



# Harnessing Red Hat AI on ROSA (AWS)

Practical Life Cycle Tips for Intelligent Applications

Yury Titov

Sr. EMEA Black Belt for Managed Cloud Services

# Introduction



## Yury Titov

- ▶ former senior EMEA Architect
- ▶ present: senior BlackBelt for Managed Cloud Services
- ▶ always: open source dude

# Agenda

- ▶ **How the AI/ML landscape is evolving: market opportunities & challenges**
- ▶ **AI Application Examples vs intelligent Application?**
- ▶ **Challenges of Operationalizing AI ?**
- ▶ **Team topologies and operationalizing models**
- ▶ **Red Hat OpenShift AI - key features and walkthrough**
- ▶ **Why application platforms? Gartner two speed architecture**
- ▶ **Where to start?**
- ▶ **Conclusion**

# How the AI/ML landscape is evolving





AI is becoming a part of our everyday lives

watsonx Code Assistant



Chat GPT



Stable Diffusion

IBM Granite Models

Llama 2

Meta AI



DALL·E 2

Gemini



GitHub  
Copilot

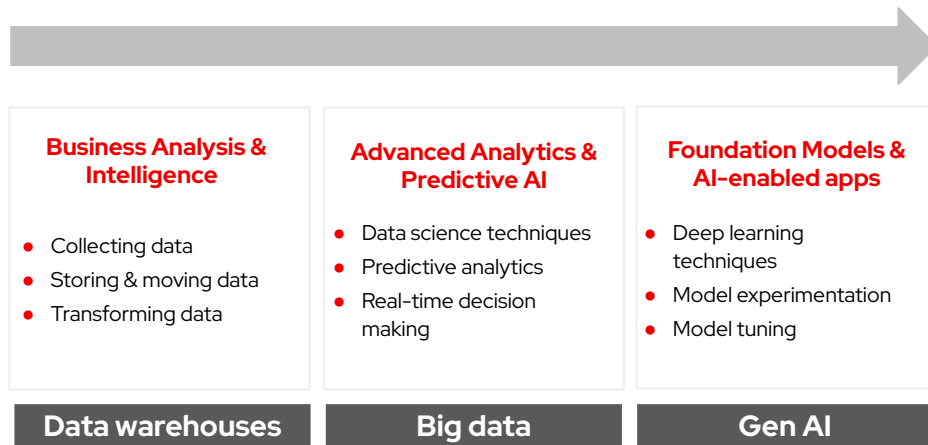


MISTRAL  
AI\_

# AI has undergone significant evolution

The evolution of AI: from Business Intelligence to Generative AI

- ▶ Predictive AI runs businesses today
- ▶ Foundation models provide a shortcut for realizing the value of AI



# Intelligent Applications?



# Examples of intelligent applications

- ▶ **Recommendation engines**

Netflix, Amazon, etc.

- ▶ **Virtual assistant**

Siri, Alexa, etc.

- ▶ **Detecting fraudulent activity**

Money laundering, spam, hacking, insurance

- ▶ **Quantifying risks and making smart decisions**

Insurance, loans

- ▶ **Pattern detection**

Images, videos: how many cars, humans, etc.

- ▶ **Analyze specialized data**

Seismic data for oil and gas

- ▶ **Teach AI to play video games**

AI opponents

- ▶ **Text analysis**

Summarization, accuracy, offensive, plagiarism detection

- ▶ **Medical**

Tumour detection

- ▶ **Customer retention**

Predict who's about to leave

# Generative AI Application Examples

- ▶ **Text Generation**

Content creation, chatbots, etc.

- ▶ **Code Generation**

Automate and supplement code development

- ▶ **Image Creation**

Create new images for art, design, games, etc.

- ▶ **Music development**

Create original music based on existing styles

- ▶ **Medical applications**

Suggest new molecules for drug development

- ▶ **Data augmentation (synthetic data)**

Create additional training data for model development

- ▶ **Anomaly detection**

Detect outliers in new data

- ▶ **Content personalization**

Personalize content like product recommendations

- ▶ **Language translation and summarization**

Translate text or summarize long passages

- ▶ **Compliance**

Analyze contracts or other documents for compliance

# Operationalize AI with Red Hat OpenShift AI



# What is Machine Learning?

Machine learning can solve business problems

Artificial Intelligence (AI)

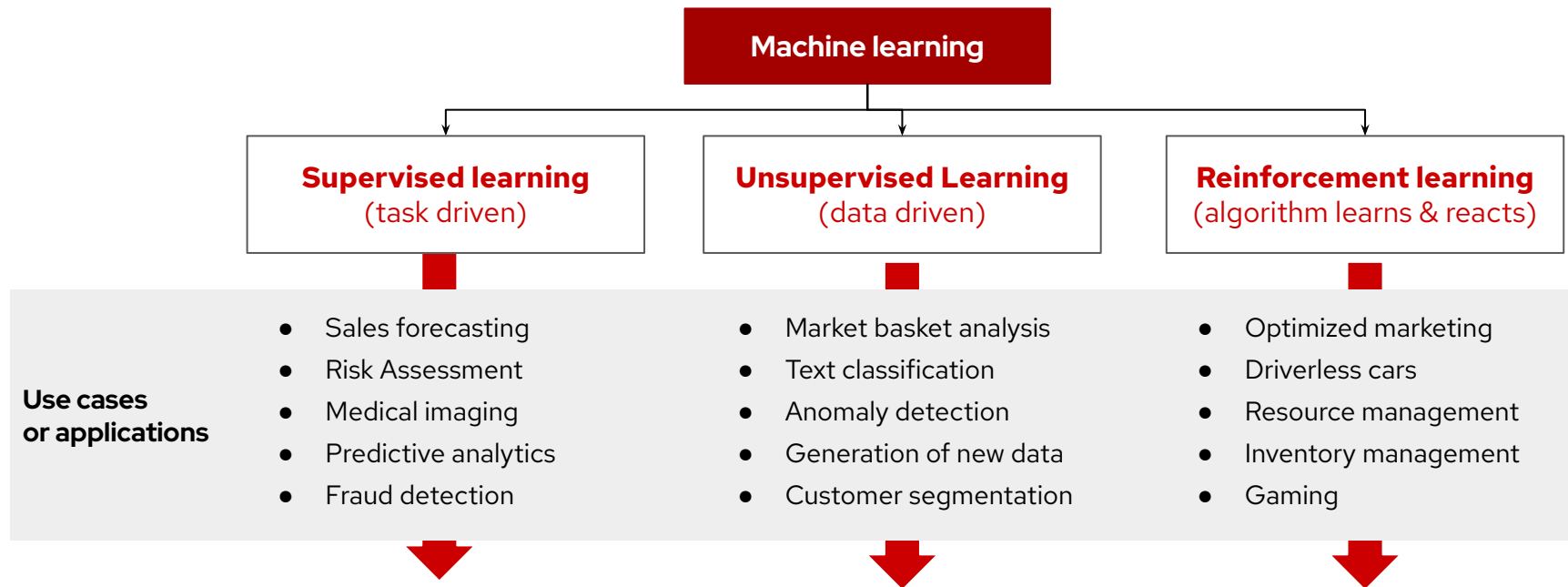
Machine Learning (ML)

**Machine Learning** is a subset of AI:

- ▶ This technique ***empowers computers to learn from data*** and improve their performance with time.
  - Training statistical models to extract knowledge and patterns from data
  - Training is done using supervised or unsupervised learning
  - ML models can make accurate predictions and decisions

# What are the different types of Machine Learning?

Each type conglomerates a variety of common algorithms



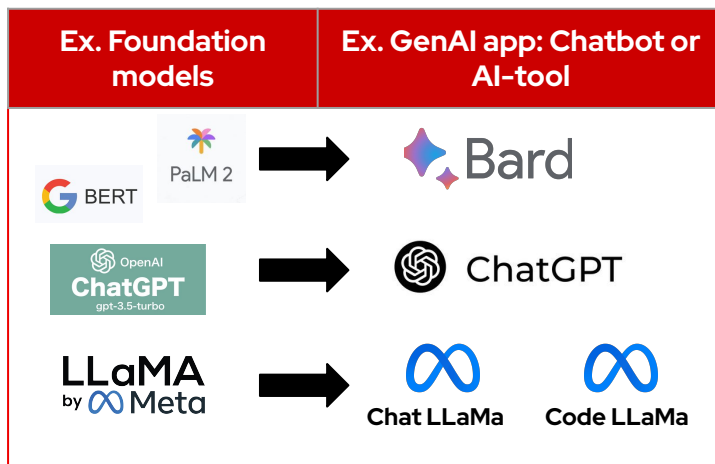


# Generative AI applications are powered by foundation models

Foundation models allow developing specialized AI-enabled applications

## Benefits of foundation models:

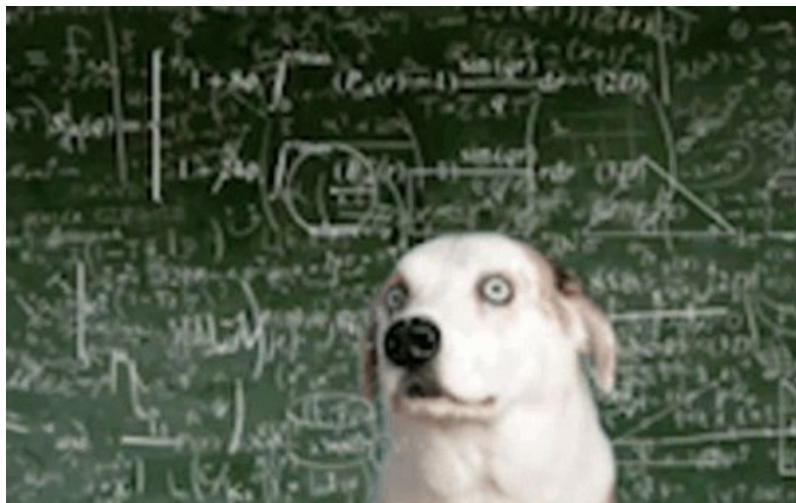
- **Time to value** - alleviates the cost of compute and people
- **Accuracy** - increases with the amount of data use during training
- **Accessibility** - makes advanced AI capabilities available to non-experts
- **Versatility** - offers support for a wide range of tasks and applications



## Most common Gen AI applications:

- Text summarization
- Text generation (including code)
- Sentiment analysis
- Classification
- Conversational questions and answers

It's not magic.  
It's math.



**All of the amazing things that AI and Generative AI can do all comes down to mathematical computation.**

- Compute intensive
- Storage intensive
- There are no small workloads
- Quota attainment

# Way to production

It all starts with some models...



# Real Life View of Technical Teams on AI\*

\*gathered from real life experience in EMEA ;)



Legacy  
Monolith



Modern  
Microservices

# Real Life View of Technical Teams on AI\*

\*gathered from real life experience in EMEA ;)



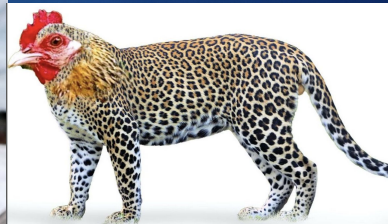
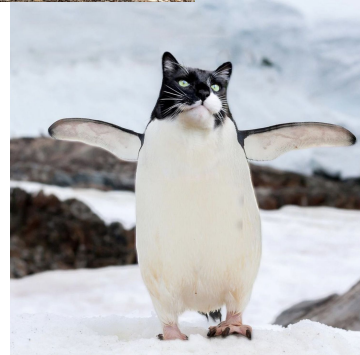
Legacy  
Monolith



Modern  
Microservices

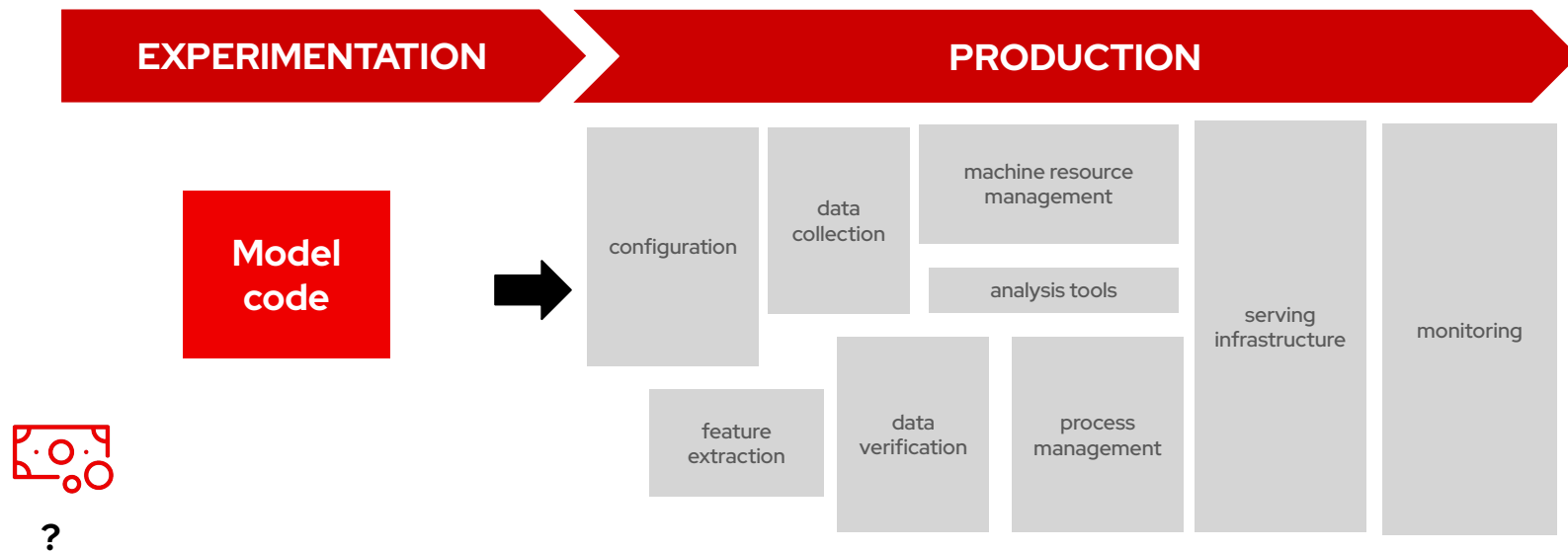


AI



# Poorly designed systems lead to failed ML projects

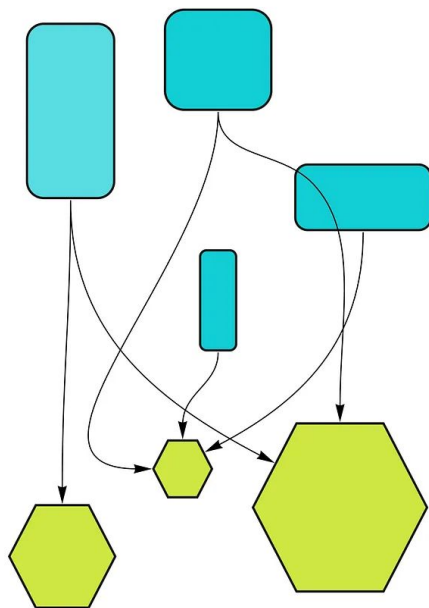
Lack of focus on end-to-end system builds technical debt



Technical debt is a barrier to production

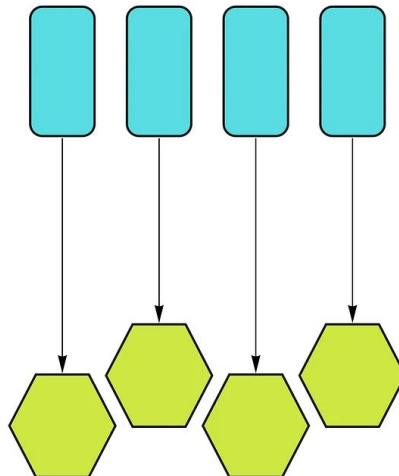
# Expected Outcome

## Typical mapping between systems and teams



## Mapping between systems and teams in a Team-First approach

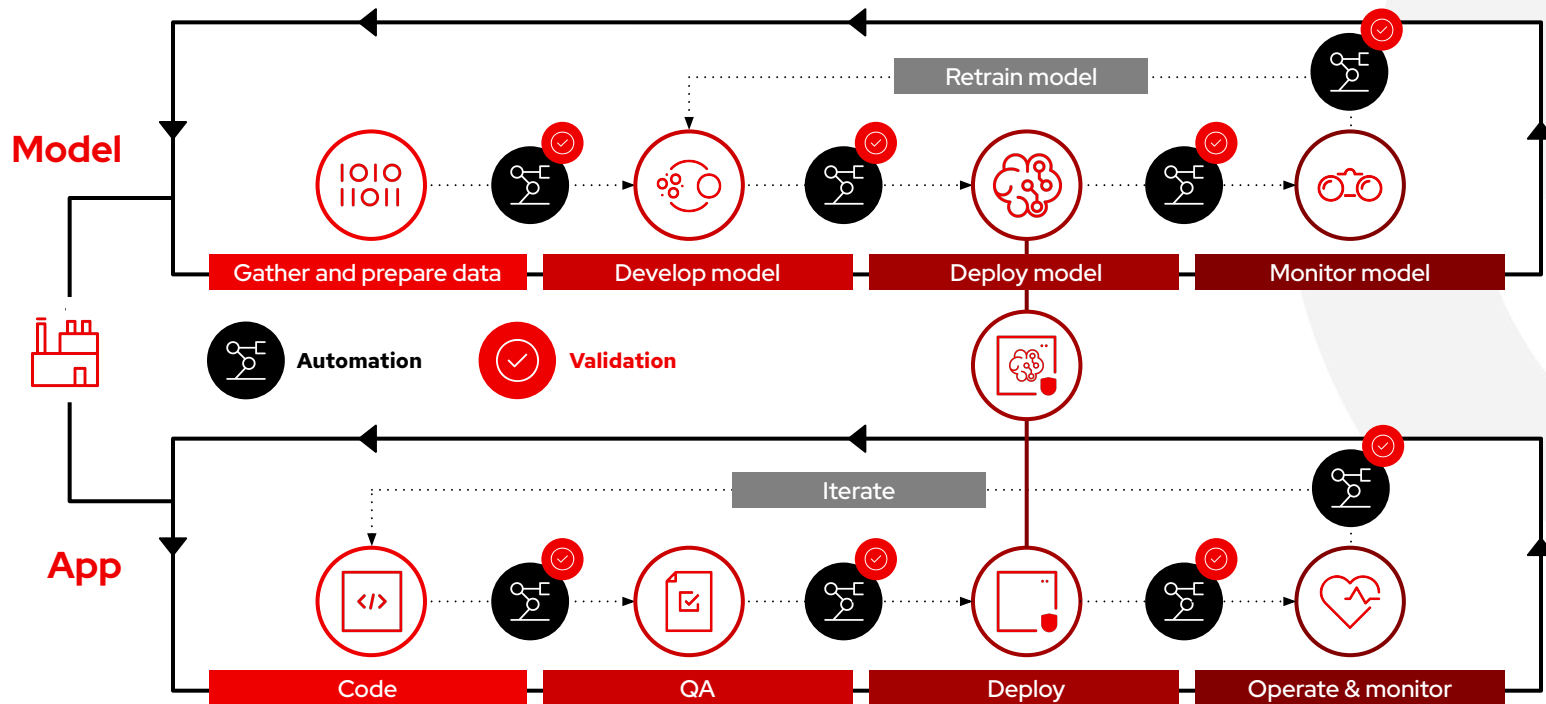
Systems



Teams



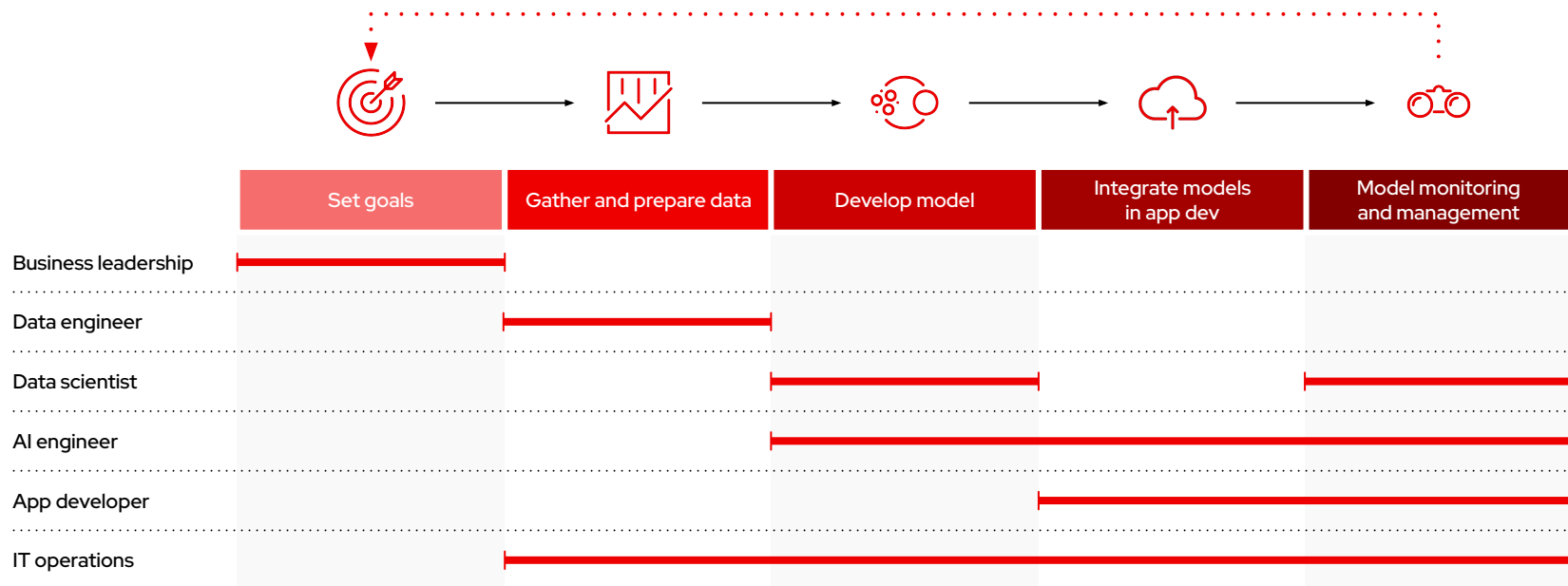
# Lifecycle for operationalizing (containerised) models





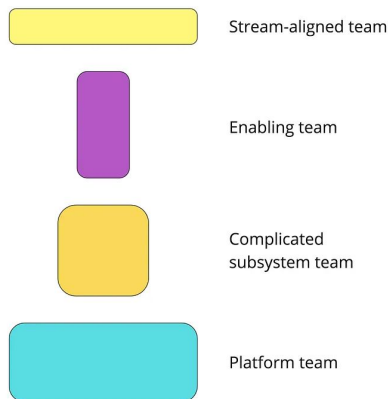
# Operationalizing AI/ML requires collaboration

Every member of your team plays a critical role in a complex process

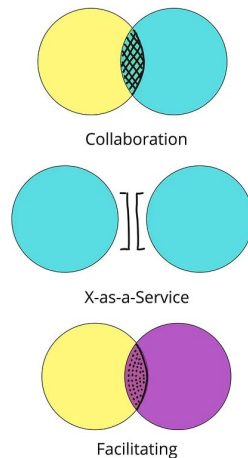


# Teams

## 4 team types



## 3 interaction modes



## 1. Stream-aligned teams

aligned to a single delivery stream, such as a product or service (what others might call a “product team” or a “feature team”).

## 2. Enabling teams

specialists in a particular domain that guide stream-aligned teams

## 3. Complicated-subsystem teams

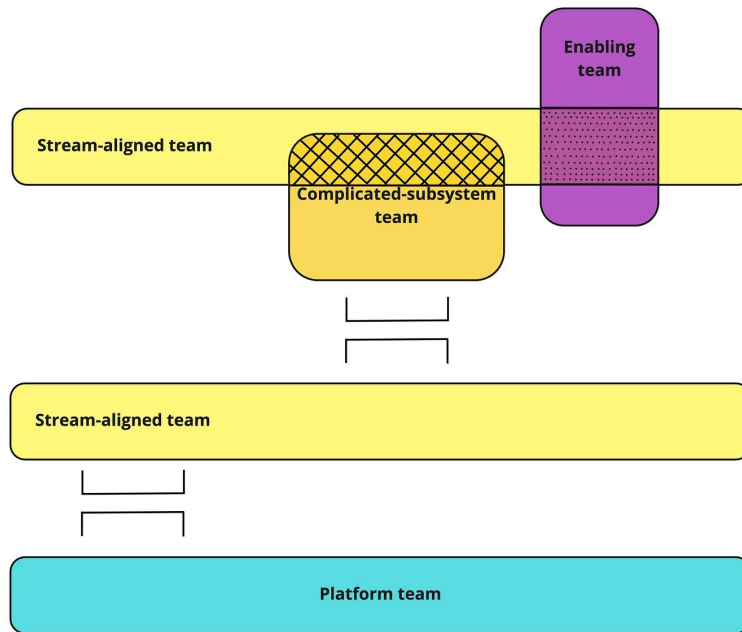
maintain a particularly complex subsystem, such as an ML model

## 4. Platform teams

provide internal services like deployment platforms or data services

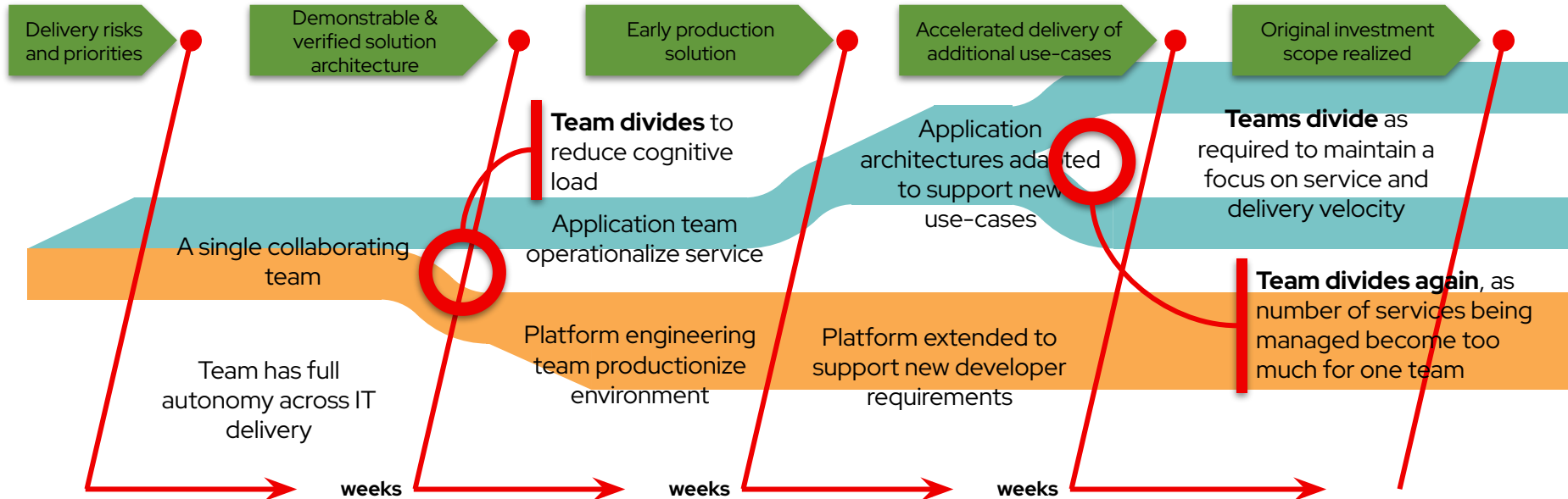


# Typical interactions between different teams



# Red Hat recommends an evolutionary approach to organisational change

Organisational change is seeded through delivery of specific services, and designed to scale as required



Team Topologies: Organizing Business and Technology Teams for Fast Flow, Pias & Skelton  
ISBN: 9781942788812

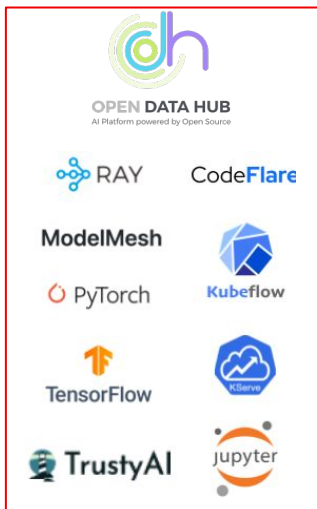
Red Hat's approached are informed by, and align with, Team Topologies

Version number here V00000

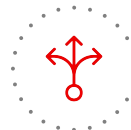


# Simplify AI adoption

Designed to increase AI adoption and enhance trust in AI initiatives

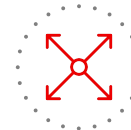


**Red Hat**  
OpenShift AI



**Flexible**

A composable platform  
for rapid dev and  
delivery of AI-enabled  
apps



**Expand**

A certified AI partner  
ecosystem for  
delivering an E2E AI/ML  
experience

# Red Hat OpenShift AI – Key features

## Model development

**Interactive, collaborative UI** for exploratory data science, and model training, tuning and serving

## 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, tuning and serving

# Build an AI platform for E2E AI lifecycle management



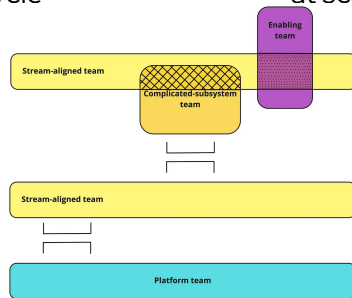
**Red Hat's**  
AI Partner Ecosystem



Trusted, comprehensive and consistent hybrid application platform for managing the entire application lifecycle

Open hybrid AI/ML platform, built on top of OpenShift, to create and deliver AI-enabled apps securely at scale across hybrid-clouds

Best-of-breed AI technologies from a certified partner ecosystem to complement or extend Red Hat's AI capabilities



## Everything in RHOAI has an OpenShift representation

Models and model servers [Deploy model](#) Single-model serving enabled

Model name ↑	Serving runtime	Inference endpoint	Status
▼ My Model ⓘ	OpenVINO Model Server		✓
Framework	onnx-1		
Model server replicas			
Model server size	Small		
	1 CPUs, 4Gi Memory requested		
	2 CPUs, 8Gi Memory limit		
Accelerator	None		
Number of accelerators	0		



Project: test ▼

[InferenceServices](#) > InferenceService details

**IS my-model**

Details [YAML](#)

```
1  apiVersion: serving.kserve.io/v1beta1
2  kind: InferenceService
3  > metadata: --
6  spec:
7    predictor:
8      maxReplicas: 2
9      minReplicas: 2
10     model:
11       modelFormat:
12         name: onnx
13         version: '1'
14       name: ''
15       resources: {}
16       runtime: my-model
17       storage:
18         key: aws-connection-abc
19         path: mymodel/v01
20 > status: --
61
```



# UI to Yaml

## GitOps

The diagram illustrates the GitOps workflow for deploying a model. It consists of three main components:

- GitHub Repository:** A screenshot of a GitHub repository showing the `model.yaml` file. The file contains the following YAML content:

```
24 ---
25 apiVersion: serving.kserve.io/v1beta1
26 kind: InferenceService
27 metadata:
28   annotations:
29     openshift.io/display-name: img-det
30     serving.kserve.io/deploymentMode: ModelServer
31   labels:
32     name: "img-det"
33     opendatahub.io/dashboard: "true"
34   name: "img-det"
35 spec:
36   predictor:
37     model:
38       modelFormat:
39         name: onnx
40         version: '1'
41       runtime: ovms
42       storage:
43         key: aws-connection-minio
44         path: accident/
```

**Terminal Window:** A terminal window shows the deployment command being executed:

```
1 apiVersion: serving.kserve.io/v1beta1
2 kind: InferenceService
3 > metadata:
4   annotations:
5     openshift.io/display-name: img-det
6     serving.kserve.io/deploymentMode: ModelServer
7   labels:
8     name: "img-det"
9     opendatahub.io/dashboard: "true"
10   name: "img-det"
11 spec:
12   predictor:
13     model:
14       modelFormat:
15         name: onnx
16         version: '1'
17       runtime: ovms
18       storage:
19         key: aws-connection-minio
20         path: mymodel/
21 > status: --
```

**Console Window:** A console window shows the deployment details for the `my-model` project. The `Models and model servers` section displays the following information:

Model name	Serving runtime	Inference endpoint	Status
My Model	OpenVINO Model Server		

The `Single-model serving enabled` toggle is turned on. The `Deploy model` button is visible.

**Workflow:** The workflow is represented by two red arrows labeled `sync`. The first arrow points from the `model.yaml` file in the GitHub repository to the terminal window, indicating the deployment command. The second arrow points from the terminal window to the console window, indicating the deployment details.

# What is Red Hat OpenShift AI (RHOAI) solving

- **MLOps**

- RHOAI helps you build out an enterprise grade AI and MLOps platform to create and deliver GenAI and predictive models by providing supported AI tooling on top of OpenShift.
- It's based on OpenShift, a container based application platform that efficiently scales to handle workload demands of AI operations and models.
- You can run your AI workloads across the hybrid cloud, including edge and disconnected environments.

- **Unified app platform**

- OpenShift supports the end-to-end application lifecycle. RHOAI extends OpenShift to AI models, getting them into to AI models and getting them into production with OpenShift best practices.
- Seamless collaboration across multiple personas including IT Ops, Data scientists and application developers by providing a unified platform.

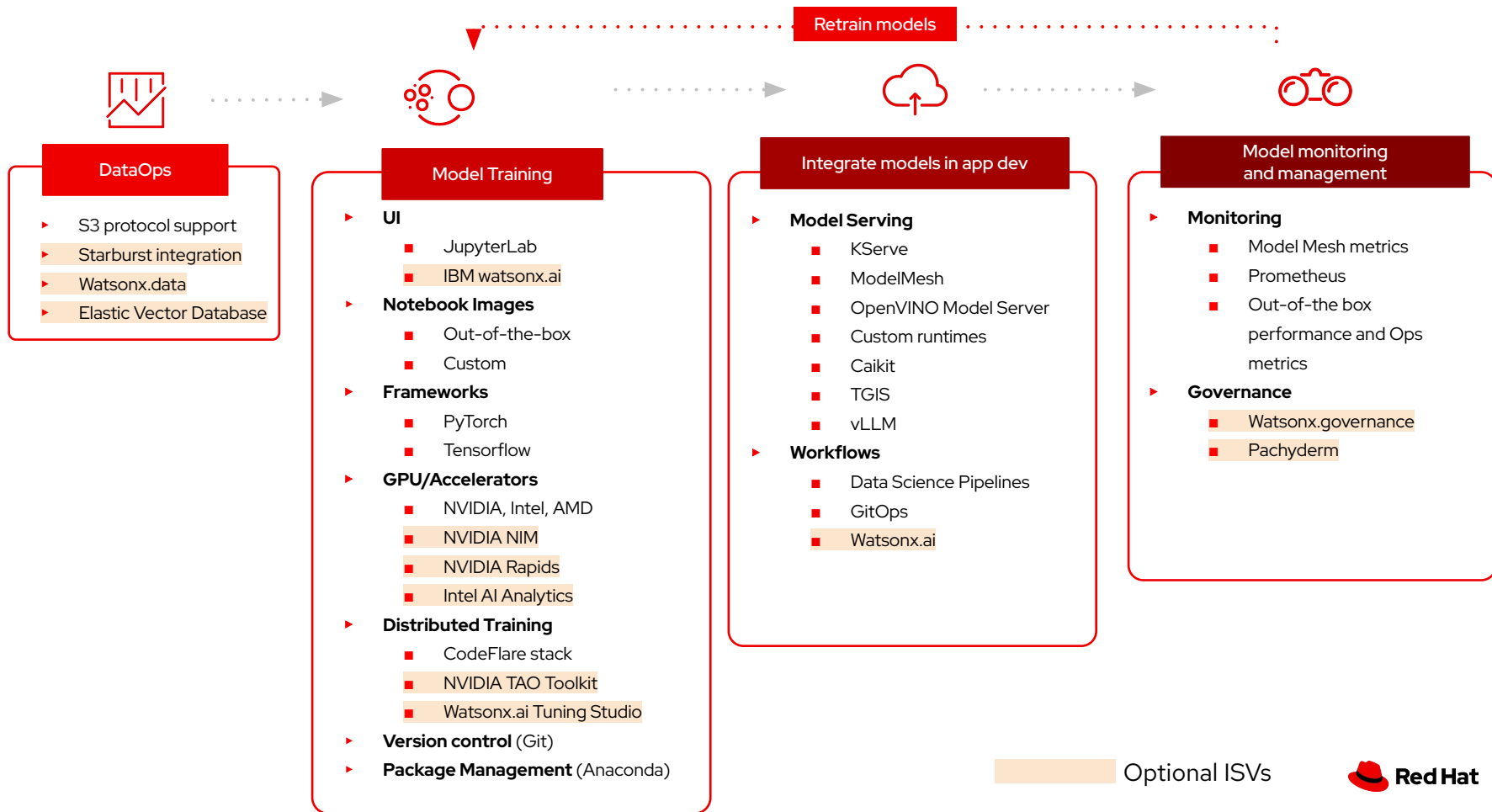
- **Extensibility**

- RHOAI is built to be modular, allowing for a customizable AI/ML stack where you can plug in partners or open source software and technologies where needed to build out an MLOps platform that fits your organization.

- **No vendor lock-in**

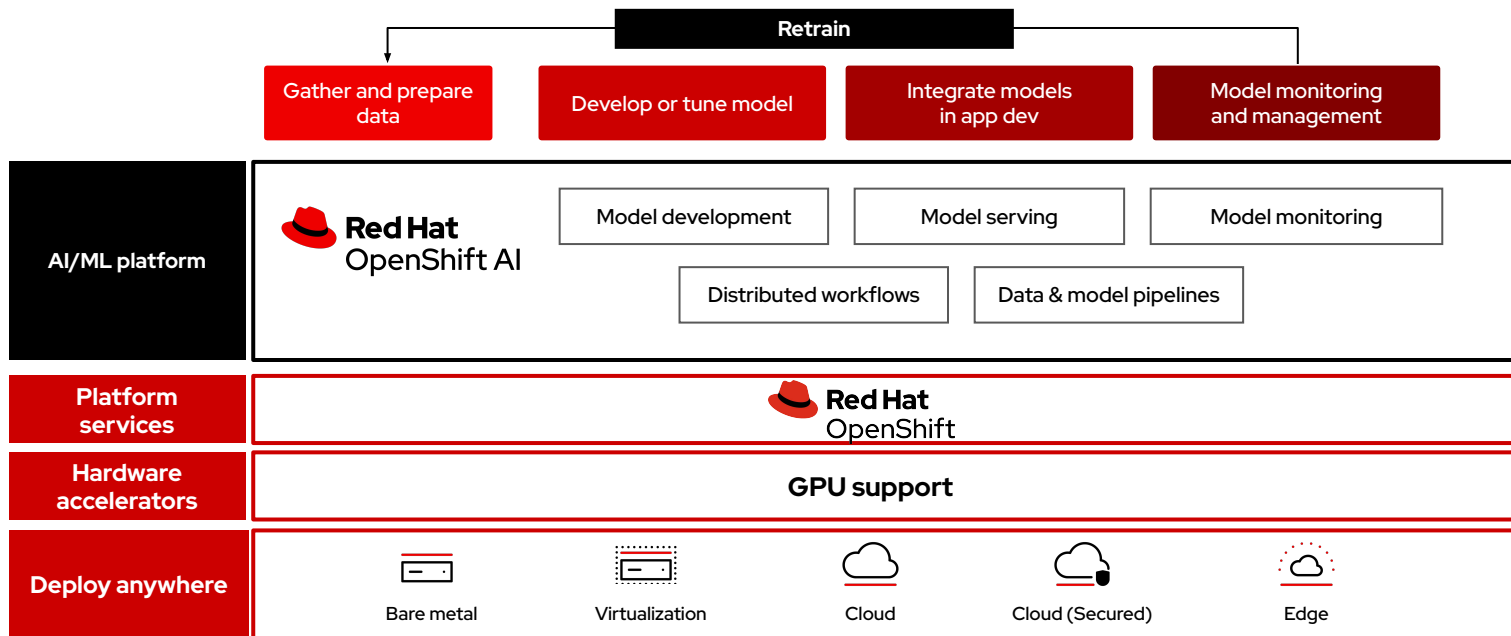
- Thanks to being modular and able to run across the hybrid cloud, you have the freedom to migrate and extend as needed, allowing you to keep up with the speed of AI innovation.

# Product (and Open Source Projects) Overview



# Red Hat OpenShift AI

Red Hat's AI/ML platform



# Red Hat OpenShift AI

Dashboard Application

Data Science Projects

Admin Features

Model Registry

Object Storage



## Model Development & Training

### Workbenches

- Minimal Python
- PyTorch
- CUDA
- Standard Data Science
- TensorFlow
- VSCode
- RStudio
- TrustyAI

CodeFlare SDK

ISV images

Custom images

### Distributed workloads

KubeRay

CodeFlare

Data and model  
Pipelines

## Model Serving

### Serving Engines

Kserve

ModelMesh

### Serving Runtimes

OVMS (built-in)

Caikit/TGIS (built-in)

Custom

## Model Monitoring

Performance metrics

Model explainers

Quality metrics

OpenShift  
Operators

OpenShift  
GitOps



OpenShift  
Pipelines



OpenShift  
ServiceMesh

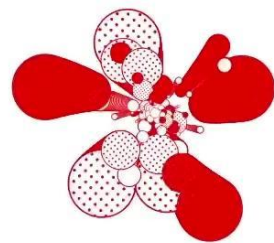


OpenShift  
Serverless



Prometheus





# Red Hat OpenShift cloud services

A turnkey application platform  
with management and support  
from Red Hat and leading cloud  
providers



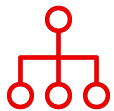
## Accelerate time to value

Quickly build, deploy, and run applications  
that scale as needed.



## Operational efficiency

Enhance operational consistency, efficiency and  
security with proactive management and support.



## Focus on innovation

Simplify operations so your teams can refocus  
on innovation, not managing infrastructure.



## Hybrid cloud flexibility

Deliver a consistent experience on premises  
and in the cloud.

# Flavors of RHOAI

Supported deployment options		
Options available	Self-managed RHOAI	Cloud Service RHOAI
Bare metal	✓	
Virtual	✓	
Private cloud	✓	
Red Hat OpenShift on AWS (ROSA)	✓	✓
Azure Red Hat OpenShift (ARO)	✓	(future)
IBM Cloud	✓	
OSD-GCP/OSD-AWS	✓	✓
Edge	(future)	



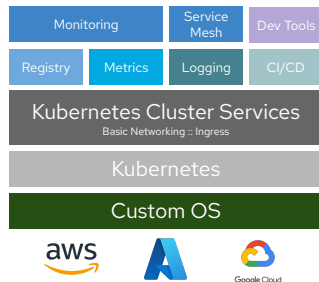
# Flavors of RHOAI

Supported deployment options		
Options available	Self-managed RHOAI	Cloud Service RHOAI
Bare metal	✓	
Virtual	✓	
Private cloud	✓	
Red Hat OpenShift on AWS (ROSA)	✓	✓
Azure Red Hat OpenShift (ARO)	✓	(future)
IBM Cloud	✓	
OSD-GCP/OSD-AWS	✓	✓
Edge	(future)	

# Build and run a platform *versus* using a turnkey cloud service



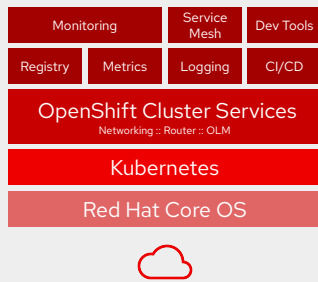
The Parts



xKS + 'native'  
services



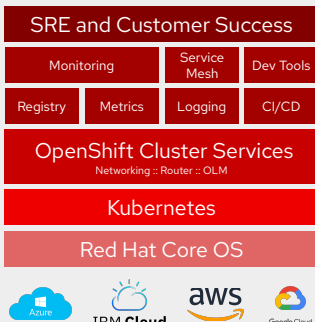
The Assembled Car



- Application Platform -  
Self-managed Red Hat OpenShift



The Car & Pit Crew



- Turnkey Application Platform -  
Red Hat OpenShift cloud services

## "Batteries Included"

... but swappable

Individual components can be swapped out

Eg.

- Using AWS CloudWatch for logging on AWS
- Use specific cloud services or ISV offerings

# Move from 24x7 operations to 9-5 innovation

## End-to-End support for your entire application platform

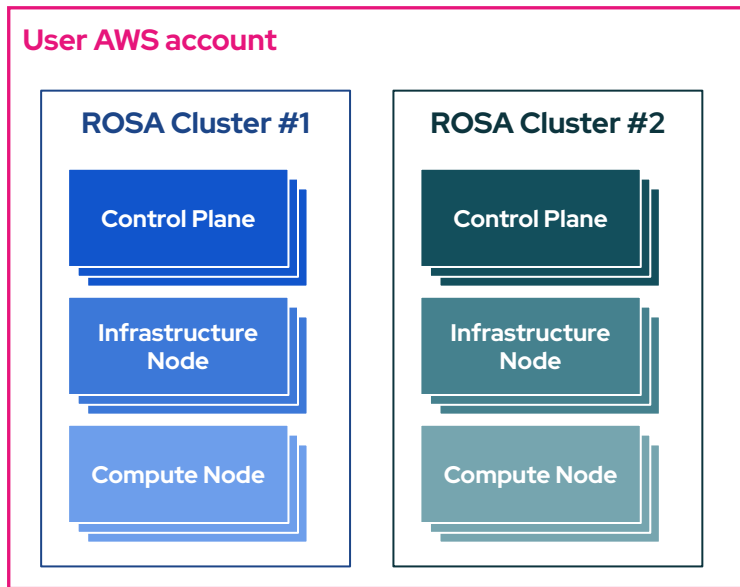


- ▶ OpenShift cloud services includes **full support for worker nodes**
  - Zero downtime upgrades,
  - proactive monitoring
  - automated patching
  - **Compliance and certifications** extend to worker nodes
- ▶ **99.95%** financially backed SLA
- ▶ **24x7 joint support** from Red Hat and cloud provider
- ▶ Automation and **Day 2 Operations** by **global SREs**

# ROSA Variants

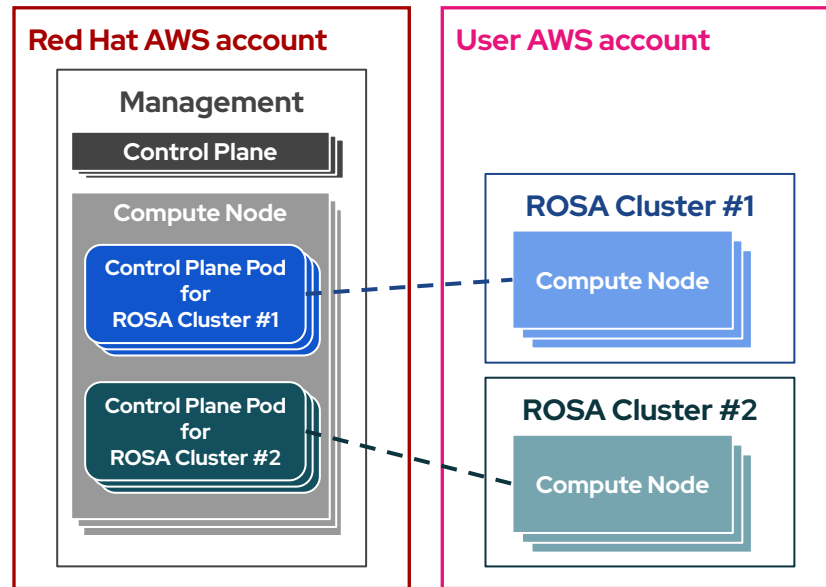
Generally Available

## ROSA Classic



1 or 3 AZs per MachinePool

## ROSA with Hosted Control Plane (HCP)



Cluster per Region (each 3 AZs)



## Reduce AWS infrastructure costs

### Minimum configuration of ROSA Classic

#### User AWS account

#### ROSA Cluster

Control Plane #1

Control Plane #2

Control Plane #3

Infrastructure Node #1

Infrastructure Node #2

Compute Node #1

Compute Node #2

Control Plane Infrastructure Node reduces "total cost of ownership" (TCO)

\$



Classic



Hosted Control Plane

~\$ /year  
cost savings

# Complexity of running your own Kubernetes Cluster

## Responsibilities

User management

Project and quota management

Application life cycle

Cluster creation

Cluster management

Monitoring and logging

Network configuration

Software and security updates

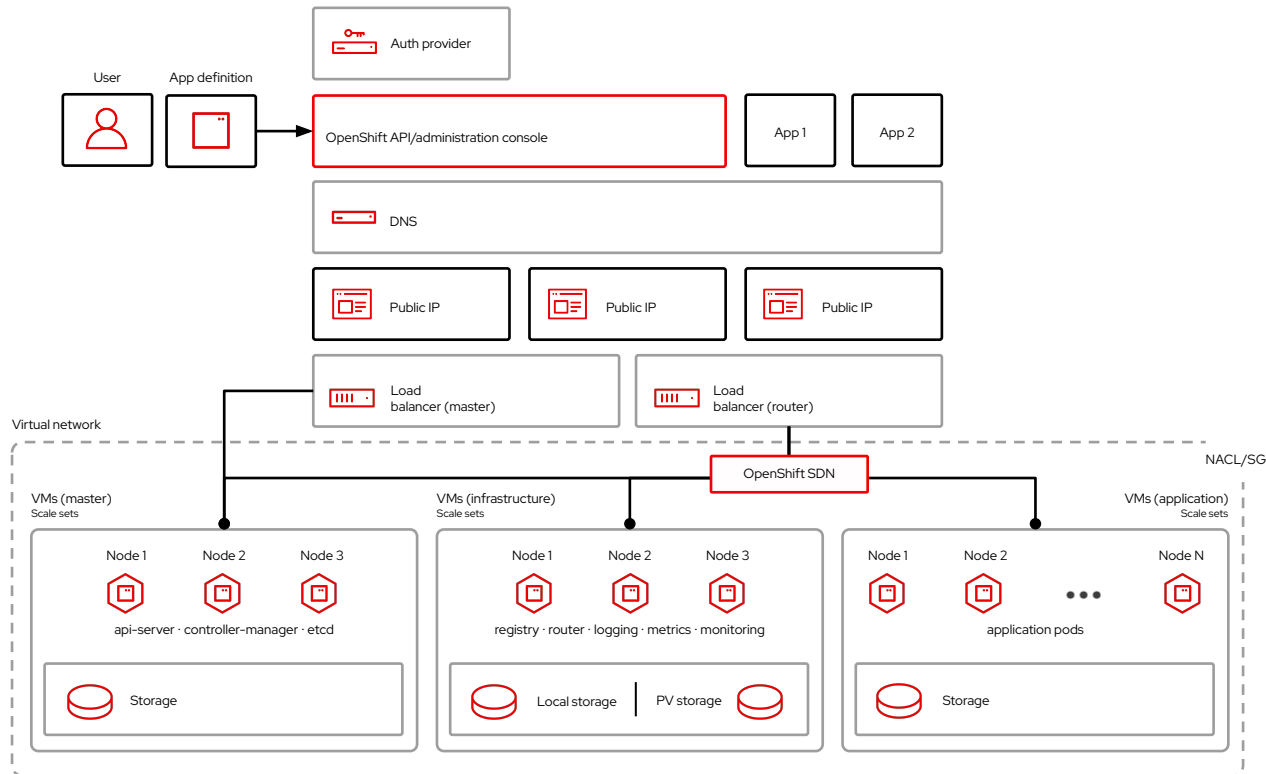
Platform support



Customer












Cloud provider & Red Hat



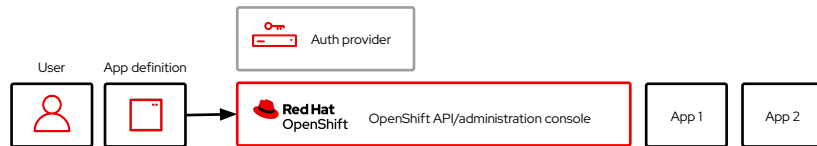
# Simplify with fully managed clusters

## Red Hat OpenShift cloud services

Responsibilities	
User management	
Project and quota management	
Application life cycle	
Cluster creation	
Cluster management	
Monitoring and logging	
Network configuration	
Software and security updates	
Platform support	

 Customer

 Cloud provider and Red Hat



### Let Red Hat & AWS...

Manage all your clusters.

Monitor and operate your VMs.

Secure your nodes.

Manage environment patches.

### You...

# Compliance

Red Hat OpenShift Service on AWS follows common industry best practices for security and controls. The certifications are outlined in the following table.

*Table 1. Security and control certifications for Red Hat OpenShift Service on AWS*

Compliance	Red Hat OpenShift Service on AWS (ROSA)	Red Hat OpenShift Service on AWS (ROSA) with hosted control planes (HCP)
HIPAA Qualified	Yes	No
ISO 27001	Yes	Yes
ISO 27017	Yes	Yes
ISO 27018	Yes	Yes
PCI DSS	Yes	Yes
SOC 1 Type 2	Yes	Yes
SOC 2 Type 2	Yes	Yes
SOC 3	Yes	Yes
FedRAMP High <sup>[1]</sup>	Yes ( <a href="#">GovCloud</a> requisite)	No


[https://docs.openshift.com/rosa/rosa\\_architecture/rosa\\_policy\\_service\\_definition/rosa-policy-process-security.html](https://docs.openshift.com/rosa/rosa_architecture/rosa_policy_service_definition/rosa-policy-process-security.html)



# Where do we start?



# Install RHOAI



Administrator

Home

Operators

OperatorHub

Installed Operators

Workloads

Serverless

Networking

Storage

Builds

Project: All Projects

## OperatorHub

Discover Operators from the Kubernetes community and Red Hat partners, curated by Red Hat. You can purchase commercial software through [Red Hat Marketplace](#) services to your developers. After installation, the Operator capabilities will appear in the [Developer Catalog](#) providing a self-service experience.

All Items

A list of comma separated categories that your operator falls under.

AI/Machine Learning

Application Runtime

Big Data

Cloud Provider

Database


Developer Tools

Development Tools

Drivers and plugins

Integration & Delivery

All Items





Red Hat

Red Hat OpenShift AI

provided by Red Hat

Operator for deployment and...

 Installed

 Red Hat

## Everything in RHOAI has an OpenShift representation

Models and model servers [Deploy model](#) Single-model serving enabled

Model name ↑	Serving runtime	Inference endpoint	Status
▼ My Model ⓘ	OpenVINO Model Server		✓
Framework	onnx-1		
Model server replicas			
Model server size	Small		
	1 CPUs, 4Gi Memory requested		
	2 CPUs, 8Gi Memory limit		
Accelerator	None		
Number of accelerators	0		



Project: test ▼

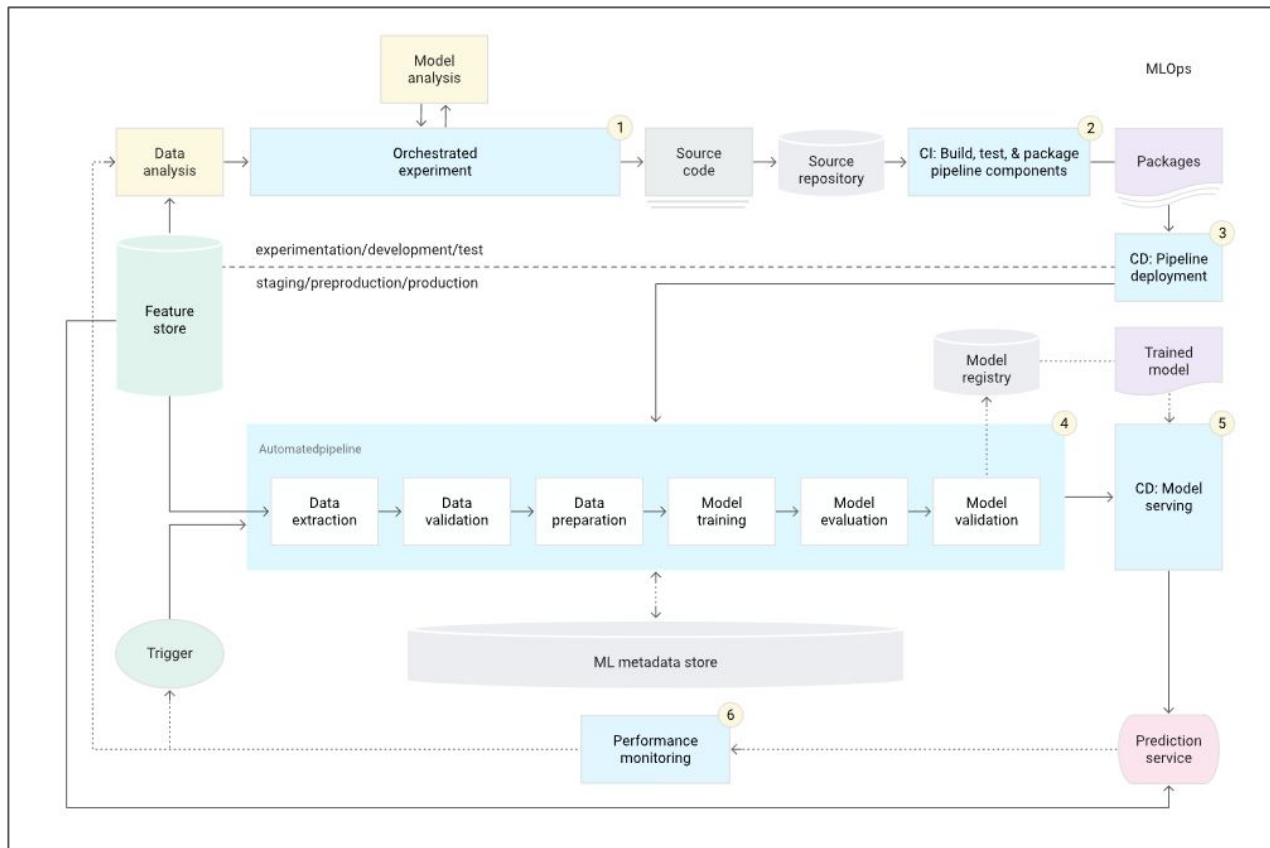
[InferenceServices](#) > InferenceService details

**IS my-model**

Details [YAML](#)

```
1  apiVersion: serving.kserve.io/v1beta1
2  kind: InferenceService
3  > metadata: --
6  spec:
7    predictor:
8      maxReplicas: 2
9      minReplicas: 2
10     model:
11       modelFormat:
12         name: onnx
13         version: '1'
14       name: ''
15       resources: {}
16       runtime: my-model
17       storage:
18         key: aws-connection-abc
19         path: mymodel/v01
20 > status: --
61
```

## Mature MLOps Flow



# Data Science Projects

Red Hat  
OpenShift AI

Applications

Enabled

Explore

Data Science Projects

Data Science Pipelines

Model Serving

Resources







Settings

Data Science Projects

View your existing projects or create new projects

Data science projects

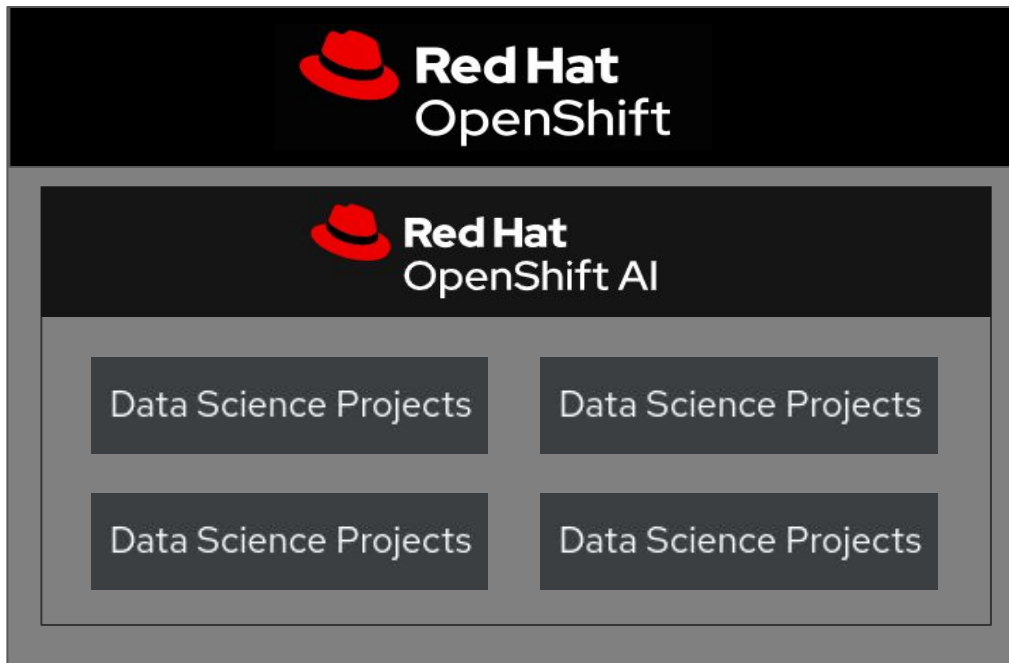
Name Find by name

Project	Created	Workbenches
Name		Name Status
 <a href="#">project-1</a>  panbalag@redhat.com	5/31/2024, 9:06:22 AM	<a href="#">Create a workbench</a> to add a custom notebook. 
 <a href="#">project-2</a>  panbalag@redhat.com	5/31/2024, 9:06:32 AM	<a href="#">Create a workbench</a> to add a custom notebook. 

Data Science projects allow users to **organize** and **manage** contents related to their AI/ML experiments in **isolation** from other projects

science project

# Data Science Projects



- Multiple data science projects.
- Isolation from other projects
- Created by admins or users
- User/Group access privileges

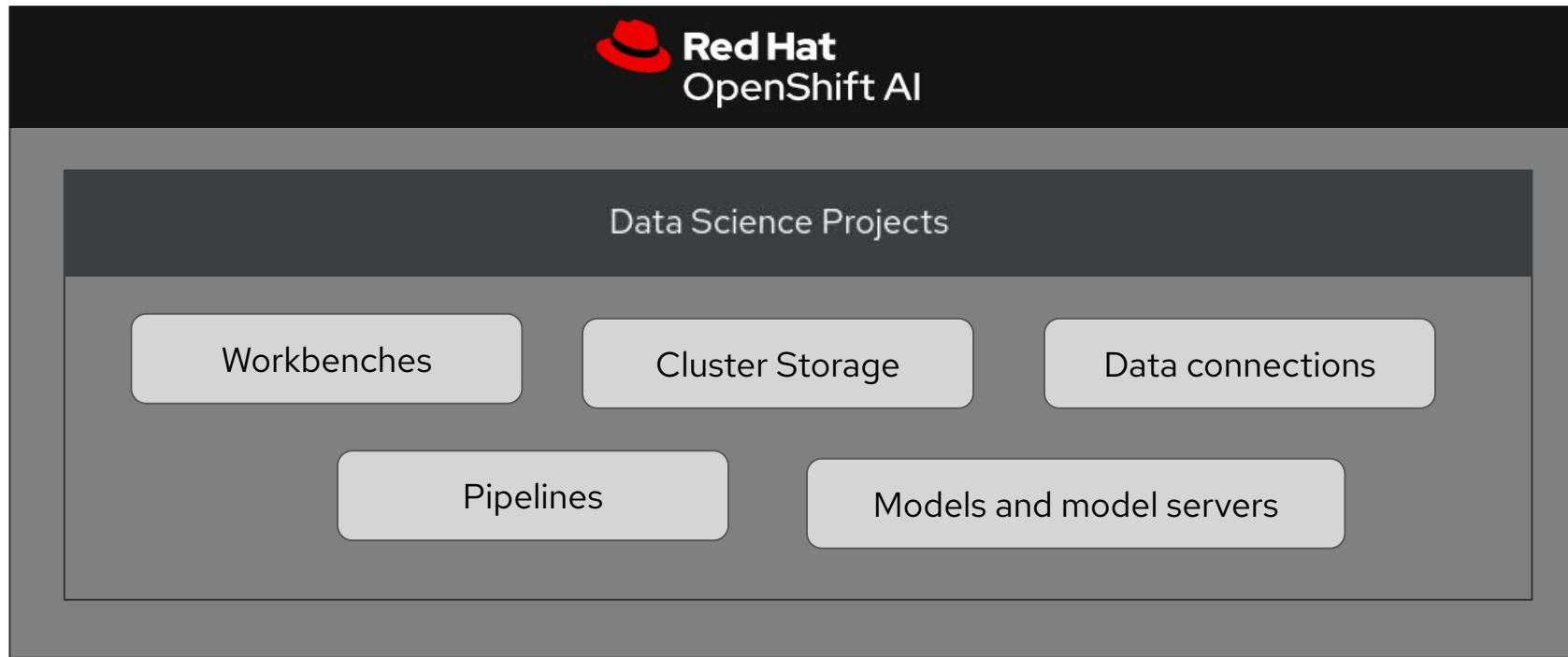
## Collaborate within a project

- Users that create a data science project
  - become an admin of that project
  - can give access to a project to any user or group
- Users with access permissions can access all resources in the project, modify them, and create new ones.
- Limiting user level access to data science projects needs to be handled at an OpenShift level at the moment

## Collaborate between projects

- Due to isolation of data science projects, resources need to be explicitly exposed in order to be shared between projects.
- A good way to do this is to have an external resource which the projects have access to.
  - Examples:
    - A git repository with shared code
    - An object storage with shared artifacts
    - A structured database with shared data

# Data Science Projects





# Data Science Projects

**Red Hat OpenShift AI**

panbalag@redhat.com

**Data Science Projects**

View your existing projects or create new projects.

Data science projects

Name Find by name

Launch Jupyter Create data science project

**Project**

Name	
project-1	panbalag@redhat.com
project-2	panbalag@redhat.com

Resource names and types are used to find your resources in OpenShift.

**Resource name** project-1

**Resource type** Project

Data science projects are 'Projects' in OpenShift identified by the label under 'Resource name'

# Data Science Projects

The screenshot displays the Red Hat OpenShift AI console interface. The left sidebar contains a navigation menu with the following items: Applications, Enabled, Explore, Data Science Projects (highlighted), Data Science Pipelines, Model Serving, and Resources. The main content area is titled "Data Science Projects" and includes a sub-header "View your existing projects or create new projects." Below this, there is a dropdown menu for "Data science projects" and a search bar labeled "Find by name". A table of projects is displayed, with the first project, "project-1", highlighted. A red box and arrow indicate the navigation path from the "Red Hat OpenShift AI" header to the "Data Science Projects" section, then to the "Projects" table, and finally to the "project-1" entry.

**Red Hat OpenShift AI**

**Data Science Projects**

View your existing projects or create new projects.

Data science projects

Name Find by name

Project

Name ↑

project-1 ?  
panbalag@redhat.com

Filter

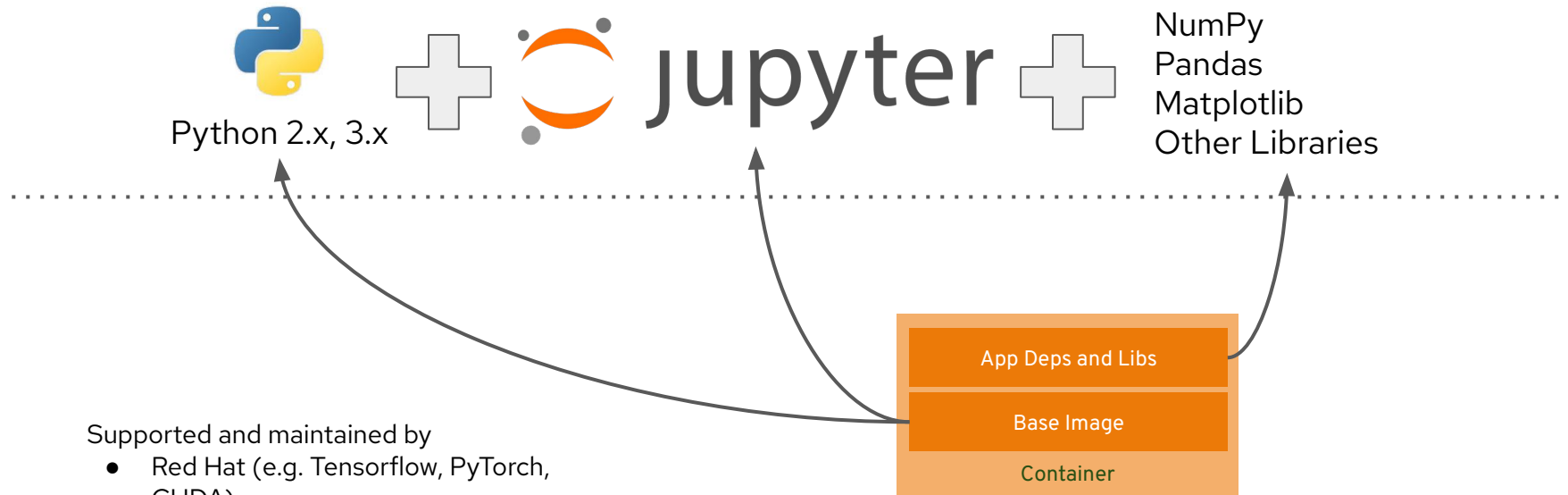
Name Search by name...

Name	Display name	Status
PR project-1	project-1	✓ Active
PR project-2	project-2	✓ Active

# Customizing Workbenches

## Base Notebook Images

Reproducible and shareable environments for building, training and serving



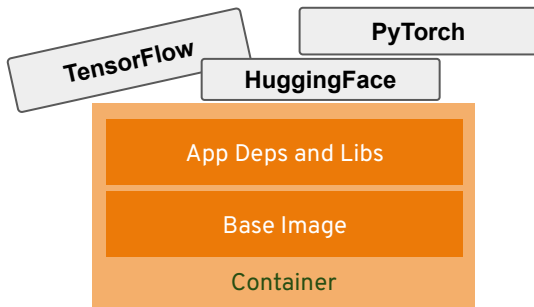
Supported and maintained by

- Red Hat (e.g. Tensorflow, PyTorch, CUDA)
- partner (Anaconda, Intel)
- you (custom notebooks)

# Customizing Workbenches

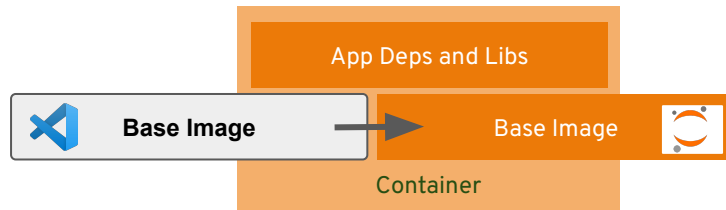
## Customizing the workbench

Adding packages on top of a good image



Just remember that they are removed when restarting the workbench\*

Creating your own custom image with all dependencies you need

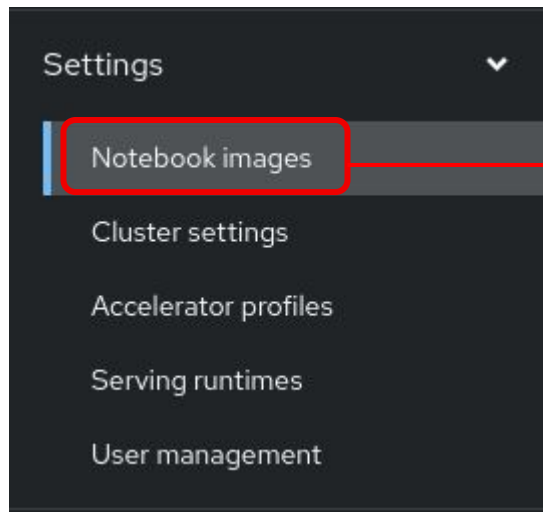


You can now version and maintain it according to your preferences

\* This is on purpose so that you can un-mess-up your environment easily if you get into dependency issues.

# Custom Notebook Image

## Import new image



### Notebook images

Manage your notebook images.

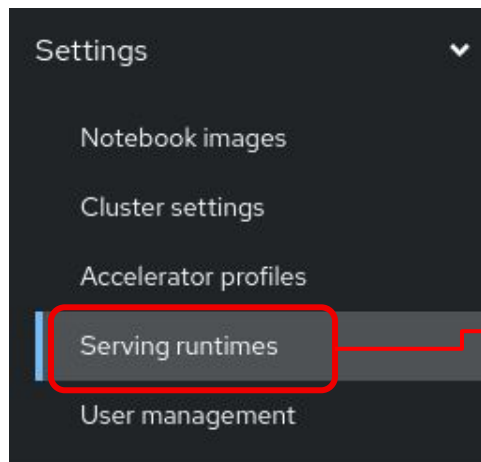
Name  Find by name

Import new image

Name	Description	Enable
> Custom RStudio ?	RStudio	<input checked="" type="checkbox"/>
> Custom VSCode ?	Custom Visual Studio Code	<input checked="" type="checkbox"/>
> Custom PyTorch ?		<input type="checkbox"/>

# Custom Serving Runtime

## Add serving runtime



### Serving runtimes

Manage your model serving runtimes.

Single-model serving enabled

Multi-model serving enabled ?

Add serving runtime

Name

⋮ vLLM ?

⋮ Triton Inference Server 24.01 ?

⋮ NVIDIA NIM ?

# Customize RHOAI Cluster

## Enable or disable components

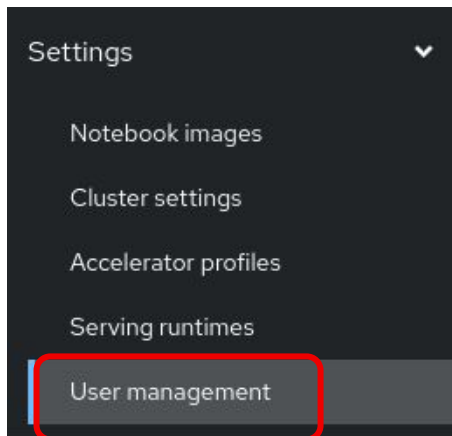
You can enable or disable RHOAI components inside of your DataScienceCluster yaml.

These are the components you can enable/disable:

- CodeFlare (for distributed training)
- Dashboard
- Data Science Pipelines
- Kserve (the component for single-model serving)
- Modelmesh serving (the component for multi-model serving)
- Ray (for distributed training)
- TrustyAI
- Workbenches

```
101 spec:
102   components:
103     codeflare:
104       devFlags: {}
105       managementState: Removed
106     dashboard:
107       devFlags: {}
108       managementState: Managed
109     datasciencepipelines:
110       devFlags: {}
111       managementState: Managed
112     kserve:
113       devFlags: {}
114       managementState: Managed
115       serving:
116         ingressGateway:
117           certificate:
118             secretName: knative-serving-cert
119             type: SelfSigned
120           managementState: Managed
121           name: knative-serving
122     modelmeshserving:
123       devFlags: {}
124       managementState: Managed
125     ray:
126       devFlags: {}
127       managementState: Removed
128     trustyai:
129       devFlags: {}
130       managementState: Removed
131     workbenches:
132       devFlags: {}
133       managementState: Managed
```

# User Management



## User management

Define OpenShift group membership for Data Science administrators and users.

### Data Science administrator groups

Select the OpenShift groups that contain all Data Science administrators.

cluster-admins ✕ dedicated-admins ✕ rhods-admins ✕

View, edit, or create groups in OpenShift under User Management

 All cluster admins are automatically assigned as Data Science administrators.

### Data Science user groups

Select the OpenShift groups that contain all Data Science users.

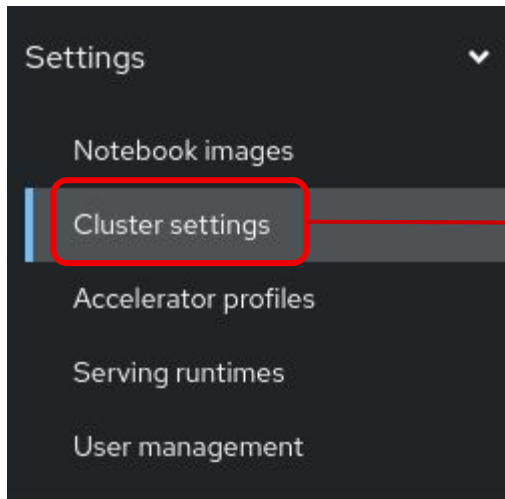
system:authenticated ✕

View, edit, or create groups in OpenShift under User Management

Save changes



# Cluster Settings



1. Model serving platforms
2. PVC size
3. Stop idle notebooks
4. Usage data collection
5. Notebook pod tolerations

# Accelerator Profile

Applications



Enabled

Explore

Data Science Projects

Data Science Pipelines



Pipelines

Runs

Model Serving

Resources

Settings



Notebook images

Cluster settings

Accelerator profiles

Serving runtimes

User management

## Accelerator profiles





Manage accelerator profile settings for users in your organization

Name



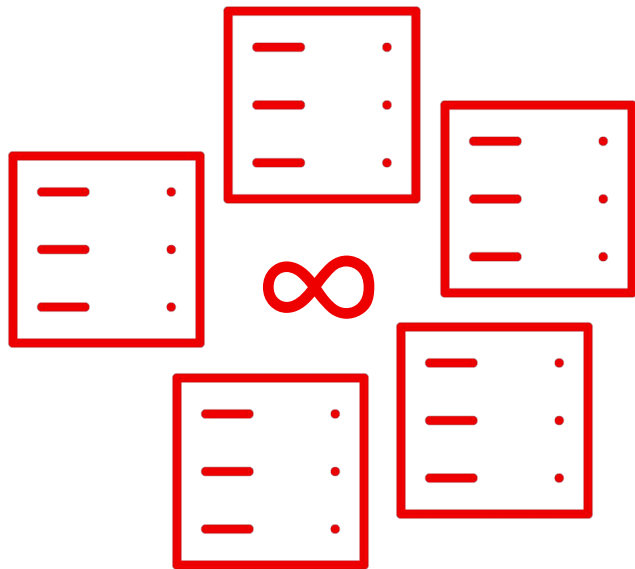
Find by name

Create accelerator profile

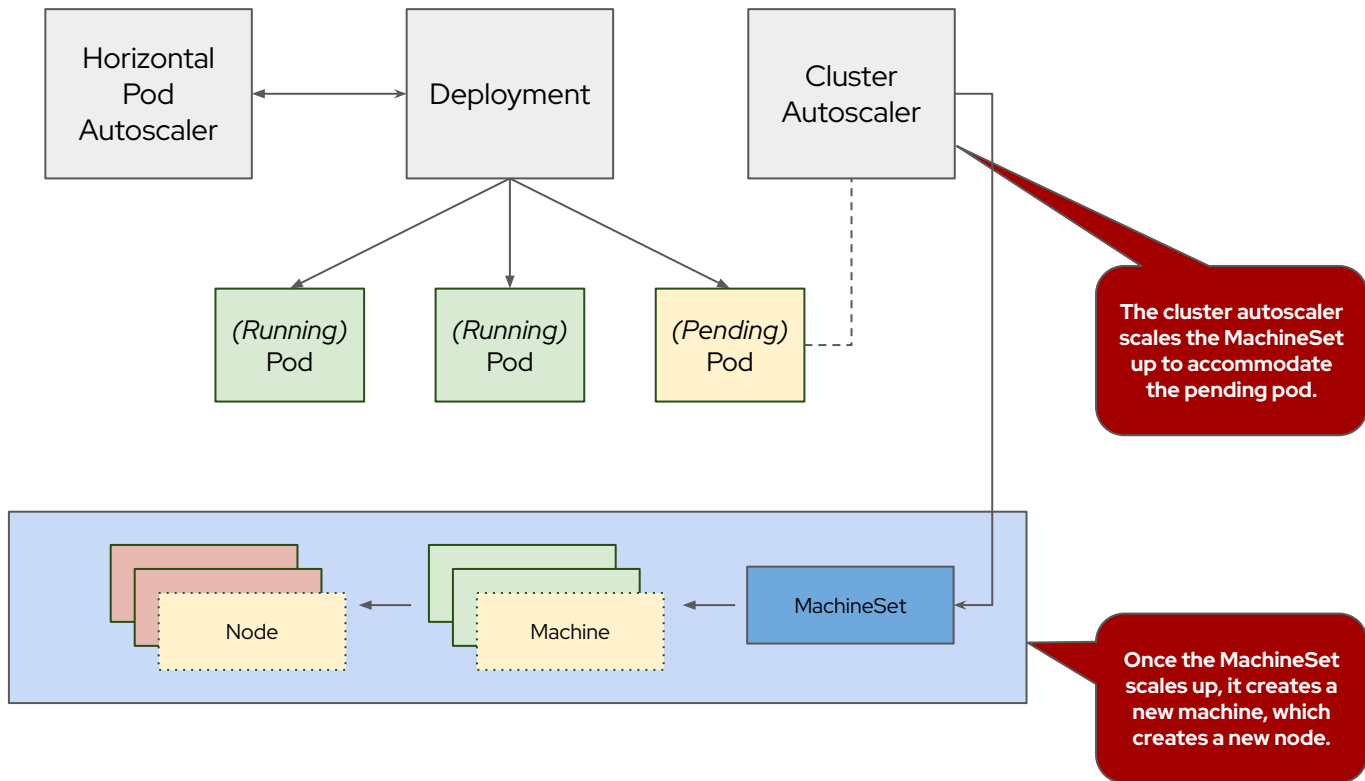
Name 	Identifier  	Enable 
<b>fractional small</b> 1/7th of a real GPU	nvidia.com/gpu-frac	<input type="checkbox"/>
<b>Habana HPU - 1st Gen Gaudi</b> This Accelerator Profile is for 1st Gen Gaudi Devices	habana.ai/gaudi	<input type="checkbox"/>
<b>Large GPU Card</b>	nvidia.com/gpu	<input type="checkbox"/>
<b>NVIDIA GPU - use sparingly</b> We have very few GPUs in this cluster. Although you can use them fo...	nvidia.com/gpu	<input checked="" type="checkbox"/>
<b>tinyGPU</b>	nvidia.com/gpu	<input type="checkbox"/>

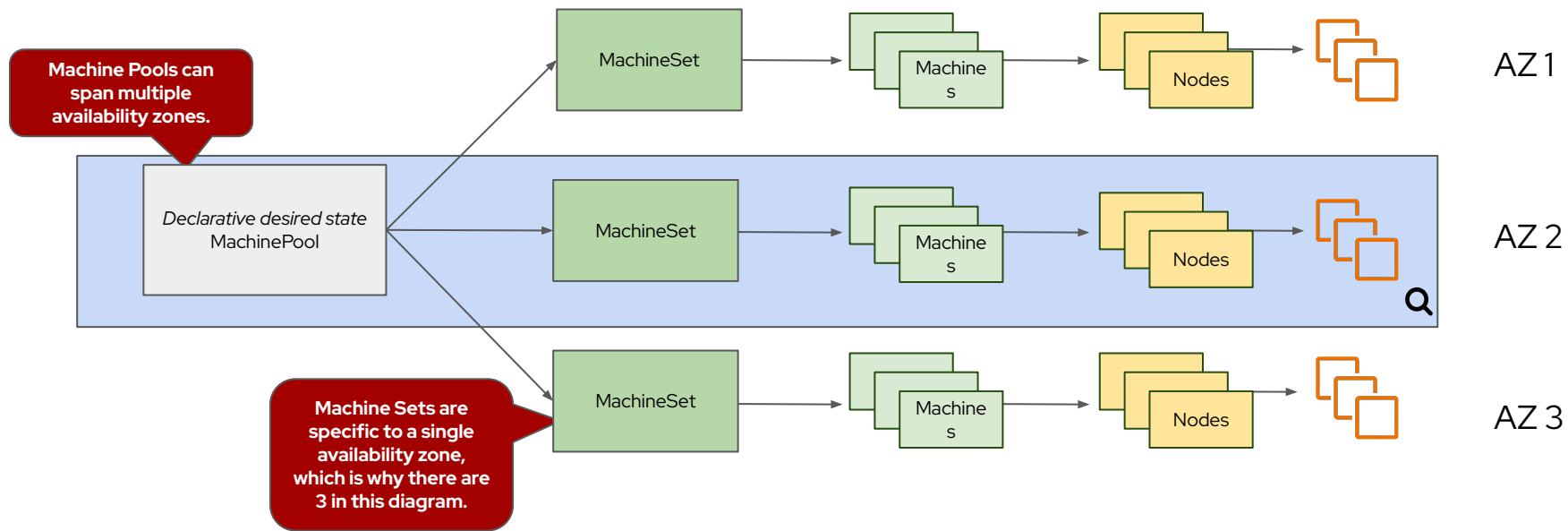
# Cluster Autoscaling

Automatically responding to cluster demand.

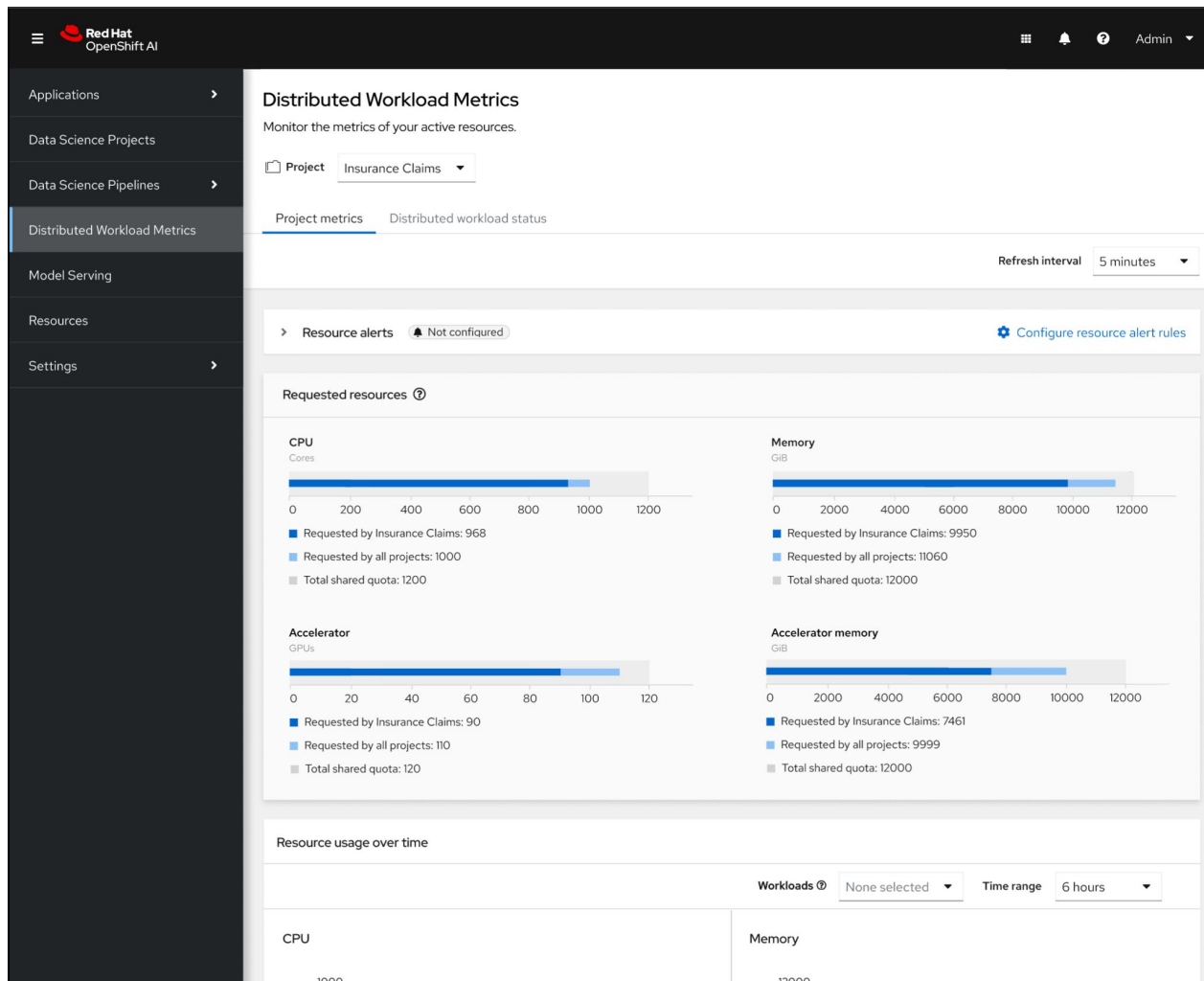


- ▶ MachinePools can be scaled to meet applications demands.
- ▶ Cluster AutoScaler will provision additional worker nodes when pods can not be scheduled due to resource constraints.
- ▶ Cluster AutoScaler will not scale beyond predefined limits.





# Monitoring



## Automatic vs Manual

### Change update approval strategy

What strategy is used for approving updates?

☐ Automatic (default)

New updates will be installed as soon as they become available.

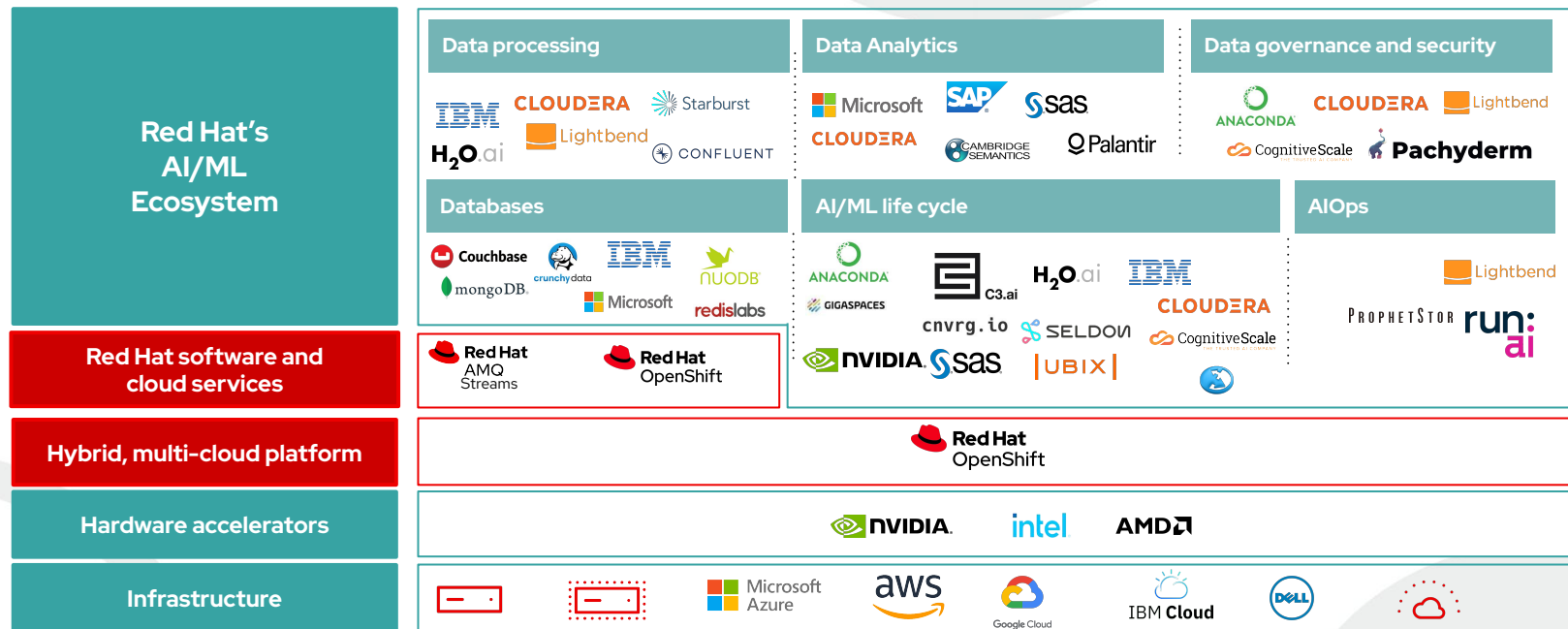
☒ Manual

New updates need to be manually approved before installation begins.

Cancel

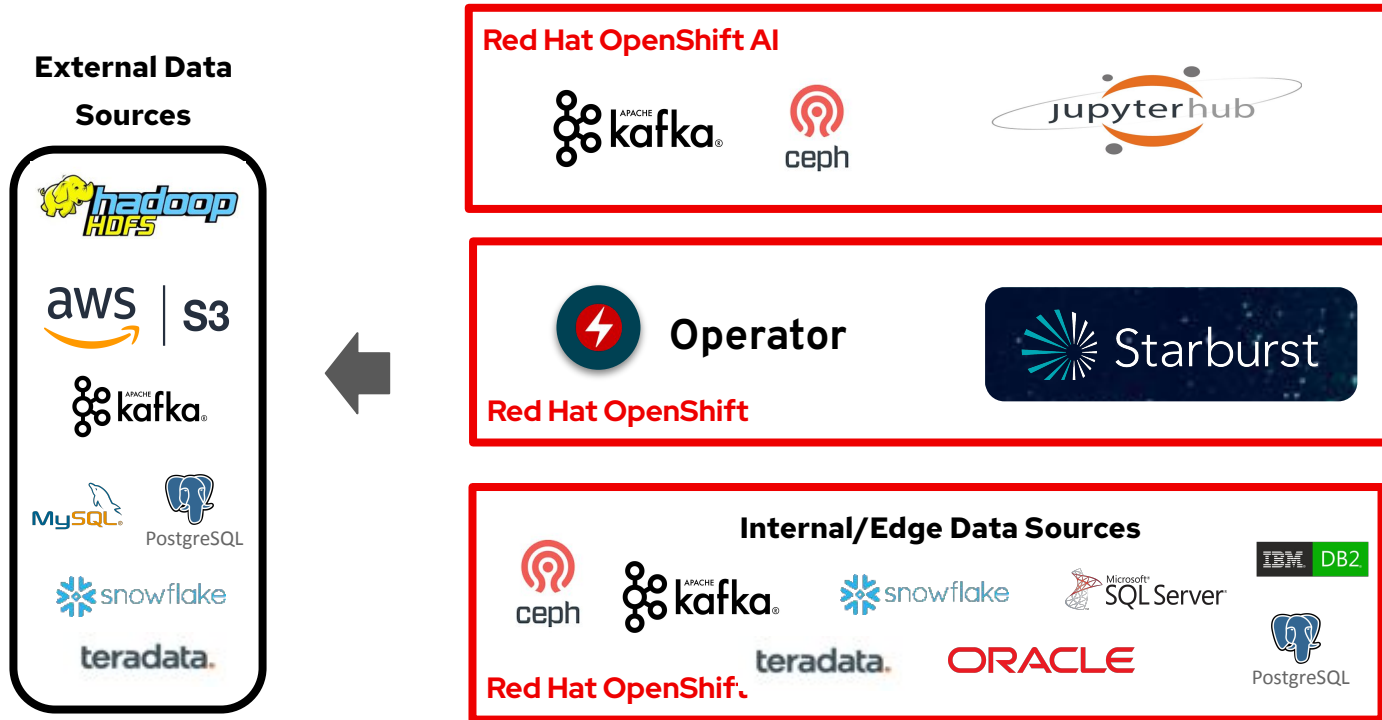
Save

## Strategic partnerships + Red Hat AI/ML offerings





# Data Acquisition and Preparation



# Operators on ROSA for AWS Services

operatorhub.io/?keyword=AWS

OperatorHub.io

Welcome to OperatorHub.io

OperatorHub.io is a new home for the Kubernetes community to share Operators. Find an existing Operator or list your own today.

49 ITEMS

CATEGORIES

- AI/Machine Learning
- Application Runtime
- Big Data
- Cloud Provider
- Database
- Developer Tools
- Drivers and plugins
- Integration & Delivery
- Logging & Tracing
- Modernization & Migration
- Monitoring
- Networking
- OpenShift Optional
- Security
- Storage
- Streaming & Messaging

PROVIDER

- ☐ Aerospike (1)
- ☐ aiven (1)
- ☐ alauda (1)
- ☐ Alibaba Cloud (1)
- ☐ Altinity (1)
- [Show 221 more](#)

CAPABILITY LEVEL

- ☐ Basic Install (182)

VIEW SORT A-Z

**Astra Trident**  
provided by NetApp, Inc.  
Trident Operator, to manage Astra Trident installations

**AWS Auth Operator**  
provided by Gennady Potapov  
Automates AWS auth-cm ConfigMap management

**aws**  
AWS Controllers for Kubernetes - Amazon ACM  
provided by Amazon, Inc.  
AWS ACM controller is a service controller for managing AWS ACM resources

**aws**  
AWS Controllers for Kubernetes - Amazon ACM PCA  
provided by Amazon, Inc.  
AWS ACM PCA controller is a service controller for managing AWS ACM PCA resources

**aws**  
AWS Controllers for Kubernetes - Amazon API Gateway v2  
provided by Amazon, Inc.  
AWS API Gateway v2 is a service controller for managing AWS API Gateway v2 resources

**aws**  
AWS Controllers for Kubernetes - Amazon Application Auto Scaling  
provided by Amazon, Inc.  
AWS Application Auto Scaling is a service controller for managing AWS Application Auto Scaling resources

**aws**  
AWS Controllers for Kubernetes - Amazon CloudFront  
provided by Amazon, Inc.  
AWS CloudFront controller is a service controller for managing AWS CloudFront resources

**aws**  
AWS Controllers for Kubernetes - Amazon CloudTrail  
provided by Amazon, Inc.  
AWS CloudTrail controller is a service controller for managing AWS CloudTrail resources

**aws**  
AWS Controllers for Kubernetes - Amazon CloudWatch  
provided by Amazon, Inc.  
AWS CloudWatch controller is a service controller for managing AWS CloudWatch resources

**aws**  
AWS Controllers for Kubernetes - Amazon CloudWatch Logs  
provided by Amazon, Inc.  
AWS CloudWatch Logs is a service controller for managing AWS CloudWatch Logs resources

**aws**  
AWS Controllers for Kubernetes - Amazon DocumentDB  
provided by Amazon, Inc.  
AWS DocumentDB controller is a service controller for managing AWS DocumentDB resources

**aws**  
AWS Controllers for Kubernetes - Amazon DynamoDB  
provided by Amazon, Inc.  
AWS DynamoDB controller is a service controller for managing AWS DynamoDB resources

**aws**  
AWS Controllers for Kubernetes - Amazon EC2  
provided by Amazon, Inc.  
AWS EC2 controller is a service controller for managing AWS EC2 resources

**aws**  
AWS Controllers for Kubernetes - Amazon ECR  
provided by Amazon, Inc.  
AWS ECR controller is a service controller for managing AWS ECR resources

**aws**  
AWS Controllers for Kubernetes - Amazon ECS  
provided by Amazon, Inc.  
AWS ECS controller is a service controller for managing AWS ECS resources

# Strategic partnerships within AI/ML ecosystem

## AI/ML life cycle



## Data governance and security



## Data processing



## Data analytics



## Databases



## AI Ops



## Infrastructure partners



## Hardware acceleration





## Training and Certification

Close skills gaps and hone your teams' Red Hat product expertise; flexible learning options to meet your business needs

*\*coming soon\**

<b>Introduction to Red Hat OpenShift AI (AI262)</b>	<b>Red Hat OpenShift AI Administration (AI263)</b>	<b>Creating Machine Learning Models with Red Hat OpenShift AI (AI264)</b>	<b>Deploying Machine Learning Models with Red Hat OpenShift AI (AI265)</b>	<b>Automation using Data Science Pipelines (AI266)</b>
Learn about the general architecture and main features of RHOAI and create some basic projects.	Get hands-on administration experience with RHOAI including install, upgrades, management and configuration.	Learn basic Machine Learning concepts and create and train Machine Learning models.	Deploy and serve Machine Learning models on RHOAI and troubleshoot deployed models.	Learn about and get hands-on experience with KubeFlow and Elyra pipelines.
<b>Developing and Deploying AI/ML Applications on Red Hat OpenShift AI (AI267):</b> 4 day instructor-led virtual training				
<b>Red Hat Certified Specialist: OpenShift AI (EX267):</b> 3 hour exam				

# Conclusion

- ▶ How the AI/ML landscape is evolving: market opportunities
- ▶ AI Application Examples vs intelligent Application?
- ▶ Challenges of Operationalizing AI ?
- ▶ Team topologies and operationalizing models
- ▶ Red Hat OpenShift AI - key features and walkthrough
- ▶ Why application platforms? Gartner two speed architecture
- ▶ Where to start?



# Thank you!

Yury Titov

**ytitov@redhat.com**



[linkedin.com/company/red-hat](https://linkedin.com/company/red-hat)



[youtube.com/user/RedHatVideos](https://youtube.com/user/RedHatVideos)



[facebook.com/redhatinc](https://facebook.com/redhatinc)



[twitter.com/RedHat](https://twitter.com/RedHat)