



# OpenShift Container Platform 4.17

## Hardware accelerators

Hardware accelerators



# OpenShift Container Platform 4.17 Hardware accelerators

---

Hardware accelerators

## Legal Notice

Copyright © Red Hat.

Except as otherwise noted below, the text of and illustrations in this documentation are licensed by Red Hat under the Creative Commons Attribution–Share Alike 3.0 Unported license . If you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, the Red Hat logo, JBoss, Hibernate, and RHCE are trademarks or registered trademarks of Red Hat, LLC. or its subsidiaries in the United States and other countries.

Linux<sup>®</sup> is the registered trademark of Linus Torvalds in the United States and other countries.

XFS is a trademark or registered trademark of Hewlett Packard Enterprise Development LP or its subsidiaries in the United States and other countries.

The OpenStack<sup>®</sup> Word Mark and OpenStack logo are trademarks or registered trademarks of the Linux Foundation, used under license.

All other trademarks are the property of their respective owners.

## Abstract

This document provides instructions for installing and configuring the GPU Operators supported by Red Hat OpenShift AI for the provided hardware acceleration capabilities for creating artificial intelligence and machine learning (AI/ML) applications.

---

## Table of Contents

<b>CHAPTER 1. ABOUT HARDWARE ACCELERATORS</b> .....	<b>3</b>
1.1. HARDWARE ACCELERATORS	4
<b>CHAPTER 2. NVIDIA GPU ARCHITECTURE</b> .....	<b>5</b>
2.1. NVIDIA GPU PREREQUISITES	5
2.2. NVIDIA GPU ENABLEMENT	5
2.2.1. GPUs and bare metal	6
2.2.2. GPUs and virtualization	7
2.2.3. GPUs and vSphere	7
2.2.4. GPUs and Red Hat KVM	7
2.2.5. GPUs and CSPs	8
2.2.6. GPUs and Red Hat Device Edge	8
2.3. GPU SHARING METHODS	9
2.3.1. CUDA streams	9
2.3.2. Time-slicing	10
2.3.3. CUDA Multi-Process Service	10
2.3.4. Multi-instance GPU	10
2.3.5. Virtualization with vGPU	11
2.4. NVIDIA GPU FEATURES FOR OPENSIFT CONTAINER PLATFORM	11
<b>CHAPTER 3. AMD GPU OPERATOR</b> .....	<b>13</b>
3.1. ABOUT THE AMD GPU OPERATOR	13
3.2. INSTALLING THE AMD GPU OPERATOR	13
3.3. TESTING THE AMD GPU OPERATOR	13

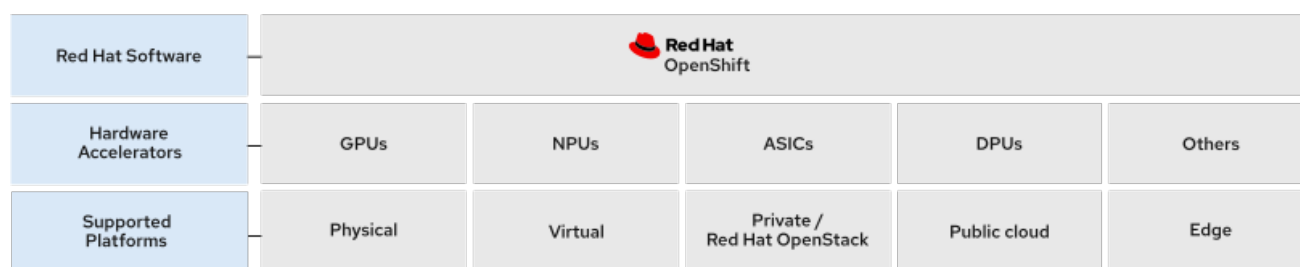


# CHAPTER 1. ABOUT HARDWARE ACCELERATORS

Specialized hardware accelerators play a key role in the emerging generative artificial intelligence and machine learning (AI/ML) industry. Specifically, hardware accelerators are essential to the training and serving of large language and other foundational models that power this new technology. Data scientists, data engineers, ML engineers, and developers can take advantage of the specialized hardware acceleration for data-intensive transformations and model development and serving. Much of that ecosystem is open source, with several contributing partners and open source foundations.

Red Hat OpenShift Container Platform provides support for cards and peripheral hardware that add processing units that comprise hardware accelerators:

- Graphical processing units (GPUs)
- Neural processing units (NPUs)
- Application-specific integrated circuits (ASICs)
- Data processing units (DPUs)



Specialized hardware accelerators provide a rich set of benefits for AI/ML development:

## One platform for all

A collaborative environment for developers, data engineers, data scientists, and DevOps

## Extended capabilities with Operators

Operators allow for bringing AI/ML capabilities to OpenShift Container Platform

## Hybrid-cloud support

On-premise support for model development, delivery, and deployment

## Support for AI/ML workloads

Model testing, iteration, integration, promotion, and serving into production as services

Red Hat provides an optimized platform to enable these specialized hardware accelerators in Red Hat Enterprise Linux (RHEL) and OpenShift Container Platform platforms at the Linux (kernel and userspace) and Kubernetes layers. To do this, Red Hat combines the proven capabilities of Red Hat OpenShift AI and Red Hat OpenShift Container Platform in a single enterprise-ready AI application platform.

Hardware Operators use the operating framework of a Kubernetes cluster to enable the required accelerator resources. You can also deploy the provided device plugin manually or as a daemon set. This plugin registers the GPU in the cluster.

Certain specialized hardware accelerators are designed to work within disconnected environments where a secure environment must be maintained for development and testing.

## 1.1. HARDWARE ACCELERATORS

Red Hat OpenShift Container Platform enables the following hardware accelerators:

- NVIDIA GPU
- AMD Instinct® GPU
- Intel® Gaudi®

### Additional resources

- [Introduction to Red Hat OpenShift AI](#)
- [NVIDIA GPU Operator on Red Hat OpenShift Container Platform](#)
- [AMD Instinct Accelerators](#)
- [Intel Gaudi AI Accelerators](#)

## CHAPTER 2. NVIDIA GPU ARCHITECTURE

NVIDIA supports the use of graphics processing unit (GPU) resources on OpenShift Container Platform. OpenShift Container Platform is a security-focused and hardened Kubernetes platform developed and supported by Red Hat for deploying and managing Kubernetes clusters at scale. OpenShift Container Platform includes enhancements to Kubernetes so that users can easily configure and use NVIDIA GPU resources to accelerate workloads.

The NVIDIA GPU Operator uses the Operator framework within OpenShift Container Platform to manage the full lifecycle of NVIDIA software components required to run GPU-accelerated workloads.

These components include the NVIDIA drivers (to enable CUDA), the Kubernetes device plugin for GPUs, the NVIDIA Container Toolkit, automatic node tagging using GPU feature discovery (GFD), DCGM-based monitoring, and others.



### NOTE

The NVIDIA GPU Operator is only supported by NVIDIA. For more information about obtaining support from NVIDIA, see [Obtaining Support from NVIDIA](#).

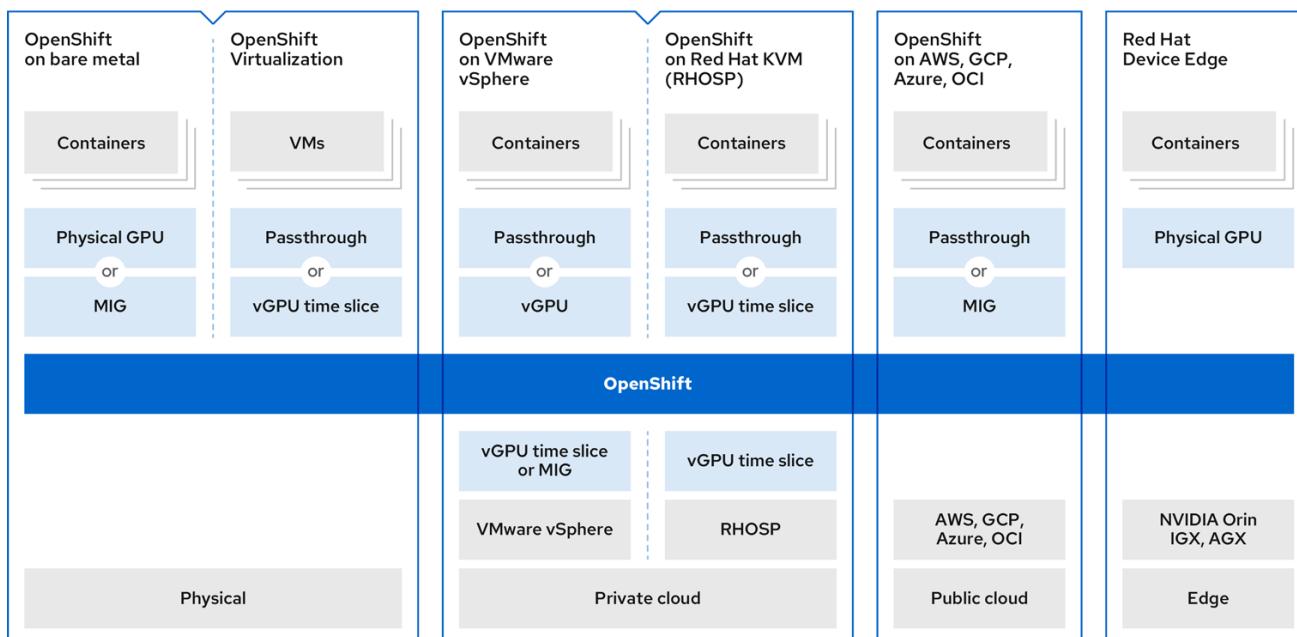
### 2.1. NVIDIA GPU PREREQUISITES

- A working OpenShift cluster with at least one GPU worker node.
- Access to the OpenShift cluster as a **cluster-admin** to perform the required steps.
- OpenShift CLI (**oc**) is installed.
- The node feature discovery (NFD) Operator is installed and a **nodefeaturediscovery** instance is created.

### 2.2. NVIDIA GPU ENABLEMENT

The following diagram shows how the GPU architecture is enabled for OpenShift:

Figure 2.1. NVIDIA GPU enablement



512\_OpenShift\_1223

**NOTE**

MIG is supported on GPUs starting with the NVIDIA Ampere generation. For a list of GPUs that support MIG, see the [NVIDIA MIG User Guide](#).

**2.2.1. GPUs and bare metal**

You can deploy OpenShift Container Platform on an NVIDIA-certified bare metal server but with some limitations:

- Control plane nodes can be CPU nodes.
- Worker nodes must be GPU nodes, provided that AI/ML workloads are executed on these worker nodes.  
In addition, the worker nodes can host one or more GPUs, but they must be of the same type. For example, a node can have two NVIDIA A100 GPUs, but a node with one A100 GPU and one T4 GPU is not supported. The NVIDIA Device Plugin for Kubernetes does not support mixing different GPU models on the same node.
- When using OpenShift, note that one or three or more servers are required. Clusters with two servers are not supported. The single server deployment is called single node openShift (SNO) and using this configuration results in a non-high availability OpenShift environment.

You can choose one of the following methods to access the containerized GPUs:

- GPU passthrough
- Multi-Instance GPU (MIG)

**Additional resources**

- [Red Hat OpenShift on Bare Metal Stack](#)

### 2.2.2. GPUs and virtualization

Many developers and enterprises are moving to containerized applications and serverless infrastructures, but there is still a lot of interest in developing and maintaining applications that run on virtual machines (VMs). Red Hat OpenShift Virtualization provides this capability, enabling enterprises to incorporate VMs into containerized workflows within clusters.

You can choose one of the following methods to connect the worker nodes to the GPUs:

- GPU passthrough to access and use GPU hardware within a virtual machine (VM).
- GPU (vGPU) time-slicing, when GPU compute capacity is not saturated by workloads.

#### Additional resources

- [NVIDIA GPU Operator with OpenShift Virtualization](#)

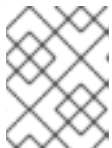
### 2.2.3. GPUs and vSphere

You can deploy OpenShift Container Platform on an NVIDIA-certified VMware vSphere server that can host different GPU types.

An NVIDIA GPU driver must be installed in the hypervisor in case vGPU instances are used by the VMs. For VMware vSphere, this host driver is provided in the form of a VIB file.

The maximum number of vGPUS that can be allocated to worker node VMs depends on the version of vSphere:

- vSphere 7.0: maximum 4 vGPU per VM
- vSphere 8.0: maximum 8 vGPU per VM



#### NOTE

vSphere 8.0 introduced support for multiple full or fractional heterogeneous profiles associated with a VM.

You can choose one of the following methods to attach the worker nodes to the GPUs:

- GPU passthrough for accessing and using GPU hardware within a virtual machine (VM)
- GPU (vGPU) time-slicing, when not all of the GPU is needed

Similar to bare metal deployments, one or three or more servers are required. Clusters with two servers are not supported.

#### Additional resources

- [OpenShift Container Platform on VMware vSphere with NVIDIA vGPUs](#)

### 2.2.4. GPUs and Red Hat KVM

You can use OpenShift Container Platform on an NVIDIA-certified kernel-based virtual machine (KVM) server.

Similar to bare-metal deployments, one or three or more servers are required. Clusters with two servers are not supported.

However, unlike bare-metal deployments, you can use different types of GPUs in the server. This is because you can assign these GPUs to different VMs that act as Kubernetes nodes. The only limitation is that a Kubernetes node must have the same set of GPU types at its own level.

You can choose one of the following methods to access the containerized GPUs:

- GPU passthrough for accessing and using GPU hardware within a virtual machine (VM)
- GPU (vGPU) time-slicing when not all of the GPU is needed

To enable the vGPU capability, a special driver must be installed at the host level. This driver is delivered as a RPM package. This host driver is not required at all for GPU passthrough allocation.

### 2.2.5. GPUs and CSPs

You can deploy OpenShift Container Platform to one of the major cloud service providers (CSPs): Amazon Web Services (AWS), Google Cloud, or Microsoft Azure.

Two modes of operation are available: a fully managed deployment and a self-managed deployment.

- In a fully managed deployment, everything is automated by Red Hat in collaboration with CSP. You can request an OpenShift instance through the CSP web console, and the cluster is automatically created and fully managed by Red Hat. You do not have to worry about node failures or errors in the environment. Red Hat is fully responsible for maintaining the uptime of the cluster. The fully managed services are available on AWS, Azure, and Google Cloud. For AWS, the OpenShift service is called ROSA (Red Hat OpenShift Service on AWS). For Azure, the service is called Azure Red Hat OpenShift. For Google Cloud, the service is called OpenShift Dedicated on Google Cloud.
- In a self-managed deployment, you are responsible for instantiating and maintaining the OpenShift cluster. Red Hat provides the OpenShift-install utility to support the deployment of the OpenShift cluster in this case. The self-managed services are available globally to all CSPs.

It is important that this compute instance is a GPU-accelerated compute instance and that the GPU type matches the list of supported GPUs from NVIDIA AI Enterprise. For example, T4, V100, and A100 are part of this list.

You can choose one of the following methods to access the containerized GPUs:

- GPU passthrough to access and use GPU hardware within a virtual machine (VM).
- GPU (vGPU) time slicing when the entire GPU is not required.

#### Additional resources

- [Red Hat Openshift in the Cloud](#)

### 2.2.6. GPUs and Red Hat Device Edge

Red Hat Device Edge provides access to MicroShift. MicroShift provides the simplicity of a single-node deployment with the functionality and services you need for resource-constrained (edge) computing. Red Hat Device Edge meets the needs of bare-metal, virtual, containerized, or Kubernetes workloads deployed in resource-constrained environments.

You can enable NVIDIA GPUs on containers in a Red Hat Device Edge environment.

You use GPU passthrough to access the containerized GPUs.

### Additional resources

- [How to accelerate workloads with NVIDIA GPUs on Red Hat Device Edge](#)

## 2.3. GPU SHARING METHODS

Red Hat and NVIDIA have developed GPU concurrency and sharing mechanisms to simplify GPU-accelerated computing on an enterprise-level OpenShift Container Platform cluster.

Applications typically have different compute requirements that can leave GPUs underutilized. Providing the right amount of compute resources for each workload is critical to reduce deployment cost and maximize GPU utilization.

Concurrency mechanisms for improving GPU utilization exist that range from programming model APIs to system software and hardware partitioning, including virtualization. The following list shows the GPU concurrency mechanisms:

- Compute Unified Device Architecture (CUDA) streams
- Time-slicing
- CUDA Multi-Process Service (MPS)
- Multi-instance GPU (MIG)
- Virtualization with vGPU

Consider the following GPU sharing suggestions when using the GPU concurrency mechanisms for different OpenShift Container Platform scenarios:

### Bare metal

vGPU is not available. Consider using MIG-enabled cards.

### VMs

vGPU is the best choice.

### Older NVIDIA cards with no MIG on bare metal

Consider using time-slicing.

### VMs with multiple GPUs and you want passthrough and vGPU

Consider using separate VMs.

### Bare metal with OpenShift Virtualization and multiple GPUs

Consider using pass-through for hosted VMs and time-slicing for containers.

### Additional resources

- [Improving GPU Utilization](#)

### 2.3.1. CUDA streams

Compute Unified Device Architecture (CUDA) is a parallel computing platform and programming model developed by NVIDIA for general computing on GPUs.

A stream is a sequence of operations that executes in issue-order on the GPU. CUDA commands are typically executed sequentially in a default stream and a task does not start until a preceding task has completed.

Asynchronous processing of operations across different streams allows for parallel execution of tasks. A task issued in one stream runs before, during, or after another task is issued into another stream. This allows the GPU to run multiple tasks simultaneously in no prescribed order, leading to improved performance.

#### Additional resources

- [Asynchronous Concurrent Execution](#)

### 2.3.2. Time-slicing

GPU time-slicing interleaves workloads scheduled on overloaded GPUs when you are running multiple CUDA applications.

You can enable time-slicing of GPUs on Kubernetes by defining a set of replicas for a GPU, each of which can be independently distributed to a pod to run workloads on. Unlike multi-instance GPU (MIG), there is no memory or fault isolation between replicas, but for some workloads this is better than not sharing at all. Internally, GPU time-slicing is used to multiplex workloads from replicas of the same underlying GPU.

You can apply a cluster-wide default configuration for time-slicing. You can also apply node-specific configurations. For example, you can apply a time-slicing configuration only to nodes with Tesla T4 GPUs and not modify nodes with other GPU models.

You can combine these two approaches by applying a cluster-wide default configuration and then labeling nodes to give those nodes a node-specific configuration.

### 2.3.3. CUDA Multi-Process Service

CUDA Multi-Process Service (MPS) allows a single GPU to use multiple CUDA processes. The processes run in parallel on the GPU, eliminating saturation of the GPU compute resources. MPS also enables concurrent execution, or overlapping, of kernel operations and memory copying from different processes to enhance utilization.

#### Additional resources

- [CUDA MPS](#)

### 2.3.4. Multi-instance GPU

Using Multi-instance GPU (MIG), you can split GPU compute units and memory into multiple MIG instances. Each of these instances represents a standalone GPU device from a system perspective and can be connected to any application, container, or virtual machine running on the node. The software that uses the GPU treats each of these MIG instances as an individual GPU.

MIG is useful when you have an application that does not require the full power of an entire GPU. The MIG feature of the new NVIDIA Ampere architecture enables you to split your hardware resources into multiple GPU instances, each of which is available to the operating system as an independent CUDA-

enabled GPU.

NVIDIA GPU Operator version 1.7.0 and higher provides MIG support for the A100 and A30 Ampere cards. These GPU instances are designed to support up to seven multiple independent CUDA applications so that they operate completely isolated with dedicated hardware resources.

#### Additional resources

- [NVIDIA Multi-Instance GPU User Guide](#)

### 2.3.5. Virtualization with vGPU

Virtual machines (VMs) can directly access a single physical GPU using NVIDIA vGPU. You can create virtual GPUs that can be shared by VMs across the enterprise and accessed by other devices.

This capability combines the power of GPU performance with the management and security benefits provided by vGPU. Additional benefits provided by vGPU includes proactive management and monitoring for your VM environment, workload balancing for mixed VDI and compute workloads, and resource sharing across multiple VMs.

#### Additional resources

- [Virtual GPUs](#)

## 2.4. NVIDIA GPU FEATURES FOR OPENSIFT CONTAINER PLATFORM

### NVIDIA Container Toolkit

NVIDIA Container Toolkit enables you to create and run GPU-accelerated containers. The toolkit includes a container runtime library and utilities to automatically configure containers to use NVIDIA GPUs.

### NVIDIA AI Enterprise

NVIDIA AI Enterprise is an end-to-end, cloud-native suite of AI and data analytics software optimized, certified, and supported with NVIDIA-Certified systems.

NVIDIA AI Enterprise includes support for Red Hat OpenShift Container Platform. The following installation methods are supported:

- OpenShift Container Platform on bare metal or VMware vSphere with GPU Passthrough.
- OpenShift Container Platform on VMware vSphere with NVIDIA vGPU.

### GPU Feature Discovery

NVIDIA GPU Feature Discovery for Kubernetes is a software component that enables you to automatically generate labels for the GPUs available on a node. GPU Feature Discovery uses node feature discovery (NFD) to perform this labeling.

The Node Feature Discovery Operator (NFD) manages the discovery of hardware features and configurations in an OpenShift Container Platform cluster by labeling nodes with hardware-specific information. NFD labels the host with node-specific attributes, such as PCI cards, kernel, OS version, and so on.

You can find the NFD Operator in the Operator Hub by searching for "Node Feature Discovery".

### NVIDIA GPU Operator with OpenShift Virtualization

Up until this point, the GPU Operator only provisioned worker nodes to run GPU-accelerated containers. Now, the GPU Operator can also be used to provision worker nodes for running GPU-accelerated virtual machines (VMs).

You can configure the GPU Operator to deploy different software components to worker nodes depending on which GPU workload is configured to run on those nodes.

### GPU Monitoring dashboard

You can install a monitoring dashboard to display GPU usage information on the cluster **Observe** page in the OpenShift Container Platform web console. GPU utilization information includes the number of available GPUs, power consumption (in watts), temperature (in degrees Celsius), utilization (in percent), and other metrics for each GPU.

### Additional resources

- [NVIDIA-Certified Systems](#)
- [NVIDIA AI Enterprise](#)
- [NVIDIA Container Toolkit](#)
- [Enabling the GPU Monitoring Dashboard](#)
- [MIG Support in OpenShift Container Platform](#)
- [Time-slicing NVIDIA GPUs in OpenShift](#)
- [Deploy GPU Operators in a disconnected or airgapped environment](#)
- [Node Feature Discovery Operator](#)

## CHAPTER 3. AMD GPU OPERATOR

AMD Instinct GPU accelerators combined with the AMD GPU Operator within your OpenShift Container Platform cluster lets you seamlessly harness computing capabilities for machine learning, Generative AI, and GPU-accelerated applications.

This documentation provides the information you need to enable, configure, and test the AMD GPU Operator. For more information, see [AMD Instinct™ Accelerators](#).

### 3.1. ABOUT THE AMD GPU OPERATOR

The hardware acceleration capabilities of the AMD GPU Operator provide enhanced performance and cost efficiency for data scientists and developers using Red Hat OpenShift AI for creating artificial intelligence and machine learning (AI/ML) applications. Accelerating specific areas of GPU functions can minimize CPU processing and memory usage, improving overall application speed, memory consumption, and bandwidth restrictions.

### 3.2. INSTALLING THE AMD GPU OPERATOR

As a cluster administrator, you can install the AMD GPU Operator by using the OpenShift CLI and the web console. This is a multi-step procedure that requires the installation of the Node Feature Discovery Operator, the Kernel Module Management Operator, and then the AMD GPU Operator. Use the following steps in succession to install the AMD community release of the Operator.

#### Next steps

1. Install the [Node Feature Discovery Operator](#).
2. Install the [Kernel Module Management Operator](#).
3. Install and configure the [AMD GPU Operator](#).

### 3.3. TESTING THE AMD GPU OPERATOR

Use the following procedure to test the ROCmInfo installation and view the logs for the AMD MI210 GPU.

#### Procedure

1. Create a YAML file that tests ROCmInfo:

```
$ cat << EOF > rocminfo.yaml
apiVersion: v1
kind: Pod
metadata:
  name: rocminfo
spec:
  containers:
  - image: docker.io/rocm/pytorch:latest
    name: rocminfo
    command: ["/bin/sh", "-c"]
    args: ["rocminfo"]
  resources:
```

```
limits:
  amd.com/gpu: 1
requests:
  amd.com/gpu: 1
restartPolicy: Never
EOF
```

2. Create the **rocminfo** pod:

```
$ oc create -f rocminfo.yaml
```

### Example output

```
apiVersion: v1
pod/rocminfo created
```

3. Check the **rocminfo** log with one MI210 GPU:

```
$ oc logs rocminfo | grep -A5 "Agent"
```

### Example output

```
HSA Agents
=====
*****
Agent 1
*****
Name:          Intel(R) Xeon(R) Gold 6330 CPU @ 2.00GHz
Uuid:          CPU-XX
Marketing Name: Intel(R) Xeon(R) Gold 6330 CPU @ 2.00GHz
Vendor Name:   CPU
--
Agent 2
*****
Name:          Intel(R) Xeon(R) Gold 6330 CPU @ 2.00GHz
Uuid:          CPU-XX
Marketing Name: Intel(R) Xeon(R) Gold 6330 CPU @ 2.00GHz
Vendor Name:   CPU
--
Agent 3
*****
Name:          gfx90a
Uuid:          GPU-024b776f768a638b
Marketing Name: AMD Instinct MI210
Vendor Name:   AMD
```

4. Delete the pod:

```
$ oc delete -f rocminfo.yaml
```

### Example output

```
pod "rocminfo" deleted
```

