# Azure Databricks

avanade

# Agenda

Azure Databricks – Introduction

Azure Databricks – Benefits

Azure Databricks -  Core Artifacts

Azure Databricks Delta

avanade

# Azure Databricks – Introduction

Apache Spark is an open source standard for advanced analytics, machine learning, AI and Big data

Databricks is extension of power of Spark with additional features

Azure Databricks Microsoft aims to help businesses work with data in a faster, easier and more collaborative way.

Using Databricks customers can accelerate innovation with one click setup, streamlined workflows and an interactive workspace that enables easy collaboration between data scientists, engineers and business analysts.
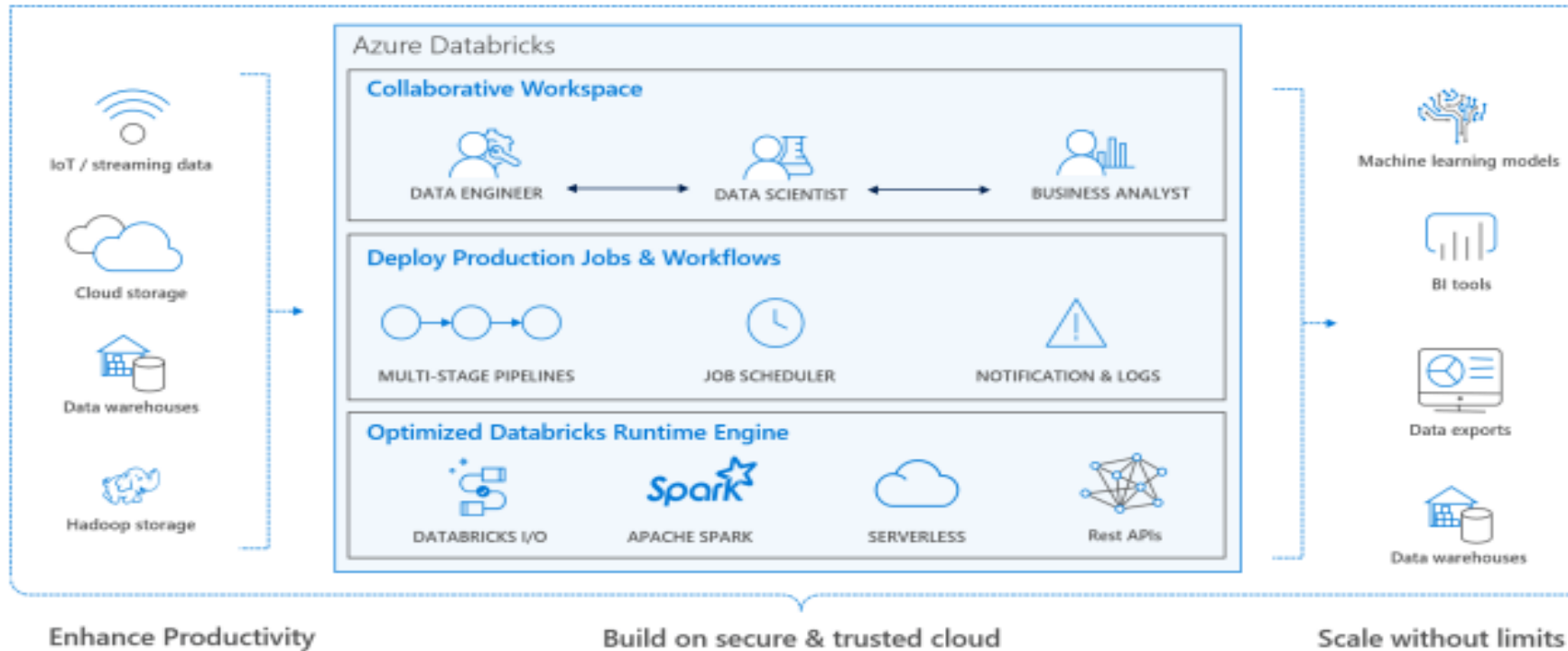
# Azure Databricks – Introduction

**Best of Databricks** + **Best of Microsoft**

Azure Databricks provides a unified platform to manage clusters for various use cases such as running production ETL pipelines, streaming analytics and ad-hoc analytics.

# Azure Databricks – Introduction

# Azure Databricks – Key benefits

Enhance productivity

❖Get started quickly by launching your spark environment with one click.

❖Integrate effortlessly with a wide variety of data stores and services such as Azure Synapse, Azure Cosmos DB, Azure Data Lake Store, Azure Blob storage, and Azure Event Hub. Add artificial intelligence (AI) capabilities instantly and share insights through rich integration with Power BI.

❖Use Databricks notebook to unify the process and instantly deploy into production. Improves collaboration amongst various teams through a unified workspace. Notebooks supports for multiple programming languages like R, Python, Scala and SQL.
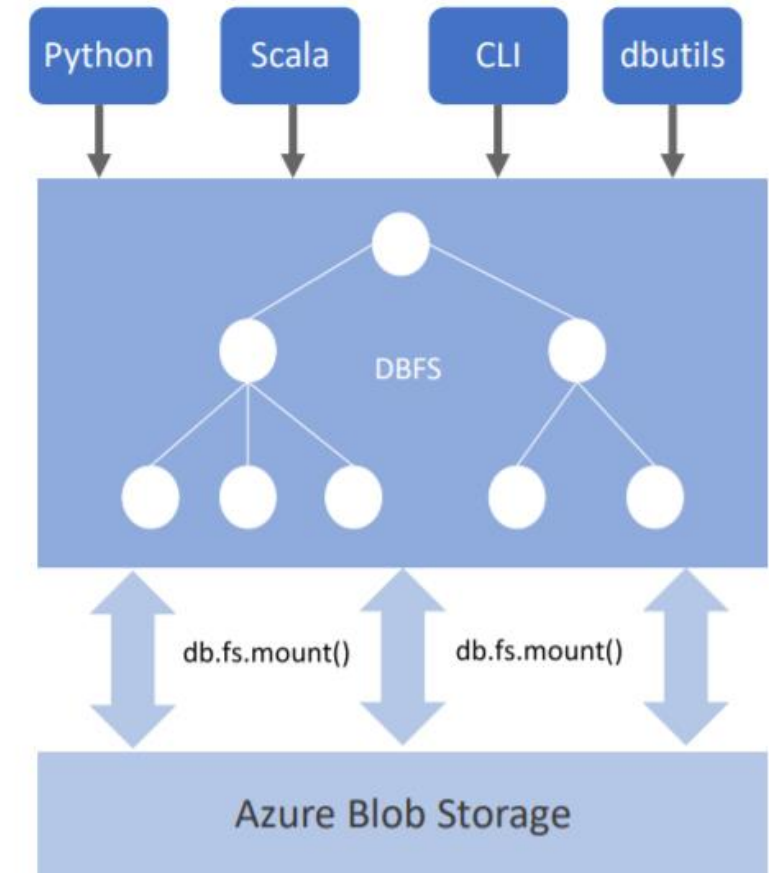
# Azure Databricks – Key benefits

## Scale without limits

❖ Operate at massive scale without limits globally. Databricks server less and highly elastic cloud service is designed to remove operational complexity.

❖ Minimize costs by Auto scaling and Auto termination of Spark Clusters.

❖ Accelerate data processing with the fastest Spark engine.

❖ High-speed connectors to Azure storage services such as Azure Blob Store and Azure Data Lake, developed together with the Microsoft teams behind these services.
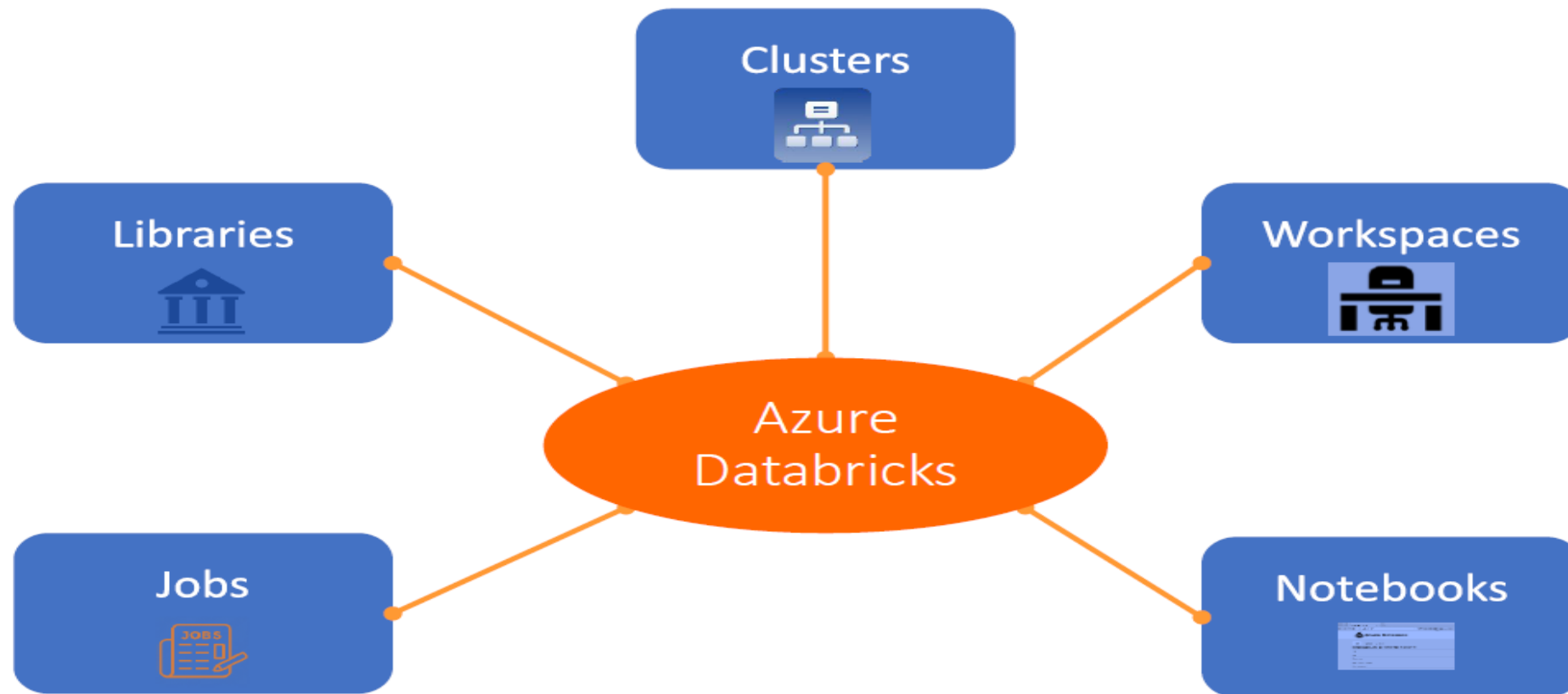
# DATABRICKS FILE SYSTEM (DBFS)

## Is a distributed File System (DBFS) that is a layer over Azure Blob Storage

- Azure Storage buckets can be mounted in DBFS so that users can directly access them without specifying the storage keys

- DBFS mounts are created using *dbutils.fs.mount()*

- Azure Storage data can be cached locally on the SSD of the worker nodes

- Available in both Python and Scala and accessible via a DBFS CLI

- Data persist in Azure Blob Storage – is not lost even after cluster termination

- Comes pre-installed on Spark clusters in Databricks

# Azure Databricks – Core Artifacts

# Azure Databricks – *Clusters*

- Azure Databricks clusters are the set of Azure Linux VMs that host the Spark Worker and Driver Nodes

- Your Spark application code (i.e. Jobs) runs on the provisioned clusters.

- Azure Databricks clusters are launched in your subscription—but are managed through the Azure Databricks portal.

- Azure Databricks provides a comprehensive set of graphical wizards to manage the complete lifecycle of clusters—from creation to termination.

- You can create two types of clusters – *Standard* and *Serverless Pool* (see next slide)

- While creating a cluster you can specify:
  - Number of nodes
  - **Autoscaling and Auto Termination policy**
  - Spark Configuration details
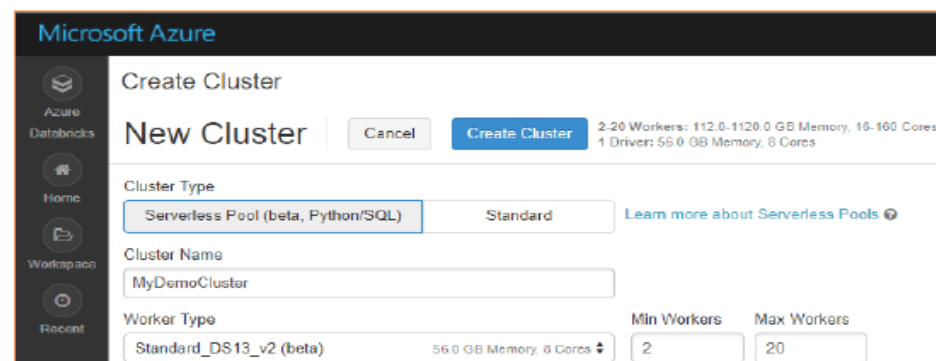  - The Azure VM instance types for the Driver and Worker Nodes

# Azure Databricks – Clusters

## SERVERLESS POOL (BETA)

### A self-managed pool of cloud resources, auto-configured for interactive Spark workloads

- You specify only the **minimum** and **maximum** number of nodes in the cluster—Azure Databricks provisions and adjusts the compute and local storage based on your usage.

- Limitation: Currently works only for SQL and Python.

**Microsoft Azure**

Create Cluster

**New Cluster** | Cancel | **Create Cluster** | 2-20 Workers: 112.0-1120.0 GB Memory, 16-160 Cores
1 Driver: 56.0 GB Memory, 8 Cores

Cluster Type

| Serverless Pool (beta, Python/SQL) | Standard | Learn more about Serverless Pools |

Cluster Name

MyDemoCluster

Worker Type

| Standard_DS13_v2 (beta) | 56.0 GB Memory, 8 Cores | Min Workers: 2 | Max Workers: 20 |

| **Benefits of Serverless Pool** | |
|---|---|
| Auto-Configuration | <ul><li>Databricks chooses the best configuration for Spark to get the best performance</li><li>Users don't need to worry about providing any of the Databricks runtime version or Spark configuration.</li><li>Databricks also chooses the best cluster parameters to save cost on infrastructure</li></ul> |
| Elasticity | <ul><li>Automatically scales the compute and local storage, independently, based on usage</li></ul> |
| Fine grained Sharing | <ul><li>Offers maximum resource utilization and minimum query latencies<ul><li>*Preemption:* Databricks proactively preempts Spark tasks from over-committed users to ensure all users get their fair share of cluster time and their jobs complete in a timely manner even when contending with dozens of other users. Uses the "Task Preemption for High Concurrency" feature of Spark in Databricks.</li><li>*Fault isolation:* Databricks sandboxes the environments belonging to different notebooks from one another.</li></ul></li></ul> |

# Azure Databricks – *Jobs*

## JOBS

Jobs are the mechanism to submit Spark application code for execution on the Databricks clusters

- Spark application code is submitted as a 'Job' for execution on Azure Databricks clusters

- Jobs execute either 'Notebooks' or 'Jars'

- Azure Databricks provide a comprehensive set of graphical tools to create, manage and monitor Jobs.

When you create a new Job you have to specify:

- The Notebook or Jar to execute

- Cluster: The cluster on which the Job execute. This could be an exiting or new cluster.

- Schedule i.e. how often the Job runs. Jobs can also be run one time right away.

When you create a new job you can optionally specify advanced options:

- Maximum number of concurrent runs of the Job

- Timeout: Jobs still running beyond the specified duration are automatically killed

- Retry Policy: Specifies if—and when—failed jobs will be retried

- Permissions: Who can do what with jobs. This allows for Job definition and management to be *securely shared* with others (see next slide)
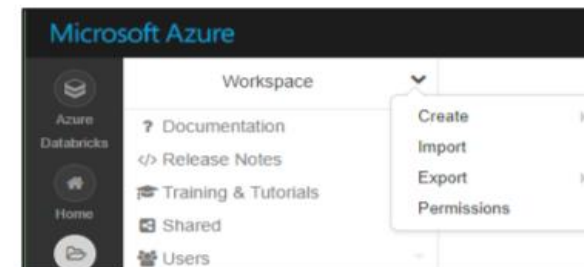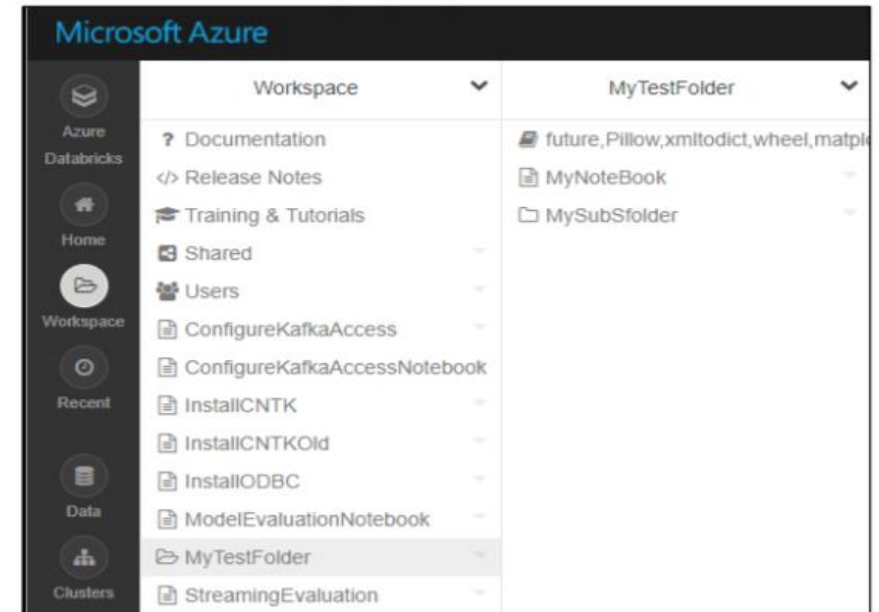
# WORKSPACES

## Workspaces enables users to organize—and share—their Notebooks, Libraries and Dashboards

- Workspaces—sort of like Directories— are a convenient way to organize an user's Notebook, Libraries and Dashboards.

- Everything in a workspace is organized into hierarchical folders. Folders can hold Libraries, Notebooks, Dashboard or more (sub) folders.
  - Icons indicate the type of the object contained in a folder

- Every user has one directory that is private and unshared.
  - By default, the workspace and all its contents are available to users.

- Fine grained access control can be defined on workspaces (next slide) to enable *secure collaboration with colleagues.*
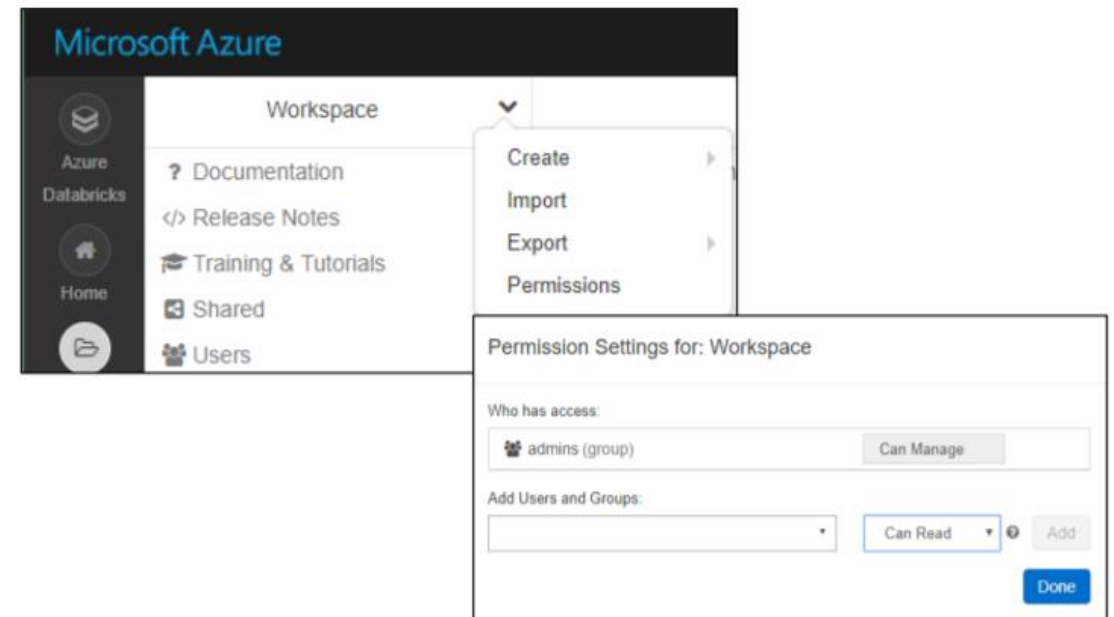
# WORKSPACE OPERATIONS

You can search the entire Databricks workspace

In the Azure Databricks Portal, via the Workspaces drop down menu, you can:

- Create Folders, Notebooks and Libraries

- Import Notebooks into the Workspace

- Export the Workspace to a database archive

- Set Permissions. You can grant 4 levels of permissions
  - Can Manage
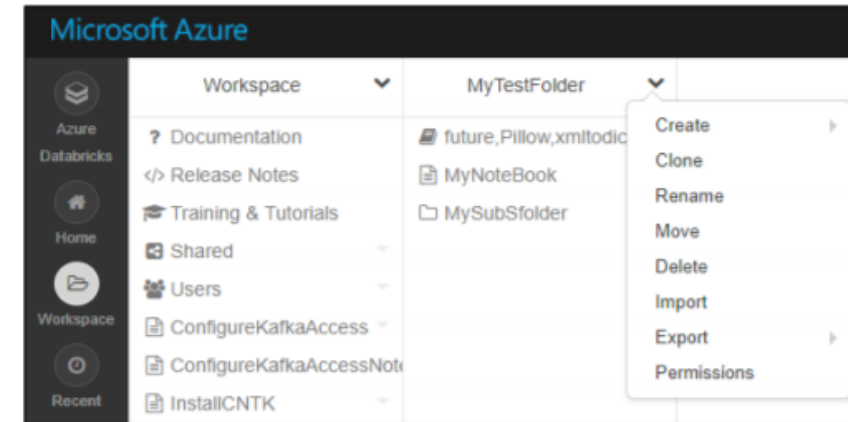  - Can Read
  - Can Edit
  - Can Run

# FOLDER OPERATIONS AND ACCESS CONTROL

In the Azure Databricks Portal, via the Folder drop down menu, you can:

- Create Folders, Notebooks and Libraries within the folder

- Clone the folder to create a deep copy of the folder

- Rename or delete the folder

- Move the folder to another location

- Export a folder to save it and its contents as a Databricks archive

- Import a saved Databricks archive into the selected folder

- Set Permissions for the folder. As with Workspaces you can set 5 levels of permissions: *No Permissions, Can Manage, Can Read, Can Edit, Can Run*



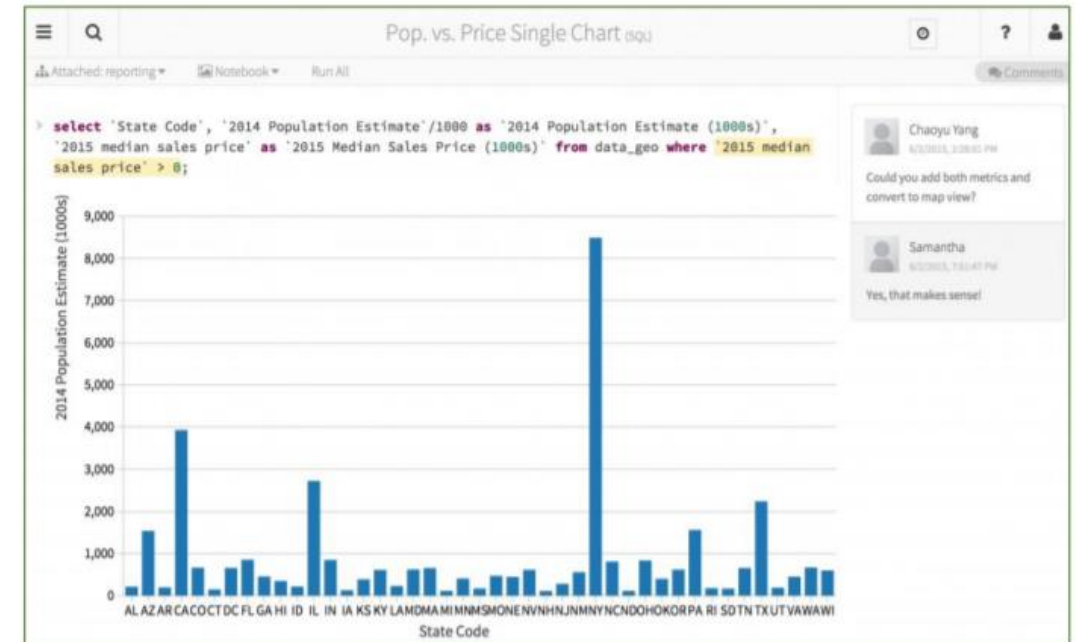| Abilities | No Permissions | Read | Run | Edit | Manage |
|---|---|---|---|---|---|
| Create items | | | | | ✅ |
| Delete items | | | | | ✅ |
| Move/rename items | | | | | ✅ |
| Change permissions | | | | | ✅ |

Abilities associated with each permission level

# AZURE DATABRICKS NOTEBOOKS OVERVIEW

## Notebooks are a popular way to develop, and run, Spark Applications

- Notebooks are not only for authoring Spark applications but can be *run/executed directly* on clusters

  - Shift+Enter

  - click the ▶ at the top right of the cell in a notebook

  - Submit via Job

- Notebooks support fine grained permissions—so they can be *securely shared* with colleagues for collaboration (see following slide for details on permissions and abilities)

- Notebooks are well-suited for prototyping, rapid development, exploration, discovery and iterative development



Notebooks typically consist of code, data, visualization, comments and notes

# MIXING LANGUAGES IN NOTEBOOKS

## You can mix multiple languages in the same notebook

Normally a notebook is associated with a specific language. However, with Azure Databricks notebooks, you can mix multiple languages in the same notebook. This is done using the language magic command:

- **%python**   Allows you to execute python code in a notebook (even if that notebook is not python)
- **%sql**   Allows you to execute sql code in a notebook (even if that notebook is not sql).
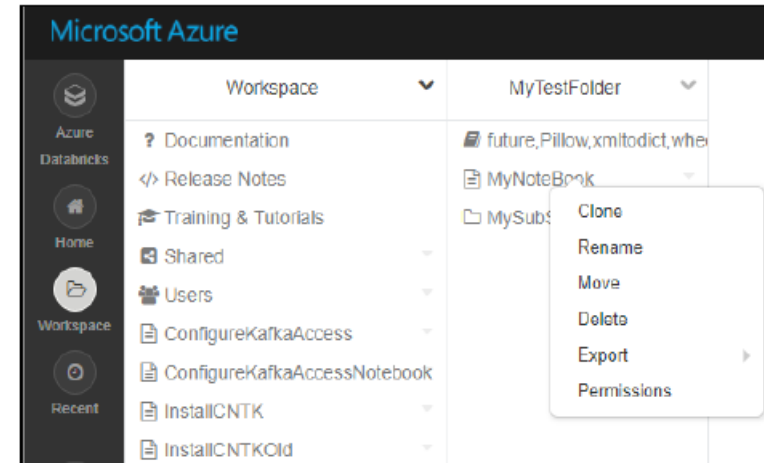- **%r**   Allows you to execute r code in a notebook (even if that notebook is not r).
- **%scala**   Allows you to execute scala code in a notebook (even if that notebook is not scala).
- **%sh**   Allows you to execute shell code in your notebook.
- **%fs**   Allows you to use Databricks Utilities - dbutils filesystem commands.
- **%md**   To include rendered markdown

# NOTEBOOK OPERATIONS AND ACCESS CONTROL

You can create a new notebook from the Workspace or the folder drop down menu (see previous slides)

From a notebook's drop down menu you can:

- Clone the notebook

- Rename or delete the notebook

- Move the notebook to another location

- Export a notebook to save it and its contents as a Databricks archive or IPython notebook or HTML or source code file.

- Set Permissions for the notebook As with Workspaces you can set 5 levels of permissions: *No Permissions, Can Manage, Can Read, Can Edit, Can Run*

- You can also set permissions from notebook UI itself by selecting the 🔒 Permissions menu option.

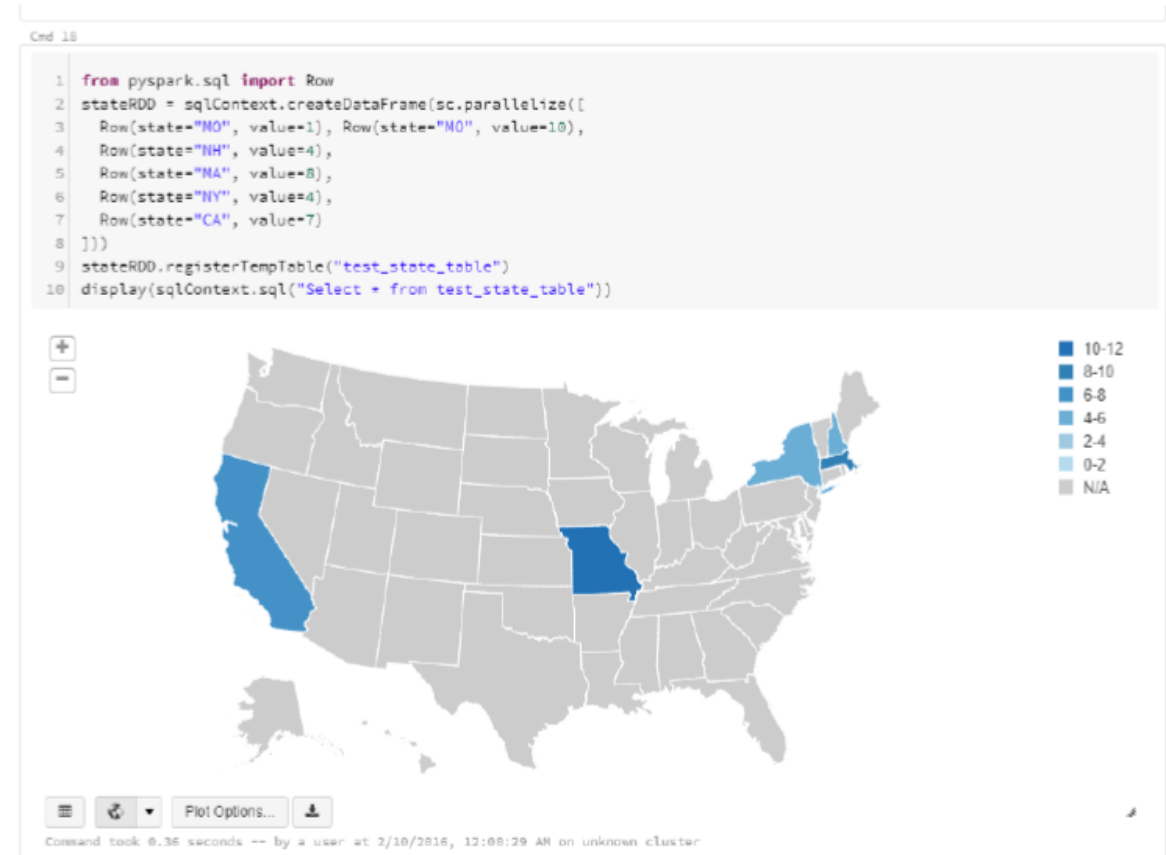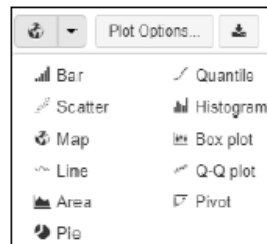| Abilities | No Permissions | Read | Run | Edit | Manage |
|---|---|---|---|---|---|
| View cells | | ✅ | ✅ | ✅ | ✅ |
| Comment | | ✅ | ✅ | ✅ | ✅ |
| Run Commands | | | ✅ | ✅ | ✅ |
| Attach/detach notebooks | | | ✅ | ✅ | ✅ |
| Edit cells | | | | ✅ | ✅ |
| Change permissions | | | | | ✅ |

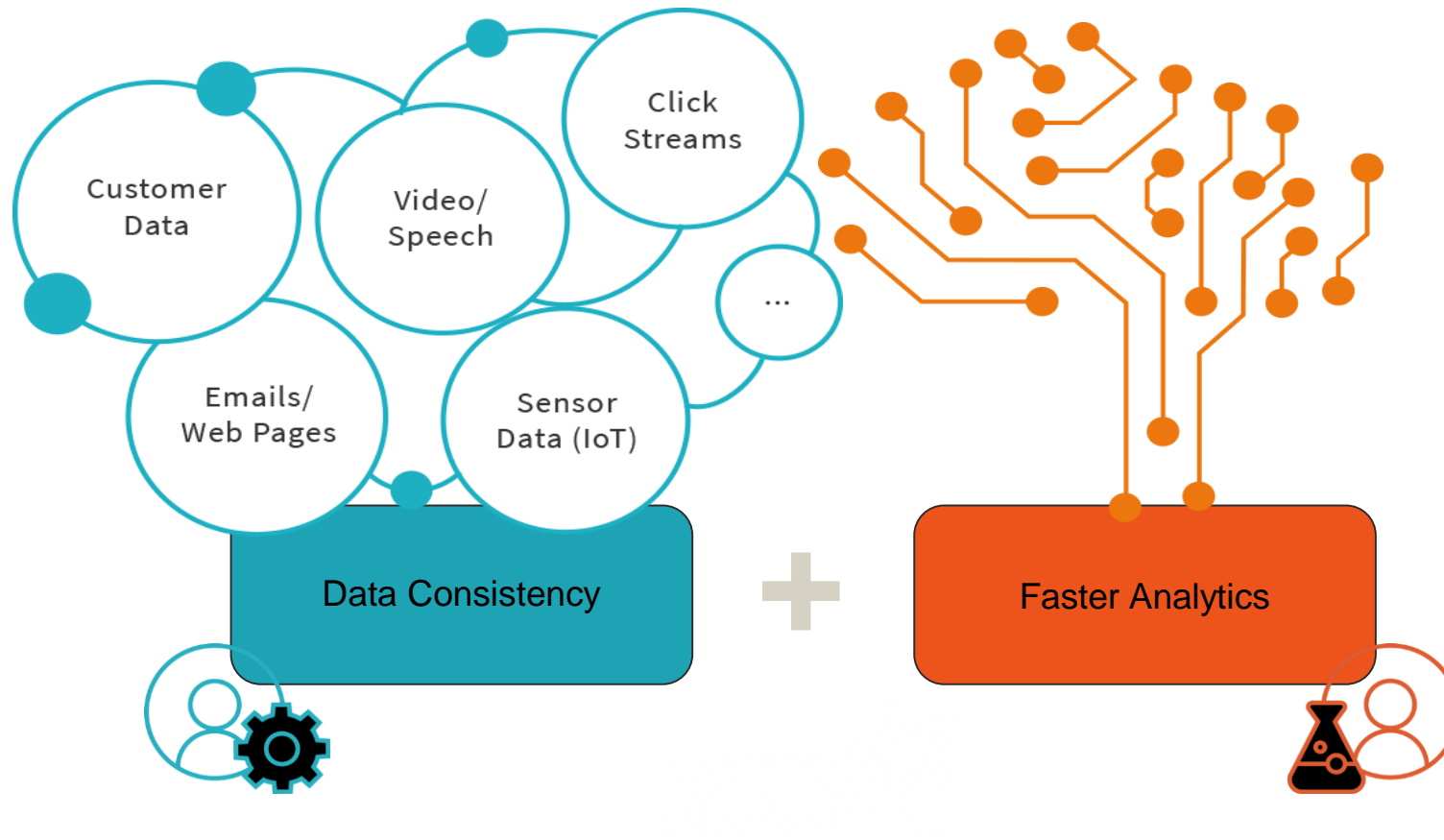Abilities associated with each permission level

# VISUALIZATION

## Azure Databricks supports a number of visualization plots out of the box

- All notebooks, *regardless of their language,* support Databricks visualizations.

- When you run the notebook the visualizations are rendered inside the notebook in-place

- The visualizations are written in HTML.

  - You can save the HTML of the entire notebook by exporting to HTML.

  - If you use Matplotlib, the plots are rendered as images so you can just right click and download the image

- You can change the plot type just by picking from the selection

# Azure Databricks Delta



Data Consistency   +   Faster Analytics

# Azure Databricks - Delta

Databricks Delta delivers a powerful transactional storage layer by harnessing the power of Apache Spark and Databricks DBFS (Data bricks file system). The core abstraction of Databricks Delta is an optimized Spark table that

- ➢ Stores data as Parquet files in DBFS.
- ➢ Maintains a *transaction log* that efficiently tracks changes to the table.
- ➢ Improves performance by organizing data and creating metadata.
    - ✓ Manages optimal file sizes.
    - ✓ Stores statistics that enable data skipping.
    - ✓ Creates and maintains indexes.
- ➢ Improves data reliability by adding data management capabilities.
    - ✓ Schema enforcement ensures data is not corrupted.
    - ✓ Upserts makes fixing bad data simple.
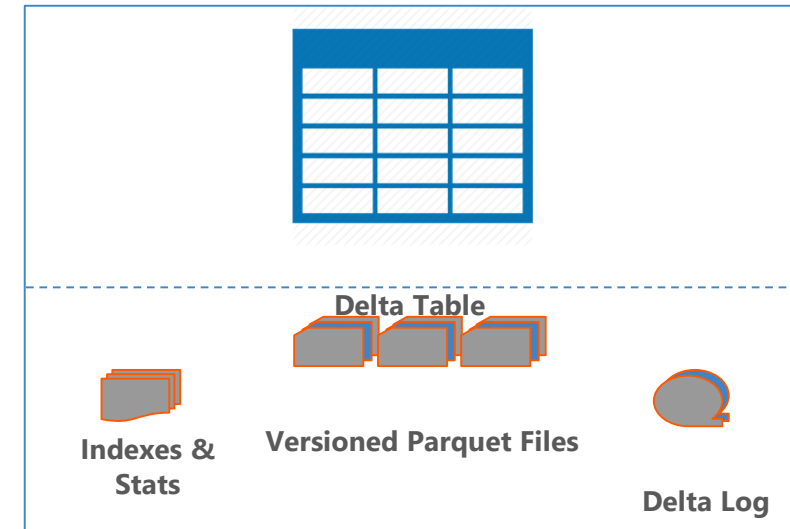
avanade

# Azure Databricks Delta Architecture

## Delta Table = Parquet + Transaction Log

Linear history of atomic changes

Optimistic Concurrency Control

Log checkpoint is stored as Parquet

**Delta Table**



Indexes & Stats

Versioned Parquet Files

Delta Table

Delta Log

# Azure Databricks Vs Spark

| | | Spark | databricks |
|---|---|:---:|:---:|
| **DATABRICKS RUNTIME**<br>*Built on Apache Spark and optimized for performance* | | ✕ | ✓ |
| **MANAGED DELTA LAKE**<br>*Reliable and Performant Data Lakes* | | ✕ | ✓ |
| **INTEGRATED WORKSPACE**<br>*Interactive Data Science and Collaboration* | | ✕ | ✓ |
| **PRODUCTION JOBS AND WORKFLOWS**<br>*Data Pipelines and Workflow Automation* | | ✕ | ✓ |
| **ENTERPRISE SECURITY**<br>*End-to-End Data Security and Compliance* | Learn More | ✕ | ✓ |
| **INTEGRATIONS**<br>*Compatible with Common Tools in the Ecosystem* | Rectangular Snip | ✕ | ✓ |
| **EXPERT SUPPORT**<br>*Unparalled Support by the Leading Committers of Apache Spark* | | ✕ | ✓ |

avanade

# Resources

[Azure Databricks Delta Guide](#)

[QuickStart notebook](#)

[Performance optimization](#)