



Connect

Agentic AI in Action

Red Hat & Intel Shaping the Future of Enterprise AI

Helsinki

October 23rd, 2025



Danai Skournetou

EMEA Strategic Sales
Intel



Andreas Bergqvist

EMEA Sales Specialist AI
Red Hat

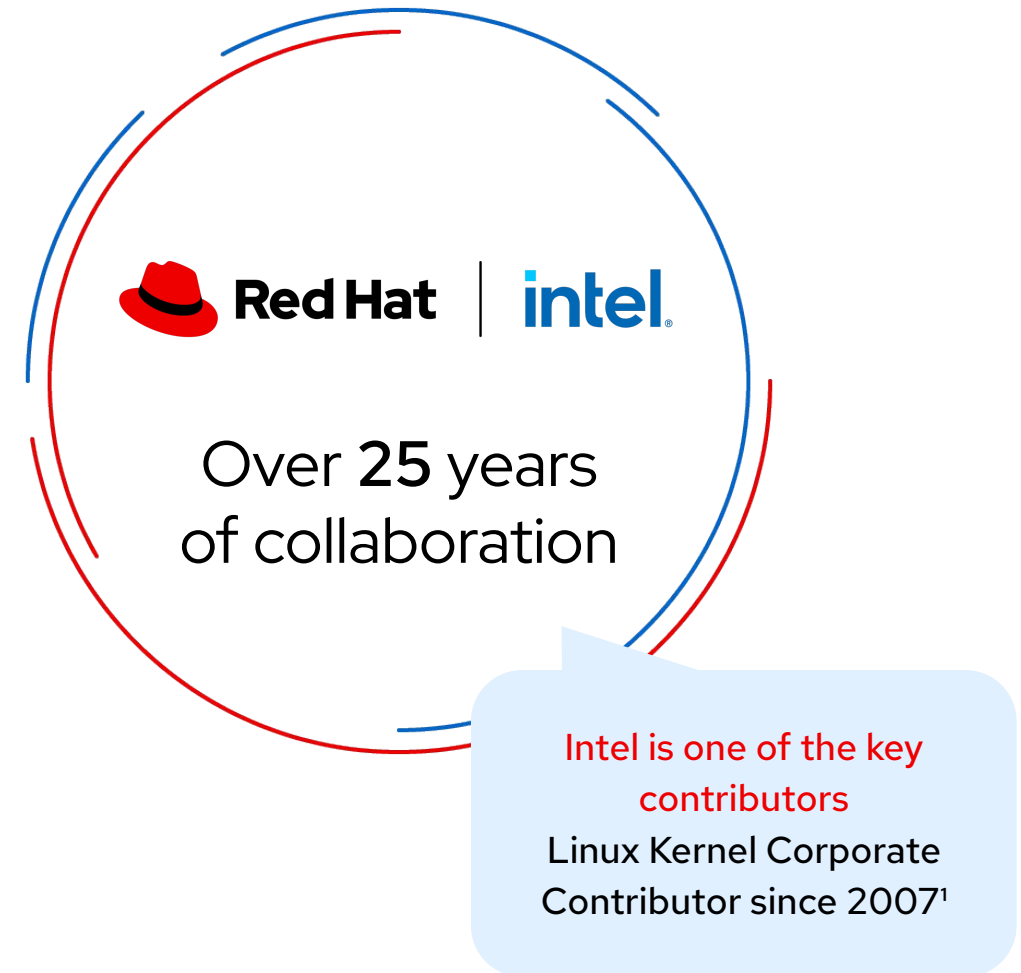


Intel – RH Partnership

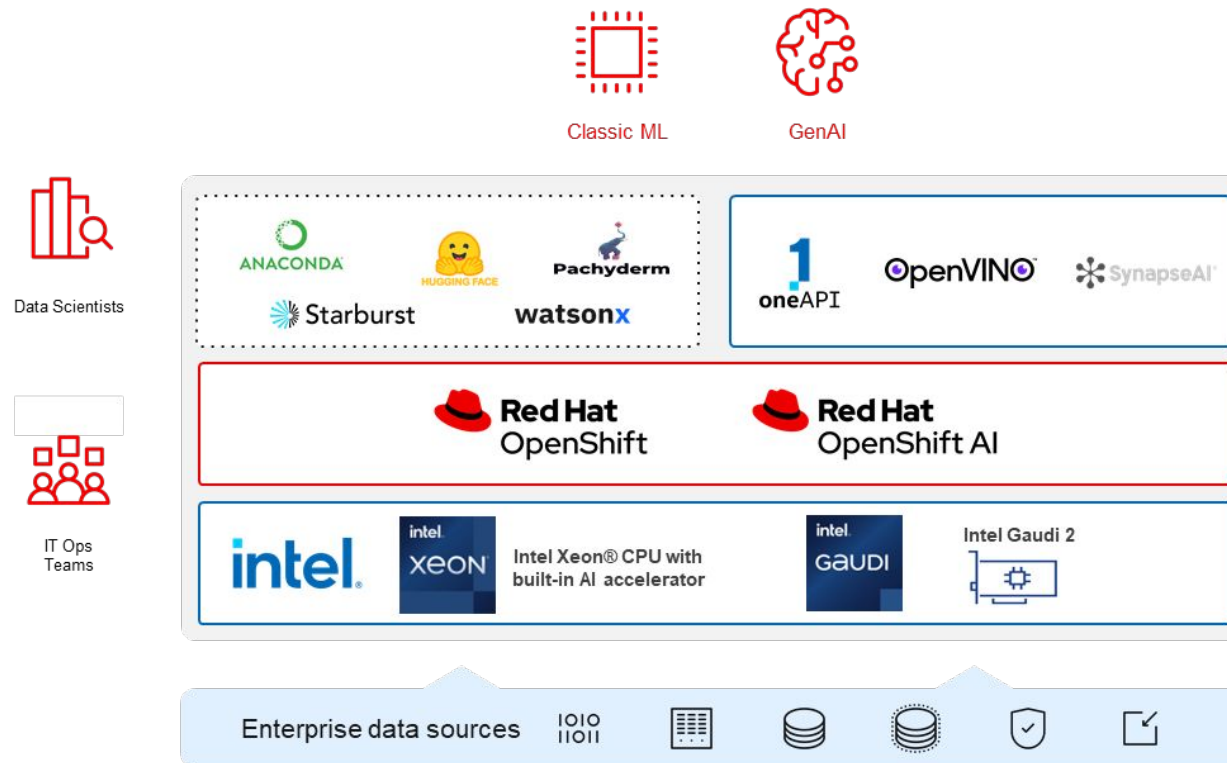
Open source software: Intel is committed

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: AI Sweden



- ▶ 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 AI applications

Intel's AI Strategy and Capabilities

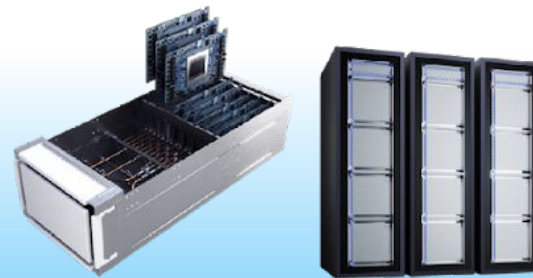
Bringing AI Everywhere

Intel's AI Strategy



AI PC Node
AI Developer Productivity & Light
Inference

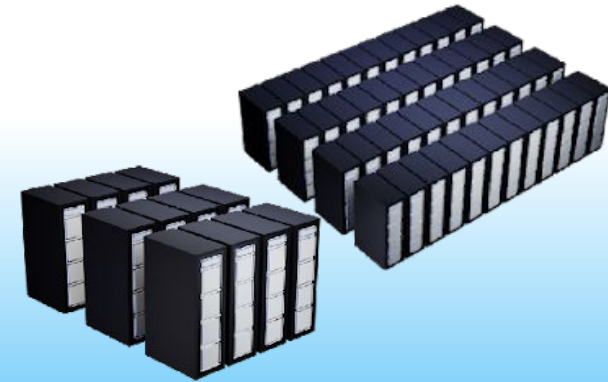
AI PC
Broadest AI SW Ecosystem



Node
Fine-tuning,
Inference

Cluster
Light Training, Tuning, Peak
Inference

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



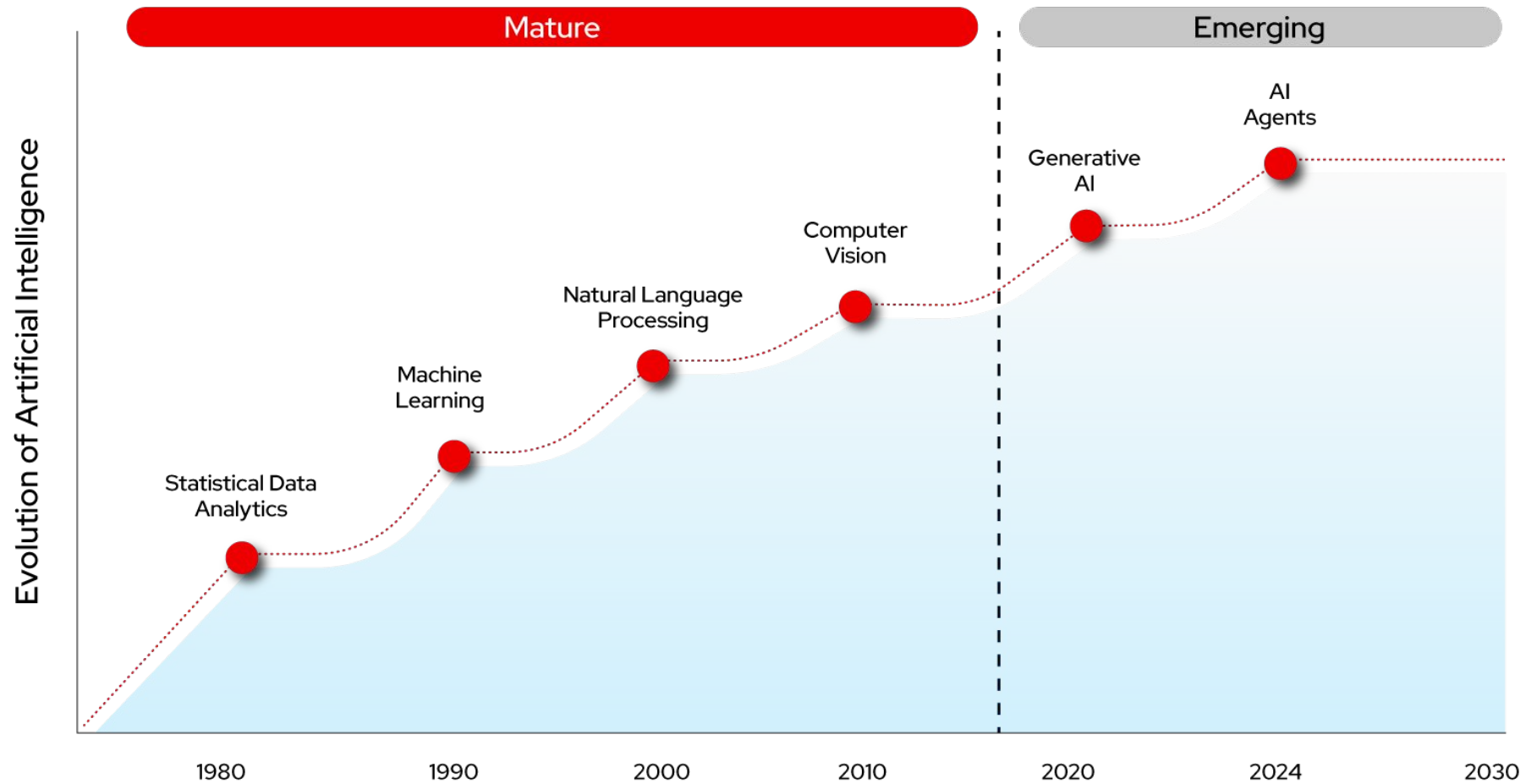
Super Cluster
Training, Tuning, Peak
Inference

Mega Cluster
Large Scale Training
& Inference

DATA CENTER AI
AI Open, Scalable Systems & Reference Arch



Evolution of AI Applications in Enterprise Use Cases



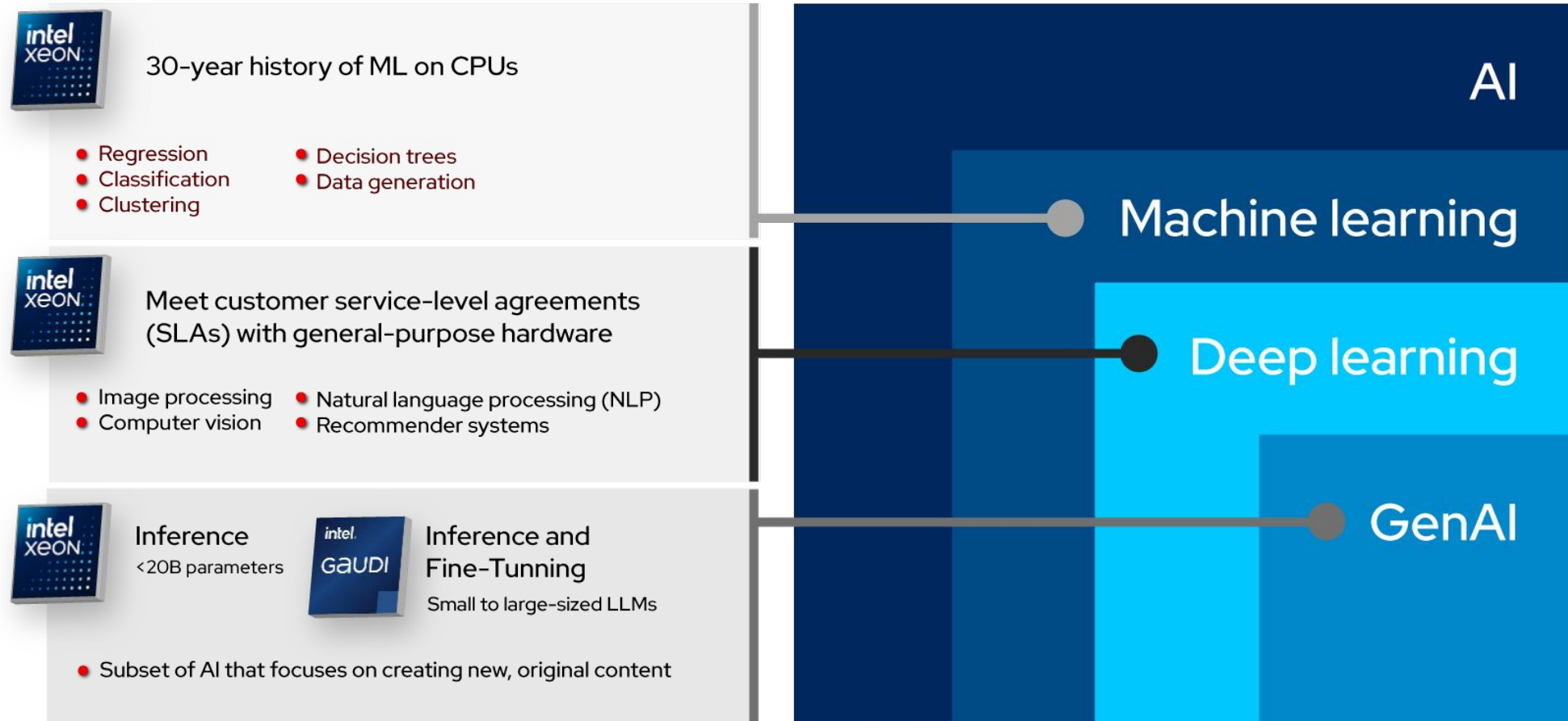
Intel's AI Strategy



- Open** Less cost, No lock in
- Innovation** AIPC to Edge to Datacenter & Cloud
- Efficient** Performance per \$ & per W leadership
- Secure** Data as your IP & Models as your IP

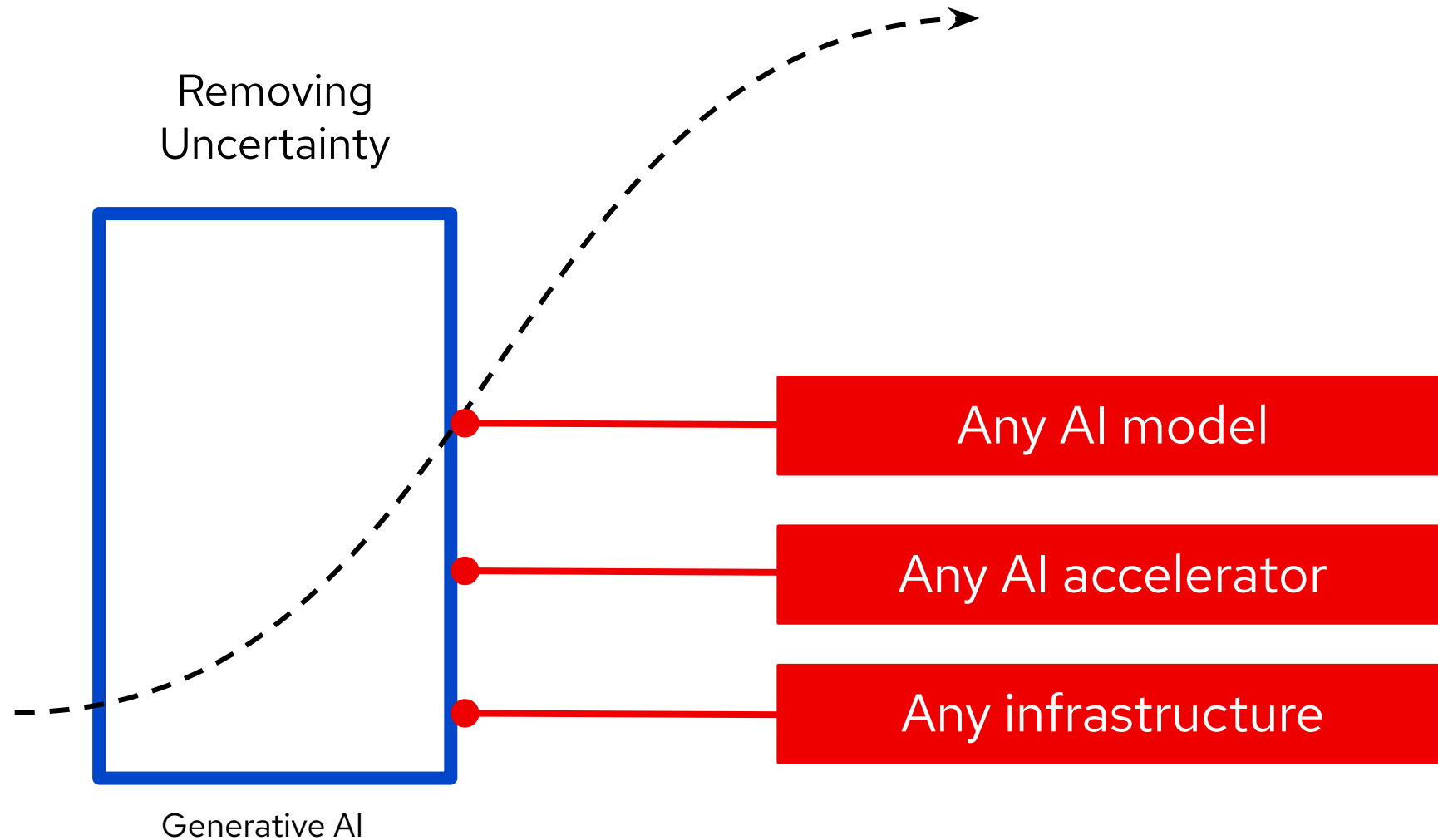
The AI Hierarchy: Mapping ML, DL, and GenAI with Intel

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



Red Hat's AI Strategy and Capabilities

Red Hat AI - Enabling AI Success





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**

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

Accelerate Agentic AI
delivery and stay at the
forefront of innovation





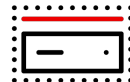
Trusted, Consistent and Comprehensive foundation



Hardware Acceleration



Physical



Virtual



Private
Cloud



Public
Cloud

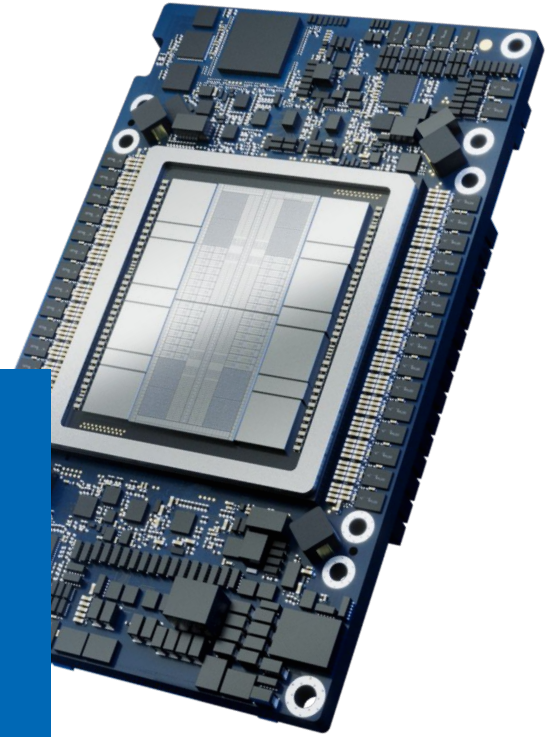


Edge

Intel Gaudi AI Accelerators

Intel® Gaudi® 3 AI Accelerator: AI 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

Up to
120%

More cost efficient
(tokens per dollar)

on Mixtral-8x7B-Instruct-v0.1
vs. leading GPU competitor
with long input and short output sizes

Up to
92%

More cost efficient
(tokens per dollar)

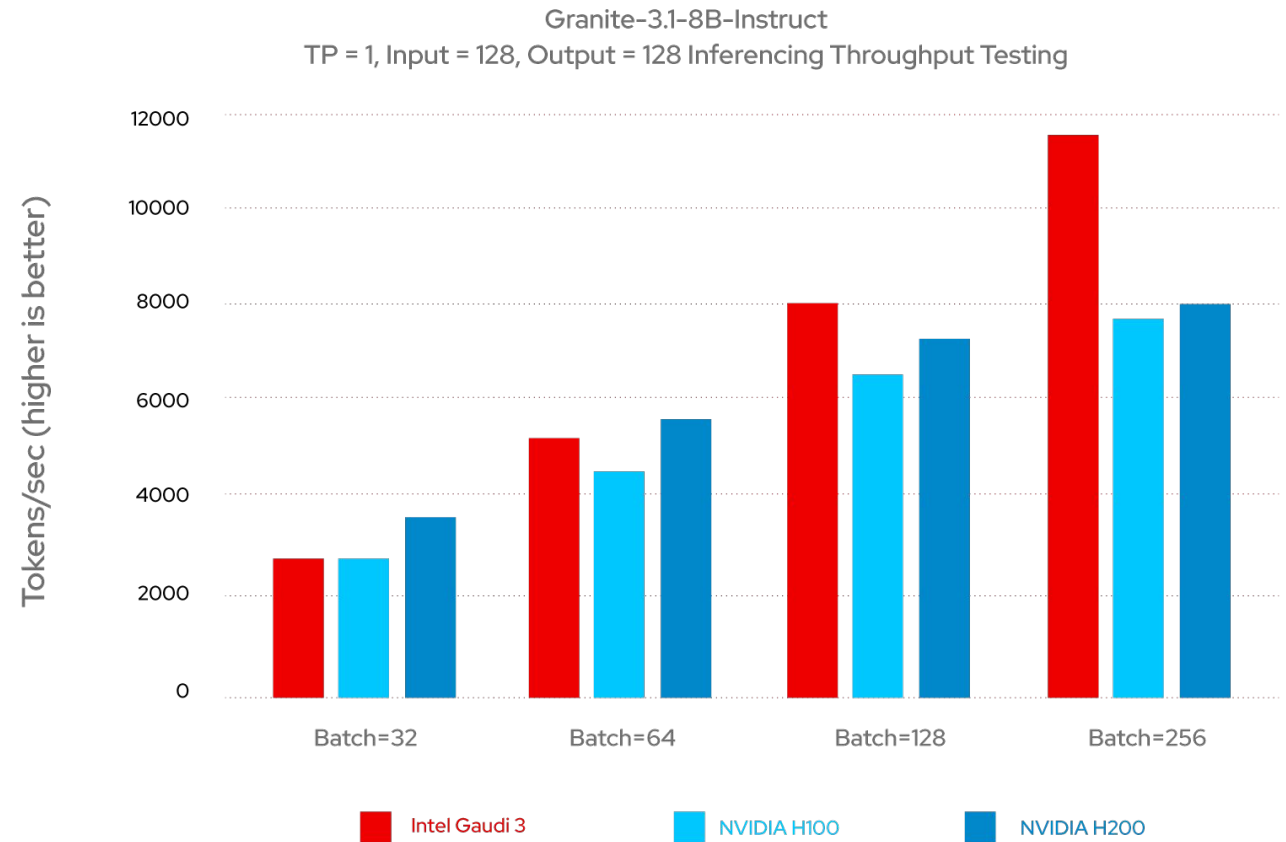
on Llama-3.1-405B-Instruct-FP8
vs. leading GPU competitor
with large context sizes



Up to **43% higher**
throughput than NVIDIA H200

Up to **52% higher**
throughput than NVIDIA H100

For lightweight AI Use Cases



*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 NVIDIA 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.

Results may vary.

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 AI Inferencing performance
vs. AMD EPYC

Unmatched
Performance

See [9G2, 9H9, 9A3] at [intel.com/processorclaims](https://www.intel.com/processorclaims): Intel Xeon 6. Results may vary

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

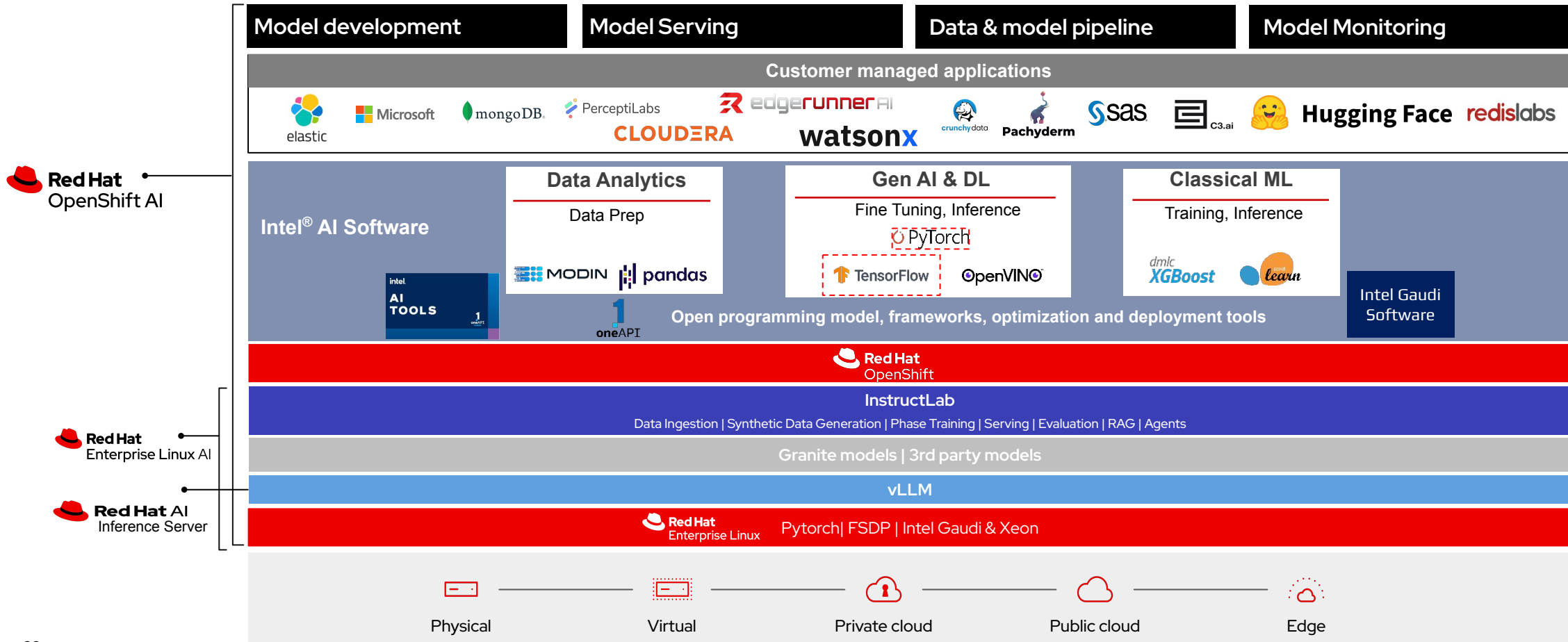
Intel adds Quantum attack protection while providing 1.89 Tb IPsec throughput.

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

Intel AI Software

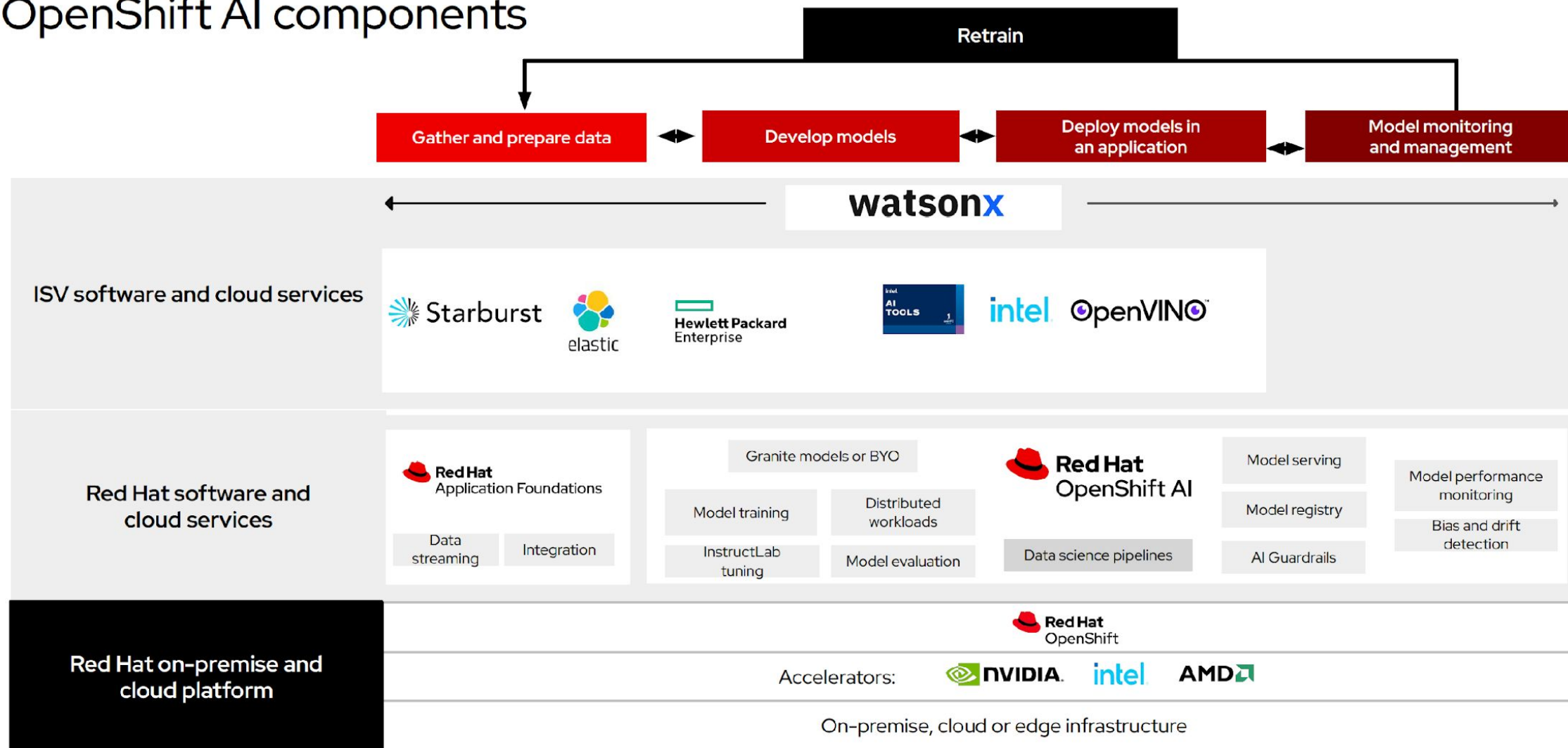
Red Hat AI with Intel AI platform

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



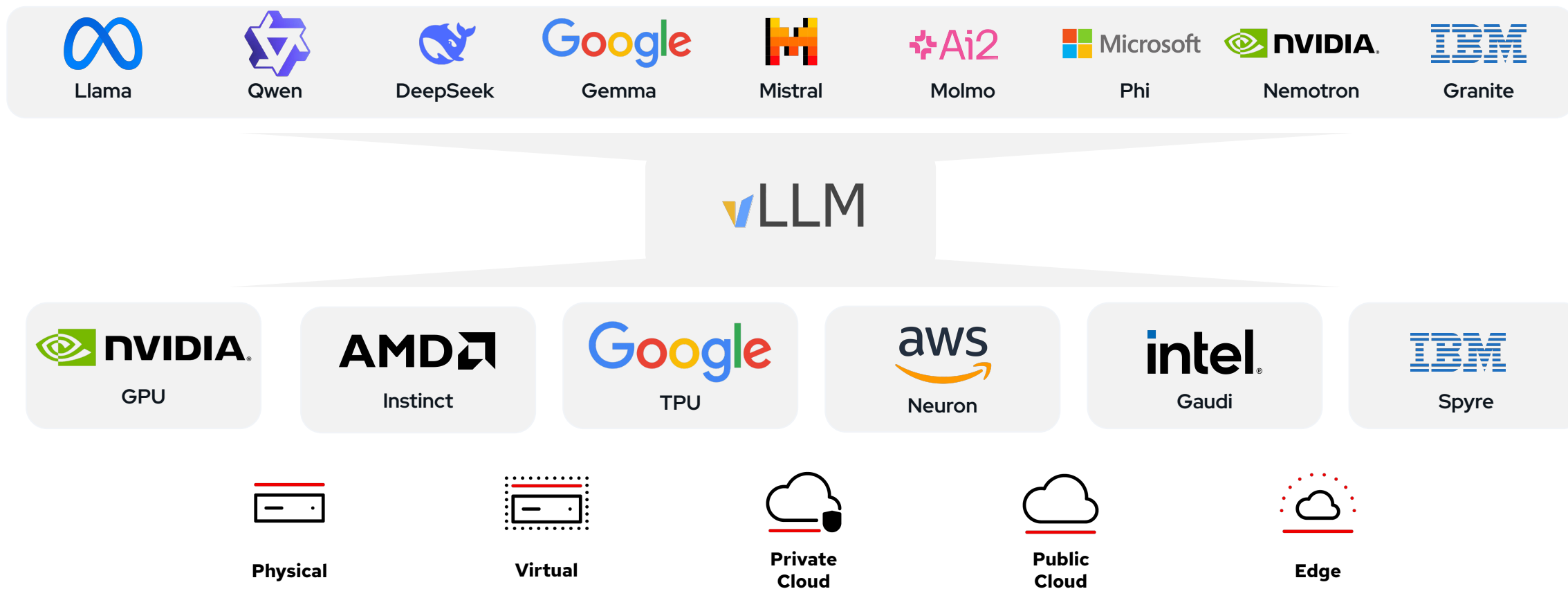
Red Hat AI Platform

OpenShift AI components



Red Hat AI the inference engine for the hybrid cloud

vLLM supports the key models on the key hardware accelerators



Red Hat AI repository on Hugging Face

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

Broad Collection of models



Llama



Qwen



Gemma



Mistral



DeepSeek



Microsoft

Phi



Molmo

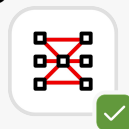


Granite



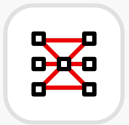
Nemotron

Validated models



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

Optimized models

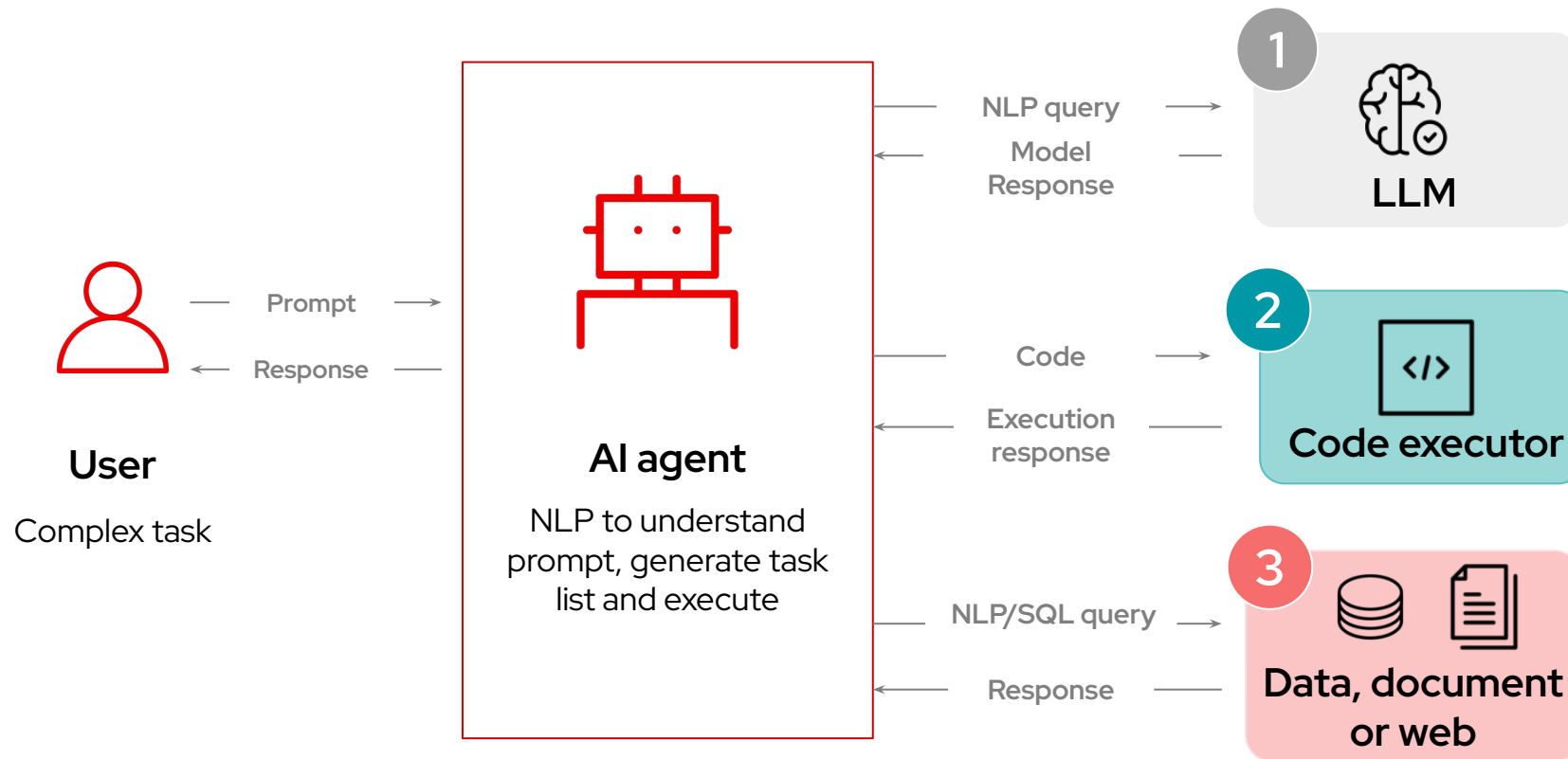


- ▶ Compressed for speed and efficiency
- ▶ Designed to run faster, use fewer resources, maintain accuracy
- ▶ Done using LLM Compressor with latest algorithms

Intro to Agentic AI

AI agents is a **system** that can take independent **actions** to achieve **goals**

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



Agentic AI Demo

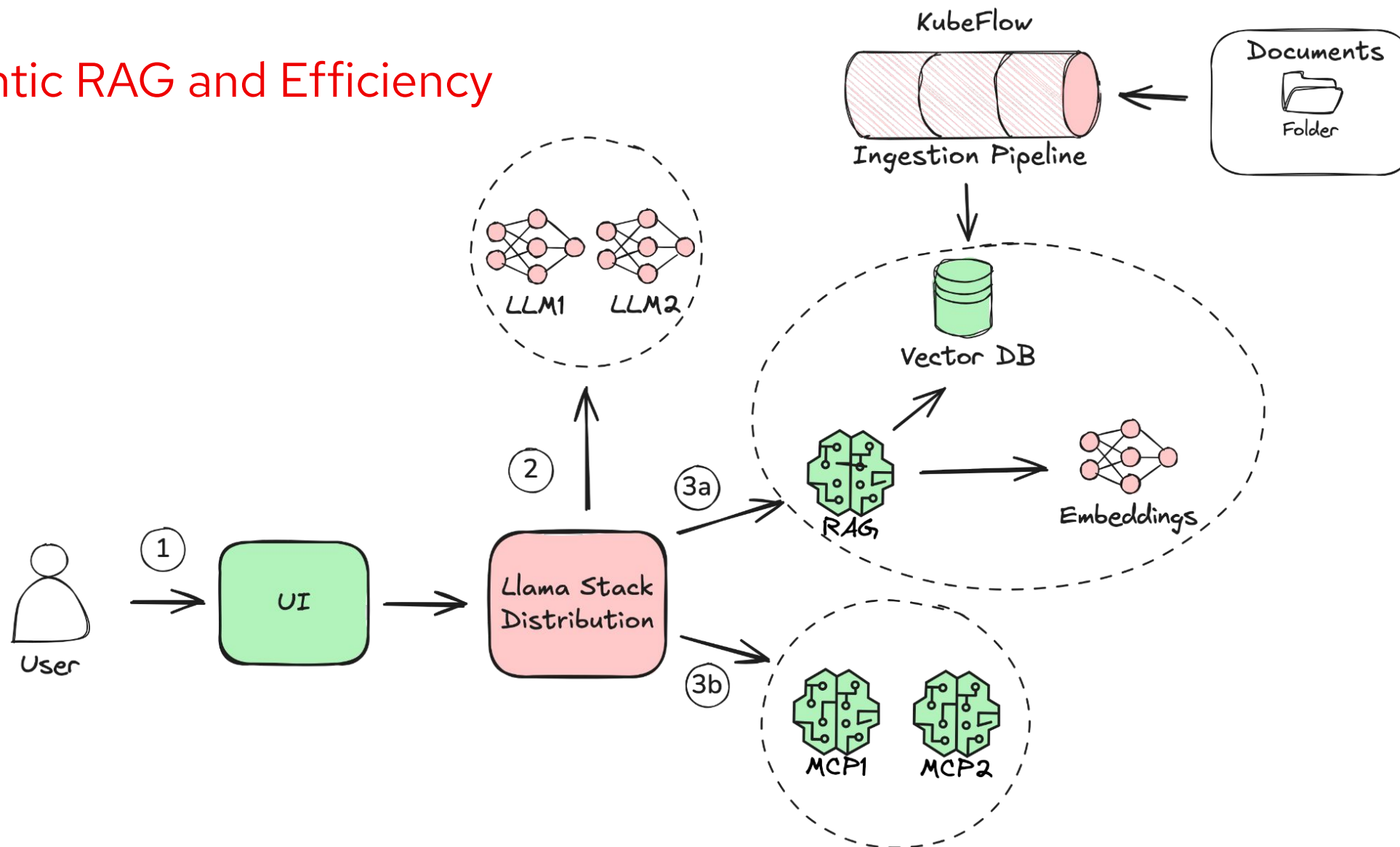
Background

You build an app helping citizens to get money for different family situations.

You want the answer to be based on the context of the person

You want the answer to be specific and with a clear action for the citizen to take.

Agentic RAG and Efficiency



KalturaCapture Edit

eligibility-mcp - Project - Workl... Red Hat OpenShift A

console-openshift-console.apps.sno...fm2aihpcsed.com/k8s/cluster/projects/eligibility-mcp/workloads?view=graph

Red Hat OpenShift

You are logged in as a temporary administrative user. Update the [cluster OAuth configuration](#) to allow others to log in.

Projects > Project details

PR eligibility-mcp Active Actions

Overview Details YAML **Workloads** RoleBindings

Application: All applications View shortcuts

Display options Filter by resource Name Find by name...

eligib...ground eligib...engine

eligibility-mcp-llamastack

eligibility-lsd

granit_dictor granit_dictor llama-...dictor

eligibility-mcp-llamastack

intel

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

Choose your GenAI or Virtualization PoC:

- ☐ Building Inference, RAG, AgenticAI, Model-as-a-Service, and other AI 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 AI 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!



Connect

Thank you



linkedin.com/company/red-hat



facebook.com/redhatinc



youtube.com/user/RedHatVideos



twitter.com/RedHat