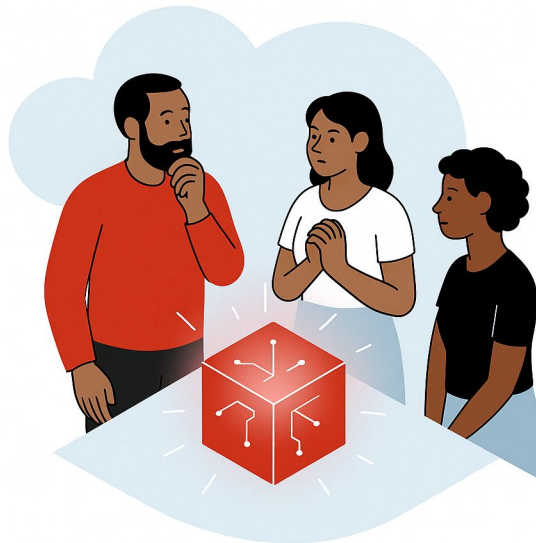# Navigating AI/ML Adoption

## From safe **experimentation** to collaborative **development** to scaled AI **integration**

**Stage 1:** Trial

How do I *evaluate* models and pick the best one for my use case?

**Stage 2:** Experiment

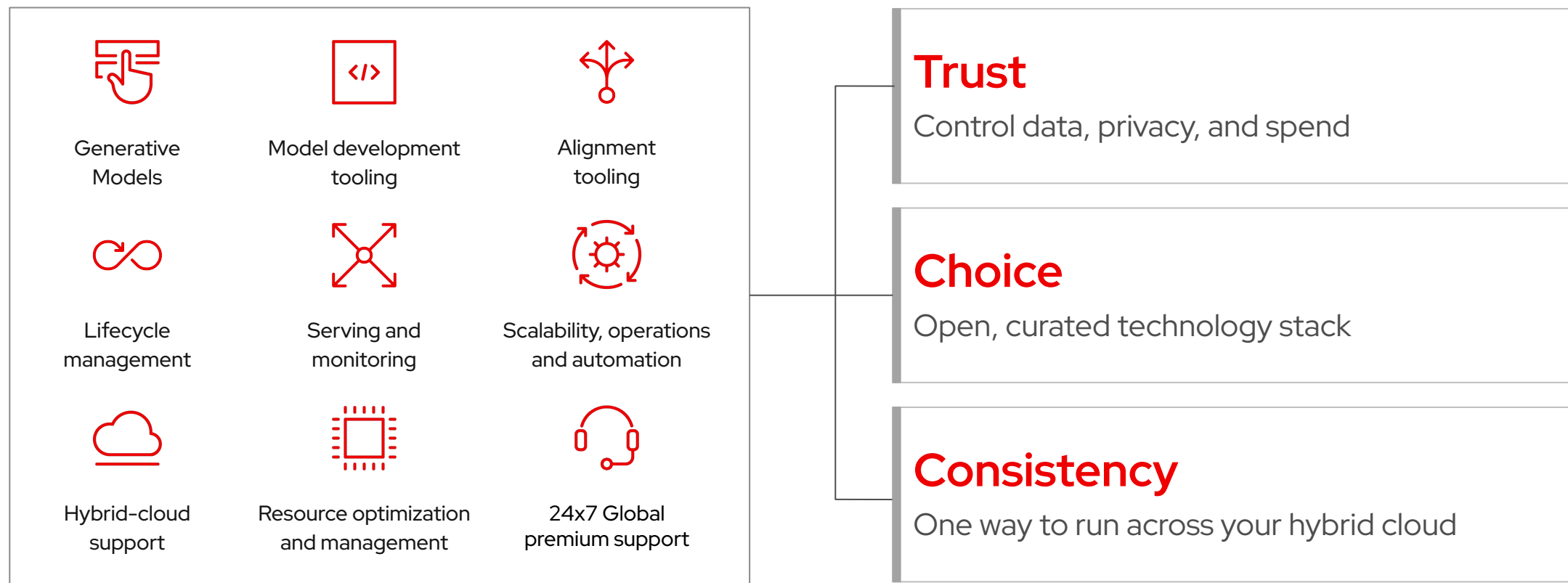How do I *build and integrate* an AI application?

**Stage 3:** Adoption & Scale

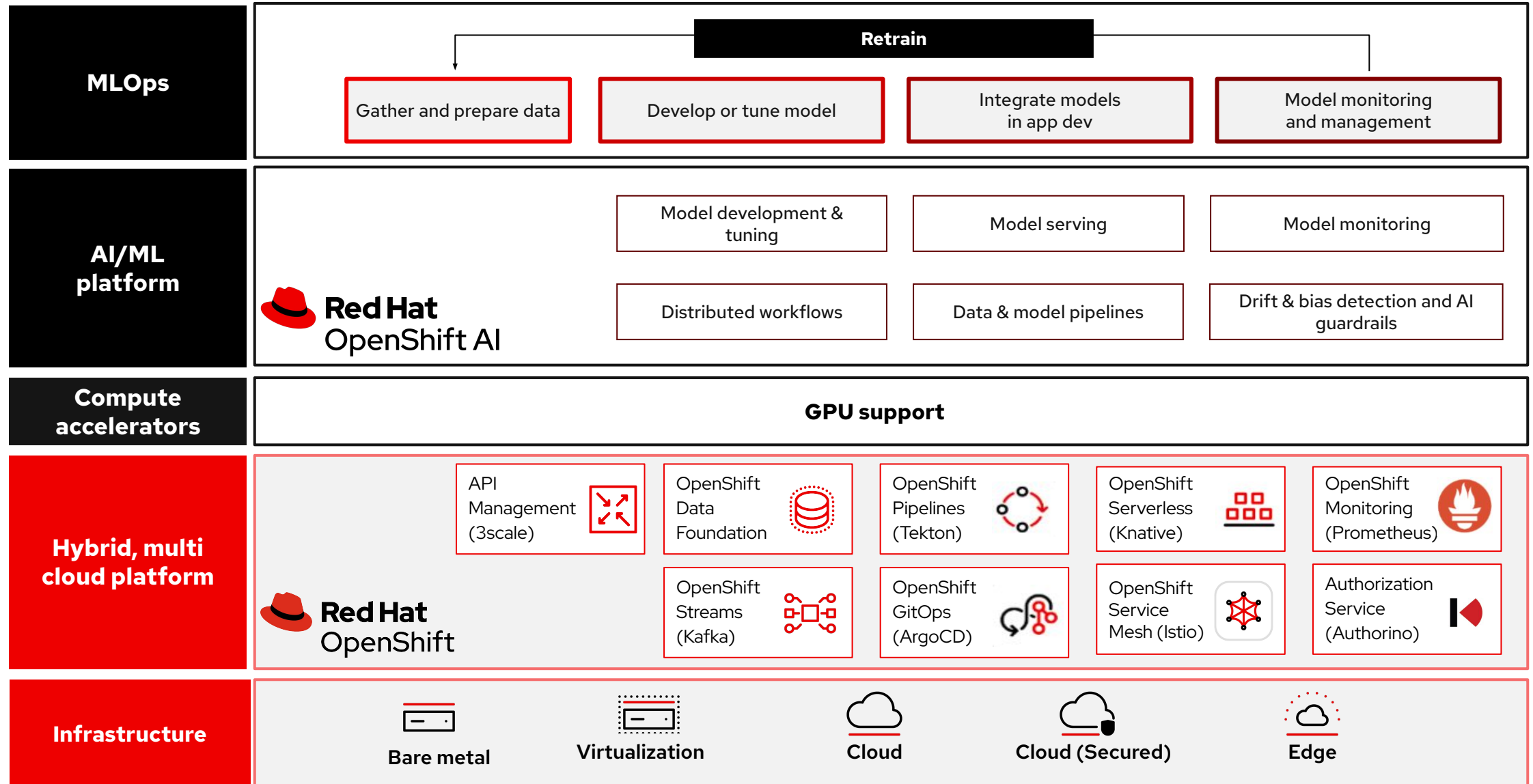How do I *deploy and operate* my AI application?

# Building AI applications requires more than just models

## Red Hat offers generative AI and MLOps capabilities for building flexible, trusted AI solutions at scale

Generative Models

Model development tooling

Alignment tooling

Lifecycle management

Serving and monitoring

Scalability, operations and automation

Hybrid-cloud support

Resource optimization and management

24x7 Global premium support

### Trust
Control data, privacy, and spend

### Choice
Open, curated technology stack

### Consistency
One way to run across your hybrid cloud

Red Hat

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

## MLOps

**Retrain**

Gather and prepare data | Develop or tune model | Integrate models in app dev | Model monitoring and management

## AI/ML platform

**Red Hat OpenShift AI**

Model development & tuning | Model serving | Model monitoring

Distributed workflows | Data & model pipelines | Drift & bias detection and AI guardrails

## Compute accelerators

**GPU support**

## Hybrid, multi cloud platform

**Red Hat OpenShift**

API Management (3scale) | OpenShift Data Foundation | OpenShift Pipelines (Tekton) | OpenShift Serverless (Knative) | OpenShift Monitoring (Prometheus)

OpenShift Streams (Kafka) | OpenShift GitOps (ArgoCD) | OpenShift Service Mesh (Istio) | Authorization Service (Authorino)

## Infrastructure

Bare metal | Virtualization | Cloud | Cloud (Secured) | Edge

**Red Hat**

# AI/ML Adoption Is a Collaborative Journey

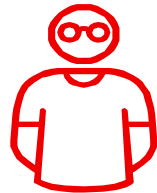Every member of your team plays a critical role—from early exploration to scaled integration

| **Business leadership** | **Data engineer** | **Data scientists** | **ML engineer** | **App developer** | **IT operations** |
|---|---|---|---|---|---|
| Define AI vision and success criteria | Prepare and pipeline data for modeling | Train, fine-tune, and evaluate models | Operationalize models with MLOps | Integrate models into business applications | Manage platform and ensure secure, scalable access |

Red Hat

# Stage: Trial
# *Evaluate Models*

**Strategic Goal**

Establish a low-risk entry point for AI adoption. *Focus on enabling initial infrastructure and secure access to pre-trained models* to support evaluation and early alignment.
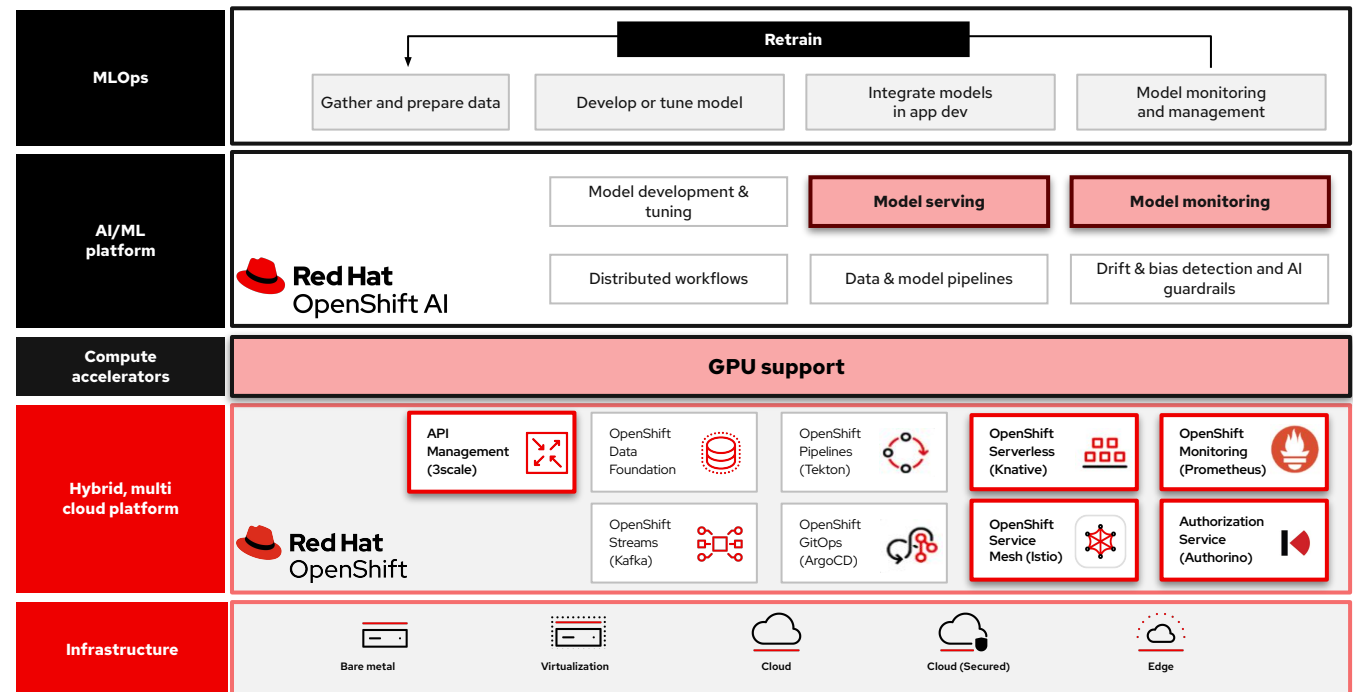
**Why This Stage Matters**

- *Establishes clarity*
- *Aligns teams*
- *Lays governance*

# *Enabling Foundational AI Capabilities*

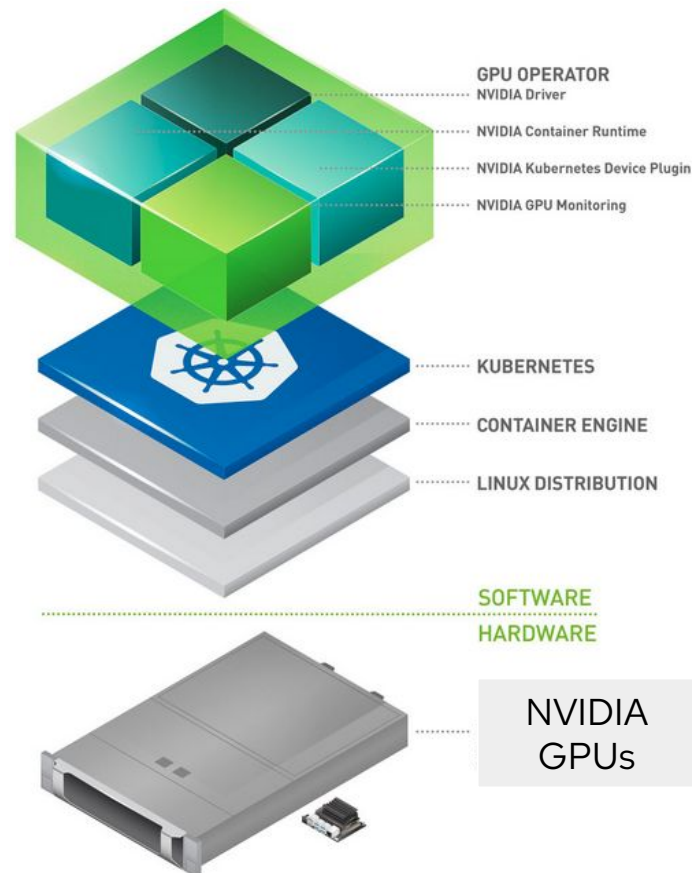## Focus on establishing secure, observable, and scalable **infrastructure** to support early experimentation

### **Platform Objectives:**

- **Provide a secure and scalable AI infrastructure**

- **Support model inference**

- **Expose models through secure APIs**
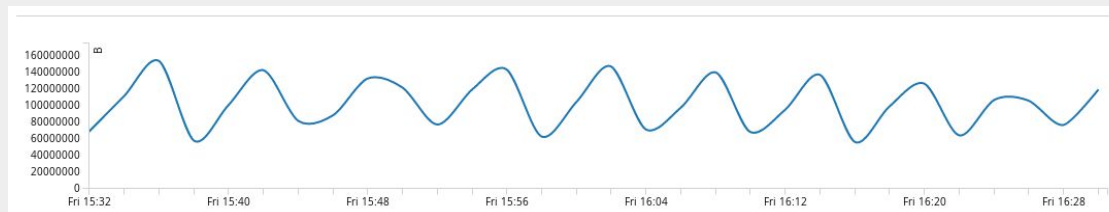
- **Provide a governed model catalog**

# NVIDIA GPU Operator on OpenShift

## Automated GPU provisioning and monitoring for scalable AI workloads



GPU OPERATOR
NVIDIA Driver
NVIDIA Container Runtime
NVIDIA Kubernetes Device Plugin
NVIDIA GPU Monitoring

KUBERNETES
CONTAINER ENGINE
LINUX DISTRIBUTION
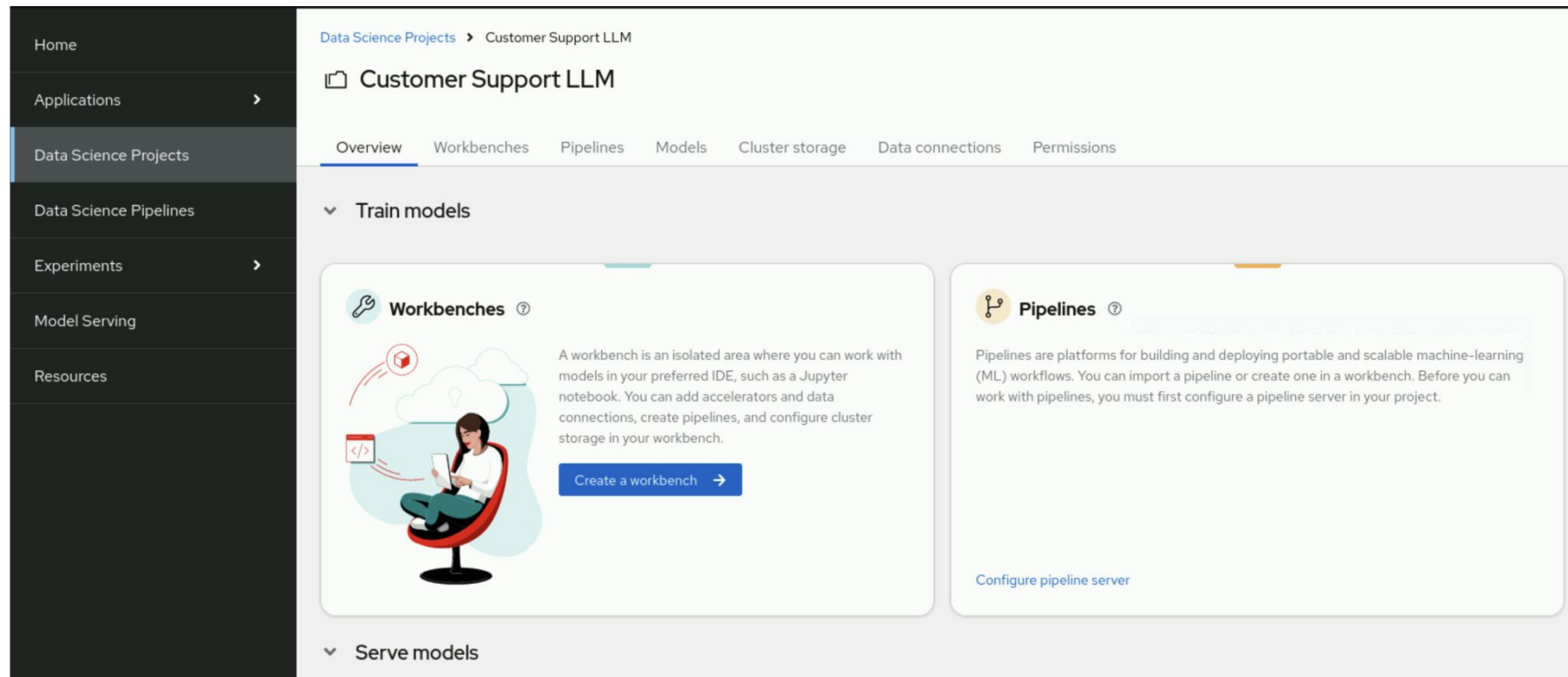
SOFTWARE
HARDWARE

NVIDIA GPUs

1. Build NVIDIA GPU driver for RHCOS

2. Expose NVIDIA GPUs as K8s extended resources
   `nvidia.com/gpu:2`

3. Advertise NVIDIA GPU features with Node labels
   ```
   nvidia.com/gpu.count=2
   nvidia.com/gpu.memory=20096
   nvidia.com/gpu.family=ampere
   nvidia.com/gpu.product=A100-PCIE-40GB-MIG-...
   ```
   …

4. Expose NVIDIA GPU metrics to OpenShift



5. Handle operator upgrades

# OpenShift AI Data Science Projects

## Unified, secure workspaces for building, managing, and collaborating on AI projects
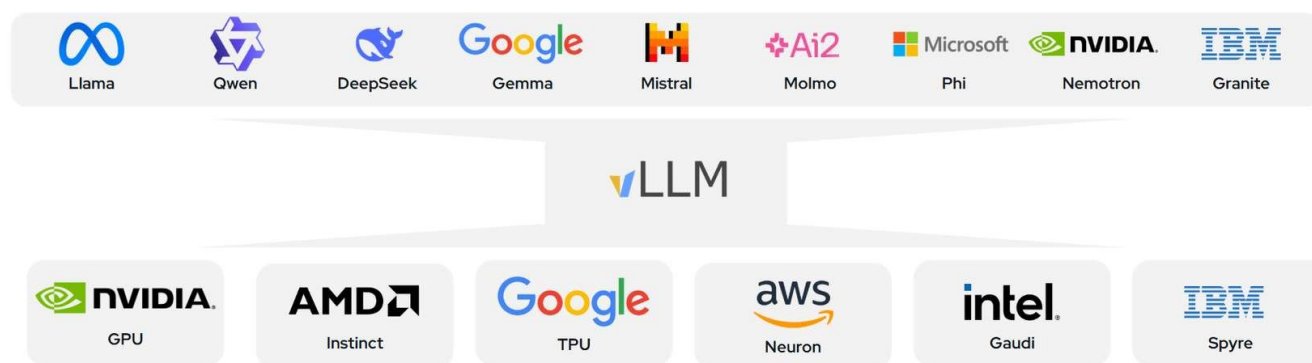


**Data Science Projects** enable teams to *organize, collaborate, and manage end-to-end AI projects*—offering workbenches, data connections, storage, pipelines, and model serving in a unified, secure environment.

**Key Features:**

- Scoped project workspaces: Kubernetes namespaces encapsulating notebooks, compute, storage, and pipelines

- Integrated development workbenches: Launch Jupyter, VS Code, or RStudio environments directly within projects

- Managed data & storage: Attach cluster storage and external data connections (e.g. S3) via UI

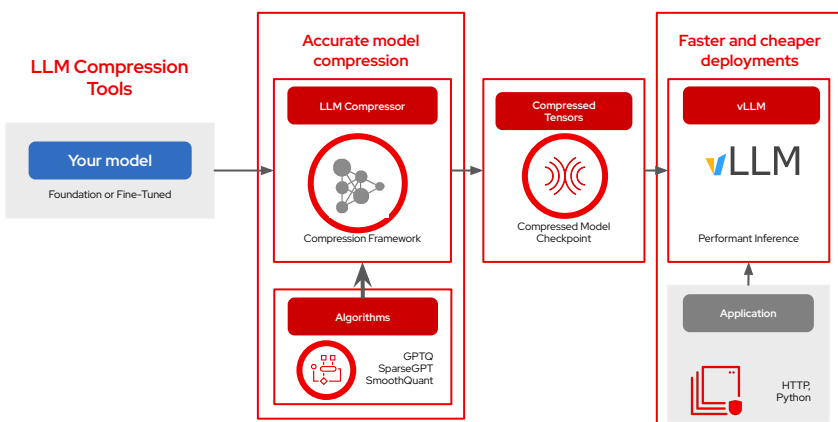Red Hat

# vLLM Inference Server

## Fast, memory-efficient, production-grade LLM serving for GenAI workloads



**vLLM Inference Server** enables efficient, production-grade LLM serving while *significantly reducing infrastructure and compute costs*—making it ideal for real-time, scalable AI applications across enterprises and research environments.

Key Features:

- **Faster response time:** Higher throughput handles more requests in less time for faster responses.

- **Efficient memory management:** Organized memory enables larger models a on existing hardware.

- **Reduced hardware costs:** Efficient resource utilization reduces required GPUs and overall infrastructure costs.

- **Designed for security and scale:** Self-hosting with vLLM strengthens data privacy, control, and scalable growth.

# OpenShift AI Model Registry

## Establishing model governance from data to deployment



**A centralized AI/ML model registry** *that tracks models from registration through deployment*, enabling collaboration, version control, governance, and streamlined MLOps workflows.

**Key Features:**

- **Central hub:** stores models, parameters, metrics, and deployment events
- **Model versioning:** register models and versions with object storage
- **Collaboration & governance:** role-based access control

DEMO

# Early Role Alignment for AI Trial

## Introducing core teams to AI through goal-setting and hands-on evaluation of pre-trained models

| Set goals | Gather and prepare data | Develop model | Integrate models in app dev | Model monitoring and management |
|---|---|---|---|---|
| **Business leadership** ▬▬▬▬ | | | | |
| Data engineer | | | | |
| **Data scientist** | | ▬▬▬▬ | | |
| ML engineer | | | | |
| App developer | | | | |
| **IT operations** | | ▬▬▬▬ | | |

▸ **Business leadership** sets initial goals and success criteria

▸ **Data scientists** evaluate pre-trained models to test feasibility and value.

▸ **IT operations** enables secure access and controlled infrastructure

Red Hat

# Laying the Groundwork for Enterprise AI

## Establishing an AI/ML foundation with secure infrastructure, and early alignment across teams

**Stage: Trial - Outcomes:**

- **Platform Enablement**
  ✅ *Foundational infrastructure is operational* (GPU provisioning, platform metrics).
  ✅ *Secure access to models and usage controls is enabled via curated catalog and secure APIs.*
  ✅ *As-a-Service capabilities activated* (GPU, MaaS) to accelerate experimentation.

- **Team Enablement**
  ✅ *Cross-functional alignment* between platform teams, data scientists, and business stakeholders.
  ✅ *AI roles and responsibilities clarified* across cross-functional teams.
  ✅ *Teams gain first hands-on AI experience* in a controlled environment.

- **Business Readiness**
  ✅ *Early success criteria* and AI use cases align with business goals.
  ✅ *Clarity on AI value* through trial workloads using pre-trained models.
  ✅ *A governance foundation* is established to prevent "shadow AI" and support safe experimentation.

Red Hat

# Stage: Experiment
## *Fine-tune Models*

**Strategic Goal**

Enable experimentation with real enterprise data. *Focus on building collaborative environments, customizing models, and defining repeatable workflows* to support scalable and secure AI development.

**Why This Stage Matters**

- *Introduces enterprise data*
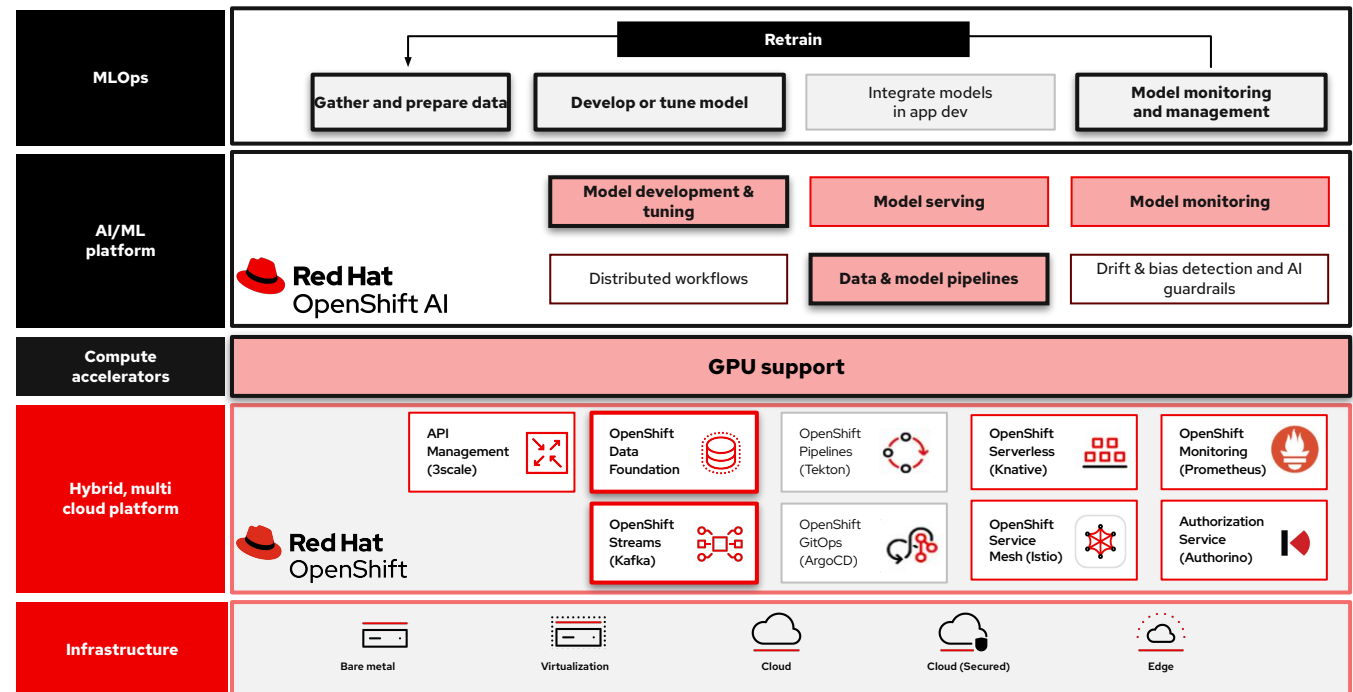- *Promotes collaboration*
- *Enables model customization*

Red Hat

# *Enabling Collaborative AI Experimentation*

## Focus on establishing secure, observable, and scalable infrastructure to support early experimentation
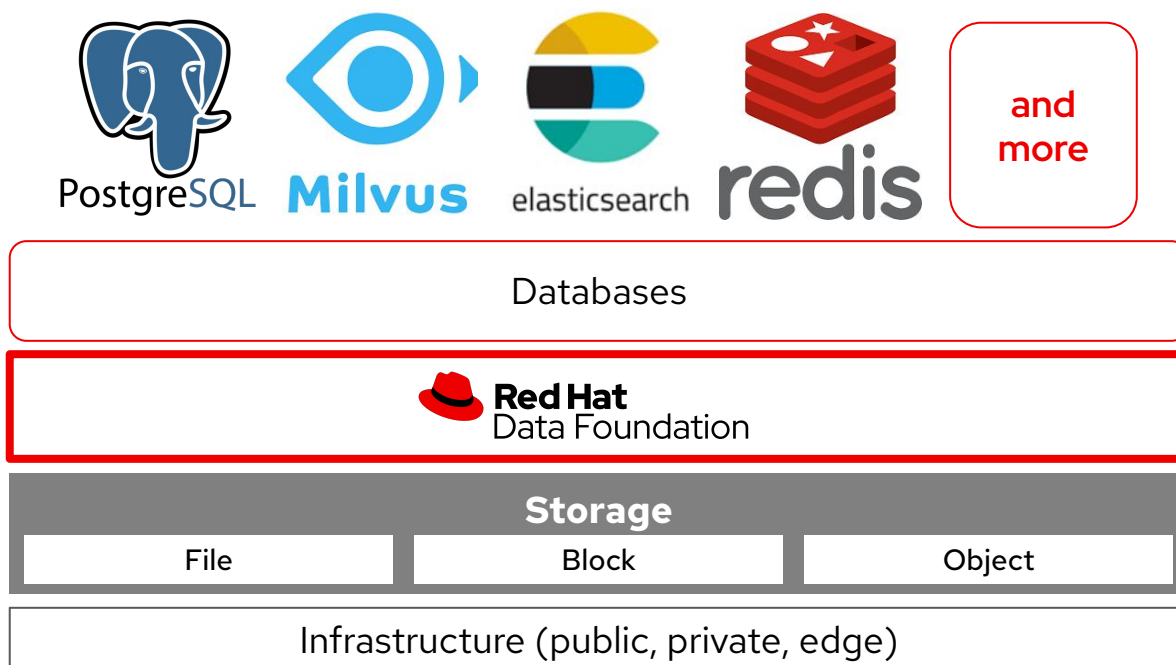
### Platform Objectives:

- Ingest, store, and prepare enterprise data

- Enable collaborative development environments

- Support model training for differentiated AI

# OpenShift Data Foundation and OpenShift Operators

## Scalable, resilient storage for AI/ML workloads with integrated vector database support



PostgreSQL

**Milvus**

elasticsearch

redis

**and more**

Databases

**Red Hat**
Data Foundation

**Storage**

| File | Block | Object |
|------|-------|--------|

Infrastructure (public, private, edge)

**OpenShift Data Foundation** provides the *scalable, resilient storage backbone for AI/ML workloads*, supporting object, block, and file storage. Combined with OpenShift Operators, it *enables seamless deployment and lifecycle management of vector databases* like:

- PostgreSQL: Stores structured data, model metadata, and vectors via `pgvector`.

- Milvus (Vector DB): High-performance embedding storage and vector search in RAG pipelines.

- Elasticsearch: Hybrid text and vector search for RAG and large-scale indexing.

- Redis: Low-latency feature serving, caching, and vector search for real-time inference workloads.

Red Hat

# Docling

## Automated document transformation for GenAI and retrieval-augmented generation (RAG)



**User-Provided Docs**

**Context-Aware Chunking**
Text, tables, figures, lists, columns

**Docling** simplifies document processing by *converting complex files into structured, AI-ready formats*—making it easier to use documents and audio in search, RAG, and analytics workflows.

Key Features:

- Multi-format parsing: supports PDF, DOCX, PPTX, HTML —even scanned or layout-complex formats

- Advanced layout & table extraction: leverages DocLayNet and TableFormer models for accurate structure, formulas, code, and tables

- Built-in Automatic Speech Recognition: transcribes audio (WAV/MP3) into structured text

# JupyterLab on OpenShift AI

## Web-based notebooks for scalable, containerized AI development



**JupyterLab** delivers a web-based notebook IDE running on Kubernetes—*empowering teams to develop, train, and deploy AI workflows* with cluster-grade compute, storage, and data integrations.

Key Features:

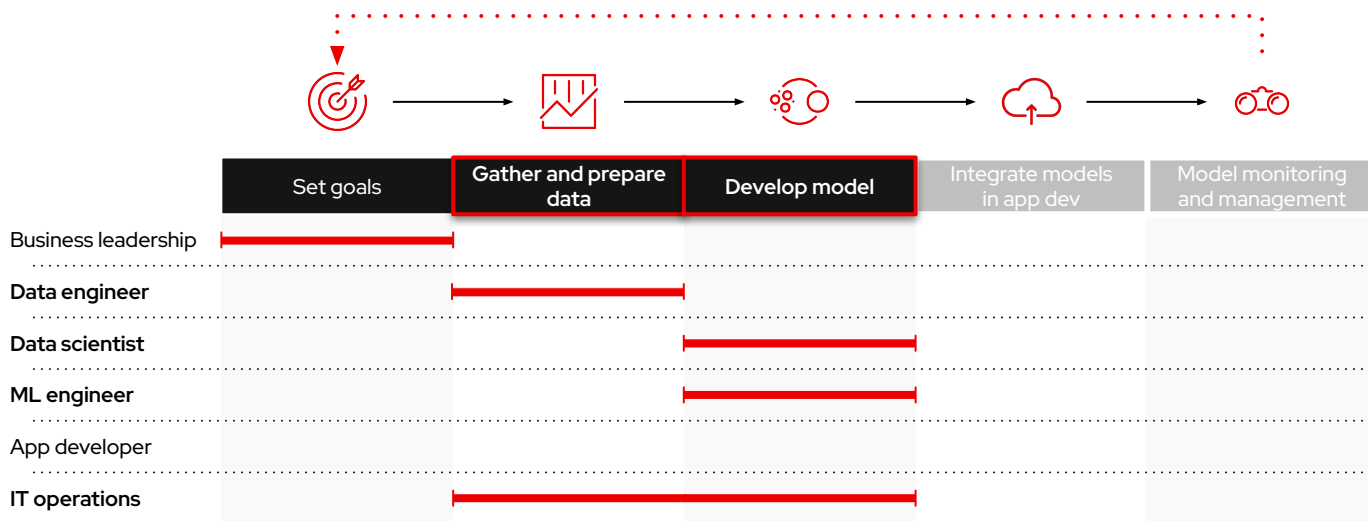- **Managed notebook server:** Launch JupyterLab containerised environments for secure and scalable development**.**

- **Integrated IDE options:** Choose from JupyterLab, VS Code, or RStudio—all containerized within projects.

- **Persistent storage & data access:** Attach cluster or S3-compatible external data sources seamlessly.

- **Cluster-powered compute:** Run code on powerful CPUs/GPUs in the cluster.

DEMO

# Role Expansion for AI Experimentation

## Enabling hands-on collaboration with enterprise data, pipelines, and model customization



| | Set goals | Gather and prepare data | Develop model | Integrate models in app dev | Model monitoring and management |
|---|---|---|---|---|---|
| Business leadership | ▬▬▬ | | | | |
| Data engineer | | ▬▬▬ | | | |
| Data scientist | | | ▬▬▬ | | |
| ML engineer | | | ▬▬▬ | | |
| App developer | | | | | |
| IT operations | | ▬▬▬▬▬▬ | | | |

▸ **Data engineers** build pipelines and prepare data for modeling.

▸ **Data scientists** customize and evaluate models with enterprise context data.

▸ **ML engineers** fine-tune models and manage experimentation workflows.

▸ **IT operations** maintains environment and storage access

Red Hat

# Accelerating Collaborative AI Development

## Establishing scalable, secure, and reusable foundations for experimentation and model customization

**Stage: Experiment - Outcomes:**

- **Platform Enablement**
  ✅ *Enterprise data pipelines operational* (ingestion, storage, reuse).
  ✅ *Training-ready environments available* to support model fine-tuning.
  ✅ *As-a-Service capabilities enabled* (Data as a Service, Data Science as a Service) to streamline experimentation.

- **Team Enablement**
  ✅ *Cross-functional collaboration extends* across data engineers, scientists, and platform teams.
  ✅ *Shared environments in place* for collaborative development.
  ✅ *Automated workflows in use* for reproducibility and governance.

- **Business Readiness**
  ✅ *Real enterprise data in use* to test business-aligned use cases.
  ✅ *Differentiated AI development enabled* via model fine-tuning and customization.
  ✅ *Compliance controls introduced* to support scale and secure use.

Red Hat

# Stage 3: Adopt
## *Scale & Integrate*

**Strategic Goal**

Operationalize AI at scale across the enterprise. *Focus on automating model delivery, embedding AI into business applications, and ensuring trust, governance, and observability* in production.

**Why This Stage Matters**
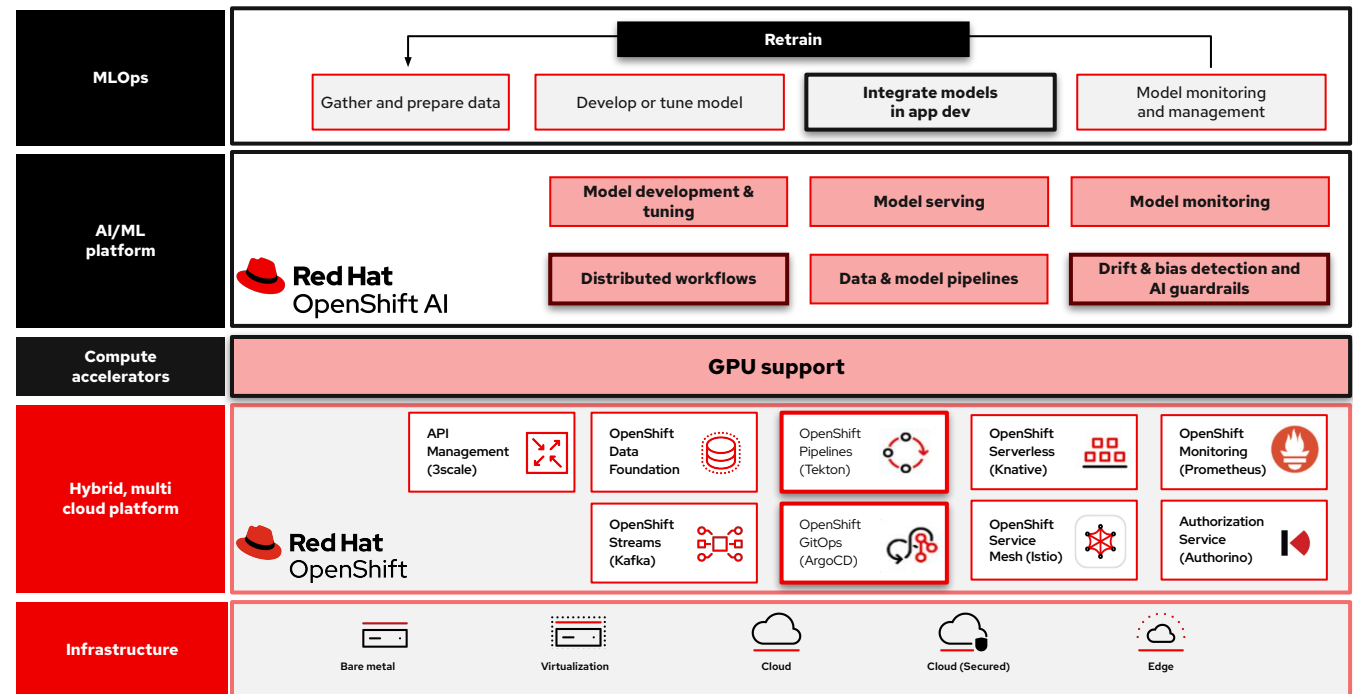
- *Integrates AI*
- *Establishes trust*
- *Sustainable AI operations*

# *Enabling Scalable and Trusted AI Integration*

## Focus on automating delivery, ensuring trust, and embedding AI into enterprise applications
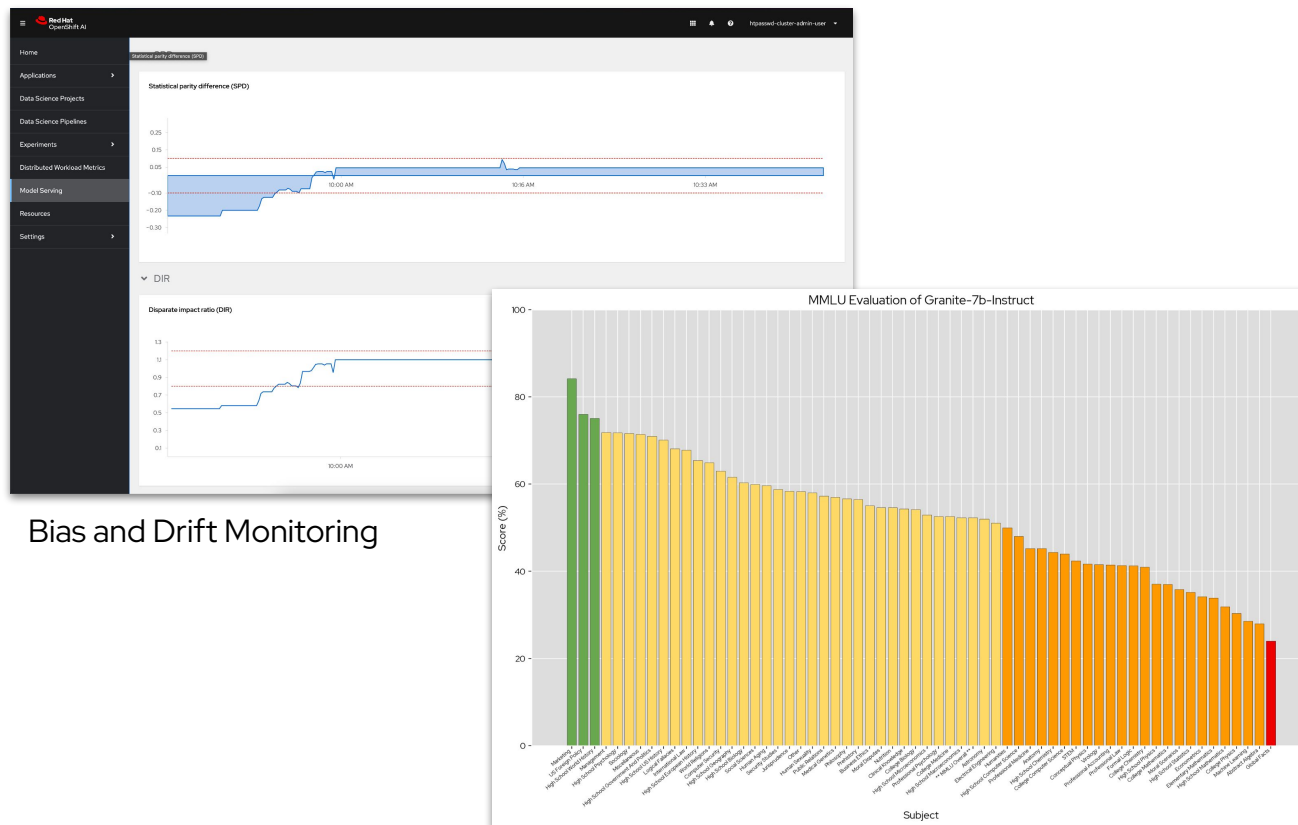
### Platform Objectives:

- **Embed trust, observability, and accountability**

- **Integrate AI models into enterprise applications**

- **Automate model delivery workflows**

# TrustyAI on OpenShift AI

## Automate monitoring, fairness checks, and safety enforcement for responsible AI



Bias and Drift Monitoring



Language Model Evaluation

**TrustyAI** a comprehensive toolkit for responsible AI offering explainability, fairness, drift monitoring, and guardrails. Integrated via an Operator, it *empowers teams to deploy transparent, auditable models with automated monitoring and safety enforcement* as part of their MLOps pipelines.

Key Features:

- **Explainability built-in:** Offers local and global model explainers (e.g., LIME, SHAP, counterfactuals) to interpret predictions.

- **Fairness & drift detection:** Continuously monitors bias and data drift metrics in production.

- **AI guardrails:** Deploys modular detectors (e.g., profanity, prompt injection) to enforce safety and policy compliance on LLM inputs/outputs.

25

# Kubeflow on OpenShift AI

## Run scalable, reproducible ML workflows across hybrid environments



**Kubeflow** a *Kubernetes-native MLOps platform that supports the full AI lifecycle*—from model development to distributed training and serving. It offers multi-tenancy, GPU autoscaling —*empowering AI/ML teams to build scalable, composable, and cloud-portable workflows*.

Key Features:

- **Modular pipeline:** Includes Kubeflow Pipelines for building, scheduling, and managing reproducible ML workflows.

- **Distributed training operators:** Scale training jobs with built-in support for TensorFlow (TFJob), PyTorch (PyTorchJob), and MPI-based workloads.

- **Hyperparameter tuning (Katib):** Automate model optimization using scalable search algorithms across distributed training jobs.

# Llama Stack on OpenShift AI

## Standardize and accelerate AI agent development with modular components

**AI/ML platform**

**Llama Stack**

| | |
|---|---|
| Datasets | Inference |
| Vector.io | Telemetry |
| Agents | Evaluation |
| Safety | Tool Calling (MCP) |

**Other Agent frameworks**

**Platform services**

- Over the air updates
- Monitoring
- Networking
- Egress
- Storage
- Log forwarding
- Authorization
- Registry
- install

**Hardware accelerators**

**Deploy anywhere**

**Llama Stack** a modular, *agent-capable AI framework* that combines LLM inference, RAG, and tool integration. It *simplifies building intelligent agents by offering standardized APIs* and seamless orchestration—*accelerating AI application development with agentic workflows* and operational tooling.

Key Features:

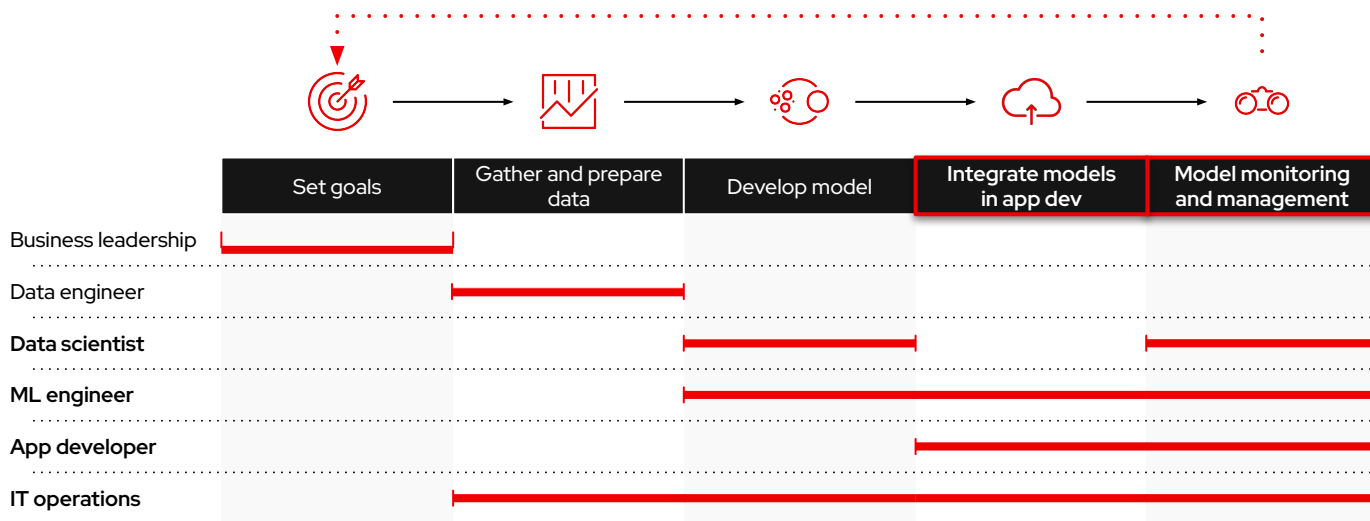- Unified API: Standard interface for inference, embeddings, RAG, and tool execution.

- MCP-based tool integration: Use the Model Context Protocol to discover and invoke external tools via APIs

- Built-in RAG support: Natively ingests documents and retrieves context using vector-DB backends

- Interactive agents: Build multi-step agents that retrieve, reason, and act with real-time context.

27

Red Hat

# DEMO

# Full Role Integration for Scaled AI Delivery

## Operationalizing models across apps with CI/CD, monitoring, and business impact in mind

| | Set goals | Gather and prepare data | Develop model | Integrate models in app dev | Model monitoring and management |
|---|---|---|---|---|---|
| Business leadership | ▬ | | | | |
| Data engineer | | ▬ | | | |
| Data scientist | | | ▬ | | ▬ |
| ML engineer | | | ▬ | | |
| App developer | | | | ▬ | |
| IT operations | | ▬ | | | |

▶ **Data scientists** validate fairness, explainability, and post-deployment performance.

▶ **ML engineers** manage MLOps pipelines for model lifecycle, from deployment to rollback.

▶ **App developers** integrate AI models into production applications using APIs, service mesh, or custom logic.

▶ **IT operations** deliver secure, compliant, and stable infrastructure for model deployment.

# Operationalizing AI Across the Enterprise

## Scaling delivery, embedding AI into business systems, and reinforcing trust in production

### Stage: Adopt – Outcomes:

- **Platform Enablement**
  - ✅ *Production platform fully operational* with automated governance controls.
  - ✅ *MLOps as a Service enabled* for model lifecycle management, promotion, and rollback.
  - ✅ *AI as a Service enabled* for scalable access and reuse.
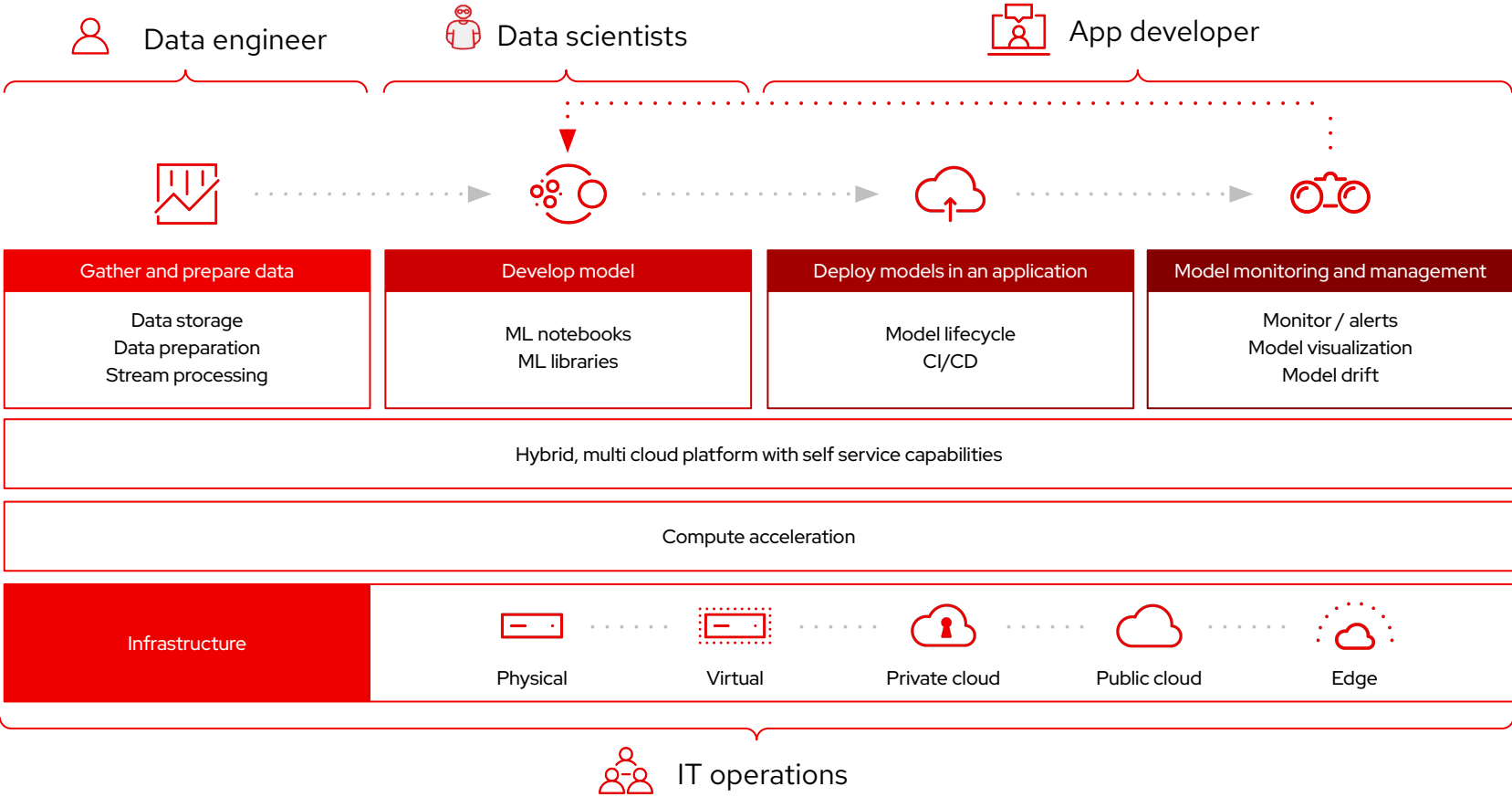
- **Team Enablement**
  - ✅ *End-to-end collaboration established* from model training to monitoring.
  - ✅ *AI knowledge extended* across AI/ML engineers, App developers, and IT operations.
  - ✅ *Governed delivery in place* with clear ownership and measurable outcomes.

- **Business Readiness**
  - ✅ *AI integrated into business applications* through secure APIs and workflows.
  - ✅ *Trust and compliance reinforced* with drift detection, explainability, and fairness checks.
  - ✅ *Scalability foundation established* for enterprise-wide AI adoption.

Red Hat

# Navigating AI/ML Adoption: What We've Achieved

## A milestone-driven journey from safe experimentation to scalable, responsible AI adoption

Data engineer    Data scientists    App developer



| Gather and prepare data | Develop model | Deploy models in an application | Model monitoring and management |
|---|---|---|---|
| Data storage<br>Data preparation<br>Stream processing | ML notebooks<br>ML libraries | Model lifecycle<br>CI/CD | Monitor / alerts<br>Model visualization<br>Model drift |

Hybrid, multi cloud platform with self service capabilities

Compute acceleration

| Infrastructure | Physical | Virtual | Private cloud | Public cloud | Edge |

IT operations

**Stage 1:** Trial

**A foundational platform**

**Stage 2:** Experiment

**An AI development environment**

**Stage 3:** Adopt & Scale

**AI models are integrated**

Red Hat

Red Hat Summit

Connect

# Thank you

## Adnan Drina

Specialist Solution Architect
AppDev & AI