# RHOAI Admin Config

**Red Hat**

# Agenda

- Users and groups
- Custom notebook images
- Custom serving runtimes
- Accelerators and GPUs
- Managing RHOAI resources

**Red Hat**

# Authentication

# OpenShift AI – Authentication

▶ OpenShift AI uses the same Authentication mechanisms used by the underlying OpenShift Cluster

▶ See [this page](#) for more details on identity providers

- htpasswd
- Keystone
- LDAP
- Basic authentication
- Request header
- GitHub or GitHub Enterprise
- GitLab
- Google
- OpenID Connect

# Users and groups

Red Hat

# OpenShift AI – Group Membership

▶ Group membership is managed at the OpenShift Level:

# OpenShift AI – Users and Admins

▶ Only RHOAI admins can see the settings panel:



**Red Hat**
OpenShift AI

- Applications ›
- Data Science Projects
- Data Science Pipelines ⌄
  - Pipelines
  - Runs
- Model Serving
- Resources
- Settings ⌄
  - Notebook image settings
  - Cluster settings
  - Accelerator profiles
  - Serving runtimes
  - User management

Versio

- ▶ User Management:
  - · Choose the OpenShift Groups that should be RHOAI admins and RHOAI users
- ▶ "system:authenticated" means "any user declared in the identity provider". i.e. All users

▶ Example configuration

**Data Science administrator groups**

Select the OpenShift groups that contain all Data Science administrators.

Team1 ✖   rhods-admins ✖

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

ℹ **All cluster admins are automatically assigned as Data Science administrators.**

**Data Science user groups**

Select the OpenShift groups that contain all Data Science users.

Team1 ✖   Team2 ✖

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

# Creating Projects

▶ OpenShift users can only create projects if they have the Project Self-Provisioner role

▶ The OpenShift admin will have configured which users (if any) have this role.

▶ To disable/enable, consult [this documentation](#).

**Data Science Projects**

View your existing projects or create new projects.

Data science projects ▼

| Name ▼ | Q Find by name | Create data science project | Launch Jupyter |

| Name ↑ | Workbench | Status | Created ↕ |
| --- | --- | --- | --- |
| DS test ⌾<br>admin | – | – | 3/2/2024, 12:22:24 PM |

Red Hat

# Overview of user types and permissions

▶ **Data scientists** - Can access and use individual components of RHOAI, such as Jupyter.

▶ **Administrators** -  Have additional permissions to perform these actions:

  ■  Configure Red Hat OpenShift AI settings

  ■  Access and manage notebook servers

- Why customize your workbench?
- How can you customize a workbench?
  - Customizing the workbench
  - Creating custom notebook images (don't dig too deep)
  - Using custom notebook images
- Best practices around custom notebook images

# Custom notebook images

Red Hat

# Custom Workbench Images

- ▶ Support
- ▶ Why
- ▶ How

# Custom Workbench Images

▶ Support

- Red Hat only supports the images that are provided by default with RHOAI

- Red Hat also supports your ability to add a custom-built image

- Red Hat does not provide support for your custom-built images

- Any issue experienced with a Custom Image should be reproduced in one of the default images

# Custom Workbench Images

▶ Why

- Customization
    - create an image with the exact packages and versions you require
- Stability
    - control over the versioning of your custom images
    - (updates, overlap, phase out)
- Experimental
    - create bleeding edge versions of images
- Flexibility
    - Jupyter vs VSCode vs R-Studio, etc....

# Custom Workbench Images

▶ How

- Look at examples
- Create custom image (on laptop or on cluster) or use pre-built examples
- Upload/Store image into Container Image repository
- Add image reference into RHOAI admin interface
- Test

# Custom Workbench Images

▶ How

- There is an infinite ways of creating Workbench Images, therefore, the process cannot be fully documented.
- Instead, we refer to examples in order to get started:
  - https://ai-on-openshift.io/odh-rhoai/custom-notebooks/
  - https://github.com/opendatahub-io-contrib/workbench-images
  - https://quay.io/repository/opendatahub-contrib/workbench-images?tab=tags

# Serving runtimes
# Default and Custom

Red Hat

# Serving Runtimes

- A Serving Runtime can be either:

  - Default (Provided by Red Hat)

  - Custom (built and deployed by the customer

- For Single-Model serving or for Multi-Model serving

- Serving Runtimes can be managed only by the RHOAI Admins

- Serving Runtimes can be enabled/disabled

- Changes to the runtime definitions do not affect existing Model Servers

# Serving Runtimes

# Models and model servers

Select the type of model serving platform to be used when deploying models in this project.



## Single-model serving platform

Each model is deployed on its own model server. This platform works well for large models or models that need dedicated resources.

**Deploy model**



## Multi-model serving platform

Multiple models can be deployed on a single-model server. This platform works well for sharing resources amongst deployed models

**Add model server**

## Serving runtime *

Select one

Caikit TGIS ServingRuntime for KServe

OpenVINO Model Server

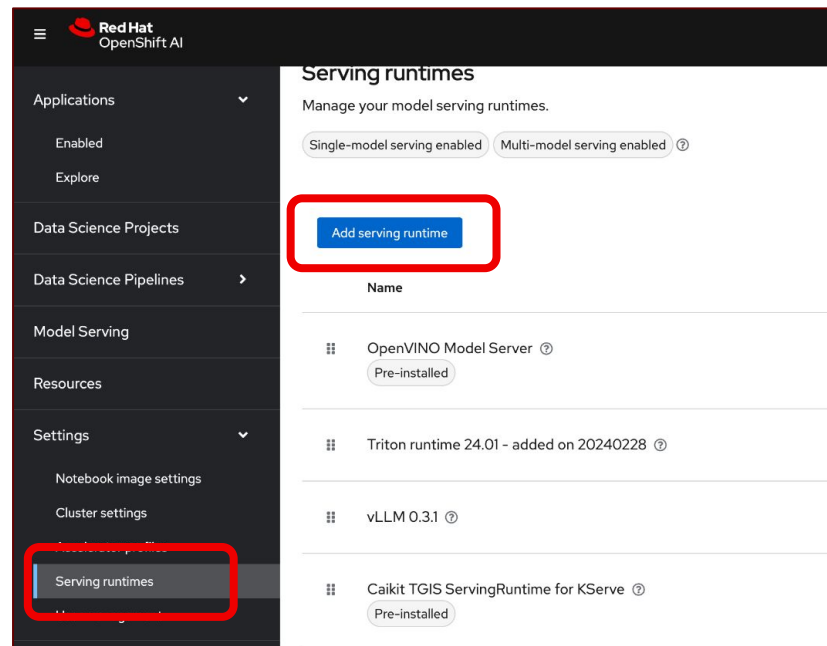TGIS Standalone ServingRuntime for KServe (gRPC)

## Serving runtime *

Select one

OpenVINO Model Server

**Red Hat**

# Adding a Custom Serving Runtimes

▶ Locate an example

- https://redhatquickcourses.github.io/rhods-deploy/rhods-deploy/1.33/chapter1/section3.html

- or

- Custom Serving Runtime (Triton) - AI on OpenShift

▶ Navigate to the right location:

# Serving Runtimes

Settings ❯ Serving runtimes ❯ Add serving runtime

## Add serving runtime

Add a new runtime that will be available for users on this cluster.

**Select the model serving platforms this runtime supports** *

| Select a value ▼ |

Single-model serving platform

Multi-model serving platform

Single-model and multi-model serving platforms

# Serving Runtimes

Applications ⌄

    Enabled

    Explore

Data Science Projects

Data Science Pipelines ❯

Model Serving

Resources

Settings ⌄

    Notebook image settings

    Cluster settings

    Accelerator profiles

    Serving runtimes

    User management

**Select the model serving platforms this runtime supports** *

Single-model and multi-model serving platforms ▾

```
 1  apiVersion: serving.kserve.io/v1alpha1
 2  kind: ServingRuntime
 3  metadata:
 4    name: triton-23.05-20230804
 5    labels:
 6      name: triton-23.05-20230804
 7    annotations:
 8      maxLoadingConcurrency: "2"
 9      openshift.io/display-name: "Triton runtime 23.05"
10  spec:
11    supportedModelFormats:
12      - name: keras
13        version: "2"
14        autoSelect: true
15      - name: onnx
16        version: "1"
17        autoSelect: true
18      - name: pytorch
19        version: "1"
20        autoSelect: true
21      - name: tensorflow
22        version: "1"
23        autoSelect: true
```

Create    Cancel

Hat

# Accelerators and GPUs

Red Hat

# Accelerator profiles

▶ Why

- There are various types of accelerators

- e.g. Nvidia GPUs, Intel Gaudi, etc...

- We need a way to surface these resources to the users

- We need to differentiate, for example, different types of GPUs, and land the pods on the right machines

▶ How

- Accelerator Profiles are managed by RHOAI admins.

# Accelerators and GPUs

▶ Accelerator Profiles

# Accelerators and GPUs

▶ Creating a new profile

## Create accelerator profile

Jump to section

**Details**

| Details |
| --- |
| **Tolerations** |

### Details

**Name** *

Habana HPU – 1st Gen Gaudi

**Identifier** * ⑦

habana.ai/gaudi

**Description**

This Accelerator Profile is for 1st Gen Gaudi Devices

**Enabled**

[toggle on]

### Tolerations    [Add toleration]

| Operator | Key | Value | Effect | Toleration Seconds | |
| --- | --- | --- | --- | --- | --- |
| Equal | habana.ai/gaudi | present | NoSchedule | – | ⋮ |

[Create accelerator profile]    Cancel

# Accelerators and GPUs

▶ Equivalent YAML

Project: redhat-ods-applications ▼

AcceleratorProfiles > AcceleratorProfile details

**AP** habana-hpu---1st-gen-gaudi

Details **YAML**

```
1   apiVersion: dashboard.opendatahub.io/v1
2   kind: AcceleratorProfile
3 > metadata: ···
30  spec:
31    description: This Accelerator Profile is for 1st Gen Gaudi Devices
32    displayName: Habana HPU – 1st Gen Gaudi
33    enabled: true
34    identifier: habana.ai/gaudi
35    tolerations:
36      - effect: NoSchedule
37        key: habana.ai/gaudi
38        operator: Equal
39        value: present
40
```

Save    Reload    Cancel

# Accelerators and GPUs

## Notebook image

**Image selection** *

| HabanaAI | ▼ |

❓ View package information

## Deployment size

**Container size**

| Tiny | ▼ |

**Accelerator**

| None | ▼ |

None

Habana HPU – 1st Gen Gaudi

This Accelerator Profile is for 1st Gen Gaudi Devices

Compatible with image

NVIDIA GPU

# Accelerators and GPUs
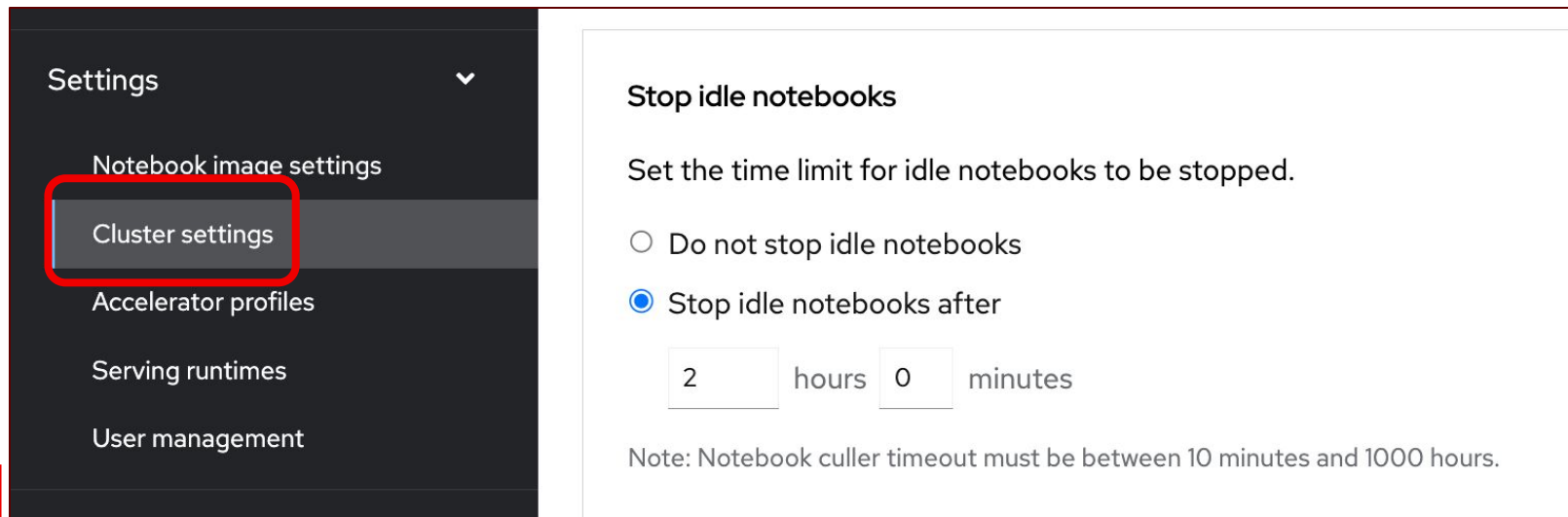
# Managing RHOAI resources

# Idle Culler and other admin Settings

Red Hat

# Idle Culler

▶ What

- The idle culler is a mechanism that will stop Notebook servers left "idle"
- "idle" in this context means that the notebook pod is running, but no browser is connected to it

▶ Why

- Notebooks/Workbenches are supposed to be interactive environments
- But if too many are left running all the time, they can consume all the cluster's resources
- Idle Culler ensures regular cleanup of leftover workbenches

# Idle Culler

▶ How

- Default value is "do not stop..."

# Toleration for notebooks

▶ Why

- If you want to dedicate some nodes to notebooks
- Very useful for auto-scaling

**Notebook pod tolerations**

☑ Add a toleration to notebook pods to allow them to be scheduled to tainted nodes

Toleration key for notebook pods: notebooksonly

The toleration key above will be applied to all notebook pods when they are created. Add a matching taint key (with any value) to the Machine Pool(s) that you want to dedicate to Notebooks.

Red Hat

# Managing Container Sizes

Red Hat

# Managing Container Sizes

- ▶ What
  - · When launching Workbenches and Model Servers, users are prompted for a Container Size (S/M/L/XL)
- ▶ RHOAI admins can control the sizes available to the users

# Managing Container Sizes

Project: redhat-ods-applications ▾

OdhDashboardConfigs ❯ OdhDashboardConfig details

**ODC** **odh-dashboard-config**

Details    **YAML**

```
106        key: notebooksonly
107      pvcSize: 20Gi
108    notebookSizes:
109    - name: Tiny
110      resources:
111        limits:
112          cpu: '1'
113          memory: 1Gi
114        requests:
115          cpu: 500m
116          memory: 1Gi
117    - name: Small
118      resources:
119        limits:
120          cpu: '2'
121          memory: 2Gi
122        requests:
123          cpu: '1'
124          memory: 2Gi
```

- Notebook sizes and Model Servers sizes are managed independently
- Model Servers Sizes can also be "custom":

**Compute resources per replica**

Model server size ⃝?

| Custom | ▾ |

**CPUs requested**

| − | 1 | + | | Cores ▾ |

**Memory requested**

| − | 4 | + | | Gi ▾ |

**CPU limit**

| − | 2 | + | | Cores ▾ |

**Memory limit**

| − | 8 | + | | Gi ▾ |

🎩 **Red Hat**

# Managing Container Sizes

▶ Note:

- When changing a size, existing Workbenches and Model Servers remain unchanged

- To avoid "unknown" sizes, leave a (deprecated) placeholder in place for a little while

# end of section

linkedin.com/company/red-hat

youtube.com/user/RedHatVideos

facebook.com/redhatinc

twitter.com/RedHat

**Red Hat**