Red Hat AI

# RH AI 3 Overview & Roadmap
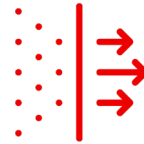
Red Hat

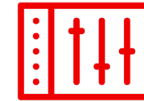# Generative AI adoption challenges



## Cost

Generative AI frontier model services are cost prohibitive at scale for many enterprise customer use cases.



## Complexity

Integrating models with private enterprise data for customer use cases is too complex for non-data scientists.



## Control

Increasing concerns with data privacy, security, and latency are compelling organizations to adopt hybrid strategies.

Trusted, Consistent and Comprehensive foundation

Hardware Acceleration

**Physical**

**Virtual**

**Private Cloud**

**Public Cloud**

**Edge**

3

**Red Hat** AI

**Accelerate the development and delivery of AI solutions** across hybrid-cloud environments

Increase efficiency with **fast, flexible and efficient inferencing**

Simplified and consistent experience for **connecting models to data**

**Accelerate Agentic AI** deployments

Flexibility and consistency when **scaling AI across the hybrid cloud**

**Red Hat**

# Introducing RH AI 3

## Flexible and Efficient Inference

- ‣ GA distributed inference (llm-d)
- ‣ New validated and optimized models
- ‣ vLLM enhancements
- ‣ LLM Compressor GA

## Agentic AI

- ‣ AI experiences: AI hub and gen AI studio
- ‣ Model Context Protocol support & MCP Server access in gen AI studio
- ‣ Llama Stack API integration

## Connecting Data to Models

- ‣ Modular and extensible approach for: data ingestion, synthetic data generation, tuning, evaluations.
- ‣ RAG enhancements & partner integrations
- ‣ Feature Store UI

## AI Platform

- ‣ Model catalog and registry GA
- ‣ Model as a Service provider enhancements and API Mgt integration
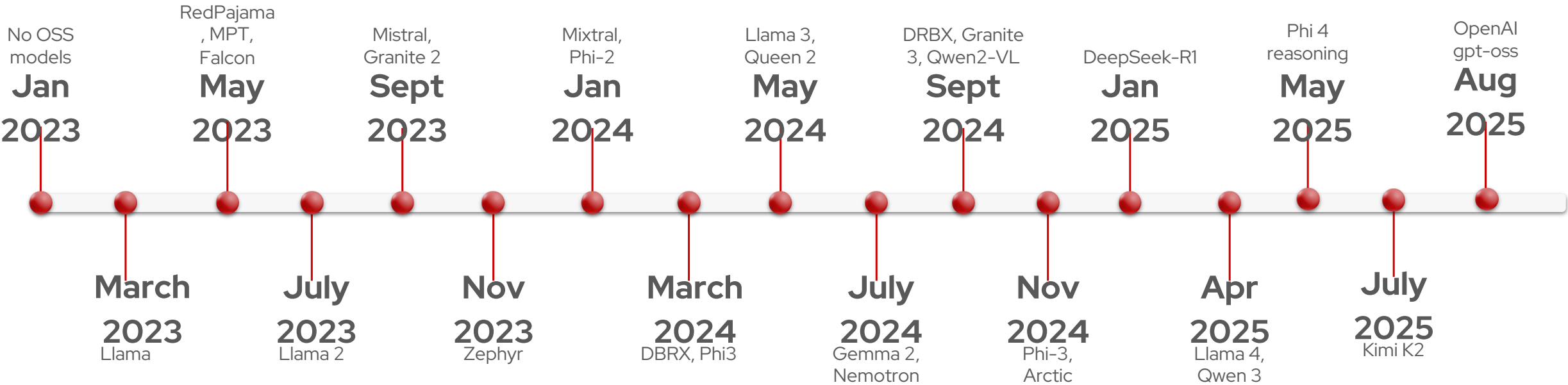- ‣ GPU as a Service enhancements

**Red Hat AI**

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

Red Hat

# Fast, flexible, and efficient inference

# Expanding choice of models

## There has been an explosion of capability from open-source over the last 2 years

| No OSS models | RedPajama, MPT, Falcon | Mistral, Granite 2 | | Mixtral, Phi-2 | | Llama 3, Queen 2 | | DRBX, Granite 3, Qwen2-VL | | DeepSeek-R1 | | Phi 4 reasoning | | OpenAI gpt-oss |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Jan 2023** | **May 2023** | **Sept 2023** | | **Jan 2024** | | **May 2024** | | **Sept 2024** | | **Jan 2025** | | **May 2025** | | **Aug 2025** |

| | **March 2023** Llama | | **July 2023** Llama 2 | | **Nov 2023** Zephyr | | **March 2024** DBRX, Phi3 | | **July 2024** Gemma 2, Nemotron | | **Nov 2024** Phi-3, Arctic | | **Apr 2025** Llama 4, Qwen 3 | | **July 2025** Kimi K2 |

Meta · Ai2 · NVIDIA · MISTRAL AI_ · Hugging Face · databricks · snowflake · IBM · Google · Microsoft · deepseek

8

Red Hat

# Expanding choice of Accelerators

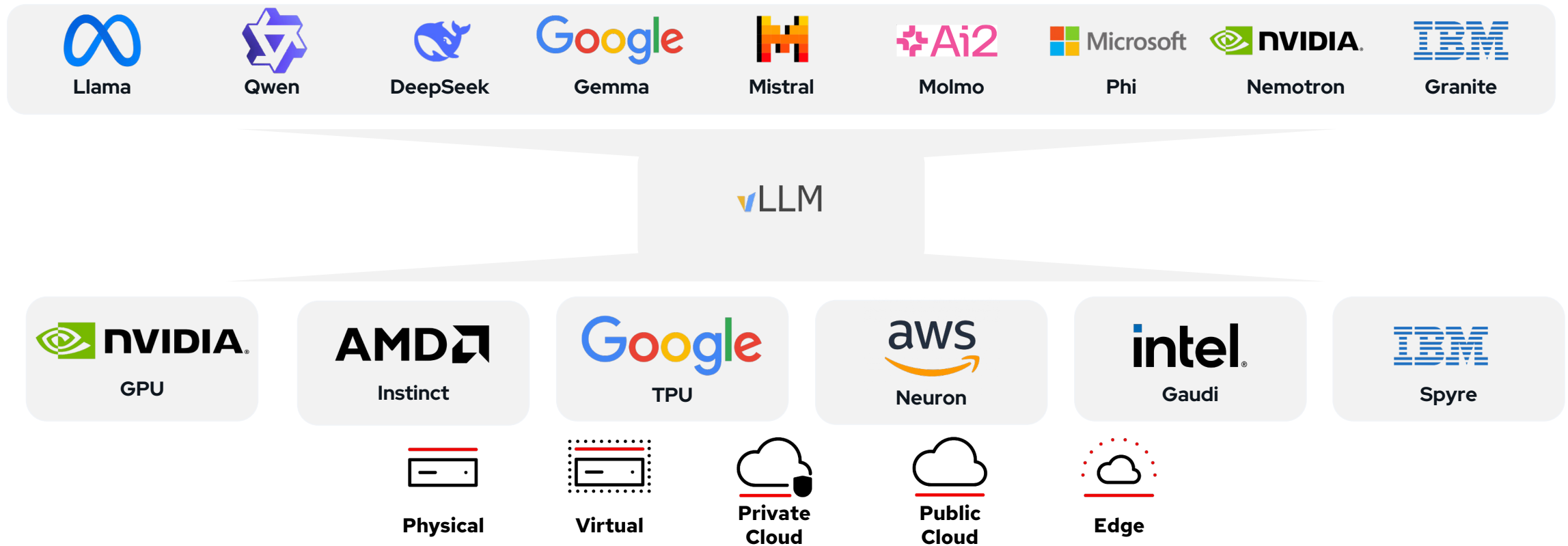## NVIDIA is still the major player but other entrants gaining ground

# Enterprise GenAI Inference Platform

## Holistic approach to optimize and operationalize deployment and scaling of open-source LLMs

# vLLM Inference Server

## vLLM v1 with enhanced performance, expanded model and hardware support



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

# Red Hat AI Model Repository

New validated and optimized models and LLM Compressor now GA

### Broad Collection of models

**Llama**  **Qwen**  **Google / Gemma**

**Mistral**  **DeepSeek**  **Microsoft / Phi**
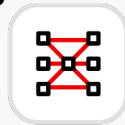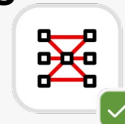
**Ai2 / Molmo**  **IBM / Granite**  **NVIDIA / Nemotron**
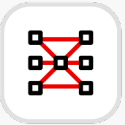
**Choice of Models**
- ‣ Transformers (Dense, MOE), Multi-modal LLMs, Embeddings Models, Hybrid / Novel Attention, Vision
- ‣ Hugging Face compatible (safe tensors), OCI-compatible containers

**Validated models**
- ‣ Tested using realistic scenarios
- ‣ Assessed for performance across a range of hardware
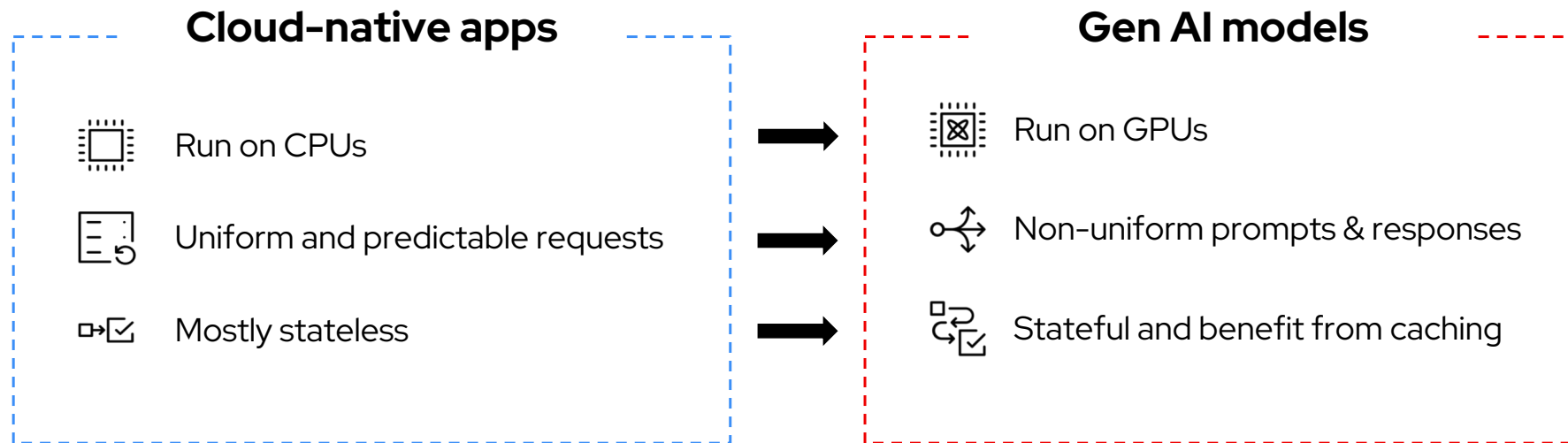- ‣ Done using GuideLLM benchmarking and LM Eval Harness

**Optimized models**
- ‣ Compressed for speed and efficiency
- ‣ Designed to run faster, use fewer resources, maintain accuracy
- ‣ Done using LLM Compressor with latest algorithms

**Red Hat**

# Overcoming the generative AI challenges

## Running LLMs efficiently

**Cloud-native apps**

Run on CPUs

Uniform and predictable requests

Mostly stateless

**Gen AI models**

Run on GPUs

Non-uniform prompts & responses

Stateful and benefit from caching

Red Hat

# Inference at scale everywhere

Distributed, scalable gen AI inference for Enterprise AI

**Red Hat** AI

Now includes llm-d

**llm-d reimagines how LLMs run on Kubernetes**

▸ Lower infrastructure spend for AI

▸ Great performance at larger scale
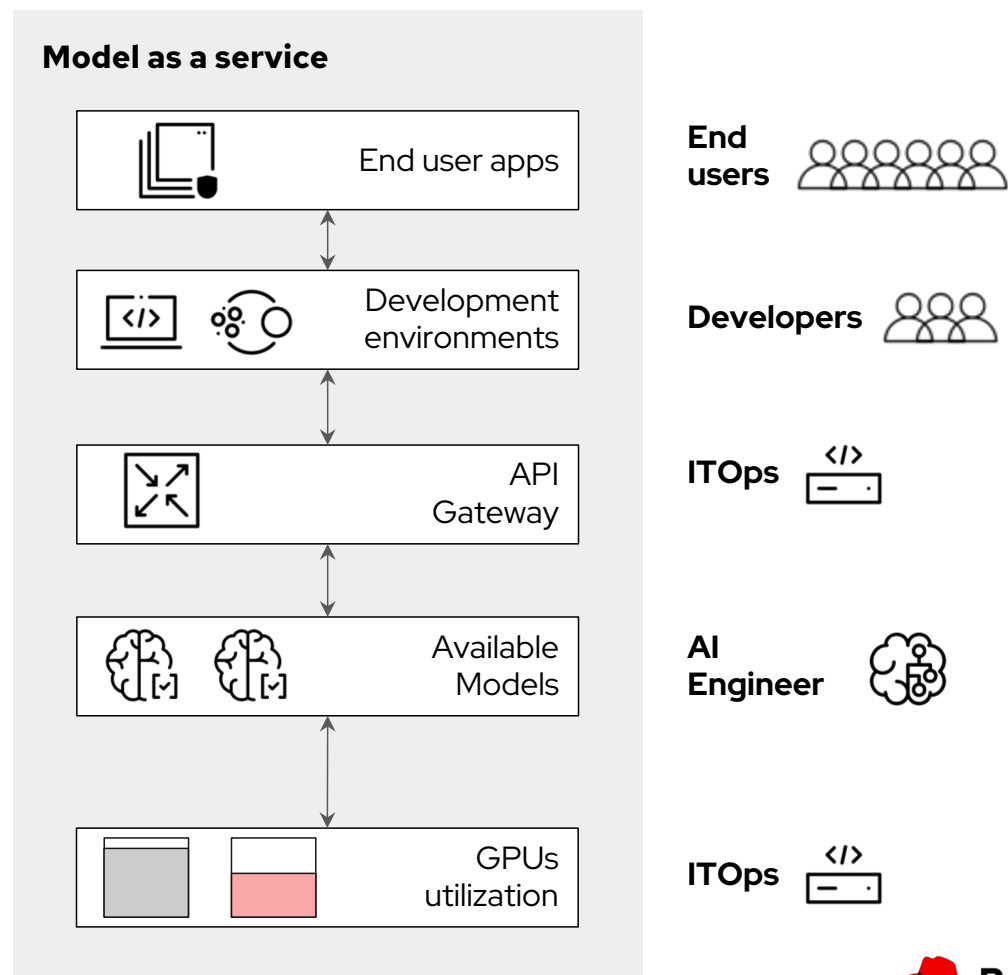
▸ Seamless scaling for unpredictable AI demand

Deliver faster, cheaper, and more manageable AI systems for enterprise production

**llm-d is GA in RHOAI 3.0**

**Red Hat**

# Model-as-a-Service in OpenShift AI (Dev Preview)

Offering AI models as the service to a larger audience

▸ IT serves common models centrally

  ■ Generative AI focus, applicable to any model

  ■ Centralized pool of hardware

  ■ Platform Engineering for AI

▸ Models available through the RHOAI console

▸ Developers consume models, build AI applications

  ■ For end users (private assistants, etc)

  ■ To improve products or services through AI

▸ Shared Resources business model keeps costs down

**Model as a service**



16

# Connecting models to data

Red Hat

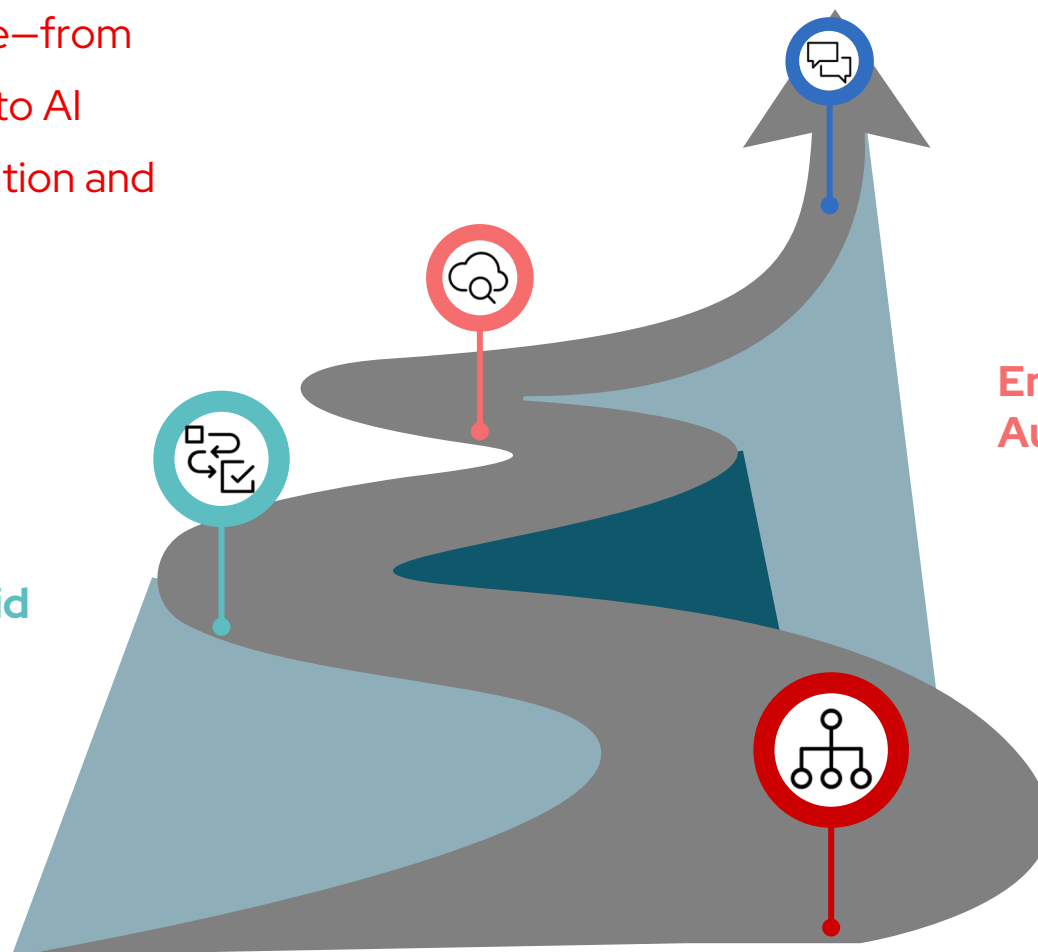# The evolution of Red Hat AI's tuning and alignment tools

Access to AI tooling that caters to different levels of AI expertise—from developers to data scientists to AI engineers—ensuring collaboration and frictionless interaction.

**Modular and extensible approach for data ingestion, synthetic data generation, model tuning, and evaluation**

**Enhance gen AI apps with Retrieval Augmented Generation (RAG)**

**Continual fine-tuning via Online Supervised Fine Tuning and rapid domain adaptation via LoRA/QLoRA**

**InstructLab allows customers to align models by adding new knowledge and skills**

18

Red Hat

# New model customization approach offers a modular extensible architecture

**Data processing**

**Synthetic data Generation hub**

**Training hub**

**Evaluations**

19

Simplifies document processing and parsing into AI-readable data for model customization and RAG applications

Generate high-quality data, with dynamic parameters, run-time visibility, and multilingual support

An algorithm-focused interface for common llm training, continual learning, and reinforcement learning techniques

Enables large language model inference evaluations.

Red Hat

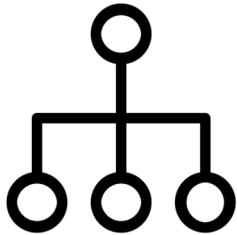# Connecting models to data

## Build customized AI solutions that address domain specific business cases

**Coming soon**

### Prompt design
*Prompt tuning and engineering*



**Design and engineer the prompts** to enhance GenAI model responses and achieve more specific and accurate outcomes.

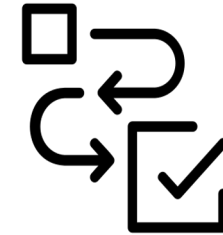**Enhanced**

### RAG
*Retrieval Augmented Generation*



**Enhance Gen AI model generated text** by retrieving relevant information from external sources, improving accuracy and depth of model's responses.
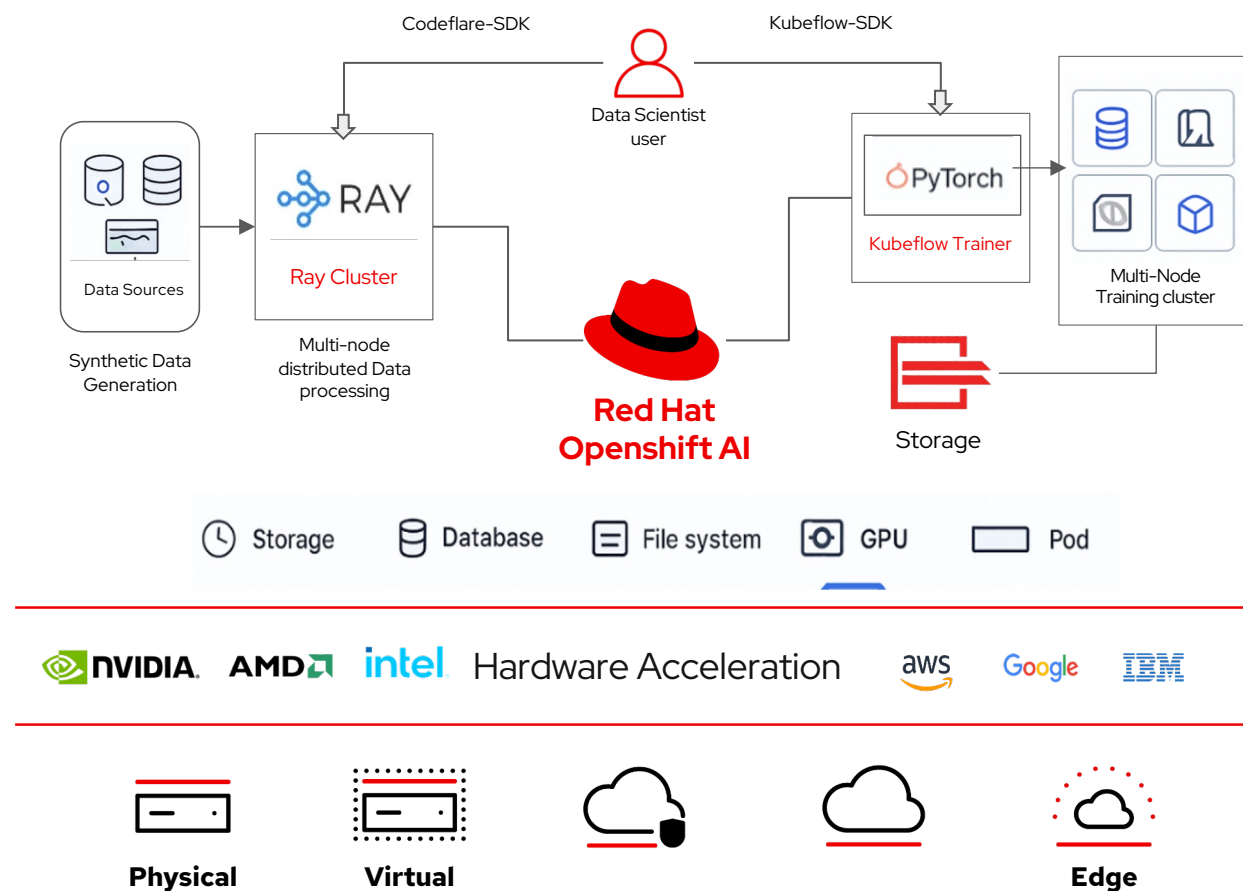
**Enhanced**

### Fine tuning
*InstructLab, LoRa and QLora*



**Adjust a pre-trained model on specific tasks or data**, improving its performance and accuracy for specialized applications without full retraining.

21

Red Hat

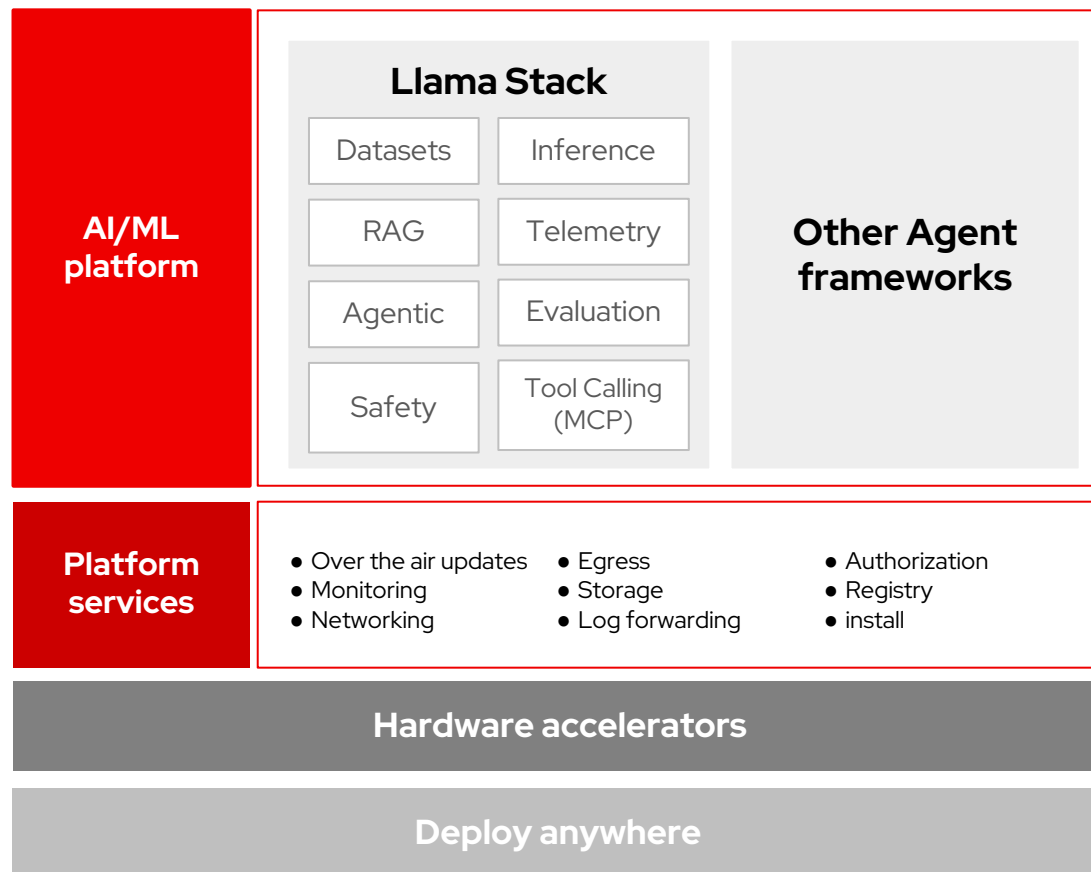# Ray + Kubeflow Training : End to End ML Pipeline

Accelerate Agentic AI

# Red Hat AI provides an agile, stable foundation to accelerate the development and deployment of AI agentic workflows.

▸ Allows running and managing agents as microservices.

▸ Simplifies production deployment by managing LLM serving and scaling.

▸ Offers native capabilities to build and manage agents with Llama Stack, and standardized communication protocols (MCP).

▸ Provides the flexibility to integrate preferred tools like LangChain and Crew AI.

# A modular approach to building AI agents

| AI/ML platform | **Llama Stack** | | **Other Agent frameworks** |
|---|---|---|---|
| | Datasets | Inference | |
| | RAG | Telemetry | |
| | Agentic | Evaluation | |
| | Safety | Tool Calling (MCP) | |

| Platform services | • Over the air updates | • Egress | • Authorization |
|---|---|---|---|
| | • Monitoring | • Storage | • Registry |
| | • Networking | • Log forwarding | • install |

**Hardware accelerators**

**Deploy anywhere**

**Red Hat AI allows to:**

▸ Build agents using **Llama Stack's native capabilities and implementations**.

▸ **Bring compatible Llama Stack implementations** to OpenShift AI.

▸ **Use your own agent framework** and selectively incorporate Llama Stack APIs.

▸ **Build with Core Primitives** and manage your own agent framework as a standard workloads.

**Red Hat**

# AI Dedicated Experiences

## Dedicated dashboard experiences provide a seamless experience to platform and AI engineers
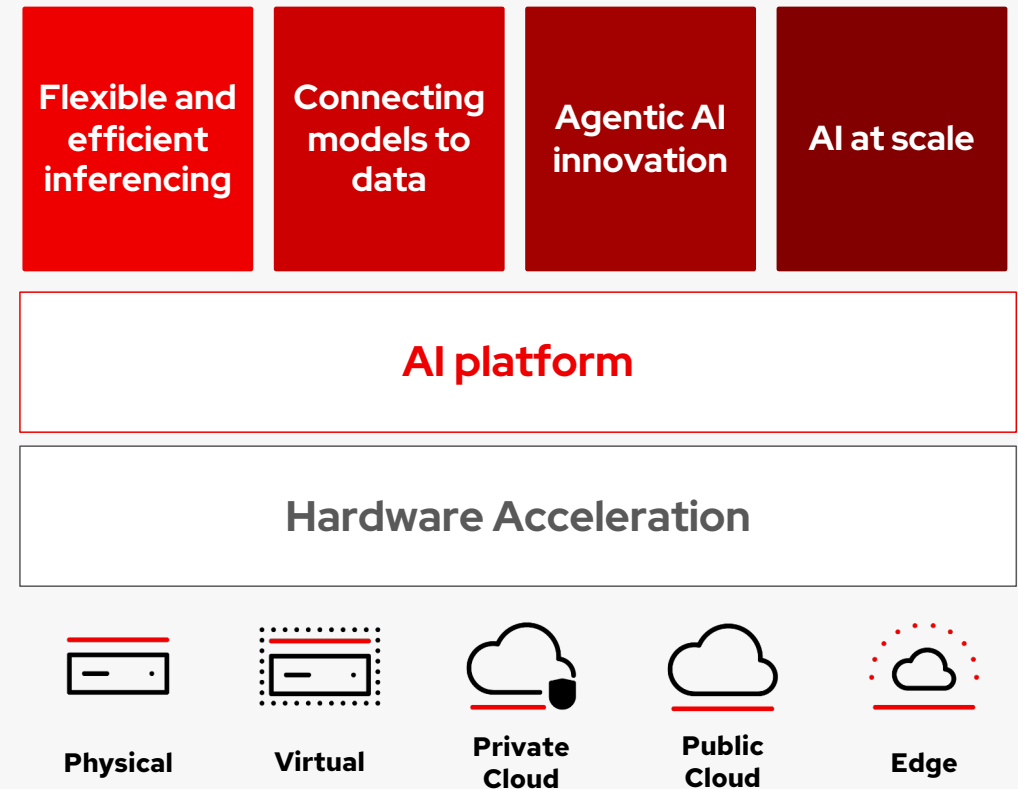
**AI hub**

**Gen AI studio**



28

# Scaling AI across the hybrid cloud

Red Hat

**Red Hat AI** provides a platform to consistently build, deploy and manage AI models, training and agentic applications at scale across the hybrid cloud

It includes:

▸ Enterprise-grade, flexible and secure AI platform

▸ Private and sovereign AI capabilities and practices

▸ GPU-as-a-Service and customer-deployable Model-as-a-Service

▸ Full MLOps and Gen AI Ops lifecycle support from experimentation to production

▸ Model safety and trust with explainability, fairness, and guardrails

| Flexible and efficient inferencing | Connecting models to data | Agentic AI innovation | AI at scale |

**AI platform**

**Hardware Acceleration**

Physical   Virtual   Private Cloud   Public Cloud   Edge

32

## Registry (Models)

OpenShift AI users now have a centralized repository within the OpenShift AI platform designed to manage the lifecycle of machine learning models via the Registry.



**Registry is GA in OpenShift AI in the 3.0 release**

# GPUaaS

▶ Enables efficient management and allocation of GPU resources for a variety of AI workloads: workbenches, training/tuning, model serving

▶ Supports both whole and fractional GPU allocation

▶ Observability tools for resource optimization and to facilitate chargeback scenarios

# Why does it matter?

– *Improves Resource Utilization*: Reclaiming idle GPUs and optimizing allocation to reduce waste

– *Supports the complete AI Lifecycle*: Handling workloads from notebooks to model serving

– *Provides Visibility*: Offering metrics for both data scientists and administrators

# Responsible AI and governance with Red Hat AI

**Telemetry API**

**Evaluation API**

**Safety API**

**Llama Stack**



## AI Monitoring

Monitors tabular model inferences with customizable metrics for bias (outcome disparities) and drift (deployment vs. training data differences)

## AI Evaluation

Perform a huge variety of evaluation tasks over LLMs to understand and quantify their knowledge, capabilities, and behaviors

## Guardrails for AI

Customizable guardrails framework to moderate interactions between users and generative AI models, ensuring secure, compliant, and efficient operations

## RHOAI 3.0

*4Q 2025*

**Inference**
- Distributed inference – llm-d GA
- Additional validated models
- LLM Compressor images GA
- MaaS (Dev preview)

**Model development & Alignment**
- Modular & extensible solns for data ingestion, SDG, training & evaluation
- RAG – OpenAI APIs, local Milvus vector DB
- RAGAS evaluation framework
- Docling for document processing (Dev Preview)
- Feature Store GA

**Agentic**
- Llama Stack APIs
- AI Hub UI & gen AI studio
- MCP support & MCP Server access in gen AI studio

**AI platform**
- Model registry & catalog GA
- Platform metrics foundation – centralized observability infrastructure
- Expand platform support – ARM / Grace Hopper, IBM Z (GA), Power (TP)
- Hardware Profiles GA

## RHOAI NEXT

*1H 2026*

**Inference**
- Llm-d enhancements – eg. autoscaling, multi modal model support
- KEDA GA including UI
- OOTB ML Server runtime for expanded predictive model support
- MaaS GA – platform admin UI, external models mgmt

**Model development & alignment**
- Enhanced experiment tracking (MLflow)
- Docling support (GA)
- Prompt engineering capabilities
- Kubeflow Trainer v2 and checkpointing capabilities
- Enhanced model customization experience

**Agentic**
- MCP Server catalog & registry
- Agent catalog & registry
- Enhanced AI Playground – eg. model comparison, guardrails
- Knowledge source creation (Gen AI Studio UI)
- Inclusion of more Llama Stack providers across modules

**AI platform**
- Enhanced native observability –  model/agent monitoring, GPUaaS
- Integration of NVIDIA Nemo Guardrails
- Support Kubeflow Spark operator
- Workbenches 2.0 – next-gen workbench experience
- RH Trusted Software Supply Chain integration
- GPU partitioning GA (based on DRA)

# Thank you

Red Hat is the world's leading provider of
enterprise open source software solutions.
Award-winning support, training, and consulting
services make
Red Hat a trusted adviser to the Fortune 500.

linkedin.com/company/red-hat

youtube.com/user/RedHatVideos

facebook.com/redhatinc

twitter.com/RedHat

Red Hat

# What's New and Next in GPU-as-a-Service (GPUaaS)

## Highlights

| RHOAI 3.0 | RHOAI NEXT |
|---|---|

*Q4 2025*

*H1 2026*

- **Ray and Codeflare** updates:
  - **Architectural Simplification** - Complete CodeFlare Operator removal for cleaner, more maintainable platform
  - CodeFlare SDK providing a more aligned user experience by focusing on **RayJob**, while providing full access to RayClusters
- **Accelerator Slicing**: for all MIG-enabled devices, via the NVIDIA operator
- **Kueue**: supporting more workloads:
  - Ray Training Jobs
  - Training Operator-based Jobs
  - Inferences Services and Workbenches

- **Kubeflow Trainer v2** and **Red Hat Build of Kueue 1.2**:
  - major technology updates for workload scheduling
  - converged API for Training Jobs
- **GPU and Workload specific Observability**: based on RHOAI observability and accompanied by data scientists- and admin-centric views
- **Model Checkpointing**: to support better workload management, like Job-preemption and resume, to increase GPU utilization
- **Dynamic Resource Allocation (DRA)**: focus on this technology for (not only) GPU partitioning, long lead-time, as a lot of OCP work is to be done!
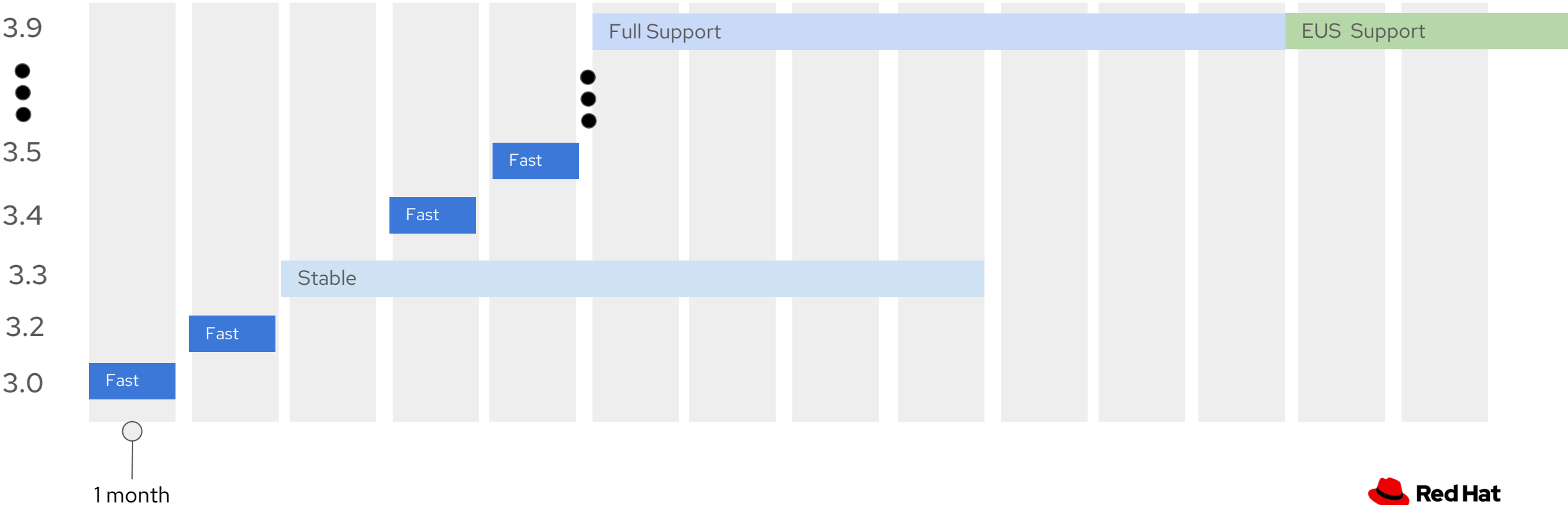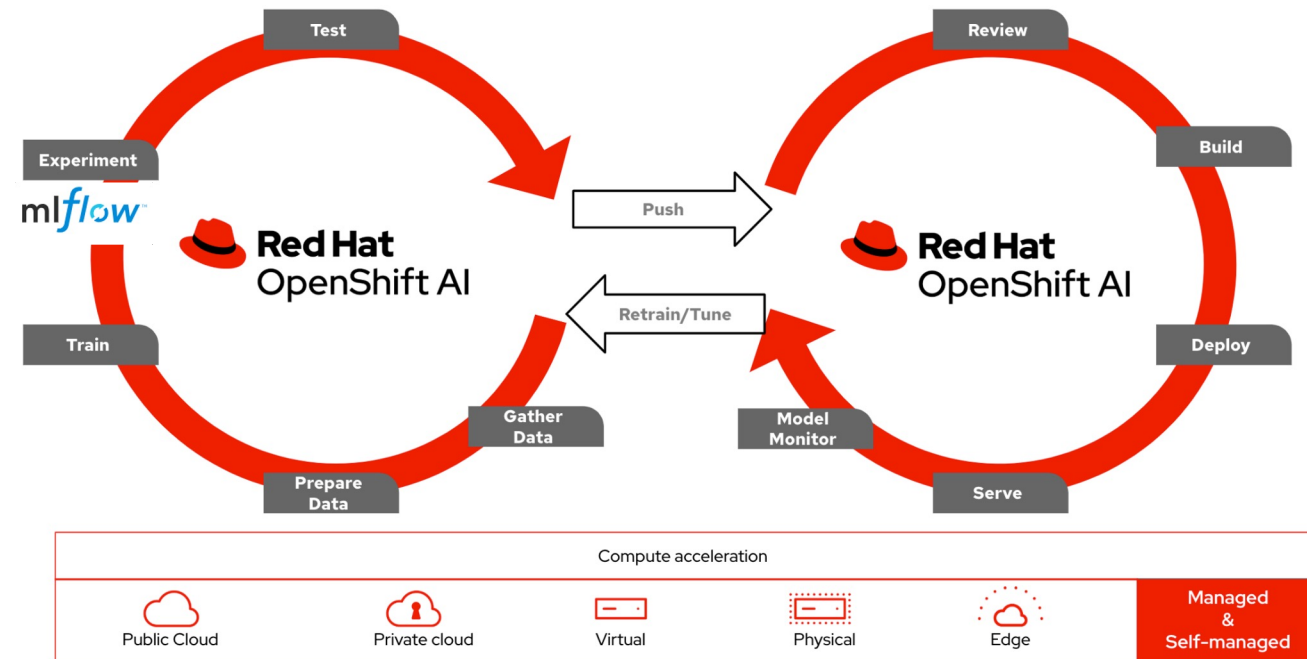
Red Hat

# What's New and What's Next in Experiment tracking

## MLFlow integration

**What's next?**

▸ Experiment tracking for more than just pipelines

▸ Integration of MLFlow on RH AI for inner loop experiment tracking

▸ More GenAI visualizations and metrics

▸ AIP and MLFlow integration

# Introducing multi-architecture

**IBM Power/ Z**

## Increasing flexibility and choice with an open source approach

▸ **In other to have a true hybrid platform, we are enhancing our release process to support multiple architectures. Starting with IBM Z and ARM. IBM Power shortly after.**

**Flexibility**

Access to cutting-edge open source innovations to keep up with a fast moving market.

**Choice**

Access to an open ecosystem of communities, technology providers, ISVs and customers.

Red Hat