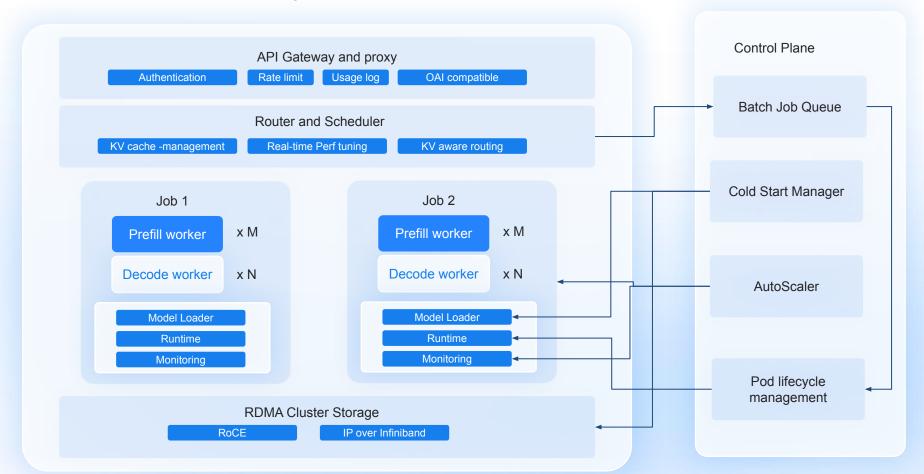
< 5

Workload supported

Due to challenges managing multiple jobs at the same time

Cluster level Autoscaling





Cluster level Autoscaling





Global auto scaling

Traffic aware
Dynamic Autoscaling

Fast model spin up with cluster preheat

Deploy same model in both Canada and US with one workflow.



Manual multi-cluster setup

days

Spin up time

Cluster, network, environment, file

transfer... you name it

>5

Manual steps

Just to get resource provisioned

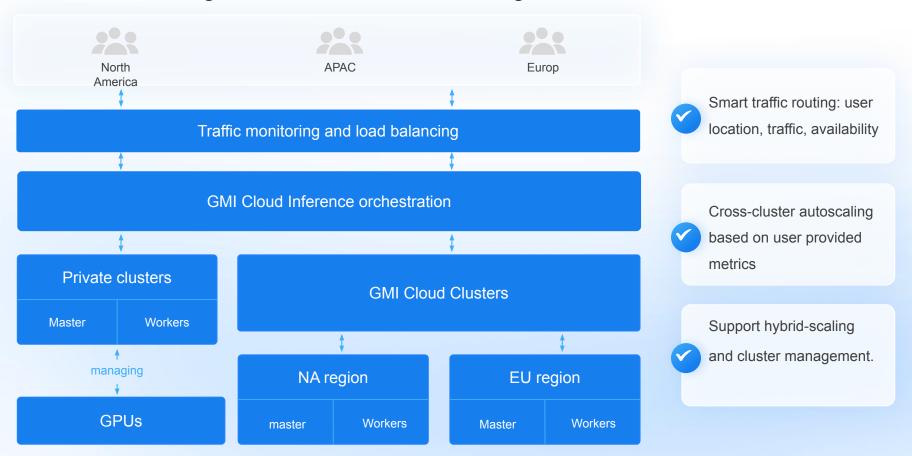
no

Support functionality

E.g. monitoring, endpoint control, etc

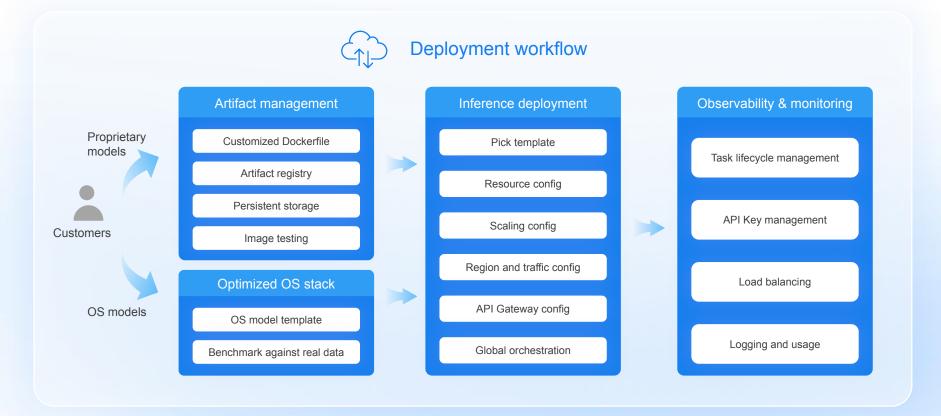
Inference Scaling Stack: Global Autoscaling





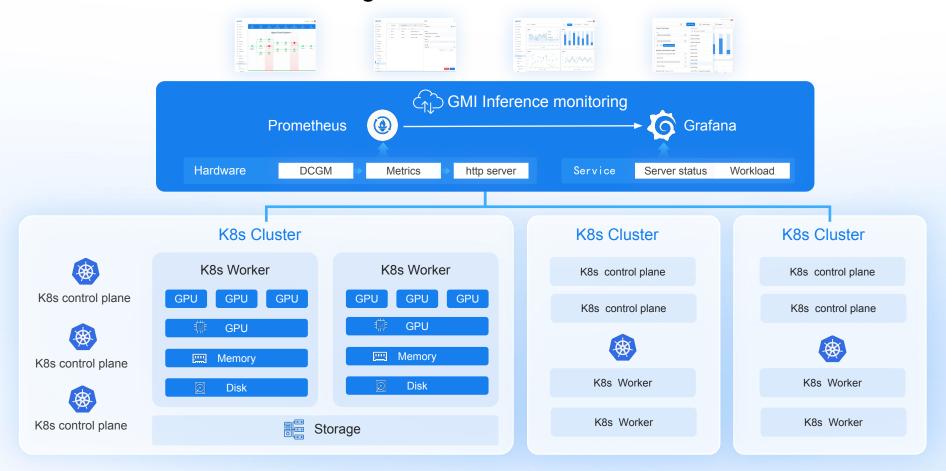
One-click Inference Engine Deployment





Global End to End monitoring





Key Takeaways Summary

Hardware Acceleration
Scale efficiently with dynamic resource management

Performant inference engine
Serve tokens with highest cost-efficiency

Workload Optimization

Balance online and offline jobs efficiently

Global Auto Scaling

Deploy and scale across multiple clusters and hybrid clouds



Get Started with Real-World Al Stacks



Try GMI Cloud Inference Engine
Get \$100 of credits with code: **INFERENCE**

Contact us at sales@gmicloud.ai