

Agentic AI in Action

Red Hat & Intel Shaping the Future of Enterprise Al

Copenhagen 9 October

2025



Joachim Aertebjerg

Head of ISV Partnerships, EMEA Intel



Kristoffer Nærland

EMEA Ecosystem Sales Specialist Al Red Hat

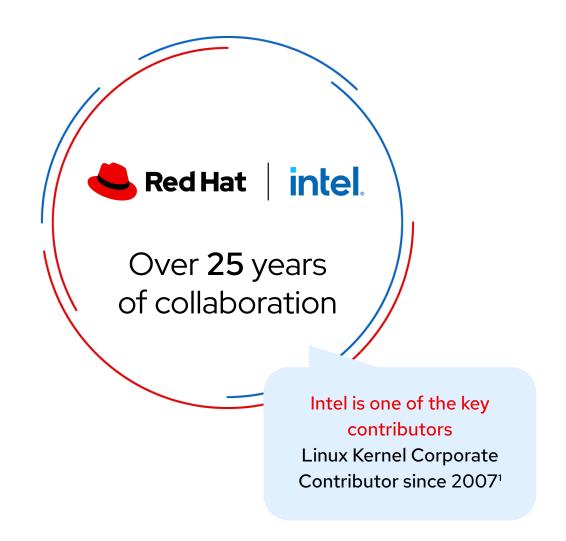


Intel - RH Partnership

Open source software: Intel is committed

Intel® has a long history with Linux®, actively participating in open source development and collaboration with the Linux community, to ensure hardware is well-supported and delivers optimal performance on Linux-based systems.

Intel contributes to more than 100 different open source projects, from the Linux kernel to cloud orchestration and plugins for Kubernetes.



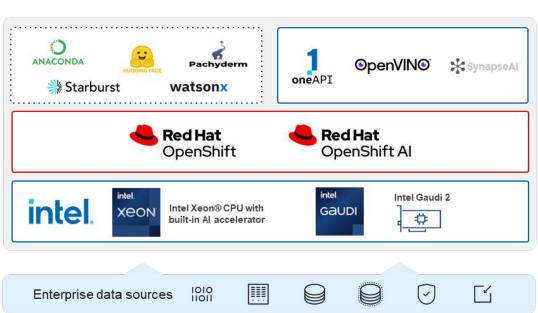
Real Customer Example: Al Sweden







Data Scientists



- Collaborating to deliver AI solutions
- Deeper, product collaboration focused on customer enablement with OpenShift AI, Intel Xeon, Gaudi 2 and the Intel AI Suite
- Testing, validation, and proof of concepts
- Receive support for building Al applications



Intel's Al Strategy and Capabilities

Bringing Al Everywhere

Intel's AI Strategy



AI PC Broadest AI SW Ecosystem



ENTERPRISE AI & EDGE AI Open Standard, "Ready to Use"



DATA CENTER AI
AI Open, Scalable Systems & Reference Arch











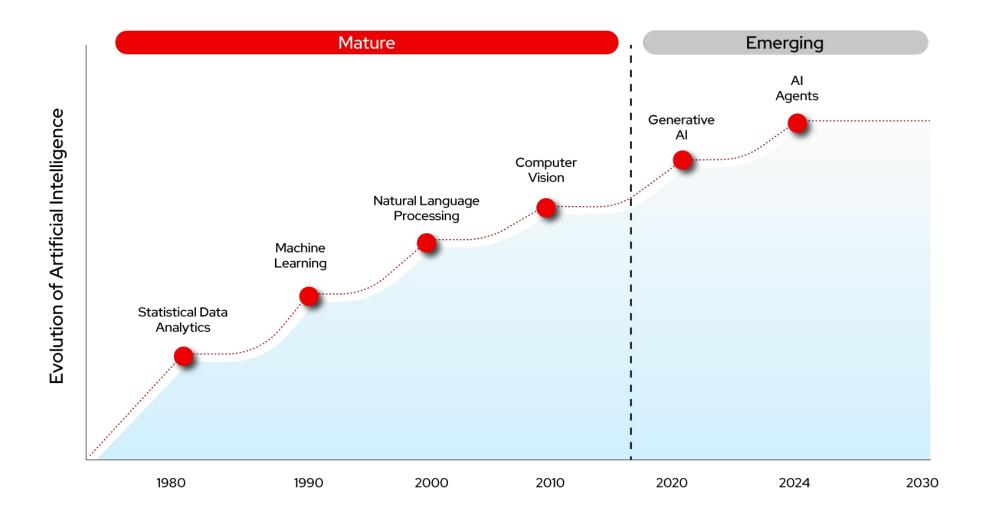




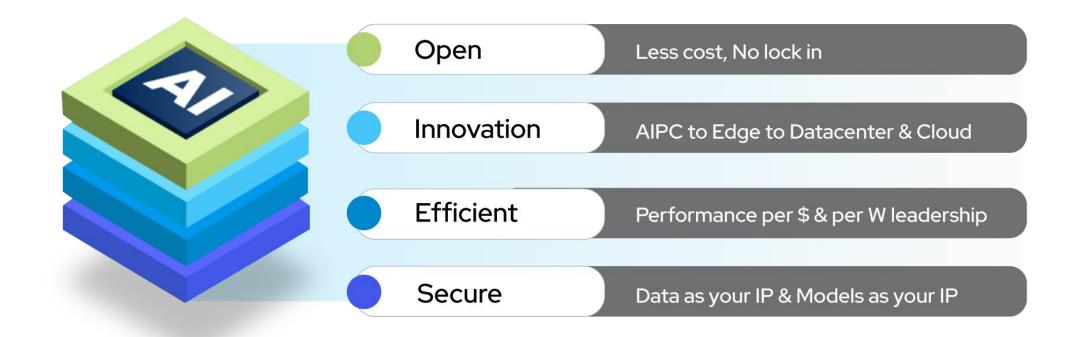




Evolution of Al Applications in Enterprise Use Cases

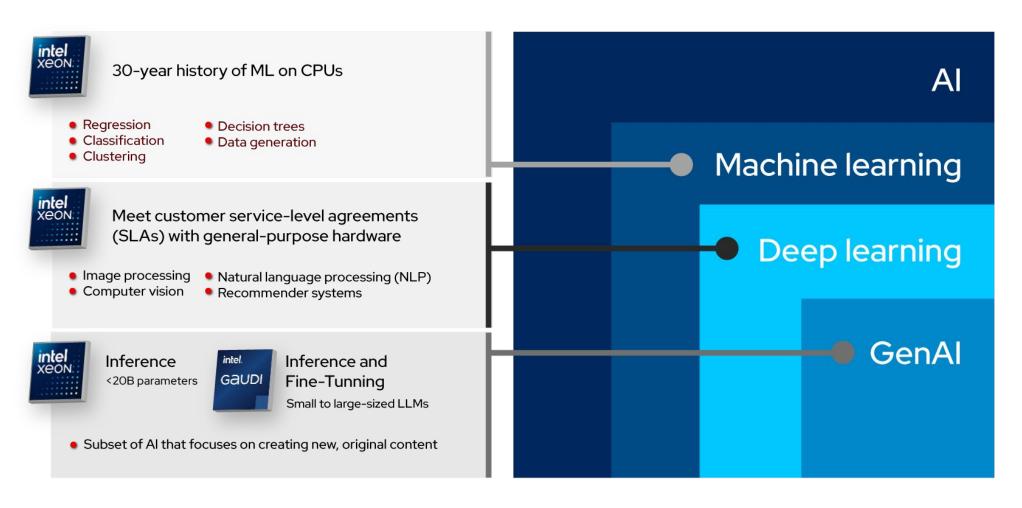


Intel's Al Strategy



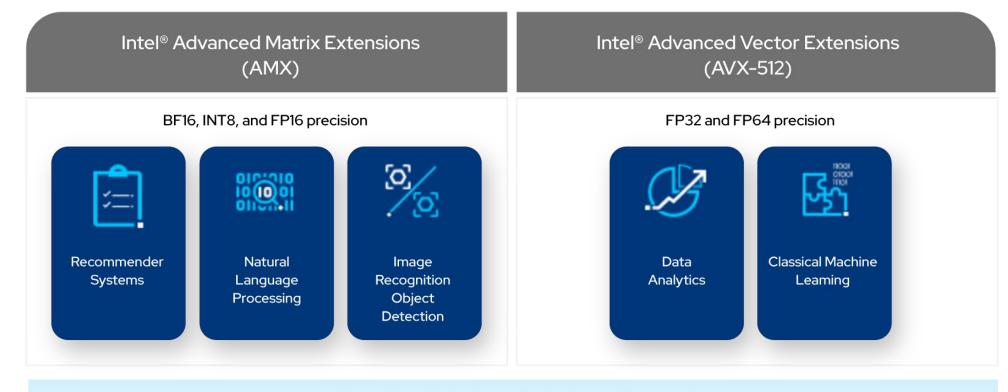
The Al Hierarchy: Mapping ML, DL, and GenAl with Intel

Discover how Intel® processors fuel AI workloads across inference, training, and next-generation GenAI applications





Intel® AMX Accelerates **DEEP LEARNING** Use Cases



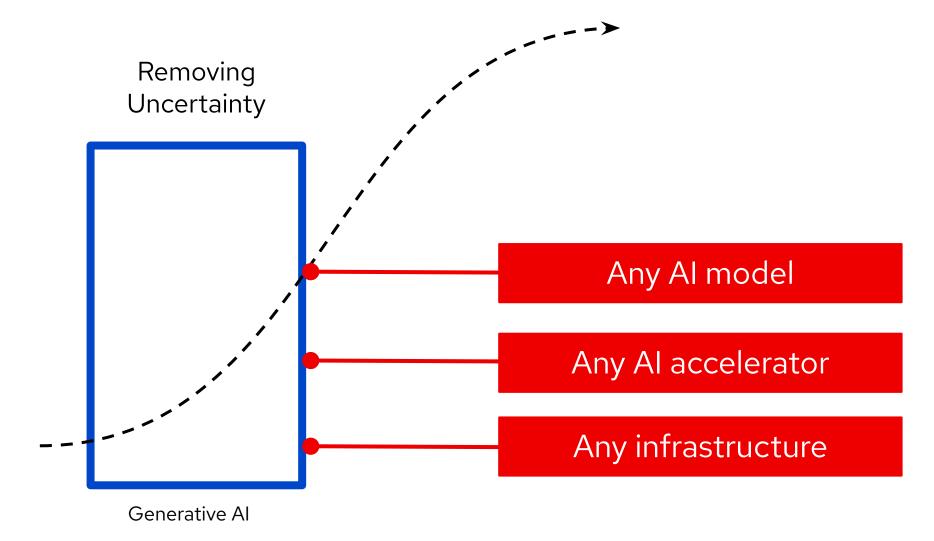
Many DL workloads are "mixed precision" and 5th Gen Xeon can seamlessly transition between AMX and AVX-512 as needed

intel ai Al Gold Deck **Public**



Red Hat's Al Strategy and Capabilities

Red Hat AI - Enabling AI Success





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

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

Simplified and consistent experience for connecting models to data

Flexibility and consistency when scaling Al across the hybrid cloud

Accelerate Agentic AI delivery and stay at the forefront of innovation











Trusted, Consistent and Comprehensive foundation





intel.

Hardware Acceleration













Virtual



Private Cloud



Public Cloud



Edge



Intel Gaudi Al Accelerators

Intel® Gaudi® 3 Al Accelerator: Al Inferencing

Price Performance Advantage

Up to

43%

Higher throughput

(tokens per second)

on IBM Granite-3.1-8B-Instruct

vs. leading GPU competitor with small context sizes

120%

More cost efficient

(tokens per dollar)

on Mixtral-8x7B-Instruct-v0.1

vs. leading GPU competitor with long input and short output sizes

More cost efficient

(tokens per dollar)

on Llama-3.1-405B-Instruct-FP8

vs. leading GPU competitor with large context sizes









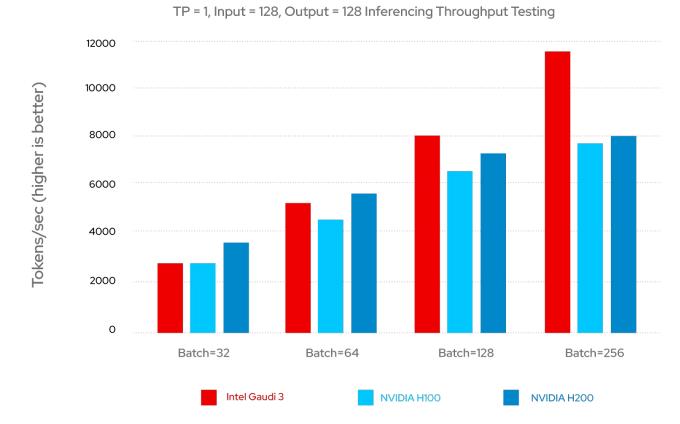
Up to 43% higher

throughput than NVIDIA H200

Up to 52% higher

throughput than NVIDIA H100

For lightweight Al Use Cases



Granite-3.1-8B-Instruct

Reported numbers are inferencing results for IBM Granite-3.1-8B-Instruct on Intel® Gaudi® 3 vs NVIDI H100 GPU and NVIDIA H200 GPU. Refer to this link for the latest published Gaudi3 performance https://www.intel.com/content/www/us/en/developer/platform/gaudi/model-performance.html

Pricing estimates based on publicly available information and Intel internal analysis.



^{*}Source: NV H100 and H200 comparisons based on Signal65 Lab Insight: Intel Gaudi 3 Accelerates AI at Scale on IBM Cloud. April 2025.

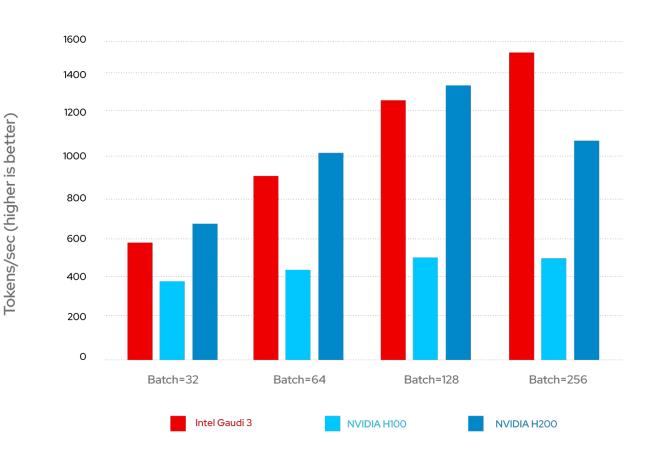
Up to 36% higher

throughput than NVIDIA H200

Up to 200% higher

throughput than NVIDIA H100

For Large Al Workloads



^{*}Source: NV H100 and H200 comparisons based on Signal65 Lab Insight: Intel Gaudi 3 Accelerates AI at Scale on IBM Cloud. April 2025.

Reported numbers are inferencing results for IBM Granite-3.1-8B-Instruct on Intel® Gaudi® 3 vs NVIDI H100 GPU and NVIDIA H200 GPU. Refer to this link for the latest published Gaudi3 performance https://www.intel.com/content/www/us/en/developer/platform/gaudi/model-performance.html

Pricing estimates based on publicly available information and Intel internal analysis.



Intel Xeon Processors



Intel® Xeon® 6 Processor

1.9x

higher performance per watt at a typical 40% server utilization vs. prior generation

> Designed for Efficiency

2.5x

higher HPC performance vs. prior generation

Significant Performance Leaps 5.5x

higher Al Inferencing performance vs. AMD EPYC

> Unmatched Performance

Resolve Customer Queries Faster with More Concurrent Users in Your LLMs and Agents

Get superior performance for batch, real-time inference, and training for small and medium language models with Intel® Xeon® processors.

Use your CPU for cost-effective model updates.



Large language models (LLMs)

Intel Xeon 6 vs. AMD EPYC Turin

Intel Xeon 6 vs. 5th Gen Intel Xeon

5th Gen Intel Xeon vs.

Llama2-7B

Up to

1.38x

higher throughput

with Intel Xeon 6980P vs. AMD EPYC 9965'

GPTJ-6B

Up to

2x

Higher performance

Intel Xeon 6980P vs. Intel Xeon 8592+2 Llama-13B

Up to

2x

Higher performance

Intel Xeon 6980P vs. Intel Xeon 8592+2 Llama2-7B

Up to

2.3x

Higher training performance

Intel Xeon 6980P vs. Intel Xeon 8592+3' 3rd Gen Intel Xeon

Llama2-13B

Up to

2.1x

real-time inference performance speedup

> 5th Gen Intel Xeon vs. 3rd Gen Intel Xeon4

Intel Confidential Computing

Confidential Computing and Post Quantum Crypto for Information & Data Security

Intel® Software Guard Extensions (Intel® SGX)

Smallest Trust Boundary - Confidential data access is restricted to attested application code

Intel® Trust Domain Extensions (Intel® TDX)

Virtual machine isolation from cloud stack, admins, and other tenants

Post Quantum

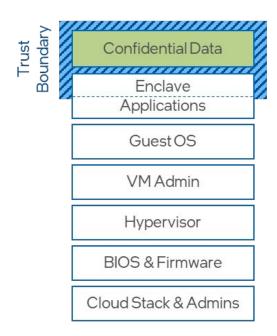
Intel adds Quantum attack protection while providing 1.89 Tb IPsec throughput.

Performant Post-Quantum Cryptography (PQC) leveraging the Intel NetSec Accelerator and Arqit SKA-Platform™ for PQC.



App Isolation

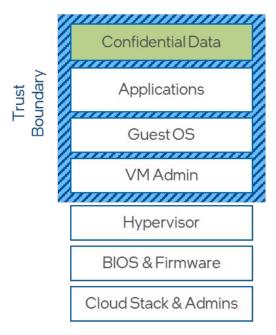
Intel®SGX



Smallest trust boundary for greatest data protection & code integrity

VM Isolation

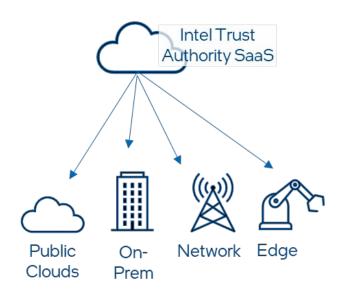
Intel® TDX



Most straightforward path to greater security for legacy apps

Trust Services

Intel® Tiber™ Trust Authority



Uniform, independent attestation of trustworthy environments

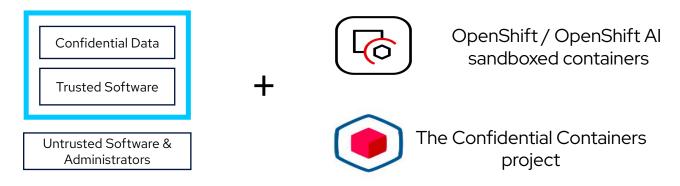
Founded on Intel's Security-First Development & Lifecycle Support



Confidential AI Helps Protect Data & Models In-Use

Utilizing Confidential Computing for Containers with Intel TDX

Hardware-Based Protection of Data In-Use With Intel Trusted Domain Extensions (TDX)



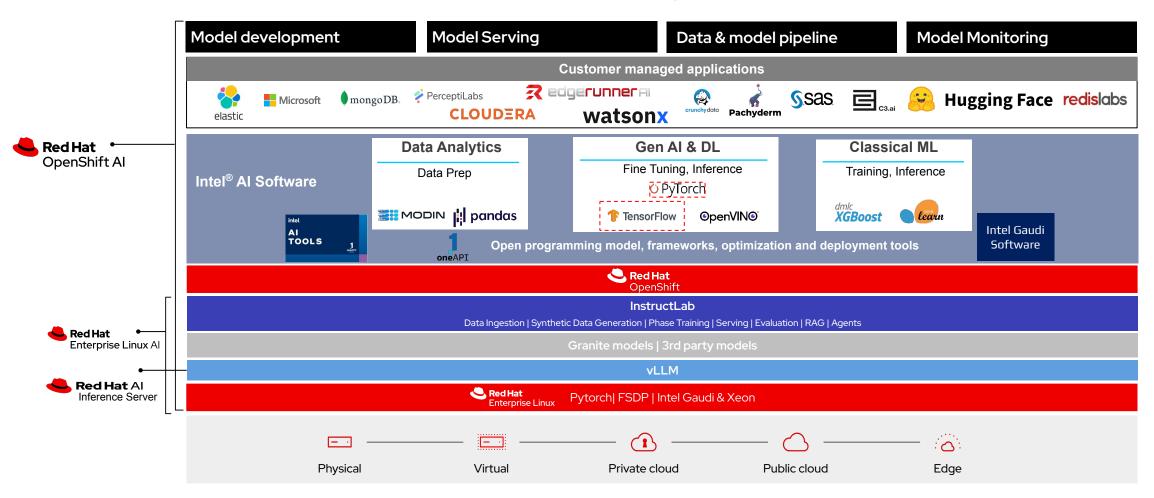
Confidential Computing is about protecting data in-use. You do not have to trust the system admins of the providers any longer.



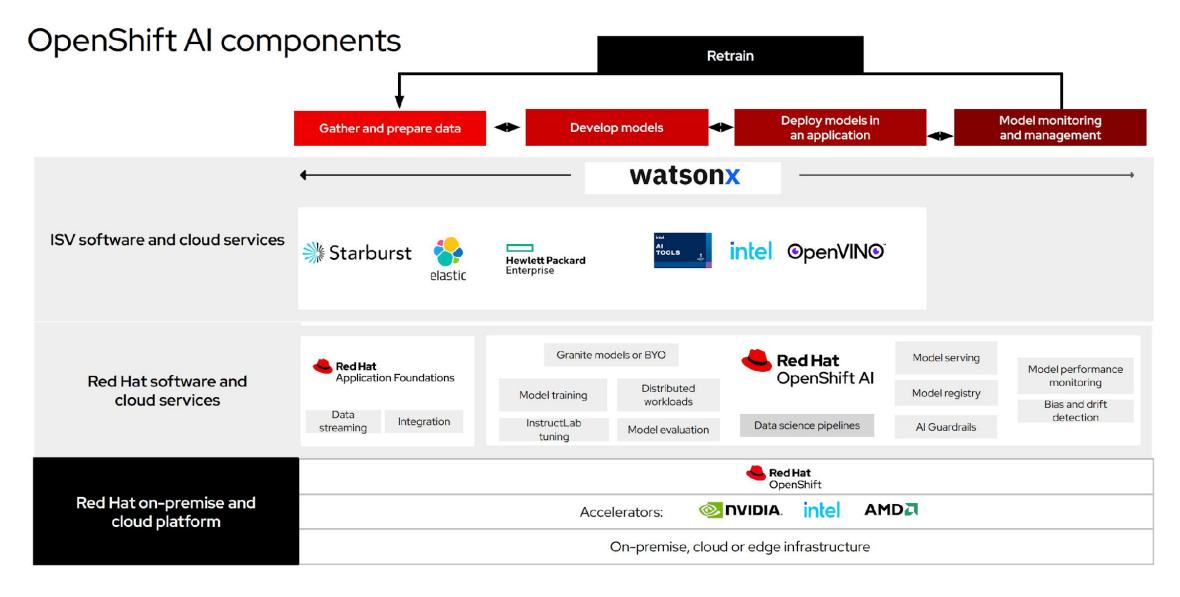
Intel Al Software

Red Hat AI with Intel AI platform

Generative AI and MLOps capabilities for building flexible, trusted AI solutions at scale



Red Hat Al Platform





Red Hat AI the inference engine for the hybrid cloud

vLLM supports the key models on the key hardware accelerators



















Molmo

Phi

Nemotron

















Spyre







Virtual



Private Cloud



Public Cloud

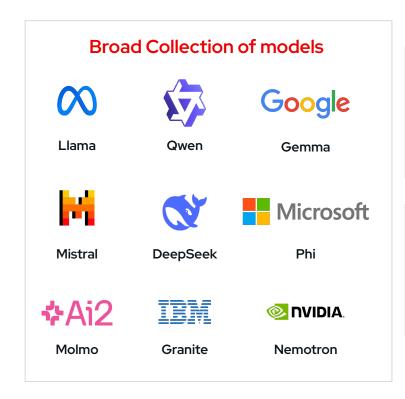


Edge



Red Hat AI repository on Hugging Face

A collection of third-party validated and optimized large language models



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



Intro to Agentic Al

... Imagine you're the CEO and need to provide coffee to your employees



Expense Starbucks for all

Easy, quick, consistent

Linear cost & complexity

Sort of OpenAl, Gemini, etc.

like...



Machine for everyone

Easy, quick, consistent

Economy of scale

Models-as-a-Service



Full time Professional Barista

Great for short time periods

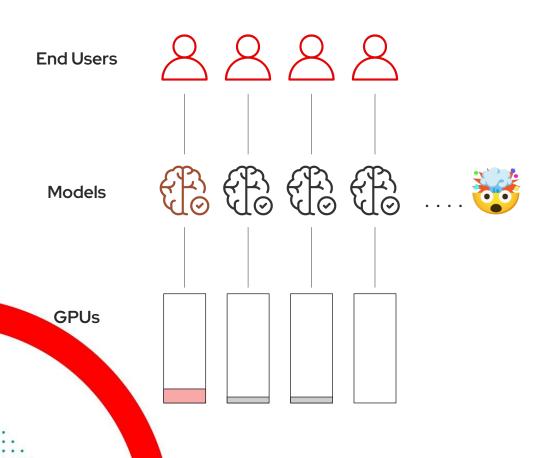
Expenses high but limited

Al Infrastructure (aka GPUs) as-a-Service



CPUs & especially GPUs

Infrastructure as a Service can be costly



Self-Service is good for plentiful resources & small teams

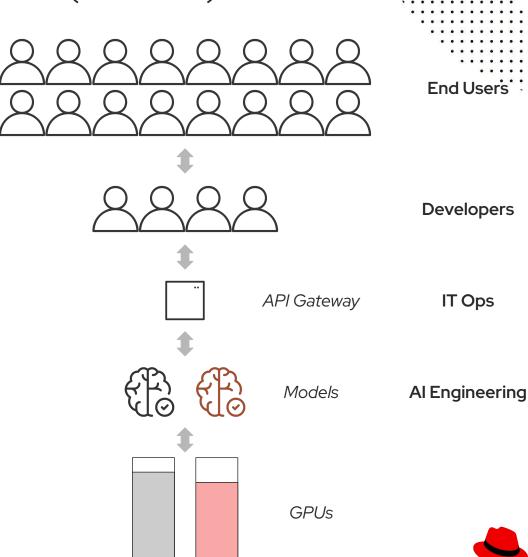
- Throwing GPUs at the problem is risky
- Few people know how to use them correctly
- Leads to duplication and underutilization
- Leads to high costs
- Most people want an LLM endpoint, not a GPU



Models as a Service (MaaS)

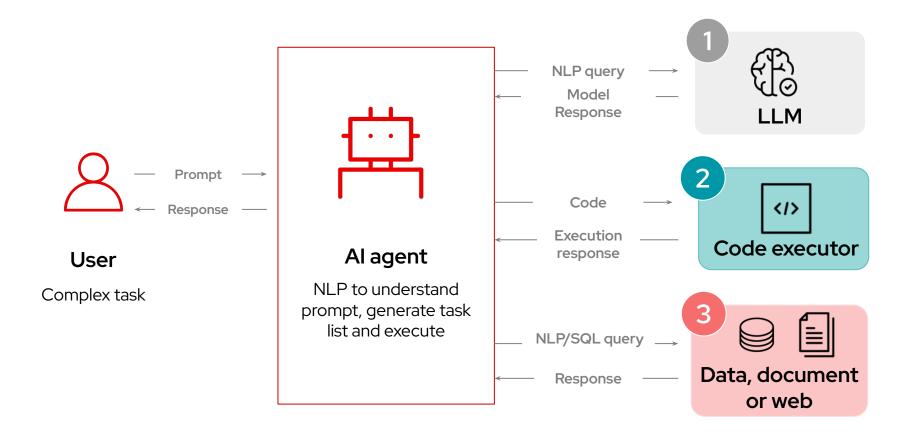
Offering AI models as the service to a larger audience

- IT serves common models centrally
 - Generative AI focus, applicable to any model
 - Centralized pool of hardware
 - Platform Engineering for Al
 - Al management (versioning, regression testing, etc)
- Models available through API Gateway
- Developers consume models, build Al applications
 - For end users (private assistants, etc)
 - To improve products or services through Al
- Shared Resources business model keeps costs down



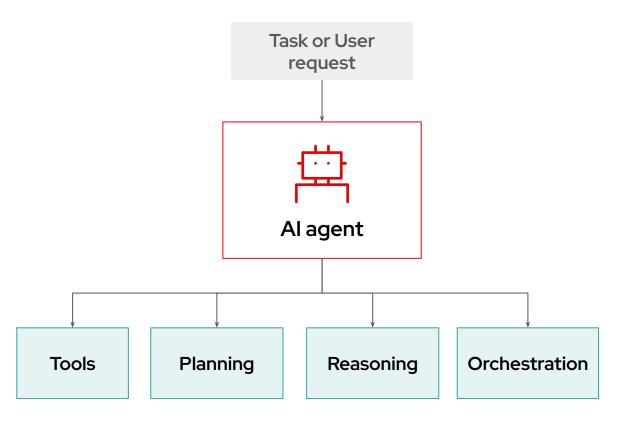
Al agents integrate models, functions & tools

Gen Al Models, Predictive Al Models, Code Functions, Search & more





The components of an Al Agent system



- ► **Tool Utilization:** Leverages external tools to gather data and perform tasks.
- Planning and Execution: Develops and executes multistep plans to achieve goals autonomously.
- Reasoning: Applies logic and contextual understanding to make informed decisions.
- Orchestration: Coordinates actions, tools, and agents to dynamically adjust and complete tasks.
- Communication protocols: enables the connections between the components.



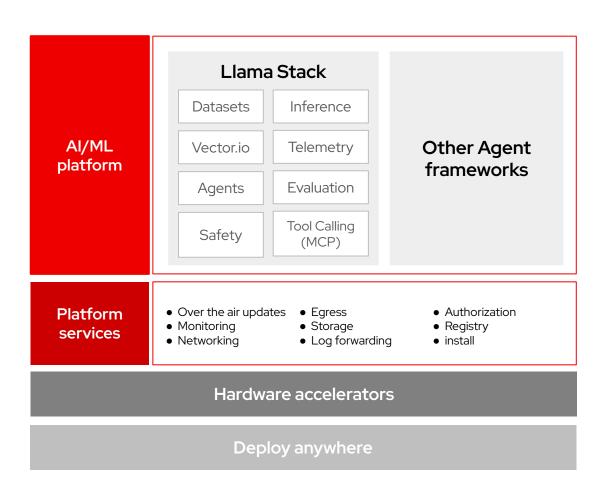
Red Hat AI provides an agile, stable foundation to accelerate the development and deployment of AI agentic workflows.

- Offers built-in agent frameworks with Llama Stack, and standardized communication protocols (MCP).
- Provides the flexibility to integrate preferred tools like LangChain and Crew AI.
- Allows running and managing agents as microservices.
- Simplifies production deployment by managing LLM serving and scaling.





A modular approach to building AI agents



Red Hat Al allows to:

- Build agents using Llama Stack's native capabilities and implementations.
- Bring compatible Llama Stack implementations to OpenShift Al.
- 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.



RAG vs RAG + MCP

Context Token-limited Session memory, multi-step

Integration Custom APIs Standardized plug-in tools

Multi-agent coordination **Agents** Hard to scale

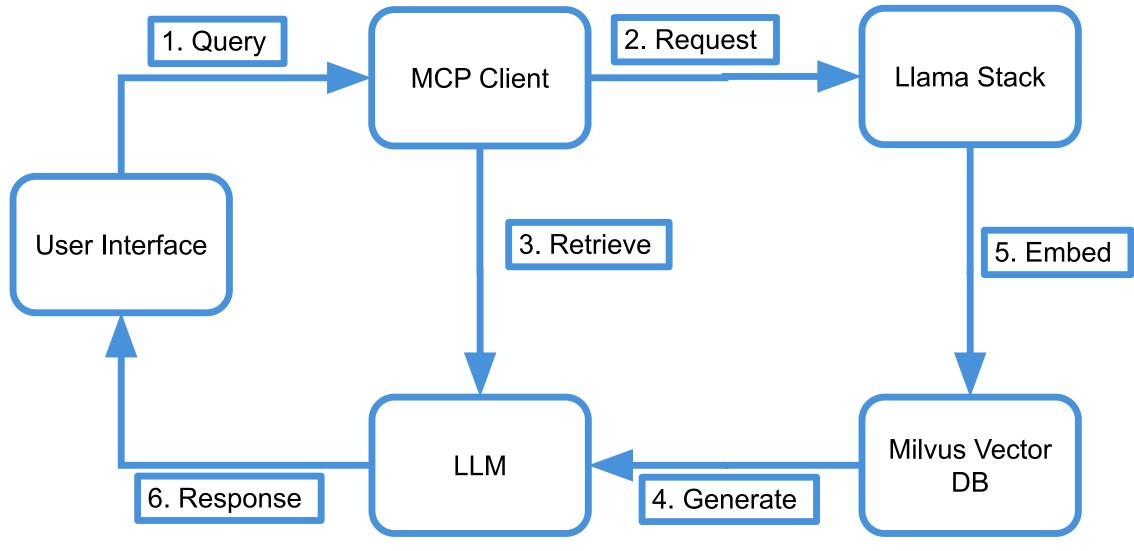
Deployment Manual, standalone Cloud-ready, replicable

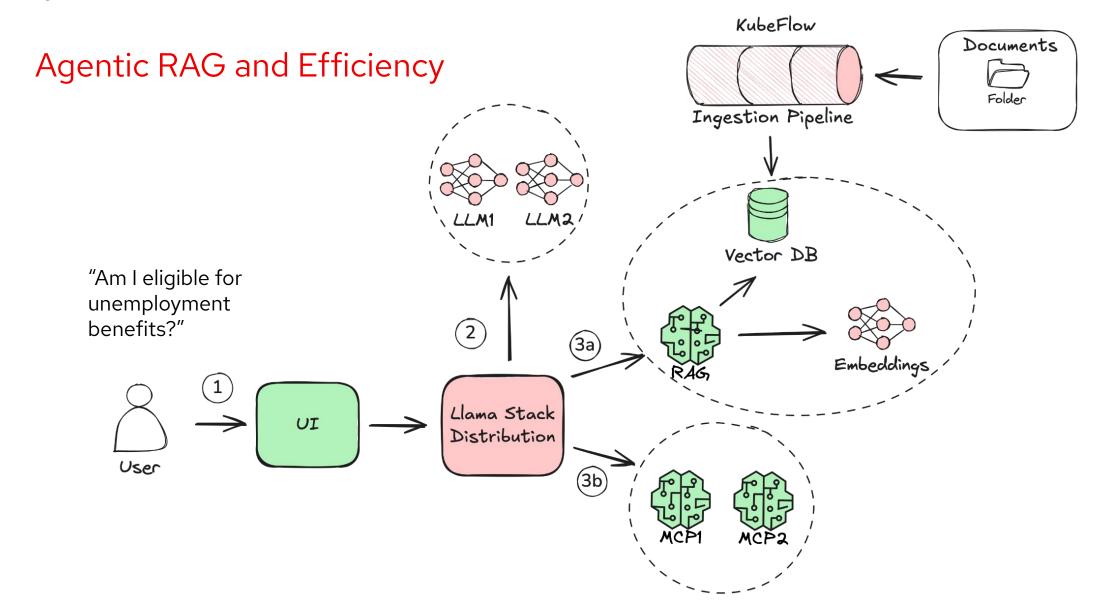
Monitoring Centralized logging Isolated errors

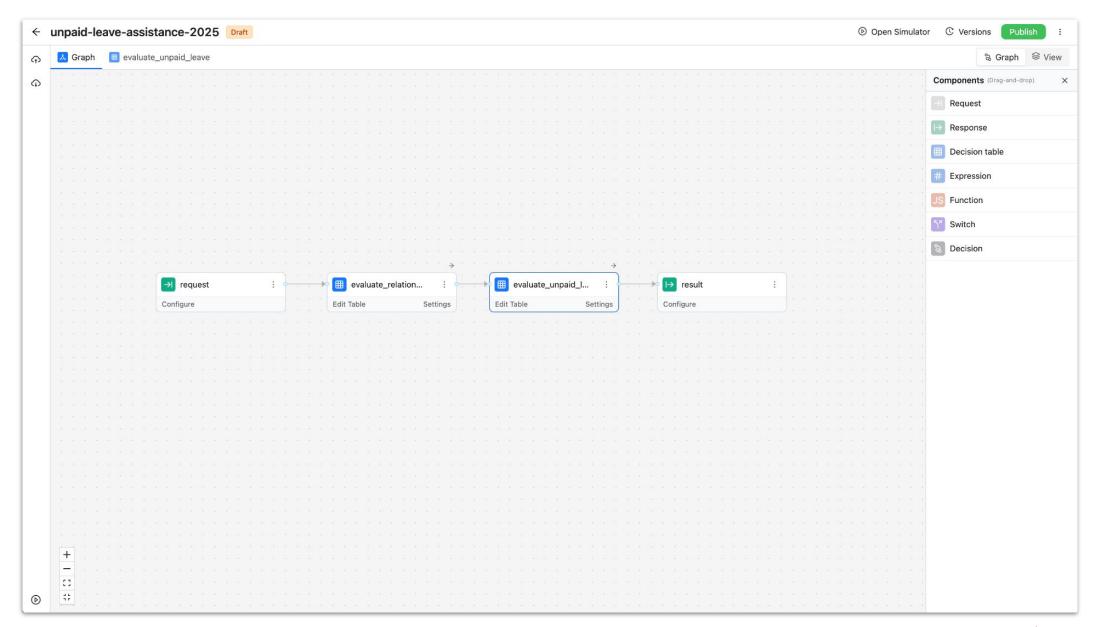
Ad hoc calls Batching, caching Resources

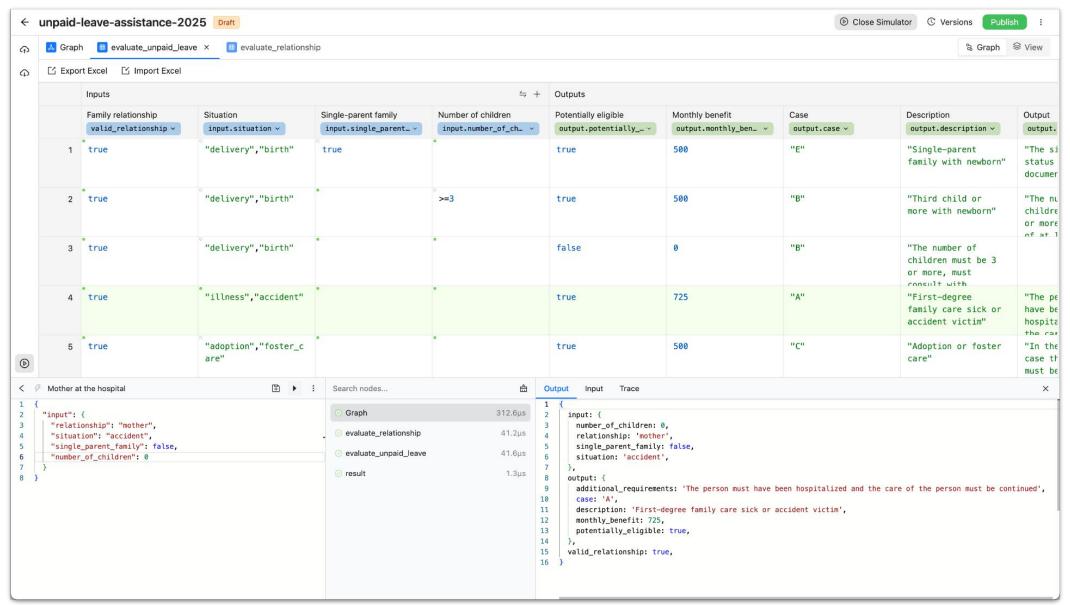
Agentic Al Demo

Agentic Al Demo Architecture

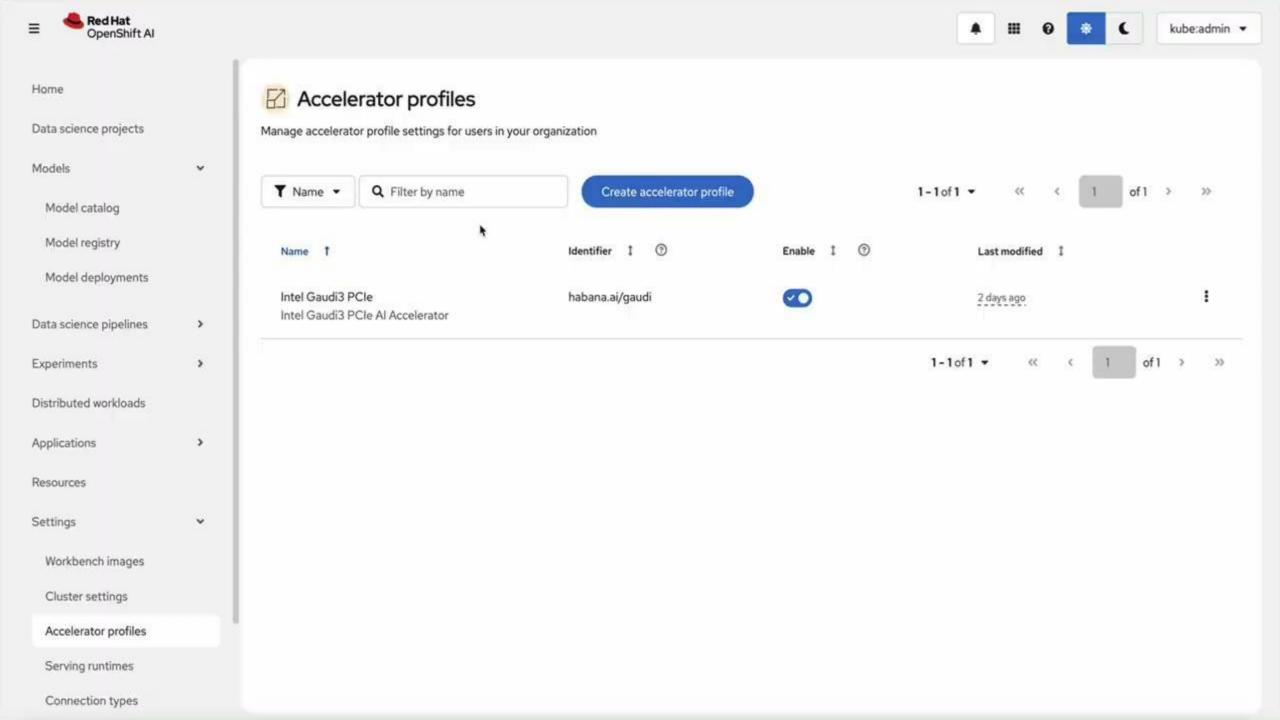


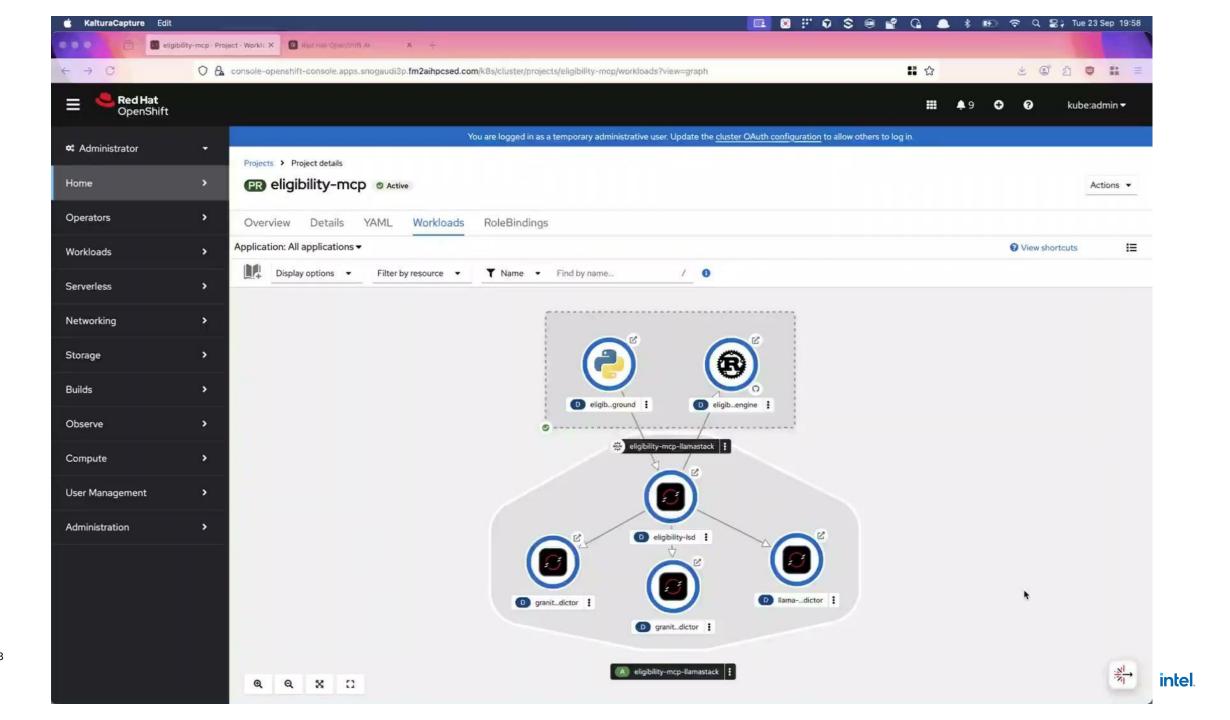


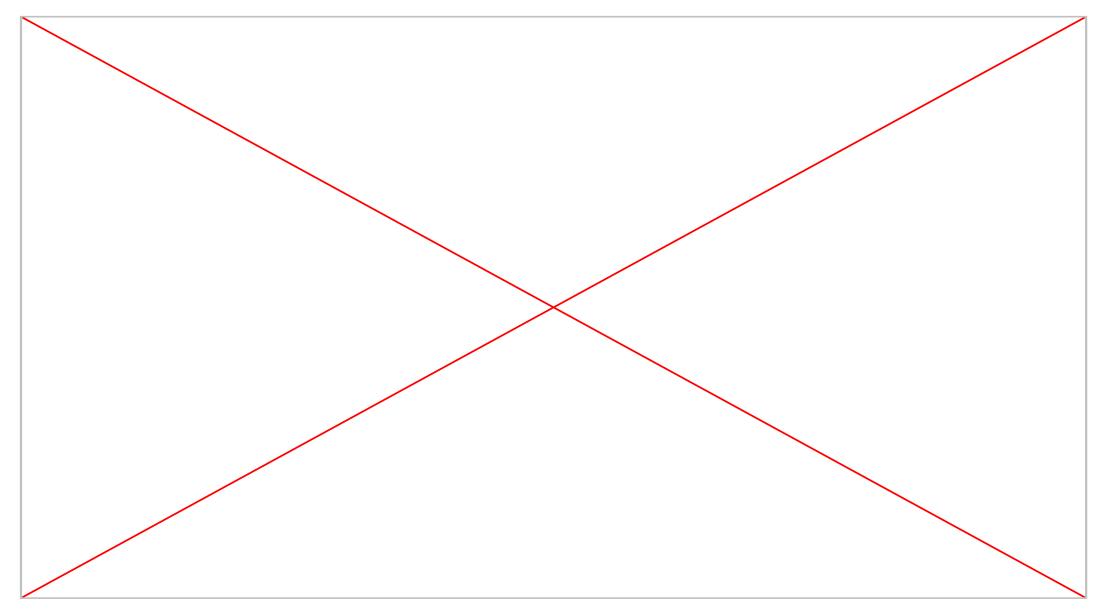












<u>Demo - OpenShift Al Model Serving with Gaudi and Xeon</u>

Demo Screenshots for Distributing the Deck (Unskip screenshot slides when generating PDF)

Agentic Al Demo - Admin Video Slides









kube:admin 🔻



Data science projects

Models

>

>

Data science pipelines

Experiments >

Distributed workloads

Applications >

Resources

Settings V

Workbench images

Cluster settings

Accelerator profiles

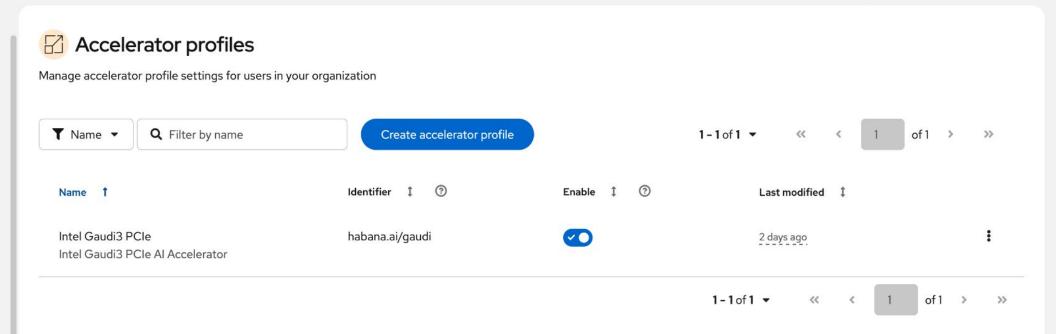
Serving runtimes

Connection types

Storage classes

Model registry settings

User management











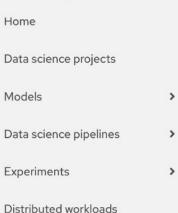








kube:admin ▼



Applications >

Resources

Settings

Workbench images

Cluster settings

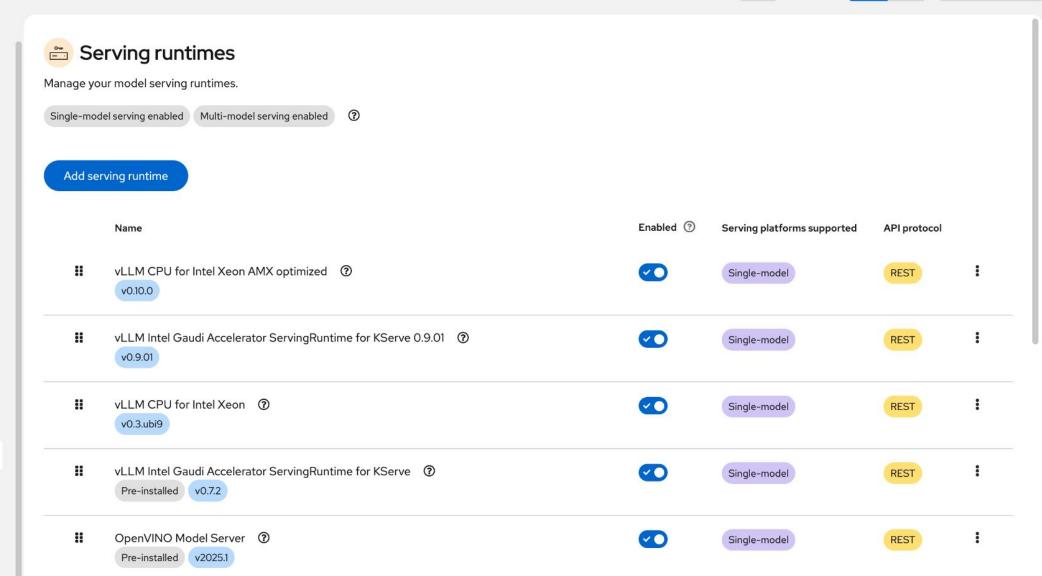
Accelerator profiles

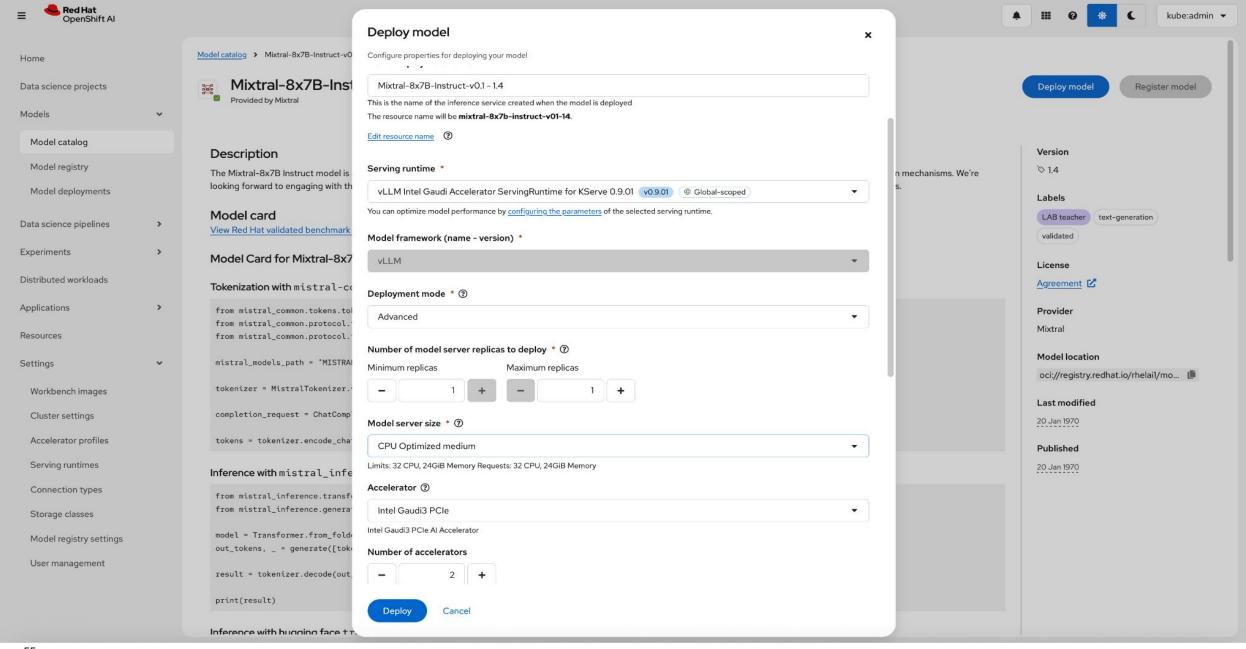
Serving runtimes

Connection types

Storage classes

Model registry settings





Workbenches Overview Pipelines Models Cluster storage Connections Permissions Settings



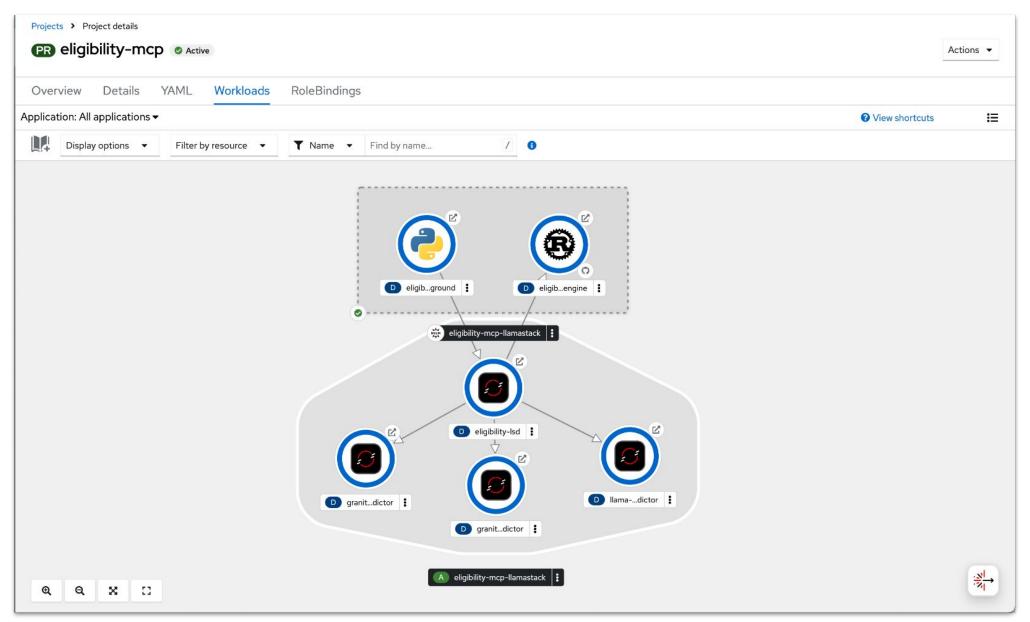
Models and model servers ②

Deploy model

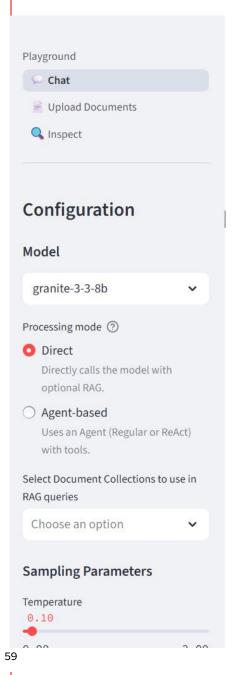
Single-model serving enabled

	Model deployment name 1	Serving runtime	Inference endpoint	API protocol	Status	
•	Granite 3.3 2B ②	vLLM CPU for Intel Xeon AMX optimized	Internal endpoint details	REST	•	:
	Framework	vLLM				
	Model server replicas	1				
	Model server size Custom 64 CPUs, 24GiB Memory requested 64 CPUs, 24GiB Memory limit					
	Accelerator Token authentication	No accelerator selected A Tokens disabled				
>	Granite 3.3 8B ②	vLLM Intel Gaudi Accelerator ServingRuntime for KServe 0.9.01	Internal endpoint details	REST	•	•

Agentic Al Demo - Agentic Al - Application Video Slides



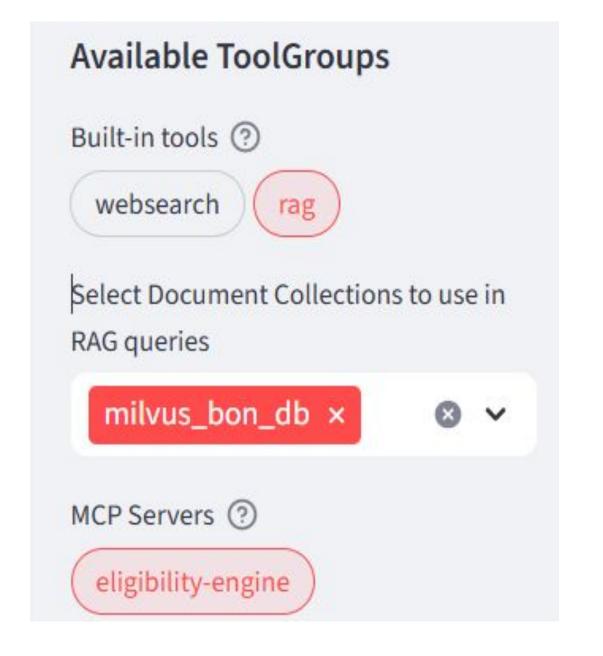


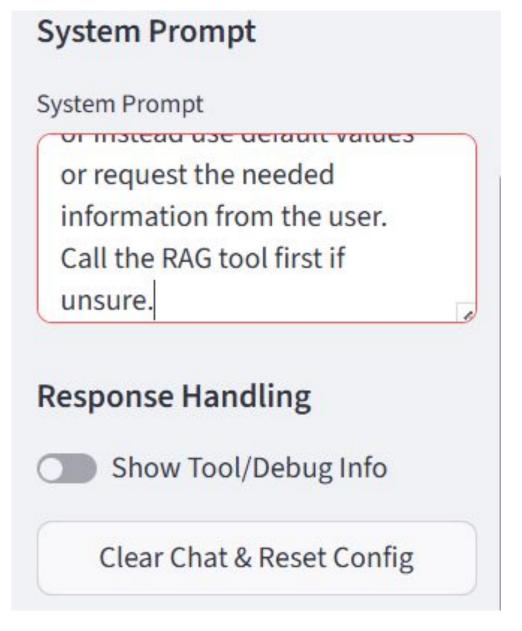




How can I help you?

Ask a question...





Chat

- How can I help you?
- My mother had an accident and she's at the hospital. I have to take care of her, can I get access to the unpaid leave aid?
- Yes, you may be eligible for unpaid leave assistance. According to the regulations, case A (Sick family care) provides 725€. To confirm your eligibility, I need to gather some details:
 - 1. Relationship to the person in need of care: 'mother'
 - 2. Situation: 'illness'
 - 3. Is your family a single-parent family? (true/false)
 - 4. Number of children involved in care (if more than one): '1'

Please provide these details so I can evaluate your case.



Q&A

Apply for a free Gaudi 3 Proof of Concept in 30 seconds

Choose your GenAl or Virtualization PoC:

- □ Building Inference, RAG, AgenticAl, Model-as-a-Service, and other Al Use Cases with Intel Gaudi and Xeon
- ☐ Optimize finetuning with intel Gaudi

Why work with Intel + Red Hat?:

Benefit from access to free highly qualified experts from Red Hat and Intel and free access to the latest hardware to build your Al use case / application.

If selected, a Intel / Red Hat representative will contact you via email.



Come visit the Intel and Red Hat booths to learn more!





Thank you



linkedin.com/company/red-hat



facebook.com/redhatinc



youtube.com/user/RedHatVideos



twitter.com/RedHat



