

## Azure Databricks

Azure Databricks is a unified analytics platform provided by Microsoft in collaboration with Databricks, a company founded by the creators of Apache Spark. It is designed to simplify big data and artificial intelligence (AI) analytics tasks. Azure Databricks integrates with Azure services and provides a collaborative environment for data scientists, data engineers, and business analysts to work together on big data and machine learning projects.

Key features of Azure Databricks include:

1. **Unified Analytics Platform:** Azure Databricks combines data engineering, data science, and business analytics in a single platform. Users can perform tasks such as data preparation, exploration, analysis, and visualization seamlessly.
2. **Apache Spark-Based:** Azure Databricks leverages Apache Spark, an open-source distributed computing framework, for processing large-scale data sets. Spark provides high-performance processing capabilities for tasks like batch processing, streaming analytics, machine learning, and graph processing.
3. **Collaborative Environment:** Azure Databricks offers a collaborative workspace where teams can work together on data projects. It includes features such as notebooks, which allow users to write and execute code in languages like Python, Scala, SQL, and R, as well as collaborate through shared notebooks.
4. **Integration with Azure Services:** Azure Databricks integrates tightly with other Azure services such as Azure Blob Storage, Azure SQL Data Warehouse, Azure Cosmos DB, Azure Data Lake Storage, and more. This integration enables seamless data ingestion, storage, and analysis workflows.
5. **Machine Learning Support:** Azure Databricks provides built-in support for machine learning workflows. Users can leverage libraries like MLLib (Spark's machine learning library) and MLflow (an open-source platform for managing the machine learning lifecycle) to build, train, deploy, and manage machine learning models.
6. **Scalability and Performance:** Azure Databricks is designed to scale dynamically to handle large volumes of data and processing workloads. It offers features such as auto-scaling and optimized performance for efficient data processing.

Overall, Azure Databricks simplifies the process of building and deploying big data and AI solutions on the Microsoft Azure cloud platform, enabling organizations to derive valuable insights from their data more effectively.

## Use cases of Azure Databricks:

Azure Databricks is utilized across various industries and domains for a wide range of use cases. Here are some common examples:

1. **Data Engineering and ETL:** Organizations use Azure Databricks to streamline their data engineering processes, including Extract, Transform, Load (ETL) tasks. It enables the efficient processing of large volumes of data from disparate sources, transforming it into a usable format for analysis and reporting. For instance, companies can use Azure

Databricks to ingest data from sources like databases, IoT devices, logs, and files, perform data cleansing, aggregation, and enrichment, and then load the processed data into data warehouses or data lakes.

2. **Real-time Stream Processing:** Azure Databricks is well-suited for real-time stream processing use cases. Organizations can leverage its integration with Apache Spark Streaming and Structured Streaming to analyze and gain insights from continuous streams of data in real-time. This is particularly useful in scenarios such as fraud detection, IoT data processing, clickstream analysis, and monitoring system health.
3. **Machine Learning and AI:** Azure Databricks provides a powerful platform for developing and deploying machine learning (ML) and artificial intelligence (AI) models at scale. Data scientists and ML engineers can use its collaborative environment to explore data, build ML models using libraries like MLlib and TensorFlow, and deploy models into production seamlessly. Use cases include predictive maintenance, customer churn prediction, recommendation systems, sentiment analysis, and image recognition.
4. **Data Analytics and Business Intelligence:** Organizations leverage Azure Databricks for advanced data analytics and business intelligence (BI) initiatives. By combining its analytics capabilities with visualization tools like Power BI and Tableau, users can gain insights from large datasets, perform ad-hoc analysis, and create interactive dashboards and reports. Use cases include market analysis, customer segmentation, financial forecasting, and operational analytics.
5. **Genomics and Healthcare Analytics:** In the healthcare and life sciences industry, Azure Databricks is used for genomics data analysis and healthcare analytics. Researchers and healthcare professionals can analyze large-scale genomics datasets to identify patterns, correlations, and genetic variations associated with diseases. Additionally, Azure Databricks can facilitate population health analytics, personalized medicine, drug discovery, and clinical trial optimization.
6. **Financial Services:** Financial institutions utilize Azure Databricks for risk management, fraud detection, algorithmic trading, and customer analytics. By processing vast amounts of financial data in real-time, organizations can identify market trends, detect anomalies, optimize trading strategies, and improve customer experiences.

In this tutorial, we're setting up an environment in Azure Databricks for data analytics and machine learning tasks. The end goal is to provide users with a unified platform where they can collaborate, analyze large datasets, build machine learning models, and derive valuable insights from their data. By creating a workspace and setting up a cluster, users can start performing various data processing and analytics tasks efficiently within Azure Databricks.

## To begin with the Lab:

### Creating Workspace

1. In your Azure Portal search for Azure Databricks and choose the service accordingly.

# Azure Databricks

Microsoft

The screenshot shows the Azure Marketplace page for 'Azure Databricks'. It features a dark icon with three stacked squares, the title 'Azure Databricks' in bold, a rating of '★ 4.5 (389 ratings)', and a 'Create' button. The page also includes a 'Plan' dropdown set to 'Azure Databricks'.

2. Then you need to choose your resource group and give your workspace a name then you have to choose the pricing tier.
3. For pricing tier, you have to choose Trial. It will provide you free access for 14 days but you will be charged based on your spark plugs.

## Create an Azure Databricks workspace

Basics Networking Encryption Security & compliance Tags Review + create

### Project Details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *	Azure Pass - Sponsorship
Resource group *	new-grp
	<a href="#">Create new</a>

### Instance Details

Workspace name *	databricksworkspace120
Region *	North Europe
Pricing Tier *	Trial (Premium - 14-Days Free DBUs)
Managed Resource Group name	Enter name for managed resource group

4. Then just move to the review page and create your workspace. Once the deployment is completed then move to resources.

**new-grp\_databricksworkspace120 | Overview**

**Your deployment is complete**

Deployment name : new-grp\_databricksworkspace120  
Subscription : Azure Pass - Sponsorship  
Resource group : new-grp

Start time : 14/5/2024, 11:29:34 am  
Correlation ID : 64c1a8ca-06b8-48ba-8454-1ec3a4edd3ba

**Deployment details**

**Next steps**

**Go to resource**

## 5. Now from the dashboard of data bricks you have to click on launch workspace.

**databricksworkspace120**

Status : Active  
Resource group : new-grp  
Location : North Europe  
Subscription : Azure Pass - Sponsorship

Managed Resource Group : databricks-rg-databricksworkspace120-yw2h6q3g3yy  
URL : <https://adb-93258239280560.azuredatabricks.net>  
Pricing Tier : Trial (Premium - 14-Days Free DBUs) (Click to change)

**Tags**  
Tags (edit) : Add tags

**Launch Workspace**

## 6. Then you will be in the data bricks workspace.

**Get started**

- Import and transform data
- Notebook
- SQL query editor
- AutoML

**Pick up where you left off**

No recent items  
Start exploring and your recently viewed items will show up here.

Popular

No popular items  
Start exploring and popular items in your workspace will show up here.

**New**

- Workspace
- Recents
- Catalog
- Workflows
- Compute

**Compute**

- SQL
- SQL Editor
- Queries
- Dashboards
- Alerts
- Query History
- SQL Warehouses

**Data Engineering**

- Job Runs
- Data Ingestion
- Delta Live Tables

**Machine Learning**

- Playground
- Experiments
- Features
- Models
- Serving

**Marketplace**

**Partner Connect**



## 1. Now to create a cluster, first and foremost, you need to click compute. Then click on create compute.

The screenshot shows the 'Compute' section of the Azure Databricks portal. On the left sidebar, the 'Compute' option is selected and highlighted with a red box. At the top right, there is a message: 'Free trial ends in 14 days. Upgrade to Premium in Azure Portal'. Below the header, there are tabs for 'All-purpose compute', 'Job compute', 'SQL warehouses', 'Vector Search', 'Pools', and 'Policies'. A search bar and filter options ('Created by', 'Only pinned') are also present. The main area displays a table with columns: State, Name, Policy, Runtime, Active mem..., Active cores, Active DBU ..., Source, Creator, Notebooks, and a gear icon. A large blue '+' button is centered above the table. Below it, a message reads: 'No compute. Create compute to run workloads from your notebooks and jobs. Learn more about best practices for compute configuration.' A 'Create compute' button is located at the bottom of this section.

2. Then you need to give your cluster a name and choose a single node.
3. After that keep everything to default or you can check the properties here. Like what is the runtime version and the node type?
4. Then just click on create compute.

The screenshot shows the 'New compute' configuration page for a 'Data Cluster'. The top navigation bar includes 'Compute > New compute >' and a 'Create compute' button. A 'Data Cluster' section is highlighted with a red box. Configuration options include:

- Policy:** Unrestricted (radio button selected)
- Access mode:** Single user access (radio button selected)
- Single user:** PULKIT KUMAR
- Performance:** Databricks runtime version (Runtime: 13.3 LTS (Scala 2.12, Spark 3.4.1)) and Use Photon Acceleration (checked)
- Node type:** Standard\_D4ds\_v5 (16 GB Memory, 4 Cores)
- Termination:** Terminate after 120 minutes of inactivity
- Tags:** Add tags (Kev, Value, Add button)

A summary panel on the right provides details:

Summary	
1 Driver	16 GB Memory, 4 Cores
Runtime	13.3.x-scala2.12
Unity Catalog Photon Standard_D4ds_v5 2 DBU/h	

5. It will take some time for the cluster to be in place.