

Accelerate your Al Journey with Red Hat Al and NVIDIA with Hitachi Unified Compute Platform and Virtual Storage Platform One

© 2025 Hitachi Vantara LLC. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including copying and recording, or stored in a database or retrieval system for commercial purposes without the express written permission of Hitachi, Ltd., Hitachi Vantara, Ltd., or Hitachi Vantara Corporation (collectively "Hitachi"). Licensee may make copies of the Materials provided that any such copy is: (i) created as an essential step in utilization of the Software as licensed and is used in no other manner; or (ii) used for archival purposes. Licensee may not make any other copies of the Materials. "Materials" mean text, data, photographs, graphics, audio, video and documents.

Hitachi reserves the right to make changes to this Material at any time without notice and assumes no responsibility for its use. The Materials contain the most current information available at the time of publication.

Some of the features described in the Materials might not be currently available. Refer to the most recent product announcement for information about feature and product availability, or contact Hitachi Vantara LLC at <u>https://support.hitachivantara.com/en_us/contact-us.html</u>.

Notice: Hitachi products and services can be ordered only under the terms and conditions of the applicable Hitachi agreements. The use of Hitachi products is governed by the terms of your agreements with Hitachi Vantara LLC.

By using this software, you agree that you are responsible for:

- **1.** Acquiring the relevant consents as may be required under local privacy laws or otherwise from authorized employees and other individuals; and
- **2.** Verifying that your data continues to be held, retrieved, deleted, or otherwise processed in accordance with relevant laws.

Notice on Export Controls. The technical data and technology inherent in this Document may be subject to U.S. export control laws, including the U.S. Export Administration Act and its associated regulations, and may be subject to export or import regulations in other countries. Reader agrees to comply strictly with all such regulations and acknowledges that Reader has the responsibility to obtain licenses to export, re-export, or import the Document and any Compliant Products.

Hitachi and Lumada are trademarks or registered trademarks of Hitachi, Ltd., in the United States and other countries.

AIX, DB2, DS6000, DS8000, Enterprise Storage Server, eServer, FICON, FlashCopy, GDPS, HyperSwap, IBM, OS/390, PowerHA, PowerPC, S/390, System z9, System z10, Tivoli, z/OS, z9, z10, z13, z14, z15, z16, z/VM, and z/VSE are registered trademarks or trademarks of International Business Machines Corporation.

Active Directory, ActiveX, Bing, Excel, Hyper-V, Internet Explorer, the Internet Explorer logo, Microsoft, Microsoft Edge, the Microsoft corporate logo, the Microsoft Edge logo, MS-DOS, Outlook, PowerPoint, SharePoint, Silverlight, SmartScreen, SQL Server, Visual Basic, Visual C++, Visual Studio, Windows, the Windows logo, Windows Azure, Windows PowerShell, Windows Server, the Windows start button, and Windows Vista are registered trademarks or trademarks of Microsoft Corporation. Microsoft product screen shots are reprinted with permission from Microsoft Corporation.

All other trademarks, service marks, and company names in this document or website are properties of their respective owners.

Copyright and license information for third-party and open source software used in Hitachi Vantara products can be found in the product documentation, at <u>https://www.hitachivantara.com/en-us/company/legal.html</u>.

Feedback

Please send comments to doc.feedback@hitachivantara.com. Include the document title and number, including the revision level (for example, -07), and refer to specific sections and paragraphs whenever possible. All comments become the property of Hitachi Vantara LLC.

Thank you!

Revision history

Changes	Date
Initial release	March 2025

Reference Architecture Guide

Red Hat AI is a portfolio of products, including both Red Hat Enterprise Linux AI and Red Hat OpenShift AI, and services that accelerate the development and deployment of AI solutions across hybrid cloud environments.

This paper provides a reference configuration of Red Hat OpenShift with the OpenShift AI feature enabled, deployed on Hitachi's integrated systems using Hitachi Unified Compute Platform leveraging VSP One Storage Platforms and NVIDIA certified data center servers. It leverages the latest capabilities and services provided by OpenShift AI and NVIDIA NIM microservices.

Red Hat OpenShift AI includes the components of Red Hat Enterprise Linux AI and extends and scales the capabilities with Kubernetes-based containerized orchestration and the ability to distribute AI workloads across multiple nodes. It is a platform that empowers enterprises to develop and deploy AI-driven solutions efficiently across hybrid cloud environments, facilitating the creation and delivery of AI-enabled applications at scale with high reliability and security. Additionally, Red Hat OpenShift AI supports model serving and hosting, making real-time AI inference and integration seamless.

A key element in the successful deployment of an AI platform is having a robust and flexible infrastructure (compute, storage, network) that can meet a wide variety of requirements in a highly dynamic environment. Hitachi infrastructure with Red Hat OpenShift AI and NVIDIA AI Enterprise provides a highly available and high-performance infrastructure for AI workloads.

Hitachi Unified Compute Platform includes a Hitachi Storage CSI driver supported with Hitachi Storage Plug-in for Containers (HSPC). Using these proven CSI (Container Storage Interface) storage integrations, you can provide persistent reliable storage for AI workloads.

This reference architecture also provides the reference design for a build-your-own AI platform using Hitachi Unified Compute Platform with Red Hat OpenShift and NVIDIA GPUs. Although a specific converged system is used as an example, this reference design can be applied to building your own AI/ML container platform.

The intended audience of this document is IT administrators, system and AI/ML architects, consultants, and sales engineers to assist in planning, designing, and implementing Hitachi Unified Compute Platform with OpenShift Container Platform solutions.

Overview

Red Hat OpenShift Container Platform is a successful container orchestration platform and is one of the container orchestration and virtualization solutions available in the Hitachi Unified Compute Platform solution portfolio. The following figure shows a high-level diagram of Red Hat OpenShift Container Platform with OpenShift AI alongside NVIDIA AI Enterprise. There is flexibility in the hypervisor that can be used. This reference guide leverages KVM that as part of the OpenShift Container Platform/OpenShift Virtualization deployment. It also integrates Hitachi Virtual Storage Platform for persistent storage, with Hitachi Storage Plug-in for Containers enabling automated storage provisioning to the Kubernetes cluster(s).



The integration of Hitachi Unified Compute Platform with Red Hat OpenShift and NVIDIA provides a unified and full-stack hybrid platform to deploy and manage all your containerized applications, virtual machines (VMs), and AI/ML workloads across a hybrid environment, reducing operational complexity and paving the way for modernization. A high-level architecture of this platform for all workloads can be seen in the following illustration.



Red Hat OpenShift AI is a powerful artificial intelligence (AI) platform for data scientists and developers of AI applications. Red Hat OpenShift AI provides the tools that let you rapidly develop, train, test, deploy and monitor AI/ML models on-premises and/or in the public cloud.

A key benefit we concluded with this reference pattern is that it helps to accelerate the time to value of AI projects. Before evaluating Red Hat AI, we observed that our customers had to spend resources to find the right combination of software packages suitable to run AI/ML workloads. This is proven to be time consuming, and many times required significant expertise to have an environment up and running to start prototyping/deploying AI/ML applications.

The very interactive and collaborative UI enables data scientists and data engineers to eliminate spending too much time on the infrastructure/DevOps processes and instead use that time to focus on developing/training, fine-tuning their models, and deploying/ integrating those models into workloads/applications in production environments.

We see significant benefits in the distributed workloads capability for more efficient model tuning. Other areas such as model monitoring and data pipelines were not extensively evaluated in this phase of paper.

This reference architecture was validated by the following use cases:

- Generative AI deployment with NVIDIA NIM microservices on OpenShift AI (on page 19)
- Model training from a pretrained model using OpenShift AI (on page 29)
- Testing a retrained model (on page 36)

Follow the steps in the Solution design, implementation, and use case validation sections to learn how to fast track and enable these AI accelerated capabilities for your business when using Hitachi Unified Compute Platform with Red Hat OpenShift Container Platform and OpenShift AI.

If deploying OpenShift AI with OpenShift Virtualization or OpenShift Container Platform, consider this reference guide a supplement to the existing reference architecture <u>Red Hat</u> <u>OpenShift Virtualization with Hitachi Unified Compute Platform and VSP One Storage</u>.

Solution components

This section outlines the components used in this reference architecture.

Hitachi Unified Compute Platform

Hitachi Unified Compute Platform is an optimized, preconfigured, and pretested integrated system or converged infrastructure platform to host workloads and platform stacks such as VMware vSphere/ VMware Cloud Foundation, Red Hat OpenShift, Kubernetes, Azure Local, Google Anthos, and bare metal. It offers a broad range of compute, network and storage components that can be scaled and configured independently to eliminate overprovisioning. You have a choice of operating environments to maximize your flexibility.

With Hitachi Unified Compute Platform, you can choose between single-rack configurations and multi-rack configurations with flexible local and VSP One Storage options across block, file, and object.

Hitachi Virtual Storage Platform One Block

VSP One Block arrives pre-configured, including Dynamic Drive Protection (DDP) groups. DDP replaces traditional RAID groups in VSP One Block storage systems, providing the resilience of RAID 6 with distributed spare space and support for an arbitrary number of drives (from 9-32 per group). Adding drives one (or more) at a time is also supported. DDP improves the resilience of the appliance by dramatically lowering rebuilding times.

Hitachi Virtual Storage Platform One File

<u>Hitachi Virtual Storage Platform One File</u> (VSP One File) is a modernized Hitachi NAS that is focused on improving customer experience, delivering on operational simplification, and enabling better agility.

The VSP One File series provides an easy and modern way to manage, monitor, and configure file systems and data across the organization.

Hitachi Virtual Storage Platform One File includes the following 3 dedicated models:

- VSP One File 32 intended for cost sensitive deployments, and/or hybrid flash deployment that have no short/medium term plan to move away from their 10GbE network infrastructure.
- VSP One File 34 targeted for technology buyers that desire the best price/performance, an all-flash deployment, and have standardized on 25GbE network connectivity or are planning to move short/medium term to 25GbE.
- VSP One File 38 a performance model and built for leading technology pioneers that require the most file performance, and highest rack network port density to address allflash use cases and applications.

For high SLA workloads, all VSP One File models support SMB3 multichannel, NFS 4.1 multi-pathing, Nconnect, and link aggregation control protocol (LACP) to improve availability and/or higher bandwidth for Windows and Linux clients beyond the limit of a single connection.

Hitachi Virtual Storage Platform One Object

VSP One Object offers massive scale, cloud capabilities, and broad protocol support. It supports all-flash configurations and policy-driven tiered storage, allowing performance and capacity to scale independently. It quickly ingests large volumes of small objects and cost-effectively stores very large objects.

Customers benefit from a diverse array of configurations and deployment options. These include traditional archive and backup solutions, as well as modern cloud-native environments, AI/ML applications, and analytics workloads. The platform is optimized for a compact form factor, fitting an entire cluster into a 4U footprint. This enables real-time analytics at the edge, eliminating the need to transfer data to a central location for analysis.

Additional hardware components

The tested solution used specific features based on the following hardware. You can use any qualified server platform defined in the Product Compatibility Guide. For more information, see <u>https://compatibility.hitachivantara.com/</u>.

Hardware	Description	Version	Quantity
Hitachi Advanced Server HA800 series (for OCP baremetal cluster)	 2 × Intel(R) Xeon(R) Gold 6454S processors 16 × 32 GB DIMM, 512 GB memory 	iLO 6 BIOS:U54 v1.30	6
	 NS204-u RAID1 for boot 		
	• SN1610E 32 Gb 2p FC HBA		
	• 2 × Intel(R) Eth E810-XXVDA2 NICs		
	 1 × NVIDIA L40s GPU (only on worker nodes) 		

Hardware	Description	Version	Quantity
Hitachi VSP One	• 1 TB cache	A3-04-01-40/00	1
Block (for OCP cluster)	 24 × 3.8 TB NVMe drives 		
	 4 × 32 Gbps Fibre Channel ports 		
Cisco Nexus	• 48-port 10/25 GbE	NXOS 10.3(4a)	2
93180YC-FX switch (leaf)	 6-port 40/100 GbE 		
Cisco Nexus	• 48-port 1 GbE	NXOS 9.3.(8)	1
92348	 4-port 1/10/25 GbE 		
	 2-port 40/100 GbE 		
Brocade G720	 48-port 16/32Gbps Fibre Channel switch 	9.1.1c	2

Software components

Software	Version
Hitachi Storage Microcode/Firmware Version	A3-04-01-40/00
Hitachi Storage Plug-in for Containers	3.15.0
Red Hat OpenShift Container Platform (OCP)	4.16.28
Red Hat OpenShift Al	2.16.1
NVIDIA GPU Operator	24.9.2

Red Hat OpenShift Container Platform

Red Hat OpenShift Container Platform (OCP) provides a single platform to build, deploy, and manage applications consistently across on-premises and hybrid cloud deployments. OCP provides the control plane and data plane within the same interface. OCP provides administrator views to deploy operators, monitor container resources, manage container health, manage users, work with operators, manage pods and deployment configurations, as well as define storage resources.

OCP also provides a developer view to allow users to deploy application resources from various pre-defined resources such as YAML files, Docker files, Catalogs, or GIT within user-defined namespaces. With OCP <code>kubectl</code>, a native binary of Kubernetes is complemented by the <code>oc</code> command, which provides further support for OCP resources, such as deployment and build configurations, routes, image streams, and tags. OCP provides a GUI and a CLI interface.

Red Hat OpenShift Al

Red Hat OpenShift AI provides an AI platform for building, training, tuning, deploying and monitoring AI-enabled applications and predictive and foundation models at scale across hybrid cloud environments. Red Hat OpenShift AI helps accelerate AI innovation, drive operational consistency, and optimize access to resources when implementing trusted AI solutions.

See the following for more details of OpenShift AI and RHEL AI:

https://www.redhat.com/en/resources/openshift-ai-overview

https://www.redhat.com/en/resources/foundation-model-platform-gen-ai-datasheet

NVIDIA GPU Operator

The NVIDIA GPU Operator uses the operator framework within Kubernetes to automate the management of all NVIDIA software components needed to provision GPUs. These components include the NVIDIA drivers (to enable CUDA), Kubernetes device plugins for GPUs, the NVIDIA Container Runtime, automatic node labeling, DCGM-based monitoring, and more.

NVIDIA AI Enterprise

NVIDIA AI Enterprise is a comprehensive, cloud-native software platform designed to accelerate data science pipelines and streamline the development and deployment of production-grade co-pilots and generative AI applications. It offers easy-to-use microservices that deliver optimized model performance, coupled with enterprise-grade security, support, and stability. This ensures a seamless transition from prototype to production for businesses that rely on AI to operate.

NVIDIA NIM for Large Language Models (LLMs)

NVIDIA NIM, part of NVIDIA AI Enterprise, is a set of easy-to-use microservices for accelerating the deployment of generative AI models across the cloud, data center, and workstations. It supports everything from open-source community models to NVIDIA AI Foundation models, as well as custom AI models.

NVIDIA NIMs make it easy for IT and DevOps teams to self-host large language models (LMMs) in their own managed environments allowing developers to quickly build enterprisegrade AI applications with minimum effort, reducing development time and cost. NIM offers the fastest path to inference with unparalleled performance.

Hitachi Storage Plug-in for Containers

Hitachi Storage Plug-in for Containers (HSPC) is a software component that contains libraries, settings, and commands that you can use to create a container to run stateful applications. It enables stateful applications to persist and maintain data after the lifecycle of the container has ended. Storage Plug-in for Containers provides persistent volumes from Hitachi Dynamic Provisioning (HDP) or Hitachi Thin Image (HTI) pools to bare metal or hybrid deployments using Fibre Channel, NVMe over Fibre Channel, or iSCSI protocols. iSCSI protocol is supported for both bare metal and virtual environments.

Storage Plug-in for Containers integrates Kubernetes or OpenShift with Hitachi storage systems using Container Storage Interface (CSI).

Kubernetes / OpenShift Cluster **Control Plane** Node Node Node csi-controller csi-node csi-node csi-node FC / **REST API** NVMe over FC / iSCSI protocols Storage System Volume Volume

The following diagram illustrates a container environment where Storage Plug-in for Containers is deployed.

Solution design

This section outlines a detailed solution example for Hitachi Unified Compute Platform with Red Hat OpenShift and NVIDIA GPUs.

Hitachi Unified Compute Platform infrastructure components

The following figure shows a high-level architecture of Hitachi Unified Compute Platform clusters used to validate the Red Hat OpenShift AI. It includes the following components:

- Hitachi Unified Compute Platform with Red Hat OpenShift Container Platform and OpenShift AI:
 - OpenShift cluster bare metal configured with six Hitachi Advanced Server nodes: three nodes for the control plane and three for worker nodes to run AI/ML workloads.
 - A Hitachi Virtual Storage Platform One storage system for persistent storage for standard Pods and VM Pods (current generation VSP 5000 series or VSP E series storage systems can be used as well).
- The following network infrastructure is used for these clusters:
 - Two Cisco 9332C spine Ethernet switches
 - Two Cisco 93180YC-FX leaf Ethernet switches
 - One Cisco 92348 management switch



The clusters used for this reference architecture were configured following the standard guidelines from Hitachi Unified Compute Platform.

As stated earlier, the solution in this reference architecture is based on Hitachi Unified Compute Platform with OpenShift and VSP/VSP One storage which supports multiple protocols with options for Fibre Channel (FC), FC-NVMe, and iSCSI. In addition, this solution can be integrated/deployed together with Hitachi VSP One File and VSP One Object to offer all the storage capabilities to meet the needs for the AI/ML workloads.

This reference design is targeted at medium and enterprise customers requiring an accelerated computing solution reducing the operational burden to develop and deploy AI model enhanced application/workload services. For even larger enterprises with extreme scalability requirements, the Hitachi iQ solution with NVIDIA HGX offers pre-built platforms with bespoke customization and delivery, combining outstanding software, infrastructure, distributed parallel file systems, and expertise in unified and scalable AI development solutions.

For building a large enterprise datacenter using the latest high-speed, low latency solution, follow the Hitachi iQ guides and reference architecture:

- Hitachi's AI Solutions Stack for Enterprise Workloads: <u>https://www.hitachivantara.com/en-us/pdf/datasheet/hitachi-iq-with-nvidia-hgx.pdf</u>
- Hitachi iQ with NVIDIA HGX Reference Architecture: <u>https://www.hitachivantara.com/en-us/gated-forms/iq-with-hgx-h100-powered-with-hcsf-hcp</u>
- More details about Hitachi's integrated systems and converged infrastructure can be found at <u>https://www.hitachivantara.com/en-us/products/integrated-systems/</u> <u>convergedinfrastructure</u>

Deploy OpenShift Container Platform

This guide does not cover the step-by-step details of how to implement OCP. Follow Red Hat OCP documentation for the setup of the cluster.

One OCP bare metal cluster with Hitachi VSP One storage and Hitachi Storage Plug-in for Containers has been configured to support the different use cases described in this guide. In addition, the worker nodes have been equipped with NVIDIA L40s GPUs.

The following are additional details about the cluster used for this validation.

<pre>[root@jputilitysrv1 ocpjpal]# oc vers: Client Version: 4.16.28 Kustomize Version: v5.0.4-0.202306011 Server Version: 4.16.28</pre>	ion 659 4 7-6ce	0bf390ce3		
Rubernetes Version: V1.29.10+6/d338/				
[root@jputilitysrv1 ocpjpai]#				
[root@jputilitysrv1 ocpjpai]# oc get n	nodes			
JAME	STATUS	ROLES	AGE	VERSION
cpjpai-master-1.ocpjpai.ocp.sce.lab	Ready	control-plane, master	22d	v1.29.10+67d3387
cpjpai-master-2.ocpjpai.ocp.sce.lab	Ready	control-plane, master	22d	v1.29.10+67d3387
cpjpai-master-3.ocpjpai.ocp.sce.lab	Ready	control-plane, master	22d	v1.29.10+67d3387
ocpjpai-worker-1.ocpjpai.ocp.sce.lab	Ready	worker	22d	v1.29.10+67d3387
cpjpai-worker-2.ocpjpai.ocp.sce.lab	Ready	worker	22d	v1.29.10+67d3387
cpjpai-worker-3.ocpjpai.ocp.sce.lab	Ready	worker	22d	v1.29.10+67d3387
[root@jputilitysrv1 ocpjpai]#				

For detailed OpenShift Container Platform installation procedures, see the Red Hat documentation at <u>https://docs.redhat.com/en/documentation/openshift_container_platform/</u><u>4.16/</u>.

Deploy Hitachi Storage Plug-in for Containers

Hitachi Storage Plug-in for Containers (HSPC) is easily deployed to OpenShift using the Operator, which can be installed from OperatorHub.

This guide does not cover the step-by-step how to install HSPC, follow the <u>https://</u> <u>docs.hitachivantara.com/v/u/en-us/adapters-and-drivers/3.15.x/mk-92adptr142</u> or the steps described in the <u>https://docs.hitachivantara.com/v/u/en-us/application-optimized-</u> <u>solutions/mk-sl-210</u>.

In summary, these steps include:

- Install Hitachi Storage Plug-in for Containers.
- Configure Secret settings to access Hitachi VSP One storage systems.

- Configure StorageClass settings.
- Configure Multipathing (Fibre Channel or iSCSI).

Specific steps for how to configure Storage Classes will be covered as part of each of the use cases in the Solution implementation and validation section.



Note: If there is a previous version of Storage Plug-in for Containers, remove it before performing the installation procedure.

HSPC and VSP Host Groups

Host groups required for Storage Plug-in for Containers are automatically created by Storage Plug-in for Containers. Storage Plug-in for Containers automatically searches host groups and iSCSI targets based on the name.

To use existing host groups, rename them according to the naming rule. For details, see *Host group and iSCSI target naming rules* in the *Hitachi Storage Plug-in for Containers Quick Reference Guide* at https://docs.hitachivantara.com/v/u/en-us/adapters-and-drivers/3.14.x/mk-92adptr142.



Note: Storage Plug-in for Containers will overwrite host mode options even if existing host groups have other host mode options.

HSPC and VSP Resource Groups

You can partition storage system resources by limiting the LDEV ID range added to the resource group for a specific Kubernetes cluster. You can also isolate impacts between Kubernetes clusters. The following requirements should be met:

- Multiple Kubernetes clusters can share one resource group.
- Storage system users must have access only to the resource group that they created. The storage system user must not have access to other resource groups.
- Create a pool from pool volumes with the resource group that you have created.
- Allocate the necessary number of undefined LDEV IDs to the resource group.
- Allocate the necessary number of undefined host group IDs to the resource group for each storage system port defined in StorageClass. The number of host group IDs must be equal to the number of hosts for all ports.

For details, see the Hitachi Storage Plug-in for Containers documentation at <u>https://</u> <u>docs.hitachivantara.com/search/all?query=Hitachi%252BStorage%252BPlug-in%252Bfor</u> %252BContainers&content-lang=en-US

Red Hat OpenShift AI with Hitachi Unified Compute Platform and VSP One Storage

Before you begin

- You must have a subscription for Red Hat OpenShift Al.
- If you are planning to use NVIDIA NIMs with OpenShift AI, you must have NVIDIA AI Enterprise subscription or be an NVIDIA Developer Program member.
- Have RedHat OpenShift cluster 4.14 or later installed on a bare metal or virtual infrastructure:
 - You must have at least 2 worker nodes, and at least one of them with an NVIDIA GPU.
- You can use any of the four installation methods (user-provisioned, installer-provisioned, assisted, agent-based installer) of an OCP cluster on bare metal.
- Use Hitachi Virtual Storage Platform as the backend storage to the OCP cluster and OpenShift Virtualization.
- Hitachi Storage Plug-in for Containers (HSPC). A storageClass defined on the OCP cluster with HSPC as the provisioner, and makes sure the storage class is configured as default. Follow the OpenShift Container Platform documentation to configure a default storage class: Changing the default storage class.



Install and configure OpenShift AI and its components

After the pre-requisites have been met, you can follow these steps to prepare and install Red Hat OpenShift AI and its components.

Before you begin

After OpenShift cluster is ready, log in to the console and install the following operators from OperatorHub:

- Node Feature Discovery Operator
 - The Node Feature Discovery Operator manages the detection of hardware features and configuration in a Kubernetes cluster by labeling the nodes with hardware-specific information. The Node Feature Discovery (NFD) will label the host with node-specific attributes, like PCI cards, kernel, or OS version, and more.
- NVIDIA GPU Operator
 - The NVIDIA GPU Operator uses the operator framework within Kubernetes to automate the management of all NVIDIA software components needed to provision and monitor GPUs. These components include the NVIDIA drivers (to enable CUDA), Kubernetes device plugin for GPUs, the NVIDIA Container Runtime, automatic node labelling and NVIDIA DCGM exporter.
 - After the GPU Operator and cluster policy have been created, you can check the status with the following command:

```
oc get pods -n nvidia-gpu-operator
```

[root@jputilitysrv1 ocpjpai]# oc get pods -n nvidia-	gpu-opera	ator		
NAME	READY	STATUS	RESTARTS	AGE
gpu-feature-discovery-z7flc	2/2	Running	0	10d
gpu-operator-5bf9695f4c-dpnm6	1/1	Running	0	10d
nvidia-container-toolkit-daemonset-jrnrh	1/1	Running	0	10d
nvidia-cuda-validator-kk79x	0/1	Completed	0	10d
nvidia-dcgm-69vcp	1/1	Running	0	10d
nvidia-dcgm-exporter-nks4v	1/1	Running	0	10d
nvidia-device-plugin-daemonset-cj5lr	2/2	Running	0	10d
nvidia-driver-daemonset-416.94.202412170927-0-wmfdj	2/2	Running	0	10d
nvidia-node-status-exporter-5b92f	1/1	Running	0	10d
nvidia-operator-validator-drdfm	1/1	Running	0	10d
[root@jputilitysrv1 ocpjpai]#				

For details how to install the NFD and GPU Operator, see the instructions at <u>NVIDIA GPU</u> <u>Operator on Red Hat OpenShift Container Platform</u>.

In addition to the NFD and GPU Operators, to support the KServe component, which is used by the single-model serving platform to serve large language models, you must install the following Operators:

- Red Hat OpenShift Serverless
- Red Hat OpenShift Service Mesh

For more details, see About the single-model serving platform.

If there is a need to add an authorization provider for the single-model serving platform, install the Red Hat Authorino (Technical Preview). For more details, see <u>Adding an authorization</u> provider for the single-model serving platform.

After these operators are installed and in a ready state, now proceed to install Red Hat OpenShift AI operator using the following procedure. This operator can be installed either from the CLI or web console. In this guide we are following the installation using the OpenShift web console.

Procedure

- 1. In the Red Hat OpenShift web console, navigate to **Operators > OperatorHub** and filter by the keyword *Red Hat OpenShift AI*.
- **2.** Select the Red Hat OpenShift AI Operator tile to open the RedHat OpenShift AI information pane.



- 3. Click Install.
- 4. Keep the default parameters. To start the installation, click Install again.
- 5. After the installation completes successfully, select Create DataScienceCluster.
- 6. Keep the default values and click Create.
- 7. Wait until the All instances tab shows the status Ready for all components.

Result

After all the components/operators are installed, you can verify all the installed operators.

Red Hat OpenShift						
🛠 Administrator	•	Project: All	Projects 👻			
Home	,	Installe	d Operators			
Operators	~	Installed Op	erators are represented by ClusterService\	/ersions within this Namespace. For more information, s	ee the Understanding Operators documentation 🗗 O	r create an Operator and
OperatorHub		Name 👻	Search by name /			
Installed Operators		Name	1	Namespace 💲	Managed Namespaces 🛛 🖡	Status
Workloads	>	#	Red Hat OpenShift Service Mesh	NS openshift-operators	All Namespaces	Succeeded
Serverless	>		2.6.5-U provided by Hed Hat, Inc.			
Networking	>		Red Hat OpenShift Serverless 1.35.0 provided by Red Hat	NS openshift-serverless	All Namespaces	Succeeded Up to date
Storage	>	CyselDati.cz	Red Hat OpenShift AI 2.16.1 provided by Red Hat, Inc.	NS redhat-ods-operator	All Namespaces	Succeeded Up to date
Builds	>		Package Conver	C acarchift aparatar lifeguela managar	n anarchift anarstar liferurle manager	Succeeded
Observe	>	(7)	0.0.1-snapshot provided by Red Hat	opensin coperator-mecycle-manager	opensint-operator-inecycle-manager	Succeeded
Compute	>	S	Node Feature Discovery Operator 4.16.0-202412170135 provided by Red Hat	NB openshift-nfd	NS openshift-nfd	Succeeded Up to date
User Management	>	LITTLE IN	Hitachi Storage Plug-in for	NS hspc	NS hspc	Succeeded
Administration	>	kiper fil het.	Containers 1.14.3 provided by Hitachi			Up to date
		PVIDIA.	NVIDIA GPU Operator 24.9.2 provided by NVIDIA Corporation	NS nvidia-gpu-operator	NS nvidia-gpu-operator	Succeeded Up to date
		I	Red Hat - Authorino (Technical Preview) 111 provided by Red Hat	NS openshift-operators	All Namespaces	Succeeded Up to date

Access OpenShift AI dashboard

Before we access OpenShift AI dashboard, let us get familiar with some components of a data science project:

- Workbenches are instances of your development and experimentation environment. They typically contain IDEs, such as JupyterLab, RStudio, and Visual Studio Code.
- A Cluster storage is a volume that persists the files and data you're working on within a workbench. A workbench has access to one or more cluster storage instances.
- Data connections contain configuration parameters that are required to connect to a data source, such as an S3 object bucket.
- Pipelines contain the Data Science pipelines that are executed within the project.
- Models and model servers allow you to quickly serve a trained model for real-time inference. You can have multiple model servers per data science project. One model server can host multiple models.

After OpenShift AI is installed, you can access the URL for your OpenShift AI console and share the URL with users to let them log in and work on their models.

Procedure

1. To access the OpenShift AI web console, click the application launcher (IIII) from the OCP console, located on the right side of the top navigator bar. Click and select Red Hat OpenShift AI from the drop-down menu.

Red Hat OpenShift				
📽 Administrator	•	Overview		Red Hat Applications
Home	>	Cluster		ned Hat Hybrid Cloud Console 🛛 🗹
Operators	~			OpenShift Self Managed Services
OperatorHub		Details View settings	Status	Red Hat OpenShift Al
Installed Operators		Cluster API address https://api.ocpjpai.ocp.sce.lab:	Cluster	Control Plane
		6443	• Operators	A Insights

2. Log in with your OpenShift credentials, and you will be presented with the OpenShift AI landing page where you can enable applications and create Data Science Projects.



All the data science projects you create will benefit from HSPC, which provides robust and high-performance persistent volumes. Additionally, HSPC together with HRPC (Hitachi Replication Plug-in for Containers) offers advanced disaster recovery (DR) protection, ensuring that your Al projects and applications maintain high availability and resilience. This means you can focus on innovation and development, knowing that your data and applications are secure and reliable.

Solution implementation and validation – use cases

This reference architecture was validated by the following use cases:

- Generative AI deployment with NVIDIA NIM microservices on OpenShift AI (on page 19)
- Model training from a pretrained model using OpenShift AI (on page 29)
- Testing a retrained model (on page 36)

Generative AI deployment with NVIDIA NIM microservices on OpenShift AI

In this test case we are going to show how to deploy generative AI applications leveraging NVIDIA NIM microservices and OpenShift AI on top of Hitachi Unified Compute Platform and VSP One infrastructure.

See the following flow diagram that showcases the high-level steps to deploy a GenAI application using NVIDIA NIM microservices with OpenShift AI.



Enable the NVIDIA NIM application

Procedure

 Before enabling the NIM application on the OpenShift dashboard, make sure to enable it in the OdhDashboardConfigs instance under Home > API Explorer in the Administrator perspective of the OpenShift console. The requirement is to set disableNIMModelServing to false.

For more details, see Editing the dashboard configuration file.

Red Hat OpenShift					
📽 Administrator	• •	Project: rec	Ihat-ods-applications		
		OdhDashboardConfigs > OdhDashboardConfig details			
Home	₩~		dh-dashboard-config		
Overview					
Projects		Details	YAML		
Search					
	ר	80 SP	et:		
		87	dashboardConfig:		
Events		88	enablement: true		
		89	disableTrustyBiasMetrics: false		
		90	disableDistributedWorkloads: false		
Operators	>	91	disableProjects: false		
		92	disableSupport: false		
		93	disableHome: false		
Workloads	>	94	disablePipelines: false		
		95	disableProjectSharing: false		
_ ·		96	disableModelServing: false		
Serverless	>	97	disableKServe: false		
		98	disableAcceleratorProfiles: false		
		99	disableNIMModelServing: false		
Networking	``	100	disableCustomServingRuntimes: false		
		101	disableModelMesh: false		

2. Gen an API key from NVIDIA. To do this, log in to the NVIDIA NGC catalog and generate an API key using the option **Generate API Key** under **profile menu > Setup**.

≡	IVIDIA . Organization	v		0	jose.perez2@hitachivanta Hitachi Vantara 11.C	ra.com 🔨
	Dashboard	Setup			● Hitachi Vantara LLC >	
8	Audit	Setup			 Contact Admin Account Settings 	*
μ.	Organization Profile	Keys/Secrets			🖉 Setup	
×	Teams				Organization	
áı	Usage	Generate API Key	Generate Personal Key		Terms of Use	
2	Users	Generate API Key Generate your own API key in order to use the NGC corrise	Personal Keys Generate your own API key in ordor to use the NGC convice		 → Sign Out 	
	Activate Subscription	through the Docker client or through NGC CLI.	through the Docker client or through NGC CLI for Secrets			
Ŷ	Subscriptions		Manager, NGC Catalog, Public API Endpoints, and			
¢	Service Keys		Private Registry.			
	External IdP					

3. From the Red Hat OpenShift AI dashboard, select the **Explore** option under Applications, and select the **NVIDIA NIM** tile.

≡ Red Hat CreenShift Al				III 🌲 🕢 cluster_admin 👻
Home You	Explore Add optional applications to your Red Hat Opens	NVIDIA NIM X by NVIDIA		
Enabled Explore	(Partner managed) O	😪 elastic 🛛 (Self-managed) O	IBM (Self-managed) •	
Data Science Projects	Anaconda Professional by Anaconda	Elasticsearch by Elastic	IBM watsonx.ai by IBM	NVIDIA NIM, part of NVIDIA AI Enterprise, is a set of easy-to-use microservices designed for secure, reliable
Data Science Pipelines	Anaconda Professional is a popular	Build transformative RAG applications,	IBM® watsonx.ai is part of the IBM	across the cloud, data center and workstations.
Experiments >	management experience that is	and address complex security threats –	together new generative AI capabilities	Supporting a wide range of Al models, including open- source community and NVIDIA AI Foundation models,
Distributed Workload Metrics	optimized for commercial use.	all with the power of Search Al.	traditional machine learning (ML) into a	remsures seamless, scalable Al interencing, on- premises or in the cloud, leveraging industry standard
Model Serving			lifecycle.	AP15.
Resources	(Self-managed)	(Red Hat managed) O	(Self-managed) .	
Settings >			IVERA	
	Intel® oneAPT AI Analytics Toolkit Containes by Intel® The AI Kit is a set of AI software tools to accelerate and-to-end data acience and analytics pipelines on Intel® architectures.	Jupyter of by Jupyter A multi-user version of the notebook designed for companies, classrooms and research labs.	NVIDIA NIM by NVIDIA NIM is a set of easy-to-use microservices designed for secure, reliable deployment of high- performance AI model inferencing.	

4. Click **Enable**, and enter the API key obtained from NVIDIA NGC catalog, and then click **Submit** to enable the NVIDIA NIM tile.

Enter N	Enter NVIDIA AI Enterprise license key		
NVIDIA AI E	nterprise license key		
		0	
This key is gi	ven to you by NVIDIA	_	
Submit	Cancel		

Result

When enabled, you can verify it by selecting the Enabled option under Applications.



Create a Data Science Project

To implement a data science workflow, you need to create a data science project. A data science project allows you and your team to organize and collaborate on resources within separated namespaces. From a project you can create multiple workbenches, each with their own Jupyter notebook environment, and each with their own data connections and cluster storage. In addition, the workbenches can share models and data with pipelines and model servers, all in a single project.

Procedure

- 1. From the OpenShift AI dashboard, navigate to the Data Science Projects tab.
- 2. Click Create Project, and enter a name and description for the project.
- 3. Click Create.

Deploy an NVIDIA NIM serving model

Inside the created project is where all the additional resources for this test case will be deployed.

Procedure

- 1. Select the project created in the previous step.
- Specify NIM as the serving model. This can be done from either the Overview tab or from the Models tab. Click Deploy model under the NVIDIA NIM model serving platform tile.
- 3. On the new window, enter the details for the NVIDIA NIM model.
 - Select the NIM model, in this case we will use the Llama-3.1-8b-instruct-latest.
 - Enter the size of the persistent volume (PVC).
 - Select the number of NVIDIA GPU accelerator cards, based on the capacity of your cluster.
 - Select the model serving size and number of replicas.

≡ CpenShift Al	
Home	Deploy model with NVIDIA NIM × Configure properties for deploying your model using an NVIDIA NIM. ×
Applications >	Project nvidia-nim-openshift-ai-demo
Data Science Projects	Model deployment name *
Data Science Pipelines	nvidia-nim-model-llama3-demo
Experiments >	NVIDIA NIM * Llama-3.1-8b-instruct - latest
Distributed Workload Metrics	NVIDIA NIM storage size
Model Serving	GiB GiB Specify the size of the cluster storage instance that will be created to store the downloaded NVIDIA NIM. Make sure your storage size is greater than the NIM size specified by NVIDIA
Resources	Number of model server replicas to deploy * ⑦
Settings >	- 1 +
	Model server size * 🕐
	Large Limits: 10 CPU, 20GiB Memory Requests: 6 CPU, 16GiB Memory
	Accelerator ③
	NVIDIA GPU 🗸
	Number of accelerators
	- 1 +
	Deploy Cancel

4. After the model is deployed, you can check the status.

Data Science Projects > nvidia-nim-openshift-ai-	-demo	
nvidia-nim-openshift-ai-	demo ®	Actions 👻
Overview Workbenches Pipelines	Models Cluster storage Connections Perm	issions Settings
B Models and model servers (2)	De	Single-model serving enabled
Model name t Servin	ng runtime Inference endpoint	API protocol Status
 nvidia-nim-model- NVIDI Ilama3-demo 	IA NIM Internal endpoint details	REST I
3		
Framework	llama-3.1-8b-instruct	
Model server replicas	1	
Model server size	Large	
	6 CPUs, 16GiB Memory requested	
	10 CPUs, 20GiB Memory limit	
Accelerator	NVIDIA GPU	
Number of accelerators	1	

5. Take note of the internal endpoint URL, this will be used in next steps.

Data Science Projects > nvidia-him	-openshift-ai-demo						
🗂 nvidia-nim-open	shift-ai-demo ୧)				Action	s 🔻
Overview Workbenches	Pipelines Models	Cluster storage	Connections	Permissions	Settings		
B Models and model s	Inference endpoints				× Single	-model serving e	nabled
Models and models	Inference endpoints url https://nvidia-nim-mode	el-llama3-demo.nvidia-	nim-openshift-ai-de	emo.svc.cluster.loca	Single	-model serving e Status	nabled

Configure and create a Workbench

After the NVIDIA NIM serving model is deployed and we have the endpoint URL, the next step is to create a workbench.

A workbench is an instance of your development and experimentation environment. They typically contain IDEs, such as JupyterLab, RStudio, and Visual Studio Code. When deploying a workbench you can select a notebook image, such as PyTorch, Standard Data Science, TensorFlow, etc., depending on the data science work that is required.

Procedure

- 1. From the **Workbenches** tab, click **Create workbench**, and then enter a name and description.
- **2.** Select one of the existing notebook images. In this case we select "Standard Data Science" and version 2024.2.

			 ٠	0
Home	Data Science Projects > nvidia-nim-ope	nshift-ai-demo 🔸 Create workbench		
Applications >	Create workbench Configure properties for your workbe	nch.		
Data Science Projects	N			
Data Science Pipelines	Jump to section 😼	Name * nvidia-nim-rag-workbench		
Experiments >	Name and description	The resource name will be nvidia-nim-rag-workbench . Edit resource name ⑦		
Distributed Workload Metrics	Notebook image	Description		
Model Serving	Deployment size		 	
Resources	Environment variables			
Settings >	Cluster storage	Notebook image Image selection		
	Data connections	Standard Data Science	 	•
	Connections	Minimal Python		Î
		CUDA	~	I
		PyTorch		ı
		TensorFlow		ı
		TrustyAl		
		code-server		
	Create workbench Cancel	ROCm		*

- **3.** Select the container image and accelerator. In this case we are using a Small container size and NVIDIA GPU.
- **4.** Enter the size of the Persistent Volume. The PVC will be created on the storage being used by the default Storage Class. For Persistent storage we are using Hitachi VSP storage and Hitachi Storage Plugin for Containers (HSPC).
- 5. Click Create workbench and wait for the workbench to be in a running state.

Data Science Projects > nvidia-nim-openshift-ai-d	emo					
nvidia-nim-openshift-ai-d	lemo ®					Actions 💌
Overview Workbenches Pipelines	Models Cluster storage	Connections	Permissions	Settings		
🖉 Workbenches ® 🖑					Crea	te workbench
Name 🌐	Notebook image	Container size	Status 🔱			
♥ nvidia-nim-rag-workbench ⑦	Standard Data Science Python v3.11	Small	Running) Stop	Open 🛃	:
Cluster storage	Packages	Limits				
 nvidia-nim-rag-workbench Max 1006/B Mount path: /opt/app-root/src 	Boto3 v1.35 Kafka-Python-ng v2.2 Kfp v2.9 Matplotlib v3.9 Numpy v2.1 Pandas v2.2 Scikit-learn v1.5 Scipy v1.14 Odh-Elyra v4.2 PyMongo v4.8 Pyodbc v5.1 Codeflare-SDK v0.23 Sklearn-onnx v1.17 Psycopg v3.2 MySQL Connector/Python v9	2 CPU, 8GiB Memory Requests 1 CPU, 8GiB Memory	y ,			

Create a notebook and verify access to the NVIDIA NIM model

The next step after creating the workbench is to deploy a notebook and work with an AI model.

Procedure

- **1.** To access the workbench interface, click the **Open** link. Log in with the OpenShift credentials.
- 2. The first window that opens after opening the workbench is a Jupiter environment.

~	→ C Solution Not secure https://nv	dia-nim-rag-workbench-nvidia-nim-openshift-ai-demo.apps.ocpjpai.ocp.sce.lab/notebook/nvidia-nim-openshift-ai 😭
()	File Edit View Run Kernel Git	Tabs Settings Help
	You are not currently in a Git repository. To use Git, navigate to a	Launcher +
0	local repository, initialize a repository here, or clone an existing repository.	Notebook
	Open the FileBrowser	
≡	Clone a Repository	Python 3.11
£		Console
₿		
88		Clone a repo
*		Enter the URI of the remote Git repository https://github.com/josep
		Download the repository Cancel Clone Pipelin Latest
		S_ Other
		Terminal Text File

3. In this window you can create a notebook from scratch, or you use git to clone an existing repo. In this case, we are cloning an existing repo which contains a notebook that is ready to use.

This is an example of how to invoke NVIDIA NIM running internally in a serving model inside the same project as the workbench.

Make sure to use the URL as shown on the NVIDIA serving model end-point. Change the value of the nim url parameter accordingly.

```
import json
import requests
import urllib3
urllib3.disable_warnings()
nim_url = "https://nvidia-nim-model-llama3-demo.nvidia-nim-openshift-ai-
demo.svc.cluster.local"
```

This cell is to query the list of available models. Keep note of this model "meta/ llama-3.1-8b-instruct" because it will be used in the next step.

```
response = requests.get(nim_url + "/v1/models", verify=False)
print(json.dumps(response.json()['data'], indent=2))
```

Example of the response:

[

```
"id": "meta/llama-3.1-8b-instruct",
    "object": "model",
    "created": 1739164868,
    "owned by": "system",
    "root": "meta/llama-3.1-8b-instruct",
    "parent": null,
    "max model len": 131072,
    "permission": [
      {
        "id": "modelperm-27b2704080674e848eece5efa2332ff2",
       "object": "model permission",
       "created": 1739164868,
       "allow_create_engine": false,
       "allow sampling": true,
       "allow logprobs": true,
       "allow search indices": false,
        "allow view": true,
       "allow fine tuning": false,
        "organization": "*",
       "group": null,
        "is blocking": false
     }
   ]
  }
]
```

The following cell contains the code required to send a chat request to the target NIM model. Make sure to use the model from previous output "meta/llama-3.1-8b-instruct". In this case we are asking a question about Hitachi VSP One SDS storage.

```
headers = \{
    "Content-Type": "application/json"
}
payload = {
    "model": "meta/llama-3.1-8b-instruct",
    "messages": [
        {
            "role": "user",
            "content": "What is Hitachi VSP One SDS?"
        }
    ],
    "temperature": 0.5,
    "top p": 1,
    "max tokens": 1024,
    "stream": False
}
response = requests.post(nim_url + "/v1/chat/completions", verify=False,
headers=headers, json=payload)
print(response.json()['choices'][0]['message']['content'])
```

File Edit View Run	Kernel Git Tabs Sett	ings Heip	
🗈 ± C		× R genai-nvidia-nim-example1.× +	
Filter files by name	Q = + %	🗌 🗂 🕨 🔳 C 🁐 Markdown 🗠 Run as Pipeline 🛈 git	Di Python 3.11 ()
/ gen-ai-demo /		Example notebook. Invokes the NVIDIA NIM	© ↑ ↓ 古 두 🕯
Name 🔺	Modified	Make sure to use the usi as shown on the NV/DIA conving model and point	
demo-coco.ipynb	2d ago	Make sure to use the un as shown on the INNDIA serving model end-point	
📃 demo1-african-w	2d ago [1]:	import json	
🔲 genai-nvidia-nim	yesterday	import requests import urllib3	
		urllib3.disable warnings()	
		<pre>nim_url = "https://nvidia-nim-model-llama3-demo.nvidia-nim-openshift-ai-demo.svc.cluster.local"</pre>	
		Fetch the list of available models.	
	[2]:	<pre>response = requests.get(nim url + "/v1/models", verify=False)</pre>	
		<pre>print(json.dumps(response.json()['data'], indent=2))</pre>	
		t j	
		{ "id": "meta/llama-3.1-8b-instruct",	
		"object": "model",	
		"owned_by": "system",	
		"root": "meta/llama-3.1-8b-instruct",	
		"max_model_len": 131072,	
		"permission": [
		"id": "modelperm-27b2704080674e848eece5efa2332ff2",	
		"object": "model_permission",	
		"allow create engine": false,	
		"allow_sampling": true,	
		"allow_logprobs": true, "allow search indices": false.	
		"allow_view": true,	
		"allow_fine_tuning": false,	
		"organization": """, "group": null.	
		"is_blocking": false	
		}	
		1	
		Send a chat request to the target model, use the model from previous output. (meta/llama-3.1-8b-instruct)	

The following screenshot shows the code snippets inside the notebook.

Here we can see the response:





As we can see from this test case, it is very easy to setup an NVIDIA NIM serving model and access it using Red Hat OpenShift AI. This integration allows development teams to

rapidly build and scale modern enterprise applications leveraging a set of easy-to-use inference microservices included as part of the NVIDIA AI Enterprise software platform.

Model training from a pretrained model using OpenShift AI

For this test case, we are going to show the process to re-train a model using a pretrained YOLO model (Ultralytics YOLO11 is a real-time object detection model) and a dataset of football players, export and reuse the trained model. The original dataset was obtained from RoboFlow and it has annotated images, split into test, train, and validation datasets.

The test case consists of the following:

- Retrain the model using the imported dataset.
- Verify that the retrained model can do object detection.
- Export the model to ONNX format (Open Neural Network Exchange).
- Upload the model to an S3 bucket to be reused in subsequent tests.
- Download and test the retrained model in a separate environment.

Create a Data Science Project

Procedure

- 1. From the OpenShift AI dashboard, navigate to the **Data Science Projects** tab.
- 2. Click **Create Project**, and enter a name and description for the project.
- 3. Click Create.

Inside the created project is where all the additional resources for this test case will be deployed.

Configure and create a Workbench

The next step is to create a workbench which is an instance of development and experimentation environment. Within a workbench we can select a notebook image for your data science work.

Procedure

- 1. Select the newly created project, then click the Workbenches tab.
- 2. Click Create workbench, and then enter a name and description.
- **3.** Select one of the existing notebook images. In this case we select "PyTorch" and version 2024.2.
- **4.** Select the container image and accelerator. In this case we are using a Large container size and NVIDIA GPU.

Jump to section	Notebook image						
	Image selection *						
Name and description	PyTorch 🔹						
Notebook image	Version selection *						
	2024.2						
Deployment size	CUDA v12.4, Python v3.11, PyTorch v2.4						
	Hover over a version to view its included packages.						
Environment variables	View package information						
Cluster storage							
Data connections	Deployment size						
Bata connectione	Container size						
Connections	Large 🔹						
	Limits: 14 CPU, 56GiB Memory Requests: 7 CPU, 56GiB Memory						
	Accelerator						
	NVIDIA GPU -						
Create workbench Cancel							

- 5. On the cluster storage enter the size of the persistent volume to be created. The PVC will be created on the storage being used by the default Storage Class. For Persistent storage we are using Hitachi VSP storage and Hitachi Storage Plugin for Containers (HSPC).
- **6.** (Optional) In addition to the cluster storage, add a connection to S3 storage. This can be used later to upload or download files from an S3 bucket.
 - **Note:** A data connection is used to store configuration parameters that are required to connect to a data source, such as an S3 object bucket. This can be used for storing your models and data. You can reuse this bucket and its connection for your notebooks and model servers.

To create data connections to your existing S3-compatible storage buckets, you need the following credential information for the storage buckets: Endpoint URL, Access key, Secret key, Region, Bucket name.

7. It can take a couple of minutes to have the workbench in a running state, and then you can access the workbench interface by clicking the **Open** link. Log in with the OpenShift credentials.

demo-training ⑦ Cluster storage G demo-training Max 1200iB Mount path: /opt/app-root/src	PyTorch CUDA VI2.4, Python v3.11, PyTorch v2.4 Packages PyTorch v2.4 Tensorboard v2.17	Large Limits	Running	Stop Open 🗹
Cluster storage demo-training Max 120GiB Mount path: /opt/app-root/src	Packages PyTorch v2.4 Tensorboard v2.17	Limits		
	Boto3 v1.35 Kafka-Python-ng v2.2 Kfp v2.9 Matplotlib v3.9 Numpy v2.1 Pandas v2.2 Scikit-learn v1.5 Scipy v1.4 Odh-Elyra v4.2 PyMongo v4.8 Pyodbc v5.1 Codeflare-SDK v0 23 Sklearn-onnx v1.17 Psycopg v3.2 MySGL Connector/Python v5.0	14 CPU, 56GiB Memory Requests 7 CPU, 56GiB Memory		
ce Projects > demo-training				Ac
w Workbenches Pipelines	Models Cluster storage	Connections	Permissions	Settings

Create a notebook to retrain a model

The next step after creating the workbench is to deploy a notebook and work with an AI model. In this section we are going to show some of these steps:

Procedure

1. The first window that opens after opening the workbench is a Jupiter environment similar to the following.

~	→ C S Not secure https://de	mo-training-demo-training.apps.ocpjpai.ocp.sce.lab/no 🛧 🔷 🖘 🖆 📔 😩	New Chrome availa
\bigcirc	File Edit View Run Kernel Git	Tabs Settings Help	
	You are not currently in a Git repository. To use Git, navigate to a	Zauncher +	
0	here, or clone an existing repository.	Notebook	
•>	Open the FileBrowser Initialize a Repository	2	
≣	Clone a Repository	Python 3.11	
		Console	
£			
ß		.	
88		Python 3.11	
*		💓 Elyra	
		Pipeline Editor Pipeline Editor	
		\$_ Other	

- **2.** While we can always write a notebook from scratch, the easiest way is to import a Git repository. In this case, we are cloning the same repository used in the previous use case.
- **3.** After the repo has been cloned, and before running the notebook to retrain a model, the first step is to import the dataset.
- 4. Import dataset:

From the cloned repo, navigate to notebooks folder and open the "download_extract_dataset" notebook. Make sure the update the URL based on the location of the dataset. In this case we are using a dataset obtained from RoboFlow.

Run the notebook and wait until the dataset is downloaded and extracted to the corresponding folder.

5. Retrain the model:

Now it is time to retrain the model. Open the "football-players-detection" notebook.

This notebook example shows how to retrain a pre-trained YOLO model using a dataset obtained from Roboflow. Also, a set of unseen images has been included to test the retrained model.

The first cell is to verify that we can access the NVIDIA GPU.

🛙 Launcher		× 5_ 100	0860000@demo-traini	in 🗙 📃 football-players-det	ection.ir× 🖪 upload_mod	el_to_s3.ipynb × +
B + %	[[] ►	■ C >>	Code 🗸 Run a	s Pipeline 🕚 git		₿ Python 3.11 ()
[1]:	!nvidia-smi					
	Fri Feb 14 1	19:38:11 2025	5			
	NVIDIA-SMI	550.144.03	Driver	Version: 550.144.03	CUDA Version: 12.4	
	GPU Name Fan Temp	Perf	Persistence-M Pwr:Usage/Cap	Bus-Id Disp.A Memory-Usage	Volatile Uncorr. ECC GPU-Util Compute M.	
				 +====================================	+=====================================	
	0 NVID] N/A 70C	IA L405 P0	On 128W / 350W	00000000:23:00.0 Off 3197MiB / 46068MiB	0 0% Default N/A	
	+			+	++	
	+				++	
	GPU GI	CI ID	PID Type Proces	ss name	GPU Memory Usage	

The next step is to clone the Ultralytics GitHub repository, install, and import YOLO.



After the required packages have been installed, the next cell is to verify the dataset. The output shows the folder for the dataset, including the test, training, and validation subfolders.



Here we start retraining the model using the imported dataset. As seen below, the first step is to load a pretrained model, in this case are using YOLO11n which is the smallest YOLO11 models. Next, we call the "train" function, change the parameters according to your test. For example, you can adjust the parameter "epoch" depending how many iterations you want to tune during the training.



From the output we can see which GPU is being used, in this case we are using NVIDIA L40 GPU. At the completion of the training, the weights will automatically be saved in

the "runs/detect/train" folder.

2 L	auno	her				\times	🖪 do	wnload_extr	act_dataset.	i 🗙 📃 footb	all-players-de	etection.ir×	+		
8	+	Ж	ē	Ċ	►	•	C 🕨	Code	✓ Run a	s Pipeline 🕚	git			÷	Python 3.11 🔘
				Ep	och	GI	PU_mem	box_loss	s cls_los	ss dfl_loss	Instances	Size			
			0,	20 18.53	/20 it/s	1	2.27G	1.462	0.763	37 1.221	. 118	640	: 100%	925/9	925 [00:49<00:0
			0:0	4<00:	90,	14.27	Class 7it/s]	Images	s Instance	es Box(P	P R	mAP50	mAP50-95):	100%	68/68 [0
							all	2176	5 1934	43 0.863	0.826	0.864	0.5		
		ſ	Opt Val Ult YOL	imize idati ralyt 011n	ng /	opt/a 8.3.3	app-roo 5 <mark>2</mark> (fused)	ot/src/gena Python-3.1 : 100 laye	root/src/g ai-demo/not 1.7 torch- ers, 2,582,	tebooks/ultra 2.4.1+cu121 932 paramete	olytics/ult olytics/runs CUDA:0 (NVII ers, 0 gradi	//detect/tra DIA L405, 45 ents, 6.3 G	ns/detect/tra in/weights/b 5589MiB) FLOPs	ain/weights/ est.pt	/best.pt, 5.5MB
			0:0	7<00:	90,	9.68	Class Bit/s]	Images	Instance	es Box(P	P R	mAP50	mAP50-95):	100%	68/68 [0
							all	2176	5 1934	43 0.863	0.826	0.864	0.5		
						Re	eferee	1008	3 128	81 0.831	0.88	0.895	0.532		
							Ball	1216	5 123	36 0.833	0.591	0.691	0.335		
						Goal	keeper	1144	117	74 0.889	0.914	0.926	0.567		
						F	Player	2074	1565	52 0.899	0.921	0.946	0.567		
		1	Spe	ed: 0	.1ms	prep	process	5, 0.3ms in	iference, @	0.0ms loss, 0	.4ms postpr	ocess per i	mage		
		- N.	Res	ults	ave	d to	/opt/a	app-root/si	∙c/genai-de	emo/notebooks	/ultralytic	s/runs/dete	ct/train		

After we have the re-trained model we can start making predictions with new images. For this purpose, we have included a set of unseen images on the test folder outside of the dataset.

Because we retrained this model with the football players dataset, the new re-trained model should be able to detect players, goalkeepers, and even referees on an image of a football (or soccer) field.

In this case we are testing with a sample image "football image1.jpg". Modify the path if needed.

Here we are loading the updated weights file to the model and passing a sample image, and as we can see the model is able to predict the images correctly, detecting goalkeeper, players, and ball.



Export the model and upload to S3 storage

Now that we have verified that the retrained model works as expected, we can export and save the model in ONNX format so it can be used for inferencing in another container image on prem or in cloud. ONNX stands for Open Neural Network Exchange which is an open standard for representing machine learning models.

Procedure

1. Export and save the model.

To export the model in ONNX format we use the following command:

Z I	aun	cher				×	\$_	100	0860000@	0dem	o-trainin \times	💌 footb	all-player:	s-detection.ir:	×	💌 upload	_mo	del_	_to_s	s3.ipy	nb >	< +
8	+	Ж	Ū	Ĉ	►	•	C	**	Code	\sim	Run as Pipe	line 🕓	git					ŧ	Ŭ I	Pytho	n 3.1	10
	I	7]	mod	el.ex	kport	- t(for	mat	dona	1 <mark>x')</mark> #.ex	kport	to ONNX					1	9 4	↑	\downarrow	÷	Ŧ	Î

The following output shows the path where the ONNX file is stored. This file can easily be copied to another location so it can be used in another container.

🛙 Launcher	×	s_ 1000	0860000@dem	o-trainin×	🖲 football-playe	rs-detection.ir×	🗏 upload_model_to_s3.ipynb 🗙 🕂	+
B + % 🗇 (" ▶ ■	C ₩	Code 🗸 🗸	Run as Pipe	eline 🛈 git		₿ Python 3.11 〇) =
ONNX: ONNX: ONNX: ONNX: N/Wei Expor Resul Predd ct/tr Valid rain/ Visua [7]: '/opt	starting e slimning w export succ ghts/best.c t complete ts saved to ct: ain/weights ate: weights/bes lize: /app-root/s	xport wi ith onnx cess onnx' (10 (10.0s) /opt/ap yolo pr /best.on yolo va t.onnx ii https:/	th onnx 1.17 slim 0.1.48. 8.75, saved .1 MB) p-root/src/ge edict task-de nx imgsz=640 l task=detect mgsz=640 datz /netron.app -demo/noteboo	0 opset 19 as '/opt/a enai-demo/r stect model t model=/op s=/opt/app- oks/ultraly	<pre> pp-root/src/gen otebooks/ultral =/opt/app-root/ t/app-root/src/ root/src/genai- tics/runs/detec</pre>	ai-demo/noteboo ytics/runs/dete src/genai-demo/ genai-demo/note demo/dataset/de	ks/ultralytics/runs/detect/trai ect/train/weights notebooks/ultralytics/runs/dete ebooks/ultralytics/runs/detect/t ita.yaml	

2. Upload the model to Hitachi VSP One Object (S3 storage).

Inside the cloned git repository is a notebook file that can be used to upload files to an S3 storage. We are using it to upload to Hitachi S3 storage so it can be used in the next use case in a different workbench.

HITACHI S3 Console				8
< genai-sea genai-sea / genai-demo /				
Browse Overview Properties			Create D	Directory Upload Object. (5)
				🗍 Delete 🕑 Destroy
Name	Size	Modified	Owner	Storage class
best.onnx	10.1 MB	02/13/2025 17:33:21		STANDARD :

Testing a retrained model

Now that we have retrained, exported the ONNX model, and placed a copy within an S3 bucket we are going to test the retrained model against another set of unseen images using a new workbench/container.

Procedure

1. Create a new workbench inside the same project. The project is called "demo-trainingretest", and we have chosen the notebook image "Standard Data Science" and no GPU.

Data Science Projects > demo-training										
🗋 den	no-training ®								Act	ions 🔻
Overview	Workbenches	Pipelines	Models	Cluster st	orage	Connecti	ons Perr	nissions	Settings	
& Worl	kbenches 🔊		Notebook image	e	Containe	er size	Status 🗍		Create wo	orkbench
> di	emo-training ③		PyTorch CUDA v12.4, Python PyTorch v2.4	v 3.11,	Large		Running	Stop	Open 🗹	:
> d	emo-training-retest 💿		Standard Data S Python v3.11	cience	Small		Running	Stop	Open 🗹	÷

2. To verify the S3 connections, click the **Connections** tab and you will see that the S3 connections are used by both workbenches.

Data Science Pro	ojects > demo-train	ning										
🖆 demo-training ®												
Overview	Workbenches	Pipelines	Models	Cluste	r storage	Connections	Permissions	Settings				
(^(o)) Conne	ections @								Add connection			
Name †		Туре	1		Model ser	ving compatibility	Connected	resources				
hitachi-s3	0	S3 co stora	ompatible o ge - v1	bject	S3 comp	atible object storage	demo-t	:				

- 3. Click Open and it will redirect to a new window with a Jupiter environment.
- 4. Use Git to clone the same repository.

To test the retrained model, use the notebook "test_retrained_model" from the cloned repository.

Here are some details of the cells of this notebook:

The first cell is to install the required packages.



Next, we get the retrained model, exported in the ONNX format, from the S3 bucket. In the output below we can see the ONNX model is being downloaded and stored locally.



The result shows that the retrained model obtained from the S3 bucket can do image detection and works as expected, detecting, and differentiating between players, goalkeeper, and referee.

Conclusion

This reference architecture also provides a reference pattern for an AI enabled platform using Hitachi Unified Compute Platform with Red Hat AI, OpenShift AI, and NVIDIA GPUs leveraging VSP One Storage Platforms and Hitachi's NVIDIA certified data center HA compute servers. It leverages the latest capabilities and services provided by Red Hat including OpenShift AI Operators, InstructLab, Granite, and NVIDIA NIM microservices

RHEL AI, along with Red Hat OpenShift AI, underpins Red Hat AI, Red Hat's portfolio of solutions that accelerate time to market and reduce the operational cost of delivering AI solutions across the hybrid cloud.

Red Hat OpenShift AI enhances RHEL AI with Kubernetes-based containerized orchestration and scalability. It is a platform that empowers enterprises to develop and deploy AI-driven solutions efficiently across hybrid cloud environments, facilitating the creation and delivery of AI-enabled applications at scale with high reliability and security. Additionally, it supports model serving and hosting, making real-time AI inference and integration seamless.

A key benefit we concluded with this reference pattern is that it helps to accelerate the time to value of AI projects. It significantly eliminates the time to find and deploy the right combination of software packages suitable to run AI/ML workloads and minimizes the expertise required to have an environment up and running to start prototyping/deploying AI/ML applications. The combination of Hitachi automated integrated system with Red Hat OpenShift Software stack will help data scientists and data engineers eliminate spending too much time on the infrastructure/DevOps processes and instead use that time to focus on developing/training, fine-tuning their models and deploying/integrating those models into workloads/applications in production environments.

With this enterprise-grade container orchestration platform, IT and DevOps teams have a set of easy-to-use tools for accelerating the deployment of AI/ML models on-premises and/or the public cloud. Organizations can self-host large language models (LMMs) in their own managed environments allowing developers to quickly build enterprise-grade AI applications with minimum effort.

Some of the key benefits of this platform include:

- Flexible and high-performance environment
- Simplified deployment, reducing the time required to retrain, test, and validate AI models
- Provides capabilities to manage the entire lifecycle of AI/ML projects

Overall, this solution provides a foundational platform for easy deployment and scalability of AI/ML workloads.

Product descriptions

This section provides information about the hardware and software components used in this solution.

Hitachi Virtual Storage Platform One Block

The Hitachi Virtual Storage Platform One Block series simplifies system setup and management through Hitachi Clear Sight and VSP One Block Administrator. Dynamic Drive Protection reduces RAID complexity, and always-on compression and deduplication enhance simplicity.

Dynamic Carbon Reduction optimizes energy usage by switching CPUs to ECO mode during low activity. Adaptive Data Reduction (ADR) is always on, enhancing efficiency and reducing the overall CO2 footprint.

Thin Image Advanced (TIA) integrates with major snapshot ecosystems, prioritizing security by defending against threats and ensuring data confidentiality. CyberArk Privileged Access Manager plugins enhance block storage system security by prioritizing data confidentiality, ensuring compliance, and actively defending against security threats.

<u>Hitachi Virtual Storage Platform One Block 20</u> includes 3 dedicated models. All have the same capacity (72 NVMe flash drives, the appliance and 2 × media trays) and they support Fibre Channel, iSCSI, and NVMe TCP connectivity. The new capabilities remove complexity: data reduction is always on, Dynamic Drive Protection removes complicated RAID setup, and Dynamic Carbon Reduction delivers real world reduction in power consumption. In addition, the models are FIPS compliant.

- VSP One Block 24 256 GB Cache + SW Advanced Data Reduction (ADR) + 24 cores
- VSP One Block 26 768GB Cache + 2x Compression Accelerator Module (CAM) + 24 cores
- VSP One Block 28 1TB Cache + 4x CAM + 64 cores

In short, the Hitachi Virtual Storage Platform One Block series combines simplicity, sustainability, and robust security features to optimize system management, energy efficiency, and data protection.

Hitachi Virtual Storage Platform One Object

VSP One Object offers massive scale, cloud capabilities, and broad protocol support. It supports all-flash configurations and policy-driven tiered storage, allowing performance and capacity to scale independently. It quickly ingests large volumes of small objects and cost-effectively stores very large objects.

Customers benefit from a diverse array of configurations and deployment options. These include traditional archive and backup solutions, as well as modern cloud-native environments, AI/ML applications, and analytics workloads. The platform is optimized for a compact form factor, fitting an entire cluster into a 4U footprint. This enables real-time analytics at the edge, eliminating the need to transfer data to a central location for analysis.

Hitachi Advanced Server

Designed to unlock the full benefits of the hybrid cloud, Hitachi Advanced Server models deliver high performance and enhanced security while reducing operational costs. Global enterprises, cloud service providers, and governments trust Hitachi servers to run bare metal, virtualized, or containerized applications. Powered by industry-leading scalable processors, Hitachi servers are ideal to deliver edge, core, and cloud IT services.

Hitachi Advanced Server supports a variety of GPUs for parallel processing which is used in a wide range of applications, including graphics and video rendering.

Hitachi servers are designed and optimized to maximize performance for VMware, Red Hat, Oracle, bare metal, virtual desktop infrastructure (VDI), SAP, analytics, high-performance computing (HPC), and other enterprise workloads.



Hitachi Vantara

Corporate Headquarters 2535 Augustine Drive Santa Clara, CA 95054 USA HitachiVantara.com/contact