

# IBM watsonx.data

## Technical Hands-on Lab

Developed, tested, and recommended for use with the IBM TechZone  
**IBM watsonx.data Development Lab 1.1.3 R1M0** environment

Kelly Schlamb ([kschlamb@ca.ibm.com](mailto:kschlamb@ca.ibm.com))  
Worldwide Technology Enablement

## Contents

1. Introducing watsonx.data .....	4
2. About this Lab .....	6
3. Getting Help .....	7
4. Prerequisites & Getting Started .....	8
4.1. Provision a watsonx.data Environment from TechZone .....	9
4.2. Accessing the watsonx.data Environment .....	13
4.3. Command Line Access .....	15
4.4. watsonx.data Infrastructure Components .....	17
4.5. Stopping and Starting watsonx.data .....	20
4.5.1. Stopping watsonx.data .....	20
4.5.2. Starting watsonx.data .....	21
5. Exploring the watsonx.data User Interface .....	23
5.1. Starting the watsonx.data User Interface .....	23
5.2. Infrastructure Manager Page .....	27
5.3. Data Manager Page .....	34
5.4. Query Workspace Page .....	44
5.5. Query History Page .....	52
5.6. Access Control Page .....	55
6. Working with Presto .....	63
6.1. Presto Command Line Interface (CLI) .....	63
6.2. Presto Web Interface .....	69
7. Working with MinIO .....	73
7.1. Introduction to Object Storage .....	73
7.2. Exploring MinIO Object Storage .....	74
8. Data Ingest .....	78
9. Federated Queries .....	84
10. Offloading Data from a Data Warehouse .....	97
11. Time Travel .....	101
11.1. Table Rollback .....	101
11.2. Time Travel Queries .....	105

12. Summary.....	111
Appendix A. Troubleshooting.....	112
watsonx.data Web Interface: .....	112
MinIO Web Interface: .....	113
Miscellaneous:.....	113
Appendix B. Acknowledgements .....	114
Appendix C. Revision History .....	115

## 1. Introducing watsonx.data

Watsonx.data is a core component of **watsonx**, IBM's enterprise-ready AI and data platform designed to multiply the impact of AI across an enterprise's business.

The watsonx platform comprises three powerful components: the **watsonx.ai** studio for new foundation models, generative AI and machine learning; the **watsonx.data** fit-for-purpose data store that provides the flexibility of a data lake with the performance of a data warehouse; plus the **watsonx.governance** toolkit, to enable AI workflows that are built with responsibility, transparency, and explainability.

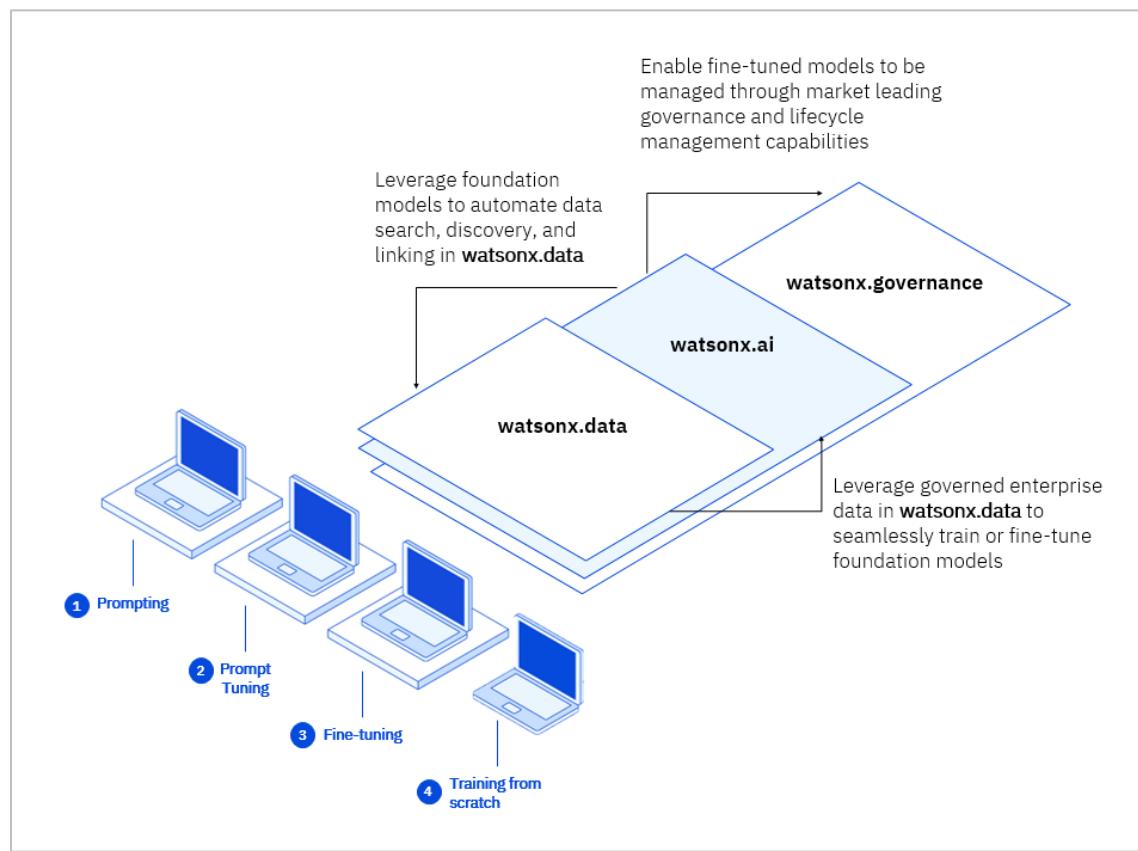


Figure 1: IBM watsonx platform

The watsonx.data component (the focus of this lab) makes it possible for enterprises to scale analytics and AI with a data store built on an open lakehouse architecture, supported by querying, governance, and open data and table formats, to access and share data. With watsonx.data, enterprises can connect to data in minutes, quickly get trusted insights, and reduce their data warehouse costs.

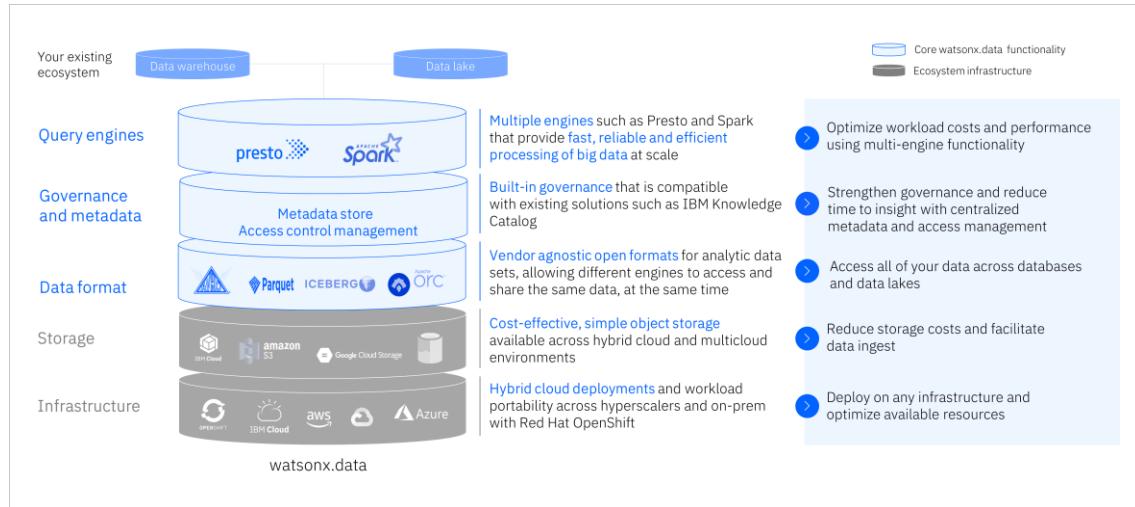


Figure 2: IBM watsonx.data's lakehouse architecture and value proposition

## 2. About this Lab

This IBM watsonx.data hands-on lab introduces you to some of the core components and capabilities of IBM watsonx.data. By completing this lab, you will gain the experience needed to demo these capabilities to clients.

Specifically, you will get hands-on experience in the following areas:

- The watsonx.data web-based user interface (UI), including infrastructure management, data management, running SQL statements, and managing user access
- The Presto web interface and the Presto command line interface (CLI)
- MinIO object storage
- Ingesting data into watsonx.data
- Creating schemas and tables
- Running queries that combine data from multiple data sources (data federation)
- Offloading tables from Db2 into watsonx.data
- Rolling back a table to a previous point in time

This lab requires that you provision an environment from IBM Technology Zone (TechZone). The image used comes pre-configured with watsonx.data Developer Edition, additional database systems including Db2 and PostgreSQL, and sample data sets. Provisioning details are provided in the [Prerequisites & Getting Started](#) section of this lab guide.

Some of the lab sections require that you copy commands or SQL statements and paste them into your browser or a command line terminal. If you have problems copying and pasting from this document, all of the commands and SQL statements can be found in this [text document](#). Download this file and copy the text from there instead.

Watsonx.data is being developed and released in an agile manner. In addition to new capabilities being added, the web interface is also likely to change over time. Therefore, the screenshots used in this lab may not always look exactly like what you see.

### 3. Getting Help

**Lab guide help:** If you require assistance in interpreting any of the steps in this lab, please post your questions to the [#data-ai-demo-feedback](#) Slack channel (IBMers only). Business Partners can request help at the [Partner Plus Support](#) website.

**TechZone environment:** If you are encountering issues regarding the TechZone environment being used in this lab, including the inability to provision an environment, please see the [TechZone Help](#) page.

**watsonx.data:** Assistance with the watsonx.data product itself is available in the [#watsonxdata-lakehouse-discussion-open-to-all-ibmers](#) Slack channel (IBMers only). Additionally, please refer to the watsonx.data documentation as needed ([SaaS](#), [software](#)).

**Additional troubleshooting:** See [Appendix A. Troubleshooting](#) at the end of this lab guide for guidance on commonly encountered issues.

## 4. Prerequisites & Getting Started

Clients can deploy watsonx.data in a number of different ways:

- As software-as-a-service (SaaS) on IBM Cloud and Amazon Web Services (AWS)
- As a cartridge with Cloud Pak for Data
- As standalone hybrid-cloud software that can be installed on Red Hat OpenShift (on-premises or in the cloud)
- As a simple, single-node Developer Edition installation

Watsonx.data is currently available in a *Standard Edition*, with an *Enterprise Edition* planned for the future. For the developer and partner community, IBM also offers an entry-level *Developer Edition*, which can be used to get familiar with the watsonx.data console and environment. The Developer Edition has the same code base as the Standard Edition, but some features are restricted, and it is not intended for production use.

This lab utilizes a pre-installed Developer Edition virtual machine (VM) environment that can be easily provisioned from IBM Technology Zone (TechZone). Provisioning instructions are included in the section that follows this one.

A number of software components are included as part of the VM. In this lab, you will be primarily interfacing with the following components:

- watsonx.data (which includes a graphical console)
- Presto (a SQL query engine used in watsonx.data; it includes its own graphical console)
- MinIO (local object storage; it includes its own graphical console)

The Developer Edition comes pre-configured with a number of watsonx.data *infrastructure components*, including a Presto query engine, three catalogs, and three object storage buckets.

## 4.1. Provision a watsonx.data Environment from TechZone

Follow the steps below to provision the watsonx.data VM lab environment from IBM Technology Zone (TechZone). This environment is available to IBM employees and IBM Business Partners, but not clients.

1. Open the **IBM watsonx.data Developer Base Image** collection in **IBM Technology Zone** at: <https://techzone.ibm.com/collection/ibm-watsonxdata-developer-base-image>. Sign in with your IBMid and accept any terms and conditions you are presented with.
2. Select the **Environments** tab in the left-side menu.

This image contains the required software and prerequisites that are needed to run the IBM watsonx.data software. This image is built on top of RHEL 9 running on a 4 VPC/16Gb /400Gb server which has been found to be sufficient for running the developer edition of IBM watsonx.data. The developer version is meant to be used on single nodes. While it uses the same code base as a production version, there are some restrictions, especially on scale. This system is meant to demonstrate how everything works, but not performance! Keep that in mind when running this system!

Note: The lab materials are available in the resources section. See <https://ibm.biz/wxd-lab> for details on the lab.

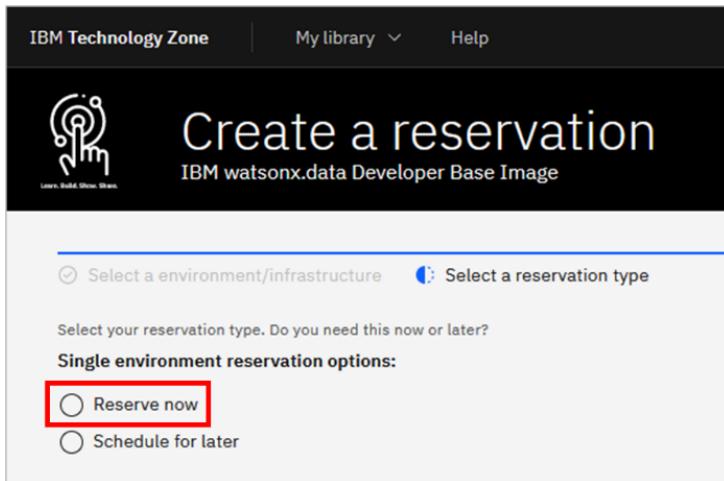
**Note:** This TechZone Collection includes its own distinct lab guide (which can be found in the **Resources** tab) section. To be clear, that lab is *not* the same as this lab that you are currently going through. There are some common elements between the two labs, so consider taking it to reinforce what you are learning in this one, but that other lab is *not* a part of the watsonx.data for Technical Sales Level 3 learning plan and does not need to be done to complete this learning plan.

3. Click the **IBM watsonx.data Development Lab - 1.1.3 R1M0** tile (note the specific version).

Environment	Date	Description
IBM watsonx.data Development Lab - 1.1.1 R1M0	Apr 21, 2024	Ibmcloud 2: any, any, any This is version 1.1.1 of the <a href="#">watsonx.data</a> development lab and should be used for
IBM watsonx.data Development Lab - 1.1.2 R1M1	Apr 21, 2024	Ibmcloud 2: any, any, any IBM watsonx.data Development Lab - 1.1.2 R1M1
IBM watsonx.data Development Lab - 1.1.3 R1M0	Apr 23, 2024	Ibmcloud 2: any, any, any IBM watsonx.data Development Lab - 1.1.3 R1M0 This is the latest version of the <a href="#">watsonx.data development lab</a> . The
IBM watsonx.data Development POC R1M2	Apr 21, 2024	Ibmcloud 2: any, any, any IBM watsonx.data Development POC R1M2 Note: Updated POC image to fix the self-signed certificate error.

**Note:** You may see multiple tiles with different, albeit similar, versions. This lab guide is based on watsonx.data v1.1.3. If you use a different version of watsonx.data then you may experience issues due to changes in product behavior or in the user interface. Use the version specified above to avoid such issues.

4. For the **reservation type**, select the **Reserve now** radio button.



IBM Technology Zone | My library | Help

## Create a reservation

IBM watsonx.data Developer Base Image

Select a environment/infrastructure | Select a reservation type

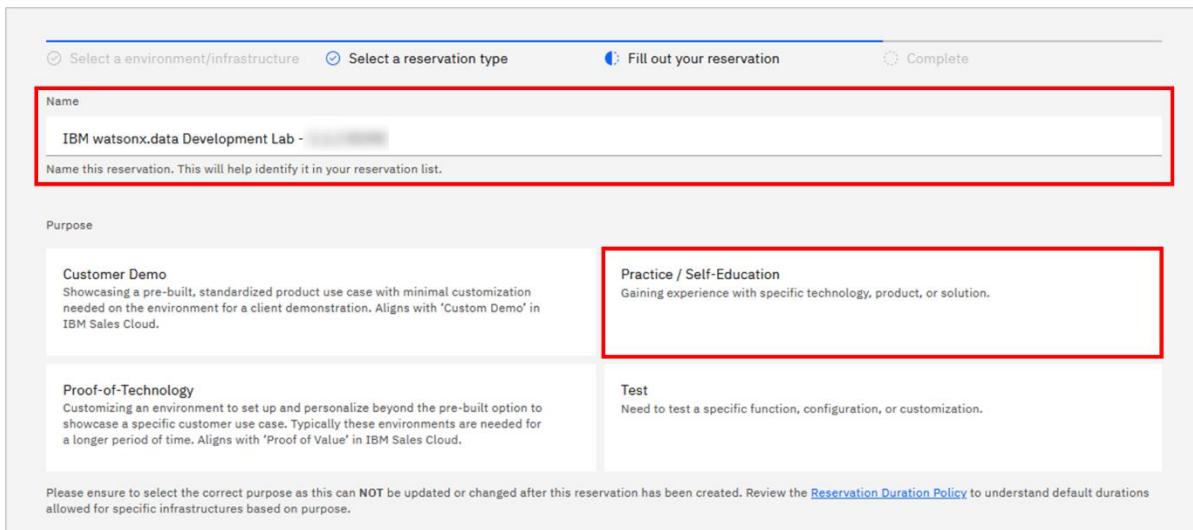
Select your reservation type. Do you need this now or later?

**Single environment reservation options:**

Reserve now

Schedule for later

5. Accept the default for the reservation **Name**, or provide a name of your choosing. For the **Purpose** of the reservation, select **Practice / Self-Education**.



Select a environment/infrastructure | Select a reservation type | Fill out your reservation | Complete

**Name**  
IBM watsonx.data Development Lab -

Name this reservation. This will help identify it in your reservation list.

**Purpose**

**Customer Demo**  
Showcasing a pre-built, standardized product use case with minimal customization needed on the environment for a client demonstration. Aligns with 'Custom Demo' in IBM Sales Cloud.

**Practice / Self-Education**  
Gaining experience with specific technology, product, or solution.

**Proof-of-Technology**  
Customizing an environment to set up and personalize beyond the pre-built option to showcase a specific customer use case. Typically these environments are needed for a longer period of time. Aligns with 'Proof of Value' in IBM Sales Cloud.

**Test**  
Need to test a specific function, configuration, or customization.

Please ensure to select the correct purpose as this can NOT be updated or changed after this reservation has been created. Review the [Reservation Duration Policy](#) to understand default durations allowed for specific infrastructures based on purpose.

- Fill in the **Purpose description** box (you may have to scroll down to see it) with the reason you are making the reservation. Then, scroll further down and select your **Preferred Geography** based on your location.

**Purpose description**

Self-education on watsonx.data to earn the watsonx.data Technical Sales Intermediate badge.

**Preferred Geography**

itzvmware - AMERICAS - any region - any datacenter

- Scroll down further. Adjust the reservation's **End date and time** as needed (by default it's two days (48 hours) from now; it can't exceed two days initially, but you can extend the reservation by two days, up to two times, for a total of 6 days, before it expires).

**End date and time**

Select a date      Select a time

07/29/2023      10:00 AM      America/Toronto

Reservation policy: Recommended 2 days, but can be reserved up to 2 days on this reservation form. Extend later for 2 days increments up to 4 days total. Max time 6 days total.

Note: There is a field to enable VPN access to the environment. It is not necessary to enable this.

- On the right-side panel, follow the links to read the **Terms & Conditions** and the **End User Security Policies** documents. Then, select the checkbox to agree to those terms. Finally, click **Submit**.

IBM watsonx.data Developer Base Image

IBM watsonx.data Development Lab - [View details](#)

Policy

Optimized by [IBM Turonomic](#)

I agree to IBM Technology Zone's [Terms & Conditions](#) and [End User Security Policies](#)

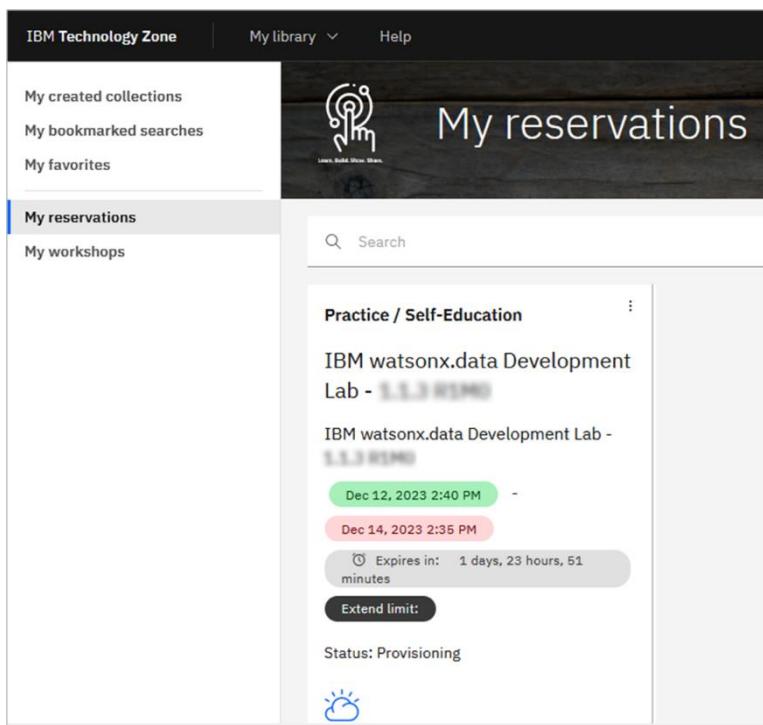
**Submit**

A message in the upper-right corner will briefly appear stating that the reservation has been created. You may also be presented with an opportunity to provide feedback on the process. Feel free to share your feedback.

Shortly after, you will receive an email from **IBM Technology Zone** acknowledging receipt of the request and you will receive another email when the provisioning is complete. Provisioning may be as quick as 15 minutes, or it may take an hour or more.

***NOTE:** Due to the massive popularity of watsonx, both for clients and for self-education, the TechZone environments occasionally reach their capacity. If you experience a problem trying to provision an instance, you can try again. However, you may have to wait until a later time when there is less demand. Please do not post provisioning-related issues to the #data-ai-demo-feedback Slack channel as the team that monitors this channel is not able to assist in these matters. A ticket can be opened with the TechZone team instead.*

Reservation status is available at <https://techzone.ibm.com/my/reservations>.



The screenshot shows the 'My reservations' section of the IBM Technology Zone interface. On the left, a sidebar lists 'My created collections', 'My bookmarked searches', 'My favorites', 'My reservations' (which is selected and highlighted in blue), and 'My workshops'. The main area is titled 'My reservations' and features a search bar. Below the search bar, a list of reserved environments is displayed under the heading 'Practice / Self-Education'. The first item in the list is 'IBM watsonx.data Development Lab - S.S.J RSMU', with a reservation period from 'Dec 12, 2023 2:40 PM' to 'Dec 14, 2023 2:35 PM'. A status message indicates 'Expires in: 1 days, 23 hours, 51 minutes'. There is a 'Extend limit:' button and a status message 'Status: Provisioning'. At the bottom of the list is a small icon of a computer monitor.

You must now wait until the environment has been provisioned before moving on to the next section. Specifically, you can continue on to the next section once you've received the **Reservation Ready on IBM Technology Zone** email from **IBM Technology Zone**.

## 4.2. Accessing the watsonx.data Environment

In this lab you will interact with the graphical web interfaces (consoles) for watsonx.data, Presto, and the MinIO Object Store. You will also run commands from a command line interface using ssh (Secure Shell).

**Note:** You do not have to use the VM Remote Console associated with your reservation (and it's recommended that you don't use it). You should be able to do everything directly from your laptop, using your laptop's terminal/command window and your own browser.

The console URLs and the ssh command are specific to your environment. They can be found in your TechZone reservation details.

1. Open the **Reservation Ready** on IBM Technology Zone email.
2. Click the **View My Reservations** button to view your TechZone reservations (you may have to login again).

The screenshot shows a web page titled "IBM Technology Zone" with a black header bar. Below the header, a blue bar displays "Status Update: Ready". The main content area contains the following information:

- Reservation Name:** IBM watsonx.data Development Lab - [\[REDACTED\]](#)
- Reservation ID:** [65b7c292d3b1ae0010b05c4d](#)
- Start Date:** 2024-01-29 15:21:00 (UTC Time)
- End Date:** 2024-03-29 15:55:30 (UTC Time)
- Collection:** <https://techzone.ibm.com/collection/6470dc59ddb55700174b6260>
- Environment Name:** IBM watsonx.data Development Lab - [\[REDACTED\]](#)

At the bottom of the page, there is a red rectangular box highlighting the **View My Reservations** button, which is located in a white box. Below the button, under "Additional Support Resources:", there is a bulleted list:

- [TechZone Status.io page](#) for System Status Updates
- Leverage our [Runbooks library](#) for documentation and troubleshooting guides
- Leverage our [Help page](#) for FAQs and other support resources
- Contact TechZone support team at [techzone.help@ibm.com](mailto:techzone.help@ibm.com)

3. The tile associated with your reservation opens and the tile should say **Status: Ready**. Click this tile.

The screenshot shows the 'My reservations' section of the IBM Technology Zone interface. On the left sidebar, 'My reservations' is selected. In the main area, a card for 'IBM watsonx.data Development Lab' is displayed. The card includes the lab name, a green button for 'Dec 12, 2023 2:40 PM', a red button for 'Dec 14, 2023 2:35 PM', a timer indicating 'Expires in: 1 days, 23 hours, 51 minutes', and a 'Extend limit:' button. At the bottom of the card, the status 'Status: Ready' is enclosed in a red box. Below the card is a small blue icon.

4. Scroll down to the **Published services** section. You will be using the **SSH for watsonx userid** command, the **Presto console**, the **MinIO console**, and the **Watsonx UI**. Make note of or bookmark these four pieces of information so that you can refer to them later.

The screenshot shows the 'Published services' section. A list of services and their URLs is provided, with several entries highlighted by red boxes and arrows:

- SSH for watsonx userid - <ssh -p 38241 watsonx@na4.services.cloud.techzone.ibm.com>
- MySQL - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 44510
- PostgreSQL Port - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 30533
- Portainer console - <http://na4.services.cloud.techzone.ibm.com:26069>
- SSH over Browser - <http://na4.services.cloud.techzone.ibm.com:47893>
- Apache Superset - <http://na4.services.cloud.techzone.ibm.com:45648>
- Presto console - <https://na4.services.cloud.techzone.ibm.com:30812>
- Presto Port - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 30812
- Jupyter Notebook - Server: [http://na4.services.cloud.techzone.ibm.com:20978/notebooks/Table\\_of\\_Contents.ipynb](http://na4.services.cloud.techzone.ibm.com:20978/notebooks/Table_of_Contents.ipynb)
- Minio Endpoint - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 23895
- Minio console - <http://na4.services.cloud.techzone.ibm.com:20636>
- Hive Thrift URL - <thrift://na4.services.cloud.techzone.ibm.com:39794>
- Hive Metastore - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 39794
- Watsonx UI - <https://na4.services.cloud.techzone.ibm.com:43196>
- Milvus Port - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 23990
- Db2 Port - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 34788
- Open Port 1 - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 43508
- Open Port 2 - Server: [na4.services.cloud.techzone.ibm.com](http://na4.services.cloud.techzone.ibm.com) Port: 48481

### 4.3. Command Line Access

Some of the activities performed in this lab require that you run commands within the VM. You will use ssh (Secure Shell) to connect from your computer into the VM and run commands.

**Note:** It is not necessary to run the commands in this section now, but you must read and familiarize yourself with the contents of this section now. The commands are shown for reference purposes, and you will need to run them whenever you are instructed later to open a new ssh session to the watsonx.data VM server. That said, feel free to try them out now if you would like.

ssh is a secured network protocol that allows command-line execution against remote servers. Your laptop should already have the ssh command included (for Windows, macOS, and Linux). If not, you will need to download an ssh client for your specific operating system (such as PuTTY on Windows or Hyper on Mac).

To ssh into your watsonx.data VM, open a terminal/command prompt window on your laptop and run the ssh command provided in your reservation details (see the **SSH for watsonx userid** line in the **Published services** section, per the instructions in section [4.2 Accessing the watsonx.data Environment](#)). The command will be in the following form (where watsonx is the username that you're connecting as and hostname/port# are specific to your environment):

```
ssh -p <port#> watsonx@<hostname>
```

The following command is an example of what your ssh command might look like. Do not run the command below; this is for example purposes only. You must use the ssh command provided in your reservation details instead, as it is specific to your environment.

```
ssh -p 99999 watsonx@geo.services.cloud.techzone.ibm.com
```

When prompted for the user's password, enter **watsonx.data** (you will not see the password as you type it).

The first time that you run the command you may receive a message stating that the authenticity of the host can't be established. In this case, respond stating that you want to proceed anyway.

The ssh command provided logs you in as the **watsonx** user. When you need to run commands as the **root** user (you will be instructed when that is the case), enter the following command to switch to the **root** user.

```
sudo su -
```

If prompted for the password of the watsonx user, enter **watsonx.data**.

If you want to close the session, you can enter the **exit** command – twice. The first one will exit you out of the root user's shell back to the watsonx user's shell. The second time will exit you out of the ssh shell completely.

## 4.4. watsonx.data Infrastructure Components

There are five infrastructure components that can be used/configured in watsonx.data:

- **Engines:** A query engine is used to run workloads against data in watsonx.data. Watsonx.data supports multiple engines, including Presto, Spark, Db2, and Netezza; this lab will use Presto as the query engine.
- **Catalogs:** Metadata catalogs are used to manage table schemas and metadata for the data residing in watsonx.data.
- **Buckets:** Watsonx.data stores data in object storage. Specifically, data is stored in buckets, which are identified storage areas within object storage, similar to file folders. AWS Simple Storage Service (S3), IBM Cloud Object Storage (COS), and MinIO object storage are supported.
- **Databases:** External databases such as IBM Db2, PostgreSQL, and MySQL can be registered and accessed by watsonx.data.
- **Services:** Other watsonx.data-related services, like the Milvus vector database.

**Note:** The term *schema* was used above – and this term will be used often in this lab. It has multiple meanings; in this context, a schema refers to the definition of a table, including the name and data type for each column. Metadata for each table is maintained in watsonx.data (specifically in the catalogs within the metadata store), including the table's schema and where its data files can be found in object storage, among other things.

This lab uses watsonx.data *Developer Edition*, which includes an embedded instance of *MinIO*, an open-source S3-compatible object store. However, you can also register and use external buckets that you have created elsewhere (such as in IBM Cloud or Amazon Web Services).

Developer Edition comes pre-configured with the following components:

- **presto-01 (engine):** This is a Presto query engine. It is used to interact with data in the lakehouse.
- **milvus000 (service):** This is an instance of the Milvus vector database.

- **iceberg\_data (catalog)**: This is an Iceberg catalog, residing within watsonx.data's embedded Hive Metastore (HMS). It manages tables that have been created with the Iceberg open table format. This catalog is associated with the iceberg-bucket object storage bucket.
- **hive\_data (catalog)**: This is a Hive catalog, also residing within the embedded HMS. This catalog is intended for use with non-Iceberg tables, where data is stored in files (such as Parquet, ORC, or CSV), but they are not using the Iceberg table format. This catalog is associated with the hive-bucket object storage bucket.
- **wxd\_system\_data (catalog)**: This is a Hive catalog, associated with the wxd-system bucket.
- **iceberg-bucket (bucket)**: This is a bucket in the embedded MinIO object store. The table data stored here is associated with the iceberg\_data catalog.
- **hive-bucket (bucket)**: This is a bucket in the embedded MinIO object store. The table data stored here is associated with the hive\_data catalog.
- **wxd-milvus (bucket)**: This is a bucket associated with the Milvus vector database.
- **wxd-system (bucket)**: This is a bucket used to hold diagnostic data such as query history and query event-related information for the Presto engine.

In this demo environment, you are limited in what you can add, edit, or delete. Because Developer Edition is being used, some tasks are prohibited. For example, you can't delete any of the pre-configured components and you can't add new query engines (Presto or other types).

When using a non-Developer Edition of watsonx.data, you are required to create/add some of these different components before you can proceed. This applies to both SaaS (the managed service) and user-managed deployments of watsonx.data.

You are likely asking yourself why the environment is set up with two different types of catalogs (Hive and Iceberg). Generally speaking, tables that you create and populate within watsonx.data should use the Iceberg open table format because they inherit beneficial capabilities like transactional consistency, schema and partition evolution, snapshots for time travel queries and rollback, and improved query efficiency. Iceberg was built to handle huge datasets and analytics workloads.

However, you may need to land existing data into watsonx.data's object storage that is in supported data file formats (like Parquet or ORC) but aren't using the Iceberg table format. In this case, you can land your data into a Hive catalog-managed bucket and then create a non-

Iceberg format table on top of it. You can then create a new Iceberg version of the table (managed by an Iceberg catalog) with the same table definition and populate it with the data from the original table. This is made easy by a common SQL statement called *CREATE TABLE AS SELECT* (or *CTAS* for short), which is supported by Presto. You will see this in action later in this lab. For large-scale conversion jobs, Spark can also be used.

## 4.5. Stopping and Starting watsonx.data

This lab environment is configured so that watsonx.data is started automatically. You should not need to perform any manual steps yourself to use watsonx.data.

However, if you find that you need to stop and restart the watsonx.data software, for whatever reason, you can follow the instructions below to do so. This is purely for reference, so *DO NOT DO THIS NOW.*

### 4.5.1. Stopping watsonx.data

1. Open an SSH session and connect to the watsonx.data VM with the **watsonx** user ID. Use the SSH command provided in your TechZone reservation details, which is in the following form:

```
ssh -p <port#> watsonx@<hostname>
```

If you receive a message stating that the authenticity of the host can't be established, respond stating that you want to proceed anyway.

2. When prompted for the password, enter **watsonx.data**.
3. Switch to the **root** user (you may be prompted to enter the password for the **watsonx** user; if so then enter **watsonx.data**).

```
sudo su -
```

4. Change the directory to the watsonx.data product binaries (or alternatively, prefix the stop and/or start commands shown later with this path).

```
cd /root/ibm-lh-dev/bin
```

5. Stop watsonx.data.

```
./stop
```

You will see output similar to the text below (this is an excerpt and not the full output).

```
Inspecting docker image lhconsole-ui before removal:  
lhconsole-ui still running. Removing lhconsole-ui...  
lhconsole-ui  
  
Inspecting docker image lhconsole-nodeclient before removal:  
lhconsole-nodeclient still running. Removing lhconsole-  
nodeclient...  
lhconsole-nodeclient-svc  
  
Inspecting docker image lhconsole-api before removal:  
lhconsole-api still running. Removing lhconsole-api...  
lhconsole-api  
...
```

6. Check the status of watsonx.data. If all of the watsonx.data components have been stopped then you will see no output to this command.

```
./status --all
```

#### 4.5.2. Starting watsonx.data

If you don't currently have a terminal open as the root user, follow Steps 1-4 from above to do so.

1. Start watsonx.data by running the following two commands.

```
export LH_RUN_MODE=diag
```

```
./start
```

It will take a few minutes for the various component containers to start. Output from the start command should look similar to what's shown below (this is an excerpt and not the full output):

```
-- starting container for ibm-lh-minio...
96ad814705c43487ea02044d8b4e5e682901caa39734258c929510b53a74f15e
-- starting container for ibm-lh-postgres...
6f20a833ae469089f63cae7c34419161e1cf0170bd24f354c4a4823e7356f981
Checking container: ibm-lh-postgres
Container ID for ibm-lh-postgres is: 6f20a833ae46
Postgres status is:
localhost:5432 - no response localhost:5432 - accepting
connections

-- starting container for ibm-lh-control-plane-prereq...
...
```

## 2. Check the status of watsonx.data:

```
./status --all
```

If watsonx.data has started successfully then you will see a number of containers running.

ibm-lh-hive-metastore	running	0.0.0.0:9083->9083/tcp, :::9083->9083/tcp
ibm-lh-minio	running	0.0.0.0:9000-9001->9000-9001/tcp, :::9000-9001->9000-9001/tcp
ibm-lh-postgres	running	0.0.0.0:5432->5432/tcp, :::5432->5432/tcp
ibm-lh-presto	running	0.0.0.0:8443->8443/tcp, :::8443->8443/tcp
lhconsole-api	running	3333/tcp, 8081/tcp
lhconsole-nodeclient-svc	running	3001/tcp
lhconsole-ui	running	7443/tcp, 0.0.0.0:9443->8443/tcp, :::9443->8443/tcp

**Note:** Even though the containers are running, the software components within those containers might still be initializing. For example, you may need to wait a few more minutes before you'll be able to bring up the watsonx.data graphical interface in a web browser.

## 5. Exploring the watsonx.data User Interface

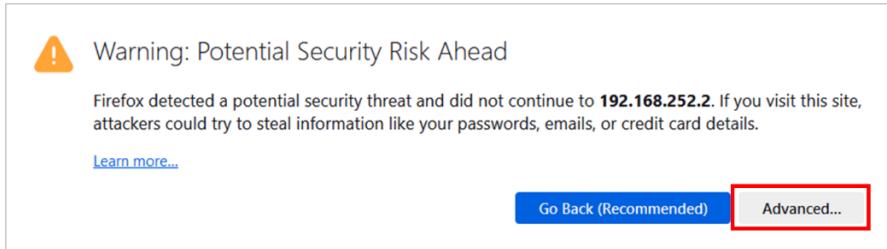
### 5.1. Starting the watsonx.data User Interface

Administration of the watsonx.data environment is primarily done with the watsonx.data user interface (also known as a console).

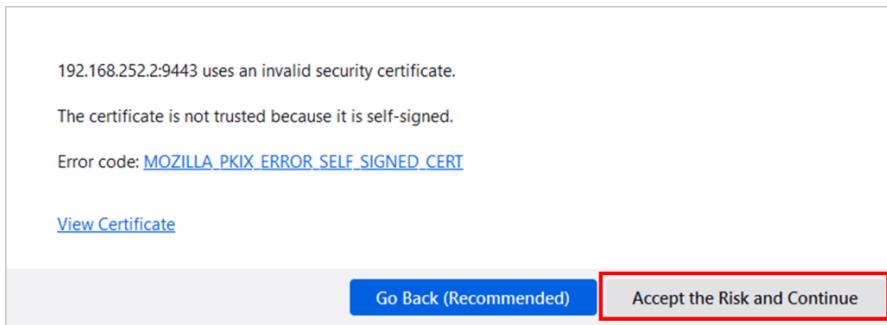
1. From your computer, open the **watsonx.data console** in your browser. The URL can be found in your TechZone reservation details (see the **Watsonx UI** line in the **Published services** section, per the instructions in section [4.2 Accessing the watsonx.data Environment](#)).

For users of the Firefox browser (steps 2 & 3):

2. You might receive a warning about a potential security risk. If so, click the **Advanced** button.

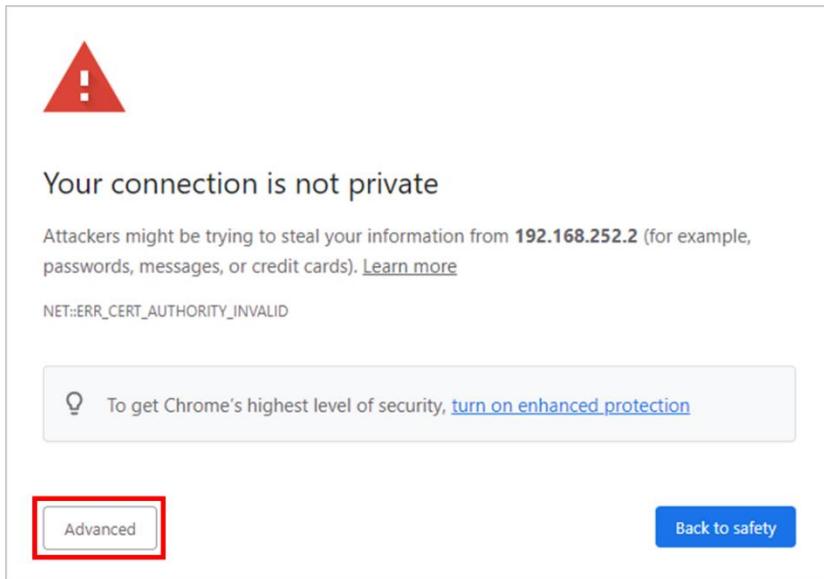


3. Scroll down the page and click the **Accept the Risk and Continue** button. Note that after doing this you will no longer see the warning when you open the console in the future.

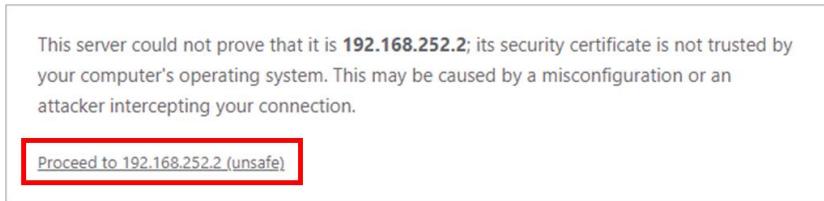


For users of the Chrome browser (steps 4 & 5):

4. You might receive a warning about the connection not being private. If so, click the **Advanced** button.



5. Scroll down the page and click the **Proceed to <URL> (unsafe)** link. Note that after doing this you will no longer see the warning when you open the console in the future.



All users should continue with the following steps:

6. In the IBM watsonx.data login screen, enter the following credentials:

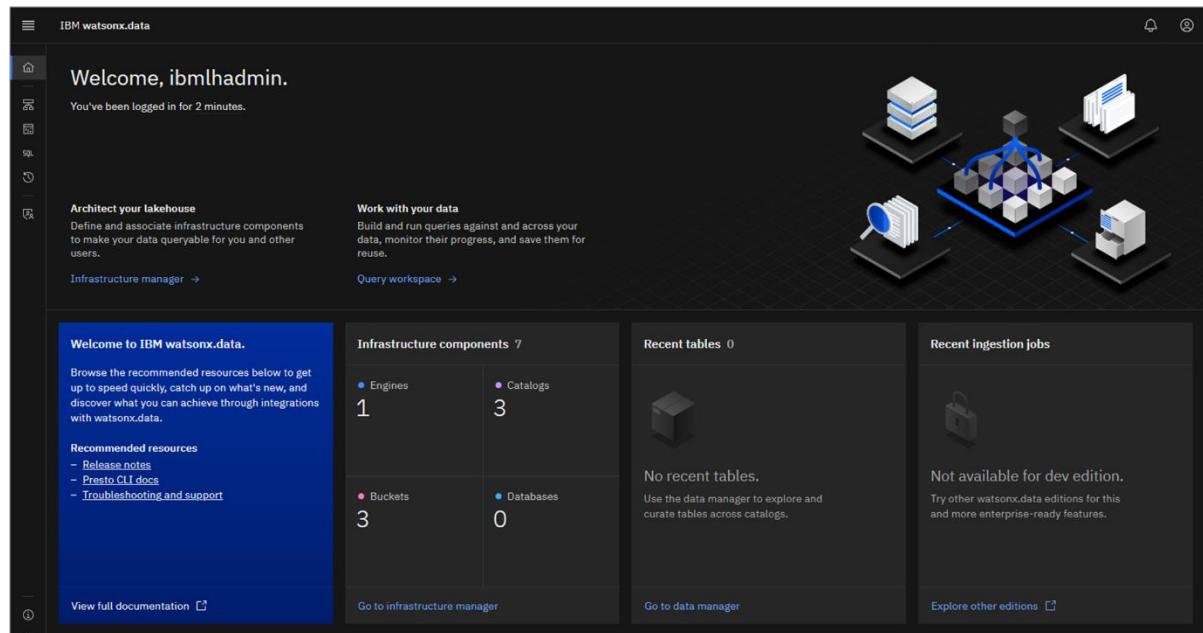
**Username:** ibmlhadmin

**Password:** password

7. Click the **Log in** button.



8. You will immediately start in the **Home** screen for watsonx.data. Scroll down and explore the contents of the page.



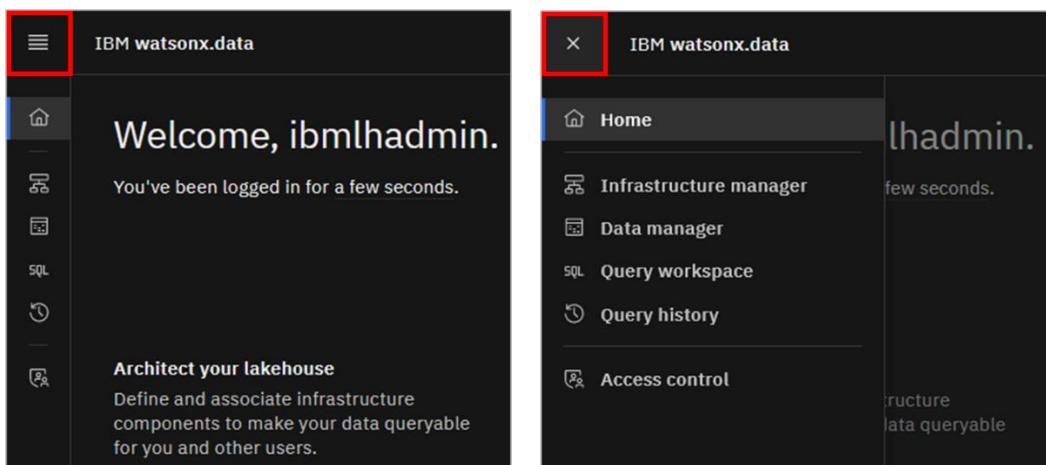
Included on the page are the following panels:

- **Welcome:** Introductory information, including a link to documentation
- **Infrastructure components:** A summary of the engines, catalogs, buckets, and databases that are registered with watsonx.data
- **Recent tables:** Tables that have recently been explored
- **Recent ingestion jobs:** Jobs that have recently moved data into watsonx.data
- **Saved worksheets:** Frequently run queries saved as worksheets, for easier reuse
- **Recent queries:** Queries that have recently been run or are in the process of being run

9. Note the left-side menu. Hover your mouse pointer over the various icons to see what actions or console pages they refer to.

-  **Home:** Returns you to the main Home screen
-  **Infrastructure manager:** Define and associate infrastructure components
-  **Data manager:** Browse schemas and tables by engine
-  **Query workspace:** Build and run queries against the data stored in the environment
-  **Query history:** Audit current and past queries across engines
-  **Access control:** Manage who can access infrastructure components and data

10. Alternatively, click on the **hamburger** icon in the upper left to expand the left-side menu such that you can see the name beside each icon. To collapse the menu back to the default view, click the **X** in the upper left.

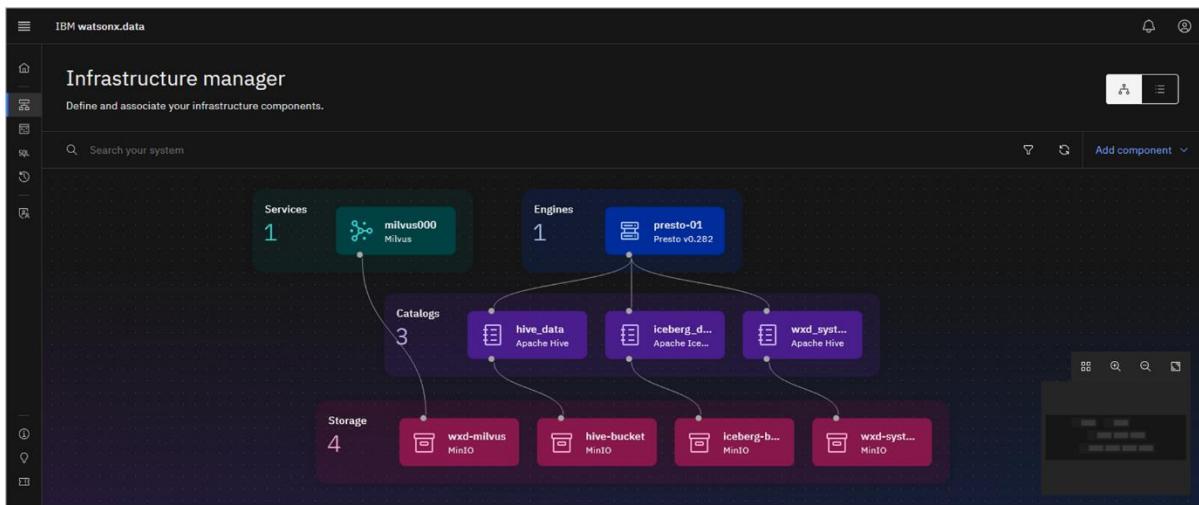


## 5.2. Infrastructure Manager Page

1. Select the **Infrastructure manager** icon (cloud) from the left-side menu.

The **Infrastructure manager** page opens with a graphical canvas view of the different infrastructure components currently defined in this watsonx.data environment: Engines (blue layer), Catalogs (purple layer), Buckets (magenta layer), and Databases (also blue, but not shown).

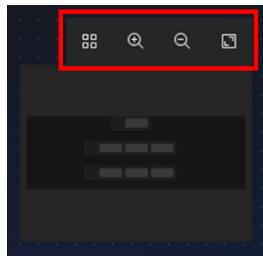
**Note:** Watsonx.data Developer Edition comes pre-configured with a Presto query engine, an instance of the Milvus vector database, three catalogs, and four object storage buckets.



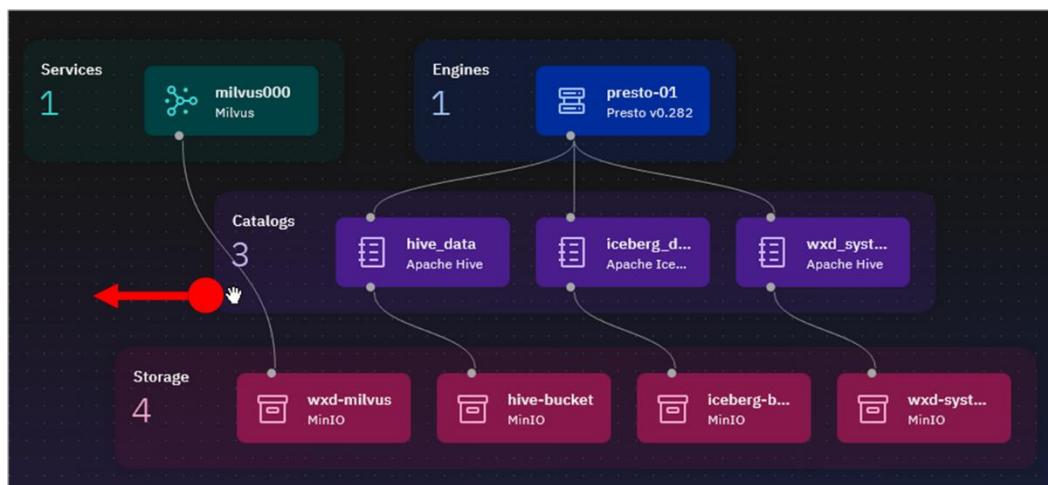
Each bucket (except in the case of the Milvus vector database service) is associated with a catalog (with a 1:1 mapping). When a bucket is added to watsonx.data, a catalog is created for it at the same time, based on input from the user. Likewise, if a database connection is added (for federation purposes), a catalog is created for that database connection as well. You will work with both object storage bucket catalogs and database catalogs later in the lab.

Each catalog is then associated with one or more engines. An engine can't access data in a bucket or a remote database unless the corresponding catalog is associated with the engine.

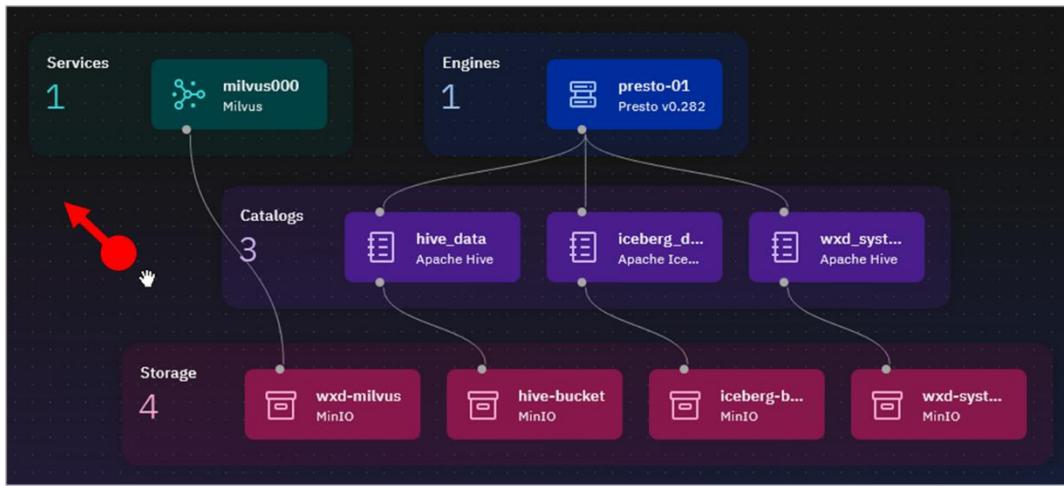
2. Note the mini-map in the lower right corner. This particular topology view is rather simple at this point, but as the number of infrastructure components grows, this control widget gives you a handy way to zoom in and out, auto-arrange the components, or fit the topology diagram to the screen. Try each of the icons to see how they affect the view.



3. Additionally, you can drag and drop the different infrastructure layers across the canvas. Click within the purple **Catalogs** area, hold the mouse button down, and move the catalogs to a different spot on the canvas.



4. Finally, you can pan across the canvas as a whole. Click somewhere in the black background of the canvas, hold the mouse button down, and move your mouse to drag the canvas.



5. In addition to the graphical topology view, infrastructure components can be listed in a table format. Click the **List view** icon in the upper-right to switch to this alternate view.

The screenshot shows the Infrastructure manager interface with the following details:

- Header:** IBM watsonx.data, Infrastructure manager, Define and associate your infrastructure components.
- Search Bar:** Search your system.
- Component Summary:** Engines 1, Services 1, Catalogs 3, Storage 4, Databases 0.
- Table View:**

Display name	Status	Type	Size	Catalogs associated
presto-01	Running	Presto v0.282	Custom	3
- Page Controls:** Items per page: 10, 1–1 of 1 items, 1 of 1 page.

A red box highlights the "List view" icon in the top right corner of the header area. Another red box highlights the component summary bar.

Tabs exist for each of **Engines**, **Services**, **Catalogs**, **Buckets**, and **Databases**. Explore the different tabs to see what information can be found there.

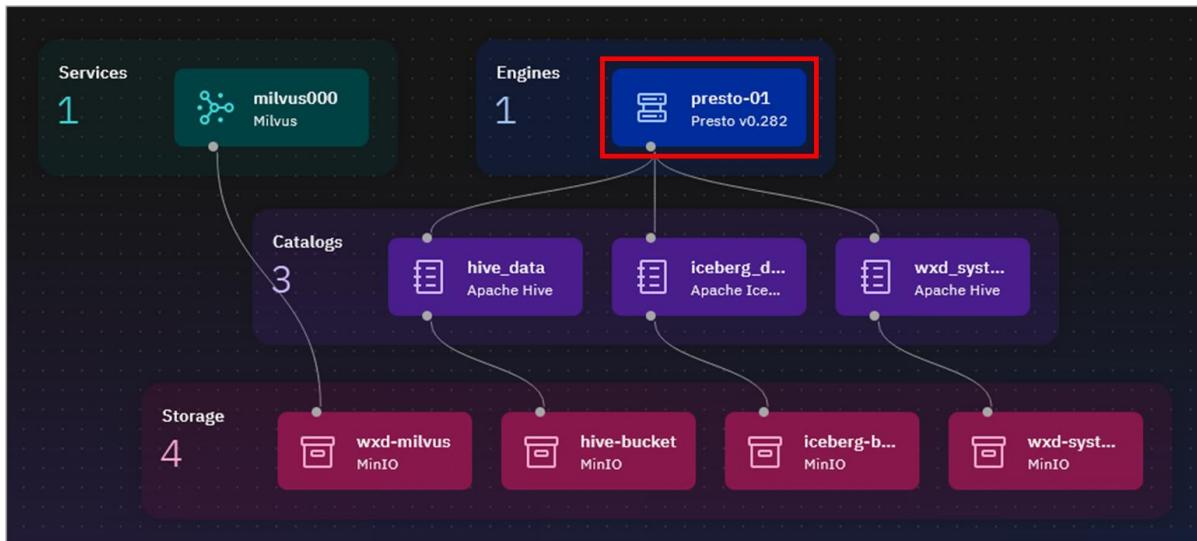
6. Click the **Topology view** icon (the icon to the left of the **List view** icon you just clicked) to switch back to the original graphical view.

The screenshot shows the 'Infrastructure manager' section of the IBM Watsonx.data interface. At the top, there's a navigation bar with icons for Home, Services, Catalogs, Storage, Databases, and a search bar. Below the navigation is a toolbar with icons for Add component, Refresh, and Help. A red box highlights the 'Topology view' icon (a grid icon) in the top right corner of the toolbar. The main area displays a table of components:

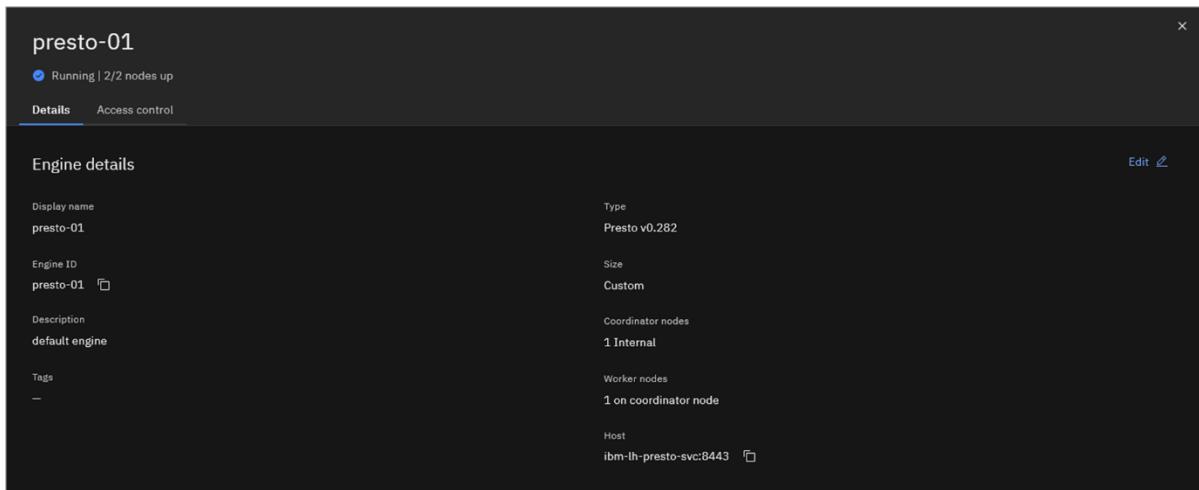
Display name	Status	Type	Size	Catalogs associated
presto-01	Running	Presto v0.282	Custom	3

At the bottom of the table are pagination controls: 'Items per page: 10', '1–1 of 1 items', and '1 of 1 page'.

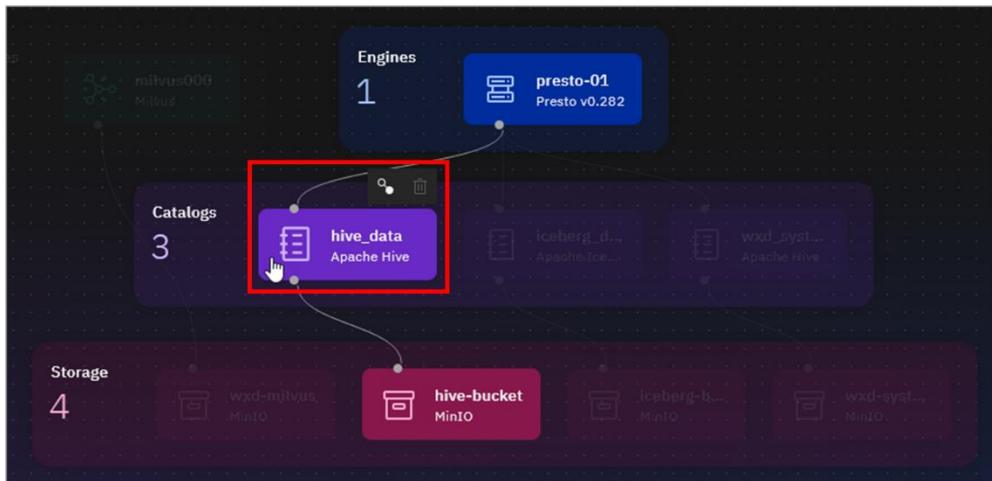
7. You can view details associated with each component. Click the **presto-01** engine tile.



Details including the Presto software version, the number of coordinator nodes, number of worker nodes, size, and host name are shown. Some of the values are editable (by clicking the **Edit** button – but don't click this now). However, as this is an instance of the watsonx.data Developer Edition, you can't do more than just change the name and description of this Presto engine.



8. Click the X in the upper right corner of the pane to return to the topology view.
9. Repeat the previous two steps for each of the catalogs and buckets, to see what information is available for them.
10. Hover your mouse pointer over the **hive\_data** catalog tile – but don't click on it.

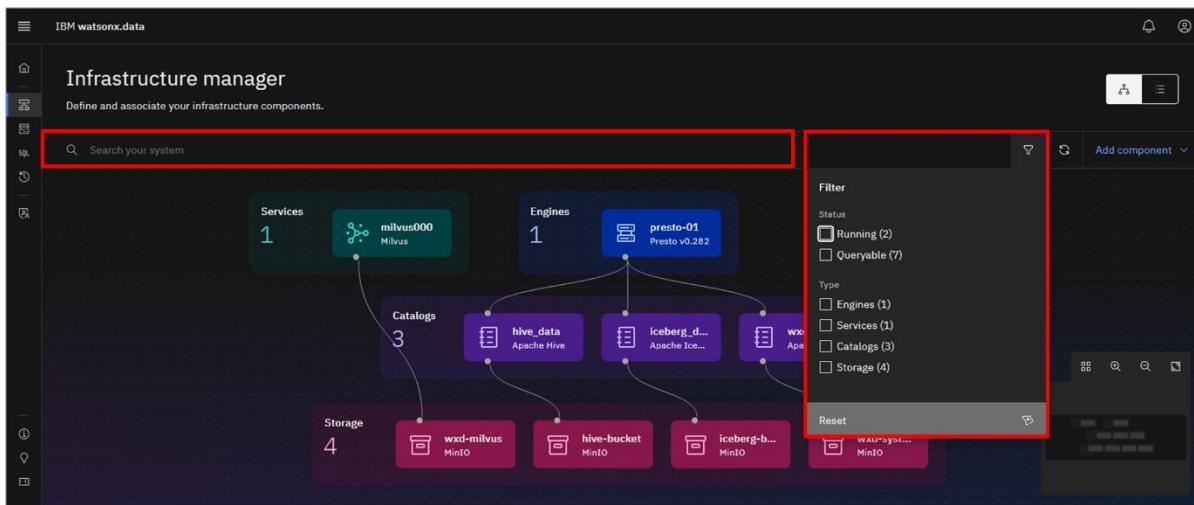


The catalog tile is highlighted, and icons appear above the tile. In this case there are two icons: **Manage associations** and **Delete**. The icons shown depend on the operations that the

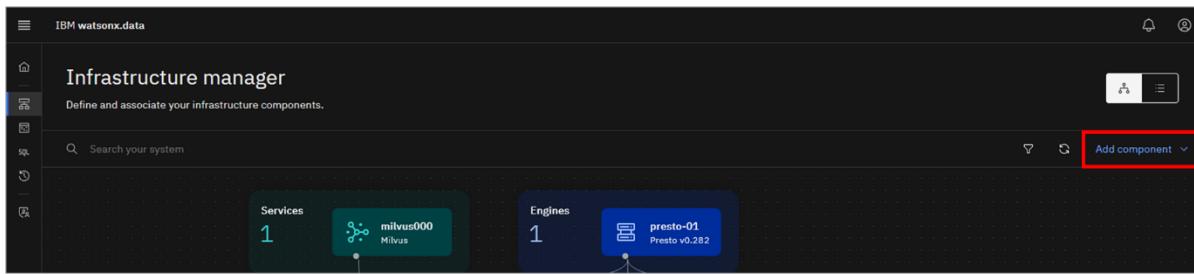
user can perform, which is based on the state of the component and the permissions that the user has on it.

**Note:** Because the **hive\_data** catalog is a default, system-managed catalog, you are not allowed to remove it and the **Delete** icon is grayed out.

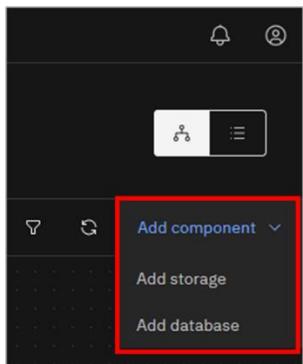
11. As the topology gets more complex, it may be difficult to find the components of interest. The console makes this easy by offering a search facility and the ability to filter what is shown (based on component type and/or the state of the component). Click the **Filter** icon to see the filter options available. Click the **Filter** icon again to close the filter options menu.



12. Click the **Add component** dropdown menu.



As this is the Developer Edition of watsonx.data, you are limited in what additional infrastructure components can be modified or added. You are not permitted to add additional engines or instances of Milvus, but you can add storage (buckets) and databases. The act of adding a new bucket or a database connection also adds an associated catalog, and so there is no explicit option for adding a catalog.



13. You will work with buckets later in this lab. For now, simply close the **Add component** dropdown menu.

### 5.3. Data Manager Page

The **Data manager** page can be used to explore and curate your data. It includes a data objects navigation pane on the left side of the page with a navigable hierarchy of *engine > catalog > schema > table*.

Earlier in this lab guide, the term *schema* was used to describe the definition of a table (including column names and types). Here, schema has a different, albeit related, meaning.

Presto (and many other database systems) organizes tables, views, and other database objects in schemas. Think of a schema as a logical collection or container of related database objects. For example, sales tables might be contained in one schema and marketing tables might be contained in another.

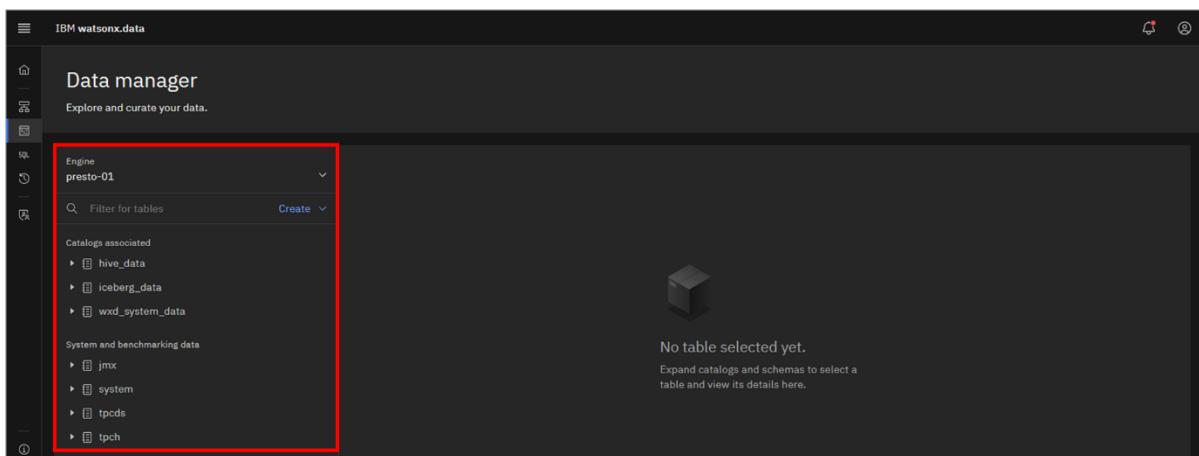
In Presto, tables can be identified in queries by specifying the catalog, schema, and table name (separated by periods). For example:

```
SELECT COUNT(*) FROM DATA_CATALOG.SALES_SCHEMA.CUSTOMER_TABLE
```

In this section you will explore the Data manager page and create your own schema and table.

1. Select the **Data manager** icon (☰) from the left-side menu.

The **Data manager** page opens with a data objects navigation pane on the left side:



The top-level navigation point is the query engine. You start by selecting an engine that is associated with the catalog and bucket you want to manage. As there is only one engine in this environment (**prestashop-01**), it is selected by default. If this was an environment with multiple Presto engines defined, you would have the choice of selecting any one of the engines you have

set up (as the administrator) or that you've been given access to (as a non-administrator).

With the engine selected, you can now navigate through the catalogs associated with the selected engine (the catalogs are listed in the **Catalogs associated** section on the left). Currently this includes the three default catalogs, but this is also where you would see any catalogs you explicitly associate with the engine.

**Note:** The web interface may automatically expand catalogs and schemas, as well as remember when catalogs and schemas were expanded previously. To simplify the screenshots in this lab, catalogs and schemas unrelated to the activity at hand are typically collapsed. So, what you see may not look exactly like what the screenshots in the lab show.

2. Expand each of the **hive\_data** and **iceberg\_data** catalogs by clicking on them.

The screenshot shows the Watsonx Data interface. At the top, it says "Engine" followed by "presto-01". Below that is a search bar with "Filter for tables" and a "Create" button. Underneath, there's a section titled "Catalogs associated" which lists three items: "hive\_data", "iceberg\_data", and "wxd\_system\_data". The "hive\_data" and "iceberg\_data" items are highlighted with red boxes around their entire row.

What do you see in these catalogs? The **iceberg\_data** catalog is empty and the **hive\_data** catalog has some schemas (with tables) in it. These pre-defined catalogs (and any catalogs you create) are empty until you create schemas and tables within them. However, some data has been added to this lab environment to provide you with interesting data to play with. These datasets are specific to this environment and are not included when you install `watsonx.data` yourself.

The screenshot shows the Watsonx Data interface with the "hive\_data" catalog expanded. The "Catalogs associated" section now shows "hive\_data" expanded, revealing three schemas: "gosalesdw", "ontime", and "taxi". The entire list of schemas is highlighted with a red box.

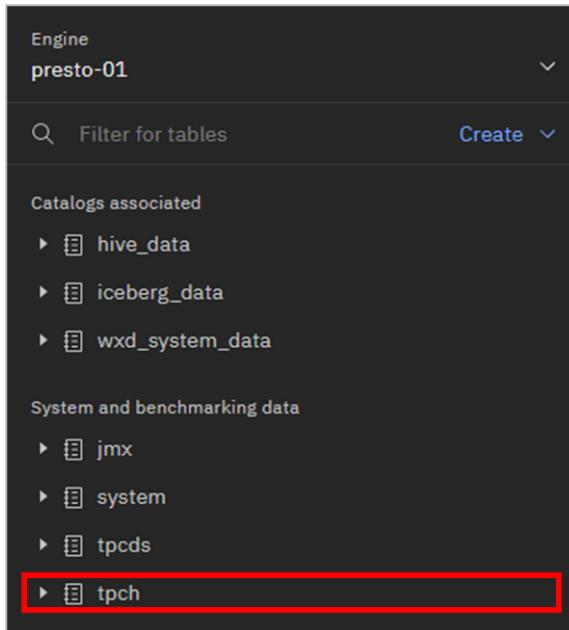
These three datasets are currently included:

1. **gosalesdw**: Sales data for the fictional Great Outdoors Company ([description](#), [schema](#)).
2. **ontime**: Airline reporting on-time performance dataset ([details](#)).
3. **taxis**: A subset of a Chicago taxi public dataset ([details](#)).

Additionally, Presto itself (and by extension watsonx.data) includes the *TPC-H* data set. TPC-H is a decision support benchmark maintained by the Transaction Processing Council (TPC). It is intended to mimic a real-world business workload involving a number of ad-hoc queries running while concurrent modifications are made to the data. The dataset that supports this benchmark includes information on customers, suppliers, orders, part numbers, and more. Datasets of different scale (size) can be generated to test workloads at different scale.

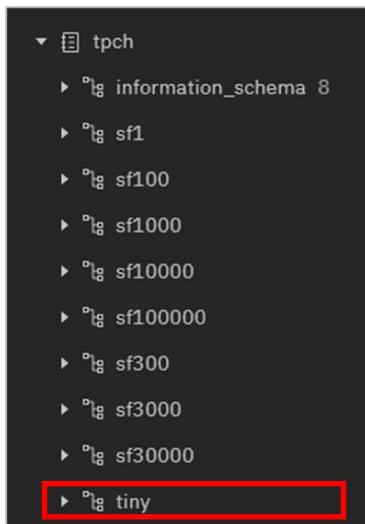
The TPC-H sample data (and other sample data) can be found in the **System and benchmarking data** section, which is below the **Catalogs associated** section.

3. Expand the **tpch** dataset by clicking on it.



Once expanded, the next level down is the schema. In the case of the TPC-H (tpch) data, each schema corresponds to a different scale factor of the dataset. The tables in the **sf100000** schema are 10x the size of the tables in the **sf10000** schema, which are 10x the size of the tables in the **sf1000** schema, and so on. The **tiny** schema is the smallest version of the dataset.

4. Expand the **tiny** schema.



5. Note the tables within this schema. Select the **customer** table. Information about this table is shown in the panel on the right. By default, the **Table schema** (table definition) tab is shown.

The screenshot shows the 'Data manager' interface for the 'tpch' catalog. On the left, a sidebar lists tables: sf3000, sf30000, tiny, lineitem, nation, orders, part, partsupp, region, and supplier. The 'customer' table under the 'tiny' schema is selected and highlighted with a red box. The main panel on the right shows the 'customer' table details. The 'Table schema' tab is active, displaying the following schema:

Name	Data type	Nullable
custkey	bigint	YES
name	varchar(25)	YES
address	varchar(40)	YES
nationkey	bigint	YES

Below the table schema, the message '1–8 of 8 items' is visible. At the bottom of the right panel, there are pagination controls: 'Items per page: 10', '1–8 of 8 items', and '1 of 1 page'.

6. Select the **Data sample** tab to see a sample of the data in the customer table.

The screenshot shows the WatsonX Data workspace interface. On the left, there is a navigation pane with a dropdown menu set to 'Engine presto-01'. Below it, a search bar says 'Search for loaded tables' and a 'Create' dropdown menu is open, with 'Create schema' highlighted. The main area shows a table named 'customer' from the 'Catalog: tpch | Schema: tiny'. The 'Data sample' tab is selected, displaying a sample of five rows from the table. The columns are labeled 'custkey', 'name', 'address', 'nationkey', 'phone', 'acctbal', and 'mktsegment'. The data includes various customer details such as names, addresses, and phone numbers, along with their respective nation keys, account balances, and market segments.

There are different ways that schemas and tables can be created in Presto. One way is through the use of SQL by running CREATE SCHEMA and CREATE TABLE SQL statements (which could be done in Presto's command line interface (CLI) or in watsonx.data's **Query workspace** page). Another approach is to use a third-party database management tool, such as [DBeaver](#).

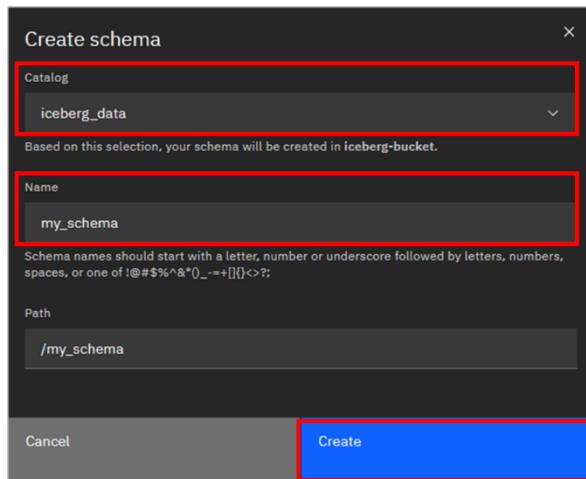
You can also use watsonx.data's **Data manager** page, which allows you to create a schema and upload a data file to define and populate it.

7. Go to the top of the left navigation pane and click the **Create** dropdown menu. Select **Create schema**.

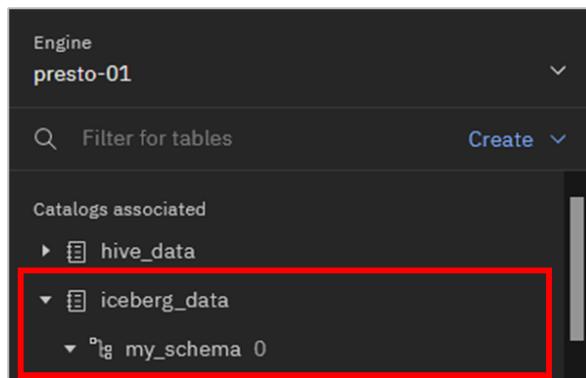
The screenshot shows the WatsonX Data workspace interface. On the left, there is a navigation pane with a dropdown menu set to 'Engine presto-01'. Below it, a search bar says 'Filter for tables' and a 'Create' dropdown menu is open, with 'Create schema' highlighted. The main area shows a section titled 'Catalogs associated' with two entries: 'hive\_data' and 'iceberg\_data'. The 'Create schema' option is highlighted with a red box.

8. In the **Create schema** pop-up window, select/enter the following information, and then click the **Create** button. The **Path** field (representing the directory folder that gets created in the object storage bucket) gets filled in automatically, based on the schema name being entered. You can override the path with your own custom path name, but leave the default as-is.

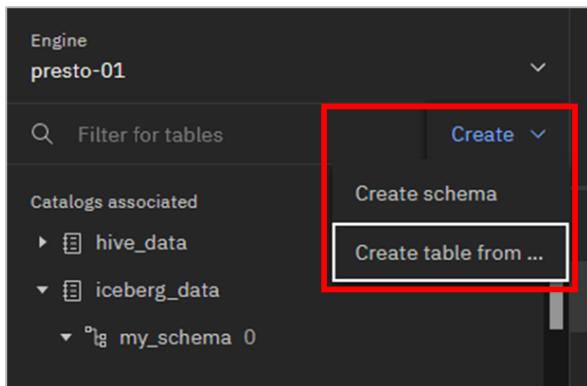
- **Catalog:** iceberg\_data
- **Name:** my\_schema



9. Expand the **iceberg\_data** catalog. The new schema should be listed (but contains no tables).



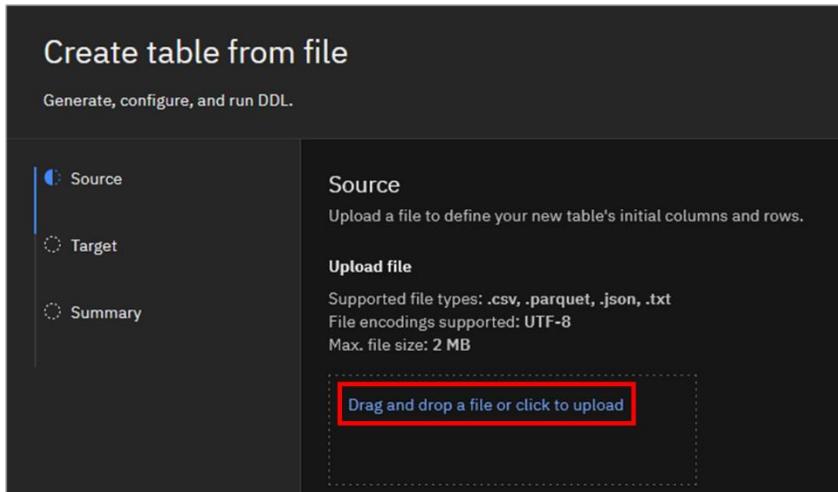
10. Click the **Create** dropdown menu again but this time select **Create table from file**.



The **Create table from file** workflow allows you to upload a small (maximum 2 MB file size) .csv, .parquet, .json, or .txt file to define and populate a new table.

11. Download the sample **cars.csv** file to your desktop (<https://ibm.box.com/v/data-cars-csv>).

12. For the **Source**, click **Drag and drop a file or click to upload**. Locate the **cars.csv** file you downloaded in the previous step and select it for upload (or simply drag and drop the file into this panel).



13. Scroll down to view a sample of the data uploaded. The schema of the table is inferred from the data in the file. Click **Next**.

Create table from file

Generate, configure, and run DDL.

Source

Target

Summary

Table schema configuration  
Confirm the data types to apply to each column.

cars.csv

Car	MPG	Cylinders	Displacement	Horsepower	Weight	Acceleration	Model Year
Chevrolet Chevelle Malibu	18	8	307	130	3504	12	70
Buick Skylark 320	15	8	350	165	3693	11.5	70
Plymouth Satellite	18	8	318	150	3436	11	70
AMC Rebel SST	16	8	304	150	3433	12	70

Cancel Back Next

14. For the **Target**, select/enter the following information (some fields are pre-populated and cannot be changed). Once filled in, click **Next**.

- **Catalog:** iceberg\_data
- **Schema:** my\_schema
- **Table name:** cars
- **Table format:** Apache Iceberg
- **Data format:** Parquet

Create table from file

Generate, configure, and run DDL.

Source

Target

Summary

Target  
Specify where your new table should be created. Some options are disabled but will be available in the future.

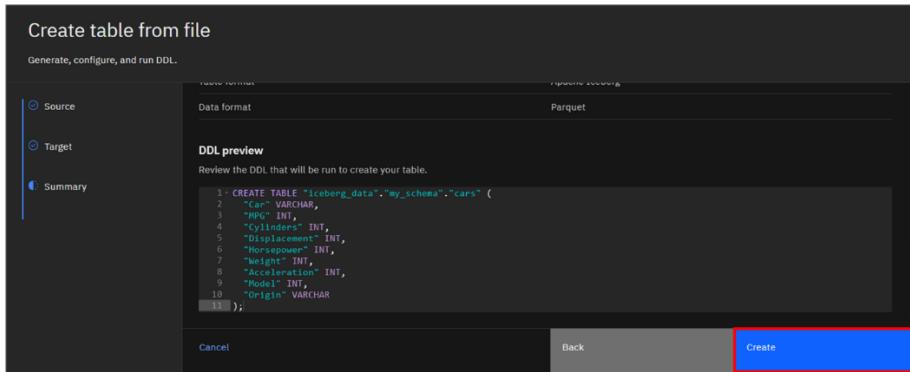
Catalog: iceberg\_data Schema: my\_schema

Table name: cars

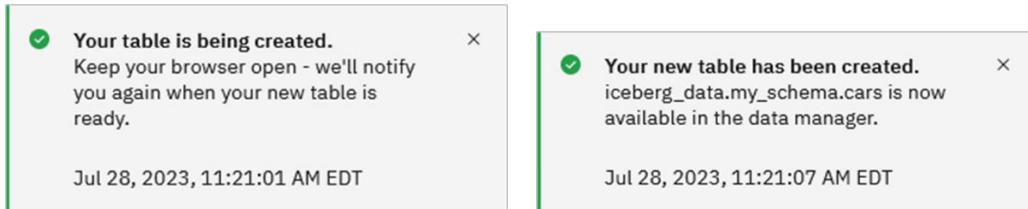
Table format: Apache Iceberg Data format: Parquet

Cancel Back Next

15. Scroll down to review the **Summary**, which includes the Data Definition Language (DDL) that will be used to create the table. You have an opportunity to alter the DDL statement if you wish, but do not change anything for this lab. Click **Create** to create the table.



You may see a pop-up message in the upper-right corner stating that the table is being created and another one after it completes. Not to worry if you miss them. Clicking on the **Notification** icon (looks like a bell) in the upper-right corner shows past notifications, including those related to creating tables – this icon will have a red dot if there is an unread notification.



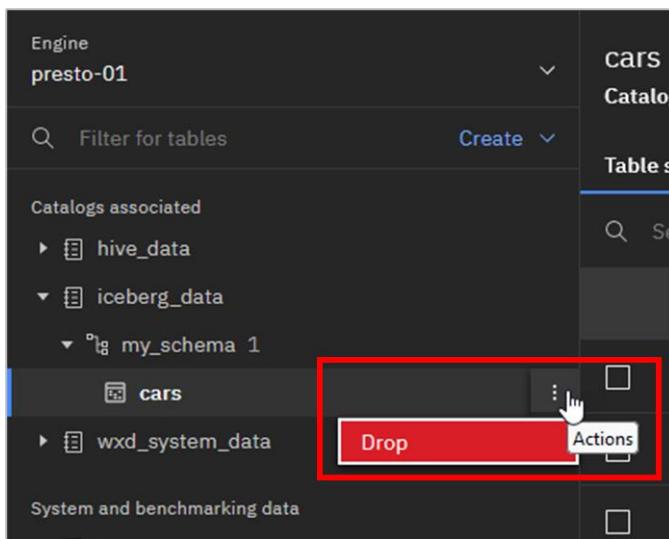
16. Navigate to your new table: `iceberg_data > my_schema > cars`.

The screenshot shows the Data manager interface. On the left, the navigation sidebar shows 'Catalogs associated' with 'hive\_data' and 'iceberg\_data'. Under 'iceberg\_data', there is a folder named 'my\_schema' containing a table named 'cars'. The main area displays the 'Table schema' for the 'cars' table. The schema table has columns for 'Name', 'Data type', and 'Nullable'. The data is as follows:

Name	Data type	Nullable
Car	varchar	YES
MPG	integer	YES
Cylinders	integer	YES
Displacement	integer	YES

At the bottom of the table, there are buttons for 'Items per page' (set to 10), '1–9 of 9 items', and navigation arrows. A blue 'Add column' button is located at the top right of the table schema table.

17. Explore the **Table schema**, **Data sample**, and **DDL** tabs. Notice a new tab called **Time travel** that you did not see with the TPC-H dataset earlier. You may view this now, but don't do anything in there. This topic will be covered later.
18. It is not immediately obvious, but there are additional menu options available for the catalogs, schemas, and tables in the navigation pane. Hover your mouse pointer at the far right of the line for the **cars** table. An **overflow menu** icon (vertical ellipses) appears. Click the **overflow menu** icon to see the option to drop the table (don't do that now). Click the **overflow menu** icon again to close it.



## 5.4. Query Workspace Page

Databases and query engines such as Presto have multiple ways that users can interact with the data. For example, there is usually an interactive command line interface (CLI) that lets users run SQL statements from a command terminal. Also, remote applications can use JDBC (Java Database Connectivity) to connect to the data store and run SQL statements.

The watsonx.data user interface includes an SQL interface for building and running SQL statements. This is called the **Query workspace**. Users can write or copy in their own SQL statements, or they can use templates to assist in building new SQL statements.

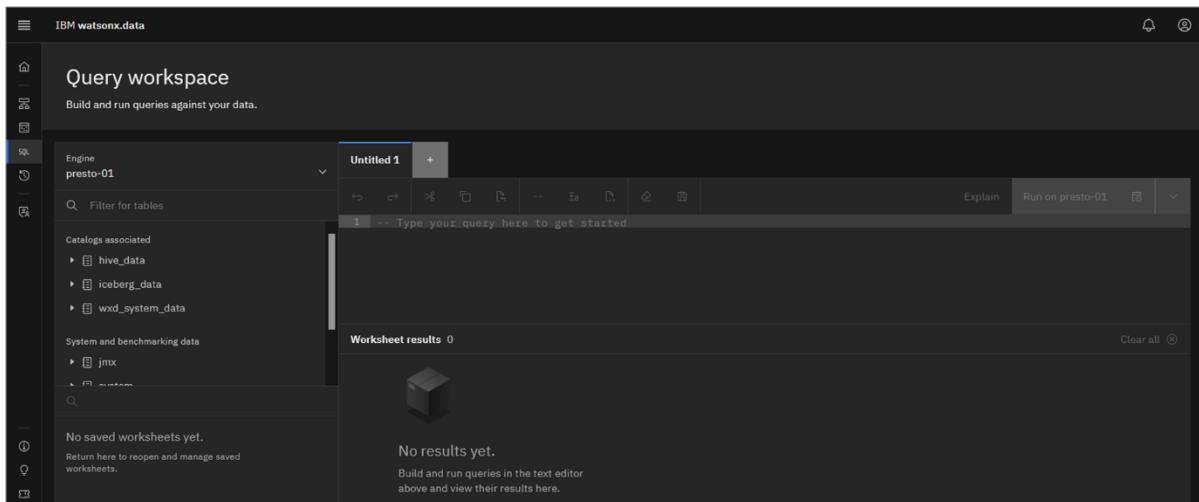
**Note:** Using Presto in watsonx.data, you can use SQL to insert new records into tables, but you aren't currently able to delete or update rows. This is currently possible with Spark SQL, though.

In addition to one-time execution of SQL statements, SQL statements that need to be run repeatedly can be saved to a *worksheet* and be run as often as needed later.

1. Select the **Query workspace** () icon from the left-side menu.

The **Query workspace** page opens with a data objects navigation pane on the left side and a SQL editor (workspace) pane on the right side.

Like in the **Data manager** page, the top of the navigation pane directs you to select an engine to use. It is this engine that will be used to run the SQL statements entered here. The only engine that exists in this lab environment is the **presto-01** Presto engine, which is selected by default.



2. Copy and paste the following text into the **SQL worksheet**. Note that the table you are about to query is identified by a 3-part name that includes the catalog, schema, and table name. Click **Run on presto-01**.

```
select car, avg(mpg) as avg_mpg from iceberg_data.my_schema.cars  
group by car order by car;
```

The screenshot shows the Presto SQL worksheet interface. At the top, there is a toolbar with various icons. Below the toolbar, a text editor window contains the SQL query: "select car, avg(mpg) as avg\_mpg from iceberg\_data.my\_schema.cars group by car order by car;". A red box highlights the "Run on presto-01" button, which is located to the right of the query text. Below the text editor is a section titled "Worksheet results 0". It features a small 3D cube icon and the message "No results yet." followed by the instruction "Build and run queries in the text editor above and view their results here.".

The query result is displayed at the bottom of the panel.

```
select car, avg(mpg) as avg_mpg from iceberg_data.my_schema.cars group by car order by car;
```

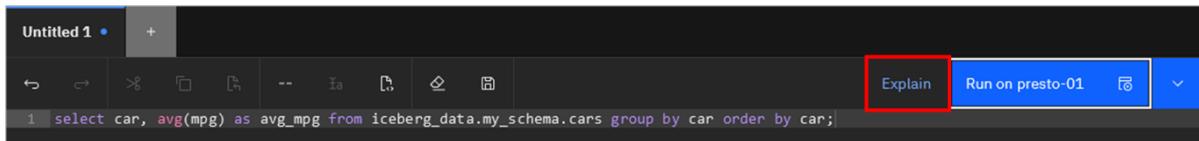
The screenshot shows the Presto SQL worksheet interface after the query has been run. The "Run on presto-01" button is highlighted with a red box. Below the results section, the "Worksheet results 1" header is visible, along with a "Clear all" button. The results table is shown with two rows: "AMC Ambassador Brougham" and "AMC Ambassador DPL", both having an "avg\_mpg" value of 13. The entire results section is highlighted with a red box.

car	avg_mpg
AMC Ambassador Brougham	13
AMC Ambassador DPL	15

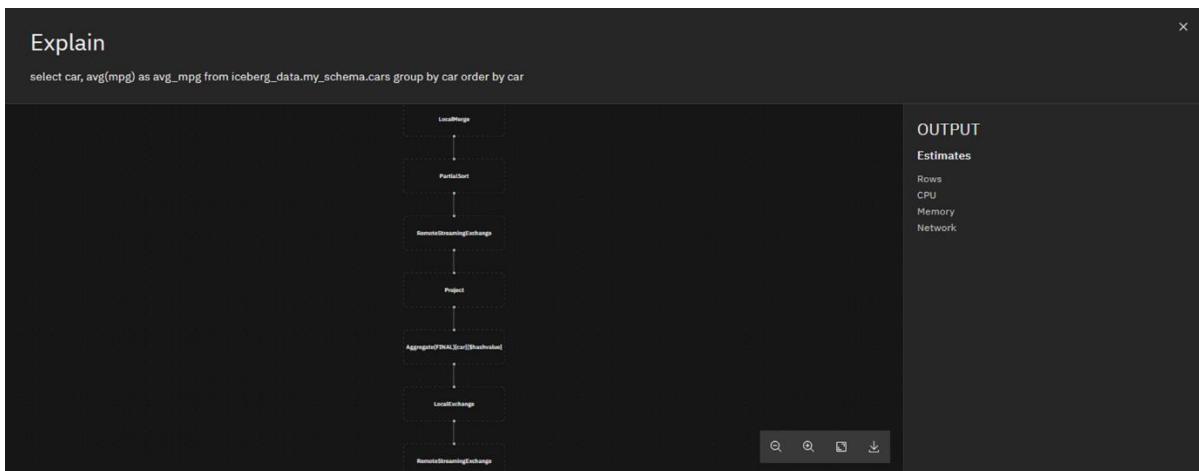
When people are developing or tuning an SQL query, they often want to understand all of the underlying work that is involved in running it – this is commonly referred to as the query's *execution plan*. Presto includes an *explain* facility that shows how Presto breaks up and distributes tasks needed to run a query. A graphical representation of a query's execution plan –

what's commonly referred to as *visual explain* – is available from within the Query workspace page. Watsonx.data uses an EXPLAIN SQL statement on the given query to create the corresponding graph, which can be used to analyze and further improve the efficiency of the query.

3. Click the **Explain** button.



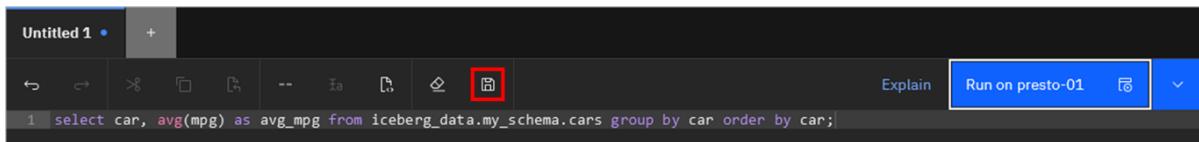
This shows a graphical representation of the query's execution plan. For this example, the query execution plan is relatively simple. Feel free to zoom in and explore the different elements of it. Clicking on a particular stage may show additional information in the pane on the right.



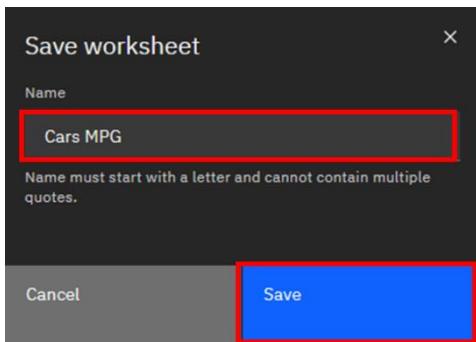
4. Click the X in the upper-right corner of the Explain window to close it.

Let's assume this is a query you will want to run time and time again in the future. To do this, you can save it as a worksheet. While this example just has a single SQL statement, imagine a scenario where you have multiple statements in the workspace editor and in a worksheet.

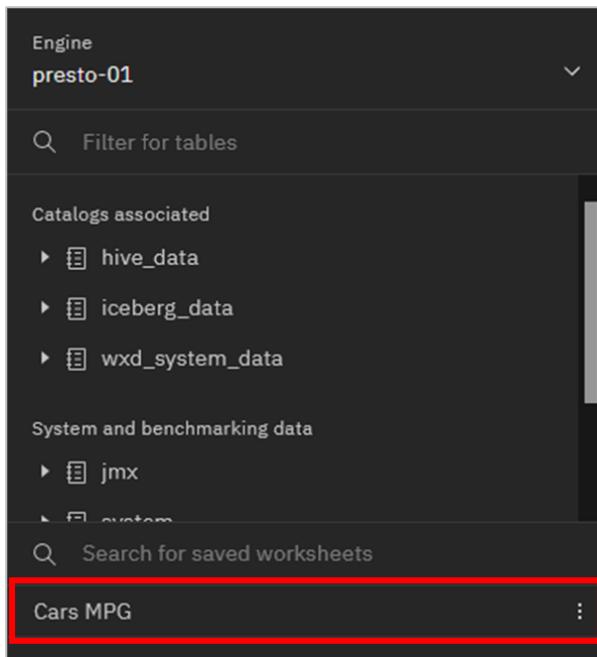
5. Click the **Save** icon (looks like a disk) in the editor menu above the SQL statement.



6. In the **Save worksheet** pop-up window, enter **Cars MPG** for the **Name** and then click **Save**.



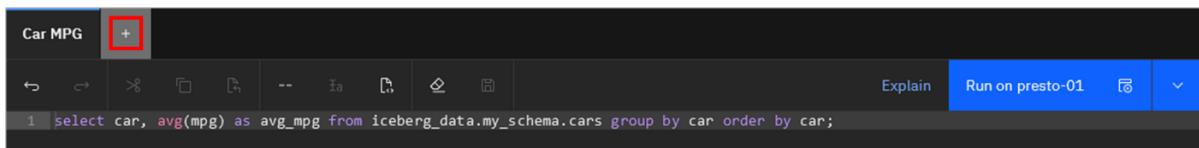
The worksheet is now displayed at the bottom of the left-side navigation pane.



The worksheet can be opened and run in the future as needed, simply by clicking on it (it will open a new SQL tab with the name of the worksheet; if you try to do it now, though, nothing will happen because the tab for this worksheet is already open).

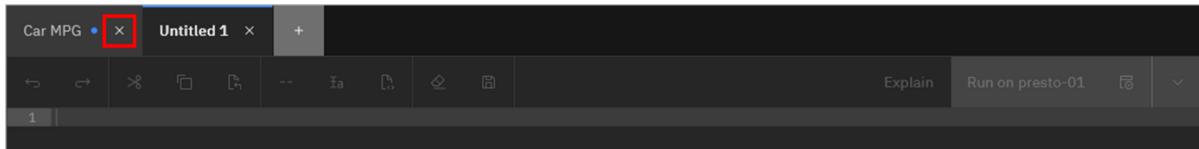
Any worksheet can be deleted at any time by clicking on the **Overflow** icon (vertical ellipses) to the right of the worksheet name and selecting **Delete**.

7. Click the **+** (New worksheet) icon at the top of the current worksheet to create a new, empty worksheet.



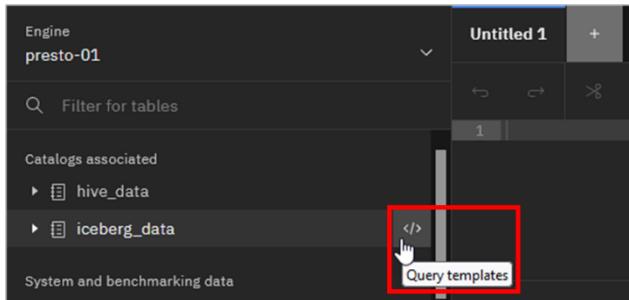
```
1 select car, avg(mpg) as avg_mpg from iceberg_data.my_schema.cars group by car order by car;
```

8. Click the **X** in the **Car MPG** tab to close that worksheet. If asked to confirm closing, click **Close**.

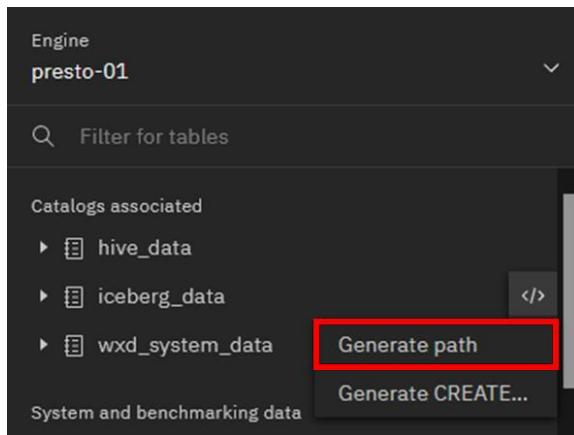


In addition to writing queries from scratch or copying and pasting queries from elsewhere, the **Query workspace** interface can also assist in generating SQL for tables that are in `watsonx.data`. Let's see what options are available here.

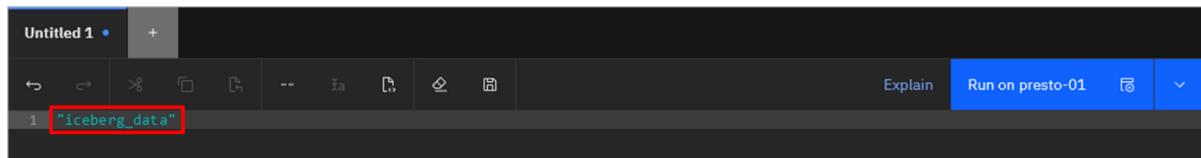
9. On the left-side navigation pane, hover your mouse pointer at the far right of the line for the `iceberg_data` catalog until the **Query templates (</>)** icon appears. When you see this icon, click it.



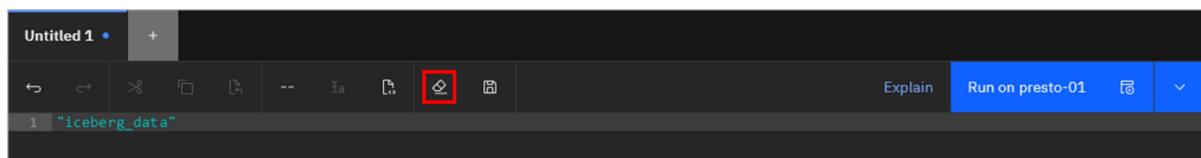
10. There are two templates listed: **Generate path** and **Generate CREATE SCHEMA**. Click **Generate path**.



Notice how the name of the catalog was entered into the worksheet. The same thing can be done for schemas, and for schemas the text that gets entered into the worksheet is in the form of “catalog”.”schema”.

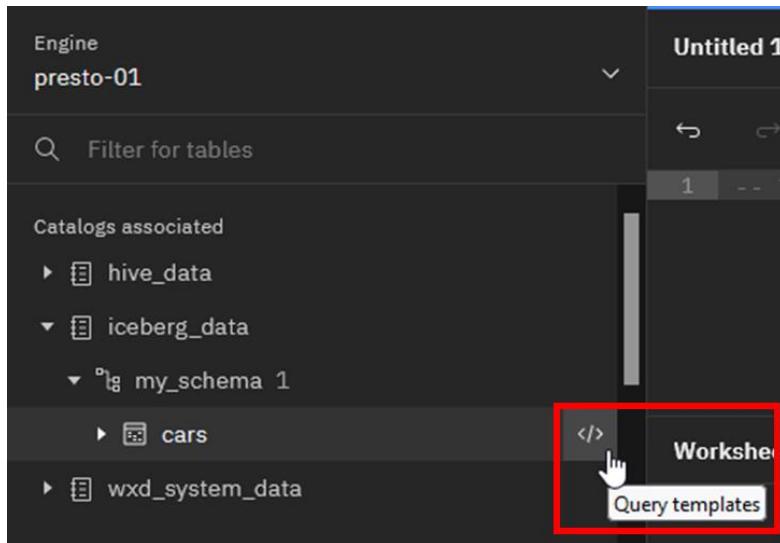


11. Click the **Clear** icon (looks like an eraser) in the menu above the SQL statement. This clears the text that was previously entered.

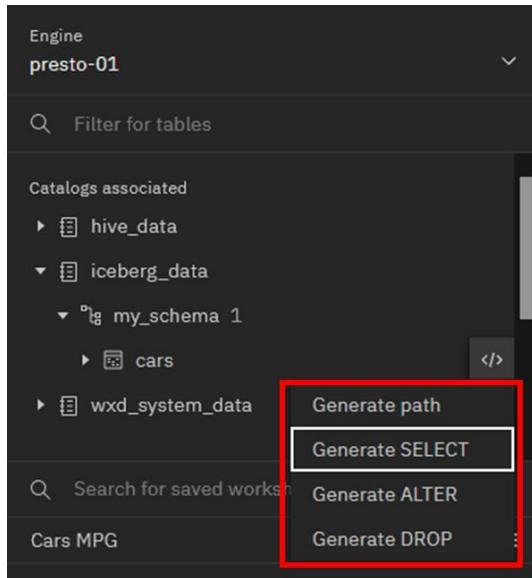


**Note:** If at any point in this lab you are instructed to copy and paste a SQL statement into a worksheet and there is already text there from a previous step you took, the implication is that you should clear out the old text first.

12. Navigate to **iceberg\_data > my\_schema > cars**. Hover your mouse pointer at the far right of the line for the **cars** table until you see the **Query templates (</>)** icon appear. When you see this icon, click it.



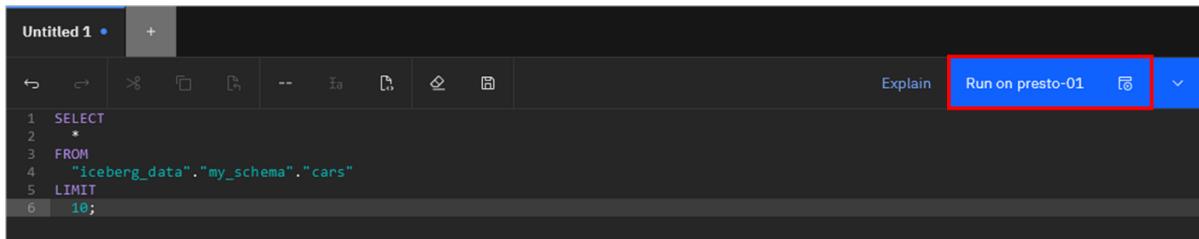
13. As you can see, more templates are available when working with tables than what you saw with the catalog:



Selecting **Generate path** generates the table name, as a 3-part name: "catalog"."schema"."table". The other three options can be used to generate SELECT, ALTER, and DROP SQL statements.

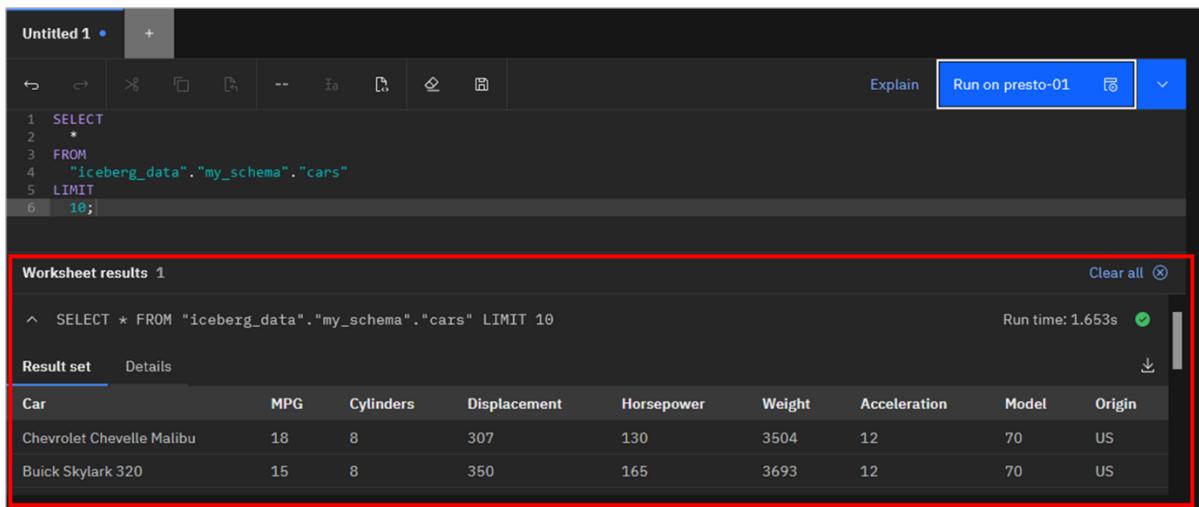
14. Click **Generate SELECT**.

15. A basic SELECT statement querying from the cars table is entered into the worksheet. This default statement returns 10 rows of the table including all columns. Click **Run on presto-01**.



```
1 SELECT
2 *
3 FROM
4 "iceberg_data"."my_schema"."cars"
5 LIMIT
6 10;
```

As before, the result of the query is shown at the bottom of the screen.



Worksheet results 1

^ SELECT \* FROM "iceberg\_data"."my\_schema"."cars" LIMIT 10 Run time: 1.653s

Result set	Details
Car	MPG Cylinders Displacement Horsepower Weight Acceleration Model Origin
Chevrolet Chevelle Malibu	18 8 307 130 3504 12 70 US
Buick Skylark 320	15 8 350 165 3693 12 70 US

## 5.5. Query History Page

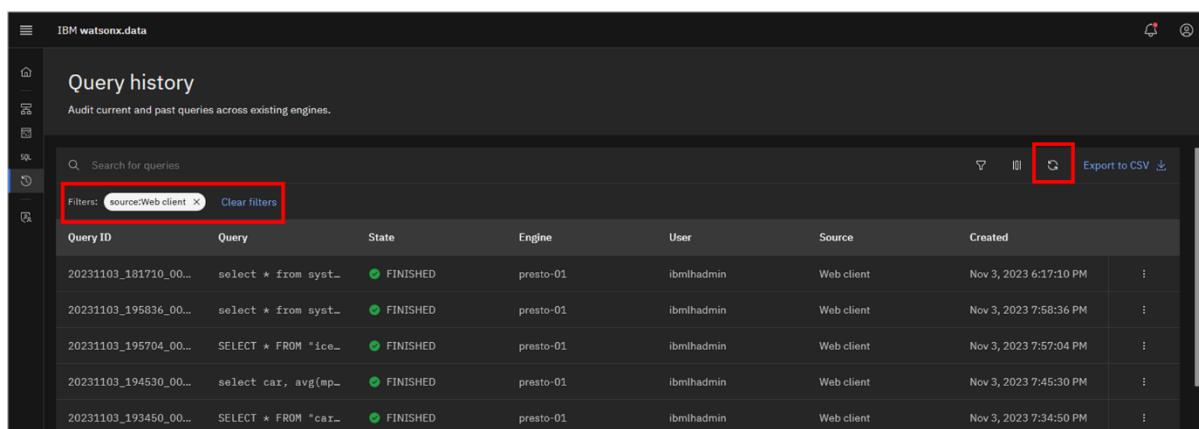
The **Query history** page lets you audit currently running queries and queries that ran in the past, across all the engines defined in the environment. This includes queries that users have explicitly run, as well as queries used in the internal management and function of the environment.

This section is intended to show some of the different ways you can interact with this page.

**Note:** What you see in your list of queries will not match the screenshots shown here.

1. Select the **Query history** (⌚) icon from the left-side menu.

The **Query history** page opens with a list of queries currently running and that ran in the past. If the list appears short or not current, click the **Refresh** icon in the upper-right. Also, by default only queries executed in (or by) the console are included. The filters can be cleared to see all queries sources, including the Presto CLI.



Query ID	Query	State	Engine	User	Source	Created
20231103_181710_00...	select * from syst...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 6:17:10 PM
20231103_195836_00...	select * from syst...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:58:36 PM
20231103_195704_00...	SELECT * FROM "ice...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:57:04 PM
20231103_194530_00...	select car, avg(mp...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:45:30 PM
20231103_193450_00...	SELECT * FROM "car...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:34:50 PM

Included (at the top) is a search bar to find specific queries of interest. The list can also be filtered by the state of the query (for example, FAILED, FINISHED, and RUNNING), the engine that ran the query, and the user that submitted the query. Information including the query text, state, and when the query was run is shown, but you can customize the columns to show more or less information.

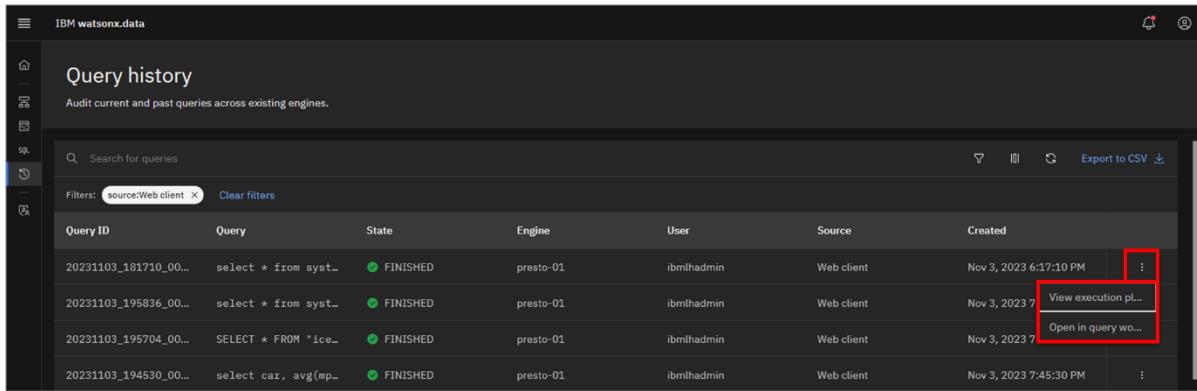
2. Click the **Filter** icon on the right (looks like a funnel). Note the filters available. Click the **Filter** icon again to close it.

The screenshot shows the 'Query history' page from IBM Watsonx.data. On the right side, there is a sidebar titled 'Filter' with several sections: State (FINISHED 32, RUNNING 1), Engine (presto-01 33), User (ibmlhadmin 33), Source (Web client 17, System 15, presto-cli 1), and Created (date range selector). A red box highlights the 'Filter' icon at the top of the sidebar.

3. Click the **Customize columns** icon (looks like a series of vertical lines; it's to the right of the **Filter** icon). Note the columns available to display. Click the **Customize columns** icon again to close it.

The screenshot shows the 'Query history' page with the 'Customize columns' sidebar open. The sidebar lists available columns: Query ID, Query, State, Engine, User, and Source. Under 'Source', 'Web client' is selected. A red box highlights the 'Customize columns' icon at the top of the sidebar.

4. Click the **overflow menu** icon (vertical ellipses) at the end of the row for any queries listed. Note the options available. You can view the explain plan for the query and you can bring up the query in the Query workspace. Click the **overflow menu** icon again to close the list.



Query ID	Query	State	Engine	User	Source	Created
20231103_181710_00...	select * from syst...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 6:17:10 PM
20231103_195836_00...	select * from syst...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:45:30 PM
20231103_195704_00...	SELECT * FROM "ice...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:45:30 PM
20231103_194530_00...	select car, avg(mp...	FINISHED	presto-01	ibmlhadmin	Web client	Nov 3, 2023 7:45:30 PM

**Note:** The query history is retrieved directly from Presto, which stores this history information in its `system.runtime.queries` table. However, this table is cleared when Presto is restarted, which means that the query history is lost whenever the engine shuts down. If users want their query history to be maintained over Presto engine restarts, they can choose to create their own query history table and populate it periodically with data from the system table.

## 5.6. Access Control Page

The **Access control** page is used to manage infrastructure access and data access policies.

Security and access control within watsonx.data are based on roles. A role is a set of privileges that control the actions that users can perform. Authorization is granted by assigning a specific role to a user, or by adding the user to a group that has been assigned one or more roles.

Access control at the infrastructural level allows permissions to be granted on the engines, catalogs, buckets, and databases. Roles for these components include Admin, Manager, User, Writer, and Reader (depending on the component).

Access to the data itself is managed through data control policies. Policies can be created to permit or deny access to schemas, tables, and columns.

User account management and access management varies between the different deployment options for watsonx.data. For instance, in the managed cloud service (SaaS), the service owner would need to invite other users to the environment and give them appropriate service access. With the standalone software, users can be added within the console's **Access control** page. In the Developer Edition, users can be added using a command line tool.

In this section you will add a new user and provide them with privileges over the infrastructure and data. Note that it's not the intention of this lab's instructions to show the results of these privileges (you will not be logging in with other users), the intention is to highlight the process of how you would assign these privileges in the first place.

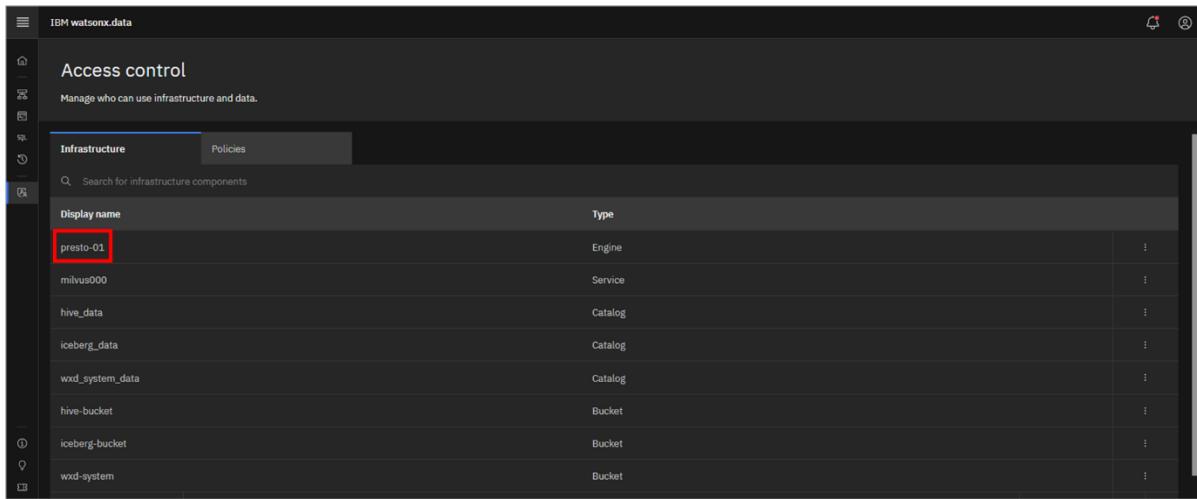
1. Open a terminal command window to the watsonx.data server as the **root** user (remember to use the `sudo` command to become the root user, or you will receive a permission denied error when running the command in the next step). (Refer to section [4.3 Command Line Access](#) for a reminder on how to open a new command/terminal window on your laptop.)
2. Run the following command to create a non-administrator user (**User** type) with the username **user1** and password **password1**. Note that this command includes the full path. If you are already in that directory, you don't need to specify the full path and can instead run the command as: `./user-mgmt add-user user1 User password1`

```
/root/ibm-lh-dev/bin/user-mgmt add-user user1 User password1
```

3. In the watsonx.data user interface running in your browser, click the Access control (key) icon from the left-side menu.

The **Access control** page opens in the **Infrastructure** tab, with the currently defined engines, services, catalogs, and buckets listed (when you add databases, they will be listed here as well).

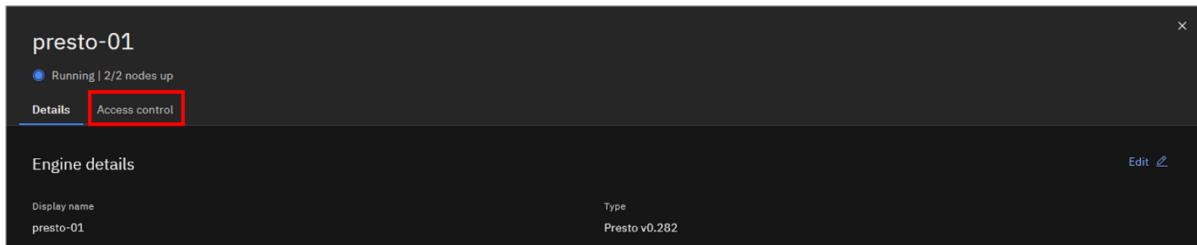
4. Click on **presto-01** (it will turn into a hyperlink when you hover over it).



The screenshot shows the 'Access control' interface with the 'Infrastructure' tab selected. A search bar at the top right contains the placeholder 'Search for infrastructure components'. Below the search bar is a table with two columns: 'Display name' and 'Type'. The table lists several entries:

Display name	Type
presto-01	Engine
milvus000	Service
hive_data	Catalog
iceberg_data	Catalog
wdx_system_data	Catalog
hive-bucket	Bucket
iceberg-bucket	Bucket
wdx-system	Bucket

5. Select the **Access control** tab.



The screenshot shows the details for the 'presto-01' engine. At the top, there is a status indicator showing 'Running | 2/2 nodes up'. Below the status are two tabs: 'Details' (selected) and 'Access control' (highlighted with a red box). Under the 'Details' tab, there is a section titled 'Engine details' containing the following information:

Display name	Type
presto-01	Presto v0.282

6. Click the blue **Add access** button on the right.

The screenshot shows the 'Access control' tab of the Presto-01 interface. It displays a table with one user ('ibmladmin') and no groups. A search bar and pagination controls are at the bottom. The 'Add access' button is highlighted with a red box.

Username	Type	Role for this component	Permission type
ibmladmin	User	Admin	Implicit

7. The **Add users** pop-up window is displayed. For **Name**, select **user1**. For **Role**, select **User**. Click **Add**.

**Note:** If user1 isn't listed then cancel the operation, refresh your browser tab (for instance, using <F5> in the Firefox browser), and repeat the steps above.

The screenshot shows the 'Add access' dialog. The 'Name' field contains 'user1' and the 'Role' dropdown is set to 'User'. The 'Add' button is highlighted with a blue box.

8. Note how the user's role has been added. Click the X in the upper-right corner to close the access controls for the Presto engine.

The screenshot shows the 'Access control' tab again. The user count has increased to 2. The newly added user 'user1' is highlighted with a red box in the list, showing it has been successfully added.

Username	Type	Role for this component	Permission type
user1	User	User	Explicit
ibmladmin	User	Admin	Implicit

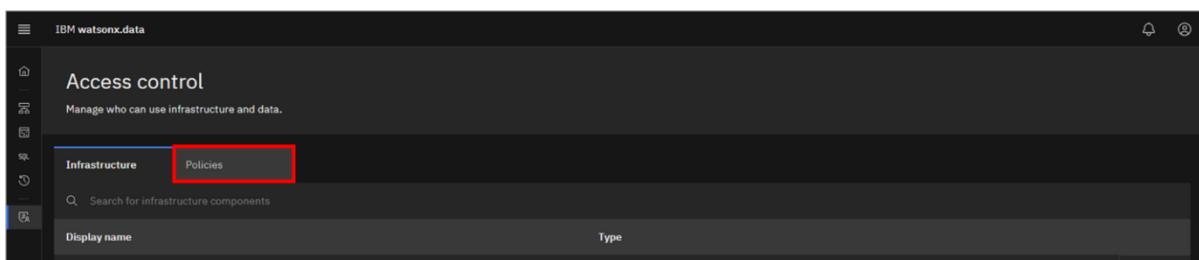
9. Optionally, repeat Steps 4 - 8 to assign the following permissions to **user1** (but rather than clicking on **presto-01**, you'll click on the component referenced below instead). (This is just intended to show the different roles that can be assigned to the other infrastructure component types.)

- **iceberg\_data** (Catalog): User
- **iceberg-bucket** (Bucket): Reader

In a real-world scenario where a user will be querying data from a table, that user will need to be given a minimum of **User access to an engine** (to be able to run the query), **User access for the catalog** associated with the data (to be able to see the schema information associated with the table), and **Reader access to the bucket** associated with the data (to be able to read the data from object storage).

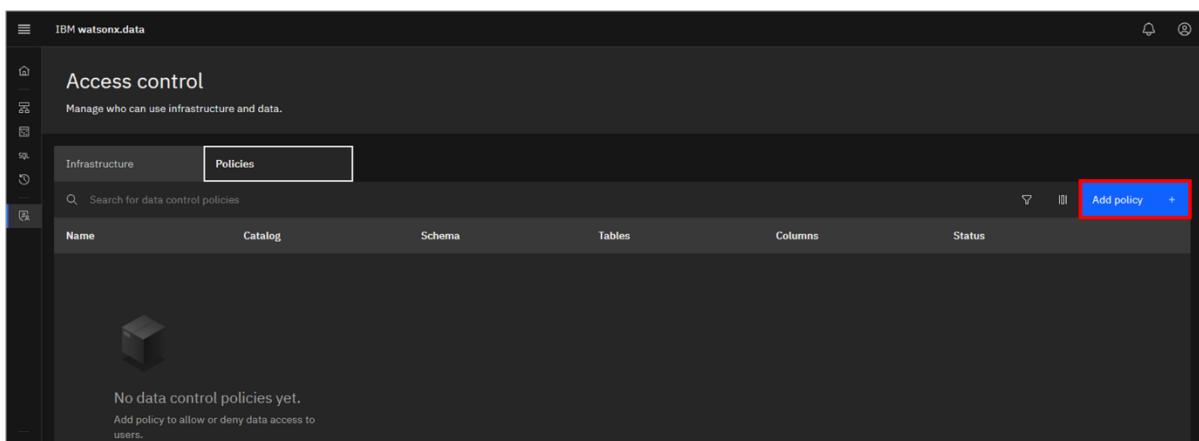
Additionally, a policy has to be created to permit the user to access the table in question.

10. Select the **Policies** tab.



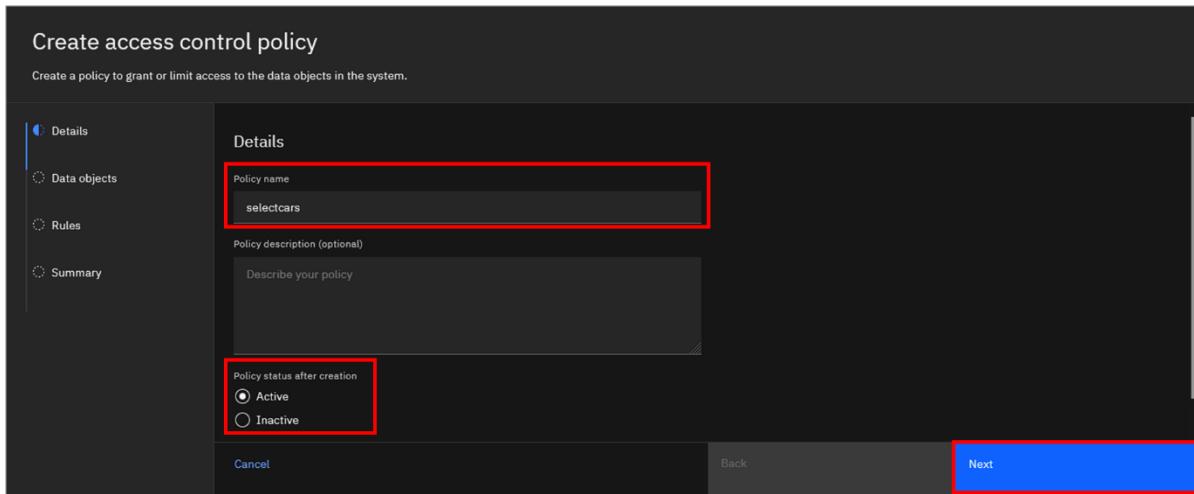
The screenshot shows the 'Access control' interface for the 'IBM watsonx.data' service. The top navigation bar includes icons for Home, Infrastructure, Policies, and Data. The 'Policies' tab is highlighted with a red box. Below the tabs is a search bar labeled 'Search for infrastructure components'. A table follows, with columns for 'Display name' and 'Type'. The table is currently empty.

11. Click the blue **Add policy** button on the right.



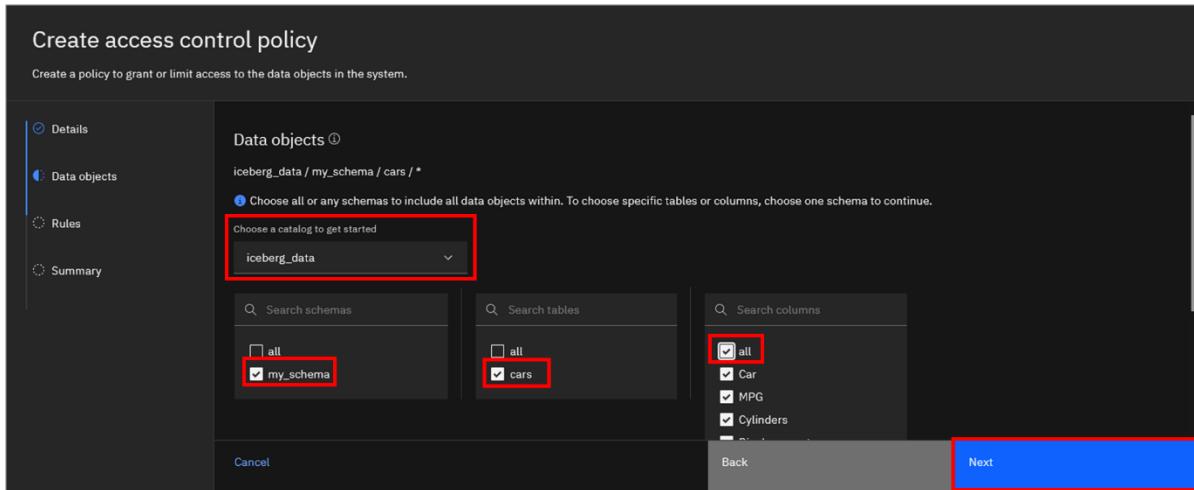
The screenshot shows the same 'Access control' interface as the previous one, but now the 'Policies' tab is selected. The 'Add policy' button, located at the top right of the main content area, is highlighted with a red box. The main content area displays a message: 'No data control policies yet.' and 'Add policy to allow or deny data access to users.'

12. Enter **selectcars** in the **Policy name** field, select **Active** for **Policy status after creation**, and then click **Next**.



13. Click on the **Choose a catalog to get started** dropdown and select **iceberg\_data**. A list of schemas is then displayed.

- Select the checkbox for the **myschema** schema.
- A list of the tables in the schema is then displayed; select the checkbox for the **cars** table.
- A list of the columns in the table is then displayed; select the checkbox for **all** columns.
- Click **Next**.



14. One or more rules can now be added to the policy. Click **Add rule** on the right.

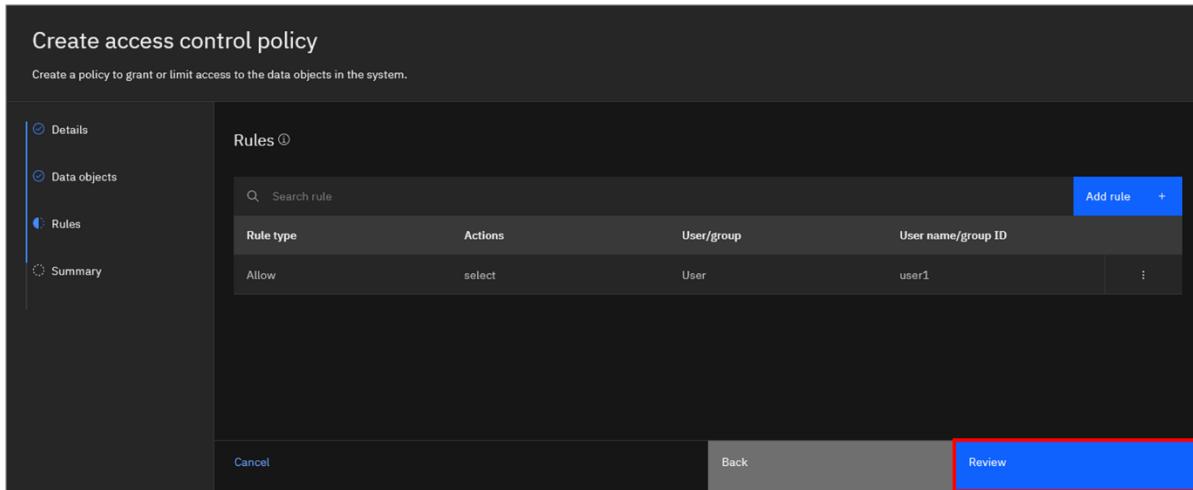
The screenshot shows the 'Create access control policy' interface. On the left, there's a sidebar with tabs: 'Details' (selected), 'Data objects', 'Rules' (which is the active tab), and 'Summary'. The main area is titled 'Rules' with a subtitle 'No rules yet.' Below this is a table header with columns: 'Rule type', 'Actions', 'User/group', and 'User name/group ID'. In the bottom right corner of the main area, there is a blue button labeled 'Add rule' with a '+' sign, which is highlighted with a red box.

15. Make the following selections and then click the **Add** button:

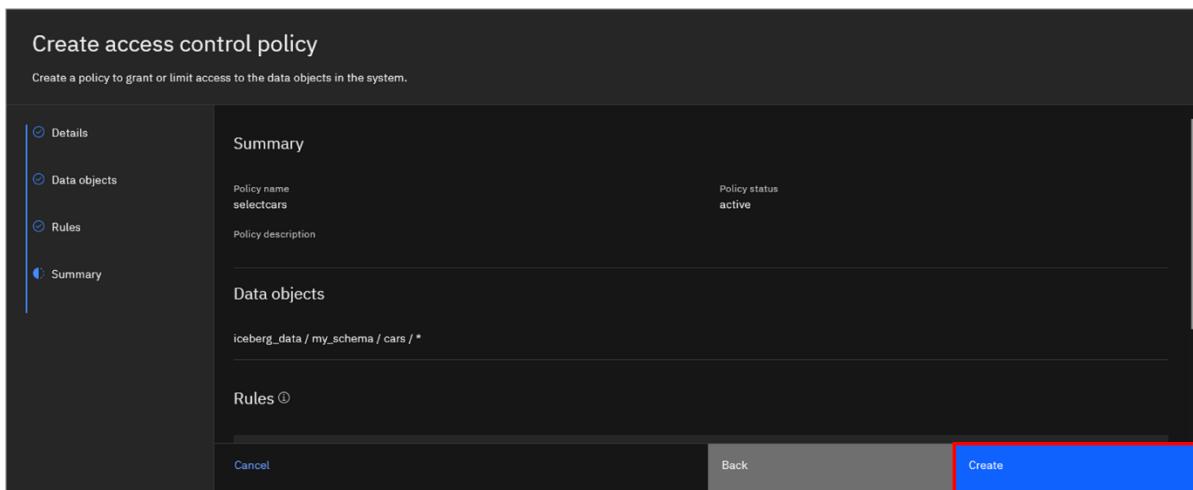
- a. For the **Rule type**, select **Allow**.
- b. For **Actions**, select the **select** checkbox.
- c. For **Users/group**, select **Users**.
- d. In the **Select user** dropdown, select **user1**.

The screenshot shows the 'Add rule' dialog box. Under 'Details', the 'Rule type' section has 'Allow' selected (radio button highlighted). In the 'Actions' section, the 'select' checkbox is checked (highlighted with a red box). Under 'Users/groups', the 'Users' radio button is selected (highlighted with a red box). In the 'Select user' dropdown, 'user1' is listed. At the bottom right of the dialog, there is a blue 'Add' button which is highlighted with a red box.

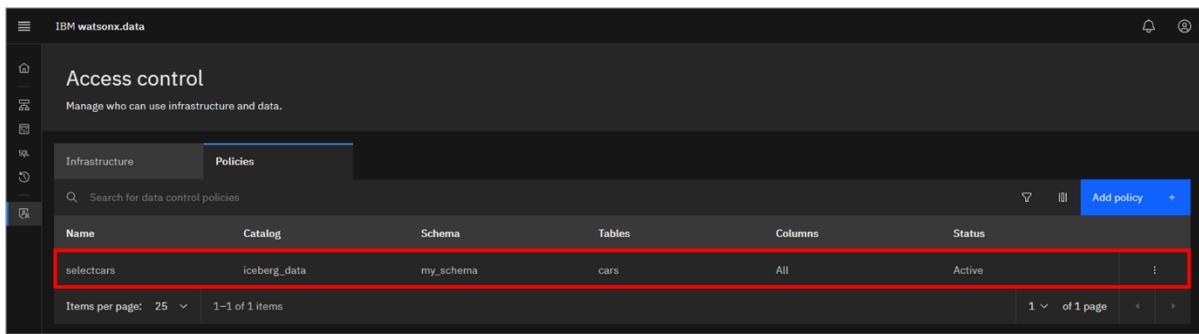
16. With the rule added to the policy, click **Review**.



17. Review the policy details you've entered and click **Create**.



The policy should now be listed. Because you specified for the policy to be active upon creation, it is in fact active (as you can see by looking at the **Status** column).



The screenshot shows the 'Access control' section of the IBM Watsonx.data interface. The 'Policies' tab is selected. A table lists a single policy:

Name	Catalog	Schema	Tables	Columns	Status
selectcars	iceberg_data	my_schema	cars	All	Active

Items per page: 25 | 1–1 of 1 items

You have just given **user1** the permissions needed to query the data within the **cars** table.

## 6. Working with Presto

Up to this point, you have seen how Presto is integrated into watsonx.data and the watsonx.data user interface. However, you can also work directly with Presto.

For example, Presto includes a terminal-based interactive command line interface (CLI) that allows you to run SQL statements. Presto also has a web-based graphical user interface for getting information on query activity in the system.

Additionally, users can connect 3<sup>rd</sup> party management and analytics tools to the Presto server. Likewise, applications can use JDBC (Java Database Connectivity) to work with data in watsonx.data, by connecting to the Presto server (this topic is beyond the scope of this lab).

### 6.1. Presto Command Line Interface (CLI)

Presto CLI is a terminal-based interactive shell that can be used to run queries. You can connect to the Presto server either through Presto CLI installed as part of the *ibm-lh-client* package or through Presto CLI installed separately.

Presto CLI is pre-installed as part of watsonx.data Developer Edition. It is started by using the `presto-cli` command (located in `/root/ibm-lh-dev/bin`).

Presto uses a 3-part name to identify tables: catalog.schema.table. These identifiers can be enclosed in double quotes as needed (for example, "mycatalog""myschema""mytable"). Note that double quotes are needed if you're using any special characters in a name, like a hyphen.

1. Open a terminal command window to the watsonx.data server as the **root** user.
2. Run the following two commands to open a Presto CLI interactive terminal.

```
cd /root/ibm-lh-dev/bin
```

```
./presto-cli
```

**Note:** In the Presto CLI interactive terminal, you can use the cursor control (arrow) keys on your keyboard to scroll through previously entered statements and edit the current statement you're about to run.

- Run the following command to list the catalogs that have been registered in Presto (Note: Presto requires that commands end with a semicolon.)

```
show catalogs;
```

The output should be similar to the text below.

```
Catalog
-----
hive_data
iceberg_data
jmx
system
tpcds
tpch
wxd_system_data
(7 rows)
```

- Run the following query. It should return a count of 1,500 rows.

```
select count(*) from tpch.tiny.customer;
```

Rather than explicitly including the catalog name and schema name for every statement, you can set the default catalog and schema for a session when you start the CLI. You will try this shortly.

You can also tell Presto which schema you want to work with by using the USE command.

- First, run the following command to see the schemas in the **tpch** catalog.

```
show schemas in tpch;
```

The output should be similar to the text below.

```
Schema
-----
information_schema
sf1
sf100
sf1000
sf10000
sf100000
sf300
sf3000
sf30000
tiny
(10 rows)
```

6. Run the following USE command to set the schema for the current session:

```
use tpch.tiny;
```

7. Run the following query. Note how the catalog and schema aren't required. Again, this should return a count of 1,500 rows.

```
select count(*) from customer;
```

Alternatively, you can provide the session catalog and schema (or just the catalog) when you start the Presto CLI. Note that even if you specify the catalog and schema when starting the Presto CLI or by running the **USE** command, you can still access tables in other schemas and catalogs. You just have to specify the full 3-part name for them (or USE the schema in question to change the default for the session).

8. Quit from the Presto CLI by running the following command.

```
quit;
```

9. Start the Presto CLI again, but this time specify the session catalog and schema as command line options.

```
./presto-cli --catalog tpch --schema tiny
```

10. Run the following query. Note how this time the catalog and schema aren't required. Again, this should return a count of 1,500 rows.

```
select count(*) from customer;
```

11. Quit from the Presto CLI by running the following command.

```
quit;
```

Next, you will create a new schema and table. When using Presto to create a schema in watsonx.data, you must specify the object storage bucket associated with the catalog.

12. Start the Presto CLI and specify that you intend to use the **iceberg\_data** catalog.

```
./presto-cli --catalog iceberg_data
```

13. Run the following SQL statement to create a new schema in the catalog.

```
create schema if not exists newschema  
with (location='s3a://iceberg-bucket/newschema');
```

14. Run the following command to see a list of the schemas in the catalog being used.

```
show schemas;
```

The **newschema** schema should be listed, along with the **my\_schema** schema created earlier.

15. Run the following SQL statements to create a new table in this schema, populate the table with some data, and then query the table.

```
create table newschema.users (id int, name varchar, age int);
```

```
insert into newschema.users values (1, 'Robert', 54);
```

```
insert into newschema.users values (2, 'Susan', 37);
```

```
select * from newschema.users order by id;
```

The output of the final SELECT statement should look similar to what is shown below.

id	name	age
1	Robert	54
2	Susan	37

(2 rows)

16. Run the following two commands to show the tables that exist in the new schema.

```
use newschema;
```

```
show tables;
```

The output of the SHOW TABLES statement should look similar to what is shown below.

Table
users

(1 row)

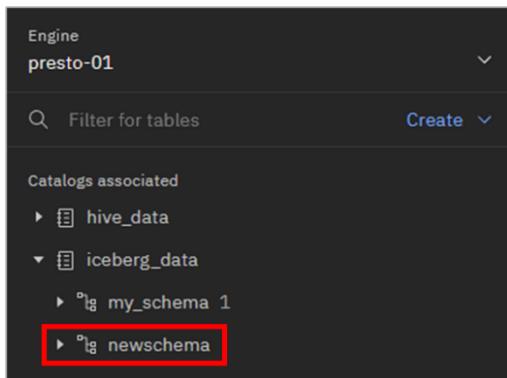
17. Quit from the Presto CLI by running the following command.

```
quit;
```

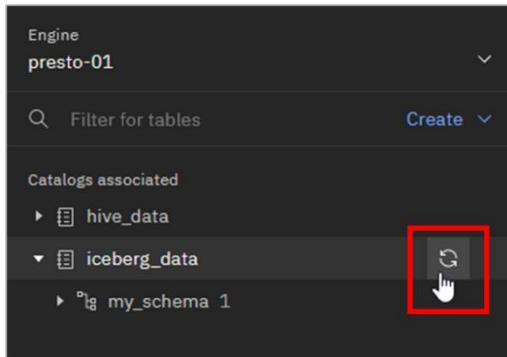
18. Open the watsonx.data user interface in a browser window (if you don't already have one open).

19. Select the **Data manager** (grid) icon from the left-side menu.

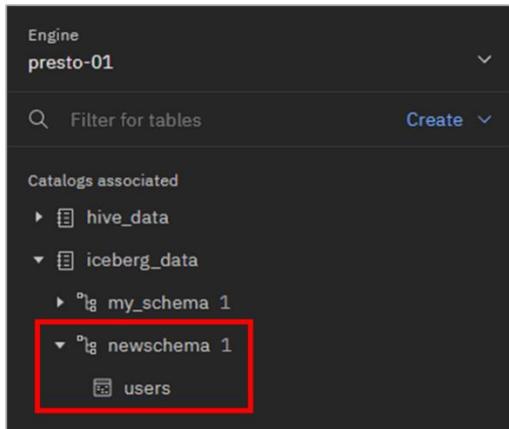
20. Expand the **iceberg\_data** catalog to see the new schema you created (**newschema**).



**Note:** If you don't see the schema listed then hover your mouse pointer over the far right of the line for the **iceberg\_data** catalog until you see the **Refresh** icon appear. Click the **Refresh** icon. You should now see the schema listed as in the above image.



21. Expand the **newschema** schema to see the new table you created (**users**).



The screenshot shows the Presto Web Interface's catalog list. At the top, it says "Engine" and "presto-01". Below that is a search bar labeled "Filter for tables" and a "Create" button. Under "Catalogs associated", there are two entries: "hive\_data" and "iceberg\_data". "iceberg\_data" is expanded, showing "my\_schema" and "newschema". "newschema" is also expanded, showing the "users" table. The "newschema" entry and its "users" table are highlighted with a red box.

The table you created through the Presto CLI is visible through the watsonx.data user interface. This is one of the benefits of having a shared metastore. In the future, any query engines that get associated with the `iceberg_data` catalog would also be able to work with this table.

## 6.2. Presto Web Interface

Presto includes its own web interface (console) for monitoring and managing Presto queries. It's a great place to get information about running queries and completed queries. This includes the query text, query state, the name of the user that ran the query, and the percentage complete if it's still running.

Each query is assigned a unique Query ID and clicking on the ID brings up a Query Details page with additional information regarding the query.

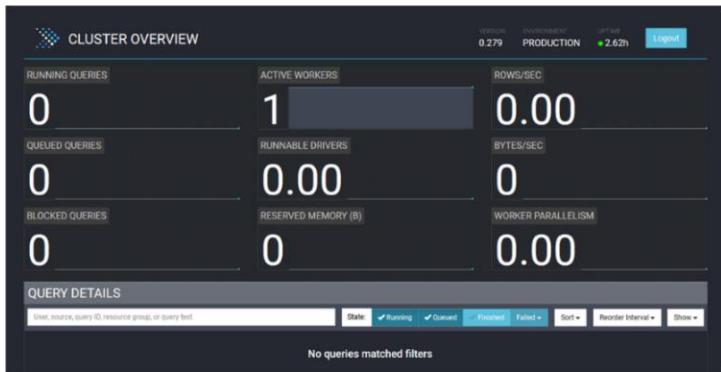
1. From your computer, open the **Presto console** in your browser. The URL can be found in your TechZone reservation details (see the **Presto console** line in the **Published services** section, per the instructions in section [4.2 Accessing the watsonx.data Environment](#)).

You may be warned about a potential security risk, in which case you can proceed as you did with the watsonx.data console earlier.

When prompted to sign in, enter the following credentials:

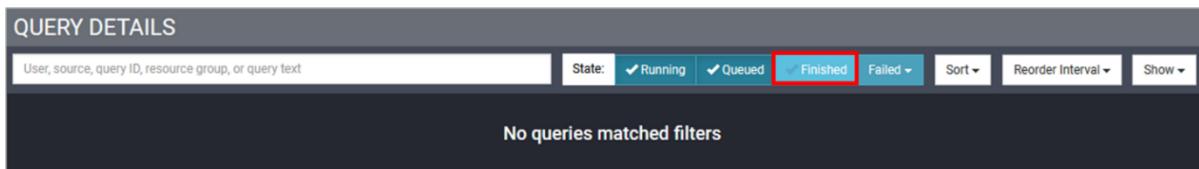
**Username:** ibmlhadmin  
**Password:** password

You immediately start on the **Cluster Overview** page. The upper portion of the page includes various metrics regarding the environment.

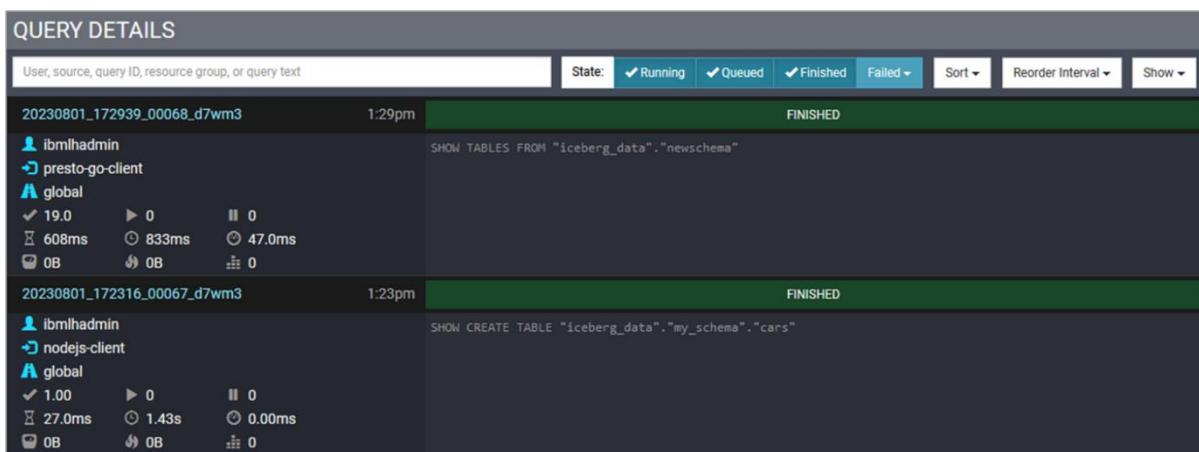


The **Query Details** section at the bottom of the page lists queries matching the **State** filter being applied. By default, only actively running queries are shown and so you probably don't see any queries listed here.

2. Select the **Finished** button to include queries that have finished running.



You should now see queries listed. You may recognize some as queries you ran, and others may have been run internally by the system.



- Click the **Query ID** link for a query that looks interesting to you.

```

20230801_172316_00065_d7wm3
1:23pm FINISHED

ibmihadmin
nodejs-client
global
✓ 17.0 ▶ 0 || 0
1.76s ○ 1.88s ○ 137ms
0B ○ 0B ─ 0

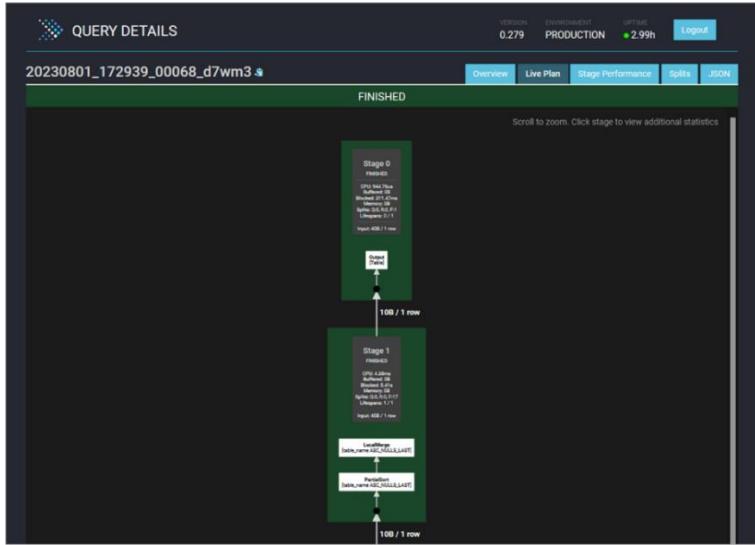
SELECT * FROM information_schema.columns WHERE
    table_catalog = 'iceberg_data' AND
    table_schema = 'my_schema' AND
    table_name = 'cars'
  
```

- This opens a new browser window with the **Query Details** page for that specific query. There are multiple tabs, and you start in the **Overview** tab by default.

Session		Execution	
User	ibmihadmin	Resource Group	global
Principal	IBMLHPrincipal [name=ibmihadmin, role=Administrator, groups=null]	Submission Time	2023-08-01 1:29pm
Source	presto-go-client	Completion Time	2023-08-01 1:29pm
Catalog	iceberg_data	Elapsed Time	833.11ms
Schema	newschema	Prerequisites Wait Time	7.82ms
Client Address	172.18.0.4	Queued Time	335.07us
Client Tags		Planning Time	242.56ms
Session Properties		Execution Time	607.79ms
Resource Estimates		Coordinator	172.18.0.6

- Scroll down the page to familiarize yourself with the information available.
- Select the **Live Plan** tab at the top of the page.

Here you can see the query's execution plan and the various steps involved in running the query (your execution plan will look different as it depends on the query you chose). This output is similar to what you see in the visual explain output for watsonx.data.



7. Click the **Presto logo** in the top-left to return to the Cluster Overview page.



8. When you're done exploring the Presto console, close the browser window (and any other Presto console windows that are still open).

## 7. Working with MinIO

### 7.1. Introduction to Object Storage

Object storage is a type of storage architecture where data is managed as *objects*. Each object typically includes the data itself, a variable amount of metadata, and a globally unique identifier. The basic unit of storage is an object and objects are organized in *buckets* (think of a bucket like a high-level directory or folder).

Object storage systems allow retention of massive amounts of structured, semi-structured, and unstructured data in which data is written once and read once or many times. Object storage is low cost, with near unlimited capacity and scalability.

Cloud object storage vendors have designed their offerings with extremely high levels of durability and reliability. The most notable provider of object storage is Amazon S3 (Simple Storage Service). Most other vendors (including IBM Cloud) offer S3 API-compatible object storage. To be clear, S3 API-compatible doesn't mean that the object storage is running in Amazon. When the industry talks of object stores being S3 API-compatible (or simply S3-compatible), it means that they support the same API set that Amazon created around their S3 object storage. In doing so, applications and tools that support S3 should work against S3-compatible object storage as well.

Clients consider characteristics such as performance, price, and workload type when deciding upon a storage format for a particular type of data or workload. One of the main reasons for the rising popularity of the lakehouse architecture is its use of low-cost object storage. Data warehouses, on the other hand, have traditionally used block storage to store data, which offers high performance, but is significantly more costly.

A major appeal of watsonx.data is that vast amounts of data can be stored in object storage at a relatively low cost, with fit-for-purpose query engines (like Presto and Spark) accessing the data concurrently. This allows for offloading existing data from a client's enterprise data warehouse (EDW), where the performance requirements and/or frequency with which the data is accessed don't justify the costs of having that data in the warehouse (keep in mind that costs aren't limited to the data storage itself; there are costs in preparing and moving data into the warehouse, additional storage costs for larger backup images, the impact of running relatively low priority workloads at the same time as higher priority workloads, and so on). Additionally, with IBM Db2 and Netezza's tight integration with watsonx.data, data that needs to live in a data warehouse can do so, while the rest of the data lives in watsonx.data. Db2 and Netezza can access the data in object storage directly (due in part to the shared metadata store) and queries can combine data in the warehouse with the data in the lakehouse. This provides clients with complete flexibility in where they store their data.

## 7.2. Exploring MinIO Object Storage

Watsonx.data Developer Edition includes a local S3-compatible object store called *MinIO*. Rather than using an external S3 object store (which is certainly possible with watsonx.data), this lab uses the local MinIO object store.

You will need MinIO credentials to log into and use the MinIO console. This includes an *access key* (username) and a *secret key* (password). These values are specific to your environment and are stored as environment variables in the Presto container.

1. Open a terminal command window to the watsonx.data server as the **root** user.
2. Run the following command to extract and display the access key (username). Make note of this and keep it in a location you can refer to later.

**Note:** This is a single command but is split over two lines in the published text below. It should copy and paste as a single line, but if it pastes as two separate lines then it will be interpreted as two separate commands and will fail. In this case, copy and paste such that this is run as a single command. Alternatively, remember that all of the commands and SQL statements can be found in this [text document](#). Download this file and copy the command from that file instead. This also applies to the command in step 3.

```
docker exec ibm-lh-presto printenv | grep LH_S3_ACCESS_KEY | sed  
's/.*=//'
```

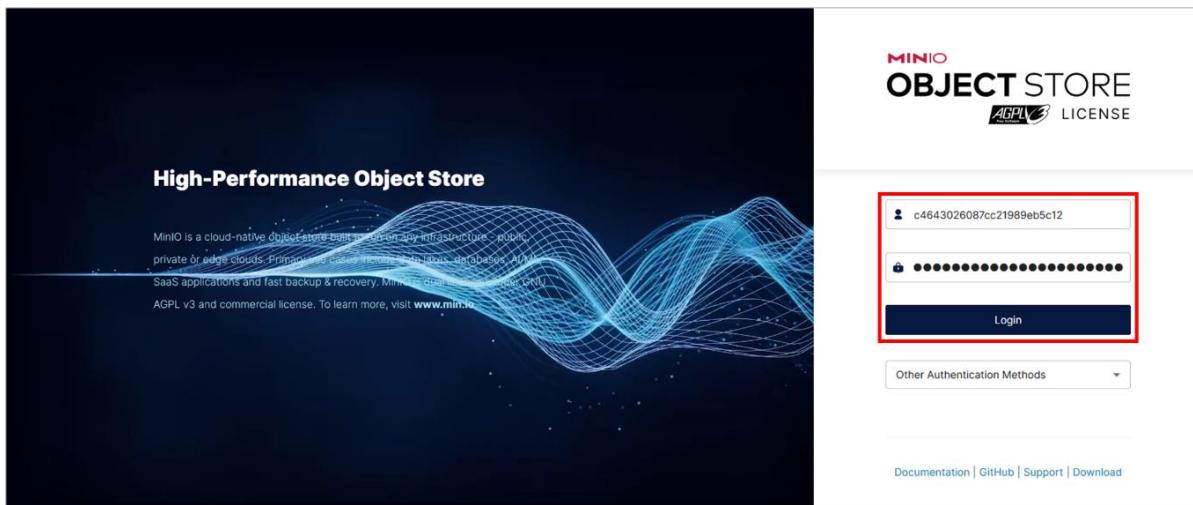
3. Run the following command to extract and display the secret key (password). Make note of this and keep it in a location you can refer to later.

```
docker exec ibm-lh-presto printenv | grep LH_S3_SECRET_KEY | sed  
's/.*=//'
```

4. Open the **MinIO console** in a new browser window. The URL can be found in your TechZone reservation details (see the **Minio console** line in the **Published services** section, per the instructions in section [4.2 Accessing the watsonx.data Environment](#)).

You may be warned about a potential security risk, in which case you can proceed as you did with the watsonx.data console earlier.

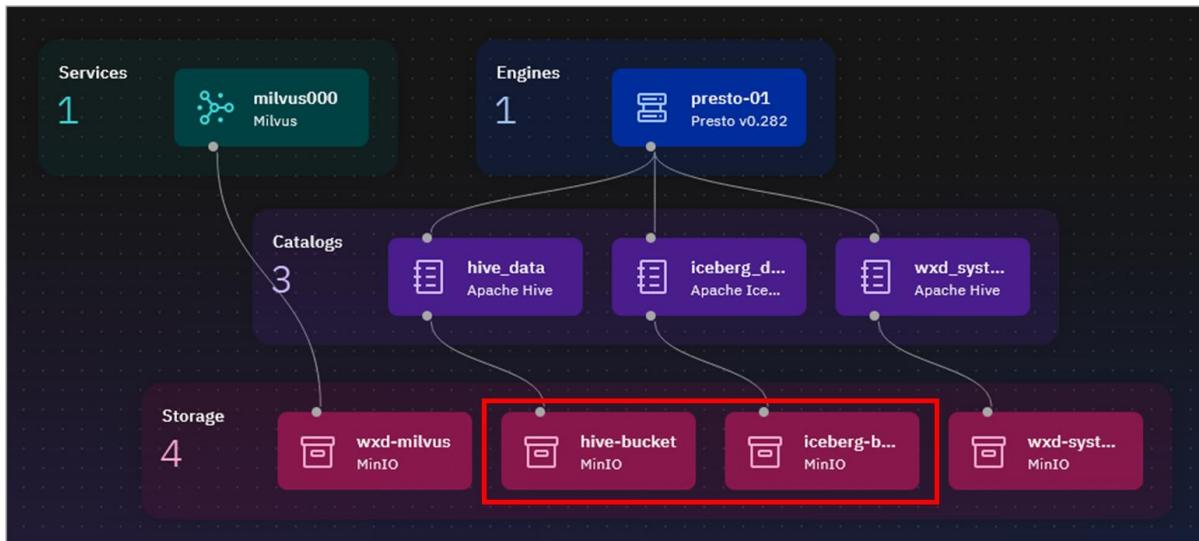
5. Refer to the MinIO credentials you copied above and enter the **access key** for the **Username** and the **secret key** for the **Password**. Click the **Login** button.



This opens the **Object Browser** screen. Various buckets already exist, but the two you will be using for table data are **hive-bucket** and **iceberg-bucket**. There are a number of objects within these two buckets already – the pre-populated tables in the hive-bucket and the tables you've created in the iceberg-bucket.

Name	Objects	Size	Access
hive-bucket	73	42.4 MiB	R/W
iceberg-bucket	16	51.1 KiB	R/W
wxd-milvus	0	0.0 B	R/W
wxd-system	2	0.0 B	R/W

- Refer back to the **Infrastructure manager** screen in the **watsonx.data** console and note how these two buckets were previously registered with **watsonx.data**.



All tables created in the **hive\_data** catalog have their files (objects) in the **hive-bucket** bucket. Likewise, all tables created in the **iceberg\_data** catalog have their files (objects) in the **iceberg-bucket** bucket.

- Go back to the MinIO console. Click the row for the **iceberg-bucket** bucket.

Name	Objects	Size	Access
hive-bucket	75	42.4 MiB	R/W
iceberg-bucket	16	49.4 KIB	R/W

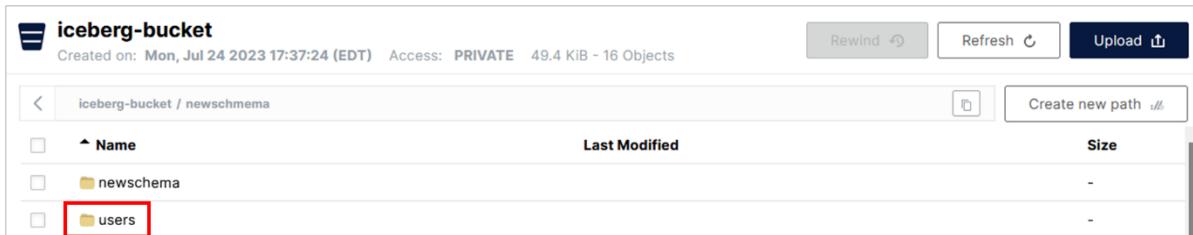
There are two folders listed that correspond to the two schemas you created earlier in this lab (both of which were created in the **iceberg\_data** catalog, which is associated with this **iceberg-bucket** bucket).

By default, **watsonx.data** uses the schema name for the folder when you create it (such as when you created the schema **my\_schema** earlier in the lab). Also recall that when you created the schema **newschema** using the Presto CLI, you specified the location as '`s3a://iceberg-bucket/newschema`', which matches what you see here.

The screenshot shows the MinIO interface with the 'iceberg-bucket' selected. It displays two objects: 'my\_schema' and 'newschema', which are highlighted by a red box.

8. Click the row for the **newschema** folder.

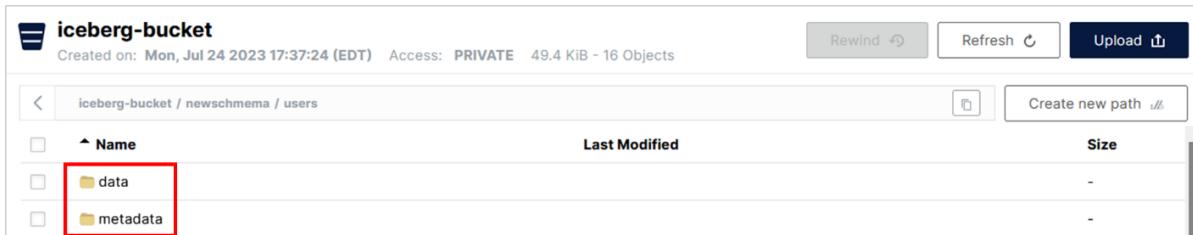
In the **newschema** folder is a sub-folder called **users**. When you created the **users** table in the Presto CLI earlier, it created this sub-folder with the same name as the table.



The screenshot shows the MinIO console interface. The top navigation bar displays the bucket name "iceberg-bucket", creation date "Mon, Jul 24 2023 17:37:24 (EDT)", access level "PRIVATE", and object count "49.4 KIB ~ 16 Objects". Below the navigation bar are buttons for "Rewind", "Refresh", and "Upload". The main area shows a file tree under the path "iceberg-bucket / newschmema". The "users" folder is highlighted with a red box. The table has columns: "Name", "Last Modified", and "Size". There are two entries: "newschema" and "users", both with size "-".

9. Click the row for the **users** folder.

There are two sub-folders within the table's folder. The **data** folder contains the Parquet file(s) holding the actual data for the table. The **metadata** folder contains a series of metadata files (of varying data file formats) used by Iceberg.



The screenshot shows the MinIO console interface. The top navigation bar displays the bucket name "iceberg-bucket", creation date "Mon, Jul 24 2023 17:37:24 (EDT)", access level "PRIVATE", and object count "49.4 KIB ~ 16 Objects". Below the navigation bar are buttons for "Rewind", "Refresh", and "Upload". The main area shows a file tree under the path "iceberg-bucket / newschmema / users". The "data" and "metadata" folders are highlighted with red boxes. The table has columns: "Name", "Last Modified", and "Size". There are two entries: "data" and "metadata", both with size "-".

10. Review the files within the **data** and **metadata** folders. You will see a mix of Parquet and Avro files.

To return back up the folder hierarchy, use the **breadcrumbs** at the top of the navigation pane. Clicking the < icon brings you to the parent of the folder you are currently in.



The screenshot shows the MinIO console interface. The top navigation bar displays the bucket name "iceberg-bucket", creation date "Mon, Jul 24 2023 17:37:24 (EDT)", access level "PRIVATE", and object count "49.4 KIB ~ 16 Objects". Below the navigation bar are buttons for "Rewind", "Refresh", and "Upload". The breadcrumb path "iceberg-bucket / newschmema / users / metadata" is highlighted with a red box. The main area shows a file tree under the path "iceberg-bucket / newschmema / users / metadata". The table has columns: "Name", "Last Modified", and "Size". There are no visible entries.

11. If you are proceeding to the next section (Data Ingest) then keep the MinIO console open. Otherwise, close the browser window.

## 8. Data Ingest

Data ingestion is the process of importing and loading data into watsonx.data. Multiple methods are available to ingest data into watsonx.data, including:

- The **Create table from file** option on the **Data manager** page (as you saw earlier in section [5.3. Data Manager Page](#)).
- The standalone, non-Developer Edition of watsonx.data includes an **Ingestion jobs** tab on the **Data manager** page, where data can be read from a source object storage bucket and ingested into a table. However, as this lab uses the Developer Edition, it is not available to use here.
- Installation of the **ibm-lh** command line tool and using commands to create a job to ingest files from S3-compatible object storage or a local file system.
- Loading the data file(s) for a table into an object storage bucket (if they don't already reside in object storage), registering the bucket with watsonx.data, and creating a table on top of the file(s) using SQL.

In this section you will use the MinIO console to upload a Parquet data file into object storage and then you will create a table on it using SQL in Presto.

1. Download the **aircraft.parquet** file from here: <https://ibm.box.com/v/ontime-aircraft-id>.
2. If you don't already have the MinIO console open, open it in a new browser window. Use the credentials (access key and secret key) you made note of earlier.

(The console URL can be found in your TechZone reservation details (see the **Minio console** line in the **Published services** section, per the instructions in section [4.2 Accessing the watsonx.data Environment](#)).

Because you are uploading a data file that is not “wrapped” in the Iceberg table format, it needs to use a Hive catalog. You will use the existing **hive\_data** catalog which is associated with the **hive-bucket** bucket.

3. Click the row for the **hive-bucket** bucket.

Name	Objects	Size	Access
hive-bucket	75	42.4 MiB	R/W
iceberg-bucket	16	49.4 KiB	R/W

- In the **hive-bucket** panel, click the **Create new path** button.

The screenshot shows the 'hive-bucket' panel. At the top, there are buttons for 'Rewind', 'Refresh', and 'Upload'. Below that is a table with columns: 'Name', 'Last Modified', and 'Size'. The table lists several folders: 'gosalesdw', 'hive\_sql', 'ontime', and 'taxi'. In the top right corner of the table area, there is a red box around the 'Create new path' button.

- In the **Choose or create a new path** pop-up window, enter **myschema2/aircraft** for the **New Folder Path**. Then, click the **Create** button.

The screenshot shows a modal dialog titled 'Choose or create a new path'. It has a 'Current Path:' label with the value '/hive-bucket'. Below it is a 'New Folder Path\*' input field containing 'myschema2/aircraft', which is also highlighted with a red box. At the bottom are two buttons: 'Clear' and 'Create', with 'Create' also highlighted with a red box.

This is going to create nested folders under the current bucket. The top-level folder (**myschema2**) should match the name of the schema you are going to create later, and the next folder down (**aircraft**) should match the name of the table you are going to create later.

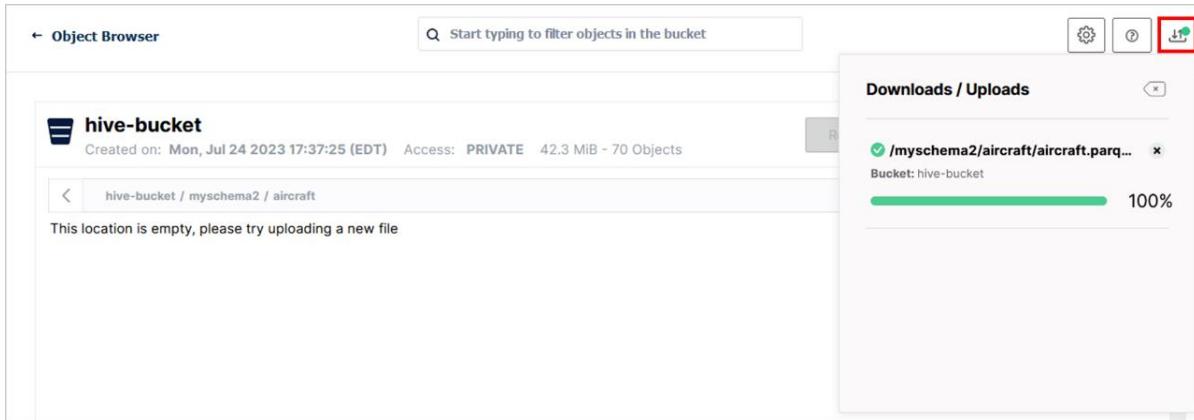
- You are instructed to upload a file. Click the **Upload** button in the upper-right and then click **Upload File**.

The screenshot shows the 'hive-bucket' panel again. The top right has an 'Upload' button, which is highlighted with a red box. Below it, there are two options: 'Upload File' (with an up arrow icon) and 'Upload Folder' (with a folder icon). Both are enclosed in a red box.

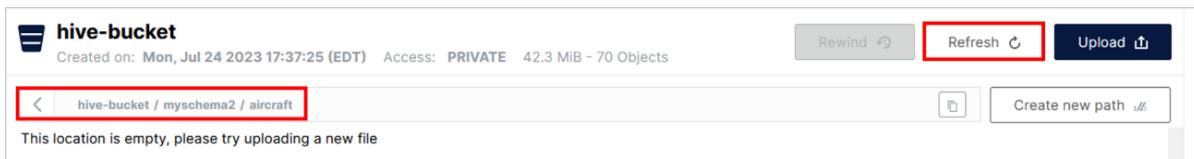
**Note:** You are also given the option of uploading a folder. If you already have a set of files in folders organized by schema and table, you could upload the entire folder (in which case you wouldn't have to create a new path for the table as you did above, as the path is already reflected in the folder being uploaded).

- Locate, select, and upload the **aircraft.parquet** file you downloaded earlier.

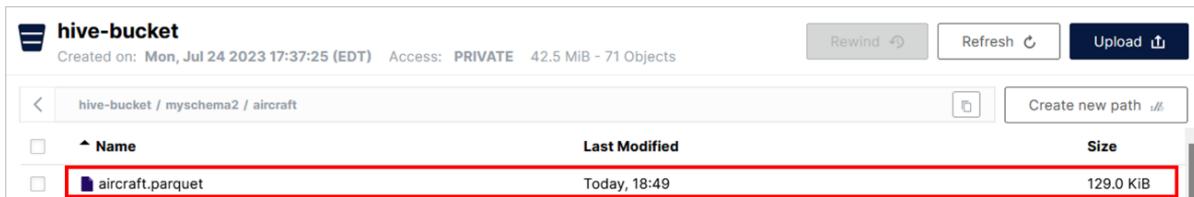
8. Click the **Downloads/Upserts** icon to close the **Downloads/Upserts** panel.



9. If the folder appears empty, click the **Refresh** button in the upper-right. If the folder is still showing as empty, navigate up one level to the parent folder (**myschema2**) and then navigate back down to this folder (**aircraft**). If this still doesn't work, refresh the browser window (for instance, using **<F5>** in the Firefox browser).



The file should now be listed, as in the image below.



10. Close the MinIO console browser window as you no longer need to use it.

You will now create a schema and a table to match the file that you uploaded. You will use the **Presto CLI**, but alternatively you could also use the **Query workspace** page in the watsonx.data web interface.

11. Open a terminal command window to the watsonx.data server as the **root** user.

12. Run the following command to open a Presto CLI interactive terminal.

```
/root/ibm-lh-dev/bin/presto-cli
```

13. Run the following SQL statement to create the **myschema2** schema (based on the catalog/bucket being used and the schema folder name you specified when uploading the file to object storage earlier; the schema name doesn't have to match, but it's easier to manage this way).

```
create schema hive_data.myschema2
with (location = 's3a://hive-bucket/myschema2');
```

14. Run the following SQL statement to create the **aircraft** table.

```
create table hive_data.myschema2.aircraft
(tail_number varchar,
 manufacturer varchar,
 model varchar)
with (format = 'Parquet',
      external_location='s3a://hive-bucket/myschema2/aircraft');
```

15. Validate that the **aircraft** table has been created correctly by querying the number of rows in it.

```
select count(*) from hive_data.myschema2.aircraft;
```

There should be 13,101 rows returned. Your table just read the data out of the file you uploaded into object storage!

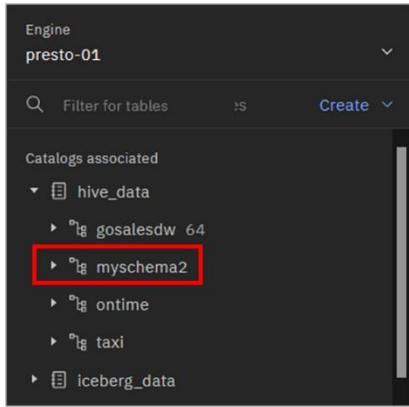
16. Quit from the Presto CLI by running the following command.

```
quit;
```

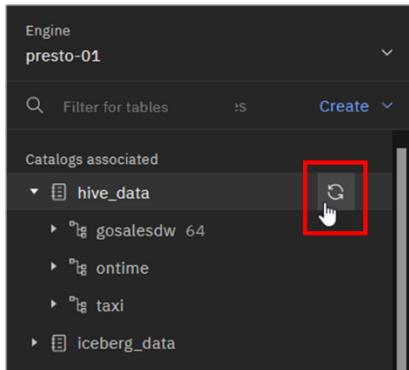
17. Open the watsonx.data user interface in a browser window (if you don't already have it open).

18. Select the **Data manager** (grid icon) from the left-side menu.

19. Expand the **hive\_data** catalog to see the new schema you created (**myschema2**).



**Note:** If you don't see the schema listed then hover your mouse pointer over the far right of the line for the **hive\_data** catalog until you see the **Refresh** icon appear. Click the **Refresh** icon. You should now see the schema listed as in the above image.



20. Expand the **myschema2** schema to see the new table you created (**aircraft**).

The screenshot shows the WatsonX Data user interface for the Presto CLI engine named 'presto-01'. In the 'Catalogs associated' section, the 'hive\_data' catalog is expanded, revealing its contents: 'gosalesdw 64', 'myschema2 1', and 'aircraft'. The 'myschema2 1' entry is highlighted with a red box, and the 'aircraft' table is also visible.

The table you just created through the Presto CLI – which is based on a data file you uploaded into object storage – is visible through the watsonx.data user interface. This is one of the benefits of having a shared metastore. In the future, any query engines that get associated with the `hive_data` catalog would also be able to work with this table.

In practice, a next step might be to convert this table to the Iceberg table format, so that it benefits from all the capabilities provided by Iceberg (such as transactional consistency, schema and partition evolution, snapshots for time travel queries and rollback, and improved query efficiency). You will see how easy it is to create a copy of a table using CREATE TABLE AS SELECT (CTAS) in section [10. Offloading Data from a Data Warehouse](#).

## 9. Federated Queries

Unlike traditional database systems, Presto doesn't have its own native database storage. Instead, Presto supports separation of compute and storage, with dozens of connectors that let Presto access data where it lives – which could be in relational databases, NoSQL databases, data warehouses, data lakes, data lakehouses, and more.

With Presto's *federated query* capability (frequently referred to as *data federation*, or simply *federation*), enterprises can combine data from their existing, diverse sources with new data that they have in watsonx.data, rapidly unlocking insights without the complexity and cost of duplicating or moving data.

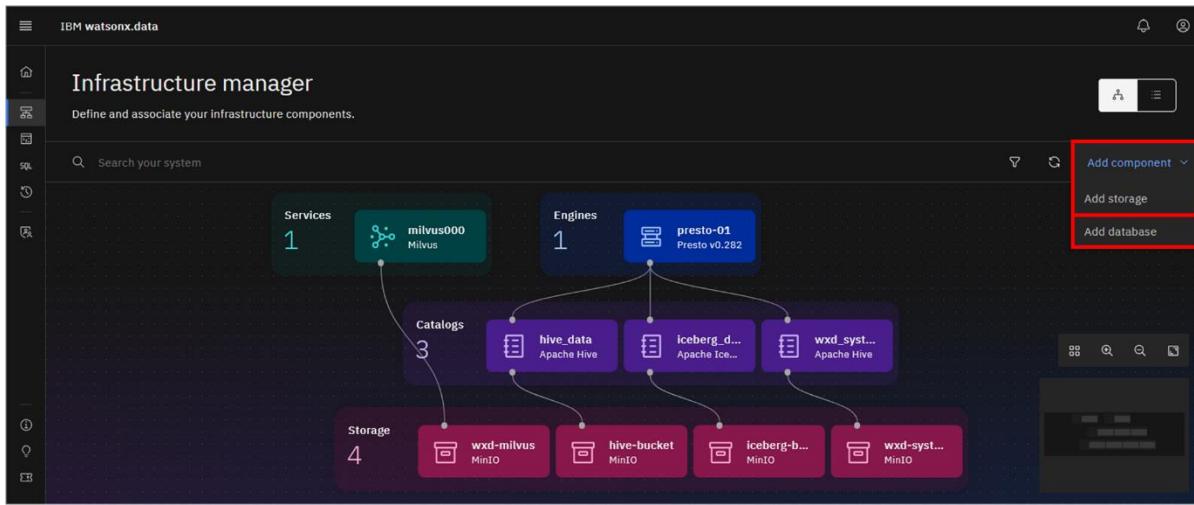
Although Presto supports a wide range of connectors, watsonx.data officially only supports a subset of them. This is because IBM wants to ensure quality, performance, and security of connectors before adding support (which may require updating connector source code to do this). More connectors will be added over time.

As of 1Q 2024, Presto in watsonx.data can currently connect to IBM Db2, Netezza, Apache Kafka, Elasticsearch, MongoDB, MySQL, PostgreSQL, SAP HANA, SingleStore, Snowflake, SQL Server, Teradata, and others through a custom connector. The most current list of connectors and the types SQL statements supported can be found [here](#). The list of supported connectors will grow over time.

In this section you will combine data from watsonx.data's object storage with data in Db2 and PostgreSQL databases. To avoid you having to provision these databases yourself, they've been installed in the same VM as watsonx.data and are pre-populated with data.

1. Select the **Infrastructure manager** icon () from the left-side menu of the **watsonx.data** console.

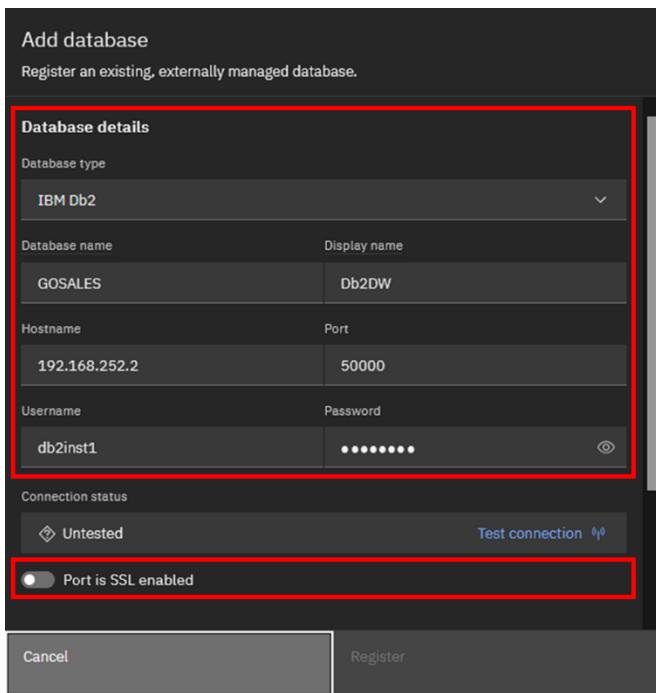
- Click the **Add component** dropdown menu at the right-side of the screen. Select **Add Database**.



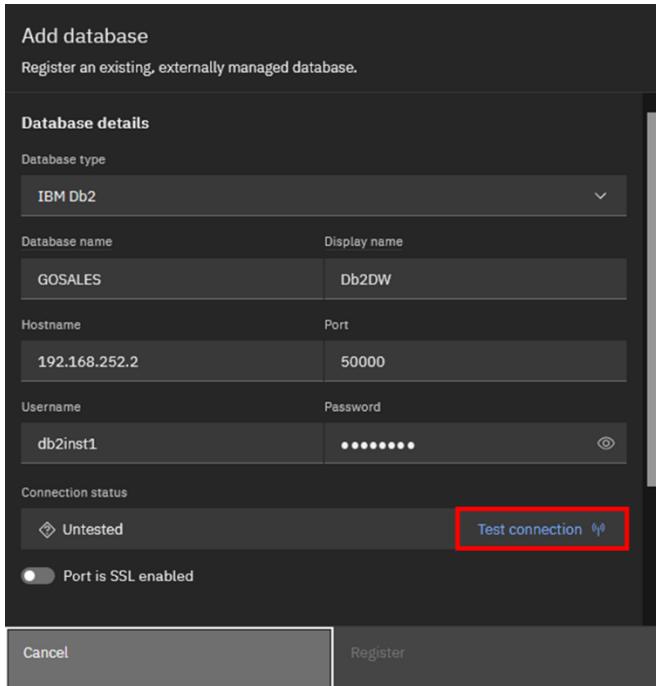
- In the **Add database** pop-up window, select/enter the following pieces of information in the **Database details** section:

- Database type:** IBM Db2
- Database name:** GOSALES
- Display name:** Db2DW
- Hostname:** 192.168.252.2
- Port:** 50000
- Username:** db2inst1
- Password:** db2inst1
- Port is SSL enabled:** <toggled off>

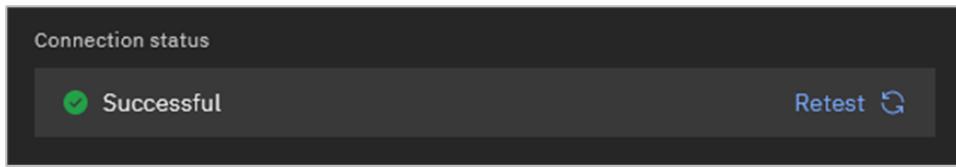
**Note:** Please verify that you have entered everything exactly as above (with no additional spaces before or after any of the values you entered). Also, ensure that the **Port is SSL enabled** toggle, which is toggled off by default, remains toggled off (or any attempt to access the Db2 database later will fail with a communication error).



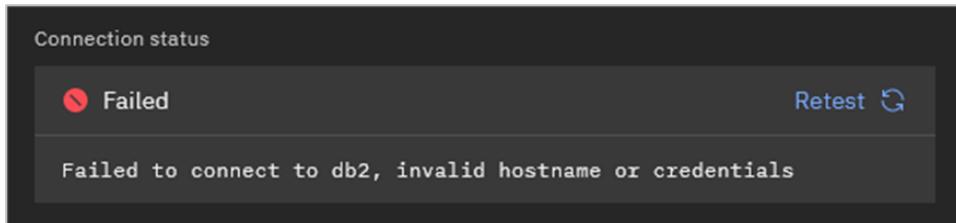
4. Click **Test connection**.



If the connection test is successful then the status indicates as such and should look like the image below:



However, if the test wasn't successful then you'll receive an error message as shown below with the reason for the failure. Re-enter all of the fields again and click **Retest** to test the connection again.



When a new bucket or database is added to watsonx.data, a new catalog is created as well. You only need to specify the name of the catalog.

5. Scroll down to the **Associated catalog** section. Enter **db2catalog** for the **Catalog name**. Then, click **Register**.

A screenshot of the "Add database" dialog. The "Associated catalog" section is highlighted with a red box, and the "Catalog name" input field contains "db2catalog" and is also highlighted with a red box. The "Register" button at the bottom right is highlighted with a blue box.

Add database

Register an existing, externally managed database.

GOSALES	Db2DW
Hostname	Port
192.168.252.2	50000
Username	Password
db2inst1	*****

Connection status

Successful Retest ⌂

Port is SSL enabled

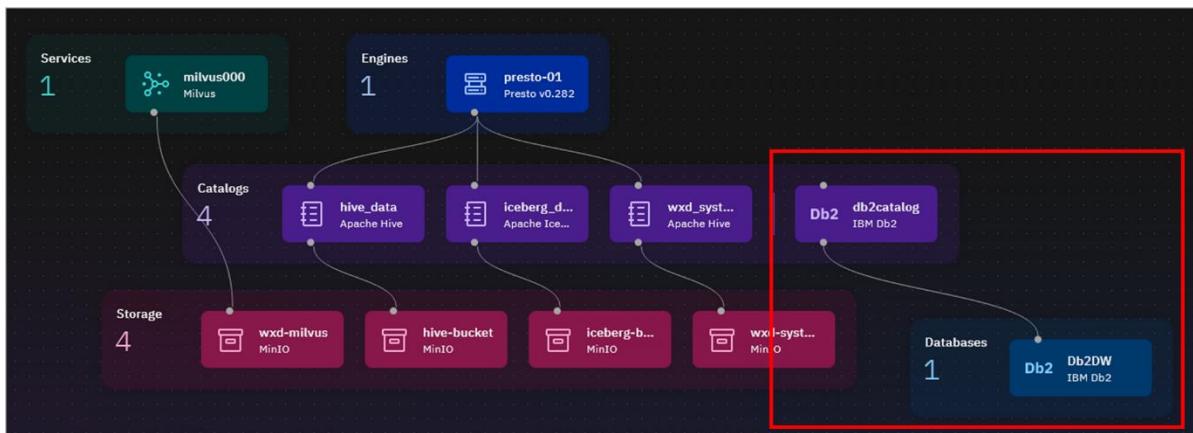
**Associated catalog** ⓘ

Catalog name

db2catalog

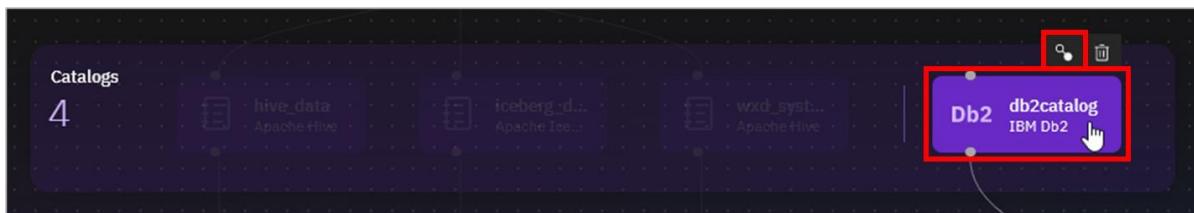
Cancel Register

The **Db2DW** database and **db2catalog** catalog have been added to watsonx.data and are now reflected in the topology view of the infrastructure components.

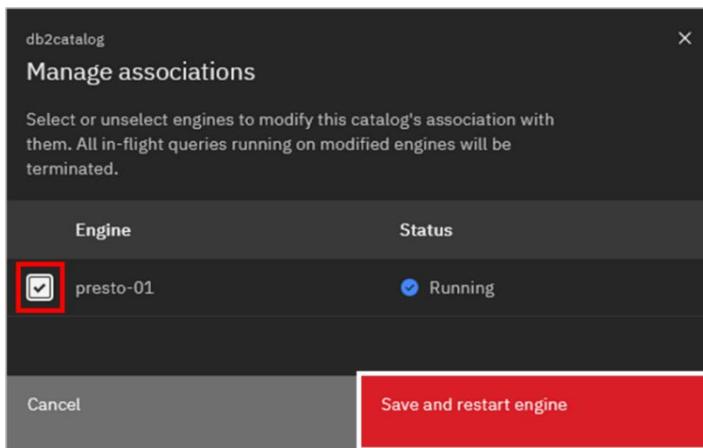


The **db2catalog** catalog is automatically associated with the **Db2DW** database, but to be able to query data from this database, the **db2catalog** catalog must also be associated with the **presto-01** engine (note how there currently is no line between **db2catalog** and **presto-01**).

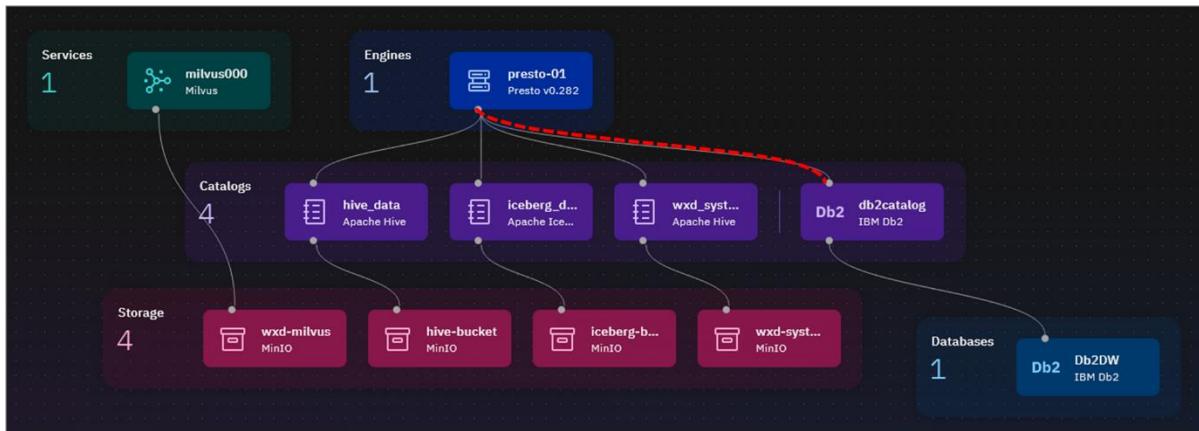
6. Hover your mouse pointer over the **db2catalog** catalog tile and the **Manage associations** icon will appear. Click the **Manage associations** icon.



7. In the **Manage associations** pop-up window, select the checkbox for the **presto-01** engine and then click **Save and restart engine**.



A line now connects **presto-01** with **db2catalog**, indicating that they're associated.



**Note:** If you ever have a need to remove a database (or bucket) and its associated catalog, you would first need to disassociate the Presto query engine with the database and its catalog. This is done by hovering over the catalog tile, clicking the **Manage associations** icon, unchecking the checkbox for the **presto-01** engine, and then clicking **Save and restart engine**. Once disassociated, you can hover over the catalog tile again such that the **Delete** (trash can) icon appears, click the **Delete** icon, and then confirm the deletion. Doing so will also remove the database. Do not do this now.

Next you will add the PostgreSQL database. The password for the PostgreSQL admin account is specific to your environment and the following steps will extract it for you.

8. Open a terminal command window to the watsonx.data server as the **root** user.
9. The administrative user for the PostgreSQL database is **admin**. Run the following command to extract and display the **password**. Copy the value shown to a location you can refer to later (be careful not to include any spaces before or after the password when you copy it).

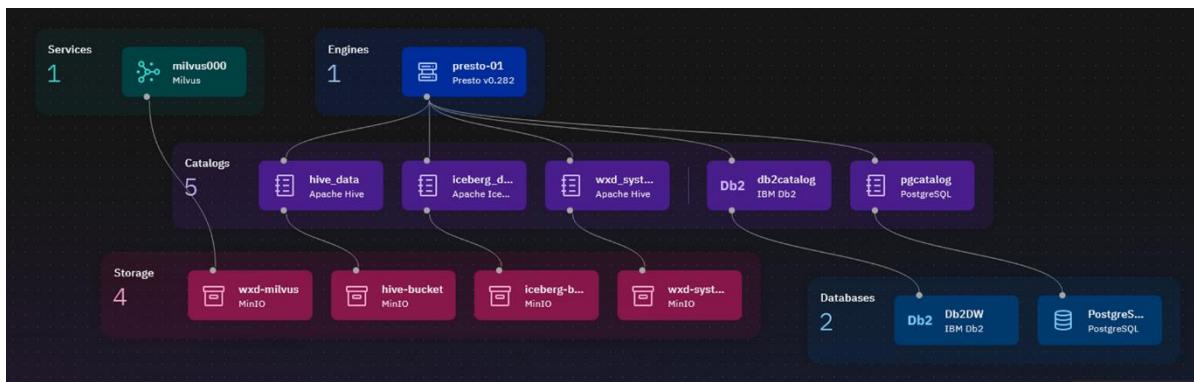
**Note:** This is a single command but is split over two lines in the published text below. It should copy and paste as a single line, but if it pastes as two separate lines then it will be interpreted as two separate commands and will fail. In this case, copy and paste such that this is run as a single command. Alternatively, remember that all of the commands and SQL statements can be found in this [text document](#). Download this file and copy the command from that file instead.

```
docker exec ibm-lh-postgres printenv | grep POSTGRES_PASSWORD | sed  
's/.*=//'
```

10. Repeat Steps 2-7 from above to add the PostgreSQL database to watsonx.data. Use the following information:

- **Database type:** PostgreSQL (under the **From Others** section)
- **Database name:** gosales
- **Display name:** PostgreSQLDB
- **Hostname:** 192.168.252.2
- **Port:** 5432
- **Username:** admin
- **Password:** <The password generated in the previous step>
- **Port is SSL enabled:** <toggled off>
- **Catalog name:** pgcatalog

With both the Db2 and PostgreSQL databases added, the topology should look like the image below.



11. Select the **Data manager** icon (grid) from the left-side menu.

12. Expand both the db2catalog and pgcatalog catalogs.

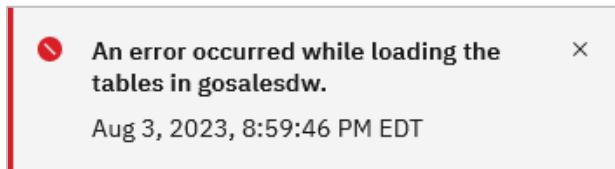
The screenshot shows the 'Catalogs associated' section of the Presto interface. It lists several catalogs: db2 catalog, hive\_data, iceberg\_data, pgcatalog, and wxd\_system\_data. The 'db2 catalog' and 'pgcatalog' entries are highlighted with red boxes.

Schema and table information are shown for both of these databases. Both include a schema with tables containing sales data for the fictional Great Outdoors company (the schema is called **gosalesdw** in PostgreSQL and **GOSALES DW** in Db2; the difference in case is simply due to how they happened to be created in each of these databases). Remember that there is also a copy of this same data in object storage, managed by the **hive\_data** catalog.

If you do not see any tables associated with the **gosalesdw (GOSALES DW)** schema in either of the catalogs, hover your mouse pointer at the far right of the line for the catalog until you see the Refresh icon appear. When you see the icon, click it. Likewise, if you didn't actually see the **gosalesdw** schema itself under the catalog, refresh the catalog in the same way.

The screenshot shows the 'Catalogs associated' section of the Presto interface. It lists several catalogs: db2 catalog, hive\_data, iceberg\_data, pgcatalog, and wxd\_system\_data. The 'db2 catalog' entry is expanded, showing its contents. A red box highlights the refresh icon (a circular arrow) next to the 'db2 catalog' entry.

**Note:** If for some reason the list of tables does not refresh, or if you see a message appear in the upper-right corner saying, “*An error occurred while loading the tables in gosalesdw*”, similar to the image below, you can still proceed. The table information is in the catalog and queries that reference these tables will still work. Please proceed with the rest of the lab.



You now have copies of the Great Outdoors company sales tables in object storage (hive\_data), Db2 (db2catalog), and PostgreSQL (pgcatalog). This isn’t likely something you’d have in a real-world scenario – after all, the benefit of being able to federate access to multiple data sources is to avoid data duplication. However, it’s being done here for the purposes of highlighting Presto’s federation capabilities.

You will now see how you can run a federated query that combines data from all three of these data sources.

13. Select the **Query workspace** ( icon) from the left-side menu.

14. Copy and paste the query below into the **SQL worksheet**. Click **Run on presto-01**.

This sample query could be used by the fictional business to determine which purchasing method is associated with the largest orders. The query accesses five tables, one of which is in PostgreSQL (**yellow**), two are in Db2 (**green**), and two are in watsonx.data's object storage (**pink**).

```
select pll.product_line_en as product,
       md.order_method_en as order_method,
       sum(sf.quantity) as total
  from
    pgcatalog.gosalesdw.sls_order_method_dim as md,
    db2catalog.GOSALES DW.SLS_PRODUCT_DIM as pd,
    db2catalog.GOSALES DW.SLS_PRODUCT_LINE_LOOKUP as pll,
    hive_data.gosalesdw.sls_product_brand_lookup as pbl,
    hive_data.gosalesdw.sls_sales_fact as sf
 where
    pd.product_key = sf.product_key
    and md.order_method_key = sf.order_method_key
    and pll.product_line_code = pd.product_line_code
    and pbl.product_brand_code = pd.product_brand_code
 group by pll.product_line_en, md.order_method_en
 order by product, order_method;
```

**Note:** If the SQL statement pastes into the worksheet as a single line, you can nicely format it by clicking the **Format worksheet** icon.



After running the query, its **Result set** is found at the bottom of the screen.

The screenshot shows the IBM Watsonx.data interface. In the center, there is a code editor window titled "Untitled 1" containing the following SQL query:

```
1 select
2     pll.product_line_en as product,
3     md.order_method_en as order_method,
4     sum(sf.quantity) as total
5 from
6     pgcatalog.gosalesdw.sls_order_method_dim as md,
7     db2catalog.gosalesdw.sls_PRODUCT_DIM as pd,
```

Below the code editor, the "Worksheet results 1" section displays the query's output:

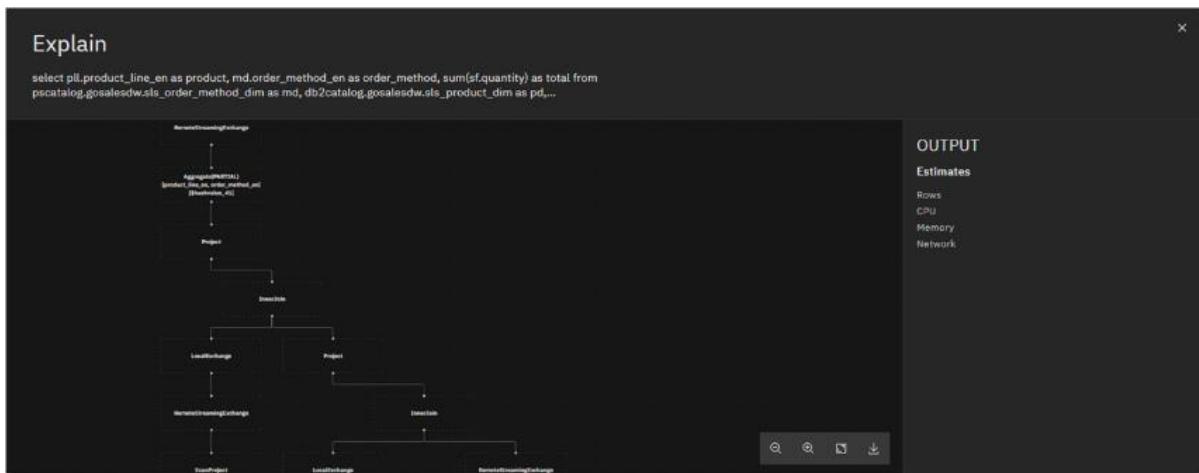
product	order_method	total
Camping Equipment	E-mail	1413084
Camping Equipment	Fax	413958
Camping Equipment	Mail	348058

The "Result set" tab is selected in the results panel. The entire results section is highlighted with a red box.

15. Click the **Explain** button in the menu bar of the worksheet.

The screenshot shows the "Untitled 1" worksheet with the "Explain" button highlighted by a red box in the top right corner of the toolbar. The code editor contains the same SQL query as the previous screenshot.

The visual explain output for this query looks a lot more interesting than the one shown earlier. If you scroll through the visual explain output, you'll see five **ScanProject** leaf nodes in the tree. These correspond to the five tables being read.



16. Click the X in the upper-right corner of the Explain window to close it.

**Optional step #1:** Here are two other queries you can try running as well, which combine data from the same three data sources.

1. The following query displays all Canadian and Mexican employees, along with their region and country. This is the kind of query that a reporting tool might generate, based on input from the user.

```
select distinct branch_region_dim.region_en region,
               branch_region_dim.country_en country,
               emp_employee_dim.employee_name employee
        from hive_data.gosalesdw.go_region_dim branch_region_dim,
             pgcatalog.gosalesdw.emp_employee_dim emp_employee_dim,
             db2catalog.GOSALES DW.GO_BRANCH_DIM go_branch_dim
       where branch_region_dim.country_en in ('Canada', 'Mexico') and
             branch_region_dim.country_code = go_branch_dim.country_code and
             emp_employee_dim.branch_code = go_branch_dim.branch_code
    order by region, country, employee;
```

2. In many businesses, departments (or organizations, as they're called in this dataset) are hierarchical, in that a department falls under another department, which in turn falls under another one. This query displays the two parent departments for a given set of departments. As in the previous query, this is the kind of query that a reporting tool might generate.

```
select gosalesdw.go_org_dim.organization_key,
       go_org_dim_1.organization_parent as org_level1_code,
       go_org_name_lookup_1.organization_name_en as org_level1_name,
       gosalesdw.go_org_dim.organization_parent as org_level2_code,
       go_org_name_lookup_2.organization_name_en as org_level2_name,
       gosalesdw.go_org_dim.organization_code as org_code,
       gosalesdw.go_org_name_lookup.organization_name_en as org_name
  from db2catalog.GOSALESDW.GO_ORG_NAME_LOOKUP go_org_name_lookup_2
    inner join
      hive_data.gosalesdw.go_org_dim
        inner join
          pgcatalog.gosalesdw.go_org_name_lookup
            on hive_data.gosalesdw.go_org_dim.organization_code =
               pgcatalog.gosalesdw.go_org_name_lookup.organization_code
            on go_org_name_lookup_2.organization_code =
               hive_data.gosalesdw.go_org_dim.organization_parent
              inner join
                pgcatalog.gosalesdw.go_org_name_lookup go_org_name_lookup_1
                  inner join
                    hive_data.gosalesdw.go_org_dim go_org_dim_1
                      on go_org_name_lookup_1.organization_code =
                         go_org_dim_1.organization_parent
                      on hive_data.gosalesdw.go_org_dim.organization_parent =
                         go_org_dim_1.organization_code
            where (hive_data.gosalesdw.go_org_dim.organization_code
                   between '1700' and '5730')
              order by org_code;
```

**Optional step #2:** Try running the two queries above from within the Presto CLI as well. You should get the same results in the Presto CLI as you did in the Query workspace.

## 10. Offloading Data from a Data Warehouse

The previous section described the federated query capability of watsonx.data and Presto. In addition to it making it easy to access data across the enterprise, federation also opens the door to significant cost-savings for clients.

Many enterprises maintain expensive data warehouses, whether they are on-premises or in the cloud. These warehouses support mission-critical workloads with low-latency and high-performance requirements that justify their high costs. However, there is also often data in these warehouses that are supporting less important or less stringent workloads. This could be infrequently accessed (cold) data being kept for audit purposes, or data used for business intelligence, reporting, and data science. It makes financial sense for this data to be offloaded to a lower cost solution but keeping it in the warehouse has just been the easier path to take. The result, though, is periodic scaling of the data warehouse to support growing workloads, which can be very expensive.

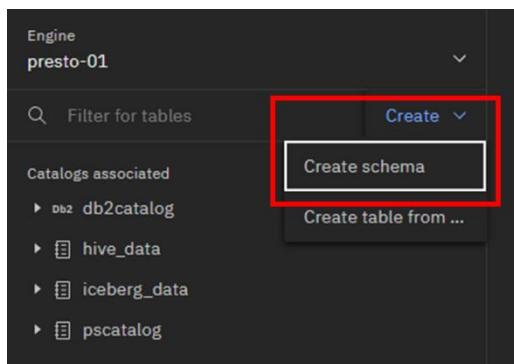
With watsonx.data, enterprises can offload data that really doesn't need to be in the warehouse to lower-cost object storage. With fit-for-purpose engines, data warehousing costs can be reduced by offloading workloads from a data warehouse to watsonx.data. Specifically, applications that need to access this data can query it through Presto (or Spark). This includes being able to combine the offloaded data with the data that remains in the warehouse. This section will show an example of how this can be done with Db2 (but the steps are equivalent for other data warehouse products).

Additionally, external engines that support the Iceberg open table format can also work directly with data in watsonx.data's object storage. For example, Db2 and Netezza have been enhanced to read from and write to the Iceberg table format (not all deployment options for Db2 and Netezza currently include this support; as of 1Q 2024 it is limited to Db2 Warehouse Gen3 on AWS, Db2 Warehouse on Red Hat OpenShift, Db2 Warehouse on Cloud Pak for Data, and Netezza Performance Server as a Service on AWS). This means that Db2 and Netezza can participate in the watsonx.data ecosystem as well, accessing the lakehouse data in object storage directly (just as Presto and Spark can). This, however, is beyond the scope of this lab.

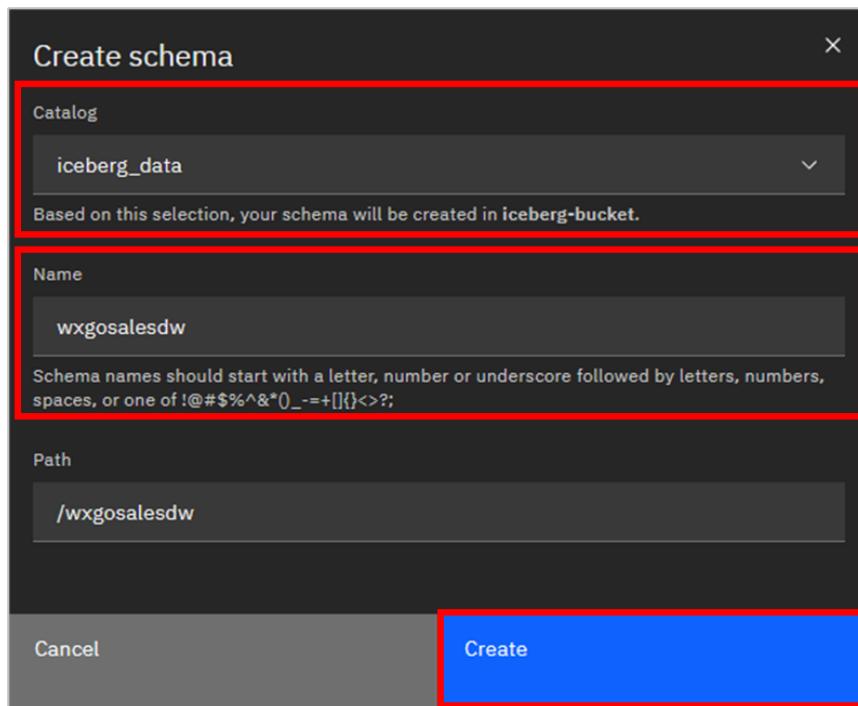
In the previous section you registered an existing Db2 environment with watsonx.data (calling it **Db2DW**) and created a catalog for it called **db2catalog**.

In this example scenario, you're going to "move" the **GOSALES DW.SLS\_SALES\_FACT** table from Db2 to watsonx.data's object storage. It will be created as an Iceberg table, in a new schema you create called **wxgosalesdw**, managed by the **iceberg\_data** catalog.

1. Select the **Data manager** (grid icon) from the left-side menu.
2. Go to the top of the navigation pane and click the **Create** dropdown menu. Select **Create schema**. (Note: These steps are using the web interface, but you can also create a schema through SQL.)



3. In the **Create schema** pop-up window, select/enter the following information, and then click the **Create** button. The **Path** field (representing the directory folder that gets created in the object storage bucket) gets filled in automatically, based on the schema name being entered. You can override the path with your own custom path name, but leave the default as-is.
  - **Catalog:** iceberg\_data
  - **Name:** wxgosalesdw



4. Expand the **iceberg\_data** catalog. The new schema should be listed.

The screenshot shows the Presto UI interface. At the top, it says "Engine" and "presto-01". Below that is a search bar labeled "Filter for tables" and a "Create" button. Under "Catalogs associated", there is a list of catalogs: "db2 db2catalog", "hive\_data", "iceberg\_data", and "pgcatalog". The "iceberg\_data" catalog is expanded, showing its contents: "my\_schema\_1", "newschema", and "wxgosalesdw". The "wxgosalesdw" schema is highlighted with a red box.

5. Select the **Query workspace** ( icon) from the left-side menu.

To create a new table with the same table definition as the original table, you can use the `CREATE TABLE AS SELECT` (CTAS) SQL statement.

6. Copy and paste the following SQL into the **SQL worksheet**. Click **Run on presto-01**.

```
create table iceberg_data.wxgosalesdw.sls_sales_fact  
as select * from db2catalog.GOSALESDW.SLS_SALES_FACT;
```

The screenshot shows the SQL worksheet titled "Untitled 1". It contains the SQL query: "create table iceberg\_data.wxgosalesdw.sls\_sales\_fact as select \* from db2catalog.GOSALESDW.SLS\_SALES\_FACT;". Below the query, the results are shown: "1 rows inserted". At the top right of the worksheet, there is a "Run on presto-01" button, which is highlighted with a red box.

The result shown at the bottom of the worksheet shows the number of rows that have been inserted from the source table to the new table (446,023 rows).

The screenshot shows the SQL worksheet results. The query was run on presto-01 and completed in 7.799 seconds. The result set shows one row with the value "446023". The "rows" header is highlighted with a red box.

rows
446023

In a real-world scenario, you would then remove the table from the data warehouse. However, to keep the GOSALES DW sample dataset intact in your environment, you won't do that here.

As a test, you can run a federated query that combines the new table in object storage with existing tables in Db2.

7. Copy and paste the following SQL into the **SQL worksheet**. Click **Run on presto-01**.

```
select pll.product_line_en as product,
       sum(sf.quantity) as total
  from
    db2catalog.GOSALES DW.SLS_PRODUCT_DIM as pd,
    db2catalog.GOSALES DW.SLS_PRODUCT_LINE_LOOKUP as pll,
    iceberg_data.wxgosalesdw.sls_sales_fact as sf
 where
    pd.product_key = sf.product_key
    and pll.product_line_code = pd.product_line_code
 group by pll.product_line_en
 order by product;
```

This business query calculates total sales for each of the high-level product lines. The output should be similar to the image below.

The screenshot shows a SQL worksheet interface. At the top, there is a command line with the query: "select pll.product\_line\_en as product, sum(sf.quantity) as total from db2catalog.GOSALES". To the right of the command line, it says "Run time: 3.502s" with a green checkmark. Below the command line, there are two tabs: "Result set" (which is selected) and "Details". The "Result set" tab displays a table with two columns: "product" and "total". The table contains five rows of data:

product	total
Camping Equipment	27301149
Golf Equipment	5113701
Mountaineering Equipment	9900091
Outdoor Protection	12014445
Personal Accessories	34907705

## 11. Time Travel

The Iceberg open table format provides a number of benefits to users, including the ability to see a table as it existed at a point in the past. This *time travel* capability is useful in a number of different ways. For example, having the ability to query historical data is useful for auditing purposes. Or if an application corrupts table data in some way, you can imagine the value in being able to quickly reverse those changes by resetting the table to a known good state.

Iceberg uses *snapshots* to support this time travel capability. A snapshot represents the state of a table at some point in time. When data is modified in a table, such as inserting, updating, or deleting records, a new snapshot is created. Existing data files aren't modified, but new data files may be created depending on the operation and metadata files are updated. There are maintenance operations that can be used to clean up older snapshots and the corresponding data files that are no longer needed.

### 11.1. Table Rollback

Watsonx.data supports the ability to roll back a table to an earlier point in time using snapshots. When a table is first created with data added, a snapshot is created that represents that specific point in time. Each subsequent transaction then affects the data in the table creates a new snapshot. In this section you will work with an existing Iceberg table that has a single snapshot associated with it. You will then insert a row to the table, which creates a new snapshot. Finally, you will rollback the table to its initial state, represented by the original snapshot. In between you will observe the state of the data in the table.

1. Select the **Query workspace** () icon from the left-side menu.
2. Copy and paste the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**.

```
create table iceberg_data.my_schema.airport_id
  as select * from hive_data.ontime.airport_id;
```

The screenshot shows a SQL worksheet interface. At the top, there is a toolbar with various icons. Below the toolbar, a red box highlights the SQL statement: `create table iceberg_data.my_schema.airport_id as select * from hive_data.onetime.airport_id;`. To the right of the statement are buttons for 'Explain' and 'Run on presto-01'. Another red box highlights the 'Run on presto-01' button. Below the statement, the text 'Worksheet results 6' is displayed. A red box highlights the 'rows' section in the result set, which shows the value '6250'. The bottom right corner shows a status message: 'Run time: 1.661s' with a green checkmark.

There should be **6,250 rows** in this table, as shown in the output above.

3. Select the **Data manager** (grid icon) from the left-side menu.
4. Navigate to the **iceberg\_data > my\_schema > airport\_id** table (if you don't see the table, refresh the schema). Then, select the **Time travel** tab.

The screenshot shows the Data manager interface. On the left, a sidebar lists databases and schemas. Under 'iceberg\_data', 'my\_schema' is expanded, and 'airport\_id' is selected, highlighted with a red box. The main panel shows the 'Time travel' tab for the 'airport\_id' table. It displays a table with the following data:

Snapshot ID	Operation	Total records	Added records	Total deleted files	Committed at
539029806241330...	append	6250	6250	0	2023-07-31T17:54...

This tab lists the snapshots associated with the table (identified by a **Snapshot ID**). At this point there is only one snapshot and there are 6,250 total records (matching what you saw earlier).

5. Select the **Query workspace** (SQL icon) again from the left-side menu.
6. Copy and paste the following SQL statement into the **SQL worksheet** and click **Run on presto-01**. This inserts one new row into the table.

```
insert into iceberg_data.my_schema.airport_id
values (10000, 'North Pole: Reindeer Field');
```

7. Copy and paste (replacing the SQL already there) the following SQL statement into the **SQL worksheet** and click **Run on presto-01**.

```
select * from iceberg_data.my_schema.airport_id
where code = 10000;
```

You should see the row you inserted in the previous step.

Result set		Details	Run time: 1.231s
code	description		
10000	North Pole: Reindeer Field		
1–1 of 1 items			1 ▾ of 1 page ▾ ▾ ▾ ▾

8. Copy and paste (replacing the SQL already there) the following SQL statement into the **SQL worksheet** and click **Run on presto-01**.

```
select count(*) from iceberg_data.my_schema.airport_id;
```

You should now see a count of 6,251 rows.

9. Select the **Data manager** (grid icon) again from the left-side menu.

10. As before, navigate to the **iceberg\_data > my\_schema > airport\_id** table. Then, select the **Time travel** tab.

airport_id						
Catalog: iceberg_data   Schema: my_schema						
Table schema		Time travel	Data sample	DDL		
Search for snapshots						
Snapshot ID	Operation	Total records	Added records	Total deleted files	Committed at	⋮
539029806241330...	append	6250	6250	0	2023-07-31T17:54...	⋮
775501768280471...	append	6251	1	0	2023-07-31T18:09...	⋮

Note how there are now two snapshots shown. If you do not see a second snapshot then refresh your browser (for instance, using **<F5>** in the Firefox browser) and repeat this step.

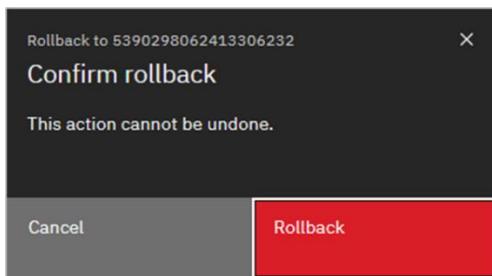
The second snapshot shows that there are 6,251 total rows, with 1 row having been added in this new version of the table.

You are now going to roll the table back to the first snapshot, representing the initial state of the table. This version of the table did *not* have the row you added.

11. Click the **overflow menu** icon (vertical ellipses) at the end of the row for the original snapshot (the first snapshot shown, with the earlier **Committed at timestamp** and **6250 Added records**). Click **Rollback**.

Snapshot ID	Operation	Total records	Added records	Total deleted files	Committed at	⋮
539029806241330...	append	6250	6250	0	2023-07-31T17:54...	
775501768280471...	append	6251	1	0	2023-07-	<b>Rollback</b>

12. In the **Confirm rollback** pop-up window, click **Rollback**.



**Note:** For reference, the following SQL statement will perform the equivalent roll back operation (DO NOT RUN THIS NOW):

```
call iceberg_data.system.rollback_to_snapshot
    ('my_schema', 'airport_id', <snapshotID>);
```

13. Select the **Query workspace** ( icon) again from the left-side menu.

14. Copy and paste the following SQL statement into the **SQL worksheet** and click **Run on presto-01**.

```
select * from iceberg_data.my_schema.airport_id
where code = 10000;
```

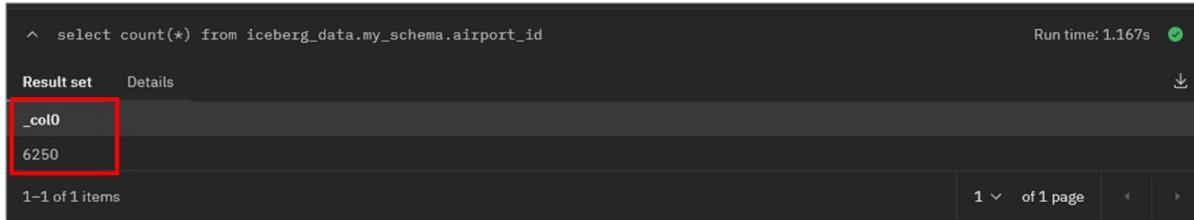
The row you added earlier is now gone!



A screenshot of a SQL worksheet interface. At the top, a query is shown: `^ select * from iceberg_data.my_schema.airport_id where code = 10000`. To the right, it says "Run time: 0s" with a green checkmark. Below the query is a table with two tabs: "Result set" (which is selected) and "Details". The "Result set" tab shows the message "0-0 of 0 items". On the far right of the table, there are navigation buttons for pages and rows, with "1 of 1 page" and "1" selected.

15. Copy and paste (replacing the SQL already there) the following SQL statement into the **SQL worksheet** and click **Run on presto-01**.

```
select count(*) from iceberg_data.my_schema.airport_id;
```



A screenshot of a SQL worksheet interface. At the top, the same query is shown: `^ select count(*) from iceberg_data.my_schema.airport_id`. To the right, it says "Run time: 1.167s" with a green checkmark. Below the query is a table with two tabs: "Result set" (selected) and "Details". The "Result set" tab shows a single row with one column labeled "\_col0" containing the value "6250". On the far right of the table, there are navigation buttons for pages and rows, with "1 of 1 page" and "1" selected.

You should see a count of 6,250 rows, which matches the original count when you first created the table. The table is back at the state it was in when it was first created, and the initial data was added!

## 11.2. Time Travel Queries

Time travel queries allow you to “look back in time” at your data. Note that this is not supported with Hive tables, only Iceberg tables (as the base support for this capability is provided by the Iceberg table format itself).

Watsonx.data supports the following query syntax options:

- `SELECT <columns> FROM <icebergTable> FOR TIMESTAMP AS OF TIMESTAMP <timestamp>;`
- `SELECT <columns> FROM <icebergTable> FOR VERSION AS OF <snapshotId>;`

In this section you will create a simple table and insert rows into it, over multiple transactions. You will then see how you can query the table as of different times using the time travel query syntax shown above.

The previous section showed how you can see a table's snapshots by looking at the **Time travel** tab for it in the **Data manager** page. This information is also available by querying the “`<tablename>$snapshots`” metadata table (where `<tablename>` is the name of the table in question).

**Note:** When you are instructed below to copy and paste a SQL statement into the **SQL worksheet**, clear any statement(s) you previously ran before running the new statement.

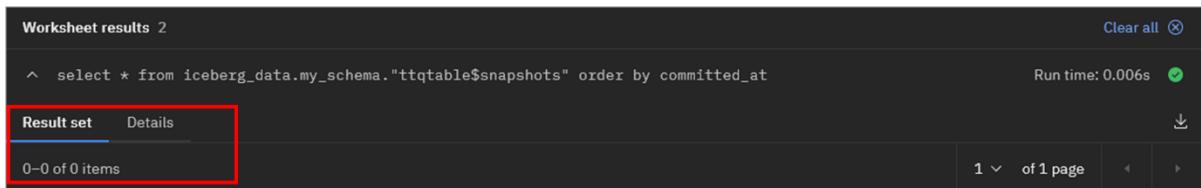
1. Select the **Query workspace** () icon from the left-side menu.
2. Create a table by copying and pasting the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**. This creates the table used in this section.

```
create table iceberg_data.my_schema.ttqtable (id int, name
char(10));
```

3. Query the list of snapshots for this table by copying and pasting the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**.

```
select * from iceberg_data.my_schema."ttqtable$snapshots"
order by committed_at;
```

The result set has zero rows. This indicates that there are currently no snapshots associated with the table (as no data has been added yet).



A screenshot of the Presto SQL worksheet interface. At the top, it says "Worksheet results 2". Below that is a query history entry: "`^ select * from iceberg_data.my_schema."ttqtable$snapshots" order by committed_at`". To the right of the query is "Run time: 0.006s" with a green checkmark. Below the query history is a results table. The table has two tabs at the top: "Result set" (which is selected and highlighted with a red border) and "Details". The main body of the table shows the message "0-0 of 0 items". At the bottom right of the results table are navigation controls: "1" and "of 1 page", and arrows for navigating through the results.

4. Copy and paste the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**. This inserts three rows into the table, all as part of the same transaction.

```
insert into iceberg_data.my_schema.ttqtable values (1, 'TV'), (2,
'Computer'), (3, 'Stereo');
```

5. Copy and paste the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**.

```
select * from iceberg_data.my_schema."ttqtable$snapshots"  
order by committed_at;
```

The result set has one row, indicating that there is one snapshot (which corresponds to the current state of the table after inserting the first set of data).

committed_at	snapshot_id	parent_id	operation	manifest_list
2024-02-01 16:40:42.548 UTC	1465650329006827451	null	append	s3a://iceberg-bucket/my_schema/ttqtable/metadata/snap-146565032900

6. Copy and paste the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**. This inserts three additional rows to the table, all as part of the same transaction.

```
insert into iceberg_data.my_schema.ttqtable values (4, 'Phone'),  
(5, 'Tablet'), (6, 'Camera');
```

7. Copy and paste the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**. This inserts three additional rows to the table, all as part of the same transaction.

```
insert into iceberg_data.my_schema.ttqtable values (7,  
'Headphones'), (8, 'Microphone'), (9, 'Watch');
```

8. Copy and paste the following SQL statement into the **SQL worksheet**. Click **Run on presto-01**.

```
select * from iceberg_data.my_schema."ttqtable$snapshots"  
order by committed_at;
```

The result set has three rows, with the snapshots corresponding to the changes made by the three insert transactions just performed. Since the result set is ordered by the **committed\_at** column (by default in ascending order), the three rows are in chronological order, with the most recent change to the table corresponding to the last row/snapshot.

committed_at	snapshot_id	parent_id	operation	manifest_list
2024-02-01 16:40:42.548 UTC	1465650329006827451	null	append	s3a://iceberg-bucket/my_schema/ttqtable/metadata/snap-1
2024-02-01 16:43:27.962 UTC	5302208642517846066	1465650329006827451	append	s3a://iceberg-bucket/my_schema/ttqtable/metadata/snap-2
2024-02-01 16:43:40.560 UTC	2901270876951439131	5302208642517846066	append	s3a://iceberg-bucket/my_schema/ttqtable/metadata/snap-3

- For each of the three snapshots in the result set, copy and paste the **committed\_at** values (including date, time, and time zone) and the **snapshot\_id** values to a location you can refer to later. They will be referred to as **Tran1Time**, **Tran1SnapshotID**, **Tran2Time**, **Tran2SnapshotID**, **Tran3Time**, and **Tran3SnapshotID** in subsequent instructions.

committed_at	snapshot_id
2024-02-01 16:40:42.548 UTC	1 1465650329006827451
2024-02-01 16:43:27.962 UTC	2 5302208642517846066
2024-02-01 16:43:40.560 UTC	3 2901270876951439131

As mentioned earlier, there are two ways to issue time travel queries. You can explicitly specify a snapshot ID or you can provide a timestamp. The snapshot ID must refer to one of the table's snapshots, but the timestamp can be any timestamp since the table was created and data added (for a given timestamp, the relevant snapshot will automatically be determined and used).

You will now run different queries using both of these methods. In each query below, replace the timestamp/snapshot ID placeholder with the corresponding value you copied from the snapshot query results earlier. Ensure that the single quotes (') are not removed from the queries when you paste the timestamp values in.

10. For each of the queries below, clear the **SQL worksheet** and copy and paste the query into it. Click **Run on presto-01**.

<b>Query 1:</b> Retrieve data from the initial state of the table, using the corresponding snapshot ID.	<b>Expected results:</b> 3 rows
---	---------------------------------

```
select * from iceberg_data.my_schema.ttqtable for version as of  
Tran1SnapshotID order by id;
```

<b>Query 2:</b> Retrieve data from the initial state of the table, using the corresponding timestamp.	<b>Expected results:</b> 3 rows
---	---------------------------------

```
select * from iceberg_data.my_schema.ttqtable for timestamp as of  
cast('Tran1Time' as timestamp with time zone) order by id;
```

<b>Query 3:</b> Retrieve data post-second transaction, using the corresponding snapshot ID.	<b>Expected results:</b> 6 rows
---	---------------------------------

```
select * from iceberg_data.my_schema.ttqtable for version as of  
Tran2SnapshotID order by id;
```

<b>Query 4:</b> Retrieve data post-third transaction, using the corresponding timestamp.	<b>Expected results:</b> 9 rows
--	---------------------------------

```
select * from iceberg_data.my_schema.ttqtable for timestamp as of  
cast('Tran3Time' as timestamp with time zone) order by id;
```

<b>Query 5:</b> Retrieve data from a point in time between the second and third transaction. (This shows that you can specify an arbitrary timestamp that doesn't have to match the timestamp of a snapshot.)	<b>Expected results:</b> 6 rows
---	---------------------------------

```
select * from iceberg_data.my_schema.ttqtable for timestamp as of  
cast('Pick-a-time-between-Tran2Time-and-Tran3Time' as timestamp  
with time zone) order by id;
```

**Query 6:** Retrieve data from a time before the table was created.

**Expected results:** Fails, as the table did not exist at this time.

```
select * from iceberg_data.my_schema.ttqtable for timestamp as of  
cast('2023-01-01 00:00:00.000 UTC' as timestamp with time zone)  
order by id;
```

This ability to run “as of” queries is a significant feature of Presto (and Spark!) in watsonx.data and it’s made possible by the Iceberg open table format. You can imagine how useful this would be when running queries for auditing and regulatory purposes.

## 12. Summary

Congratulations on completing this lab! You gained hands-on experience in the following areas of watsonx.data:

- The watsonx.data web-based user interface, including infrastructure management, data management, running SQL statements, and managing user access
- The Presto web interface and the Presto command line interface (CLI)
- MinIO object storage
- Ingesting data into watsonx.data
- Creating schemas and tables
- Running queries that combine data from multiple data sources
- Offloading tables from Db2 into watsonx.data
- Rolling back a table to a previous point in time
- Running time travel queries

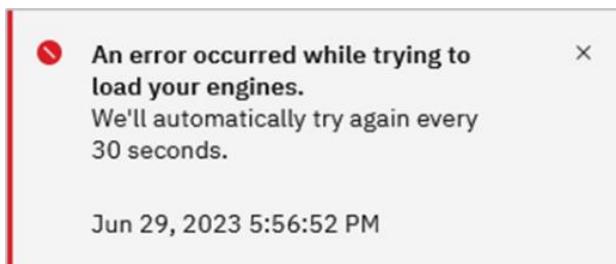
As a result, you should feel much more confident in your ability to demonstrate these capabilities to your clients, as well as be able to discuss the business value of watsonx.data with them.

## Appendix A. Troubleshooting

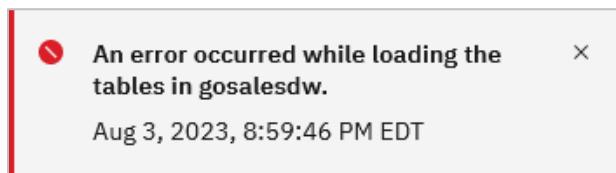
Unexpected errors, while rare, may occur. Resolutions to some of the more common issues encountered by users are included below:

### watsonx.data Web Interface:

- **The watsonx.data web interface won't start:** If you cannot bring up the watsonx.data web interface (console) in your browser (unable to connect, connection has timed out, and so on) then the environment is likely still starting up. Wait at least 10 minutes and try again. If you are still not able to bring up the web interface after this time, then you can try restarting watsonx.data as described in section [4.5. Stopping and Starting watsonx.data](#).
- **Engines, catalogs, buckets, or databases cannot be loaded:** If you see console errors stating that watsonx.data's engines, catalogs, buckets, or databases cannot be loaded (similar to the image below), this is likely due to watsonx.data still being started. Wait a minute and refresh the browser window.



- **Table or schema creation fails:** It is infrequent, but when performing an action within the console, such as creating a schema or table, you might encounter an error. The action may have in fact completed, but if not then repeat the action.
- **Loading gosalesdw tables fails:** You may see a message appear in the upper-right corner saying, “*An error occurred while loading the tables in gosalesdw.*” This occurs when watsonx.data is unable to list the GOSALES tables under db2catalog in the Data management page. It can be safely ignored; queries that access these tables will still work.



## MinIO Web Interface:

- **The MinIO web interface won't open in your browser:** If you cannot start the MinIO web interface but the watsonx.data web interface starts fine, then something in the environment is preventing http (non-https) connections from working. Please restart watsonx.data as described in section [4.5. Stopping and Starting watsonx.data](#) (ensure that LH\_RUN\_MODE=diag is set prior to starting watsonx.data).

## Miscellaneous:

- **Command/SQL statement fails:** If you are instructed to copy and paste a command or SQL statement from this lab guide and the command/statement fails, it may be because it was copied incorrectly. Try to copy and paste the command/statement again. Alternatively, all of the commands and SQL statements can also be found in this [text document](#). Download this file and copy the text from the file instead.

## Appendix B. Acknowledgements

A big “Thank you!” is owed to George Baklarz, who created the TechZone image used in this lab. There was a tremendous amount of work involved in installing and configuring all the software and populating the various data sources included in the image. George (and his lab material) was also a great source of information on the “how to’s” of watsonx.data and of the lab environment.

## Appendix C. Revision History

Date	Changes
-	Original version.
2023-09-09	Removed WireGuard VPN dependency and instructions.
2023-09-12	Minor updates.
2023-09-30	Updated query used in testing offloaded Db2 data and other minor updates.
2023-10-03	Updated Slack channel and clarified user-mgmt command.
2023-10-19	Db2 server password has expired. Added workaround steps.
2023-10-25	Clarified which environment tile to provision.
2023-11-01	Clarified purpose of section 4.3 (command line access).
2023-11-06	Updated lab guide to watsonx.data 1.0.3 TechZone image.
2023-12-13	Updated lab guide to watsonx.data 1.1.0 TechZone image.
2023-12-15	Added instructions for removing a database and its catalog.
2024-01-08	Clarified ssh instructions.
2024-01-30	Updated lab guide to watsonx.data 1.1.1 TechZone image.
2024-02-01	Added time travel queries section.
2024-04-26	Updated lab guide to watsonx.data 1.1.3 TechZone image.