



HPE Reference Architecture for HPE Ezmeral ML Ops on Kubernetes

Providing DevOps Speed and Agility for Machine Learning

CONTENTS

Executive summary.....	4
Solution overview.....	4
Solution components.....	6
HPE Ezmeral Runtime.....	6
Kubernetes features on HPE Ezmeral ML Ops.....	8
HPE Ezmeral ML Ops.....	9
Best practices.....	11
HPE EPA Platform configuration for HPE Ezmeral Runtime.....	11
HPE and NVIDIA GPUs.....	12
Storage.....	13
Networking.....	14
Kubeflow with HPE Ezmeral ML Ops.....	15
HPE Ezmeral ML Ops use cases.....	17
Software components.....	17
Hardware components.....	18
Use case NYC taxi rides.....	18
Use case on Pima Indian's diabetes prediction.....	34
Ezmeral ML Ops - in action with use case (Spark operator).....	43
Spark operator with K8s.....	43
Spark operator use case.....	45
Model training from KubeDirector Notebook using Spark with Livy.....	47
Ezmeral ML Ops – Experiment tracking with MLflow.....	54
Use case workflow.....	54
Monitoring.....	66
Kubernetes administrator.....	67
Kubernetes tenant/project administrator.....	68
Istio and Prometheus.....	68
Summary.....	70
Appendix A: Kubeflow and tests of use cases.....	71
Kubeflow components use cases GitHub issue summarization - Training with Jupyter.....	73
GitHub issue summarization – Serving with Seldon.....	77
Training with TensorFlow (Financial series).....	78
Serving a TensorFlow model with KFServing (Financial series).....	81
Training a PyTorch model (PyTorch MNIST).....	83
Sample Pipeline in the pipelines interface.....	84
Running a pipeline in Jupyter Notebook.....	87
Katib Hyperparameter Tuning.....	90
Argo workflows.....	96



Reference Architecture

ML metadata	97
Appendix B: HPE ML Ops KDapp.....	99
Centos/Ubuntu	99
MLflow	100
NVIDIA: TensorFlow (NGC).....	106
TensorFlow + Jupyter	107
Appendix C: Install and configure HPE Ezmeral Runtime.....	109
Resources and additional links	110



EXECUTIVE SUMMARY

Enterprises across all industries are embarking on a hybrid cloud journey for the development and deployment of their data-driven analytics and AI/ML applications. The continuous integration/continuous deployment (CI/CD) workflows, collectively referred to as DevOps, have become ubiquitous for all software development today. On the machine learning front, data scientists still spend a significant amount of time and effort moving projects from development to production. Model versioning is still manual, making it hard to update models in production. Code sharing is manual; data is copied onto local storage leading to the variability of results between environments. There is also a lack of standardization on the tools and frameworks used, which makes it tedious and time-consuming to deploy models across all environments.

HPE Ezmeral ML Ops includes the same capabilities and functionality as the HPE Ezmeral Runtime while also providing DevOps-like agility to enterprise machine learning. With HPE Ezmeral ML Ops, enterprises can implement CI/CD workflows and standardize their ML pipelines. The HPE Ezmeral ML Ops software platform supports every stage of the machine learning lifecycle — from supporting sandbox experimentation with the choice of ML/DL frameworks and integrating with model and code repositories to deploying and tracking models in production.

HPE Ezmeral ML Ops gives data scientists and developers the ability to quickly and easily build and train machine learning models. It allows data scientists to manage and track models built on any platform and deploy them into a scalable and secure production environment. Using HPE Ezmeral ML Ops, data scientists can spin up containerized environments for distributed data processing, Machine Learning (ML), or Deep Learning (DL) in minutes rather than weeks. It provides data science teams the flexibility to run their ML/DL workloads either on-premises, in multiple public clouds, or a hybrid model and respond to dynamic business requirements in a variety of use cases.

With HPE Ezmeral ML Ops, Hewlett Packard Enterprise is making it easier for organizations to deliver a flexible and secure multitenant architecture, with the agility, flexibility, and performance needed to address today's evolving workload and application requirements. Its deployment on HPE hardware can be done using pre-tested and optimized HPE Apollo building blocks on-premises, as well as in hybrid IT architectures and a multi-cloud model.

SOLUTION OVERVIEW

HPE Ezmeral Runtime is a unified container software platform that is built on open-source Kubernetes and designed for both cloud-native applications and non-cloud-native applications running on any infrastructure either on-premises, in multiple public clouds, in a hybrid model, or at the edge. With HPE Ezmeral Runtime, container deployment and operations can be simplified at scale. HPE Ezmeral Runtime best practices and automation can help streamline operations and improve SLAs. Hewlett Packard Enterprise delivers highly automated playbooks for Day 0 deployments combined with best practices and configuration automation to set up container HA, backup/restore, security validation, and monitoring to minimize manual overheads. HPE Ezmeral Container Platform includes KubeDirector—an open-source Kubernetes-based controller that can be used to deploy non-cloud-native apps. The HPE Ezmeral Runtime provides an App Store of curated, prebuilt images for a wide range of applications including machine learning (ML), analytics, IoT/edge, and CI/CD.

In this Reference Architecture, we discuss HPE Ezmeral ML Ops on Kubernetes and its components including the HPE Ezmeral Runtime. This is a workload-optimized platform to serve the needs of DevOps teams, CI/CD workflow integration, application modernization, and hybrid cloud solutions for the enterprise. This solution provides a cloud-like experience to customers from edge to core to cloud.

The HPE Apollo fit-for-purpose built server, storage, and networking hardware is the foundation for an infrastructure that provides both rapid deployment and scaling while delivering the highest levels of performance, quality, and availability. This solution also showcases how to modernize a legacy application using KubeDirector.

The combination of HPE Ezmeral Runtime (with pre-integrated HPE Ezmeral Data Fabric), and HPE Apollo servers deliver a composable architecture that rapidly deploys containers supporting the latest application frameworks. Ultimately, this results in faster digital transformation for the business. With services from HPE Pointnext and HPE GreenLake, the customer decides whether to purchase hardware upfront or move to a pay-as-you-go consumption model.



Figure 1 shows the architecture of the HPE Ezmeral Runtime.

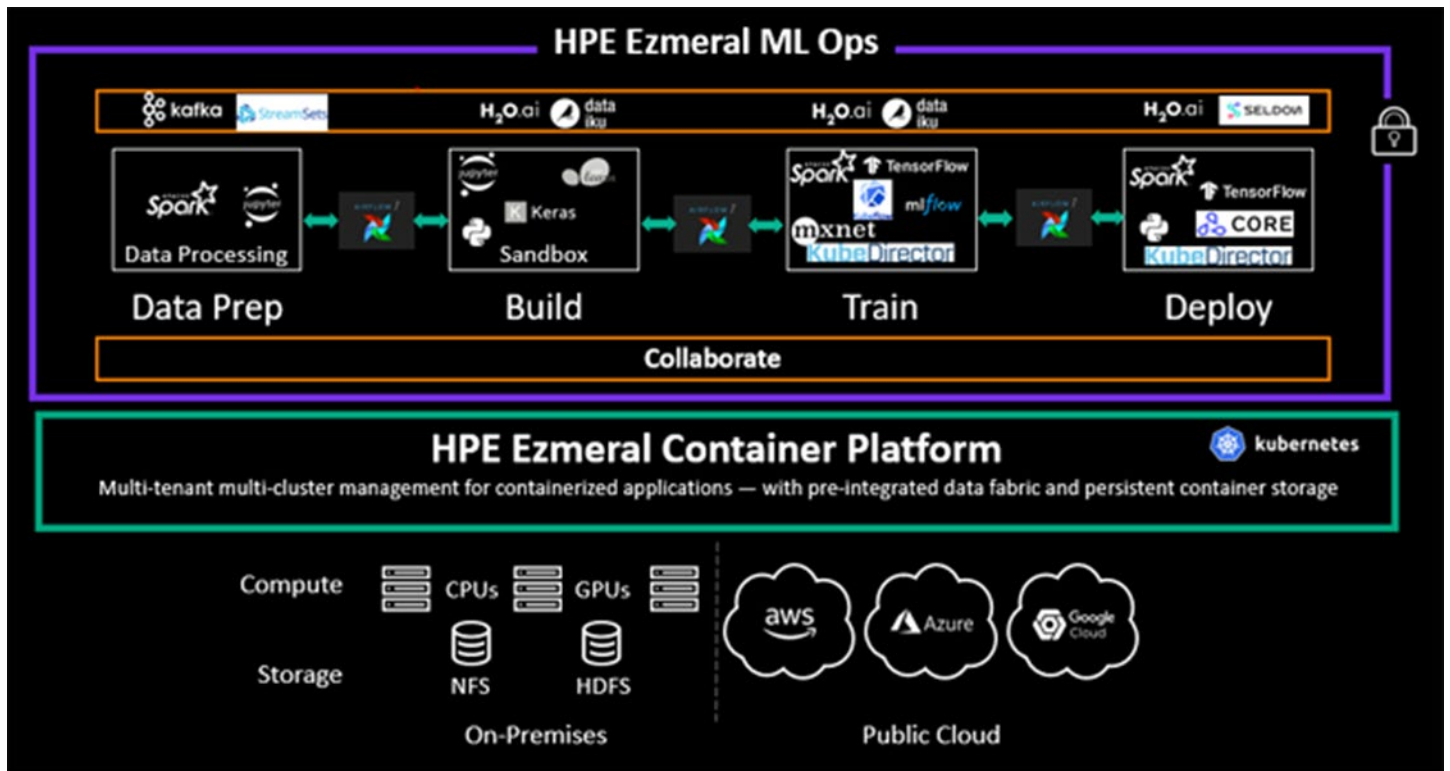


FIGURE 1. HPE Ezmeral ML Ops Architecture

HPE Ezmeral ML Ops brings the power of Kubernetes pods and Docker containers to the entire machine learning lifecycle to allow customers to build, train, deploy, and monitor machine learning (ML) and deep learning (DL) models. It supports sandbox development (notebooks), distributed training, and the deployment and monitoring of trained models in production. Project repository, source control, and model registry features allow seamless collaboration.

Some of the specific features supplied at each stage of the machine learning lifecycle include:

Table 1 shows the steps of the ML Lifecycle.

TABLE 1. ML Lifecycle

Steps	Description	Users
Data Preparation	<ul style="list-style-type: none"> Select data Preprocess the data by formatting, cleaning, and sampling it Transform data by scaling, decomposing, and aggregating it Use for fast data ingest and movement 	Data Analyst
Model Building	<ul style="list-style-type: none"> Containerized sandbox environments Choice of ML/DL tools, interface, and frameworks Secure access to shared data Data Scientists can now quickly spin up environments using their preferred data science tools 	Data Scientist, App Developers
Model Training	<ul style="list-style-type: none"> Containerized, distributed ML/DL environment Auto-scaling capabilities Prepackaged images for Python, Spark, and TensorFlow 	Data Scientist, Data Engineer, DevOps Engineer
Model Deployment	<ul style="list-style-type: none"> Support multiple runtime engines for handling scoring logic (e.g. Python, R, etc.) Can deploy distributed ML/DL environments such as TensorFlow, Caffe2, H2O, BigDL, and SparkMLlib REST endpoints with token-based authorization 	Data Scientist, DevOps Engineer



Steps	Description	Users
	<ul style="list-style-type: none"> Autoscaling and load balancing Integration with the model registry allows data scientists to track model versions and seamlessly update models 	

The following are some additional features of ML Lifecycle in HPE Ezmeral ML Ops:

- **Model Monitoring**

- Track, measure, and report model performance
- Save and inspect inputs and outputs for each scoring request
- Third-party integrations track accuracy and interpretability

- **Collaboration**

- CI/CD workflows with code, model, and project repositories
- Integration with GitHub and Bitbucket for project/code repository
- Storing multiple models (multiple versions with metadata) for various runtime engines
- A/B or Canary testing to validate the model
- NFS based project repository that eases collaboration

- **Security and control**

- Secure multi-tenancy with integration to enterprise authentication mechanisms
- Multitenancy and data isolation to ensure logical separation between each project, group, or department within the organization
- Enterprise security and authentication mechanisms such as LDAP, Active Directory, and Kerberos
- Share the same infrastructure and access the same data sources for AI/ML and Big Data analytics workloads

- **Hybrid deployment**

- On-premises, public cloud, or hybrid
- Run on-premises on any infrastructure (including in multiple data centers)
- Supports multiple public clouds (Amazon® Web Services, Google® Cloud Platform, or Microsoft® Azure)
- Provides a hybrid model for effective utilization of resources and lower operating costs

In addition to the new features and business benefits delivered through the HPE Ezmeral ML Ops Software, the underlying functionality delivered previously via the HPE Ezmeral Container Platform will continue to be part of the overall HPE Ezmeral Runtime framework and integrated with the new HPE Ezmeral Machine Learning Ops Software.

This Reference Architecture describes our solution testing performed in February 2021.

Document purpose: This Reference Architecture provides an overview of the deployment of HPE Ezmeral Runtime on servers – HPE Apollo 6500, HPE Apollo 2000. Also, it provides deployment steps of Kubeflow for pipeline management and HPE Ezmeral Data Fabric (formerly 'MapR') based Spark Operator.

SOLUTION COMPONENTS

HPE Ezmeral Runtime

HPE Ezmeral Runtime installs as a software layer between the underlying server infrastructure and the Big Data distribution, AI/ML libraries, and applications. The use of containers is completely transparent, and HPE Ezmeral Runtime customers benefit from greater agility and bare-metal performance due to the lightweight nature of containers. They can leverage the flexibility of containers to simplify the development of DevOps, CI/CD pipelines, and applications across hybrid cloud deployments.



Key features

- **Multi-cluster Kubernetes management:** Fast, easy deployment, management and monitoring of multiple clusters with either out-of-the-box or default configuration for networking, load balancing, and storage. This permits the user to run and manage different versions of Kubernetes simultaneously, and seamlessly supports the in-place upgrades.
- **Enterprise-ready persistent container storage:** Fully managed, integrated, scale-out, and edge-ready persistent storage with the HPE Ezmeral Data Fabric. This Data Fabric, along with DataTap and FS Mount functionality provides connectivity to data without copying the data locally.
- **100% open-source Cloud Native Computing Foundation (CNCF) Kubernetes:** With innovations such as KubeDirector—an open-source Kubernetes-based controller to deploy non-cloud-native, stateful apps. HPE Ezmeral Runtime is a CNCF-certified Kubernetes distribution.
- **One-click provisioning:** Pre-packaged App Store with curated, prebuilt images for a wide range of applications including machine learning (ML), analytics, IoT/edge, CI/CD, and other modern apps. The pre-bundled contents of the apps include Helm Charts, Operators, YAML configuration files, and KubeDirector scripts.
- **Simplified installation and upgrade workflows:** This includes installation on bare metal, Virtual Machines, and cloud instances.
- **Flexible multi-cluster, multitenant control plane:** Deploy multiple open-source K8s clusters and manage cloud K8s clusters (Example: GKE, EKS) from an HPE Ezmeral Runtime control plane, without vendor lock-in or modification to native K8s.
- **KubeDirector:** The first and only K8s custom controller that deploys non-cloud native, monolithic distributed stateful applications (Example: CDH, HDP, Confluent, Bring your own app).
- **Streamlined access to K8s clusters and services for end-users:** Gateway hosts isolate the HPE Ezmeral Runtime control plane and K8s hosts from the user network. This uniquely provides load balancing to multi-master K8s cluster(s) and routes to K8s services exposed via Node Ports and Ingress Controllers.
- **Bare-metal performance:** HPE Ezmeral Runtime provides storage I/O optimizations to deliver data to applications without the penalties commonly associated with virtualization or containerization. The compute cores and RAM in each host are pooled and then partitioned into virtual resource groups based on tenant requirements.
- **Self-Service Environments:** Users can get up and run quickly with HPE Ezmeral Runtime Elastic Plane functionality. New containerized environments are provisioned on-demand with just a few mouse clicks—whether they're transient for development and testing, or long-running for a production workload. Data scientists and analysts can now quickly respond to dynamic business requirements for a variety of use cases ranging from deep learning with AI Frameworks like TensorFlow to analytical SQL workloads running on Hadoop. Flexibility for Tools of Choice: The HPE Ezmeral Runtime offers pre-integrated container images, including many of the most common AI and Big Data tools, ready-to-run versions of major Hadoop distributions, such as Cloudera (CDH), Hortonworks (HDP), and MapR (CDP). Also includes recent versions of Spark standalone as well as Kafka and Cassandra.
- **Compute and Storage Separation:** HPE Ezmeral Runtime disconnects analytical processing from data storage, giving users the ability to independently scale compute and storage based on the needs of the workloads.
- **Data Access from Any Storage:** With HPE Ezmeral Runtime DataTap capability, users can access data from any shared storage system (including HDFS as well as NFS) or cloud storage (e.g. Amazon S3). It is unnecessary need to make multiple copies of data or move data before running an analysis. Sensitive data can stay in a secure storage system with enterprise-grade data governance without the costs and risks of creating and maintaining multiple copies or moving large-scale data.

HPE Ezmeral Runtime goes beyond Hadoop and Spark support by leveraging the inherent infrastructure portability and flexibility of containers to support distributed AI for both ML and DL use cases. The separation of compute and storage for Big Data and ML/DL workloads is one of the key concepts behind this flexibility. Organizations can deploy multiple containerized compute clusters for different workflows (e.g. Spark, Kafka, or TensorFlow) while sharing access to a common data lake. This also enables hybrid and multi-cloud HPE Ezmeral Runtime deployments, with the ability to mix and match on- and/or off-premises compute and storage resources to suit each workload. Furthermore, compute resources can be quickly and easily scaled and optimized independently of data storage, thereby increasing flexibility and improving resource utilization while eliminating data duplication and reducing cost.



Kubernetes features on HPE Ezmeral ML Ops

Data engineers, ML architects, and others can spin up containerized Kubernetes environments on scalable compute clusters with their choice of machine learning tools and frameworks for Big Data, AI, and/or ML use cases. Some of the key features of Kubernetes on HPE Ezmeral ML Ops include:

- **Software installation:** either on physical or virtual hosts located in a hybrid environment.
- **Storage medium:** a pre-integrated persistent container storage system known as the HPE Ezmeral Data Fabric.
- **DataTaps and FS Mounts:** access to existing data sources, with no need to copy data back and forth. See [DataTaps](#) and [FS Mounts](#).
- **Multitenant, multi-cluster management:** use open-source Kubernetes orchestration to run a variety of databases, analytics, AI/ML, CI/CD pipeline, and other applications.
- **Big Data Kubernetes tenants:** HPE Ezmeral Runtime can deploy applications with KubeDirector, or onboard Kubectl deployed applications, from the built-in Kubernetes Applications screen.
- **KubeDirector:** KubeDirector custom resource comes pre-installed with the HPE Ezmeral Runtime. The set of applications that can be automatically launched into a cluster is found by accessing a Kubernetes tenant and then clicking the Applications tab. See the Applications article at this [link](#).

Kubernetes architecture with HPE Ezmeral Runtime

This diagram depicts the physical Kubernetes cluster architecture within the HPE Ezmeral Runtime. For details on Kubernetes physical architecture, see the following [link](#).

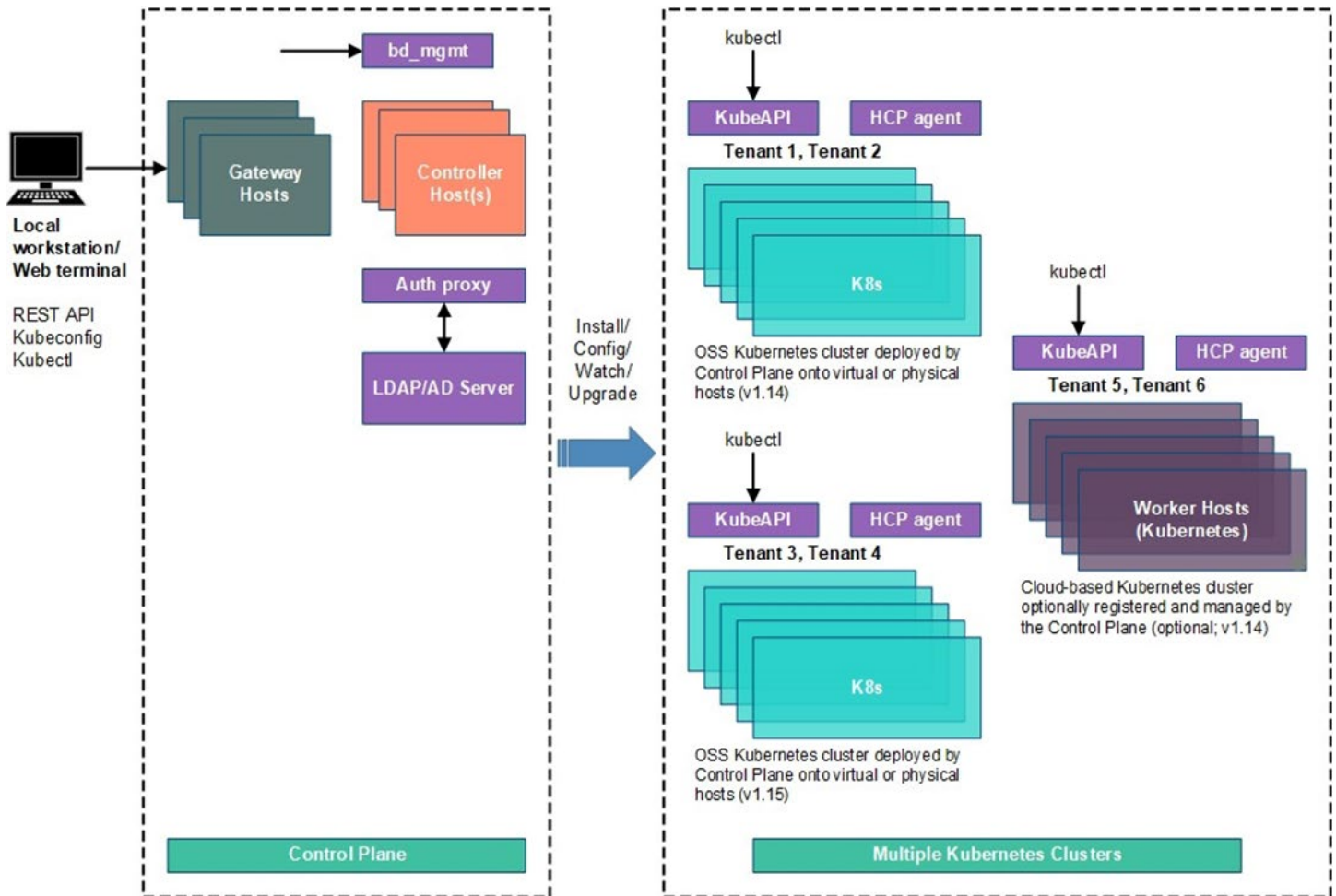


FIGURE 2. Kubernetes physical architecture in HPE Ezmeral Runtime



HPE Ezmeral ML Ops

With the HPE Ezmeral ML Ops solution, data science teams involved in the ML/DL model lifecycle can benefit from the industry's most comprehensive operationalization and lifecycle management solution for enterprise AI.

Figure 3 shows these features in the ML/DL lifecycle causal relationship. For further details, see [HPE Ezmeral Container Platform 5.3 Documentation](#).

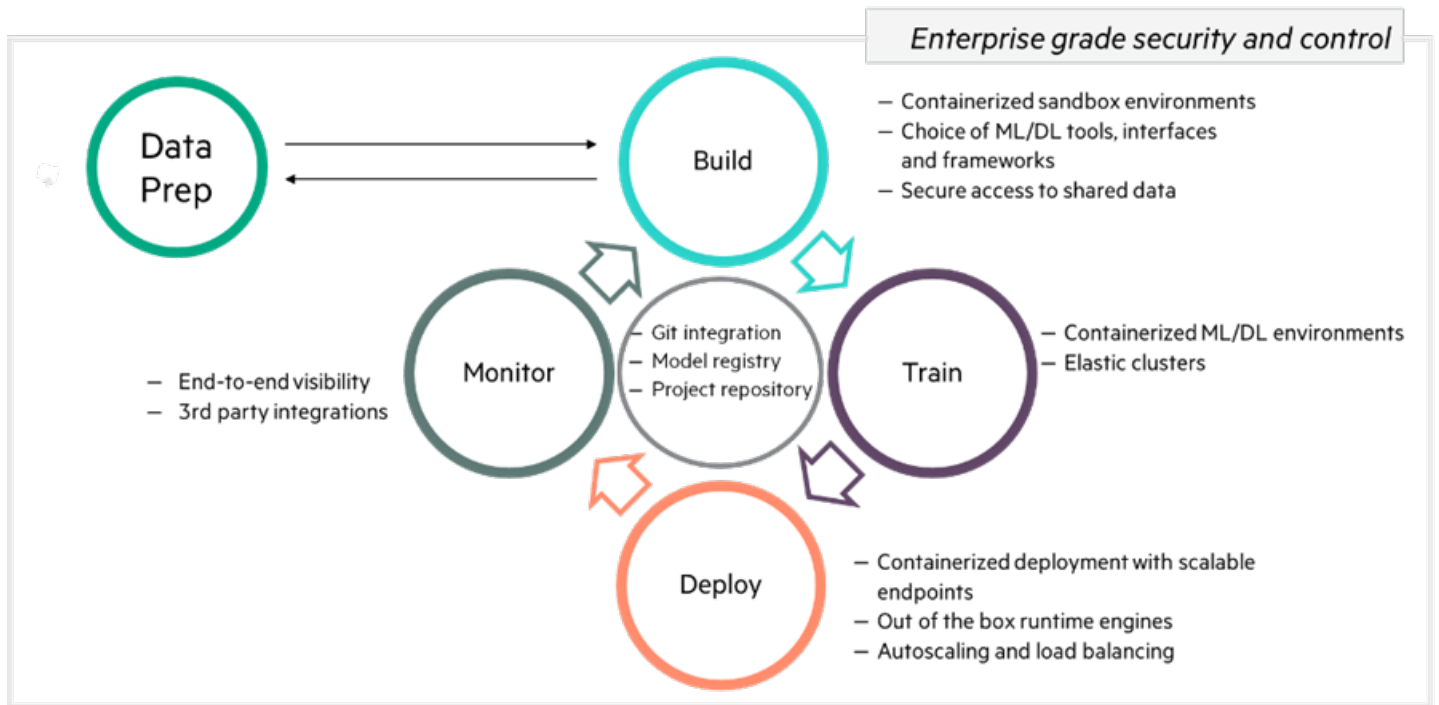


FIGURE 3. ML/DL lifecycle

HPE Ezmeral ML Ops delivers the following AI/ML features in addition to the HPE Ezmeral Runtime functionality:

- Leverage the power of containers to create a complex machine learning and deep learning stacks that include distributed TensorFlow, Apache Spark on Yarn with Kerberos, H2O, and Python ML and DL toolkits.
- Spin-up distributed, scalable, machine learning, and deep learning training environments on-premises, public cloud, or in a hybrid model.
- Support for a variety of programming languages and open-source tools designed to support even the most complex ML pipelines. For example, start with data pre-processing in Spark with Scala, followed by model development with TensorFlow on GPUs, and finally model deployment on CPUs with TensorFlow runtime.
- The model registry stores models and versions created within HPE Ezmeral ML Ops, as well as those created using different tools/platforms.
- Improves the reliability and reproducibility of machine learning projects in a shared project repository (GitHub).
- Enables the deployment of models in production with secure, scalable, highly available endpoint deployment with out-of-the-box auto-scaling, and load balancing.
- Enables out-of-the-box application images to be rapidly deployed in containerized environments – sandbox, distributed training, or serving (inference).
- Enables the creation of custom application images with any combination of tools, library packages, and frameworks.

For additional information, see [HPE Ezmeral ML Ops](#).



Kubeflow for Pipeline Management

Kubeflow is an open-source project designed to make machine learning workflows on Kubernetes simple, portable, and scalable. Kubeflow is sponsored by Google and inspired by TensorFlow Extended, or TFX, the company's internal machine learning platform. Originally intended to simply allow TensorFlow users to run training jobs on their Kubernetes clusters, the project now integrates a broad set of tools in support of many steps in an end-to-end machine learning process.

Kubeflow includes components for:

- Launching Jupyter Notebooks
- Building ML Pipelines
- Training Models
- Tracking Experiment Metadata
- Hyperparameter Tuning
- Serving Models
- Monitoring
- Continuous integration and deployment for ML

Spark within HPE Ezmeral Runtime

Spark is a data processing framework that can rapidly perform processing tasks on massive data sets and can also allocate data processing tasks across multiple nodes. Spark can work in stand-alone mode or together with other distributed computing tools. These features are the key to Big Data and ML/DL, which require significant computing power to crunch through large data stores. It also helps the developers by getting rid of some of the programming burdens with an easy-to-use API that abstracts much of the monotonous work of distributed computing and big data processing.

Spark has become one of the key big data distributed processing frameworks in the world. It can be deployed in a variety of ways and provides native bindings for the Java, Scala, Python, and R programming languages. In addition, it supports SQL, streaming data, machine learning, and graph processing.

Another Spark advantage is increased performance using its in-memory data engine. It can run tasks up to several orders of magnitude faster than MapReduce in certain situations. Furthermore, it offers the developer-friendly Spark API which hides much of the complexity of a distributed processing engine behind simple method calls.

The tutorial in the [Spark Operator with K8s](#) section describes how to set up and execute a Spark framework within the HPE Ezmeral Runtime which allows users to run Spark workloads on Kubernetes clusters. A user may instantiate the Spark framework instance on a dynamic cluster using the Spark operator. The Spark operator is a Kubernetes custom resource that is installed in a tenant namespace to support the on-demand deployment of Spark Executor pods. These pods are deleted once job execution completes.

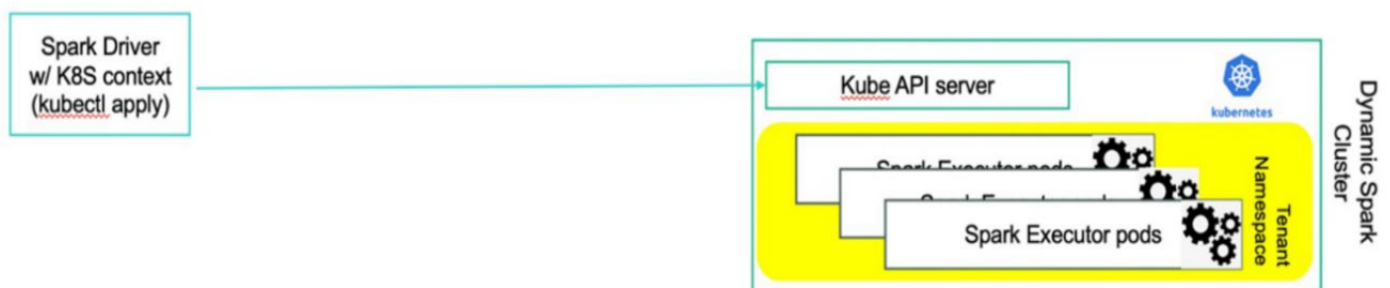


FIGURE 4. Spark Operator on Kubernetes



BEST PRACTICES

Deploying the HPE Ezmeral Runtime on the HPE Elastic Platform for Analytics (EPA) platform provides great flexibility in deploying workloads and managing resource growth, by decoupling storage from compute. This section is intended to provide high-level guidance and best practices for deploying HPE Ezmeral Runtime, and HPE Ezmeral ML Ops solution with the HPE Elastic Platform for Analytics.

HPE EPA Platform configuration for HPE Ezmeral Runtime

HPE Ezmeral Runtime uses four distinct host types, as shown in Table 2, with the recommended HPE EPA server model.

TABLE 2. HPE Ezmeral Runtime Host Types (Intel)

Host Type	HPE Server Model
Primary Controller, Shadow Controller, Arbiter	HPE Apollo 2000 Needs 3 x XL170r for controller nodes with HA
Kubernetes– Master/Worker	HPE Apollo 2000 (4 x HPE ProLiant XL170r) Recommend minimum 3 x HPE ProLiant XL170r for Kubernetes Worker/Compute nodes
Kubernetes – Compute with GPUs	<ul style="list-style-type: none"> • HPE Apollo 2000 with up to 2 x HPE ProLiant XL190r each with 1 or 2 GPUs • HPE Apollo 6500 with HPE ProLiant XL270d with 4 or 8 GPUs • HPE ProLiant DL380 with up to 4 GPUs
Kubernetes – Gateway	HPE Apollo 2000 Recommend 2 x HPE ProLiant XL170r

TABLE 3. HPE Container Platform Host Types (AMD)

Host Type	HPE Server Model
Primary Controller, Shadow Controller, Arbiter	HPE Apollo 2000 Gen10 Plus Needs 3 x HPE XL225n Gen10 Plus for controller nodes with HA
Kubernetes– Master/Worker	4 x HPE ProLiant DL325 Gen10 Plus Recommend minimum 3 x HPE ProLiant DL385 Gen10 Plus for Kubernetes Worker/Compute nodes
Kubernetes – Compute with GPUs	HPE Apollo 6500 (4 x HPE ProLiant DL675d Gen10 Plus)
Kubernetes – Gateway	HPE Apollo 2000 Gen10 Plus Recommend 2 x HPE XL225n Gen10 Plus

HPE Apollo 2000 Gen10 Plus Compute Servers

The HPE Apollo 2000 Gen10 Plus System is a shared infrastructure chassis with flexible support for up to four (4) HPE ProLiant XL225n Gen10 Plus servers (AMD) or up to four (4) HPE ProLiant XL220n Gen10 Plus servers (Intel®) or two (2) XL290n Gen10 Plus servers (Intel), helping increase rack space density. Server nodes can be serviced without impacting the operation of other nodes in the same chassis for increased server up-time. It delivers the flexibility to tailor the system to the precise needs of demanding high-performance computing (HPC) workloads with the right compute, flexible I/O, and storage options. The system can be deployed with a single server, leaving room to scale as customers' needs grow, bringing the power of supercomputing to data centers of any size. It is ideal for HPC applications in industry verticals like manufacturing, oil and gas, life sciences, and financial services.

- **HPE ProLiant XL170r Gen10 Server:** For compute-intensive workloads, HPE ProLiant XL170r delivers four servers in a single 2U chassis. Each HPE ProLiant XL170r server is serviced individually without impacting the operation of other servers sharing the same chassis to provide increased server uptime.
- **HPE ProLiant XL225n Gen10 Plus Server:** 1U Node Configure-to-order Server supports the full stack of 2nd generation AMD® EPYC™ 7000 Series processors.

For more information, see [HPE Apollo 2000 servers](#).

HPE Apollo 6500 Gen10 Plus GPU System

Built for the exascale era, the HPE Apollo 6500 Gen10 Plus System accelerates performance with NVIDIA® HGX A100 Tensor Core GPUs and AMD Instinct™ MI100 with Infinity Fabric™ accelerators to take on some of the most complex HPC and AI workloads. This purpose-built platform provides enhanced performance with premier graphics processing units (GPU), fast GPU interconnect, high-bandwidth fabric, and configurable



GPU topology, providing rock-solid reliability, availability, and serviceability (RAS). Configure with single or dual processor options for a better balance of processor cores, memory, and I/O. Improve system flexibility with support for 4, 8, 10, or 16 GPUs and a broad selection of operating systems and options, all within a customized design to reduce costs, improve reliability, and provide leading serviceability.

- **HPE ProLiant XL675d Gen10 Plus:** It is a dual-processor system for the NVIDIA HGX A100 8-GPU or AMD Instinct with 8 to 10 double-wide or 16 single-wide PCIe accelerators.

For more detailed information, see [HPE Apollo 6500 Gen10 Plus system](#).

HPE and NVIDIA GPUs

HPE ProLiant servers offer NVIDIA accelerators for high-performance computation for deep learning, high-performance computing (HPC) workloads, or graphics. The NVIDIA accelerators for HPE ProLiant servers seamlessly integrate GPU computing with select HPE server families. Designed for power-efficient, high-performance supercomputing, NVIDIA accelerators deliver dramatically higher application acceleration than a CPU-only approach for a range of deep learning, scientific, and commercial applications. The thousands of NVIDIA CUDA® cores of each accelerator allow it to divide large computing or graphics tasks into thousands of smaller tasks that can be run concurrently, thus enabling much faster simulations and improved graphics fidelity for extremely demanding 3D models.

For more detailed information, see [HPE AI and deep learning](#).

HPE Intelligent System Tuning (IST)

Available in HPE ProLiant Gen10 servers, HPE Intelligent System Tuning is a new set of revolutionary capabilities that deliver higher levels of performance, agility, and control to the server environment. With these groundbreaking new features, we can:

- Dynamically tune the servers' performance to match the needs of each workload
- Drive real cost savings
- Radically improve server performance

HPE ProLiant Gen10 Servers offer a UEFI configuration option to help customers tune their BIOS settings by using the known workload-based tuning profiles developed by the HPE performance engineering team. The default BIOS settings on HPE servers provide a balance between performance and power efficiency. For workloads running on the HPE ProLiant XL190r Gen10 and HPE Apollo 6500 servers with GPUs, the recommended workload profile is Graphic Processing which disables power management and virtualization to optimize the bandwidth between I/O and memory.

For more information about how to tune an HPE ProLiant Gen10 server using the workload profiles, refer to the [UEFI workload-based Performance Tuning Guide](#) for HPE ProLiant Gen10 servers.

TABLE 4. HPE Ezmeral Runtime host CPU and memory recommendations (Intel)

Host Type	Deployment Size	Memory	Processor
Ezmeral - Controller, Shadow Controller, Arbiter	Starter	192 GB	2 x Intel® Xeon® Gold 5215 - 10C 2.5 GHz
	Medium	384-768 GB	2 x Intel® Xeon® Gold 6242 - 16C 2.8 GHz
K8s - Master/Compute	Starter	192 GB	2 x Intel® Xeon® Gold 5215 - 10C 2.5 GHz
	Medium	384-768 GB	2 x Intel® Xeon® Gold 6226 - 12C 2.7 GHz
K8s - GPU	Starter	384 GB	2 x Intel® Xeon® Gold 5215 - 10C 2.5 GHz
	Medium	384-768 GB	2 x Intel® Xeon® Gold 6226 - 12C 2.7 GHz
	Large	384-768 GB	2 x Intel® Xeon® Gold 6242 - 16C 2.8 GHz
K8s - Gateway	All	192 GB	2 x Intel® Xeon® Gold 5215 - 10C 2.5 GHz

TABLE 5 HPE Ezmeral Runtime host CPU and memory recommendations (AMD)

Host Type	Deployment Size	Memory	Processor
Ezmeral - Controller, Shadow Controller, Arbiter	Starter	192 GB	AMD EPYC 7742 - 64C 2.25GHz
	Medium	384-768 GB	AMD EPYC 7742 - 64C 2.25GHz



Host Type	Deployment Size	Memory	Processor
K8s – Master/Compute	Starter	192 GB	AMD EPYC 7262 - 8C 3.2GHz
	Medium	384-768 GB	AMD EPYC 7262 - 8C 3.2GHz
K8s – GPU	Starter	384 GB	AMD EPYC 7542 - 32C 2.9GHz
	Medium	384-768 GB	AMD EPYC 7542 - 32C 2.9GHz
	Large	384-768 GB	AMD EPYC 7542 - 32C 2.9GHz
K8s – Gateway	All	192 GB	AMD EPYC 7742 - 64C 2.25GHz

To assist in sizing an HPE Ezmeral Runtime cluster, Hewlett Packard Enterprise has developed a [sizing tool](#).

NOTE

HPE Apollo Gen10 Plus systems support a variety of flexible memory configurations. But for optimal performance, it is recommended to balance the total memory capacity across all installed processors and make use of all six memory channels per CPU with up to two DIMM slots per channel.

Storage

A typical Kubernetes environment may have pods frequently coming and going. Large Kubernetes environments, such as a public cloud, may handle pools of systems where new hosts are added to support pod and cluster placement. In the HPE Ezmeral Runtime, a Data Fabric cluster is a Kubernetes Custom Resource that functions as a storage cluster providing access to PVCs, tenant storage, shares, and other storage needs. In a Data Fabric cluster:

- The hosts (called nodes) commit considerable disk resources that may include NVMe and enterprise-class SSDs
- The Data Fabric cluster can be deployed on a small number of nodes
- Unlike a typical k8s environment, pods are not deleted frequently
- The Data Fabric cluster must account for host resource profiles to guarantee core pod availability

HPE Ezmeral Runtime includes native support for HPE Ezmeral Data Fabric. This automates many manual steps and allows the creation of Data Fabric clusters like that used for creating Compute Kubernetes clusters (see [Creating a New Data Fabric Cluster](#) and [Creating a New Kubernetes Cluster](#)). Each Data Fabric cluster resides on nodes. See [Kubernetes Worker Installation Overview](#) and [Kubernetes Data Fabric Node Installation Overview](#).

Ephemeral Storage (Node Storage)

Ephemeral Storage is built from the local storage in each host and is used for the disk volumes that back the local storage for each virtual node. Using SEDs (Self-Encrypting Drives) will ensure that any data written to node storage is encrypted on write and decrypted on read by the OS. A tenant can optionally be assigned a quota for how much storage the nodes in that tenant can consume.



FIGURE 5. Storage

Virtual nodes/containers running on public cloud VMs (such as AWS EC2) utilize storage within the instance (such as AWS Elastic Block Storage, or EBS) as node storage.

Persistent Storage

Deploying a persistent data fabric is supported on the local disks within the hosts. This local storage can serve as either HDFS storage or as persistent volumes for Kubernetes clusters. Persistent volumes for Kubernetes stateful clusters are seamlessly available either from the native persistent data fabric or Nimble Storage using the storage interface driver (CSI) that is deployed during cluster creation.

Tenant Storage for Local Data Access

HPE Ezmeral Runtime can deploy a Data Fabric (MapR) file system on local disks for Tenant Storage within the servers running the K8s services. The DataTap interface then surfaces the physical locations of the Tenant Storage data blocks to the containers that make up the virtual cluster. This allows the Big Data task scheduling software running within the containers to route Big Data tasks to the containers running on the physical servers where copies of the required data blocks reside. This behavior mimics bare-metal Big Data deployments, thereby preserving the performance advantages of data locality without losing the flexibility and agility of a container-based virtualized compute platform. It also allows Big Data datasets to persist beyond the lifespan of a given Big Data cluster.

Operating system storage

For all host types, the recommended storage for the operating system is two 960 GB SSDs in a RAID 1 configuration.

HPE Ezmeral Runtime storage recommendations

Table 4 lists the recommended minimum storage configuration for each host type.

TABLE 5. Storage recommendations for HPE ProLiant XL170r, HPE ProLiant XL190r, and HPE ProLiant XL270d

Host Type	K8s – Storage Type	Storage Recommendation
Ezmeral – Controller, Shadow Controller, Arbiter	OS	2 x 960 GB SSD configured as RAID 1
	Ephemeral Storage	3 x 6.4 TB mixed-use SSD.
	Local Data Fabric	0-1 x 2 TB SATA 7.2 SFF HDD.
K8s – Master/ Compute	OS	2x 960 GB SSD configured as RAID 1
	Ephemeral Storage	3 x 6.4 TB mixed-use SSD.
	Local Data Fabric	1 x 2 TB SATA 7.2 SFF HDD.
K8s – Gateway	OS	2 x 960 GB SSD configured as RAID 1

NOTE

The HPE ProLiant XL170r and HPE ProLiant XL190r have six drive bays. Two are used for the OS drives and four are available for node storage and local tenant storage. For additional high availability, it is recommended that the node storage disks be configured as RAID 5. Choose the appropriate size and number of disks based on node storage and local tenant storage space requirements.

Networking

HPE Ezmeral Runtime operates on two networks, as shown in Figure 6 below.

The two networks are laid out as follows:

- Network for the Controller, Worker, and Gateway hosts: This network must be both routable and part of the organization's network that is managed by that organization's IT department.



- Network for virtual nodes (containers): HPE Ezmeral Runtime creates and manages this network, which can be either public (routable) or private (non-routable). For Kubernetes, Canal is used as the Container Network (CN) Network Provider. The container network is typically private (non-routable) instead of public (routable). See detailed description of [Public \(Routable\) Virtual Node Network](#) and [Private \(Non-Routable\) Virtual Node Network](#).

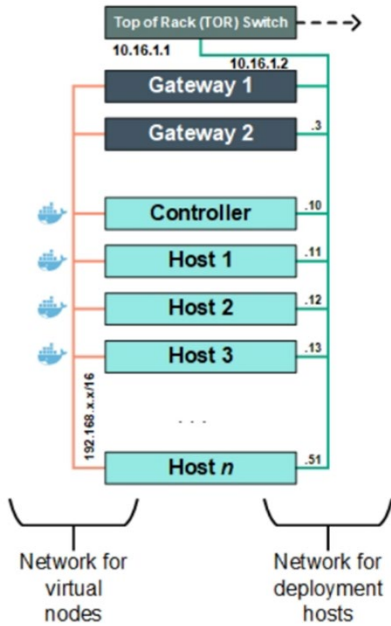


FIGURE 6. Networking layout

We recommend deploying 25GB Ethernet adapters on all the hosts. Table 6 shows the recommended networking hardware for the clusters.

TABLE 6. HPE Ezmeral Runtime networking recommendations

Host Type	Network Recommendation
All host types	HPE Eth 10/25GB 2P 640FLR-SFP28 Ethernet Adapter

KUBEFLOW WITH HPE EZMERAL ML OPS

Kubeflow conceptual overview

HPE Ezmeral Machine Learning Ops (HPE Ezmeral ML Ops) 5.3 upgraded to Kubeflow version 1.2.

Kubeflow is an open-source platform that makes the deployment of machine learning workflows in Kubernetes simple, portable, and scalable. Kubeflow deploys a suite of machine learning (ML) applications to Kubernetes that multiple users can securely access.

Kubeflow uses the Istio Gateway and LDAP authentication with the Dex authentication service to authenticate and authorize user access. Each user is assigned a profile based on a Kubernetes namespace. The profile provides an isolated view of Kubeflow for each user.

A Kubernetes administrator can install Kubeflow in environments where the computer network can access the internet, as well as in air-gapped environments where the network is isolated from outside networks.

At a high level, the execution of a pipeline proceeds as follows:

- **Python SDK:** Language use in the creation of components or pipelines using the Kubeflow Pipelines.
- **DSL Compiler:** Convert the pipeline's Python code into a static configuration (YAML).
- **Pipeline Service:** Call during the creation of a pipeline running from the static configuration (YAML).



- **Kubernetes Resources:** Allocate when Pipeline Service calls the Kubernetes API server to create the necessary resources (CRDs) to run the pipeline.
- **Orchestration Controllers:** A set of orchestration controllers execute the containers needed to complete the pipeline. An example controller is the [Argo Workflow](#) (demonstrated below) controller, which orchestrates task-driven workflows.
- **Artifact Storage:** The Pods store two kinds of data:
 - **Metadata:** Experiments, jobs, pipeline runs, and single scalar metrics. Metric data is aggregated to sort and filter. Kubeflow Pipelines stores the metadata in a MySQL database.
 - **Artifacts:** Pipeline packages, views, and large-scale metrics (time series). Use large-scale metrics to debug a pipeline run or investigate an individual run's performance. Kubeflow Pipelines stores the artifacts in an artifact store like a [Minio server](#) or [Cloud Storage](#).
- The MySQL database and the Minio server are both backed by the Kubernetes [Persistent Volume](#) subsystem.
- **Persistence Agent and ML Metadata:** The Pipeline Persistence Agent (PPA) watches the Kubernetes resources created by the Pipeline Service and persists the state of these resources in the ML Metadata Service. The PPA records the set of containers along with their inputs and outputs.
- **Pipeline Web Server:** The Pipeline web server gathers data from various services to display relevant views: the list of pipelines currently running, the history of pipeline execution, the list of data artifacts, debugging information about individual pipeline runs, execution status of the individual pipeline runs.

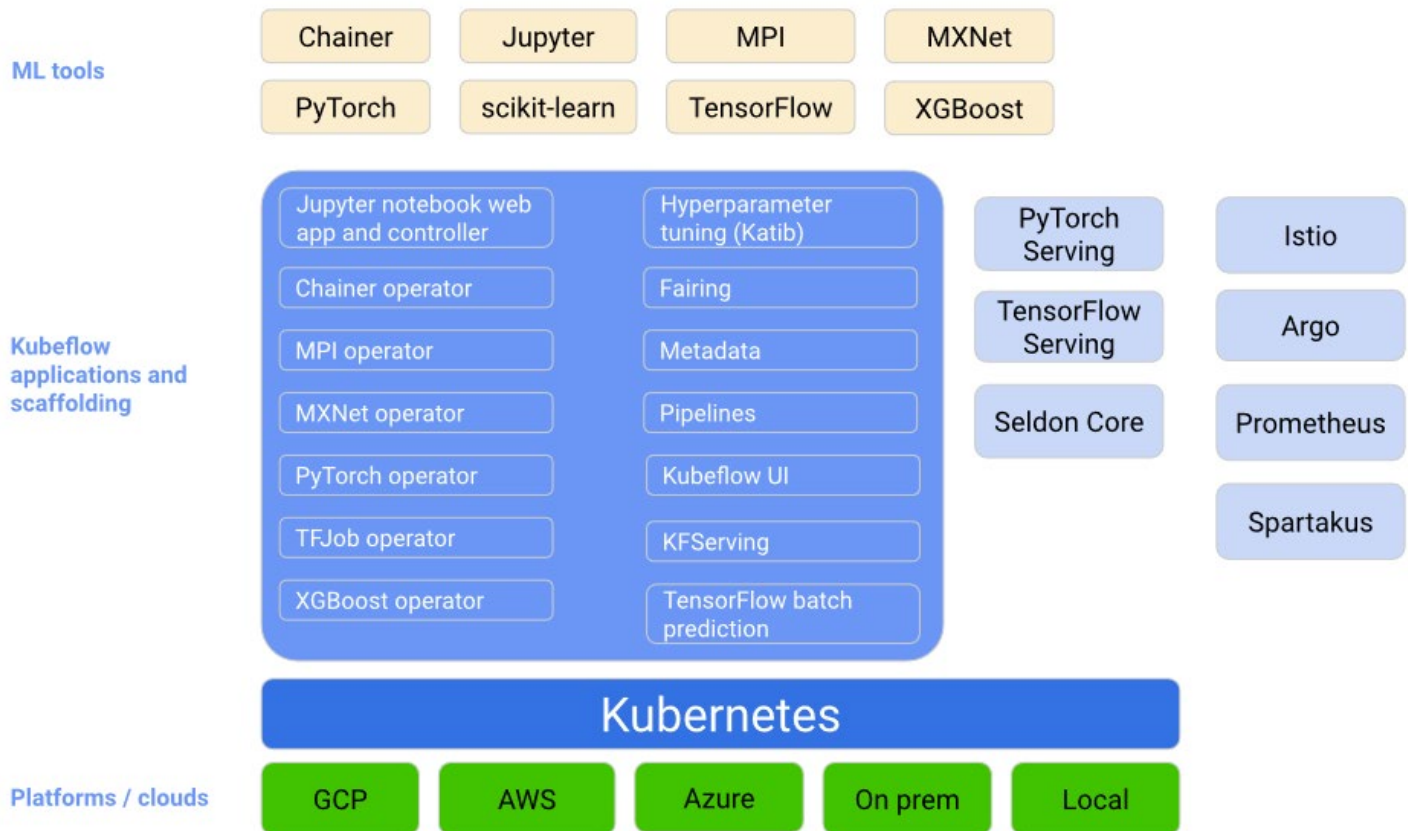


FIGURE 7. Kubeflow is a platform for components of ML systems on Kubernetes

Katib hyperparameter tuning

Katib is a Kubernetes-based system for hyperparameter tuning and neural architecture search. Katib supports many ML frameworks, including TensorFlow, MXNet, PyTorch, XGBoost, and others. It can be used to submit experiments and monitor results.



Argo Workflow

Argo Workflows is an open-source container-native workflow engine for orchestrating parallel jobs on Kubernetes. Argo Workflows is implemented as a Kubernetes CRD (Custom Resource Definition). Argo Workflows can run on any Kubernetes cluster within the HPE Ezmeral Runtime and orchestrate highly parallel jobs on Kubernetes.

Istio Prometheus

Prometheus is an open-source monitoring system and time series database. Prometheus can be used with Istio to record metrics that track the health of Istio and applications within the service mesh. Metrics can be visualized metrics using tools like Grafana and Kiali. The HPE Ezmeral Runtime supports the deployment of Prometheus

HPE EZMERAL ML OPS USE CASES

Much like pre-DevOps software development, data science organizations still spend a significant amount of time and effort moving projects from development to production. Model version control and code sharing are manual, and there is a lack of standardization on tools and frameworks thus making it tedious and time-consuming to productize machine learning models.

HPE Ezmeral Machine Learning Ops (HPE Ezmeral ML Ops) extends the capabilities of the platform and brings DevOps-like agility to enterprise machine learning. With the HPE Ezmeral ML Ops, enterprises can implement DevOps processes to standardize their ML workflows.

HPE Ezmeral ML Ops provides data science teams with a platform for their end-to-end data science needs with the flexibility to run their machine learning or deep learning (DL) workloads on-premises, in multiple public clouds, or in a hybrid model and respond to dynamic business requirements in a variety of use cases.

To complete the use case using the default datasets, 100GB of space will be required.

NOTE

The examples and documentation provided in this section are meant to supplement, not replace the HPE Ezmeral Runtime manuals.

Software components

This section describes the software versions utilized in the solution as well as notes any special installation or configuration requirements.

Table 7 lists the specific versions of the software used in this solution.

TABLE 7. Software Versions

Component	Versions
HPE Ezmeral Runtime	5.3.1
CentOS	CentOS Linux release 7.7.1908 Kernel 3.10.0-1062.el7.x86_64
389 Directory Server	1.15.1-37

Active Directory/LDAP: HPE Ezmeral ML Ops requires Active Directory (AD)/LDAP user authentication and supports the member of the attribute. The use cases use the following AD/LDAP users which can be substituted with existing LDAP users.

TABLE 8. AD/LDAP users

Projects	User	Member Of	Role
ML Ops/Kubeflow/Spark	Fraudadmin	Fraud	Admin
ML Ops/Kubeflow/Spark	fraudanalyst	Fraud	Member
ML Ops/Kubeflow/Spark	imageadmin	Image	Admin
ML Ops/Kubeflow/Spark	imageanalyst	Image	Member



Hardware components

The cluster was configured with AD/LDAP authentication, and platform high availability as shown in Figure 8.

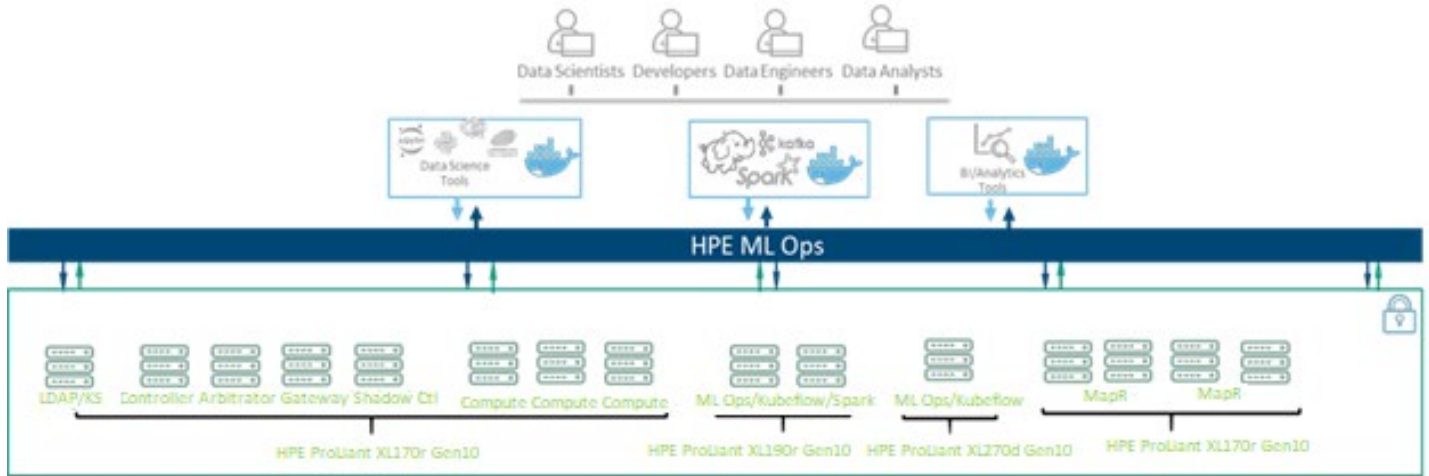


FIGURE 8. HPE Ezmeral ML Ops cluster

A complete list of the hardware components and service configuration is listed in Table 9. This is an example of how an HPE Ezmeral Runtime cluster could be deployed. The number of servers and types of servers will vary based on the workload.

TABLE 9. Hardware details

Qty	Host Service	Server Type	GPU	CPU Cores	Memory	OS Storage	Node Storage	Tenant Storage	K8s/EPIC Containers
HPE Ezmeral Runtime Cluster									
3	Ezmeral Controller	HPE ProLiant XL225n Gen10 Plus	0	8	128 GB	1.2 TB	1.2 TB	None	N
1	Ezmeral Gateway	HPE ProLiant XL225n Gen10 Plus	0	8	128 GB	1.2 TB	None	None	N
1	K8s Master/Compute	HPE ProLiant DL325 Gen10 Plus	0	32	256 GB	1.2 TB	1.2 TB	None	Y
2	K8s Compute GPU	HPE ProLiant XL675d Gen10 Plus	4	2 x 32	256 GB	1.2 TB	1.2 TB	None	Y
1	K8s Compute	HPE ProLiant DL385 Gen10 Plus	0	2 x 32	256 GB	1.2 TB	1.2 TB	None	Y
Supporting Services									
1	MIT Kerberos LDAP 389 Directory Server NTP Time Service DNS Name Server NFS Storage	HPE ProLiant DL360 Gen10		24	196 GB		1.2 TB		10 TB

Use case NYC taxi rides

The use case provided demonstrates how HPE ML Ops provides an end-to-end solution for the complete lifecycle to build, train, deploy, and monitor ML and DL models in multitenant enterprise environments. The use case demonstrates the NYC taxi ride prediction using TensorFlow.



The dataset contains a sample of approximately 375,000 NYC taxi rides from January-June 2019. Pickup and drop-off locations are specified as location ID numbers.

Sample data can be found at https://github.com/bluedatainc/solutions/tree/master/MLops/examples/NYCTaxi/Taxi_Datasets.

NOTE

The examples and documentation provided in this section are meant to supplement, not replace the manuals.

The AI/ML workflow allows the user to build, train, and deploy a model and then send API requests to that model to make predictions. This workflow consists of three high-level steps that must be performed by users with different roles in the following order:

- [Kubernetes Administrator](#)
- [LDAP/AD Administrator](#) (For Jupyter Notebook KDapp use)
- [Project Administrator](#)
- [Project Member](#) (Data Scientist)

Persona: Kubernetes Administrator

Verify that HPE Ezmeral Runtime is licensed for at least the number of CPU cores that will be used for the new Kubernetes cluster. HPE Ezmeral ML Ops requires a separate license for each of the CPU cores that will be used in AI/ML projects. Open the system settings screen, and then select the License tab to verify the number of CPU cores licensed for HPE Ezmeral ML Ops as shown in Figure 9.

License Summary

Name	Expiration Date	Status	Details
HPE Ezmeral Container Platform	LatestExpiration: Wed Sep 08 2021 NextExpiration: Wed Sep 08 2021	valid	UsedCapacity: 120 TotalCapacity: Unlimited
HPE Ezmeral Machine Learning Ops	LatestExpiration: Wed Sep 08 2021 NextExpiration: Wed Sep 08 2021	valid	UsedCapacity: 32 TotalCapacity: Unlimited
HPE Ezmeral External Cluster Management	LatestExpiration: Wed Sep 08 2021 NextExpiration: Wed Sep 08 2021	valid	UsedCapacity: 0 TotalCapacity: Unlimited

License(s)

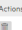
Name	Expiration Date	License Key	Details	Actions
HPE Container Platform and Machine Learning Ops Instant-on	Wed Sep 08 2021	WEE D9HA PPTA DEHQ UGDS HCTH FVIL K9PL B88H M21U 228A QHTL L7DZ LW4C GRFA LVVT DQKJ EPVW T59U N4GC 9KDC J88H L08B SPVY KX4Z WR4H HCFY TH3G DRVY XAVR PIVR BPTS XPTC LL4U R4V4 VB8B F0Z3 QW57 BPTT 4VLR HMY5 CBUD 2E47 HPE Container Platform and Machine Learning Ops Instant-on	Start: Wed Dec 31 1969 Capacity: unlimited Feature: HPE Ezmeral Container Platform and Machine Learning Ops Instant-on Evaluation: true DeviceID: any	

FIGURE 9. License Summary



Configure LDAP/AD authentication, If required.

User Authentication

Verify

Enable Multi Domain

Directory Server

Security Protocol

Service Locations

Bind Type

User Attribute

Group Attribute

Base DN

Bind DN (Optional)

Bind Password (Optional)

Verify Peer

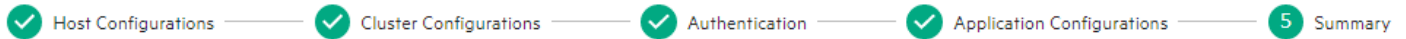
Enable SAML SSO

Submit

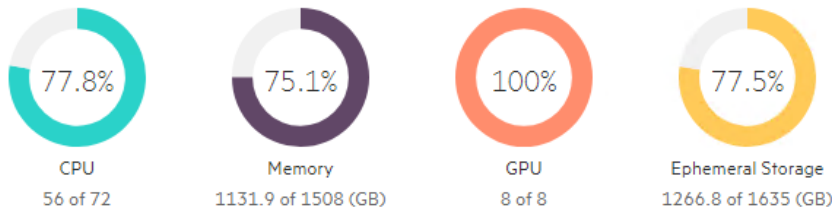
FIGURE 10. User Authentication



Create a Kubernetes cluster. Be sure to provide LDAP server information in Step 2: Authentication screen; LDAP must be configured to run HPE Ezmeral ML Ops in a Kubernetes cluster.



Cluster Size



Selected Hosts

Update

Host	Role	CPU	Memory (GB)	GPU	Ephemeral (GB)	Persistent (GB)	Disks
ezam-04.perflab.hp.com	master	16	377.5	0	372.6	745.2	Ephemeral Disks: /dev/sdb Persistent Disks: /dev/sdf
ezam-09.perflab.hp.com	worker	40	754.4	8	894.3	894.2	Ephemeral Disks: /dev/sdb Persistent Disks: /dev/sdc

Previous

Submit

FIGURE 11. Kubernetes cluster creation

Assign at least one user to be a Kubernetes Administrator for the Kubernetes cluster just created.

User Assignment



FIGURE 12. Admin Assignment to K8s cluster

Execute the following commands on the Kubernetes cluster master hosts to create LDAP/AD secret labels before creating tenants:

```
kubectl config set-context --current --namespace=hpecp
kubectl label secrets hpecp-ext-auth-secret "kubernetes.hpe.com/resource-tenant-visibility"="True"
```

```
[root@ezam-04 ~]# kubectl config set-context --current --namespace=hpecp
Context "kubernetes-admin@k8s-6" modified.
[root@ezam-04 ~]# kubectl label secrets hpecp-ext-auth-secret "kubernetes.hpe.com/resource-tenant-visibility"="True"
error: 'kubernetes.hpe.com/resource-tenant-visibility' already has a value (true), and --overwrite is false
[root@ezam-04 ~]#
```



■ Create a new Kubernetes AI/ML project, being sure to:

- a. Check the AI/ML Project checkbox.

Tenant Name

Tenant Description

K8s Cluster

Adopt Existing Namespace (Optional)

Specified Namespace Name (Optional)

Is Namespace Owner (Optional)

Map Services To Gateway (Optional)

ML Ops Project

Quotas External Authentication

Maximum Cores (License Available Capacity unlimited)

Maximum Memory (GB)

Maximum Ephemeral Storage (GB)

GPU Devices

Maximum Tenant Storage(GB)

Maximum Persistent Storage (GB)

Submit

FIGURE 13. Tenant Creation

- b. Enter the external LDAP/AD user group in the External Authentication tab.

Quotas External Authentication

External User Groups (Optional)

Member

Submit

FIGURE 14. External Authentication



Assign at least one user to be a Kubernetes Project Administrator for the newly created project.

User Assignment



FIGURE 15. Admin Assignment to the tenant

LDAP/AD Administrator (For Jupyter Notebook KDapp use)

If the environment will include the ability to use the Jupyter Notebook KubeDirector application (KDapp), LDAP server group settings must be changed to include all members of the group. For details, see [HPE Ezmeral Container Platform 5.3 Documentation](#).

Persona: Kubernetes project administrator

Assign at least one user to the new project.

User Assignment

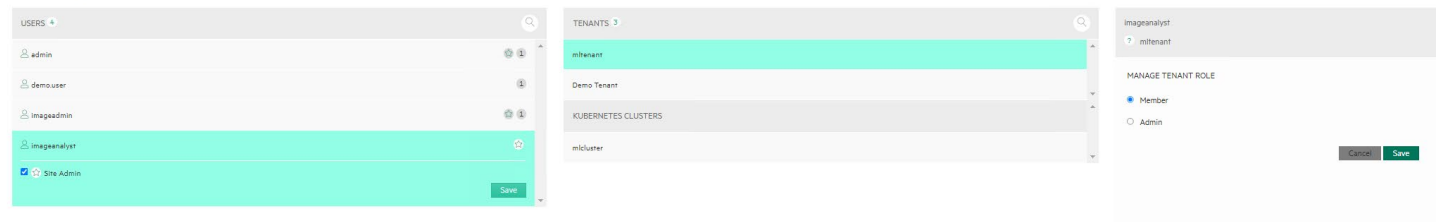


FIGURE 16. Member assignment to the tenant

NOTE

All AI/ML project users (Project Members and Project Administrators) must be LDAP/AD users. They cannot be authenticated using local authentication.



■ Configure one or more global source control configurations.

Create Source Control

Label

Configuration Name* ⓘ

Description ⓘ

Source Control Configuration

Source Control Type* ⓘ

Repository URL* ⓘ

Authentication Type* ⓘ

Branch ⓘ

Working Directory ⓘ

Proxy Protocol ⓘ

Proxy Host ⓘ

Proxy Port ⓘ

Username ⓘ

Email ⓘ

Token ⓘ

Submit

FIGURE 17. Global Source Control Configurations

Persona: Kubernetes project member (Data Scientist)

■ Log out of the web interface, and then log back in as the user that was created or assigned in step 1 of the Kubernetes Administrator workflow described above.



Configure at least one individual source control repository. Be sure to copy the name of the secret created for this source control.

Create Source Control

Label

Configuration Name* ⓘ indsrcctrl

Description ⓘ individual src ctrl

Global Configuration* ⓘ mltestglobalsrc ▼

Source Control Configuration

Source Control Type* ⓘ Github ▼

Repository URL* ⓘ https://github.com/hansha-sharma/solutions.git

Authentication Type* ⓘ Token ▼

Branch* ⓘ master

Working Directory ⓘ

Proxy Protocol ⓘ https

Proxy Host ⓘ hostname.com

Proxy Port ⓘ 00

Username* ⓘ hansha-sharma

Email* ⓘ hansha.sharma@hpe.com

Token* ⓘ

Submit

FIGURE 18. Individual Source Control Configuration



Access the Kubernetes Training screen, and then onboard the necessary training applications.

Create Training

Cluster Detail

Name*

Description

RunTime Image*

Enable DataTap

Node Roles

LoadBalancer

Instances

CPU

Memory (GB)

GPU

RESTServer

Instances

CPU

Memory (GB)

GPU

controller

Instances

CPU

Memory (GB)

GPU

FIGURE 19. Training Cluster Creation



Access the Kubernetes Notebooks screen, and then launch the notebook application. Also, attach the training cluster.

Create Notebook

Cluster Detail

Name*

Description

RunTime Image*

Enable DataTap

Node Roles

controller

Instances

CPU

Memory (GB)

GPU

Persistent Storage Size (GB)

Persistent Storage Class

Associate with Training Environments

FIGURE 20. Notebook Cluster Creation

Once the notebook status appears as configured, it is possible to view the notebook applications and access the application endpoints.

Select the Jupyterlab endpoint from the Notebook Endpoints tab.

Notebooks

Applications [Notebook Endpoints](#)

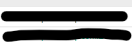
Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
nb-controller-prigm-0	controller	KubeDirectorApp: ID: jupyter-notebook Name: Jupyter Notebook with ML toolkits	nb	SSH Jupyter Notebook	22 8000		NodePort

FIGURE 21. Accessing Jupyterlab



Log into Jupyterlab using AD/LDAP credentials, and then launch a Python 3 notebook.

FIGURE 22. Sign in to Jupyterlab

Enter the `%kubeRefresh` magic in the Python cell and provide the logged-in member's password when prompted. Kubectl commands can now be entered from the notebook.

```
[6]: %kubeRefresh
.....
kubeconfig set for user imageanalyst
```

Training dataset to prepare the model. See the HPE Ezmeral Container Platform 5.3 documentation for [Notebook Magic Functions](#).

```
%%mltraining
print("Importing libraries")
import pandas as pd
import numpy as np
from scipy import stats
import math
import os
import datetime
import xgboost as xgb
import pickle

import matplotlib.pyplot as plt

# Start time
print("Start time: ", datetime.datetime.now())

# Project repo path function
def ProjectRepo(path):
    ProjectRepo = "/bd-fs-mnt/project_repo"
    return str(ProjectRepo + '/' + path)

print("Reading in data")
# Reading in dataset table
dbName = "pqyellowtaxi"
df = pd.read_csv(ProjectRepo['/data/demodata.csv'])

# Reading in latitude/longitude coordinate lookup table
```



```

lookupDbName = "pqlookup"
dflook = pd.read_csv(ProjectRepo['/data/lookup-ipyheader.csv'])
print("Done reading in data")

# merging dataset and lookup tables on latitudes/coordinates
df = pd.merge(df, dflook[[lookupDbName + '.location_i', lookupDbName + '.long', lookupDbName + '.lat']],
             how='left', left_on=dbName + '.polocationid', right_on=lookupDbName + '.location_i')
df.rename(columns = {(lookupDbName + '.long'):(dbName + '.startstationlongitude')}, inplace = True)
df.rename(columns = {(lookupDbName + '.lat'):(dbName + '.startstationlatitude')}, inplace = True)
df = pd.merge(df, dflook[[lookupDbName + '.location_i', lookupDbName + '.long', lookupDbName + '.lat']],
             how='left', left_on=dbName + '.dolocationid', right_on=lookupDbName + '.location_i')
df.rename(columns = {(lookupDbName + '.long'):(dbName + '.endstationlongitude')}, inplace = True)
df.rename(columns = {(lookupDbName + '.lat'):(dbName + '.endstationlatitude')}, inplace = True)

def fullName(colName):
    return dbName + '.' + colName

# convert string to datetime
df[fullName('tpep_pickup_datetime')] = pd.to_datetime(df[fullName('tpep_pickup_datetime')])
df[fullName('tpep_dropoff_datetime')] = pd.to_datetime(df[fullName('tpep_dropoff_datetime')])
df[fullName('duration')] = [df[fullName('tpep_dropoff_datetime')] -
df[fullName('tpep_pickup_datetime')]].dt.total_seconds()

# feature engineering
df[fullName("weekday")] = [df[fullName('tpep_pickup_datetime')].dt.dayofweek < 5].astype(float)
df[fullName("hour")] = df[fullName('tpep_pickup_datetime')].dt.hour
df[fullName("work")] = [df[fullName('weekday')] == 1] & [df[fullName("hour")] >= 8] & [df[fullName("hour")] < 18]
df[fullName("month")] = df[fullName('tpep_pickup_datetime')].dt.month
# convert month to a categorical feature using one-hot encoding
df = pd.get_dummies(df, columns=[fullName("month")])

# Filter dataset to rides under 3 hours and under 150 miles to remove outliers
df = df[df[fullName('duration')] > 20]
df = df[df[fullName('duration')] < 10800]
df = df[df[fullName('trip_distance')] > 0]
df = df[df[fullName('trip_distance')] < 150]

# drop null rows
df = df.dropna(how='any',axis=0)

# select columns to be used as features
cols = [fullName('work'), fullName('startstationlatitude'), fullName('startstationlongitude'),
fullName('endstationlatitude'), fullName('endstationlongitude'), fullName('trip_distance'), fullName('weekday'),
fullName('hour')]
cols.extend([fullName('month_' + str(x)) for x in range(1, 7)])
cols.append(fullName('duration'))
dataset = df[cols]

X = dataset.iloc[:, 0:(len(cols) - 1)].values
y = dataset.iloc[:, (len(cols) - 1)].values
X = X.copy()
y = y.copy()
del dataset
del df

print("Done cleaning data")

```



```

print("Training...")

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

xgbr = xgb.XGBRegressor(objective = 'reg:squarederror', colsample_bytree = 1, subsample = 1, learning_rate = 0.15,
booster = "gbtree", max_depth = 15, eta = 0.5, eval_metric = "rmse", tree_method='gpu-hist', gpu_id=0)
print("num train elements: " + str(len(X_train)))
print("Train start time: ", datetime.datetime.now())
xgbr.fit(X_train, y_train)
print("Train end time: ", datetime.datetime.now())
y_pred = xgbr.predict(X_test)
y_pred = y_pred.clip(min=0)

from sklearn import metrics
from sklearn.metrics import mean_squared_log_error
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
print('Root Mean Squared Log Error:', np.sqrt(mean_squared_log_error(y_test, y_pred)))
print()

print("Saving model")
pickle.dump(xgbr, open( ProjectRepo('models/') + "XGB.pickle.dat", "wb"))

# from xgboost import plot_importance
# plot_importance(xgbr, max_num_features=10) # top 10 most important features
# plt.show()

# Finish time
print("End time: ", datetime.datetime.now())

```

The output looks like this:

```

Importing libraries
Start time: 2021-08-24 16:50:31.977260
Reading in data
Done reading in data
Done cleaning data
Training...
num train elements: 264325
Train start time: 2021-08-24 16:50:35.330399
Train end time: 2021-08-24 16:51:22.312780
Mean Absolute Error: 179.24821523937294
Mean Squared Error: 83067.06130576036
Root Mean Squared Error: 288.21356891333267
Root Mean Squared Log Error: 0.3093430779701812

Saving model
End time: 2021-08-24 16:51:30.038040

```

Access the Model Management screen, and then click the Register New Model button to open the Register New Model screen.



Register the serialized model, being sure to include the Model Version and Path to Model Repo.

Register Model

Label

Name* ?	<input type="text" value="ezmlops-nyctaxi"/>
Description ?	<input type="text"/>
Model Store Type ?	<input type="text" value="Ezmeral Project Repository"/>
Model Version* ?	<input type="text" value="1.1"/>
Path to Model Repo* ?	<input type="text" value="repo://project_repo/models/XGB.pickle.dat"/> <input type="button" value="Browse"/>
Path to Scoring Script ?	<input type="text" value="repo://project_repo/models/XGB_Scoring.py"/> <input type="button" value="Browse"/>
Trained on Environment ?	<input type="text"/>

FIGURE 23. Registering Model



Access the Kubernetes Model Serving screen and then onboard the necessary training applications.

Create Model Serving

Cluster Detail

Name*

Description

Model Serving Engine*

Select Model*

Enable DataTap

Node Roles

LoadBalancer

Instances	<input type="text" value="1"/>
CPU	<input type="text" value="2"/>
Memory (GB)	<input type="text" value="4"/>
GPU	<input type="text" value="0"/>

RETServer

Instances	<input type="text" value="1"/>
CPU	<input type="text" value="2"/>
Memory (GB)	<input type="text" value="4"/>
GPU	<input type="text" value="0"/>

FIGURE 24. Launching Deployment Cluster



Now use the deployed model. To make predictions:

■ Create a Postman API call that is formulated as follows:

- a. Prefix: Either http:// or https://, as appropriate.
- b. Body: The access point from the ModelServingLoadbalancer of the Load Balancer role in the Ezmeral Serving Endpoints tab.

Model Serving

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
nyctaxi-loadbalancer-mc9vk-0	LoadBalancer	KubeDirectorApp: ID: deployment-engine Name: ML Inference	nyctaxi	Model serving request balancer stats API Server Model Serving LoadBalancer Model Serving Path	8081 10001 32700	exam-01.perflab.hp.com:10020 exam-01.perflab.hp.com:10025 [Auth Token] exam-01.perflab.hp.com:10027 [Auth Token] / <<model_name>>/<<model_version>>/predict	NodePort
nyctaxi-restserver-nvxpri-0	RETSerVer	KubeDirectorApp: ID: deployment-engine Name: ML Inference	nyctaxi	SSH API Server Model Serving Path	22 10001	exam-01.perflab.hp.com:10021 exam-01.perflab.hp.com:10026 [Auth Token] / <<model_name>>/<<model_version>>/predict	NodePort

FIGURE 25. Deployment Endpoints

- c. Suffix: Registered model name and version number, in the format/model_name_registered/version_number.

Model Management

Model Name	Model Version	Description	Model Store Type	Details	Actions
ezmlops-nyctaxi	1.1		Ezmeral Project Repository	Created At: Tue Aug 24, 2021 13:24:16 Created By: imagesadmin Model: repo://project_repo/models/XGB_pickle.dat Scoring Script: repo://project_repo/models/XGB_Scoring.py Tenant Namespace: mihenant	Register New Model 🔍 🗑️

FIGURE 26 Model Management

- d. Ending: /predict.

■ Verify that the finished API call looks similar to the example below:

m24wn02.bluedata:10038/ezmlops-nyctaxi/1.1/predict

■ In the Deployment Endpoints tab, click the Copy Auth Token link for the Load Balancer role.

Model Serving

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
nyctaxi-loadbalancer-mc9vk-0	LoadBalancer	KubeDirectorApp: ID: deployment-engine Name: ML Inference	nyctaxi	Model serving request balancer stats API Server Model Serving LoadBalancer Model Serving Path	8081 10001 32700	exam-01.perflab.hp.com:10020 exam-01.perflab.hp.com:10025 [Auth Token] exam-01.perflab.hp.com:10027 [Auth Token] / <<model_name>>/<<model_version>>/predict	NodePort
nyctaxi-restserver-nvxpri-0	RETSerVer	KubeDirectorApp: ID: deployment-engine Name: ML Inference	nyctaxi	SSH API Server Model Serving Path	22 10001	exam-01.perflab.hp.com:10021 exam-01.perflab.hp.com:10026 [Auth Token] / <<model_name>>/<<model_version>>/predict	NodePort

FIGURE 27. Auth Token

■ Launch Postman, and then enter the following information:

- a. Request URL: The URL created in Step 2.
- b. X-AUTH-TOKEN: Auth token copied in Step 3.

■ Use the raw body data to send queries.

```
{
  "use_scoring": true,
  "scoring_args": {
    "work": 0,
    "start_latitude": 40.57689727,
    "start_longitude": -73.99047356,
  }
}
```



```

    "end_latitude": 40.72058154,
    "end_longitude": -73.99740673,
    "distance": 8,
    "weekday": 1,
    "hour": 9,
    "month_1": 0,
    "month_2": 1,
    "month_3": 0,
    "month_4": 0,
    "month_5": 0,
    "month_6": 0
  }
}

```

Click **Send** to get a prediction.

The screenshot shows a Postman REST client interface. The top part displays a POST request to `http://ezam-01.perflab.hp.com:10027/ezmllops-nyctaxi/1.1/predict` with a header `X-Auth-Token: 0c35f3af9c5be243a750277c8b49a7b4`. The request body is a JSON object with fields for work status, start/end coordinates, distance, weekday, hour, and months. The response body shows a log entry with a warning about no visible GPU and a prediction of 2899.4675 seconds. The second part of the screenshot shows the same request with the 'Body' tab selected, displaying the raw JSON input and the response body with the prediction result.

FIGURE 28. Taxi Ride Prediction

See the HPE Ezmeral Container Platform 5.3 documentation for [Getting started with AI and ML in Kubernetes](#) for more details.

Use case on Pima Indian's diabetes prediction

This section explains model building, model registry, and model serving for Diabetes prediction using end-to-end ML Ops components on HPE Ezmeral Runtime.

- Dataset: <https://www.kaggle.com/uciml/pima-indians-diabetes-database>
- Attempting to classify whether or not a patient has diabetes based on some diagnostic measurements
- All patients in the dataset are females at least 21 years old of Pima Indian heritage



Follow the below step-by-step procedure to implement the use case on HPE Ezmeral Runtime.

■ Create an ML Ops tenant.


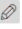

<input type="checkbox"/> mlop02	ML ops 2	ml02	Namespace: mlop02 Cores: 100 Memory: No Quota Ephemeral Storage: No Quota GPU Devices: No Quota Persistent Storage: No Quota ML Ops Project: True	  
---------------------------------	----------	------	---	---

FIGURE 29. ML Ops Tenant

■ Create a training cluster inside the ML Ops tenant.

HPE Ezmeral Container Platform (Acting as Tenant Admin)
mlop02 | admin

Dashboard

Users 0

Project Repository

Source Control

Model Registry

Training

Model Serving

DataTaps 1

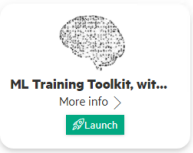
FsMounts 1

Applications

Notebooks

Training

Applications Training Endpoints



ML Training Toolkit, wit...
More info >
[Launch](#)

Running Applications

[Create Training](#)


<input type="checkbox"/>	Name	Description	Serving Framework	Created At	Role Configurations	Status	Member Status	Actions
<input type="checkbox"/>	ttl01		Kubedirector	Thu Dec 02 2021 01:18:30	LoadBalancer (1) RETServer (1) controller (1)	● configured		

FIGURE 30. Training Cluster



Create a notebook and attach a training cluster to it.

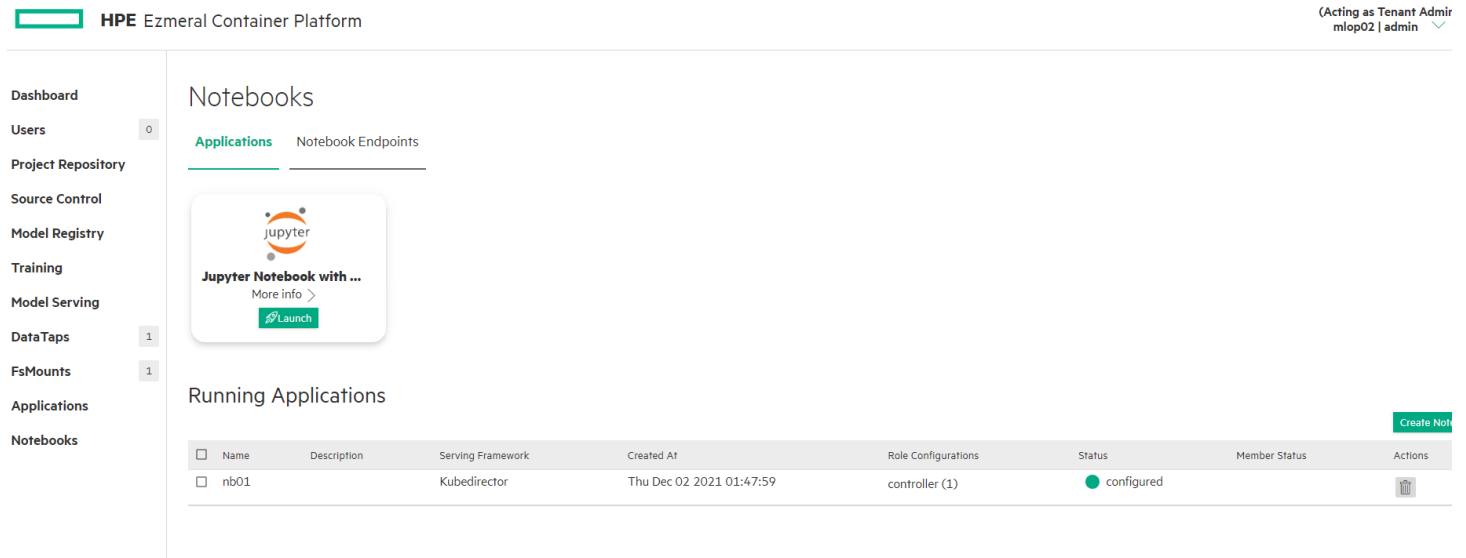


FIGURE 31. Creating notebook

Upload [dataset](#) to project Repository at project_repo/data/pima_Indians/ and scoring script at project_repo/code/Tensorflow/Diabetes_Scoring.py (can be found in the notebook default examples directory at examples/tensorflow/diabetes_prediction/ Diabetes_Scoring.py).

Also, create a new directory “Diabetes_Prediction” under project_repo/models.

The project_repo structure is as shown below.

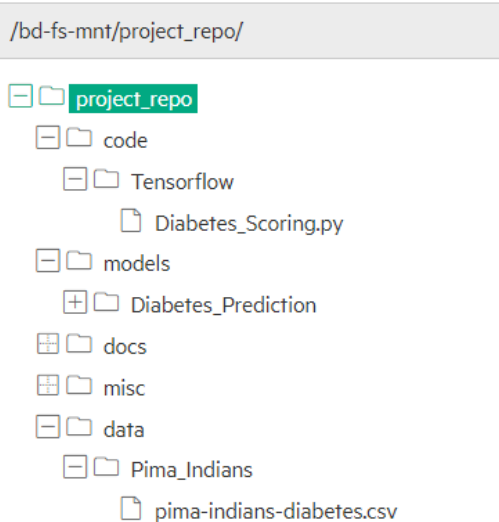


FIGURE 32. Project repo



■ Login to notebook and open the training cluster notebook at /examples/tensorflow/diabetes_prediction/Diabetes_Prediction-k8s.ipynb.

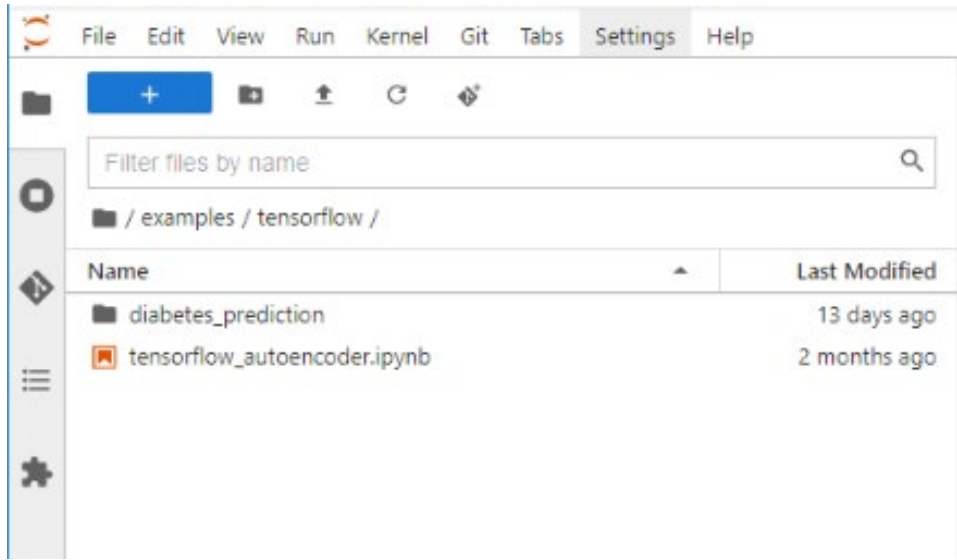


FIGURE 33. Training Cluster notebook

■ Setting up the environment and Data Preprocessing.

```
import numpy
import os
import pandas as pd
import tensorflow as tf

## Set the project repo
def ProjectRepo(path):
    ProjectRepo = '/bd-fs-mnt/project_repo/' + os.popen('bdvcli --get cluster.project_repo').read().rstrip()
    # print(ProjectRepo)
    return ProjectRepo + '/' + path
print(ProjectRepo('data/Pima_Indians/pima-indians-diabetes.csv'))
## Load the dataset
dataset = pd.read_csv(ProjectRepo('data/Pima_Indians/pima-indians-diabetes.csv'), delimiter=",")
dataset.columns = [
    "NumTimesPrg", "PlGlcConc", "BloodP",
    "SkinThick", "TwoHourSerIns", "BMI",
    "DiPedFunc", "Age", "HasDiabetes"]
```

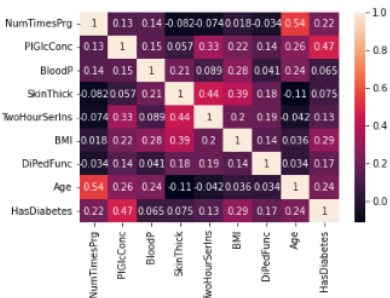


Visualizing correlation of variables with a heatmap and plotting histogram.

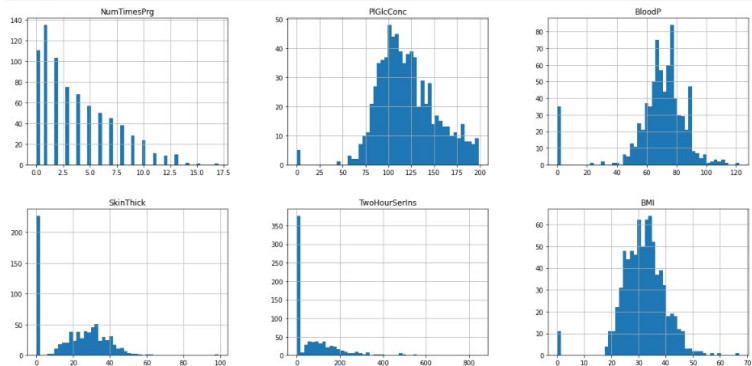
```
%matplotlib inline
import seaborn as sns
sns.heatmap(corr, annot = True)
import matplotlib.pyplot as plt
dataset.hist(bins=50, figsize=(20, 15))
plt.show()
```

```
[3]: %matplotlib inline
import seaborn as sns
sns.heatmap(corr, annot = True)
```

```
[3]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff0242f0e48>
```



```
[4]: import matplotlib.pyplot as plt
dataset.hist(bins=50, figsize=(20, 15))
plt.show()
```



Model development (Part 1)

Attempting the first model with XGB.

```
# First XGBoost model for Pima Indians dataset
from numpy import loadtxt
from xgboost import XGBClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import pickle

# load data
dataset = loadtxt('ProjectRepo\data/Pima_Indians/pima-indians-diabetes.csv', delimiter=',', skiprows=1)

# split data into X and y
X_train = dataset[:,0:8]
y_train = dataset[:,8]

# fit model no training data
model = XGBClassifier()
model.fit(X_train, y_train)

#
print(model.get_xgb_params())
```

```
[10:02:23] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
{'objective': 'binary:logistic', 'base_score': 0.5, 'booster': 'gbtree', 'colsample_bynode': 1, 'colsample_bylevel': 1, 'colsample_bytree': 1, 'gamma': 0, 'gpu_id': -1, 'interaction_constraints': '', 'learning_rate': 0.3000000012, 'max_delta_step': 0, 'max_depth': 6, 'min_child_weight': 1, 'monotone_constraints': '()', 'n_jobs': 64, 'num_parallel_tree': 1, 'random_state': 0, 'reg_alpha': 0, 'reg_lambda': 1, 'scale_pos_weight': 1, 'subsample': 1, 'tree_method': 'exact', 'validate_parameters': 1, 'verbosity': None}
```

```
[10]: %attachments
```

```
Training Cluster  ML Engine
-----
tc101             python
```

Model development (Part 2)

The second model uses Keras with a remote training cluster.

- Save model and prepare for TensorFlow Serving



```
%%tcl01

import numpy
import os
import pandas as pd
import tensorflow as tf
from numpy import loadtxt
from keras.models import Sequential
from keras.layers import Dense

with tf.device("/device:CPU:0"):

    ## Set the project repo
    def ProjectRepo(path):
        ProjectRepo = '/bd-fs-mnt/project_repo/' + os.popen('bdvcli --get cluster.project_repo').read().rstrip()
        print(ProjectRepo)
        return ProjectRepo + '/' + path

    ## Load the dataset
    print("Loading data")
    dataset = loadtxt(ProjectRepo['data/Pima_Indians/pima-indians-diabetes.csv'], delimiter=",", skiprows=1)
    dataset.shape

    # Split into input [X] and output [y] variables
    X = dataset[:,0:8]
    y = dataset[:,8]

    # Define the keras model
    print("Building model")
    model = Sequential()
    model.add(Dense(12, input_dim=8, activation='relu'))
    model.add(Dense(8, activation='relu'))
    model.add(Dense(1, activation='sigmoid'))

    # Compile the keras model
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

    # Fit the keras model on the dataset
    print("Training model")
    model.fit(X, y, epochs=150, batch_size=10, verbose=0)

    # Evaluate the keras model
    _, accuracy = model.evaluate(X, y, verbose=0)
    print('Accuracy: %.2f' % (accuracy*100))

    # Make class predictions with the model
    predictions = model.predict_classes(X)

    # Summarize the first 3 cases
    for i in range(3):
        print('%s => %d [expected %d]' % (X[i].tolist(), predictions[i], y[i]))

    # Save model weights and architecture together
    print("Saving model")
    model.save(ProjectRepo['models/Diabetes_Prediction/db_remote.h5'])

    # Evaluate the keras model
    _, accuracy = model.evaluate(X, y, verbose=0)
    print('Accuracy: %.2f' % (accuracy*100))

    # Make class predictions with the model
```



```

predictions = model.predict_classes(X)

# Summarize the first 5 cases
for i in range(5):
    print('%s => %d [expected %d]' % (X[i].tolist(), predictions[i], y[i]))

# Prepare TF Serving
print("Preparing for TF Serving")
MODEL_VERSION = 1
tf.keras.backend.set_learning_phase(0)
model = tf.keras.models.load_model(ProjectRepo('models/Diabetes_Prediction/db_remote.h5'))
export_path = ProjectRepo('models/Diabetes_Prediction/' + str(MODEL_VERSION))
tf.keras.models.save_model(model, export_path)

# Summarize model.
model.summary()
print("Done")

```

```

# Summarize model.
model.summary()
print("Done")

```

History URL: <http://tcl01-restserver-c6fgl-0.tcl01ms2q9.mlop02.svc.cluster.local:10001/history/4>

See the logs and Job Status as **Finished**.

```
[24]: %logs --url http://tcl01-restserver-c6fgl-0.tcl01ms2q9.mlop02.svc.cluster.local:10001/history/4
```

```

Job Status: Finished
Loading data
/bd-fs-mnt/project_repo/
Building model
Training model
Accuracy: 75.52
[6.0, 148.0, 72.0, 35.0, 0.0, 33.6, 0.627, 50.0] => 1 (expected 1)
[1.0, 85.0, 66.0, 29.0, 0.0, 26.6, 0.351, 31.0] => 0 (expected 0)
[8.0, 183.0, 64.0, 0.0, 0.0, 23.3, 0.672, 32.0] => 1 (expected 1)
Saving model
/bd-fs-mnt/project_repo/
Accuracy: 75.52
[6.0, 148.0, 72.0, 35.0, 0.0, 33.6, 0.627, 50.0] => 1 (expected 1)
[1.0, 85.0, 66.0, 29.0, 0.0, 26.6, 0.351, 31.0] => 0 (expected 0)
[8.0, 183.0, 64.0, 0.0, 0.0, 23.3, 0.672, 32.0] => 1 (expected 1)
[1.0, 89.0, 66.0, 23.0, 94.0, 28.1, 0.167, 21.0] => 0 (expected 0)
[0.0, 137.0, 40.0, 35.0, 168.0, 43.1, 2.288, 33.0] => 1 (expected 1)
Preparing for TF Serving
/bd-fs-mnt/project_repo/
/bd-fs-mnt/project_repo/
Model: "sequential_1"

Layer (type)                Output Shape              Param #
-----
dense_1 (Dense)              (None, 12)                108
-----
dense_2 (Dense)              (None, 8)                 104
-----
dense_3 (Dense)              (None, 1)                 9
-----
Total params: 221
Trainable params: 221
Non-trainable params: 0
Done

```

FIGURE 34. Job Status Logs



■ For the model registry, go to tenant UI and register the model.

Register Model

Label

Name* ?	<input type="text" value="diabetesmodel"/>
Description ?	<input type="text"/>
Model Store Type ?	<input type="text" value="Ezmeral Project Repository"/>
Model Version* ?	<input type="text" value="1.0"/>
Path to Model Repo* ?	<input type="text" value="//project_repo/models/Diabetes_Prediction/db_remote.h5"/> <input type="button" value="Browse"/>
Path to Scoring Script ?	<input type="text" value="//project_repo/code/Tensorflow/Diabetes_Scoring.py"/> <input type="button" value="Browse"/>
Trained on Environment ?	<input type="text"/>

FIGURE 35. Model Registry



Create Model Serving.

Create Model Serving

Cluster Detail

Name*

Description

Model Serving Engine*

Select Model*

Enable DataTap

Node Roles

LoadBalancer

Instances

CPU

Memory (GB)

GPU

RESTServer

Instances

CPU

Memory (GB)

GPU

FIGURE 36. Creating Model Serving



Once Model Serving is done, obtain the endpoint and Auth token from the Model Serving section.

Model Serving

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
diabetermodel1-loadbalancer-lk6hf-0	LoadBalancer	KubeDirectorApp: ID: deployment-engine Name: ML Inferencing	diabetermodel1	Model serving request balancer stats API Server Model Serving LoadBalancer Model Serving Path	8081 10001 32700	m24wn02.bluedata:10048 m24wn02.bluedata:10050 [Auth Token] m24wn02.bluedata:10051 [Auth Token] /<<model_name>>/<<model_version>>/predict	NodePort
diabetermodel1-restserver-5xq6m-0	RETSerVer	KubeDirectorApp: ID: deployment-engine Name: ML Inferencing	diabetermodel1	SSH API Server Model Serving Path	22 10001	m24wn02.bluedata:10047 m24wn02.bluedata:10049 [Auth Token] /<<model_name>>/<<model_version>>/predict	NodePort

FIGURE 37. Model Serving

Prediction can be done by making a POST request via curl or postman to the Model Serving endpoint.

```
curl -X POST -H "Content-Type:application/json" -H "X-Auth-Token:<Auth Token>" -d '{"use_scoring": true, "scoring_args": {"NumPreg":1.0, "Glucose": 85.0, "BloodPressure": 66.0, "SkinThick": 29.0, "Insulin": 0.0, "BMI": 26.6, "DiabetesPedFunc": 0.351, "Age": 35.0 } }' http://<model-serving-Loadbalancer:port>/<modelname>/<version>/predict
```

```
For example : curl -X POST -H "Content-Type:application/json" -H "X-Auth-Token:fb3f6521015978c2f5d58530a42ce720" -d '{"use_scoring": true, "scoring_args": {"NumPreg":1.0, "Glucose": 85.0, "BloodPressure": 66.0, "SkinThick": 29.0, "Insulin": 0.0, "BMI": 26.6, "DiabetesPedFunc": 0.351, "Age": 35.0 } }' http://m24wn02.bluedata:10050/diabetesmodel/1.0/predict
```

```
[root@m24wn13 ~]# curl -X POST -H "Content-Type:application/json" -H "X-Auth-Token:fb3f6521015978c2f5d58530a42ce720" -d '{"use_scoring": true, "scoring_args": {"NumPreg":1.0, "Glucose": 85.0, "BloodPressure": 66.0, "SkinThick": 29.0, "Insulin": 0.0, "BMI": 26.6, "DiabetesPedFunc": 0.351, "Age": 35.0 } }' http://m24wn02.bluedata:10050/diabetesmodel/1.0/predict
{"input":{"NumPreg":1.0, "Glucose": 85.0, "BloodPressure": 66.0, "SkinThick": 29.0, "Insulin": 0.0, "BMI": 26.6, "DiabetesPedFunc": 0.351, "Age": 35.0}, "log_url":"http://diabetermodel1-loadbalancer-lk6hf-0.diabetermodel14g65t.mlop02.svc.cluster.local:10001/logs/15", "node":"diabetermodel1-loadbalancer-lk6hf-0.diabetermodel14g65t.mlop02.svc.cluster.local", "output":"Chances of having diabetes: 0%\n\n", "pid":"990", "qid":"15", "request_url":"http://diabetermodel1-loadbalancer-lk6hf-0.diabetermodel14g65t.mlop02.svc.cluster.local:10001/history/15", "status":"Finished"}
```

This concludes our testing for ML Ops using the Training cluster and KubeDirector Notebook.

EZMERAL ML OPS - IN ACTION WITH USE CASE (SPARK OPERATOR)

Spark operator with K8s

Spark is the cluster computing framework for big data processing. It can also be used for distributed data processing on different machines. It can be used for batch, stream, and interactive data processing. Various file systems can be loaded in Spark. The Spark API is available in Scala, Python, Java, and R. Spark provides libraries such as Spark SQL (for structured data processing), MLlib (for scalable and easy ML), GraphX (for iterative graph computation within a single system), and Spark Streaming (to process real-time data from various sources).



The Kubernetes Operator for Apache Spark included in the HPE Ezmeral Runtime makes running Spark jobs easy. The Spark operator is a custom Kubernetes resource that is configured in a tenant namespace to allow on-demand cluster deployment with Spark Executor pods. For managing Spark jobs, it uses declarative specifications. It dynamically creates the specified number of driver and executor pods during the execution and then deletes the executor pods and leaves the driver pod in the completed state when the job finishes successfully. The driver pod does not consume any Kubernetes resources in this state, and the logs can be viewed to see execution details or results.

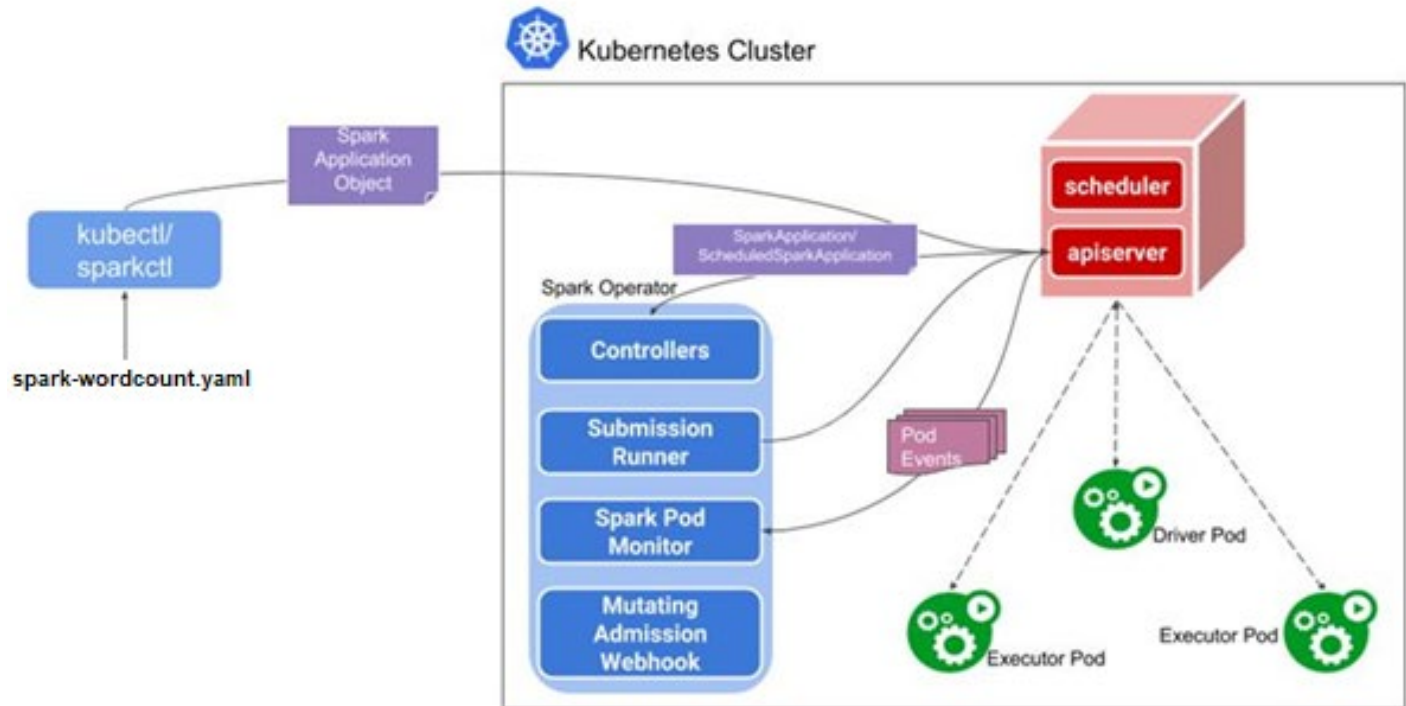


FIGURE 38. Spark operator architecture on K8s

Prerequisites

- HPE Ezmeral Runtime
- Platform administrators can access the web interface
- Root access to the controller host

System requirements

- For supported Kubernetes versions, see [Kubernetes Version Requirements](#).
- For issues and workarounds, see [Issues and Workarounds](#).
- The following resources must be available to install Kubeflow: h
- Minimum number of nodes for compute cluster: 2 (1 primary, 1 secondary)
- Minimum core and memory resources required:
 - CPU Cores: 36
 - Memory (GB): 160

AD/LDAP authentication requirements

- The Kubernetes cluster where Kubeflow will be installed must have AD/LDAP user authentication configured. The AD/LDAP user authentication configuration is posted as a secret in the cluster.
- For information about setting AD/LDAP user authentication configuration, see [Authentication, in Creating a New Kubernetes Cluster](#).



Preparing the environment

To prepare the environment:

- Log in to the web interface as a Kubernetes Administrator.
- Create a Kubernetes cluster, being sure that this cluster meets or exceeds the prerequisites described above.
- Check the 'Enable Spark Operator'.

Edit Kubernetes Cluster

Progress bar: Host Configurations (✓) — Cluster Configurations (✓) — Authentication (✓) — **4 Application Configurations** — 5 Summary

Select from the list of compute applications

- Enable Spark operator
- Istio
- Enable Kubeflow
- Enable Airflow

Policy Settings

Policies are not available. Create one first.

Previous Next

FIGURE 39. Application configurations

■ This bootstraps a Spark operator and creates the following:

- a. Integrated HPE Ezmeral Data Fabric Container Storage Interface (CSI) and associated service accounts.
- b. Spark namespaces, service accounts, cluster roles, and role bindings, Spark applications Custom Resource Definition (CRD), Spark operator, and compute templates for Spark.
- c. Auto Ticket Generator Service, roles, and role bindings.

■ Create a new Kubernetes tenant. Do not assign any quotas when creating this tenant. See Figure 13 for tenant creation.

■ Assign a Kubernetes Cluster Administrator user to the cluster created in Step 2. See Figure 12 for assigning users to the cluster.

■ Assign a Tenant Administrator user to the tenant created in Step 4. See Figure 15 for admin assignment to the tenant.

■ Log out of the web interface when this process is complete.

■ HPE Ezmeral Runtime automatically creates the following:

- a. Spark operator namespace.
- b. User secrets for the tenant members (necessary for running the Spark workload with the Spark operator).
- c. The Role-Based Access Control (RBAC) for Spark resources is also configured for the AD/LDAP tenant members.

Spark operator use case

Following are the steps to run a Spark job using the Kubernetes Web Terminal as a local user.

■ Log in to the web interface as the Site Administrator, and add/assign Tenant users. See Figure 16 for members assigned to the tenant.

■ In the FS Mounts screen, click the TenantShare link in the Name column of the table to open the FS Mount Browser screen. This screen functions identically to the DataTap browser screen.



■ Create the data subdirectory in the TenantShare filesystem mount, and then either create a text file or [download](#) this example as wordcount.txt.

■ Upload the wordcount.txt file to the data subdirectory. In HPE Ezmeral Runtime, the location to this subdirectory is /hcp/tenant-
<tenant_id>/fsmount/data/wordcount.txt.

```
$ /hcp/tenant-14/fsmount/data/wordcount.txt
```

■ Copy the following text, and then save it as `spark-wc.yaml`.

```
apiVersion: "sparkoperator.hpe.com/v1beta2"
kind: SparkApplication
metadata:
  name: spark-wordcount-secure
  namespace: mltenant
spec:
  #sparkConf:
  # Note: If you are executing the application as a K8 user that MapR can verify,
  #       you do not need to specify a spark.mapr.user.secret
  #spark.mapr.user.secret: spark-user-secret
  # Note: You do not need to specify a spark.eventLog.dir
  #       it will be auto-generated with the pattern "maprfs:///apps/spark/<namespace>"
  #spark.eventLog.dir: "maprfs:///apps/spark/sampletenant"
  type: Java
  sparkVersion: 2.4.4
  mode: cluster
  image: gcr.io/mapr-252711/spark-2.4.4:202009090453C
  imagePullPolicy: Always
  mainClass: org.apache.spark.examples.JavaWordCount
  mainApplicationFile: "local:///opt/mapr/spark/spark-2.4.4/examples/jars/spark-examples_2.11-2.4.4.6-mapr-630.jar"
  restartPolicy:
    type: Never
  arguments:
  - maprfs:///hcp/tenant-14/fsmount/data/wordcount.txt
  imagePullSecrets:
  - imagepull
  driver:
    cores: 1
    coreLimit: "1000m"
    memory: "512m"
    labels:
      version: 2.4.4
    # Note: You do not need to specify a serviceAccount
    #       it will be auto-generated referencing the pre-existing "hpe-<namespace>"
    #serviceAccount: hpe-sampletenant
  executor:
    cores: 1
    coreLimit: "1000m"
    instances: 2
    memory: "512m"
    labels:
      version: 2.4.4
```

■ Use the Kubernetes Web Terminal to edit `spark-wordcount.yaml` by updating the namespace and input file name and path for wordcount.txt.

■ Execute the Spark wordcount job by executing the following command:

```
$ kubectl apply -f /bd-fs-mnt/TenantShare/apps/spark-wordcount.yaml -n mltenant
```



```
k8suser@kdss-5zd6r-0: /bd-fs-mnt/TenantShare/apps/spark$ kubectl apply -f spark-wordcount.yaml -n mltenant
sparkapplication.sparkoperator.hpe.com/spark-wordcount-secure created
```

Check the pods running within the tenant namespace by executing the following command:

```
$ kubectl get pods -n mltenant
```

```
k8suser@kdss-5zd6r-0: /bd-fs-mnt/TenantShare/apps/spark$ kubectl get pods -n mltenant
NAME                                READY   STATUS    RESTARTS   AGE
nbone-controller-t9vtj-0            1/1    Running   0           4h13m
nyctaxi-loadbalancer-mc9vk-0        1/1    Running   0           3h50m
nyctaxi-restserver-nvxpn-0         1/1    Running   0           3h50m
spark-pi-secure-driver              0/1    Completed 0           36m
spark-wordcount-secure-driver       0/1    Completed 0           97s
sparkhs-5c5bb8c6b9-4wqsz           1/1    Running   4           25h
tenantcli-0                         1/1    Running   0           25h
```

Check the job status of the job by executing the following command:

```
$ kubectl logs spark-wordcount-secure-driver -follow
```

```
programs: 2
watch: 1
Wikipedia: 1
were: 3
length: 5
and: 23
subgenre: 1
being: 1
measures: 1
21/08/24 22:14:58 INFO server.AbstractConnector: Stopped Spark@5dbf5634{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
21/08/24 22:14:58 INFO ui.SparkUI: Stopped Spark web UI at http://spark-wordcount-secure-1629843256056-driver-svc.mltenant.svc:4040
21/08/24 22:14:58 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
21/08/24 22:14:58 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
21/08/24 22:14:58 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)
21/08/24 22:14:58 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
21/08/24 22:14:58 INFO memory.MemoryStore: MemoryStore cleared
21/08/24 22:14:58 INFO storage.BlockManager: BlockManager stopped
21/08/24 22:14:58 INFO storage.BlockManagerMaster: BlockManagerMaster stopped
21/08/24 22:14:58 INFO scheduler.OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
21/08/24 22:14:58 INFO spark.SparkContext: Successfully stopped SparkContext
```

Validate that the job has been completed by checking the status of the pod:

```
$ kubectl get pods -n mltenant
```

Model training from KubeDirector Notebook using Spark with Livy

Following are the steps to run a Spark job with Livy operator REST API from KubeDirector Notebook in the Kubernetes by AD/LDAP user.

- Log in to the web interface as the Site Administrator, add/assign Tenant users. See Figure 16 for members assigned to the tenant.
- In the DataTaps screen, click the TenantStorage link in the Name column of the table to open the DataTap Browser screen.
- Create the data subdirectory in the TenantStorage filesystem, and then upload or download it as `train.csv`. Refer to this example [train.csv](#).



Upload the train.csv file to the data subdirectory. In HPE Ezmeral Runtime, the location to this subdirectory is /hcp/tenant-<tenant_id>/data/train.csv.

HPE Ezmeral Container Platform

FIGURE 40. DataTaps

Fetch Livy endpoint.

HPE Ezmeral Container Platform

(Acting as Tenant Admin)
mlops1 | admin

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
livy-http				http	8998	m24wn02.bluedata:10015	NodePort
miflow-controller-gfsm-0	controller	KubeDirectorApp: ID: miflow Name: MLFlow server	miflow	SSH	22	m24wn02.bluedata:10026	NodePort
				MySQL	3306	m24wn02.bluedata:10027	
				Minio s3 bucket	9000	m24wn02.bluedata:10028	
				MLFlow Server	5000	m24wn02.bluedata:10029	
spark-ui-proxy				http	80	m24wn02.bluedata:10014	NodePort
sparkhs-svc				http	18480	m24wn02.bluedata:10012	NodePort
sparkts-svc				http	4440	m24wn02.bluedata:10013	NodePort
				spark-thrift	2304	m24wn02.bluedata:10016	

FIGURE 41. Service Endpoints



Launch KubeDirector Notebook.

HPE Ezmeral Container Platform

FIGURE 42. Creating Notebook

Access JupyterHub Notebook using the service endpoint.

HPE Ezmeral Container Platform

(Acting as Tenant Admin)
mlops1 | admin

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
mifftest-controller-rvppf-0	controller	KubeDirectorApp: ID: jupyter-notebook Name: Jupyter Notebook with ML toolkits	mifftest	SSH Jupyter Notebook	22 8000	m24wn02.bluedata:10017 m24wn02.bluedata:10018	NodePort
sparktest-controller-hbmv6-0	controller	KubeDirectorApp: ID: jupyter-notebook Name: Jupyter Notebook with ML toolkits	sparktest	SSH Jupyter Notebook	22 8000	m24wn02.bluedata:10024 m24wn02.bluedata:10025	NodePort

FIGURE 43. Notebook Endpoints



Launch PySpark kernel to configure and create Livy session.

Configure the Spark Livy Session Context

Modify Spark parameters for the session. In our examples we are including the DTAP Libraries and we specify an image that we want to use for our livy sessions using the spark cluster.

```
[1]: %configure -f
{"driverCores": 1, "executorCores": 2,
 "conf": {"spark.kubernetes.container.image": "gcr.io/mapr-252711/spark-2.4.7:202106220630P141", "spark.hadoop.fs.dtap.impl": "com.bluedata.hadoop.bdfs.Bdfs", "spark.hadoop.fs.AbstractFileSystem.dtap.impl": "co

Current session configs: {'driverCores': 1, 'executorCores': 2, 'conf': {'spark.kubernetes.container.image': 'gcr.io/mapr-252711/spark-2.4.7:202106220630P141', 'spark.hadoop.fs.dtap.impl':
'com.bluedata.hadoop.bdfs.Bdfs', 'spark.hadoop.fs.AbstractFileSystem.dtap.impl': 'com.bluedata.hadoop.bdfs.BdAbstractFS', 'spark.hadoop.fs.dtap.impl.disable.cache': 'false',
'spark.driver.extraClassPath': 'local:///opt/bdfs/bluedata-dtap.jar', 'spark.executor.extraClassPath': 'local:///opt/bdfs/bluedata-dtap.jar', 'spark.kubernetes.driver.label.hpecp.hpe.com/dtap':
'hadoop2', 'spark.kubernetes.executor.label.hpecp.hpe.com/dtap': 'hadoop2'}, 'kind': 'pyspark'}
No active sessions.
```

Create a Spark Context

Getting a Spark Context can take a few seconds as the livy containers are started in the background. The Startup Time depends on the speed between the k8s nodes and the image repository

```
[*]: print("hi")
Enter Livy endpoint (e.g., http://<internal-livy-session-url>:<port>) : https://172.30.226.52:10015
Enter your password: *****
```

Create a Spark Context

Getting a Spark Context can take a few seconds as the livy containers are started in the background. The Startup Time depends on the speed between the k8s nodes and the image repository

```
[2]: print("hi")
Enter Livy endpoint (e.g., http://<internal-livy-session-url>:<port>) : https://172.30.226.52:10015
Enter your password: .....
WARNING: Restart the kernel if you entered something wrong or if the kernel fails.
Starting Spark application
```

Spark Session ID	Kind	State	Spark UI	Driver log	User	Current session?
6	pyspark	idle			imageadmin	✓

```
SparkSession available as 'spark'.
hi
```

Copy this code to Notebook and execute it to read data from DataTap to Spark Data Frame.

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import isnull, when, count, col
from pyspark.ml.feature import StringIndexer
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.classification import RandomForestClassifier
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

#spark = SparkSession.builder.appName("titanic").getOrCreate()

# read locally
sdf = spark.read.format("csv").option('header', 'true').load("dtap://TenantStorage/data/train.csv") # data from
https://www.kaggle.com/c/titanic/data

# or read from mapr file system
# titanic_sdf = [spark.read.format("csv").option('header', 'true').load("maprfs:///exthcp/tenant-
5/fsmount/repo/data/train.csv")]

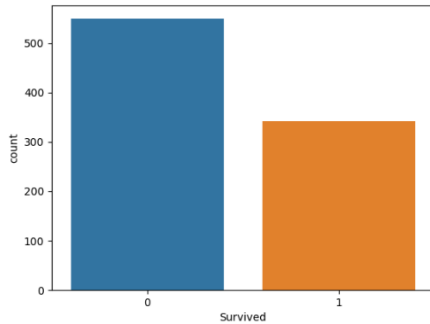
#Converting Spark DataFrame To Pandas DataFrame for exploratory data analysis
pdf = sdf.toPandas()
```



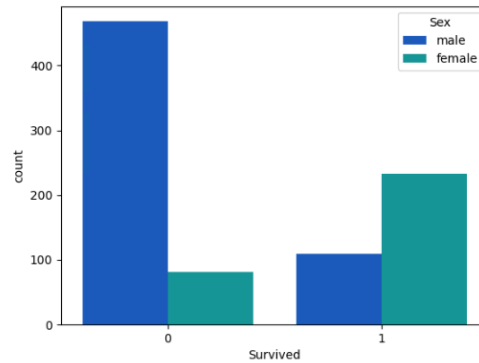
Explore Data analysis.

Exploratory data analysis

```
[18]: sns.countplot(x="Survived", data=pdf)
      %matplotlib plt
```



```
[22]: plt.clf()
      sns.countplot(x="Survived", hue="Sex", data=pdf, palette='winter')
      %matplotlib plt
```



Transform and Train model.

```
data = sdf.select(col('Survived').cast('float'), col('Pclass').cast('float'),
                  col('Sex'), col('Age').cast('float'),
                  col('Fare').cast('float'), col('Embarked'))

data.select([count(when(isnull(c), c)).alias(c) for c in data.columns]).show()

data = dataset.replace('?', None).dropna(how='any')

data = StringIndexer(
    inputCol='Sex',
    outputCol='Gender',
    handleInvalid='keep').fit(data).transform(data)
data = StringIndexer(
    inputCol='Embarked',
    outputCol='Boarded',
    handleInvalid='keep').fit(data).transform(data)

data = data.drop('Sex')
data = data.drop('Embarked')
data.show()

required_features = ['Pclass', 'Age', 'Fare', 'Gender', 'Boarded']
assembler = VectorAssembler(inputCols=required_features, outputCol='features')
transformed_data = assembler.transform(data)
(training_data, test_data) = transformed_data.randomSplit([0.8, 0.2])
rf = RandomForestClassifier(labelCol='Survived',
                           featuresCol='features',
                           maxDepth=5)
model = rf.fit(training_data)
```



```

+-----+
|Survived|Pclass|Sex|Age|Fare|Embarked|
+-----+
|      0|      0| 0|177|  0|      2|
+-----+

+-----+
|Survived|Pclass| Age|  Fare|Gender|Boarded|
+-----+
|      0.0|  3.0|22.0|  7.25|  0.0|  0.0|
|      1.0|  1.0|38.0|71.2833|  1.0|  1.0|
|      1.0|  3.0|26.0|  7.925|  1.0|  0.0|
|      1.0|  1.0|35.0|  53.1|  1.0|  0.0|
|      0.0|  3.0|35.0|  8.05|  0.0|  0.0|
|      0.0|  1.0|54.0|51.8625|  0.0|  0.0|
|      0.0|  3.0| 2.0|21.075|  0.0|  0.0|
|      1.0|  3.0|27.0|11.1333|  1.0|  0.0|
|      1.0|  2.0|14.0|30.0708|  1.0|  1.0|
|      1.0|  3.0| 4.0| 16.7|  1.0|  0.0|
|      1.0|  1.0|58.0| 26.55|  1.0|  0.0|
|      0.0|  3.0|20.0|  8.05|  0.0|  0.0|
|      0.0|  3.0|39.0|31.275|  0.0|  0.0|
|      0.0|  3.0|14.0| 7.8542|  1.0|  0.0|
|      1.0|  2.0|55.0| 16.0|  1.0|  0.0|
|      0.0|  3.0| 2.0|29.125|  0.0|  2.0|
|      0.0|  3.0|31.0| 18.0|  1.0|  0.0|
|      0.0|  2.0|35.0| 26.0|  0.0|  0.0|
|      1.0|  2.0|34.0| 13.0|  0.0|  0.0|
|      1.0|  3.0|15.0| 8.0292|  1.0|  2.0|
+-----+
only showing top 20 rows

```

Validate and save model.

```

predictions = model.transform(test_data)
evaluator = MulticlassClassificationEvaluator(
    labelCol='Survived',
    predictionCol='prediction',
    metricName='accuracy')
accuracy = evaluator.evaluate(predictions)
print('=====' * 10)
print('Test Accuracy = ', accuracy)
print('=====' * 10)
# saving locally
model_location = "maprfs:///hcp/tenant-8/dco/data/sparkml-titanic4"

# saving on mapr file system
# model_location = "maprfs:///exthcp/tenant-5/fsmount/repo/models/sparkml-titanic"

model.save(model_location)
print("Model Save to location : {}".format(model_location))

```



```
=====
Test Accuracy = 0.7755102040816326
=====
```

Saving the Model

```
[17]: # saving locally
model_location = "maprfs:///hcp/tenant-8/dco/data/sparkml-titanic4"

# saving on mapr file system
# model_location = "maprfs:///exthcp/tenant-5/fsmount/repo/models/sparkml-titanic"

model.save(model_location)
print("Model Save to location : {}".format(model_location))
```

Model Save to location : maprfs:///hcp/tenant-8/dco/data/sparkml-titanic4

HPE Ezmeral Container Platform

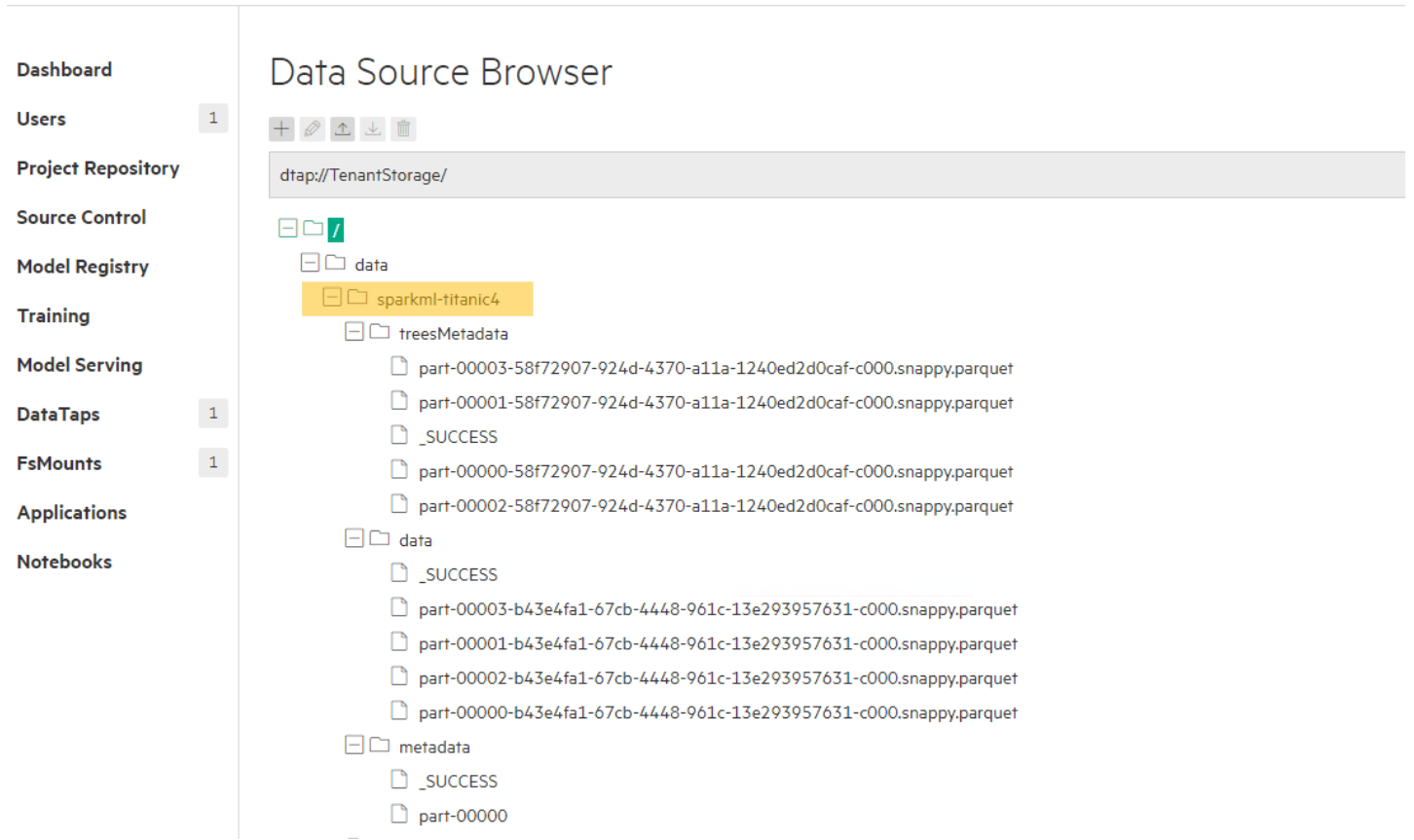


FIGURE 44. Data Source Browser



EZMERAL ML OPS – EXPERIMENT TRACKING WITH MLFLOW

The use case contains dataset preprocessing, model training and evaluation, model tuning via MLflow tracking, and finding the best-trained model.

Goal: Predict `rented_bikes` (count per hour) based on weather and time information.

Prerequisite

Dataset: Bike Sharing [Dataset](#)

Use case workflow

Following are the steps to implement the use case and run the experiment into the MLflow model server.

■ Initialize the web terminal.

■ After login to HPE Ezmeral Container Platform 5.3 and in MLOps Tenant, initialize the web terminal.

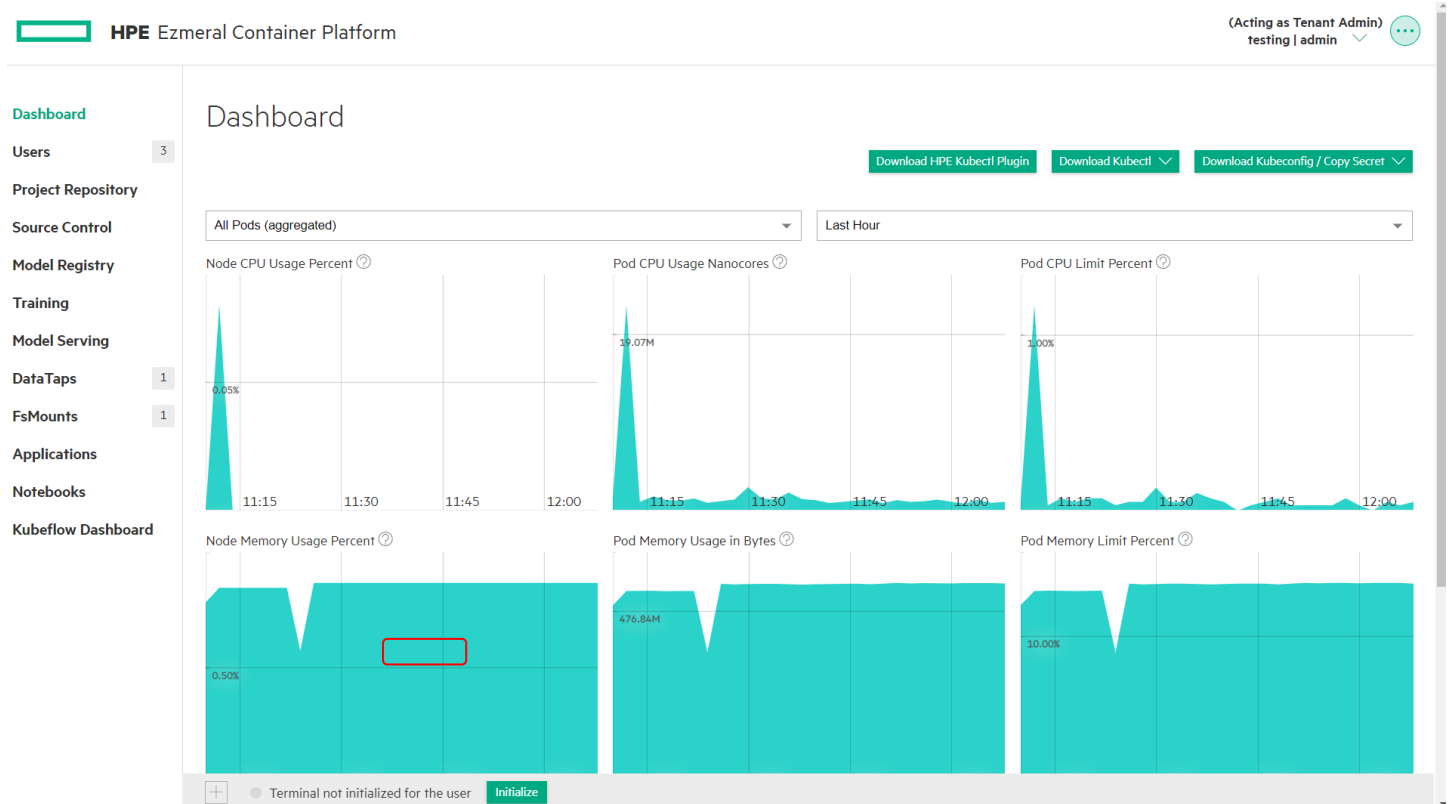


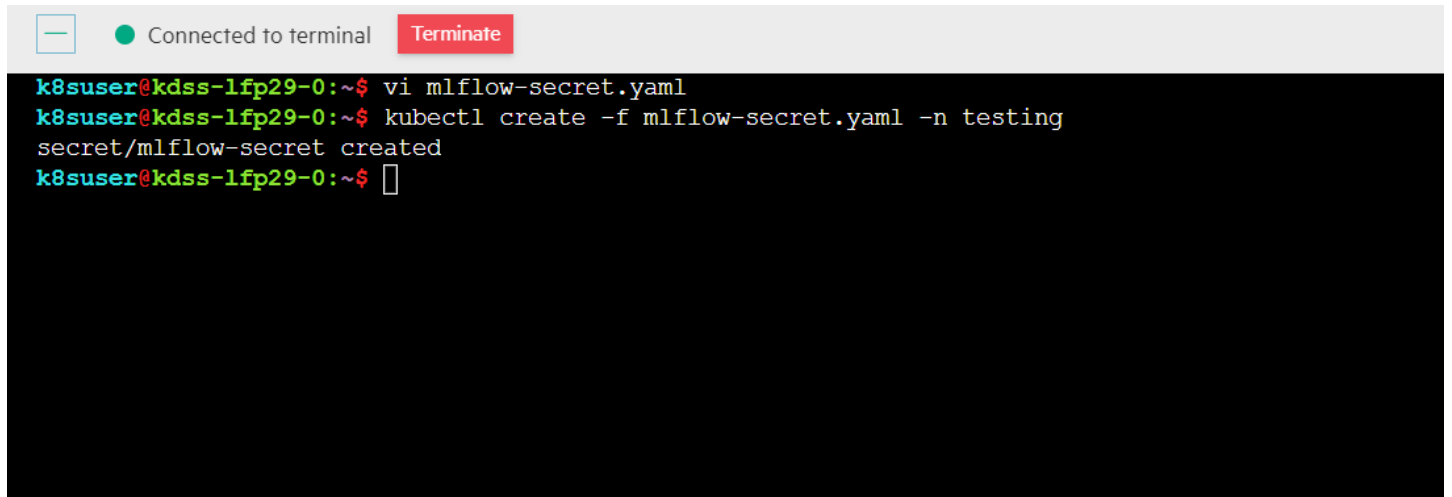
FIGURE 45. Dashboard

■ Create Secret. Once the web terminal is connected, create a file `mlflow-secret.yaml`.

```
apiVersion: v1
kind: Secret
type: Opaque
metadata:
  name: mlflow-secret
  labels:
    kubedirector.hpe.com/secretType: mlflow
data:
  AWS_ACCESS_KEY_ID: YWRtaW4= #admin
  AWS_SECRET_ACCESS_KEY: YWRtaW4xMjM= #admin123
  MLFLOW_ARTIFACT_ROOT: czM6Ly9tbGZsb3c= #s3://mlflow
```



```
$ kubectl create -f mlflow-secret.yaml -n testing
```



Launch MLflow Server. Go to the Application Tab and launch MLflow Server.

The screenshot shows the HPE Ezmeral Container Platform dashboard. At the top, it says 'HPE Ezmeral Container Platform' and '(Acting as Tenant Admin) testing | admin'. The main section is titled 'Kubernetes Applications' and has tabs for 'KubeDirector', 'Kubectl', 'Service Endpoints', and 'Virtual Endpoints'. Under 'KubeDirector', there are several application tiles: 'CentOS 8.0', 'ELK Stack 7.7.1 v1', 'MLFlow server', 'NVIDIA:TensorFlow(NGC)', and 'TensorFlow + Jupyter'. Each tile has a 'Launch' button. Below these is an 'Ubuntu 18.04' tile. A 'Running Applications' section at the bottom shows a table with columns: Name, Description, Serving Framework, Created At, Role Configurations, Status, Member Status, and Actions. The table is currently empty with the message 'Sorry, no matching records found'. A 'Create KubeDirector' button is in the top right of this section. At the bottom of the dashboard, there is a terminal window showing 'Connected to terminal' and a 'Terminate' button.

FIGURE 46. Kubernetes Applications



We need to attach the secret created in step 3 in YAML of MLflow Server before launching it.

HPE Ezmeral Container Platform (Acting as Tenant Admin) testing | admin

- Dashboard
- Users 3
- Project Repository
- Source Control
- Model Registry
- Training
- Model Serving
- DataTaps 1
- FsMounts 1
- Applications
- Notebooks
- Kubeflow Dashboard

Create KubeDirector

Edit your Yaml and hit submit to launch app

```

1 ---
2 apiVersion: "kubedirector.hpe.com/v1beta1"
3 kind: "KubeDirectorCluster"
4 metadata:
5   name: "mlflow"
6   namespace: "testing"
7   labels:
8     description: "mlflow-application"
9 spec:
10  app: "mlflow"
11  namingScheme: "CrNameRole"
12  appCatalog: "local"
13  connections:
14    secrets:
15      - mlflow-secret
16  roles:
17    -
18      id: "controller"
19      members: 1
20      resources:
21        requests:
22          cpu: "2"
23          memory: "4Gi"
24          nvidia.com/gpu: "0"
25        limits:
26          cpu: "2"
27          memory: "4Gi"
28          nvidia.com/gpu: "0"
29      #Note: "if the application is based on hadoop3 e.g. using StreamCapabilities interface, then change the below dtap label to 'hadoop3', otherwise for most
30      #podLabels:
31      #hpecp.hpe.com/dtap: "hadoop2"
32

```

+ Connected to terminal Terminate

FIGURE 47. Create/edit KubeDirector YAML

Accessing Minio and MLflow Server.

Kubernetes Applications

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
kf-dashboard-import-b5ntq				80	80	hpecp-531-vm2.rcc.local:10028	NodePort
mlflow-controller-jf5mk-0	controller	KubeDirectorApp: ID: mlflow Name: MLFlow server	mlflow	SSH MySQL Minio s3 bucket MLFlow Server	22 3306 9000 5000	hpecp-531-vm2.rcc.local:10030 hpecp-531-vm2.rcc.local:10031 hpecp-531-vm2.rcc.local:10032 hpecp-531-vm2.rcc.local:10033	NodePort

FIGURE 48. Service Endpoints



On the MinIO Browser page, create a bucket name `mlflow` which is configured in `mlflow-secret.yaml`.

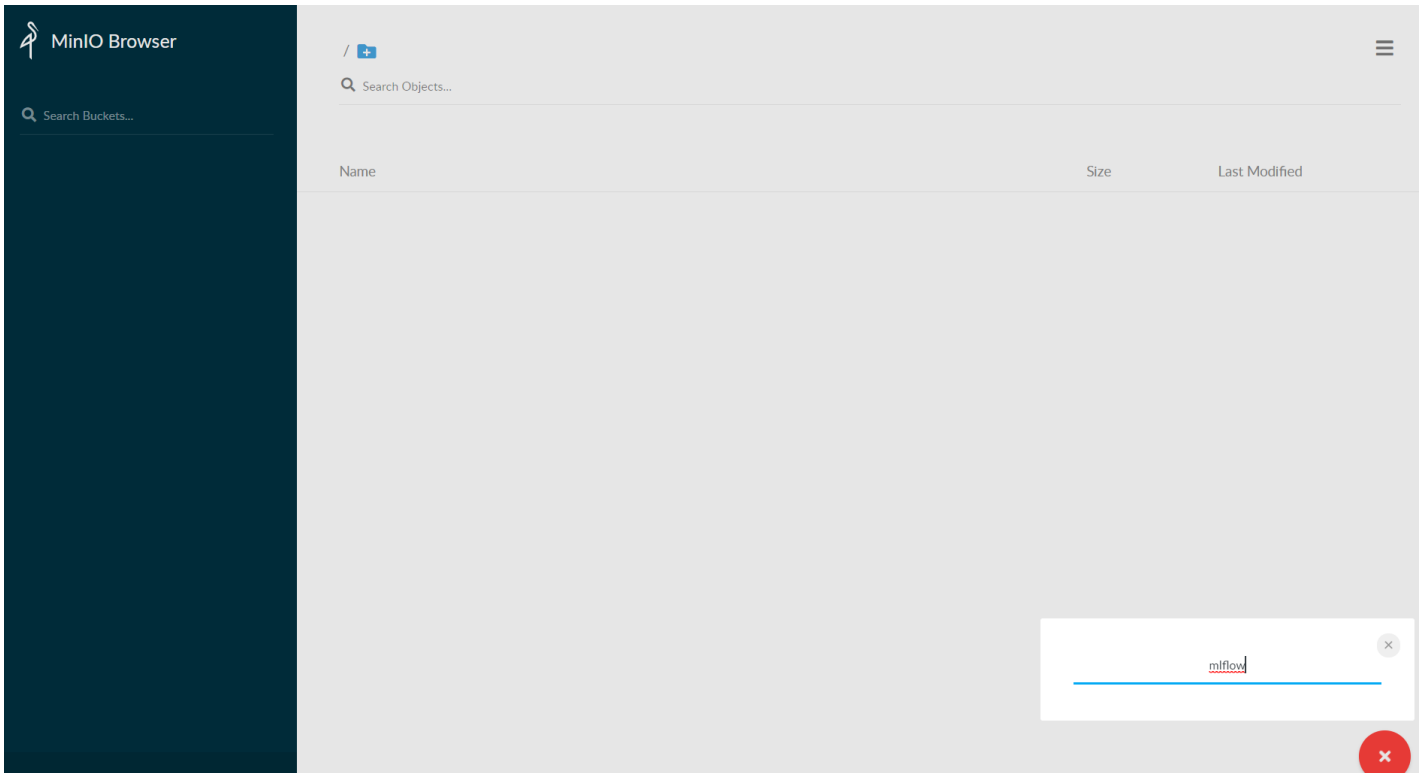


FIGURE 49. MinIO Browser

MLflow Server

Open the MLflow Server window.

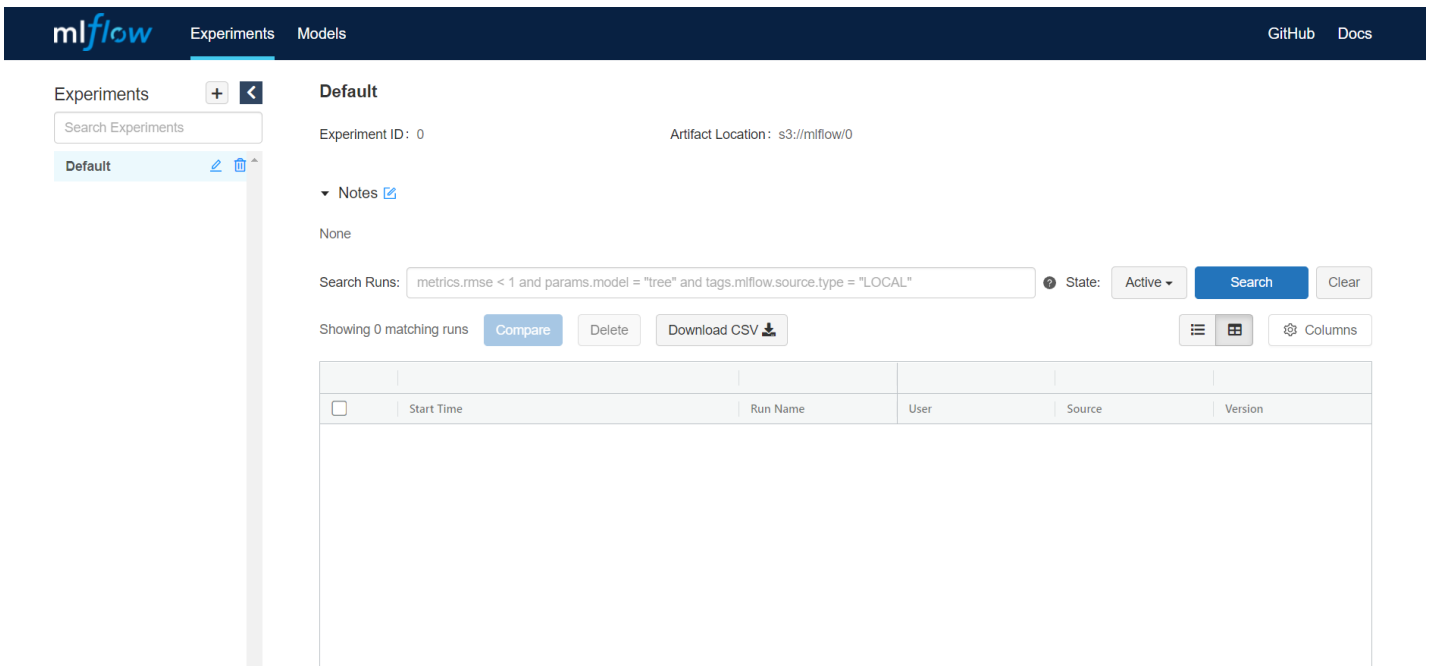


FIGURE 50. MLflow Server



Deploy KubeDirector Notebook. Go to the Notebook Tab and launch Jupyter Notebook.

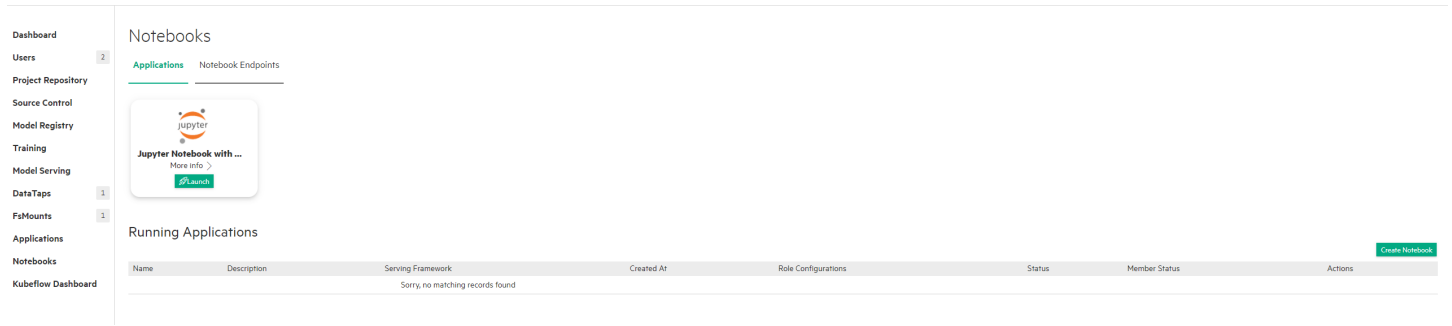


FIGURE 51. Launch Jupyter Notebook

Attach 'clusters' and 'secret' in Launch YAML of the Jupyter Notebook.

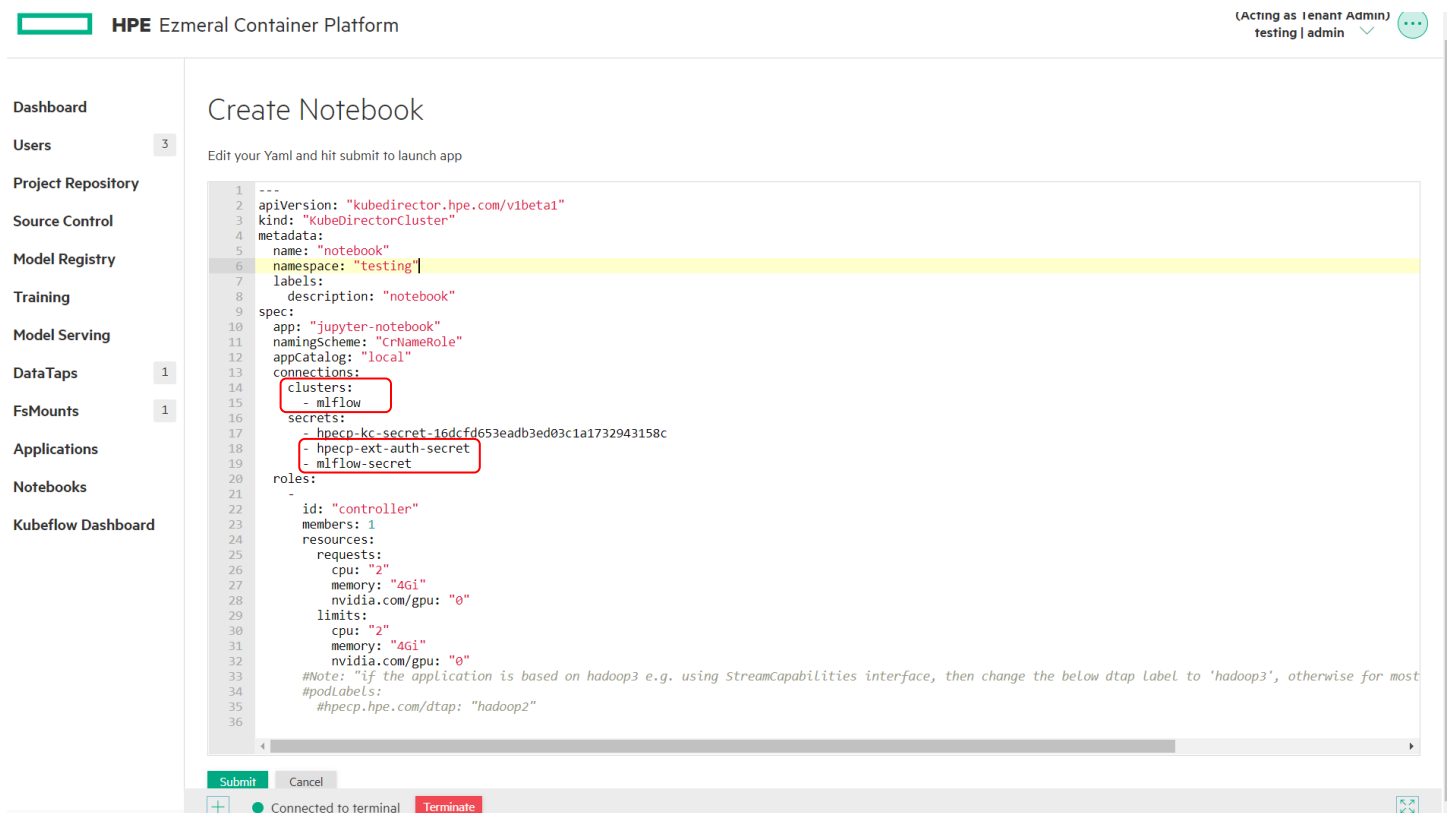


FIGURE 52. Edit Launch yaml

In Figure 52, **miflow** in clusters and **miflow-secret** in secrets are attached.

- **miflow** is the application name which deployed
- **miflow-secret** is the secret created



Launching Notebook with MLflow cluster attached from HPE Ezmeral UI has known issue, so notebook to be created from the terminal. Copy the above YAML file and create a file at the terminal and apply using kubectl, as shown below (workaround).

```

k8suser@kdss-1fp29-0:~$ ls
mlflow-secret.yaml notebook.yaml
k8suser@kdss-1fp29-0:~$ kubectl create -f notebook.yaml -n testing
kubedirectorcluster.kubedirector.hpe.com/notebook created
k8suser@kdss-1fp29-0:~$
    
```

Access the Notebook from the Notebook Endpoints tab.

HPE Ezmeral Container Platform (Acting as Tenant Admin) testing | admin

Notebooks

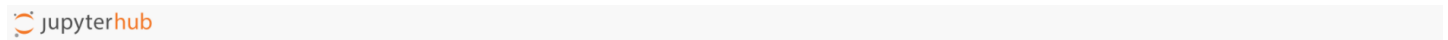
Applications **Notebook Endpoints**

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
notebook-controller-rv4vl-0	controller	KubeDirectorApp: ID: jupyter-notebook Name: Jupyter Notebook with ML toolkits	notebook	SSH Jupyter Notebook	22 8000	hpecp-531-vm2.rcc.local:10037 hpecp-531-vm2.rcc.local:10038	NodePort

FIGURE 53. Notebook Endpoints

Upload Notebook File and Dataset to run the experiment.

Login to the Jupyterhub.



Sign in

Warning: JupyterHub seems to be served over an unsecured HTTP connection. We strongly recommend enabling HTTPS for JupyterHub.

Username:

Password:

[Sign in](#)

FIGURE 54. Login to Jupyterhub



After login to Jupyter Notebook upload the care directory Bike-Sharing-MLFlow-UseCase and upload Bike-Sharing-MLFlow.ipynb, day.csv, and hour.csv from <https://github.com/SANDataHPE/MLFlow-Examples/tree/main/Bike-Sharing-MLFlow-Example>.

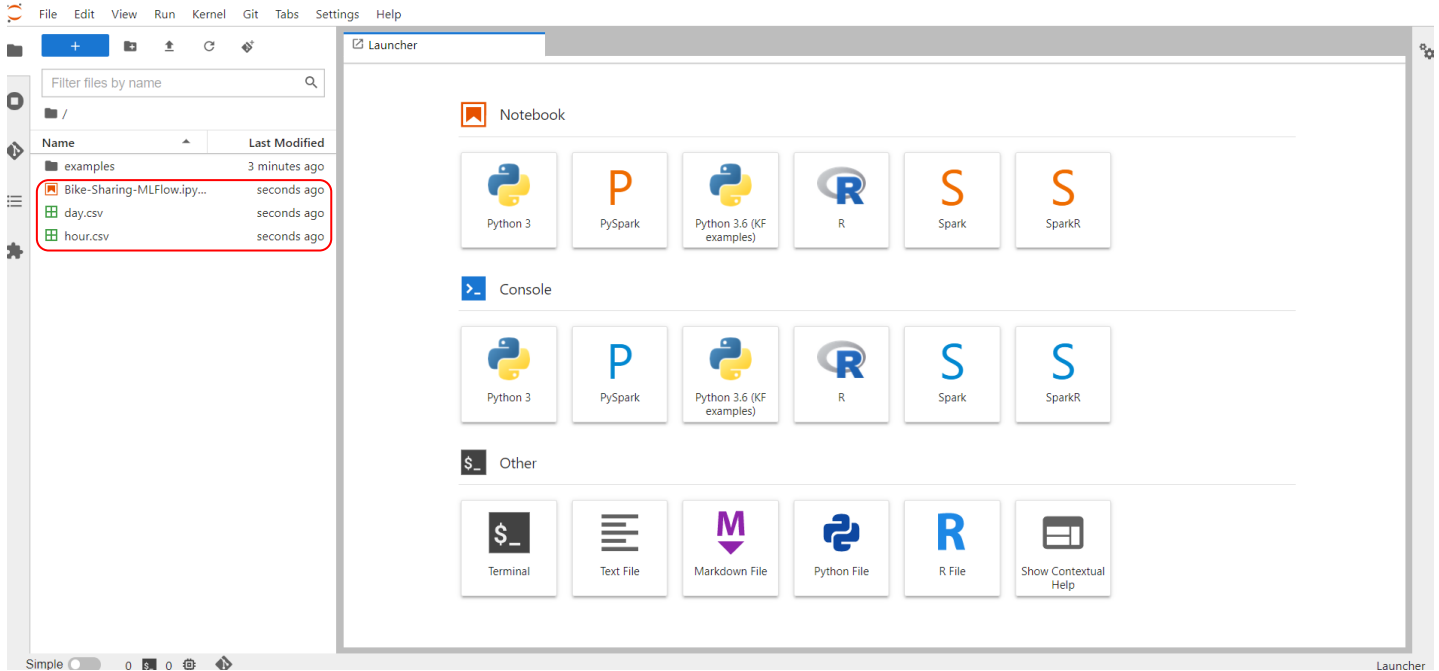


FIGURE 55. Jupyter Notebook

Setup 'ML Flow' magics.

- Add the user (which we've logged in to notebook with) to the platform and give it member access to the tenant.
- Run `%kuberRefresh` in one of the cells.

Setup ML Flow Magics

Load ML Flow Environment variable

```
[3]: %kuberRefresh
Kubeconfig is available to use
```

```
[1]: %loadMlflow
Backend configured
```

Set ML Flow Experiment Name

```
[2]: %Setexp --name bikesharetesting
```

- Load MLflow server and the secret to our notebook. Run `%loadmlflow`.



Execute cell in Notebook file.

- a. Each notebook cell contains a comment about the task which we are performing in each cell. Read the markdown before executing the cell. Follow the instruction written in the notebook before executing each cell.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import mlflow
import mlflow.sklearn
from mlflow import log_metric, log_param, log_artifact

from sklearn.ensemble import GradientBoostingRegressor
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import KFold, cross_val_score, train_test_split
from sklearn.inspection import permutation_importance
from mlflow.models.signature import infer_signature
from sklearn import tree

from pydotplus import graph_from_dot_data
import graphviz
from IPython.display import Image

import itertools

plt.style.use("fivethirtyeight")
pd.plotting.register_matplotlib_converters()

import warnings
warnings.filterwarnings('ignore')
bike_sharing = pd.read_csv("hour.csv")
bike_sharing
# remove unused columns
bike_sharing.drop(columns=["instant", "dteday", "registered", "casual"], inplace=True)

# use better names
bike_sharing.rename(
    columns={
        "yr": "year",
        "mnth": "month",
        "hr": "hour_of_day",
        "holiday": "is_holiday",
        "workingday": "is_workingday",
        "weathersit": "weather_situation",
        "temp": "temperature",
        "atemp": "feels_like_temperature",
        "hum": "humidity",
        "cnt": "rented_bikes",
    },
    inplace=True,
)

# show samples
bike_sharing
hour_of_day_agg = bike_sharing.groupby(["hour_of_day"])["rented_bikes"].sum()

hour_of_day_agg.plot(
    kind="line",
    title="Total rented bikes by hour of day",
    xticks=hour_of_day_agg.index,
```



```

    figsize=(15, 10),
]
# Split the dataset randomly into 70% for training and 30% for testing.
X = bike_sharing.drop("rented_bikes", axis=1)
y = bike_sharing.rented_bikes
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, test_size=0.3, random_state=42)

print(f"Training samples: {X_train.size}")
print(f"Test samples: {X_test.size}")
def rmse(y, y_pred):
    return np.sqrt(mean_squared_error(y, y_pred))

def rmse_score(y, y_pred):
    score = rmse(y, y_pred)
    print("RMSE score: {:.4f}".format(score))
    return score

def rmsle_cv(model, X_train, y_train):
    kf = KFold(n_splits=3, shuffle=True, random_state=42).get_n_splits(X_train.values)
    # Evaluate a score by cross-validation
    rmse = np.sqrt(-cross_val_score(model, X_train.values, y_train, scoring="neg_mean_squared_error", cv=kf))
    return rmse

def rmse_cv_score(model, X_train, y_train):
    score = rmsle_cv(model, X_train, y_train)
    print("Cross-Validation RMSE score: {:.4f} [std = {:.4f}]"
          .format(score.mean(), score.std()))
    return score

def model_feature_importance(model):
    feature_importance = pd.DataFrame(
        model.feature_importances_,
        index=X_train.columns,
        columns=["Importance"],
    )

    # sort by importance
    feature_importance.sort_values(by="Importance", ascending=False, inplace=True)

    # plot
    plt.figure(figsize=(12, 8))
    sns.barplot(
        data=feature_importance.reset_index(),
        y="index",
        x="Importance",
    ).set_title("Feature Importance")
    # save image
    plt.savefig("model_artifacts/feature_importance.png", bbox_inches='tight')

def model_permutation_importance(model):
    p_importance = permutation_importance(model, X_test, y_test, random_state=42, n_jobs=-1)

    # sort by importance
    sorted_idx = p_importance.importances_mean.argsort()[::-1]
    p_importance = pd.DataFrame(
        data=p_importance.importances[sorted_idx].T,
        columns=X_train.columns[sorted_idx]
    )

    # plot
    plt.figure(figsize=(12, 8))
    sns.barplot(

```



```
        data=p_importance,
        orient="h"
    ).set_title("Permutation Importance")

    # save image
    plt.savefig("model_artifacts/permutation_importance.png", bbox_inches="tight")

def model_tree_visualization(model):
    # generate visualization
    tree_dot_data = tree.export_graphviz(
        decision_tree=model.estimators_[0, 0], # Get the first tree,
        label="all",
        feature_names=X_train.columns,
        filled=True,
        rounded=True,
        proportion=True,
        impurity=False,
        precision=1,
    )

    # save image
    graph_from_dot_data(tree_dot_data).write_png("model_artifacts/Decision_Tree_Visualization.png")

    # show tree
    return graphviz.Source(tree_dot_data)
# Track params and metrics
def log_mlflow_run(model, signature):
    # Auto-logging for scikit-learn estimators
    # mlflow.sklearn.autolog()

    # log estimator_name name
    name = model.__class__.__name__
    mlflow.set_tag("estimator_name", name)

    # log input features
    mlflow.set_tag("features", str(X_train.columns.values.tolist()))

    # Log tracked parameters only
    mlflow.log_params({key: model.get_params()[key] for key in parameters})

    mlflow.log_metrics({
        'RMSE_CV': score_cv.mean(),
        'RMSE': score,
    })

    # log training loss
    for s in model.train_score_:
        mlflow.log_metric("Train Loss", s)

    # Save model to artifacts
    mlflow.sklearn.log_model(model, "model", signature=signature)

    # log charts
    mlflow.log_artifacts("model_artifacts")

    # misc
    # Log all model parameters
    mlflow.log_params(model.get_params())
    mlflow.log_param("Training size", X_test.size)
    mlflow.log_param("Test size", y_test.size)
# GBRT (Gradient Boosted Regression Tree) scikit-learn implementation
```



```
model_class = GradientBoostingRegressor
parameters = {
    "learning_rate": [0.1, 0.05, 0.01],
    "max_depth": [4, 5, 6],
    # "verbose": True,
}
# generate parameters combinations
params_keys = parameters.keys()
params_values = [
    parameters[key] if isinstance(parameters[key], list) else [parameters[key]]
    for key in params_keys
]
runs_parameters = [
    dict(zip(params_keys, combination)) for combination in itertools.product(*params_values)
]

# training loop
for i, run_parameters in enumerate(runs_parameters):
    print(f"Run {i}: {run_parameters}")

    # mlflow: stop active runs if any
    if mlflow.active_run():
        mlflow.end_run()
    # mlflow: track run
    mlflow.start_run(run_name=f"Run {i}")

    # create model instance
    model = model_class(**run_parameters)

    # train
    model.fit(X_train, y_train)

    # get evaluations scores
    score = rmse_score(y_test, model.predict(X_test))
    score_cv = rmse_cv_score(model, X_train, y_train)

    # generate charts
    model_feature_importance(model)
    plt.close()
    model_permutation_importance(model)
    plt.close()
    model_tree_visualization(model)

    # get model signature
    signature = infer_signature(model_input=X_train, model_output=model.predict(X_train))

    # mlflow: log metrics
    log_mlflow_run(model, signature)

    # mlflow: end tracking
    mlflow.end_run()
    print("")

best_run_df = mlflow.search_runs(order_by=['metrics.RMSE_CV ASC'], max_results=1)
if len(best_run_df.index) == 0:
    raise Exception(f"Found no runs for experiment '{experiment_name}'")

best_run = mlflow.get_run(best_run_df.at[0, 'run_id'])
best_model_uri = f"{best_run.info.artifact_uri}/model"
best_model = mlflow.sklearn.load_model(best_model_uri)
# print best run info
```




```

print("Best run info:")
print(f"Run id: {best_run.info.run_id}")
print(f"Run parameters: {best_run.data.params}")
print(f"Run score: RMSE_CV = {:.4f}".format(best_run.data.metrics['RMSE_CV']))
print(f"Run model URI: {best_model_uri}")
model_feature_importance(best_model)
model_permutation_importance(best_model)
model_tree_visualization(best_model)
test_predictions = X_test.copy()
# real output (rented_bikes) from test dataset
test_predictions["rented_bikes"] = y_test

# add "predicted_rented_bikes" from test dataset
test_predictions["predicted_rented_bikes"] = best_model.predict(X_test).astype(int)

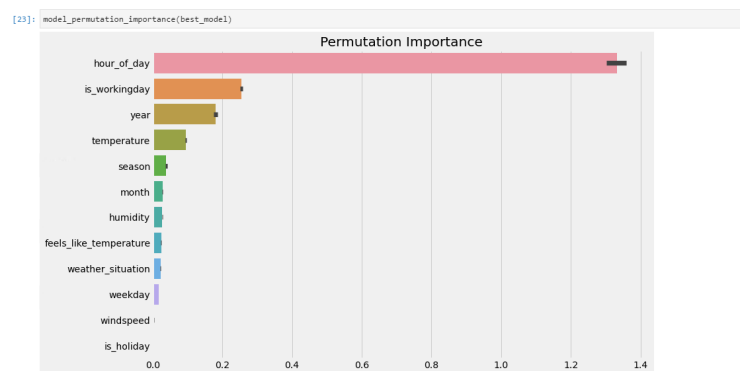
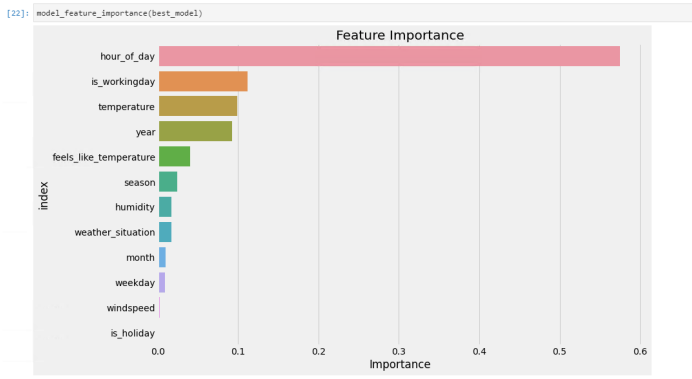
# show results
test_predictions
# plot truth vs prediction values
test_predictions.plot(
    kind="scatter",
    x="rented_bikes",
    y="predicted_rented_bikes",
    title="Rented bikes vs predicted rented bikes",
    figsize=(15, 15),
)

```

NOTE

We need to install all the required packages for the use case as shown below, if already installed this step can be skipped.

Visualize the best model.



Model Artifact in MLflow Server.

a. For the model artifact, we can access the MLflow Server.

The screenshot shows the MLflow Experiments page for an experiment named 'bike-sharing-use-case'. The interface includes a search bar for experiments, a list of experiments (Default and bike-sharing-use-case), and a search filter: 'metrics.rmse < 1 and params.model = "tree" and tags.mlflow.source.type = "LOCAL"'. The state is set to 'Active'. Below the search bar, there are buttons for 'Compare', 'Delete', and 'Download CSV'. A table displays 10 matching runs with the following columns: Start Time, Run Name, User, Source, Version, Test size, Training size, alpha, RMSE, RMSE_CV, Train Loss, estimator_n, and features. The runs are ordered by start time from latest to earliest.

Start Time	Run Name	User	Source	Version	Test size	Training size	alpha	RMSE	RMSE_CV	Train Loss	estimator_n	features
2021-12-09 13:41:18	Run 8	dev1	ipykerne	-	5214	62568	0.9	106.3	109.4	11827.7	Gradien...	[season...
2021-12-09 13:41:08	Run 7	dev1	ipykerne	-	5214	62568	0.9	112.5	116.1	13372.3	Gradien...	[season...
2021-12-09 13:40:59	Run 6	dev1	ipykerne	-	5214	62568	0.9	120.2	123.9	15396.3	Gradien...	[season...
2021-12-09 13:40:47	Run 5	dev1	ipykerne	-	5214	62568	0.9	46.3	49.83	1935.5	Gradien...	[season...
2021-12-09 13:40:38	Run 4	dev1	ipykerne	-	5214	62568	0.9	53.05	55.98	2916.2	Gradien...	[season...
2021-12-09 13:40:29	Run 3	dev1	ipykerne	-	5214	62568	0.9	63.15	67.82	4317.3	Gradien...	[season...
2021-12-09 13:40:18	Run 2	dev1	ipykerne	-	5214	62568	0.9	41.9	44.97	1323	Gradien...	[season...
2021-12-09 13:40:04	Run 1	dev1	ipykerne	-	5214	62568	0.9	44.69	48.18	1861.2	Gradien...	[season...
2021-12-09 13:39:49	Run 0	dev1	ipykerne	-	5214	62568	0.9	51.99	56.55	2765.7	Gradien...	[season...
2021-12-09 13:37:34	-	dev1	ipykerne	-	-	-	-	-	-	-	-	-

FIGURE 56. MLflow Experiments

MONITORING

There are different levels of monitoring with HPE Ezmeral ML Ops.

- At the platform level, the HPE Ezmeral Runtime provides dashboards that allow users to monitor resource utilization
- At the application level such as Kubeflow, containerized monitoring services such as Istio Prometheus are provided



Kubernetes administrator

Kubernetes users who have access to the site admin tenant can view the platform administrator dashboard which presents a high-level overview of the Kubernetes activity. Figure 57 shows the Usage tab displays usage information on a per-tenant basis. Refer to the HPE Ezmeral Container Platform 5.3 documentation, [Dashboard – Kubernetes Administrator](#) for more details.

Beginning with HPE Ezmeral Runtime 5.3, Dashboard views show additional GPU usage for Tesla-class or Quadro-class GPU families. For site administrators, the **Dashboard → Usage** tab shows the GPU devices used system-wide, while for Cluster Admin the stats would show usage cluster-wide. Finally, for tenant members, the resource usage statistics provide pod-specific GPU device information.

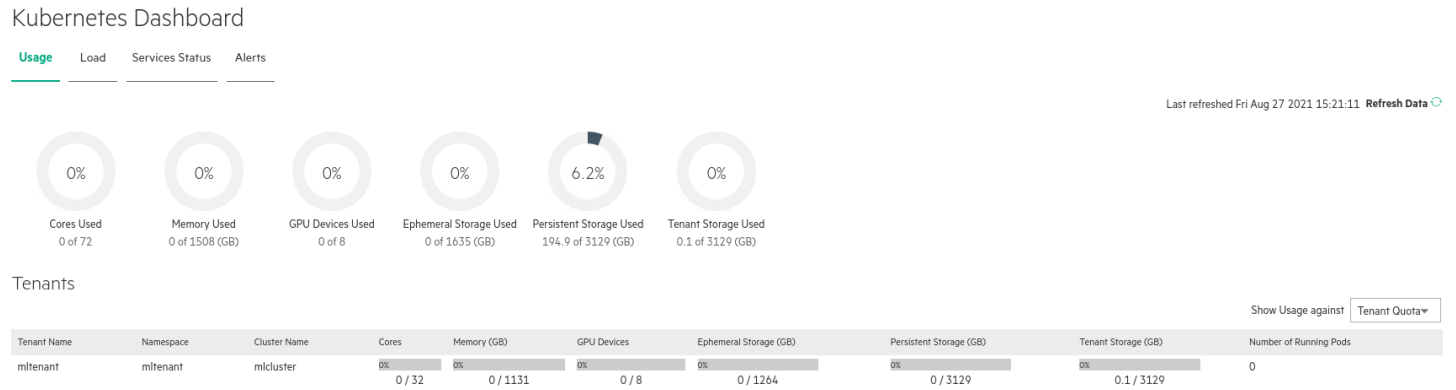


FIGURE 57. Platform administrator usage dashboard

Figure 59 shows the Load tab displays load statistics for the on-premises CPU, GPU, memory, and network resources within the K8s platform.

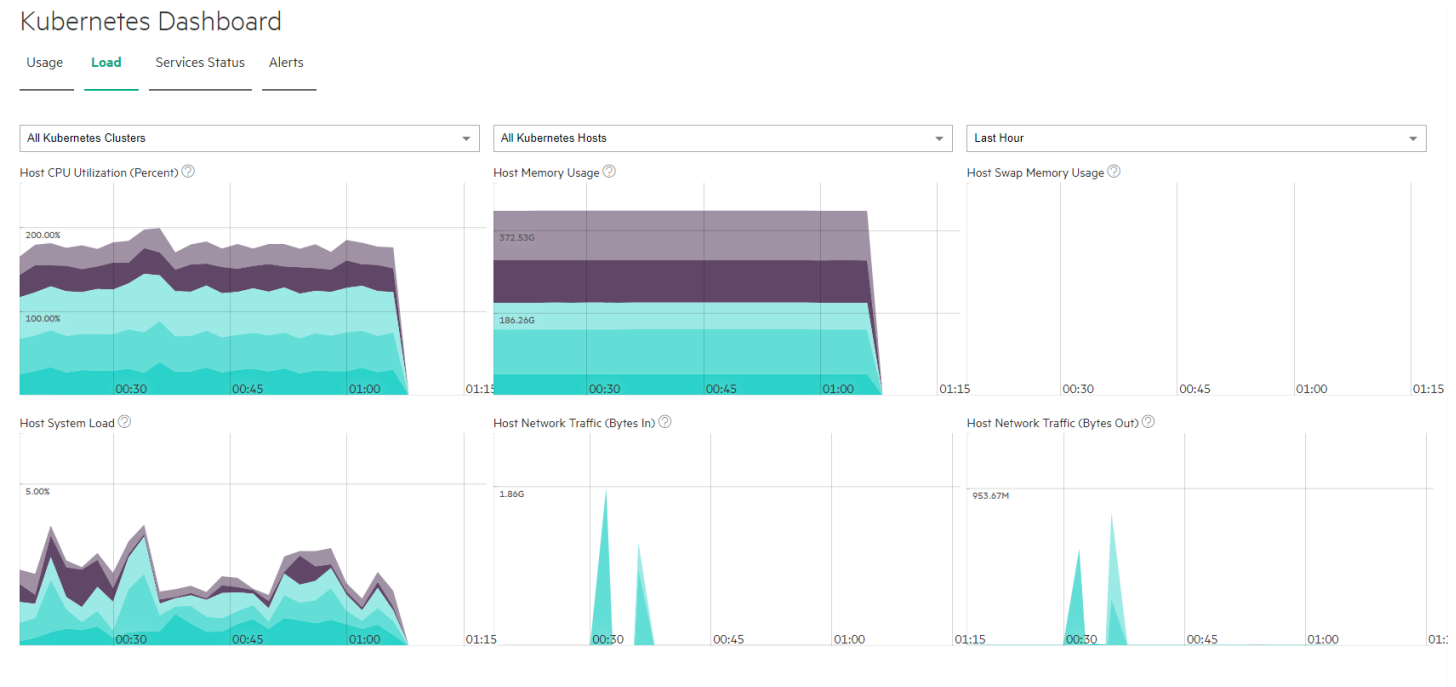


FIGURE 58. Platform administrator load dashboard



Kubernetes tenant/project administrator

HPE Ezmeral Runtime users who are logged into a Kubernetes tenant/project with the tenant/project administrator role can access the Kubernetes tenant/project administrator Dashboard. A tenant/project admin has access to the main menu and can view login details, alerts, etc. For more details, see the [Toolbar & Main Menu - Kubernetes Tenant Admin](#) page.

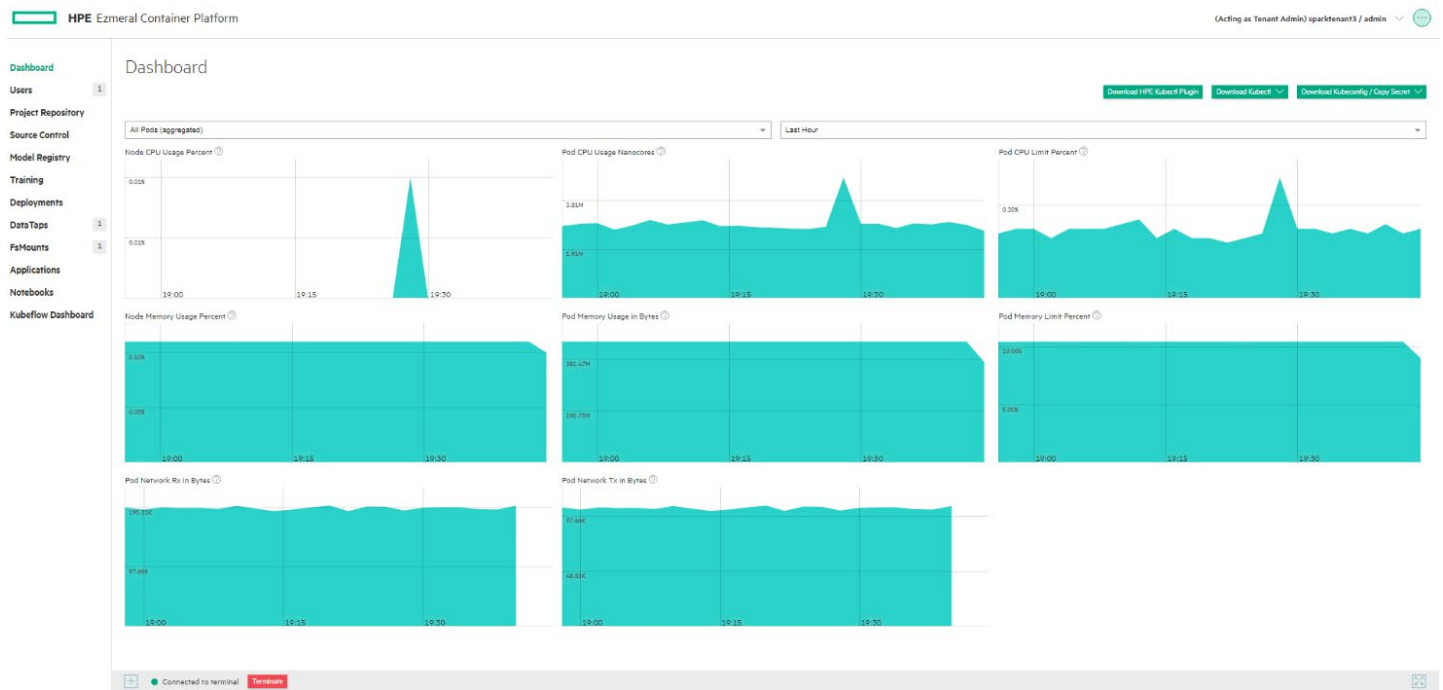


FIGURE 59. Tenant/project admin Dashboard

Istio and Prometheus

Prometheus is used for monitoring and Istio is used for managing network communication in the form of a service mesh.

Istio Prometheus use case

■ Create tunnel for 9090 port via SSH.

```
$ ssh -L 9090:localhost:9090 root@<master node>
```

```
k8suser@kdss-p9tw8-0:~$
k8suser@kdss-p9tw8-0:~$ ssh -L 9090:localhost:9090 root@172.24.2.4
The authenticity of host '172.24.2.4 (172.24.2.4)' can't be established.
ECDSA key fingerprint is SHA256:1W8mZtT7R8FzqUr2yLpHRp17Vm6UilbYeiYZeUBCy2k.
ECDSA key fingerprint is MD5:3d:36:80:3f:58:2c:c0:f3:5b:f3:65:db:e7:4f:a6:8b.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added '172.24.2.4' (ECDSA) to the list of known hosts.
root@172.24.2.4's password:
bind: Cannot assign requested address
Last login: Mon Aug 23 15:36:53 2021 from ilomxq231019r.perflab.hp.com
[root@ezam-04 ~]# kubectl port-forward svc/prometheus -n istio-system 9090:9090
Forwarding from 127.0.0.1:9090 -> 9090
Forwarding from [::1]:9090 -> 9090
```

Enable port forwarding.

```
$ kubectl port-forward svc/prometheus -n istio-system 9090:9090
```

On the master node, start up the Firefox, and go to \$ http://127.0.0.1:9090/.

Select **Status** → **Targets**, and then verify that all targets have been discovered and their statuses are being monitored.

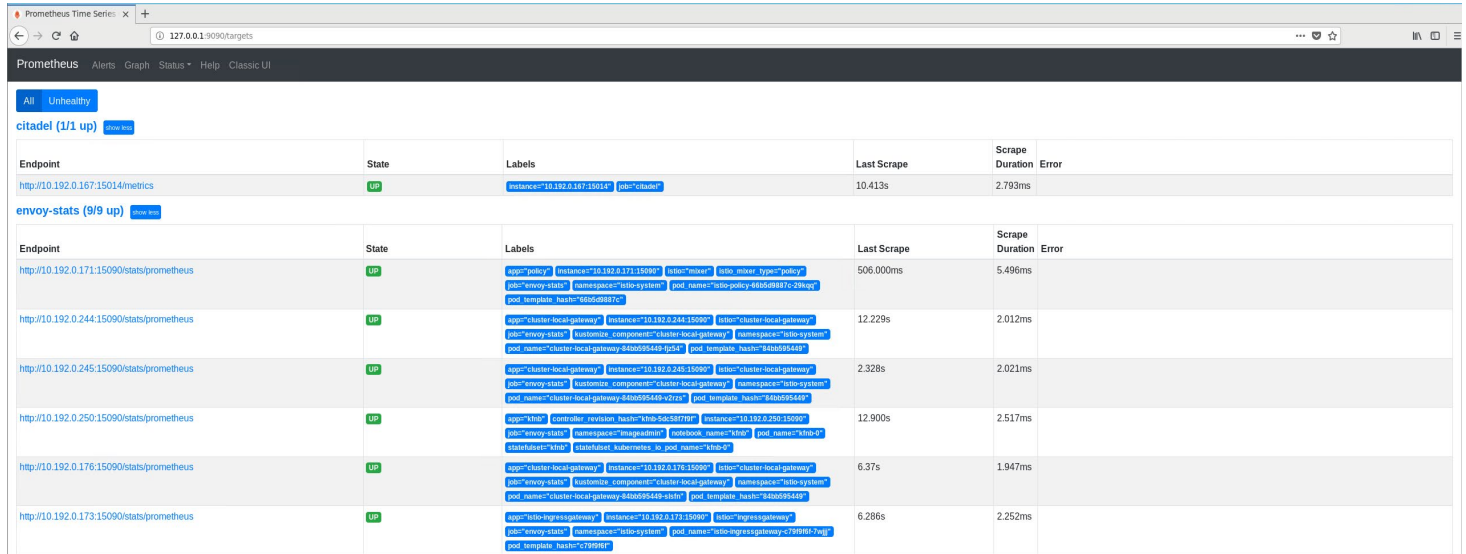


FIGURE 60. Targets

Select **Status** → **Service Discovery**, and then verify that all services have been discovered and their statuses are being monitored.

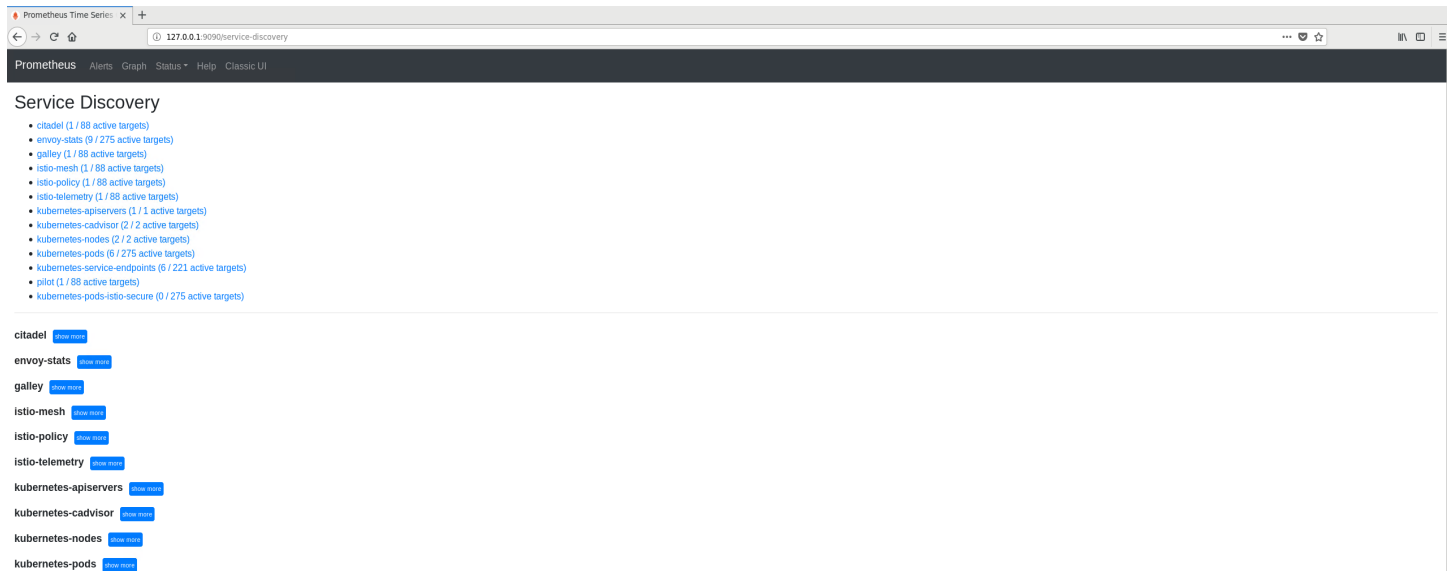


FIGURE 61. Service Discovery



Navigate to Graph.

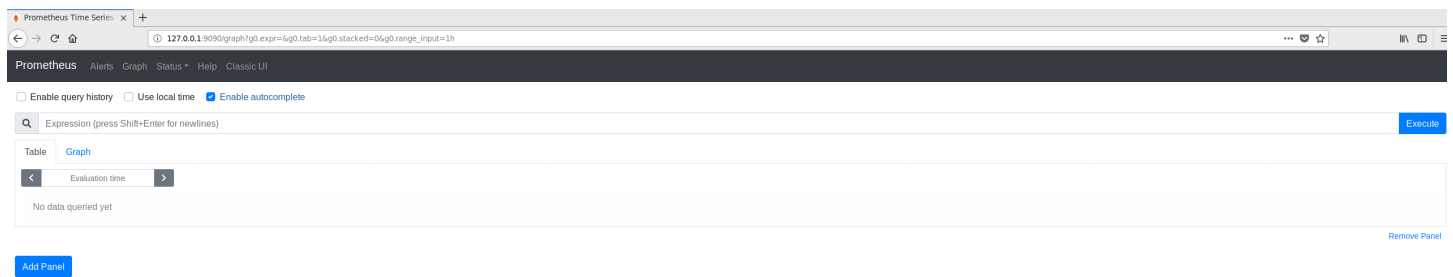


FIGURE 62. Graph

Choose any command you want and select it.

Click the Execute button and observe the Graph tab.

SUMMARY

Enterprises across all industries are embarking on a hybrid cloud journey for the development and deployment of their data-driven analytics and AI/ML applications. The continuous integration/continuous deployment (CI/CD) workflows, collectively referred to as DevOps, have become ubiquitous for software development today. On the machine learning front, data scientists still spend a significant amount of time and effort moving projects from development to production. Model version control is still largely manual, making it hard to update models in production. Code sharing is manual; data copied onto local storage leads to the variability of results between environments. There is a lack of standardization tools and frameworks, which makes it tedious and time-consuming to ensure the accuracy of predictions across all environments.

HPE Ezmeral Machine Learning Ops Software (HPE Ezmeral ML Ops) includes the capabilities and functionality of the HPE Ezmeral Runtime while also providing DevOps like agility to enterprise machine learning. With HPE Ezmeral ML Ops, enterprises can implement CI/CD workflows and standardize their ML pipelines. The HPE Ezmeral ML Ops software platform supports every stage of the machine learning lifecycle — supporting sandbox experimentation with a choice of ML/DL frameworks, integrations with model and code repositories, to deploying and tracking models in production.

HPE Ezmeral ML Ops gives data scientists and developers the ability to quickly and easily build and train machine learning models. HPE Ezmeral ML Ops allows data scientists to manage and track models built on any platform and deploy them into a scalable and secure production environment. Using HPE Ezmeral ML Ops, data scientists can spin-up containerized environments for distributed data processing, Machine Learning (ML), or Deep Learning (DL) in minutes rather than weeks. HPE Ezmeral ML Ops provides data science teams the flexibility to run their ML/DL workloads either on-premises, in multiple public clouds, or in a hybrid model and respond to dynamic business requirements in a variety of use cases.

With HPE Ezmeral ML Ops Software, Hewlett Packard Enterprise is making it easy for organizations to deliver a flexible and secure multitenant architecture, with the agility, flexibility, and performance needed to address evolving workload and application requirements. HPE Ezmeral ML Ops is deployed using pre-tested and optimized HPE Apollo building blocks on-premises, as well as in hybrid IT architectures and multi-cloud models.

Companies are driving digital transformation and investing in innovation to remain competitive. They are looking to deploy modern apps faster and simplify the production environment in a hybrid cloud architecture. They may have a mandate to move their application portfolio to the cloud or containers. Many organizations are still struggling to achieve these goals due to a lack of time and expertise. This Reference Architecture showcases the “lift-and-shift” application modernization use case which allows organizations to accelerate time-to-value by building a workable infrastructure the first time and every time.

With the HPE Ezmeral Runtime, enterprises now have a unified Kubernetes-based software solution for DevOps, CI/CD workflow, application modernization across hybrid cloud architecture, streamlining deployment, and operation with consistent orchestration and management. The platform acts as the control plane for container management and provides persistent container storage across multiple versions of open-source Kubernetes for container orchestration. The solution delivers a simpler, more scalable approach to modernizing applications. This is achieved using a scalable, code-driven container solution that, once assembled, can be configured within hours. This eliminates the complexities associated



with implementing a K8s container platform across an enterprise data center and provides the automation of hardware and software configuration to quickly provision and deploy a containerized environment at scale.

The solution provides customers with greater efficiency, higher utilization, and bare-metal performance by “collapsing the stack” and eliminating the need for virtualization. Developers have secured on-demand access to their environment. They can develop apps and release code faster, with the portability of containers to build once and deploy anywhere. IT teams can manage multiple Kubernetes clusters with multitenant container isolation and data access, for any workload, from edge to core to cloud. The benefits of containers, beyond cloud-native microservices-architected stateless applications, can be extended by providing the ability to containerize monolithic stateful analytic applications with persistent data.

The combination of HPE Ezmeral Runtime paired with HPE Apollo compute and HPE Nimble Storage delivers a composable architecture that can rapidly deploy modern containers supporting the new application framework. Ultimately, this results in faster digital transformation for businesses by helping organizations drastically increase the velocity of application development and accelerate innovation. This Reference Architecture provides an overview of an enterprise-grade solution that helps organizations increase agility, simplify operations, and deliver a cloud-like experience while offering a compelling return on investment.

APPENDIX A: KUBEFLOW AND TESTS OF USE CASES

For Kubeflow installation, see the HPE Ezmeral Container Platform 5.2 documentation [Kubeflow installation](#) page. For HPE Ezmeral Runtime and uninstallation, see [Uninstalling Kubeflow](#) page.

We have used [kubeflow_tutorials.zip](#), for Kubeflow use case testing which contains sample files for all of the included Kubeflow tutorials. The testing was done in a non-air-gapped environment.



Figure A1 shows Kubeflow architecture.

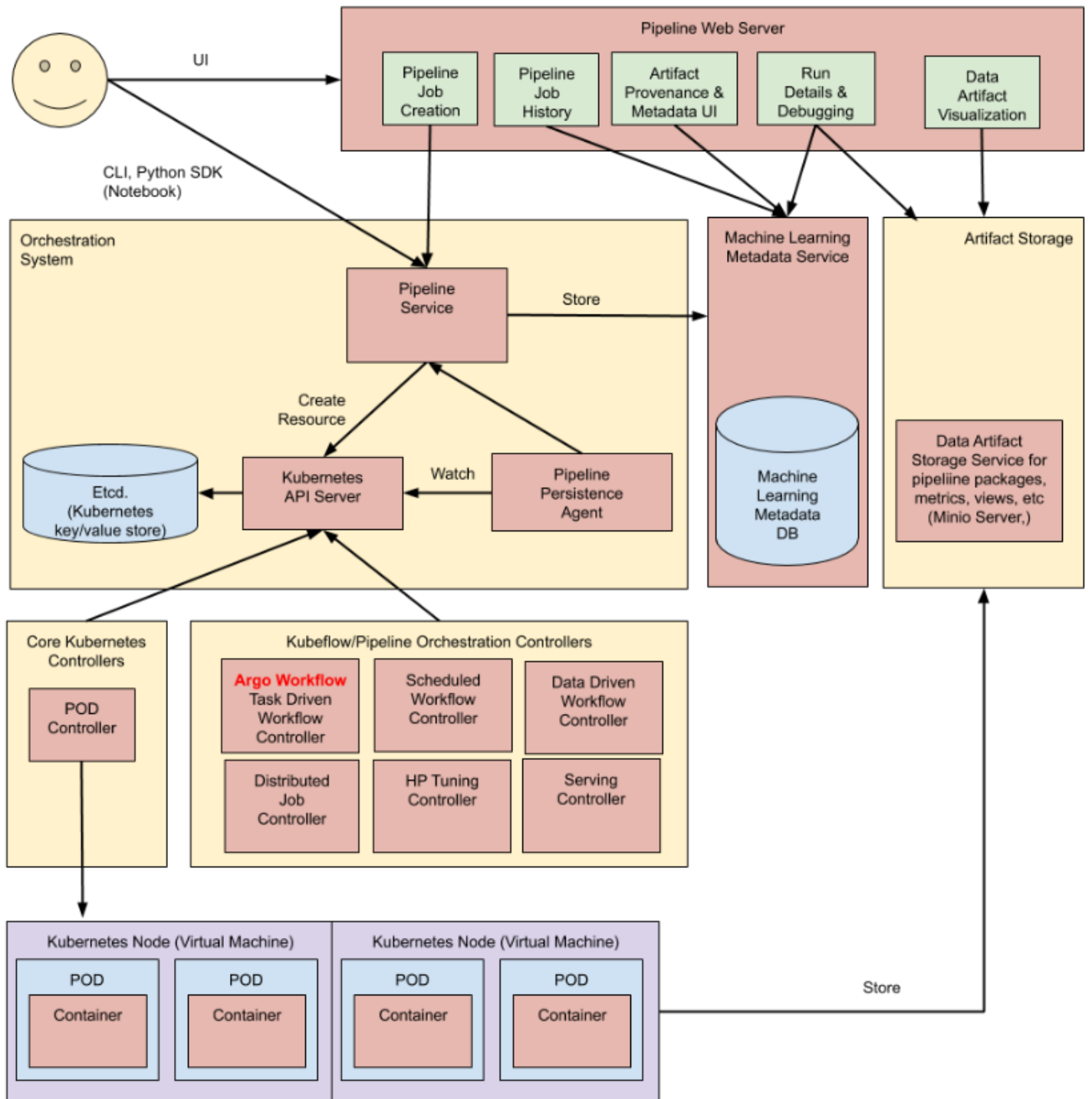


FIGURE A1. Kubeflow architectural diagram



KubeFlow components

The following table lists the components that Kubeflow can deploy.

TABLE A1. Kubeflow Components (Kubeflow Operator version 1.2)

Component	Version
Argo	2.3.0
Dex	2.22.0
Istio	1.3.1
Grafana	6.0.2
Jupyter Web Application	1.0.0
Katib	v1beta1
Kiali	1.4.0
Kfserving	0.3.0
Kubeflow Dashboard	1.0.0
ML Metadata	0.21.1
Notebook Controller	kf-ecp-5.3.0
Pipelines	1.0.4
PyTorch	1.0.0
Seldon	1.4.0
Spartakus	1.1.0
TensorFlow	1.0.0

Kubeflow components use cases GitHub issue summarization - Training with Jupyter

The steps followed in this [Tutorial: GitHub Issue Summarization](#) can be found in the official BlueData Documentation.

To begin with the Kubeflow examples:

■ Log in to HPE Ezmeral Runtime.

Sign in to

HPE Ezmeral Container Platform

Username

Password

Sign in

FIGURE A2. HPE Ezmeral Logging Page



Select the Tenant and navigate to Kubeflow Dashboard.

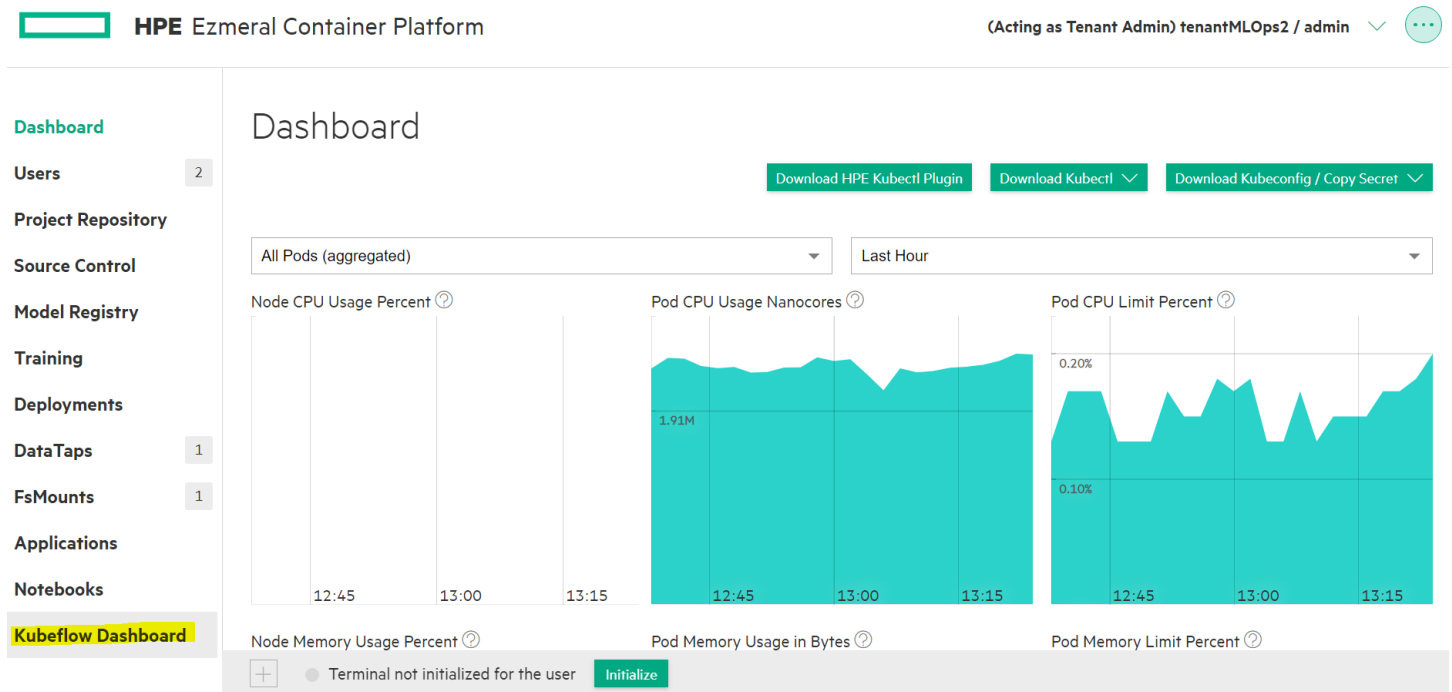


FIGURE A3. Navigating Kubeflow Dashboard

A list of Notebook Servers in Kubeflow Dashboard appears.

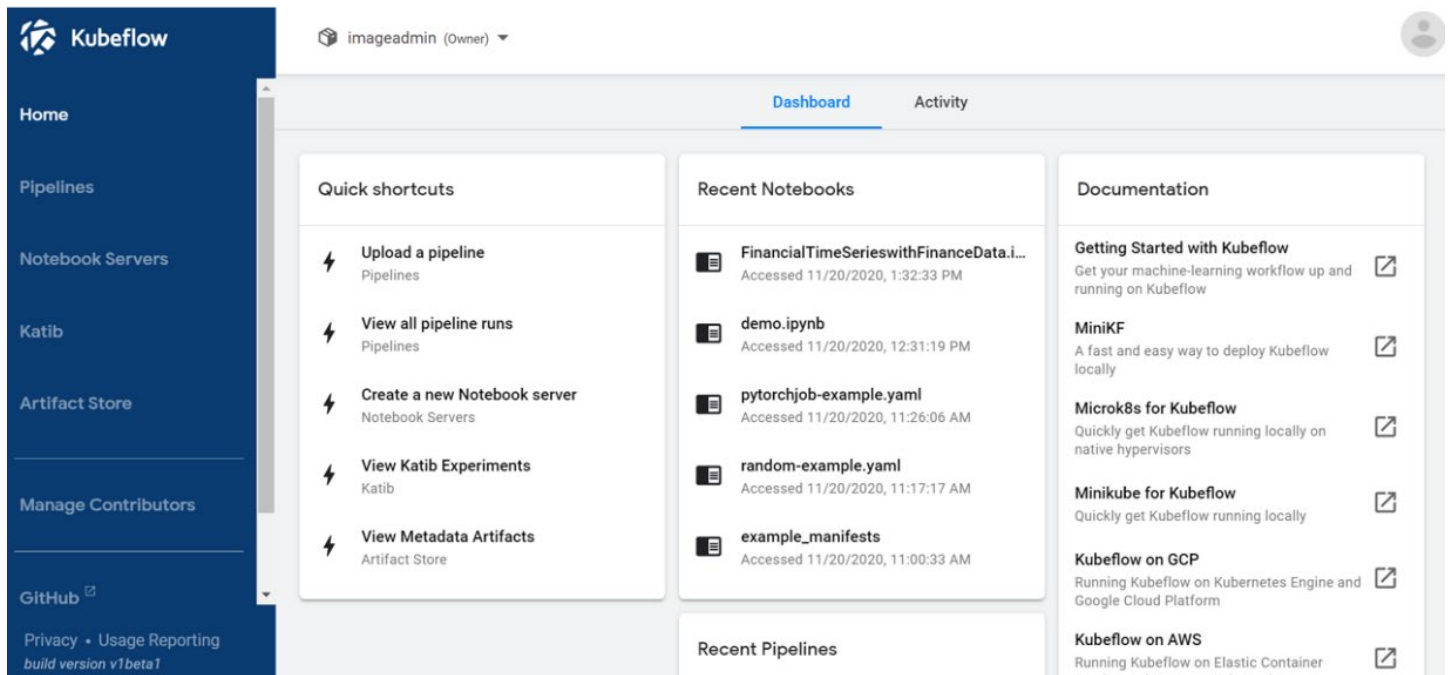


FIGURE A4. NB Servers in Kubeflow Dashboard



Select Notebook Servers and create one with the following instructions:

- In the Data Volume section, select ReadWriteMany.
- Change the Mount Path to a shorter name, such as /data.
- Leave the Workspace Volume as-is.
- Spawn the notebook by clicking Launch at the end of the page.

Name

Specify the name of the Notebook Server and the Namespace it will belong to.

Name	Namespace
nb-kf	imageadmin

Image

A starter Jupyter Docker Image with a baseline deployment and typical ML packages.

Custom Image

Image

gcr.io/mapr-252711/kf-ecp-5.3.0/tensorflow-2.1.0-notebook-cpu:latestC

CPU / RAM

Specify the total amount of CPU and RAM reserved by your Notebook Server. For CPU-intensive workloads, you can choose more than 1 CPU (e.g. 1.5).

CPU	Memory
10	16Gi

Workspace Volume

Configure the Volume to be mounted as your personal Workspace.

Don't use Persistent Storage for User's home

Type	Name	Size	Mode	Mount Point
New	workspace- nb-kf	10Gi	ReadWriteOnce	/home/jovyan

Data Volumes

Configure the Volumes to be mounted as your Datasets.

[+ ADD VOLUME](#)

Type	Name	Size	Mode	Mount Point
New	nb-kf-vol-1	10Gi	ReadWriteMany	/data

Configurations

Extra layers of configurations that will be applied to the new Notebook. (e.g. Insert credentials as Secrets, set Environment Variables.)

Configurations

FIGURE A5. NB Servers



- e. Connect to the Notebook Server created.

Notebook Servers

[+ NEW SERVER](#)

Status	Name	Age	Image	CPU	Memory	Volumes		
✓	imageadmin	3 days ago	tensorflow-1.15.2-notebook-cpu:1.0.0	5	10Gi		⋮	CONNECT 

FIGURE A6. NB Servers in Ready State

Open a new terminal from the Jupyter Hub and follow the steps below:

- a. Download the [kubeflow_tutorials.zip](#) file, containing sample files for all of the included Kubeflow tutorials.

```
$ pwd
/home/jovyan

$ wget kubeflow_tutorials.zip
```

- b. In a non-air-gapped environment only, execute the following commands to create the mapr-image-pull secrets and patch the Notebook. This must be done before pulling the Jupyter Notebook image for this tutorial.

```
$ kubectl apply -f imagepull-secrets.yaml
secret/mapr-imagepull-secrets created

$ kubectl patch serviceaccount default-editor -p '{"imagePullSecrets": [{"name": "mapr-imagepull-secrets"}]}'
serviceaccount/default-editor patched
```

- c. Connect to the Notebook, then open a new terminal, and then clone the Kubeflow examples repo:

```
$ git clone https://github.com/mapr/kubeflow-examples.git
```

- d. Return to the Jupyter folder list, and then open the file:

```
$ kubeflow-examples/github_issue_summarization/notebooks/Training.ipynb
```



The screenshot shows the Jupyter Notebook interface. At the top, there is a 'Quit' button. Below it, there are tabs for 'Files', 'Running', and 'Clusters'. A toolbar contains buttons for 'Duplicate', 'Rename', 'Move', 'Download', 'View', 'Edit', and a trash icon. On the right side of the toolbar, there are buttons for 'Upload', 'New', and a refresh icon. The main area is a file browser showing the path: /example_manifests / secrets / kubeflow-examples / github_issue_summarization / notebooks. The file browser has columns for 'Name', 'Last Modified', and 'File size'. The files listed are: '..' (seconds ago), 'test_data' (28 minutes ago), and 'Training.ipynb' (28 minutes ago, 15.9 kB). The 'Training.ipynb' file is selected with a blue checkmark.

FIGURE A7. Training Notebook



- e. In the Set path for the data_dir cell, change the data dir to:

```
%env DATA_DIR=/data
```

```
In [ ]: # Set path for data dir
%env DATA_DIR=/data
```

- f. In the pre-process data for the Deep Learning cell, comment out the magic function:

```
%%time
```

```
In [ ]: %%time
# Clean, tokenize, and apply padding / truncating such that each document length = 70
# also, retain only the top 8,000 words in the vocabulary and set the remaining words
# to 1 which will become common index for rare words
body_pp = processor(keep_n=8000, padding_maxlen=70)
train_body_vecs = body_pp.fit_transform(train_body_raw)
```

- g. To execute each of the cells in the Notebook, either:
- I. Click the Rerun button to run all the steps.
 - II. Click the Run button to manually run each step one at a time.
- h. After training completes, use the Notebook terminal to copy the files to the MapR file system:

```
$ cd kubeflow-examples/github_issue_summarization/notebooks/
$ cp *.h5 /data/
$ cp *.dpkl /data/
```

```
$ cd kubeflow-examples/github_issue_summarization/notebooks/
$ ls
body_pp.dpkl          seq2seq_model_tutorial.h5  train_title_vecs.npy
Dockerfile           seq2seq_utils.py          tutorial_seq2seq.epoch01-val7.19677.hdf5
Dockerfile.estimator test_data                  tutorial_seq2seq.epoch02-val6.43952.hdf5
environment_seldon_rest title_pp.dpkl             tutorial_seq2seq.epoch03-val5.95483.hdf5
github-issues-data   train_body_vecs.npy      tutorial_seq2seq.epoch04-val5.61074.hdf5
IssueSummarization.py trainer.py                 tutorial_seq2seq.epoch05-val5.37961.hdf5
Makefile             Training.ipynb            tutorial_seq2seq.epoch06-val5.19002.hdf5
__pycache__         train.py                  tutorial_seq2seq.epoch07-val5.10789.hdf5
requirements.txt     train_test.py             tutorial_seq2seq.log
$ cp *.h5 /data/
$ cp *.dpkl /data/
```

GitHub issue summarization – Serving with Seldon

Download the [kubeflow_tutorials.zip](#) file, which contains sample files for all of the included Kubeflow tutorials.

This tutorial uses the following image `idzikovsky/sandbox:seldon-issuesum`.

Apply the deployment by executing the following command.

```
$ kubectl apply -f seldon-issue-sum-deployment.yaml
Seldondeployment.machinelearning.seldon.io/issue-summarization created
```

Verify that the Seldon deployment was created.

```
$ kubectl get sdep
NAME          AGE
issue-summarization 2m54s
```



Verify that the pods are running.

```
$ kubectl get pods | grep classifier
issue-summarization-example-0-classifier-59fbc99c8-9zt6      2/2      Running 0      7m8s
```

Connect to the Notebook and upload the file `seldon-request.py` to send a sample request to the server model.

1 / example_manifests / scripts		Name ↓	Last Modified	File size
..			seconds ago	
<input type="checkbox"/>	<code>kfserving-request.py</code>		7 days ago	1.54 kB
<input checked="" type="checkbox"/>	<code>seldon-request.py</code>		7 days ago	1.29 kB

In the Notebook terminal, install the following Python dependencies.

```
$ pip install requests lxml --user
```

Execute `seldon-request.py` with the following options.

```
$ python seldon-request.py http://istio-ingressgateway.istio-system.svc.cluster.local:80 imageadmin hp123456
imageadmin issue-summarization
{"data":{"names":["t:0"],"ndarray":[["add support for for"]]},"meta":{}}
```

To delete the deployment, execute the following commands.

```
$ kubectl get sdep
NAME          AGE
issue-summarization 18m

$ kubectl delete sdep issue-summarization
Seldondeployment.machinelearning.seldon.io "issue-summarization" deleted
```

Training with TensorFlow (Financial series)

Before beginning this tutorial, download the [kubeflow_tutorials.zip](#) file, which contains sample files for all of the included Kubeflow tutorials.

Step 1: Mount the MapRFS Directory

To mount the MapRFS directory:

Obtain `pvc-tf-training-fin-series.yaml` from the zip file mentioned above for the Persistent Volume Claim (PVC).

Apply the `.yaml` file to create the PVC:

```
$ kubectl apply -f pvc-tf-training-fin-series.yaml
```

```
$ kubectl apply -f financial-series-tfjob.yaml
tfjob.kubeflow.org/trainingjob configured
$
```



Verify that the PVC was created and is in the bound state:

```
$ kubectl get pvc
```

```
$ kubectl get pvc
NAME                STATUS  VOLUME                                     CAPACITY  ACCESS MODES  STORAGECLASS
data                Bound   mapr-pv-e3686806-4252-404d-8f9a-0b00f85b7013  10Gi      RWX           hcp-mapr-cluster
imageadmin-vol-1   Bound   mapr-pv-d1c80ef1-1f69-46fe-a42c-d1fe93055d39  10Gi      RWX           hcp-mapr-cluster
kfnotebook1-vol-1  Bound   mapr-pv-a03cd16f-815a-4745-b1e7-14f87760c549  10Gi      RWO           hcp-mapr-cluster
pvcpy              Bound   mapr-pv-9002ba6a-13e3-4cd7-8622-25430d381ed2  5Gi       RWX           hcp-mapr-cluster
pvctf              Bound   mapr-pv-2feaa310-19f0-48ff-94c3-10e2a9b1a498  5Gi       RWX           hcp-mapr-cluster
voval-vol-1        Bound   mapr-pv-76d19390-61ce-4a1c-b8c6-b148bc092963  10Gi      RWO           hcp-mapr-cluster
```

Step 2: Exploration phase

To complete the exploration phase:

Open the Kubeflow web interface.

Follow the procedure based on the example described [here](#) (link opens an external website in a new browser tab/window) to spawn the Notebook. Do not create a Data Volume; simply select New Workspace and leave the remaining options set to their default values.

Connect to the Notebook server, and then open a terminal.

Execute the following command:

```
$ curl https://raw.githubusercontent.com/kubeflow/examples/master/financial_time_series/Financial%20Time%20Series%20with%20Finance%20Data.ipynb --output FinancialTimeSerieswithFinanceData.ipynb
```

```
$ curl https://raw.githubusercontent.com/kubeflow/examples/master/financial_time_series/Financial%20Time%20Series%20with%20Finance%20Data.ipynb --output FinancialTimeSerieswithFinanceData.ipynb
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload  Total   Spent    Left   Speed
100 47939  100 47939    0     0  178k    0  --:--:-- --:--:-- --:--:-- 178k
$
```

Open the uploaded Notebook.

	Name	Last Modified	File size
<input type="checkbox"/>	example_manifests	6 days ago	
<input type="checkbox"/>	kf-tutorials	9 days ago	
<input type="checkbox"/>	nat	8 days ago	
<input type="checkbox"/>	pipelines	12 days ago	
<input type="checkbox"/>	scripts	9 days ago	
<input type="checkbox"/>	secrets	9 days ago	
<input type="checkbox"/>	demo.ipynb	12 days ago	18.9 kB
<input type="checkbox"/>	financialSeries-tf-Training-working-example.ipynb	Running 9 days ago	3.13 MB
<input type="checkbox"/>	FinancialTimeSerieswithFinanceData.ipynb	15 minutes ago	47.9 kB
<input type="checkbox"/>	kfservingtest.py	9 days ago	1.63 kB

FIGURE A8. Financial Time Series with Finance Data Notebook

Walkthrough the Notebook step by step to better understand the problem and suggested solution(s).



Step 3: Training phase

To complete the training phase:

■ If required, use the already-pushed image `nskopiuk/sandbox:tensorflowimage-finseries`, else skip to step 3 if not.

■ Build and push the image.

```
$ git clone https://github.com/mapr/private-kubeflow-examples.git cd private-kubeflow-examples/financial_time_series/tensorflow_model export TRAIN_PATH= <your-docker-username>/sandbox:tensorflowimage-finseries docker build -t $TRAIN_PATH . docker push $TRAIN_PATH
```

■ Obtain the secrets file `user-gcp-sa` from the zip file mentioned above (required for this example).

```
$ kubectl create secret generic user-gcp-sa --from-file user-gcp-sa.json
```

■ Apply the TF-training job using the `.yaml` file from the zip file mentioned above.

```
$ kubectl apply -f financial-series-tfjob.yaml
```

```
$ kubectl apply -f financial-series-tfjob.yaml
tfjob.kubeflow.org/trainingjob configured
$
```

■ Verify that `user-gcp-sa` is mounted.

```
$ kubectl exec -it trainingjob-ps-0 -- /bin/sh
# ls /auth/
user-gcp-sa.json
```

■ Verify that the TF-job `trainingjob` was successfully created.

```
$ kubectl get tfjobs
```

```
$ kubectl get tfjobs
NAME          STATE          AGE
trainingjob  Succeeded     8d
$
```

■ Verify that pods were created, run, and completed.

```
$ kubectl get pods | grep trainingjob
```

```
$ kubectl get pods | grep trainingjob
trainingjob-ps-0          0/1    Completed    0    8d
trainingjob-worker-0     0/1    Completed    0    8d
$
```



Check the logs to see the training job description.

```
$ kubectl logs trainingjob-ps-0 tensorflow
```

```
$ kubectl logs trainingjob-ps-0 tensorflow
INFO:root:getting the ML model...
INFO:root:getting the data...
INFO:root:generating training data...
INFO:root:defining the training objective...
2020-11-23 22:03:46.400338: I tensorflow/core/platform/cpu_feature_guard.cc:140] Your CPU supports instructions that this TensorFlow binary was not compiled to use: SSE4.1 SSE4.2 AVX AVX2 AVX512F FMA
INFO:root:training the model...
INFO:root:training took 19.19 sec
INFO:root:validating model on test set...
INFO:root:Exporting model for tensorflow-serving...
INFO:tensorflow:Assets added to graph.
INFO:tensorflow:Assets added to graph.
INFO:tensorflow:No assets to write.
INFO:tensorflow:No assets to write.
INFO:tensorflow:SavedModel written to: b'model/1/saved_model.pb'
INFO:tensorflow:SavedModel written to: b'model/1/saved_model.pb'
INFO:root:copy files to /data/model/1
5000 0.5607639
10000 0.5755208
15000 0.5946181
20000 0.6145833
25000 0.6302083
30000 0.6449653
Precision = 0.9142857142857143
Recall = 0.2222222222222222
F1 Score = 0.35754189944134074
Accuracy = 0.6006944444444444
$
```

Serving a TensorFlow model with KFServing (Financial series)

Before beginning this tutorial, download the [kubeflow_tutorials.zip](#) file, which contains sample files for all of the included Kubeflow tutorials.

Step 1: Create and apply the YAML file

To complete the tutorial:

Obtain the serving YAML file from the zip file mentioned above.

Apply the file.

```
$ kubectl apply -f financial-series-serving.yaml
```

```
$ kubectl apply -f financial-series-serving.yaml
inferenceservice.serving.kubeflow.org/finance-sample unchanged
$
```

Verify that the inference service, revision, and relevant pods have been created.

K SVC:

```
$ kubectl get ksvc
```

```
$ kubectl get ksvc
NAME                                URL                                                                                               LATESTCREATED
finance-sample-predictor-default    http://finance-sample-predictor-default.imageadmin.example.com  finance-sample-predictor
-default-mhjwg                      finance-sample-predictor-default-mhjwg                        True
```



Revision:

```
$ kubectl get revision
```

```
$ kubectl get revision
NAME                                CONFIG NAME                                K8S SERVICE NAME                                GENERAT
ION  READY  REASON
finance-sample-predictor-default-mhjpg  finance-sample-predictor-default  finance-sample-predictor-default-mhjpg  1
True
```

Inference Services:

```
$ kubectl get inferencservices
```

```
$ kubectl get inferencservices
NAME                                URL                                READY  DEFAULT TRAFFIC  CANARY TR
AFFIC  AGE
finance-sample  http://finance-sample.imageadmin.example.com/v1/models/finance-sample  True  100
5d16h
```

Pods:

```
$ kubectl get pods | grep finance-sample
```

```
$ kubectl get pods | grep finance-sample
finance-sample-predictor-default-mhjpg-deployment-7549f689ftk  2/2  Running  0  5d16h
```

Verify that virtual services have been created.

```
$ kubectl get virtalservices | grep finance-sample
```

```
$ kubectl get virtalservices | grep finance-sample
finance-sample  [kubeflow-gateway.kubeflow knative-serving/cluster-local-gateway]  [finance-sample.
imageadmin.example.com finance-sample.imageadmin.svc.cluster.local]
5d16h
finance-sample-predictor-default  [knative-serving/cluster-local-gateway kubeflow/kubeflow-gateway]  [finance-sample-
predictor-default.imageadmin finance-sample-predictor-default.imageadmin.example.com finance-sample-predictor-default.imagea
dmin.svc finance-sample-predictor-default.imageadmin.svc.cluster.local]  5d16h
finance-sample-predictor-default-mesh  [mesh]  [finance-sample-
predictor-default.imageadmin finance-sample-predictor-default.imageadmin.svc finance-sample-predictor-default.imageadmin.svc
.cluster.local]
5d16h
```

Step 2: Perform inferences against the served model

To send a request to the model:

Obtain kfserving-request.py from the zip file mentioned above.

Install the following Python dependencies.

```
$ pip install requests lxml -user
```

Launch kfserving-request.py with the following options.

```
python kfserving-request.py <base_url> <login> <password> <profile_name>
```

For example:

```
$ python kfserving-request.py http://cv-hcp-lb1.qa.lab:10046 imageadmin 12341234 imageadmin
```

To send requests from the Jupyter Notebook terminal, use the ingressgateway address (istio-ingressgateway.svc.cluster.local). For example:

```
$ python kfserving-request.py http://istio-ingressgateway.istio-system.svc.cluster.local:80 imageadmin
12341234 imageadmin
```

The output will be similar to the following:

```
200
```

```
{'predictions': [{'model-version': '1', 'prediction': 0}]}
```



Training a PyTorch model (PyTorch MNIST)

Download the [kubeflow_tutorials.zip](#) file, which contains sample files for all of the included Kubeflow tutorials, if desired.

1. Obtain `vim pytorch-mnist-ddp-cpu.yaml` from the zip file mentioned above.

```
$ vim pytorch-mnist-ddp-cpu.yaml
```

■ Create the PyTorch Job.

```
$ kubectl apply -f pytorch-mnist-ddp-cpu.yaml
pytorchjob.kubeflow.org/pytorch-mnist-ddp-cpu created
```

■ Verify that the PyTorch job was created.

```
$ kubectl get pytorchjobs
NAME          STATE      AGE
pytorch-mnist-ddp-cpu  Succeeded  108s
```

■ Verify that relevant pods have been created.

```
$ kubectl get pods -l pytorch-job-name=pytorch-mnist-ddp-cpu
NAME                                READY STATUS  RESTARTS  AGE
pytorch-mnist-ddp-cpu-master-0      0/1   Completed  0         10m
pytorch-mnist-ddp-cpu-worker-0      0/1   Completed  0         10m
pytorch-mnist-ddp-cpu-worker-1      0/1   Completed  0         10m
pytorch-mnist-ddp-cpu-worker-2      0/1   Completed  0         10m
```

■ Inspect the logs to observe PyTorch training progress:

```
$ PODNAME=$(kubectl get pods -l pytorch-job-name=pytorch-mnist-ddp-cpu,pytorch-replica-type=master,pytorch-replica-index=0 -o name)
kubectl logs -f ${PODNAME}$
```

■ Verify the PyTorch Job pod statuses, and wait until all pods show the status **Completed**.

```
$ kubectl get pods -l pytorch-job-name=pytorch-mnist-ddp-cpu
```

■ Check the logs again to verify that the output contains the following information.

```
$ kubectl logs -f ${PODNAME}$
Using CUDA
Using distributed PyTorch with gloo backend
Downloading http://yann.lecun.com/exdb/mnist/train-images-idx3-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/train-labels-idx1-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz
Processing...
Done!
Train Epoch: 1 [58880/60000 (98%)]   loss=0.2060
Train Epoch: 1 [59520/60000 (99%)]   loss=0.0644

accuracy=0.9664
```

■ Use the `describe` command to check PyTorch job status.

```
$ kubectl describe pytorchjobs pytorch-mnist-ddp-cpu
```



The output should look similar to the following:

Type	Reason	Age	From	Message
Normal	SuccessfulCreatePod	16m	pytorch-operator	Created pod: pytorch-mnist-ddp-cpu-worker-0
Normal	SuccessfulCreatePod	16m	pytorch-operator	Created pod: pytorch-mnist-ddp-cpu-worker-1
Normal	SuccessfulCreatePod	16m	pytorch-operator	Created pod: pytorch-mnist-ddp-cpu-worker-2
Normal	SuccessfulCreatePod	16m	pytorch-operator	Created pod: pytorch-mnist-ddp-cpu-master-0
Normal	SuccessfulCreateService	16m	pytorch-operator	Created service: pytorch-mnist-ddp-cpu-master-0
Normal	PyTorchJobSucceeded	5m57s	pytorch-operator	PyTorchJob pytorch-mnist-ddp-cpu is successfully completed.

Additional information

After training is complete, the files will be located in the MapR file system.

```
$ ssh root@<controller_node_ip_address>
bdmapr --root bash
hadoop fs -ls <volumePath>/pytorch/model
```

Result:

Found 1 items
-rw-r--r-- 3 root 88548 2020-07-15 09:36 /mapr-csi/k8s-10--siaitqgucc/pytorch/model/model_cpu.dat

Sample Pipeline in the pipelines interface

Refer to an example from the [Kubeflow documentation](#) (The link opens an external website in a new browser tab/window).

Open the Kubeflow dashboard (see [Accessing the Kubeflow Dashboard](#)), and then select Pipelines.

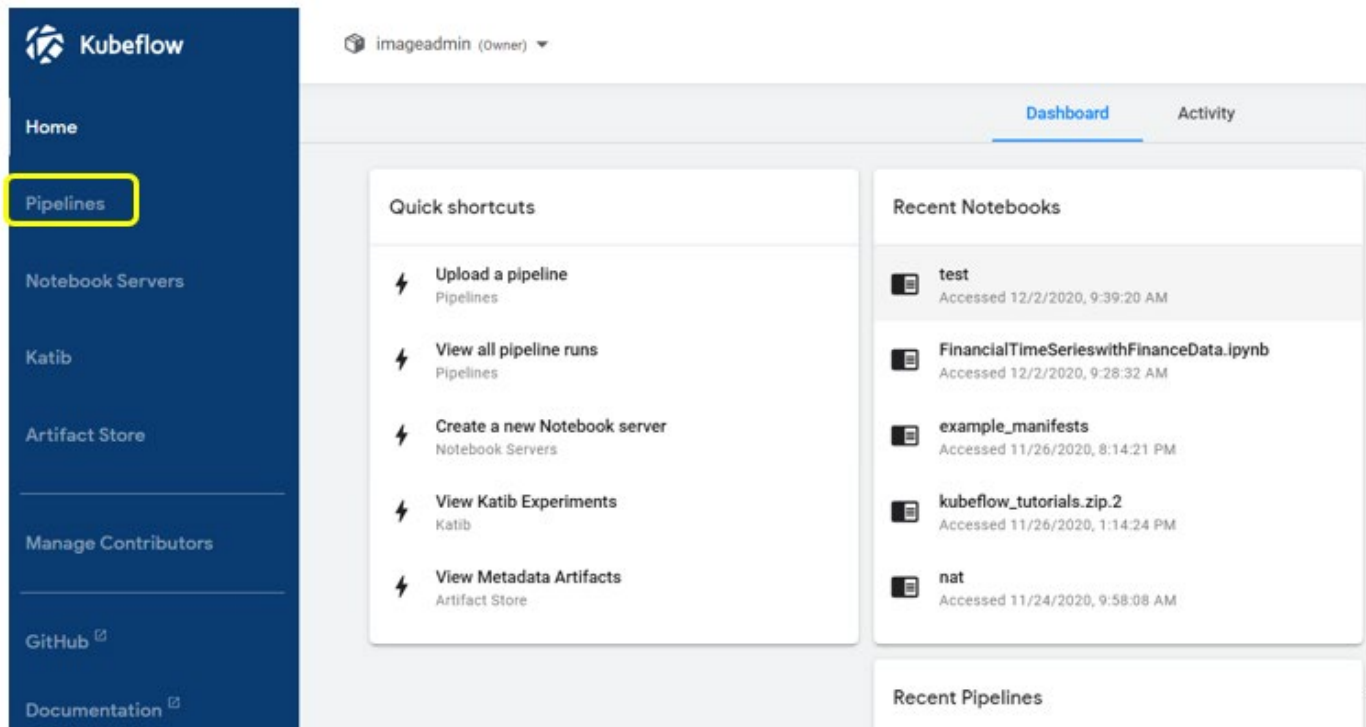


FIGURE A9. Kubeflow Pipeline

Click the sample name [Tutorial] DSL- Control Structures.



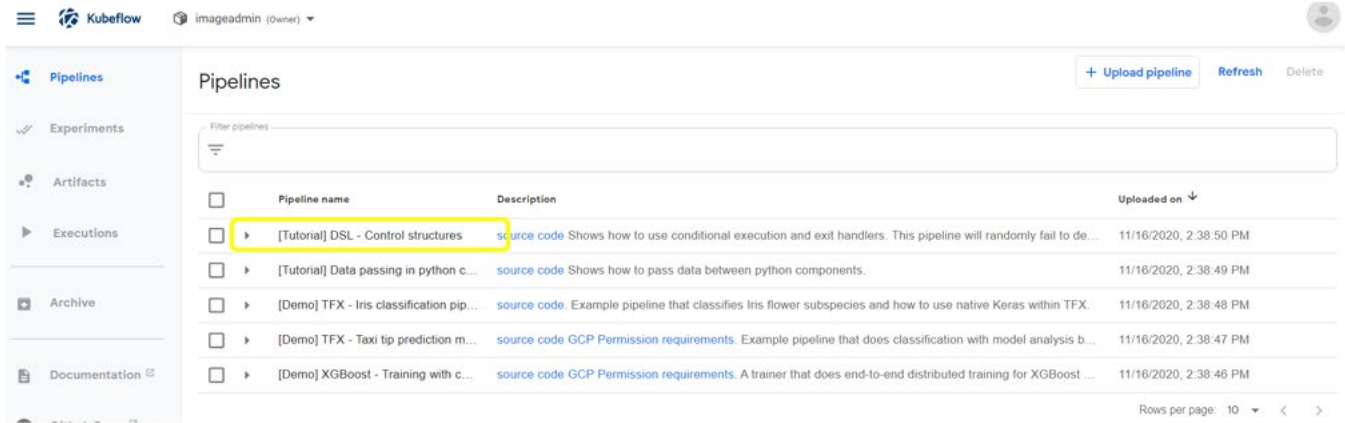


FIGURE A10. Kubeflow Pipeline Dashboard

Click **Experiments**, and then follow the on-screen prompts.

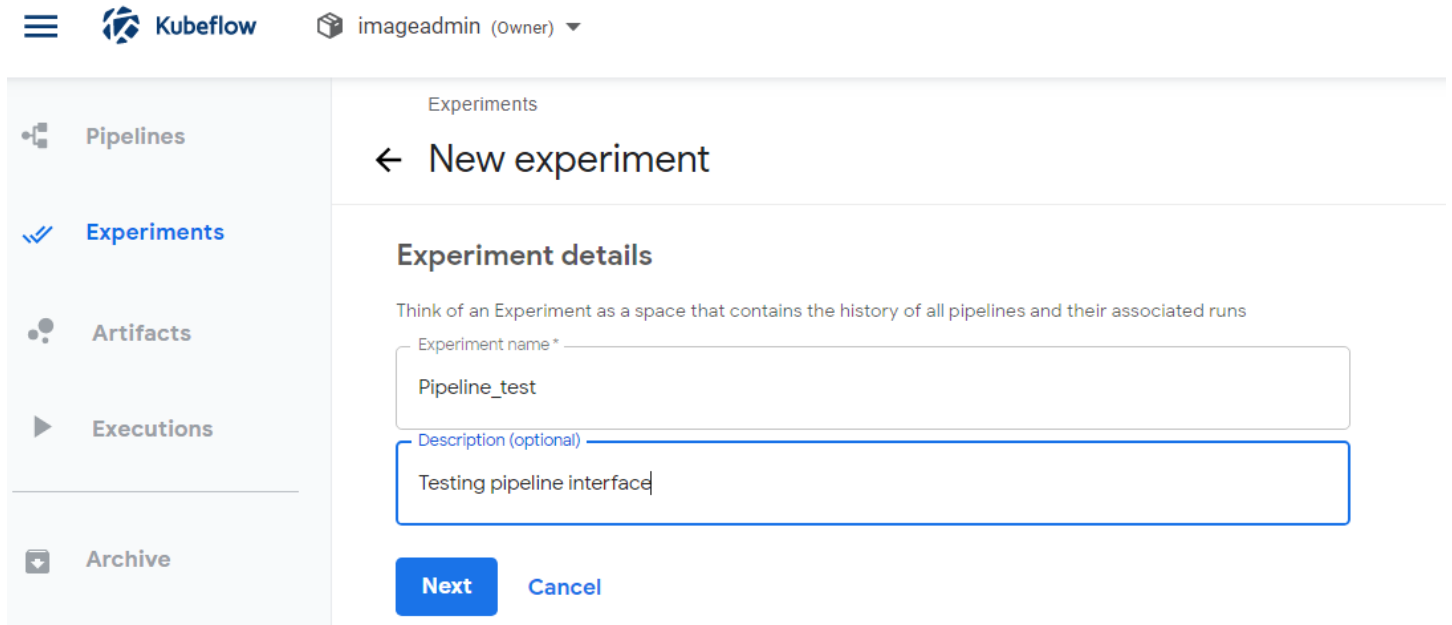


FIGURE A11. Create new Experiment



Enter the Run details and click **Start**.

Run details

Pipeline *
[Tutorial] DSL - Control structures [Choose](#)

Pipeline Version *
[Tutorial] DSL - Control structures [Choose](#)

Run name *
Run of [Tutorial] DSL - Control structures (b07c3) Run of [Tutorial] DSL - Control structures (b07c3)

Description (optional)

This run will be associated with the following experiment

Experiment *
Pipeline_test [Choose](#)

This run will use the following Kubernetes service account. [?](#)

Service Account (Optional)

Run Type

One-off Recurring

Run parameters

This pipeline has no parameters

[Start](#) [Skip this step](#)

FIGURE A12. Create Run

Select the run that was created in step 4 on the Experiments dashboard.

Experiments Pipeline_test [Refresh](#) [Archive](#)

Recurring run configs
0 active [Manage](#)

Experiment description [?](#)
Testing pipeline interface

Runs [+ Create run](#) [+ Create recurring run](#) [Compare runs](#) [Clone run](#) [Archive](#)

Filter runs

<input type="checkbox"/>	Run name	Status	Duration	Pipeline Version	Recurring Run	Start time ↓
<input type="checkbox"/>	Run of [Tutorial] DSL - Control structures (b07c3)	?	-	[Tutorial] DSL - Control structures	-	12/3/2020, 12:41:41 AM

Rows per page: 10 < >

FIGURE A13. Select Run on Experiments Dashboard



Explore the graph and other aspects of the run by clicking the graph components and other interface elements.

Experiments > Pipeline_test

← ! Run of [Tutorial] DSL - Control structures (b07c3)

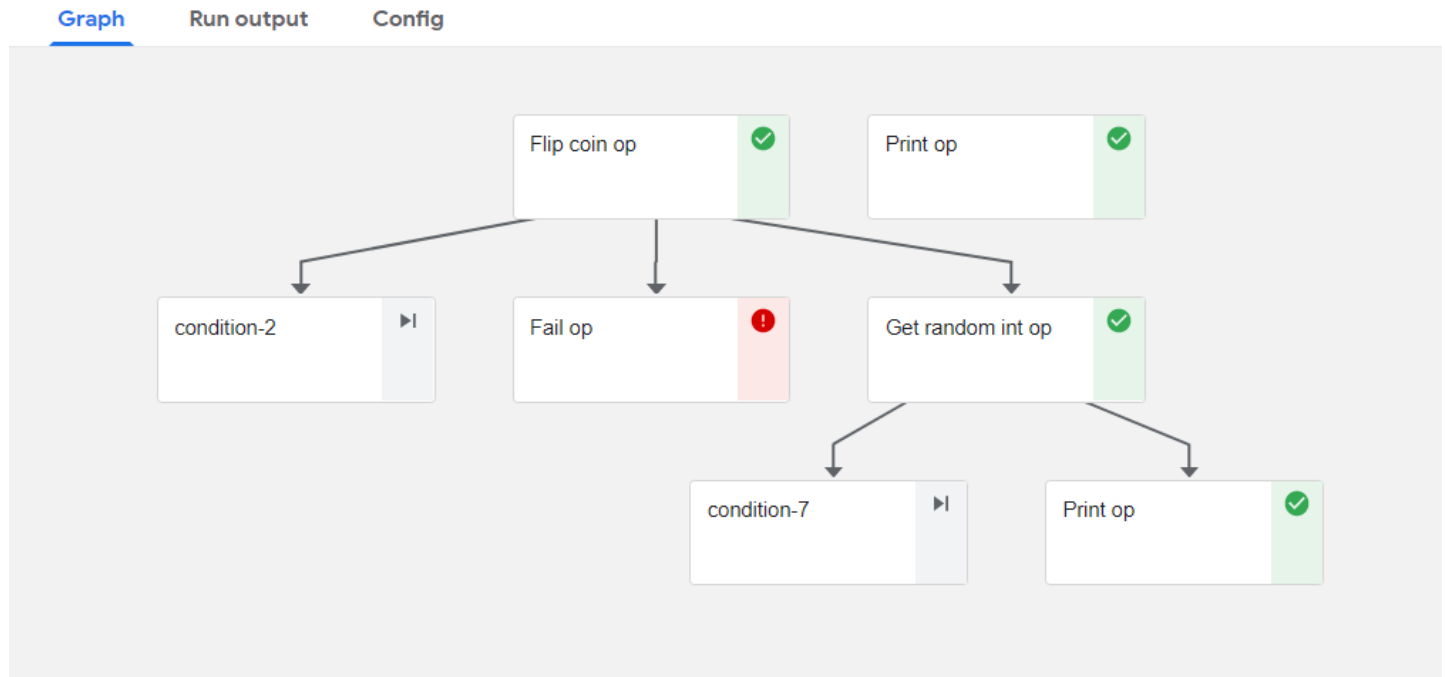


FIGURE A14. Run Graph

Running a pipeline in Jupyter Notebook

Before beginning this tutorial, download the [kubeflow_tutorials.zip](#) file, which contains sample files for all of the included Kubeflow tutorials.

■ Create a Jupyter Notebook.

■ Connect to the Notebook, and then click **New Terminal**.

■ Clone the kubeflow/pipelines repo:

```
$ git clone https://github.com/kubeflow/pipelines.git
```



Return to the Files tab, and then open the Notebook:

```
$ pipelines/samples/core/lightweight_component.ipynb
```

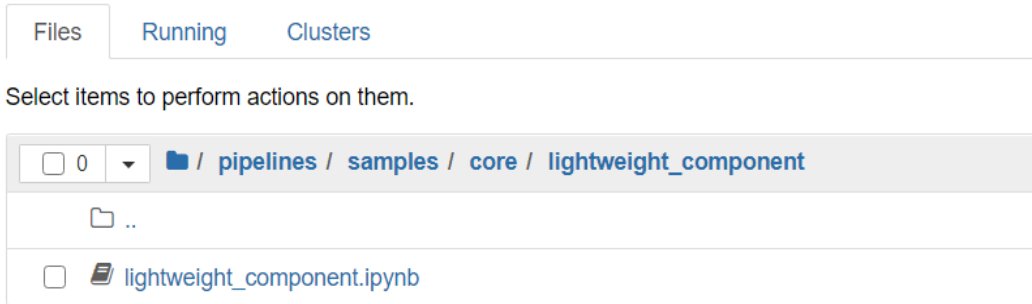


FIGURE A15. Lightweight Component Notebook

Execute each cell in the Notebook until it is finished.

Follow the Notebook link to the created experiment in the Pipelines interface.



Follow the Notebook link to the created run in the Pipelines interface.

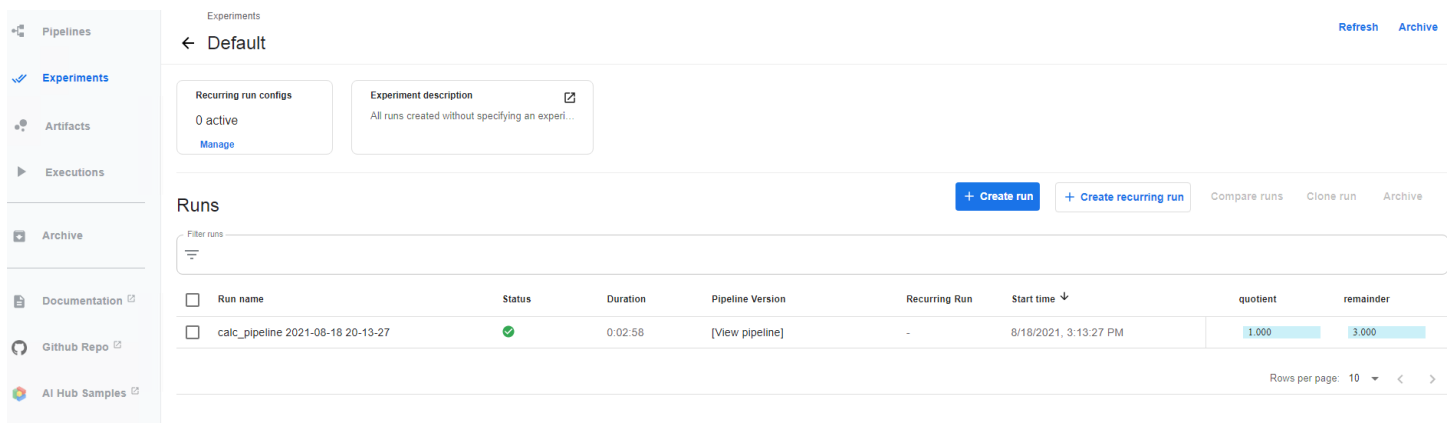


FIGURE A16. Calc Pipeline Run



Experiments > Default

← ✔ calc_pipeline 2021-08-18 20-13-27

Graph Run output Config

Run details

Status	Succeeded
Description	
Created at	8/18/2021, 3:13:26 PM
Started at	8/18/2021, 3:13:27 PM
Finished at	8/18/2021, 3:16:25 PM
Duration	0:02:58

Run parameters

a	7
b	8
c	17

FIGURE A17. Run Details

Here are some general pipeline viewing steps from the Pipelines Dashboard interface:

■ Open the **Experiments** page in the Pipelines dashboard.

FIGURE A18. Experiments Dashboard



In the **All experiments** tab, expand the Default group, and then view the pipeline graph and details per step by clicking the appropriate (view pipeline) link.

In the **All runs** tab, click the name of the run to view the **Graph**, **Run output**, and **Config** tabs.

The screenshot displays the Kubeflow Pipelines interface. On the left is a navigation sidebar with options like Pipelines, Experiments, Artifacts, Executions, Archive, Documentation, Github Repo, and AI Hub Samples. The main area shows the 'Graph' tab for a pipeline run named 'calc_pipeline 2021-08-18 20-13-27'. The graph consists of three steps: 'Add', 'My divmod', and another 'Add' step, all marked with green checkmarks. To the right, the 'Run output' tab is open, showing details for the run 'calculation-pipeline-2p65r-2935576819'. It includes sections for Input parameters (add-Output: 11.0, b: 8), Input artifacts, Output parameters (my-divmod-quotient: 1.0), and Output artifacts (mipipeline-ui-metadata, mipipeline-metrics, my-divmod-quotient, my-divmod-remainder), each with a corresponding artifact URL.

FIGURE A19. Pipeline Run Graph

Katib Hyperparameter Tuning

Example 1: TensorFlow

Download the [Kubeflow tutorials.zip](#) file which contains sample files for all of the included Kubeflow tutorials.

Edit the `tensorflow-example.yaml` to put the following on the pod template.

```
metadata:
  annotations:
    sidecar.istio.io/inject: "false"
```

Deploy the example.

```
$ kubectl apply -f tensorflow-example.yaml
```

Open the Kubeflow Dashboard, and then select Katib.



Click the left menu button, and then go to **HP** → **Monitor**.

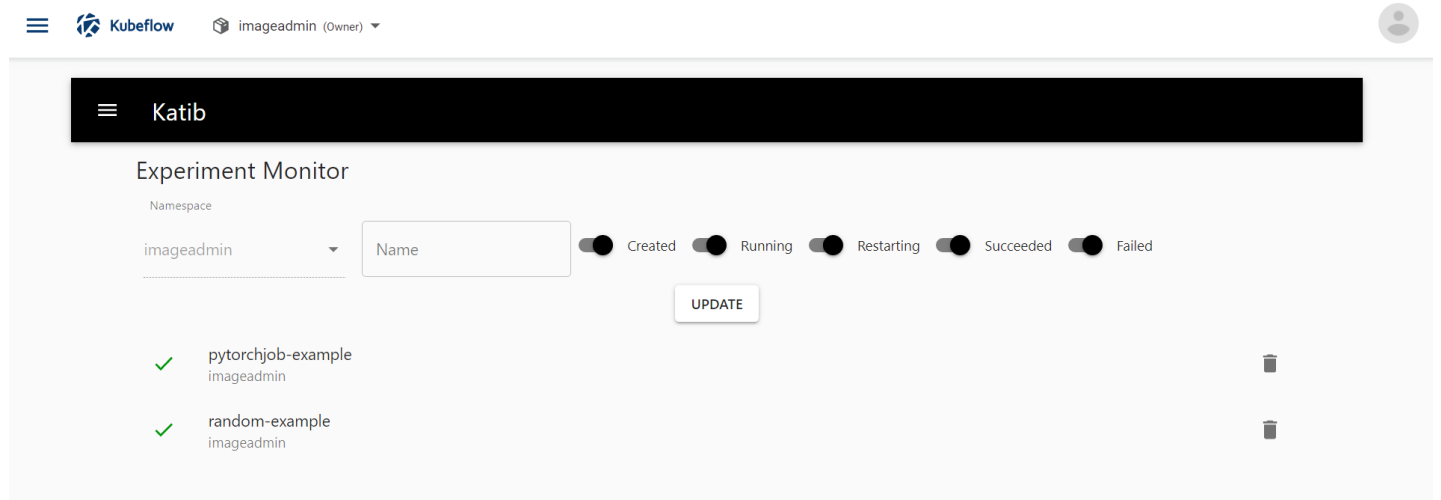


FIGURE A20. Katib HP Monitor

Click the experiment name, and then observe the running trials.

Check the experiment status:

```
$ kubectl get experiment
```

```
$ kubectl get experiment
NAME                STATUS      AGE
pytorchjob-example  Succeeded  11d
random-example      Succeeded  11d
$
```



Check the experiment trials.

```
$ kubectl get trial
```

```
$
$ kubectl get trial
NAME                                TYPE          STATUS    AGE
pytorchjob-example-2flzmzb9        Succeeded    True     11d
pytorchjob-example-4qhbxml9        Succeeded    True     11d
pytorchjob-example-55gj2kpx        Succeeded    True     11d
pytorchjob-example-55w2j4x6        Succeeded    True     11d
pytorchjob-example-58f9qtwc        Succeeded    True     11d
pytorchjob-example-5gwd5srs        Succeeded    True     11d
pytorchjob-example-7xsvgggz        Succeeded    True     11d
pytorchjob-example-dqtmjp8l        Succeeded    True     11d
pytorchjob-example-hkjg2hzm        Succeeded    True     11d
pytorchjob-example-vtd47rnv        Succeeded    True     11d
pytorchjob-example-zxlzvb6j        Succeeded    True     11d
pytorchjob-example-zzsglc9r        Succeeded    True     11d
random-example-25gr9zqw             Succeeded    True     11d
random-example-8cmpkkk9             Succeeded    True     11d
random-example-c5h242cr             Succeeded    True     12d
random-example-d9lsc7jl             Succeeded    True     11d
random-example-dmwcldxb             Succeeded    True     11d
random-example-gmb2ncgq             Succeeded    True     11d
random-example-kdp6pqjp             Succeeded    True     12d
random-example-lptj9mkv             Succeeded    True     11d
random-example-lpv84kvs             Succeeded    True     12d
random-example-m9kzdnmw             Succeeded    True     11d
random-example-nkkrz4nr             Succeeded    True     11d
random-example-sp6hhs8s             Succeeded    True     11d
$
```

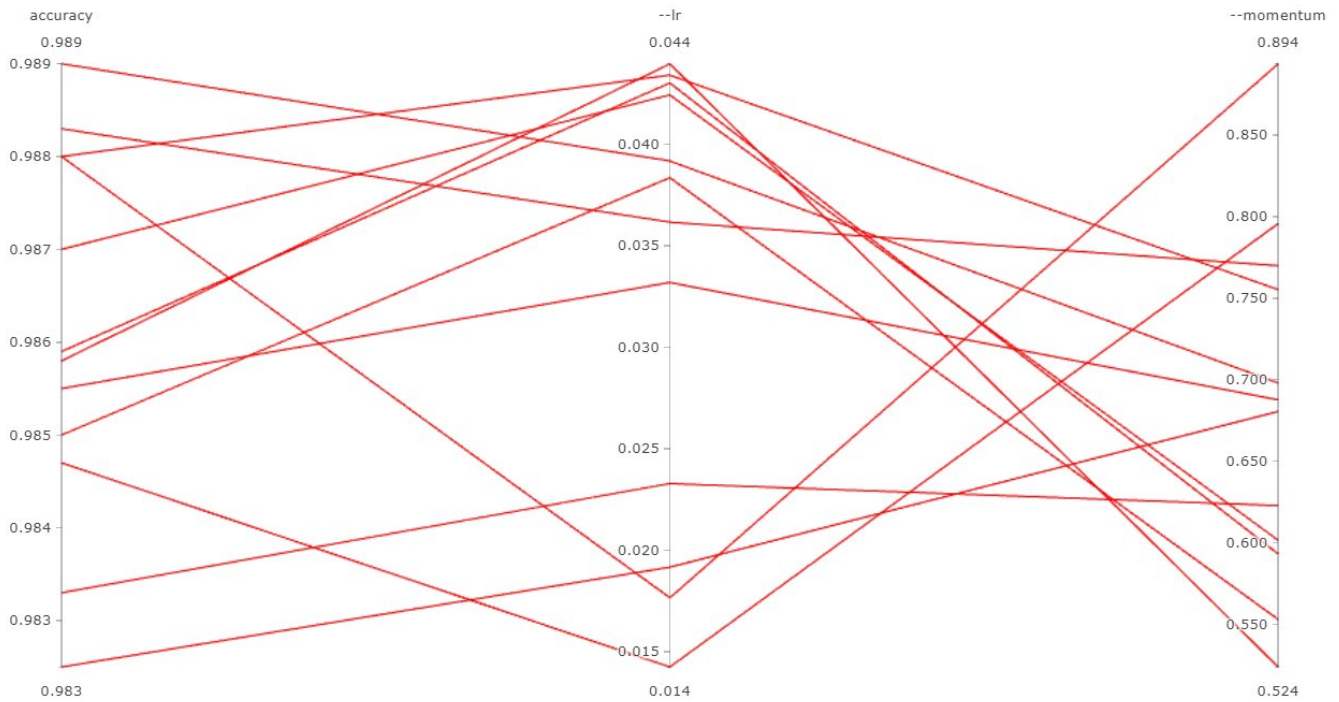


FIGURE A21. Katib HP Tuning



Example 2: Random algorithm

The following hyperparameters can be tuned:

```
--lr - learning rate  
--num-layers - Number of layers in the neural networks  
--optimizer
```

To launch an experiment using the random algorithm example:

■ Download the [kubeflow-tutorials.zip](#) (link opens an external website in a new browser tab/window).

■ Edit `random-example.yaml` to put the following on the pod template.

```
metadata:  
  annotations:  
    sidecar.istio.io/inject: "false"
```

■ Deploy the example:

```
$kubectl apply -f random-example.yaml
```

This example embeds the hyperparameters as arguments. Hyperparameters can be embedded in other ways (e.g. by using environment variables) by using the template defined in the TrialTemplate, GoTemplate, and RawTemplate section of the yaml file. The template uses the [Go template format](#) (link opens an external website in a new browser tab/window).

This example randomly generates the following hyperparameters:

```
--lr - Learning rate (type: double).  
--num-layers - Number of layers in the neural network (type: integer).  
--optimizer - Optimizer (type: categorical).
```



Check the experiment status:

```
$ kubectl describe experiment random-example
```

```
$ kubectl describe experiment random-example
Name:          random-example
Namespace:     imageadmin
Labels:        controller-tools.k8s.io=1.0
Annotations:   kubectl.kubernetes.io/last-applied-configuration:
                {"apiVersion":"kubeflow.org/v1alpha3","kind":"Experiment","metadata":{"annotations":{},"labels":{"controller
-tools.k8s.io":"1.0"},"name":"...
API Version:   kubeflow.org/v1alpha3
Kind:          Experiment
Metadata:
  Creation Timestamp:  2020-11-20T19:17:36Z
  Finalizers:
    update-prometheus-metrics
  Generation:         1
  Resource Version:   3303458
  Self Link:          /apis/kubeflow.org/v1alpha3/namespaces/imageadmin/experiments/random-example
  UID:                819d66c8-a977-46e9-8bfb-daf61c74e38
Spec:
  Algorithm:
    Algorithm Name:    random
    Max Failed Trial Count: 3
    Max Trial Count:   12
    Metrics Collector Spec:
      Collector:
        Kind: StdOut
    Objective:
      Additional Metric Names:
        Train-accuracy
      Goal: 0.99
      Objective Metric Name: Validation-accuracy
      Type: maximize
    Parallel Trial Count: 3
    Parameters:
      Feasible Space:
        Max: 0.03
        Min: 0.01
        Name: --lr
        Parameter Type: double
      Feasible Space:
        Max: 5
        Min: 2
        Name: --num-layers
        Parameter Type: int
      List:
        sgd
        adam
        ftrl
      Name: --optimizer
      Parameter Type: categorical
    Resume Policy: LongRunning
    Trial Template:
      Go Template:
        Raw Template: apiVersion: batch/v1
kind: Job
metadata:
  name: {{.Trial}}
  namespace: {{.Namespace}}
spec:
  template:
    spec:
      containers:
      - name: {{.Trial}}
        image: docker.io/kubeflow/katib/mxnet-mnist
        command:
```

Example 3: PyTorch

1. Download the [Kubeflow tutorials zip file](#) which contains sample files for all of the included Kubeflow tutorials.
2. Edit the YAML to point to put the following on the pod template.

```
metadata:D
  annotations:
    sidecar.istio.io/inject: "false"
```

3. Deploy the example.

```
$ kubectl apply -f pytorch-example.yaml
```

Go to the Katib page.

Click the Menu button, and then select **HP** → **Monitor**.



Click the **experiment** name, and then observe the trials running.

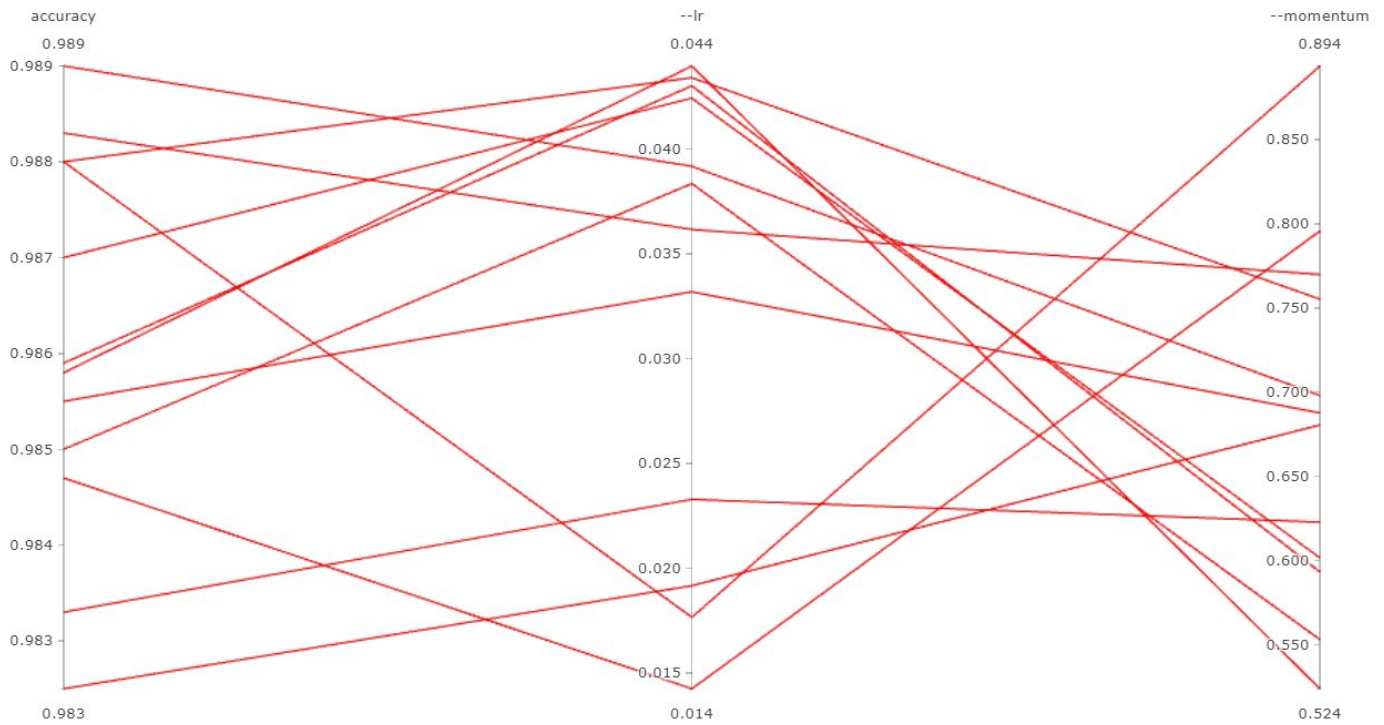


FIGURE A22. Pytorch Job Experiment

Check the experiment status.

```
$ kubectl get experiment
```

Use the following command to check the trials of the experiment.

```
$ kubectl get trial
```

Sample Katib commands

To check experiment results via the kubectl CLI.

- List experiments

```
$ kubectl get experiment
```

- Check experiment result

```
$ kubectl get experiment random-example -o yaml
```

- List trials

```
$ kubectl get trials
```

- Check trial detail

```
$ kubectl get trials random-experiment-24lgqghm -o yaml
```



To check the status using the interface:

- Go to the Katib page.
- Click the Menu button, and then select **HP** → **Monitor**.
- Click the experiment name and observe the built experiment graph after all the trials have succeeded.

Argo workflows

Download the [kubeflow_tutorials.zip](#) file, which contains sample files for all of the included Kubeflow tutorials.

This article provides the following two examples:

- [Simple workflow](#)
- [Parallel execution workflow](#)

Simple workflow

To complete the simple workflow:

- Create and apply the Argo workflow in the profile namespace.
- Obtain `argo-hello-world.yaml` from the zip file mentioned above.

```
$ kubectl apply -f argo-hello-world.yaml
workflow.argoproj.io/hello-world created
```

- Verify that the workflow was created:

```
$ kubectl get wf
NAME          AGE
hello-world   41m
parallelism-nested-dag 12d
```

- Verify that the related `hello-world` pod was created and is running.

```
$ kubectl get pods hello-world
NAME      READY STATUS  RESTARTS AGE
hello-world 0/2   Completed 0       49m
```

- Open the Argo interface by navigating to:

```
http://<kubeflow\_url>/argo/
http://ezam-01.perflab.hp.com:10053/?ns=imageadmin/argo/
```

- To remove the workflow:

```
$ kubectl delete wf hello-world
workflow.argoproj.io "hello-world" deleted
```

```
$ kubectl delete wf hello-world
workflow.argoproj.io "hello-world" deleted
```

Parallel execution workflow

To complete the parallel execution workflow:

- Create and apply the nested Argo workflow in the profile namespace.
- Obtain `argo-parallel-nested.yaml` from the zip file mentioned above.

```
$ kubectl apply -f argo-parallel-nested.yaml
workflow.argoproj.io/parallelism-nested-dag configured
```

- Verify that pods were created, as per the template (Observe the `.yaml` file the template).



Open the Argo interface by navigating to:

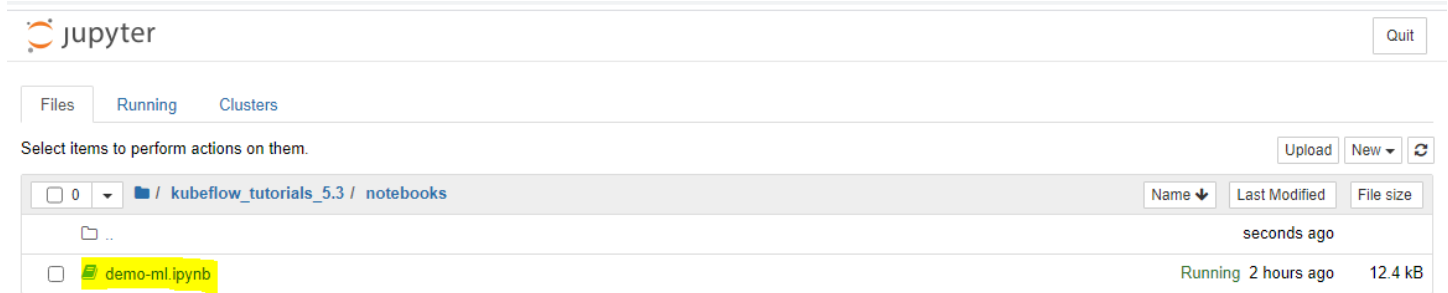
```
http://<kubeflow_url>/argo/
http://ezam-01.peeflab.hp.com:10053/?ns=imageadmin/argo/
```

ML metadata

Before beginning this tutorial, download the [Kubeflow tutorials zip file](#) which contains sample files for all of the included Kubeflow tutorials.

Create a Jupyter notebook server with any of the default images.

Connect to the created notebook server and upload the following notebook: demo-ml.ipynb.

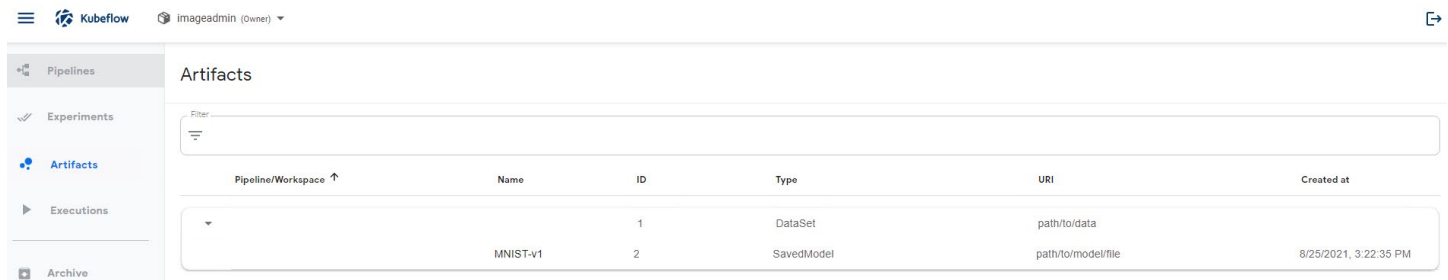


The screenshot shows the JupyterLab interface. At the top left is the Jupyter logo and a 'Quit' button. Below the logo are tabs for 'Files', 'Running', and 'Clusters'. A message says 'Select items to perform actions on them.' with 'Upload', 'New', and a refresh icon. The file browser shows the path '/ kubeflow_tutorials_5.3 / notebooks'. A table lists files with columns for Name, Last Modified, and File size. The file 'demo-ml.ipynb' is highlighted in yellow and shows 'Running 2 hours ago' and '12.4 kB'.

Name	Last Modified	File size
..	seconds ago	
demo-ml.ipynb	Running 2 hours ago	12.4 kB

FIGURE 63. Demo-ml

Run the notebook step by step, and observe the result on the **Pipelines** → **Artifacts** page in the Kubeflow UI.



The screenshot shows the Kubeflow UI. The top navigation bar includes the Kubeflow logo and the user 'imageadmin (Owner)'. The left sidebar has a menu with 'Pipelines', 'Experiments', 'Artifacts', 'Executions', and 'Archive'. The main content area is titled 'Artifacts' and has a search filter. Below the filter is a table with columns: Pipeline/Workspace, Name, ID, Type, URI, and Created at.

Pipeline/Workspace	Name	ID	Type	URI	Created at
		1	DataSet	path/to/data	
	MNIST-v1	2	SavedModel	path/to/model/file	8/25/2021, 3:22:35 PM

FIGURE A23. Artifacts



Click the name of each item to view detailed information. Click the Execution tab on the left and see the details.

The screenshot shows the Kubeflow interface. At the top, there is a navigation bar with the Kubeflow logo and the user name 'imageadmin (Owner)'. Below this is a sidebar with several menu items: Pipelines, Experiments, Artifacts, Executions (highlighted in blue), Archive, Documentation, Github Repo, and AI Hub Samples. The main content area is titled 'Executions' and features a back arrow. Underneath, it displays 'Type: Trainer' and 'Properties' with a table showing 'name: My Execution' and 'state: COMPLETED'. There are also sections for 'Custom Properties', 'Declared Inputs', and 'Declared Outputs', each with a corresponding table.

name	state
My Execution	COMPLETED

Artifact ID	Name	Type	URI
1		DataSet	path/to/data

Artifact ID	Name	Type	URI
2	MNIST-v1	SavedModel	path/to/model/file

FIGURE A24. Executions



APPENDIX B: HPE ML OPS KDAPP

Centos/Ubuntu

In the Applications screen, launch Centos/Ubuntu app with the required resources screen, launch Centos/Ubuntu app with the required resources.

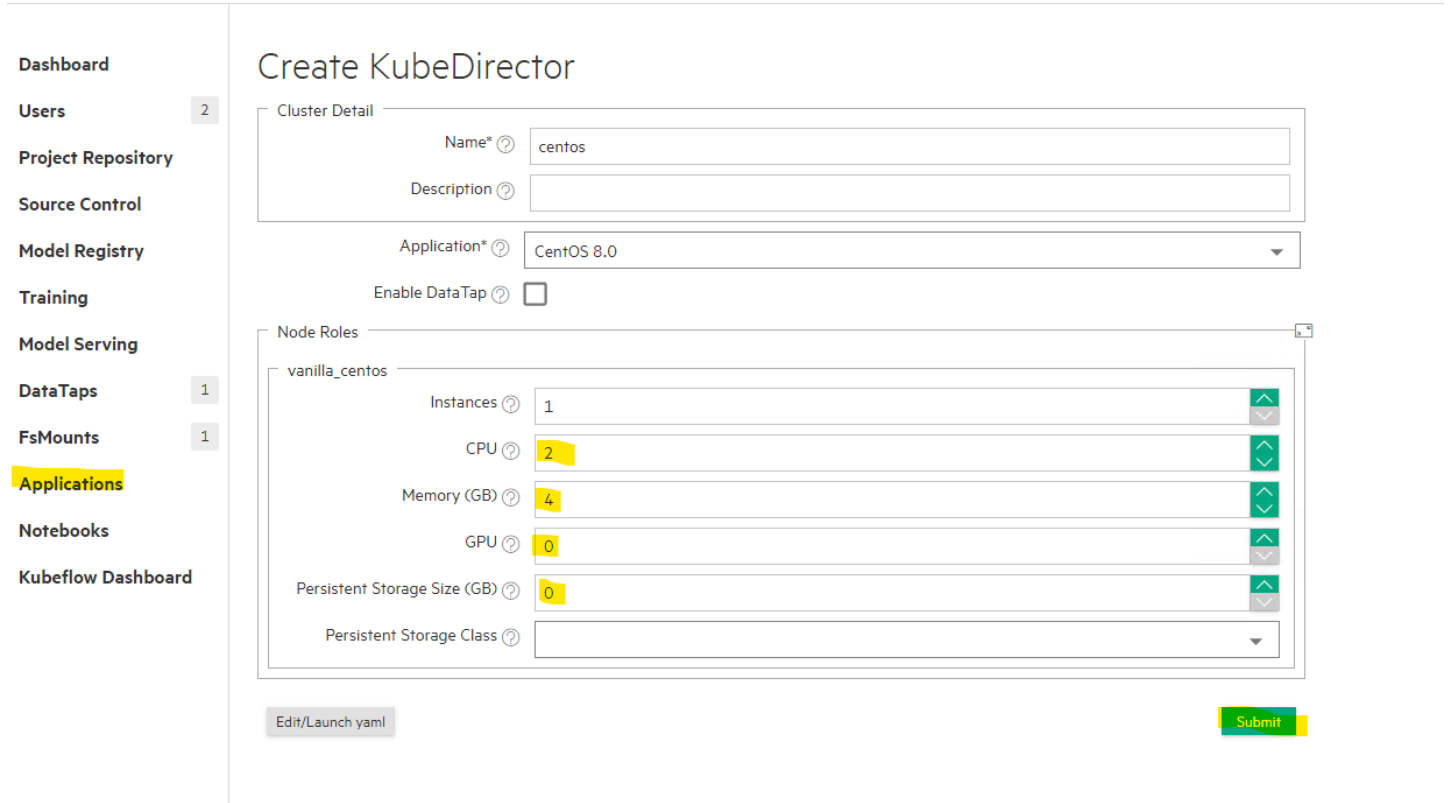


FIGURE B1. Launching KDAPP CentOS

ssh-keygen in webterm and copy id_rsa.pub.

```
k8suser@kds-bwk28-0:~$ ssh-keygen
Generating public/private rsa key pair.
Enter file in which to save the key (/home/k8suser/.ssh/id_rsa):
Created directory '/home/k8suser/.ssh'.
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /home/k8suser/.ssh/id_rsa.
Your public key has been saved in /home/k8suser/.ssh/id_rsa.pub.
The key fingerprint is:
SHA256:an0kRtzhkPUCxszZ7B+m14AcZMAqL8/n0FtTpr2a k8suser@kds-bwk28-0
The key's randomart image is:
+----[RSA 2048]-----+
| .+*o..o
|oo=ooo
|=.o.o...
|+.o..o
|o..B + + S
|.o 8 + = .
|..= o o .
|.o .+BO .
| .oo.
+----[SHA256]-----+
k8suser@kds-bwk28-0:~$
k8suser@kds-bwk28-0:~$ cat .ssh/id_rsa.pub
ssh-rsa AAAAB3NzaC1yc2EAAAADAQABAAQDuu/ZzMrbv1v8LYbFN10tX6utdhMzJIFRAw6scavGIYlYK1U5pKzP4pAMGTnk1jV0zJ6Txb/5uATCD5uBRZjs5jQcIrm2tKK3Zq0/hK-1aorQ1ANODY5F7XLNnpU6CK6h1CmGvN1ySa1/Mv/OT5d5uWeER9oqiKnMTu6uChjvYnd4pznOmMsXGStH9wSB6w/zpBxpR3jn2lBlnYGqfjwhYK9y1YFMs+8mzstCmR8z0j0nOU2LYrtJKoy9EA5r3QBS1cEVMRQjkr/Lj1k15ST4f7wFEBK1Pdc/Et0+Tle/RBvN11kpbH/s65TcH3gt6nS0ogPIPI9mT4xfYR0JER5 k8suser@kds-bwk28-0
k8suser@kds-bwk28-0:~$
```



Exec into the centos pod and run ssh-keygen.

```
k8suser@kds-4dzrd-0:~$
k8suser@kds-4dzrd-0:~$ kubectl exec -it centos-vanilla-centos-znsk6-0 -- bash
[root@centos-vanilla-centos-znsk6-0 /]# ssh-keygen
Generating public/private rsa key pair.
Enter file in which to save the key (/root/.ssh/id_rsa): c
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in c.
Your public key has been saved in c.pub.
The key fingerprint is:
SHA256:n/3GTa5DbTCvFSpfxDUKxCKz9J9VbbPA+liq27EWDTY root@centos-vanilla-centos-znsk6-0
The key's randomart image is:
+---[RSA 3072]-----+
|
|   o..  .
|  + . o o ++
| . = . . * . =
| . . . Eo*
| S . + @ = .
| . * = . + *
| o = + o B
| . . + B o
| . o + o o
+----[SHA256]-----+
[root@centos-vanilla-centos-znsk6-0 /]#
```

Vi /root/.ssh/authorized_keys, paste the id_rsa.pub that was copied, and exit.

```
[root@centos-vanilla-centos-zptfk-0 /]# vi root/.ssh/authorized_keys
[root@centos-vanilla-centos-zptfk-0 /]# cat root/.ssh/authorized_keys
ssh-tsa AAAAB3NzaC1yc2EAAAADAQABAAQDuu/2rMrb1v8Lbfn0tX6utdhMzJIFRAw65cavG1YlYKIU5pKz4pAWGTnKj1v0zJ6Txb/SuATCD5uBR2ja5J0cIrm2tKK3g0/hK+1aorQ1ANODYSF7XLWnpU6CK6h1CmGvNly1Sa1/Mv/OT5d5uWeER9oq1KnMTu6u6Cn1jvYND4pznOmMxGStB9w
SB6w/zpBXprj2lB1sYqgfjwhyK91YFM+8mzsTcMr8z0j0n0U2LYrJK0y9A5r3QBS1rEVMRQjkr/Ljk15ST4f7WFEbK1Pdc/Etc+TLe/RBVn1lIkpH/s6S7EH3g76nSGOgFIIP19mT4XYROJFK5 k8suser@kds-4dzrd-0
[root@centos-vanilla-centos-zptfk-0 /]#
```

ssh to the Access Points.

Kubernetes Applications

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
centos-vanilla-centos-zptfk-0	vanilla-centos	KubeDirectorApp: ID: centos8x Name: CentOS 8.0	centos	SSH	22	ezam-01.perflab.hp.com:10026	NodePort
kf-dashboard-import-sonbv					80	ezam-01.perflab.hp.com:10019	NodePort
lily-http				http	8998	ezam-01.perflab.hp.com:10023	NodePort
spark-ui-proxy				http	80	ezam-01.perflab.hp.com:10022	NodePort
sparkhe-svc				http	18480	ezam-01.perflab.hp.com:10020	NodePort
sparkts-svc				http	4440	ezam-01.perflab.hp.com:10021	NodePort
				spark-thrift	2304	ezam-01.perflab.hp.com:10024	NodePort

FIGURE B2. CentOS Access Points

```
k8suser@kds-bwkz8-0:~$
k8suser@kds-bwkz8-0:~$ ssh root@ezam-01.perflab.hp.com -p 10026
The authenticity of host '[ezam-01.perflab.hp.com]:10026 ([172.24.2.1]:10026)' can't be established.
ECDSA key fingerprint is SHA256:DKJi3QsAESBEAHwMOK41V4TJG1WakTK4c9kkFmsH8Ik.
ECDSA key fingerprint is MD5:c8:b0:18:40:f9:e4:6f:31:d0:2d:50:4d:64:48:90:91.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added '[ezam-01.perflab.hp.com]:10026, [172.24.2.1]:10026' (ECDSA) to the list of known hosts.
[root@centos-vanilla-centos-zptfk-0 ~]# ls
anaconda-ks.cfg anaconda-post.log original-ks.cfg
[root@centos-vanilla-centos-zptfk-0 ~]# pwd
/root
[root@centos-vanilla-centos-zptfk-0 ~]#
```

MLflow

MLflow is an open-source platform to manage the machine learning lifecycle, including experimentation, reproducibility, deployment, and a central model registry. For MLflow integration in the HPE Ezmeral Runtime details, see [MLflow for Model Management](#). See [MLflow Configuration and Deployment](#) for the process to execute one run of MLflow from training to deployment. The user can clone the MLflow notebook used in this tutorial from <https://github.com/pcao11/mlflow.git>.



Generate and apply the MLflow Secret (See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for the details).

```

k8suser@kdss-r9bk8-0:~$ vi mlflow-sc.yaml
k8suser@kdss-r9bk8-0:~$ cat mlflow-sc.yaml
apiVersion: v1
kind: Secret
stringData:
  MLFLOW_ARTIFACT_ROOT: s3://hanshabucket #s3://mlflow
  AWS_ACCESS_KEY_ID: AKIAW2IT4GWHVDDWS25E #myusername
  AWS_SECRET_ACCESS_KEY: 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 #mypassword
  #MLFLOW_BACKEND_STORE: xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx=
                        #mysql://mlflowusr:Password^12@192.0.2.0:3306/mlflowdb
  MLFLOW_S3_ENDPOINT_URL: https://s3.us-east-2.amazonaws.com #https://s3.us-east-2.amazonaws.com
  #MLFLOW_S3_ENDPOINT_URL: https://hanshabucket.s3.us-east-2.amazonaws.com #https://s3.us-east-2.amazonaws.com
  #xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx==
                        #http://myserver.example.com:10007
  #AWS_DEFAULT_REGION: us-east-2
metadata:
  name: mlflow-sc
  labels:
    kubedirector.hpe.com/secretType: mlflow
type: Opaque
  #podLabels:
    #hpecp.hpe.com/dtap: "hadoop2"

k8suser@kdss-r9bk8-0:~$ kubectl apply -f mlflow-sc.yaml
secret/mlflow-sc created
k8suser@kdss-r9bk8-0:~$

```

Create an MLflow App Instance by attaching the mlflow secret (See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for the details).

Create KubeDirector

Edit your Yaml and hit submit to launch app

```

1 ---
2 apiVersion: "kubedirector.hpe.com/v1beta1"
3 kind: "KubeDirectorCluster"
4 metadata:
5   name: "mlflow"
6   namespace: "amdtenant"
7   labels:
8     description: ""
9 spec:
10  app: "mlflow"
11  namingScheme: "CrNameRole"
12  appCatalog: "local"
13  connections:
14    secrets:
15      - mlflow-sc
16  roles:
17    -
18      id: "controller"
19      members: 1
20      resources:
21        requests:
22          cpu: "2"
23          memory: "4Gi"
24          nvidia.com/gpu: "0"
25        limits:
26          cpu: "2"
27          memory: "4Gi"
28          nvidia.com/gpu: "0"
29      #Note: "if the application is based on hadoop3 e.g. using StreamCapabilities interface, then change the below dtap Label to 'hadoop3', otherwise for most applications use the default 'hadoop2'"
30      #podLabels:
31        #hpecp.hpe.com/dtap: "hadoop2"
32

```

Submit Cancel

FIGURE B3. Launching MLflow App Instance



■ Create a Training Application Instance. See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for details (optional).

Create Training

Cluster Detail

Name*

Description

RunTime Image* ML Training Toolkit, with GPU support

Enable DataTap

Node Roles

LoadBalancer

Instances

CPU

Memory (GB)

GPU

RETServer

Instances

CPU

Memory (GB)

GPU

controller

Instances

CPU

Memory (GB)

GPU

Edit/Launch yaml
Submit

FIGURE B4. Create Training Cluster

■ Launch and configure the Notebook (See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for more details).

Create Notebook

Edit your Yaml and hit submit to launch app

```

1 ---
2 apiVersion: "kubedirector.hpe.com/v1beta1"
3 kind: "KubeDirectorCluster"
4 metadata:
5   name: "nbmlflow"
6   namespace: "amdenant2"
7   labels:
8     description: ""
9 spec:
10  app: "jupyter-notebook"
11  namingScheme: "CrNameRole"
12  appCatalog: "local"
13  connections:
14    clusters:
15      - trcpu
16      - mlflow
17    secrets:
18      - hpecp-kc-secret-d69f53e21bb68f1b0f0b5f7088726
19      - hpecp-ext-auth-secret
20      - mlflow-sc
21  roles:
22    -
23      id: "controller"
24      members: 1
25      resources:
26        requests:
27          cpu: "2"
28          memory: "4Gi"
29          nvidia.com/gpu: "0"
30        limits:
31          cpu: "2"
32          memory: "4Gi"
33          nvidia.com/gpu: "0"
34      #Note: If the application is based on hadoop3 e.g. using StreamCapabilities interface, then change the below dtap label to 'hadoop3', otherwise for most applications use the default 'hadoop2'
35      #podLabels:
36        #hpecp.hpe.com/dtap: "hadoop2"
37

```

Submit
Cancel

FIGURE B5. Create MLflow Notebook



■ Enable kubectl to run MLflow backend (mandatory). Set the experiment name, Train, and track models.

Set your password (mandatory)

```
[1]: PASSWORD = "hp123456" # use your password
```

Enable kubectl to run MLflow backend (mandatory)

```
[2]: %kubeRefresh --pwd $PASSWORD
```

```
kubeconfig set for user imageadmin
```

```
[3]: # This magic sets the environmental variables required for mlflow in backend.
%loadMlflow
```

```
Backend configured
```

Set your experiment name

```
[4]: # Magic function '%Setexp' replaces the two lines below.
#mlflow.set_experiment('demoexp')
#mlflow.set_tag('mlflow.user', 'chris')
%Setexp --name demoexp
```

■ Observe experiments, runs, metrics, parameters, dependency, and trained models in the MLflow UI.

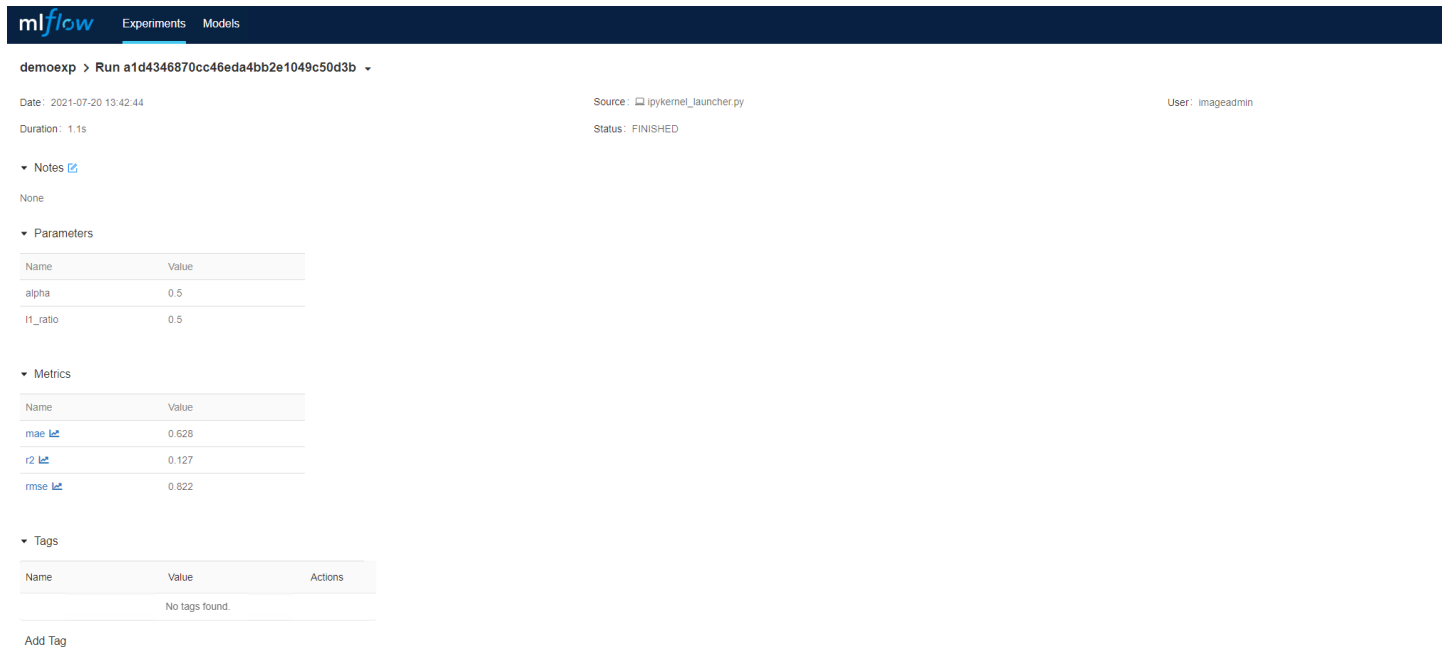


FIGURE B6. MLflow Experiment, Run



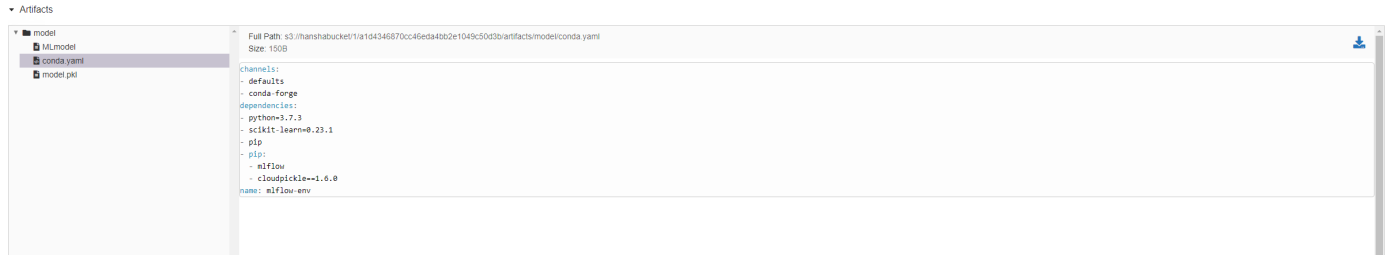


FIGURE B7. Dependency and Trained Model

Register Model for MLflow (See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for details).

Register Model

Label

Name*

Description

Model Store Type

Model Artifact URL*

FIGURE B8. Register Model



Create an MLflow Model Serving (See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for details).

Create Model Serving

Cluster Detail

Name*

Description

Model Serving Engine*

Select Model*

Resources

Instances

CPU

Memory (GB)

FIGURE B8. MLflow Model Serving

Copy Access Points and Auth Token.

Model Serving

Name	Access Points	Auth Token
mlflowserving	http://ezam-01.perflab.hp.com:10021/seldon/amdtenant2/model/mlflowservingapi/v1.0/predictions	Auth Token

Making Prediction Calls (See the [HPE Ezmeral Container Platform 5.3 Documentation](#) for details).

```
[root@ezam-09 ~]#
[root@ezam-09 ~]#
[root@ezam-09 ~]# curl -k --cookie "authservice_session=MYyHjpwMjgM3x0d3BTKWQlVFRk1VWVEV0U4MfZVZrdrVr1Zi2B5S13qITN3MfFmsxTNGQlNrvk5Vwz10U1NVCERBkuzpGovrdFVl3wEUFJavakv0Sp1ZbaUdCOYL1shvz2G-StuceSp_wcMFTT" -X POST -H "Content-Type: application/json" -d '{"data":{"names": ["fixed acidity", "volatile acidity", "citric acid", "residual sugar", "chlorides", "free sulfur dioxide", "total sulfur dioxide", "density", "ph", "sulphates", "alcohol"]}, "meta":{"names": ["data"], "ndarray": [5, 65599999229192]}, "meta": {}}
```



NVIDIA: TensorFlow (NGC)

■ Create an NVIDIA TensorFlow (NGC) App Instance.

▭ HPE Ezmeral Container Platform

Create KubeDirector

Cluster Detail

Name*

Description

Application*

Enable DataTap

Node Roles

tensorflow

Instances

CPU

Memory (GB)

GPU

Edit/Launch yaml Submit

FIGURE B9. Launching NVIDIA TensorFlow (NGC) App Instance

■ Login inside the pod to run TensorFlow jobs with GPU.

```
k8suser@kds-bwkz8-0:~$
k8suser@kds-bwkz8-0:~$ kubectl exec -it ngc-tensorflow-gvz94-0 -- bash
root@ngc-tensorflow-gvz94-0:/workspace# vi tftest.py
root@ngc-tensorflow-gvz94-0:/workspace# python3 tftest.py
2021-08-16 18:28:24.644481: I tensorflow/stream_executor/platform/default/dso_loader.cc:49] Successfully opened dynamic library libcudart.so.11.0
WARNING:tensorflow:Deprecation warnings have been disabled. Set TF_ENABLE_DEPRECATION_WARNINGS=1 to re-enable them.
2021-08-16 18:28:28.447813: I tensorflow/core/platform/profile_utils/cpu_utils.cc:94] CPU Frequency: 2500000000 Hz
```



TensorFlow + Jupyter

Create a TensorFlow + Jupyter App Instance.

HPE Ezmeral Container Platform

Dashboard

- Users 2
- Project Repository
- Source Control
- Model Registry
- Training
- Model Serving
- DataTaps 1
- FsMounts 1
- Applications
- Notebooks
- Kubeflow Dashboard

Create KubeDirector

Cluster Detail

Name*

Description

Application*

Enable DataTap

Node Roles

controller

- Instances
- CPU
- Memory (GB)
- GPU

FIGURE B10. Launching TensorFlow + Jupyter App Instance

Click the **access points** of the created App Instance.

Kubernetes Applications

KubeDirector Kubectl **Service Endpoints** Virtual Endpoints

Kubernetes Service Name	Role	Details	KubeDirector Cluster	Services	Ports	Access Points	Service Type
centos-vanilla-centos-qptfk-0	vanilla_centos	KubeDirectorApp: ID: centos8x Name: CentOS 8.0	centos	SSH	22	ezam-01.perflab.hp.com:10026	NodePort
kf-dashboard-import-l6n8v				80	80	ezam-01.perflab.hp.com:10019	NodePort
ivy-http				http	8998	ezam-01.perflab.hp.com:10023	NodePort
spark-ui-proxy				http	90	ezam-01.perflab.hp.com:10022	NodePort
sparkui-svc				http	18480	ezam-01.perflab.hp.com:10020	NodePort
sparkui-svc				http	4440	ezam-01.perflab.hp.com:10021	NodePort
tfnb-controller-7qz29-0	controller	KubeDirectorApp: ID: tensorflow-gpu-jupyter Name: TensorFlow + Jupyter	tfnb	Jupyter Notebook	8888	ezam-01.perflab.hp.com:10024	NodePort

FIGURE B11. Access Points of TensorFlow + Jupyter



Exec into the TensorFlow +Jupyter pod and run `jupyter notebook list` and copy the token.

```
k8suser@kds-bwkz8-0:~$ kubectl exec -it tfnb-controller-9cg29-0 -- bash
_____
 /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
/_/_/  /_/_/  /_/_/  /_/_/  /_/_/
/_/_/  /_/_/  /_/_/  /_/_/  /_/_/
/_/_/  /_/_/  /_/_/  /_/_/  /_/_/
/_/_/  /_/_/  /_/_/  /_/_/  /_/_/

WARNING: You are running this container as root, which can cause new files in
mounted volumes to be created as the root user on your host machine.

To avoid this, run the container by specifying your user's userid:

$ docker run -u $(id -u):$(id -g) args...

root@tfnb-controller-9cg29-0:/tf# jupyter notebook list
Currently running servers:
http://0.0.0.0:8888/?token=57590ab553e56308eca408bb89f025d3dbdd71f6c3f352d2 :: /tf
root@tfnb-controller-9cg29-0:/tf#
```

FIGURE B12. Copy the token

Paste the copied token and login to TensorFlow + Jupyter notebook.

The screenshot shows a web browser window with the URL `10.209.18.1:10034/login?next=%2Ftree%3F`. The page is the Jupyter Notebook login interface. At the top right, the Jupyter logo is visible. Below it, there is a "Password or token:" label followed by a text input field containing a long alphanumeric token and a "Log in" button. Underneath, a section titled "Token authentication is enabled" explains that a password is not required if a token is provided. It also lists the command `jupyter notebook list` and shows the output from the previous figure as an example of the command's output. A "Setup a Password" section follows, with a sub-header "You can also setup a password by entering your token and a new password on the fields below:". This section contains two input fields: "Token" and "New Password", and a "Log in and set new password" button.

FIGURE B13. Login to TensorFlow + Jupyter Notebook



Now launch the notebook and run TensorFlow with GPU

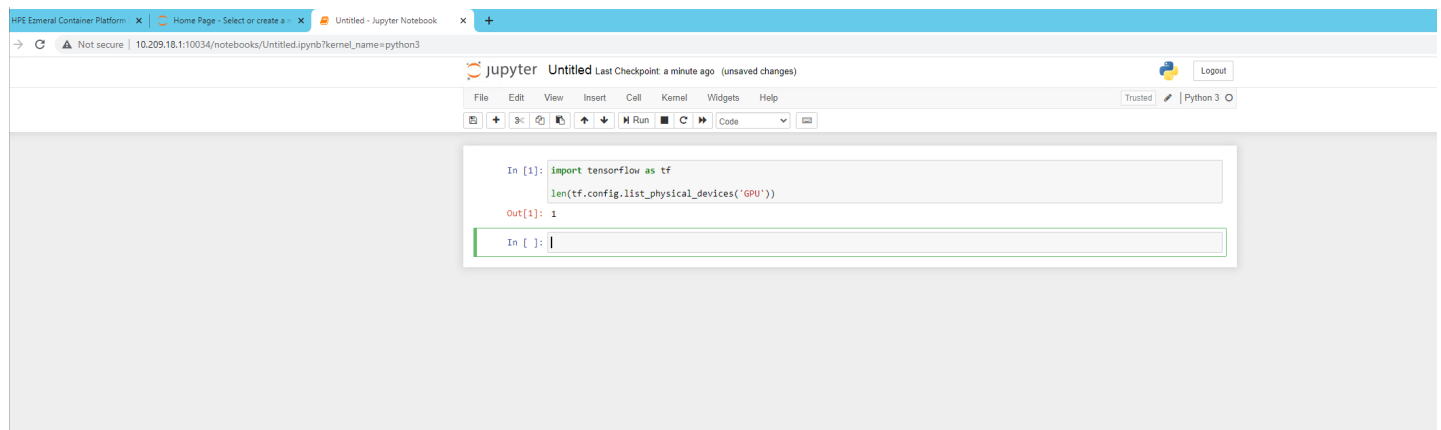


FIGURE B14. Python3 Notebook

APPENDIX C: INSTALL AND CONFIGURE HPE EZMERAL RUNTIME

Follow the steps as outlined to install and configure the HPE Ezmeral Runtime. This section assumes that all the prerequisites mentioned in the earlier sections were followed. See the [Standard Installation](#) procedure.



Reference Architecture

RESOURCES AND ADDITIONAL LINKS

HPE Reference Architectures, hpe.com/info/ra

HPE Servers, hpe.com/servers

HPE Storage, hpe.com/storage

HPE Networking, hpe.com/networking

HPE Technology Consulting Services, hpe.com/us/en/services/consulting.html

HPE Ezmeral Machine Learning Ops, <https://buy.hpe.com/us/en/enterprise-solutions/artificial-intelligence-analytics/artificial-intelligence-analytics/artificial-intelligence-analytics/hpe-ezmeral-machine-learning-ops/p/1011947349>

Operationalization for the Machine Learning Lifecycle live demonstration, <https://hpedemoportal.api.ext.hpe.com/DemoPortal/api/DocContent/GetDocByToken/26858daf-c14b-49f0-ae8f-60d3e7423e66>

HPE Ezmeral ML Ops, <https://assets.ext.hpe.com/is/content/hpedam/documents/a50000000-0999/a50000137/a50000137enw.pdf>

Kubeflow Introduction, <https://www.kubeflow.org/docs/pipelines/overview/pipelines-overview/>

<https://v0-5.kubeflow.org/docs/use-cases/>

Broadcom AIOps, <https://www.broadcom.com/sw-tech-blogs/aiops-blog/what-is-prometheus#:~:text=Prometheus%20is%20a%20time-series%20streaming%20data%20tool.%20It,make%20that%20data%20available%20for%20processing%20and%20analysis>

HPE Ezmeral Container Platform 5.3 Documentation, <https://docs.containerplatform.hpe.com/53/index.html>

To help us improve our documents, please provide feedback at hpe.com/contact/feedback.

© Copyright 2022- 2023 Hewlett Packard Enterprise Development LP. The information contained herein is subject to change without notice. The only warranties for Hewlett Packard Enterprise products and services are set forth in the express warranty statements accompanying such products and services. Nothing herein should be construed as constituting an additional warranty. Hewlett Packard Enterprise shall not be liable for technical or editorial errors or omissions contained herein.

Intel and Xeon are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries. NVIDIA, the NVIDIA logo, are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Microsoft and Windows are trademarks of Microsoft Corporation in the United States and/or other countries. © 2012 Google Inc. All rights reserved. Google and the Google Logo are registered trademarks of Google Inc. All third-party marks are property of their respective owners.