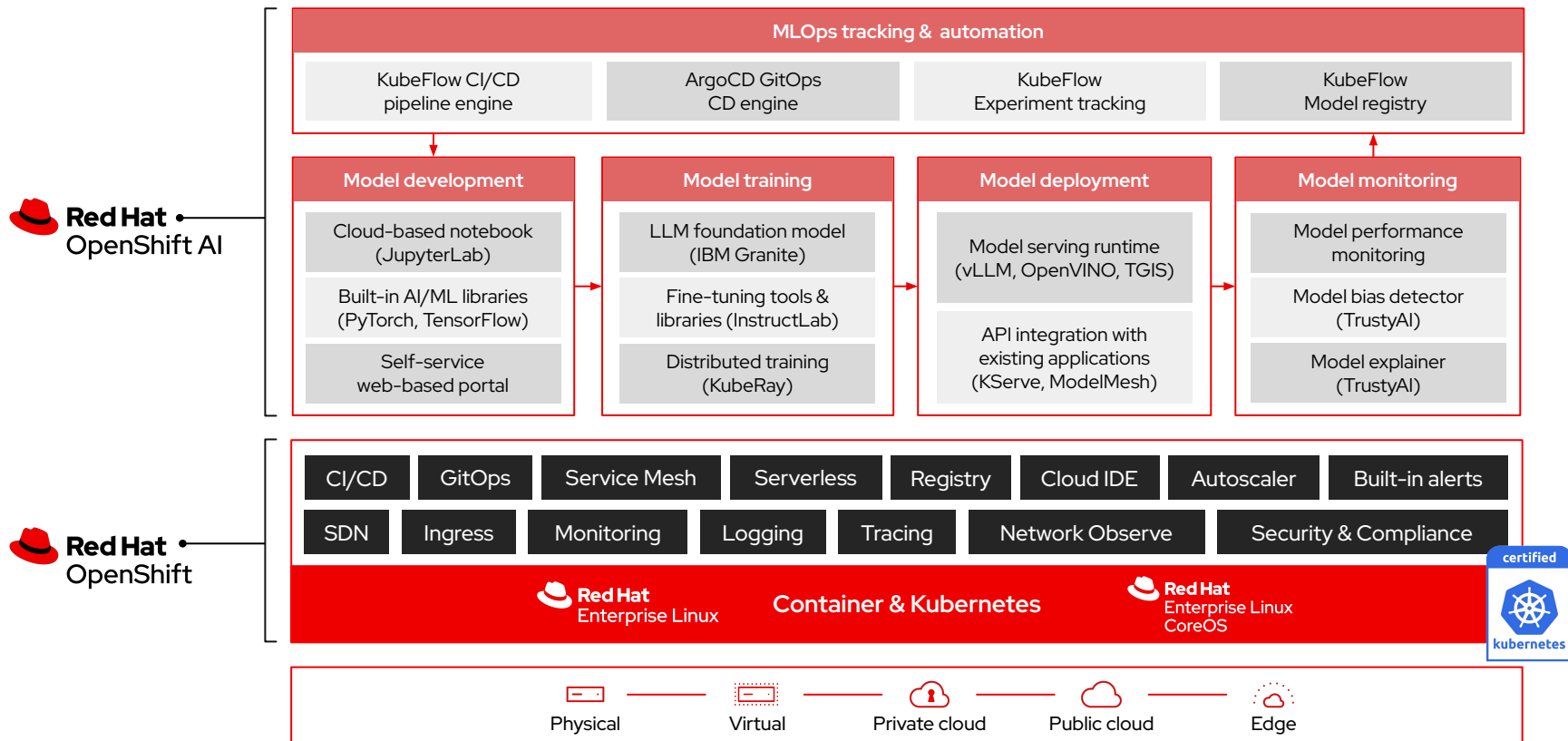


Red Hat OpenShift AI



OpenShift AI: the **MLOps** platform for **data scientists** and **operation teams** across **hybrid clouds**



Red Hat OpenShift AI

Dashboard Application

Data Science Projects

Admin Features

Model Registry

Object Storage



Model Development, Training & Tuning

Workbenches

- Minimal Python
- PyTorch
- CUDA
- Standard Data Science
- TensorFlow
- VS Code
- RStudio
- TrustyAI

CodeFlare SDK

ISV images

Custom images

Distributed workloads

KubeRay

Kueue

CodeFlare

Models

Granite Models

Ecosystem models

Data and model Pipelines

Model Serving

Serving Engines

Kserve

ModelMesh

Serving Runtimes

OVMS

vLLM, Caikit/TGIS

Custom

Model Monitoring

Performance metrics

Operations metrics

Quality metrics

OpenShift Operators

OpenShift GitOps



OpenShift Pipelines



OpenShift ServiceMesh



OpenShift Serverless



Prometheus



Red Hat OpenShift AI - Key features

Model development

Interactive, collaborative UI for **seamless access** AI/ML tooling, libraries, frameworks, etc.

Model serving

Model serving routing for **deploying models to production** environments

Model monitoring

Centralized monitoring for **tracking models performance and accuracy**

Data & model pipelines

Visual editor for **creating and automating** data science pipelines

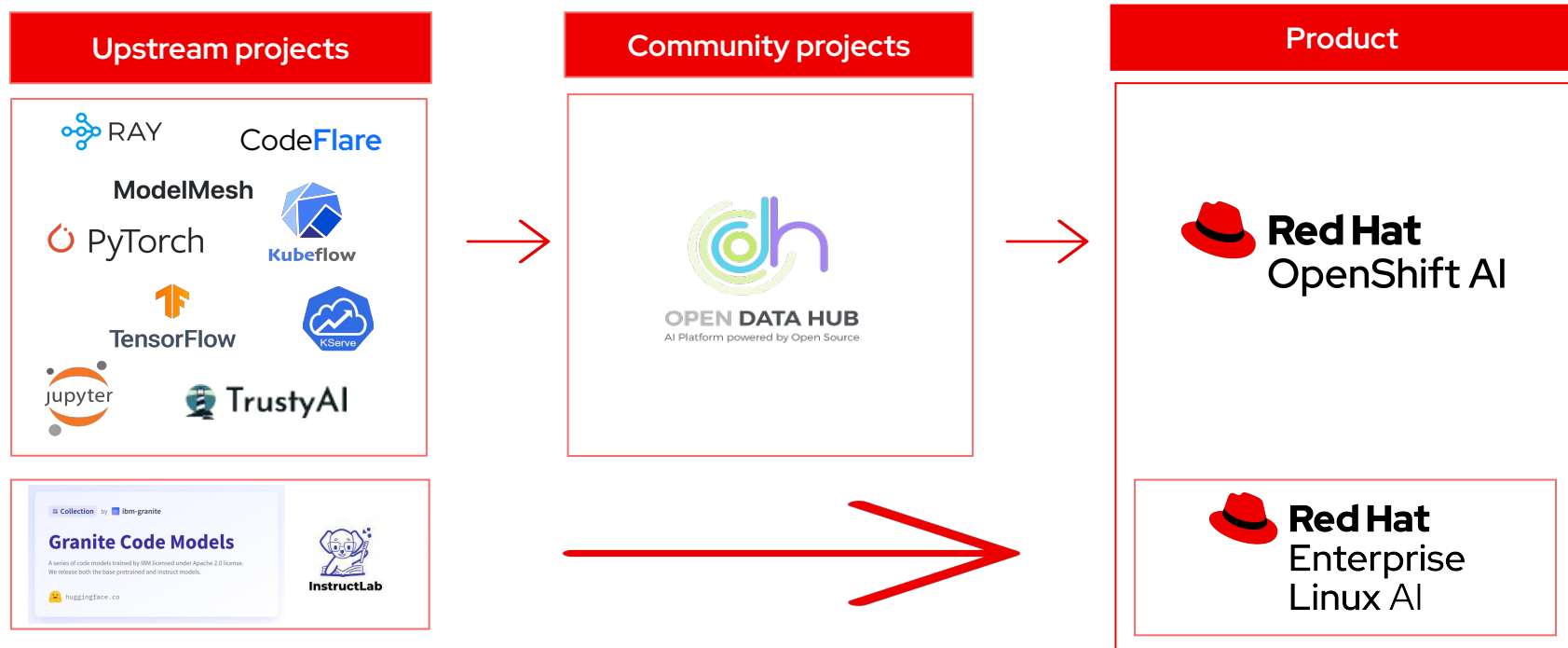
Distributed workloads

Seamless experience for **efficient data processing, model training, and tuning**

Trust & AI guardrails

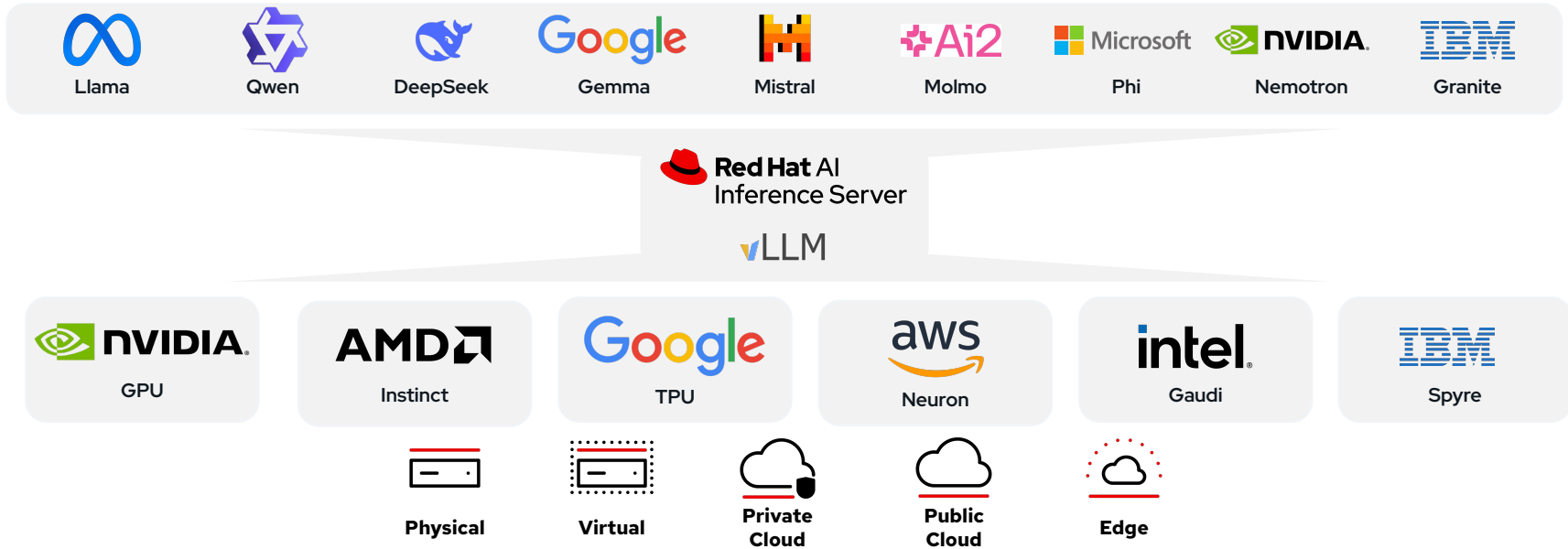
Improve LLM accuracy, performance, latency and **transparency**

Red Hat's AI/ML engineering is 100% open source



Red Hat AI Inference Server

vLLM connects model creators to accelerated hardware providers



Single platform to run any model, on any accelerator, on any cloud

Use case #1: Cost reduction by model optimization



Model quantization by Neural Magic

Leader in LLM serving, now acquired by Red Hat



Quantized by
Neural Magic



DeepSeek-R1-Distill-Llama-70B

DeepSeek-R1-Distill-Llama-8B

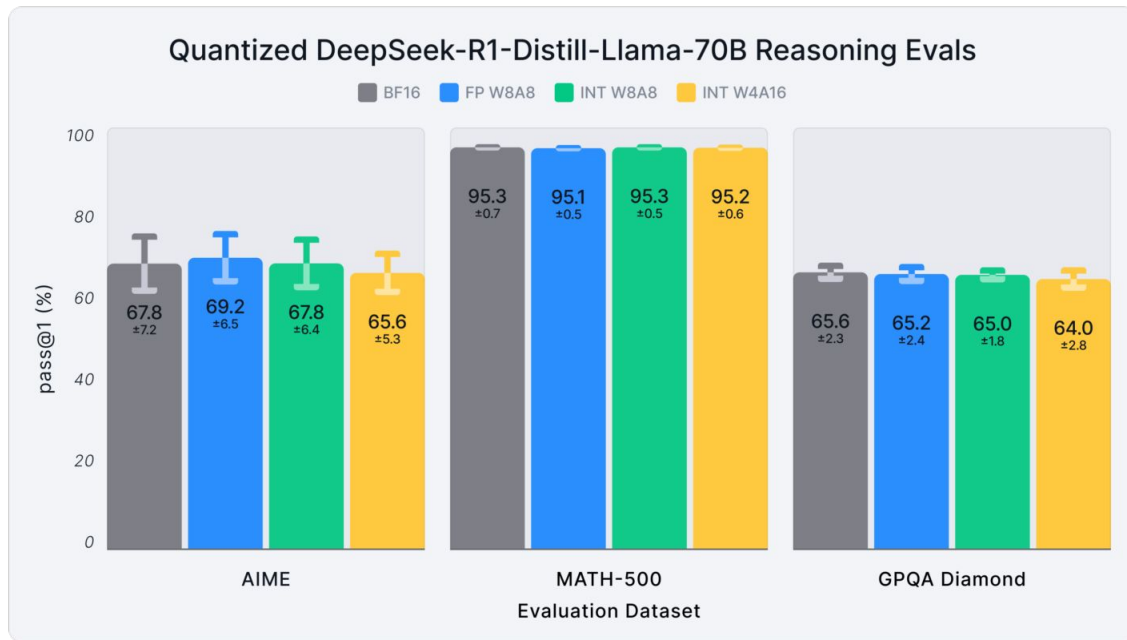
DeepSeek-R1-Distill-Qwen-32B

DeepSeek-R1-Distill-Qwen-14B

DeepSeek-R1-Distill-Qwen-7B

DeepSeek-R1-Distill-Qwen-1.5B

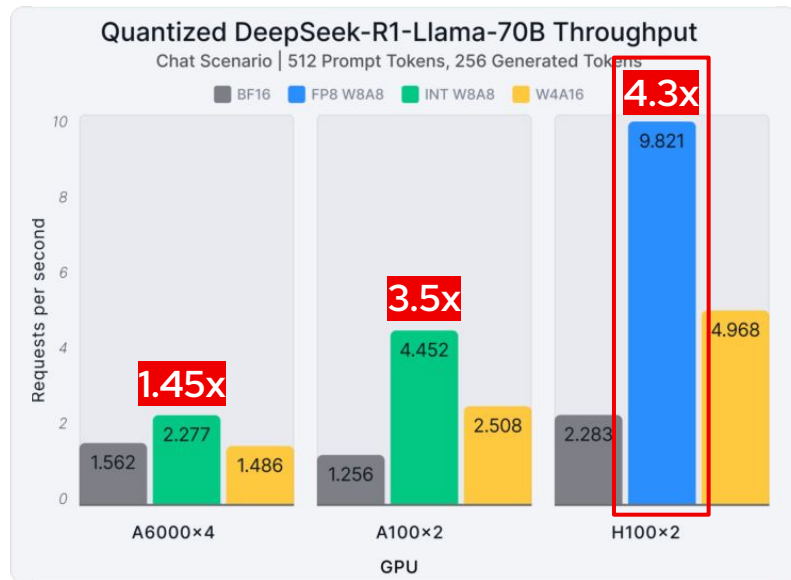
Fewer GPUs and \$\$\$, yet **same level of model accuracy**



- **Nearly 100% accuracy recovery** in W8A8 compression (i.e. 50% of model compressed)
- Only around 2.67% accuracy loss in W4A16 compression (i.e. 75% of model weight compressed)

Neural Magic reduces DeepSeek deployment \$\$\$

Use fewer GPU cards and resources, yet achieve the same = Save cost!



LLM serving throughput

↑430%

- Powered by **Neural Magic**, company acquired by Red Hat
- Advanced **quantization** (i.e. model compression) techniques

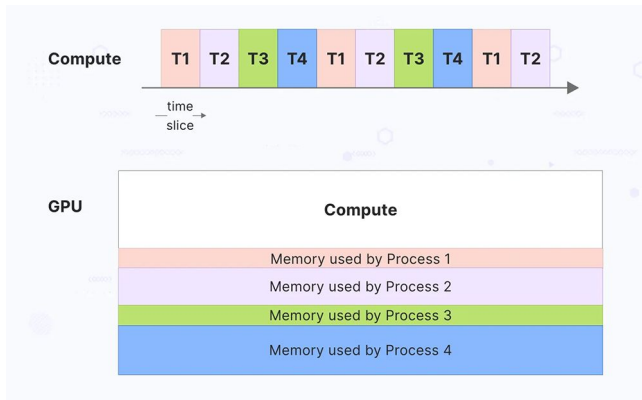
Use case #2: GPU sharing



GPU sharings across teams on OpenShift AI

Leverage NVIDIA time-slicing and MIG using **NVIDIA GPU Operator**

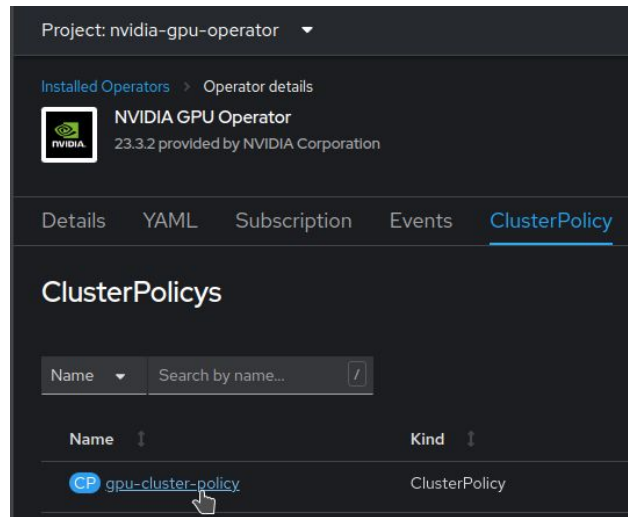
Time-slicing



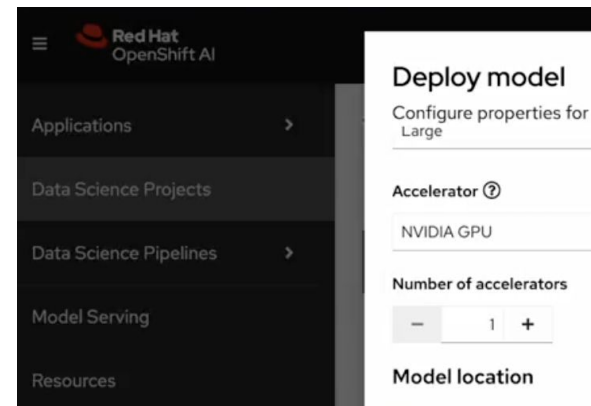
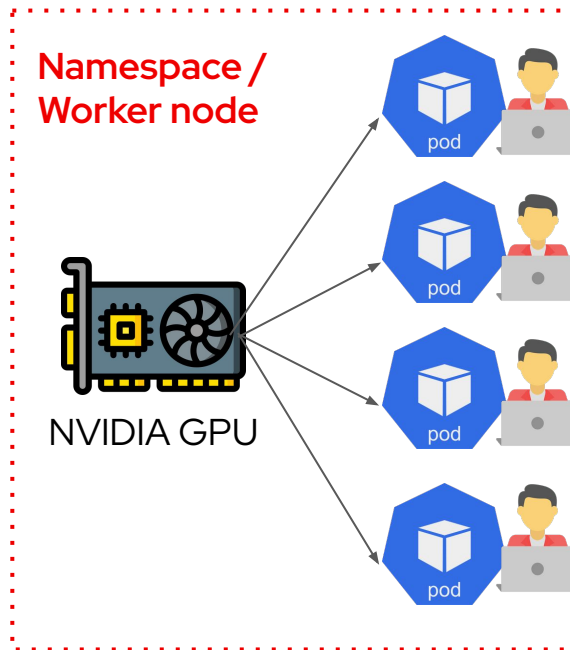
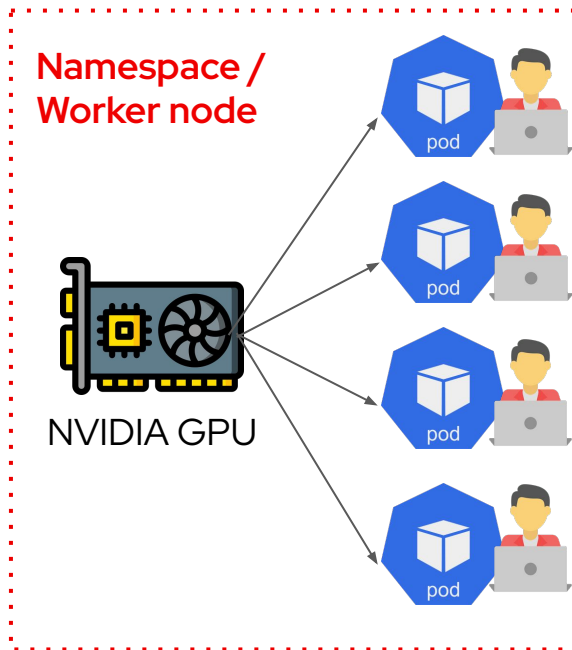
Multi-instance GPU (MIG)



- 1 x 7g.40gb
or
- 2 x 3g.20gb
or
- 3 x 2g.10gb
or
- 7 x 1g.5gb



Results of GPU sharing on OpenShift AI

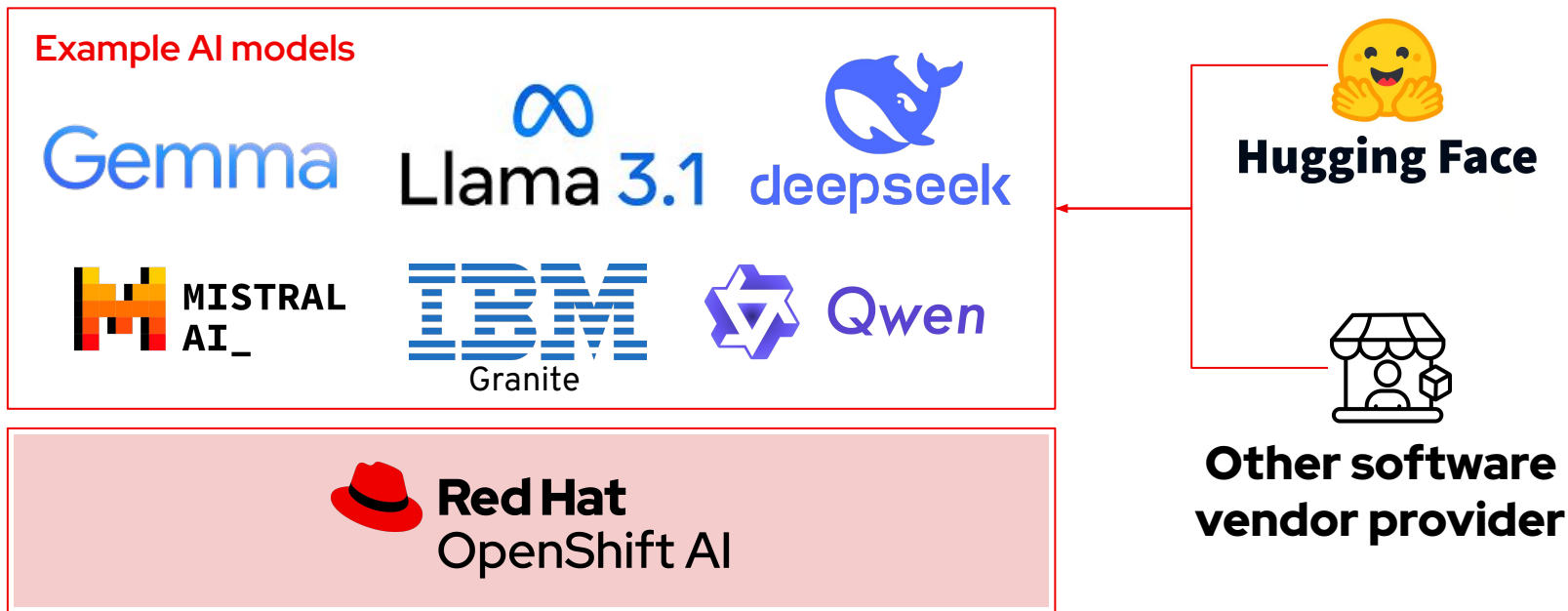


Use case #3: Diverse models support



Centralize AI model deployment on OpenShift AI

Bring whatever AI models to the platform



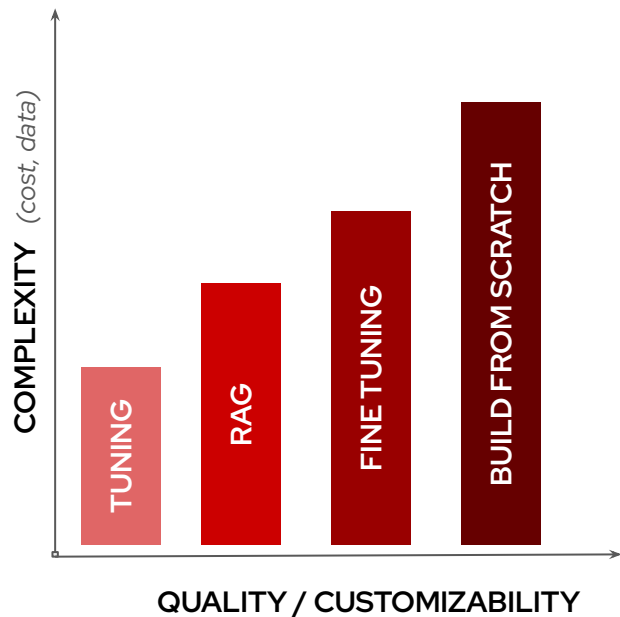
Note: Hugging Face and 3rd-party models are provided and supported by 3rd-party model providers

Use case #4: Integrate LLM with your knowledge



How to build and extend AI for your own business?

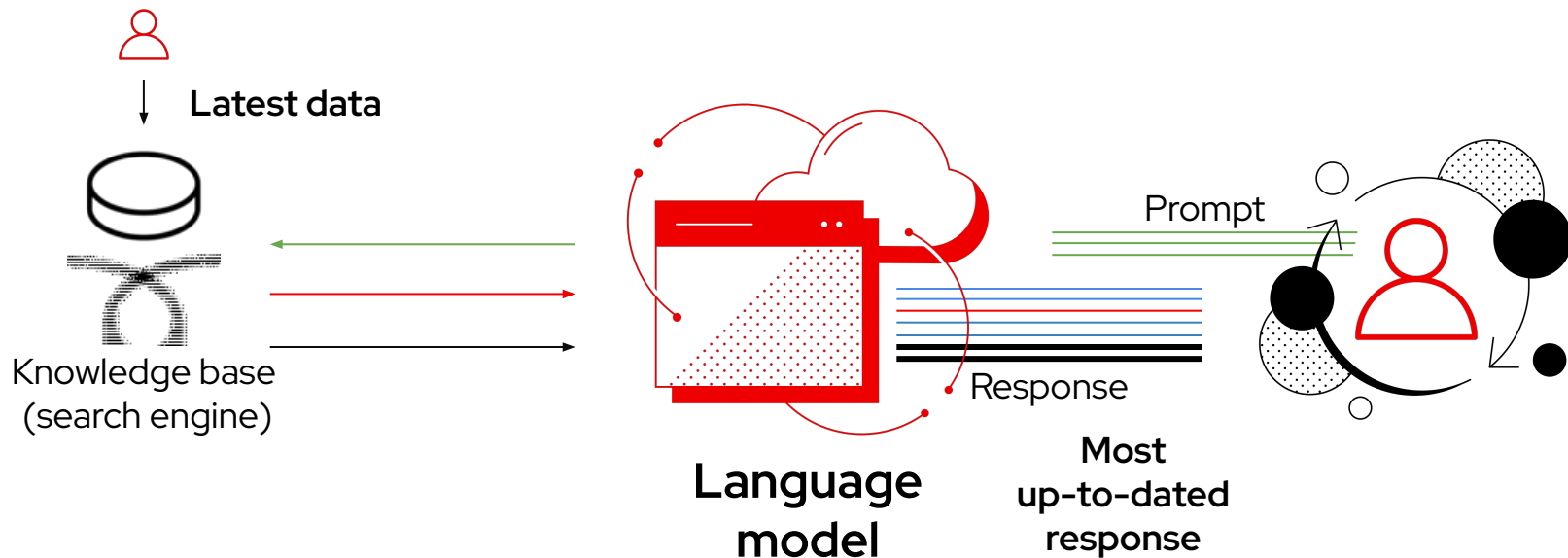
Pick the right way based on cost and customizability



- ▶ **Prompt tuning:** Prepend prompt with data.
- ▶ **Retrieval augmented generation (RAG):** return facts from external sources.
- ▶ **Fine tuning:** Retrain small part of the model with your data.
- ▶ **Training from scratch:** only make sense for predictive AI.

Retrieval-augmented generation (RAG)

Augment knowledge base with new data



Use case #5: Fine-tuning small language models





What is IBM Granite 3.1?

Pre-trained foundation model **provided and co-supported by IBM and Red Hat**

TRUE open sourced AI model

- **8B chat model, up to 32B code** SLM
- Apache 2.0 license
- Open weights on Hugging Face^[1]
- **Open training data set**^[2]

Multi-lingual support

- **English**, German, Spanish, French, Japanese, Portuguese, Arabic, Czech, Italian, Korean, Dutch and **Chinese**

Designed for Enterprise

- Up to **128K token context length**
- Language, code, agentic AI tasks.....

IP indemnification assurance

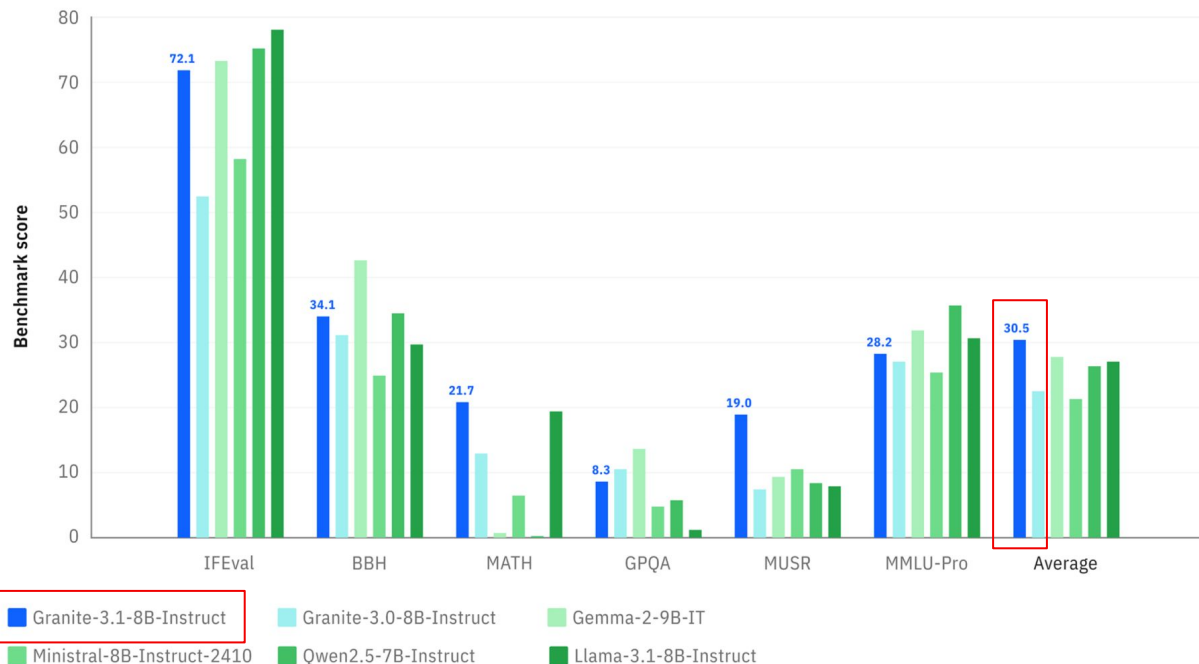
- **Protect customers from legal liability** if Granite infringes on others' IP

^[1] <https://huggingface.co/ibm-granite/granite-3.1-8b-instruct>

^[2] <https://github.com/ibm-granite/granite-3.0-language-models/blob/main/paper.pdf>

Granite can even **outperform** SOTA^[1] language model

Hugging Face OpenLLM Leaderboard benchmarks within the same class of model size

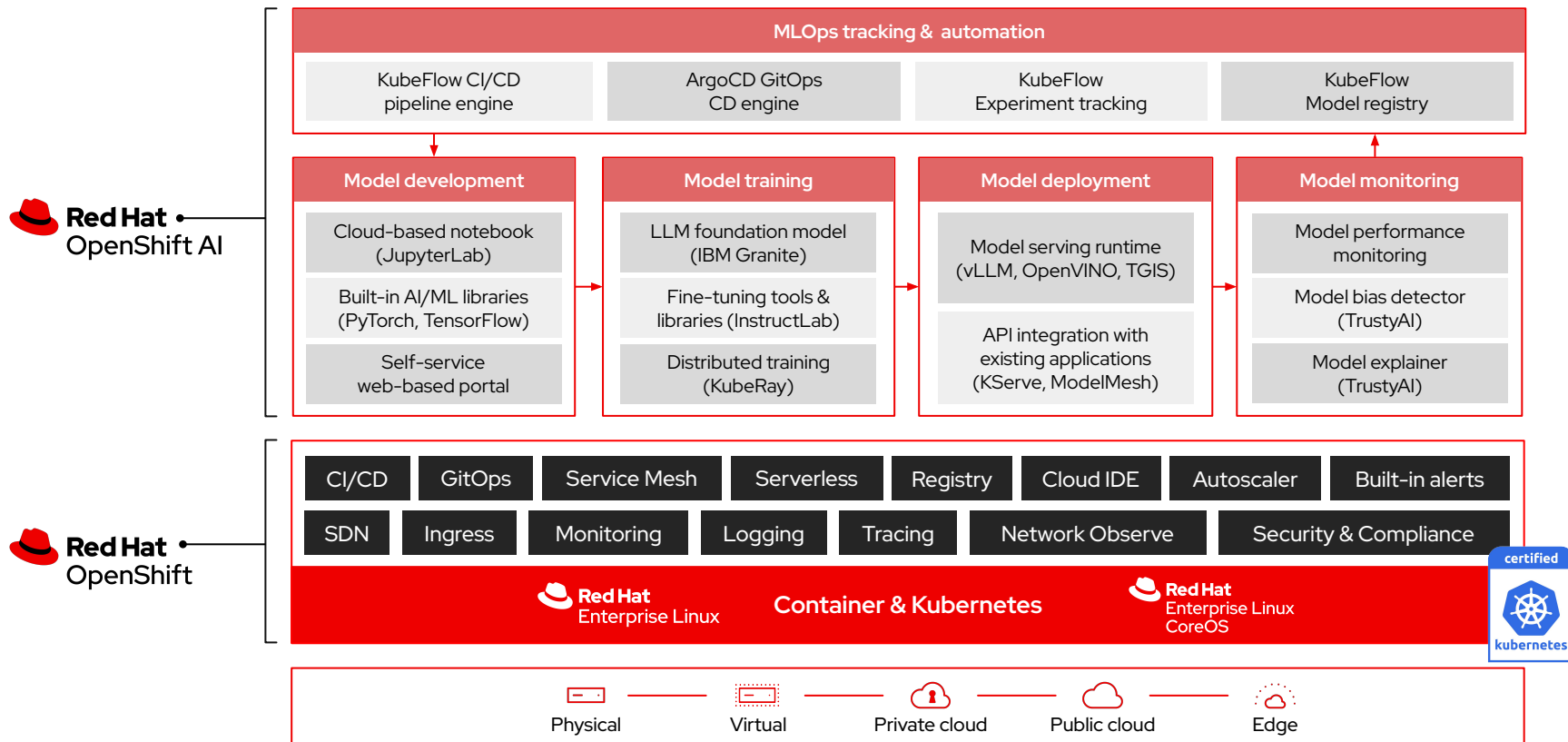


[1] State-of-the-Art

Use case #6: Streamline AI/ML with MLOps, across clouds



OpenShift AI: the **MLOps** platform for **data scientists** and **operation teams** across **hybrid clouds**

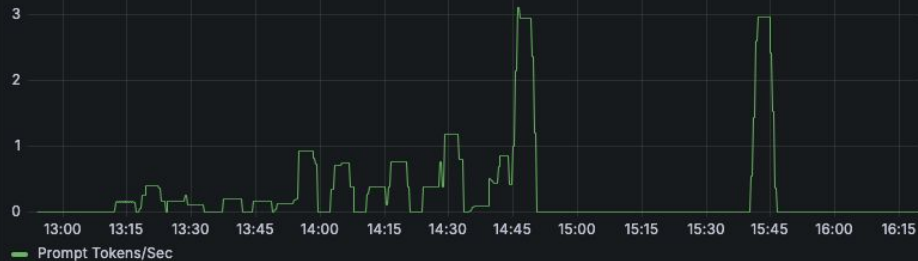




vLLM Generation Token Throughput



vLLM Prompt Token Throughput



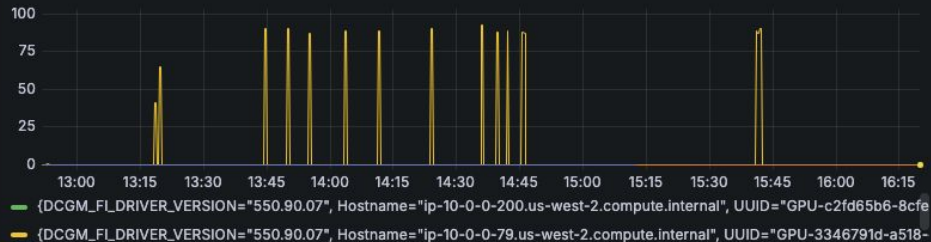
E2E Request Latency



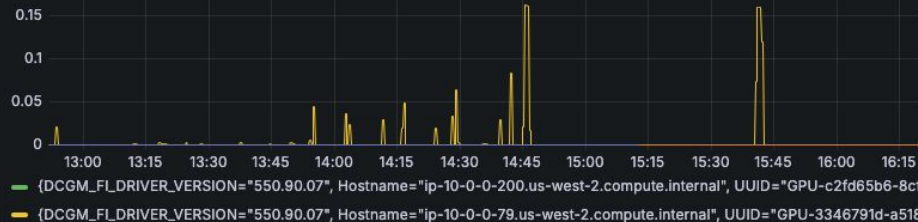
Scheduler State

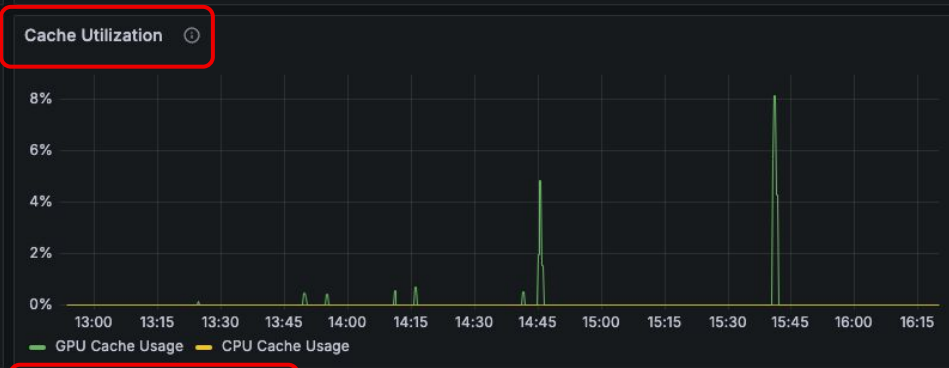
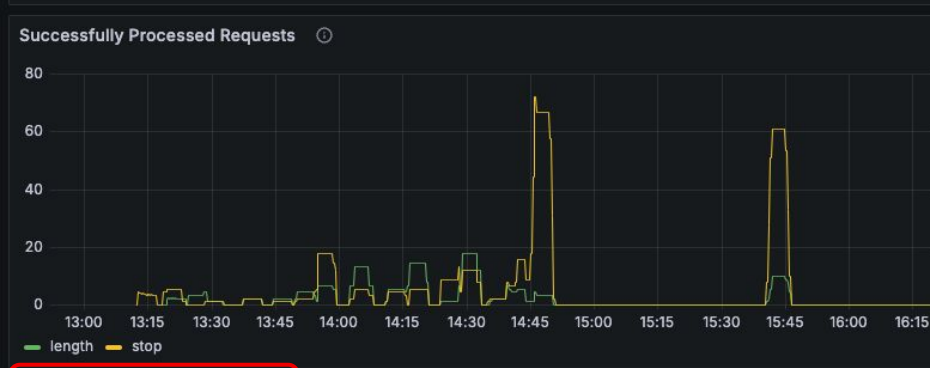
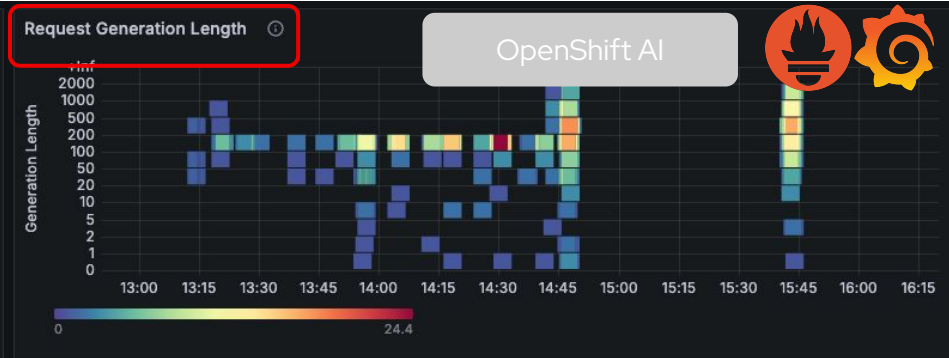


GPU Utilization



GPU Tensor Core Utilization





Values Red Hat can bring in AI



Provide governance

→ **Central place** to deploy and use AI models

→ **Standardize outputs** using approved AI models

→ **Tools to allow you shaping** your own AI for business

→ **Free to choose** Red Hat-provided or 3rd party models

Increase flexibility and control

Reduce privacy and legal risk

Run AI models **privately**,
no data is ever leaked out

IP indemnification
assurance from Red Hat

Run SLM over LLM, scale AI models horizontally

Provide AI model **inference distribution and optimization**

Reduce cost

