

Azure Databricks Assignment 6

Pratik Wani

- Notes

- to store data without performing analysis on data set hierarchical namespace disabled.
- To set it to DataLake Gen2 enabled hierarchical namespace option.

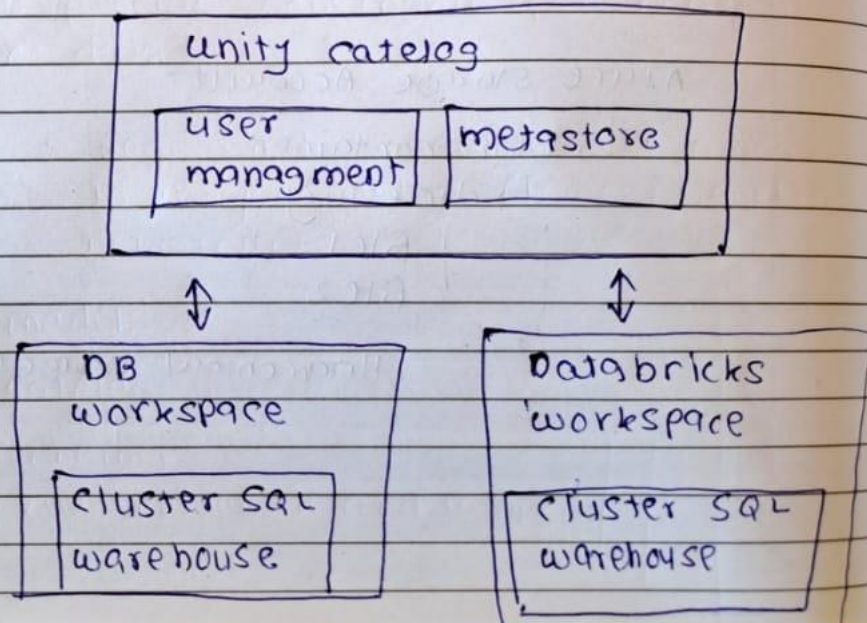
* stage of processing big data

Azure databricks

Day 6

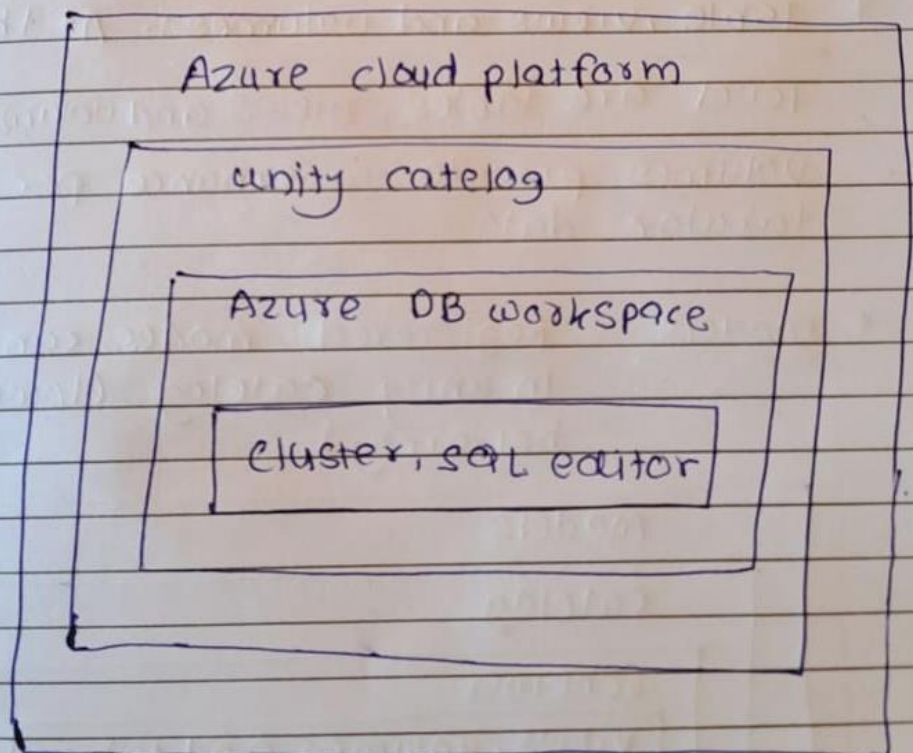
* unity catalog :-

- Provides centralised access control, auditing, linkage and data discovery capabilities



* key features:-

- 1) standards :- complaint security model based on ANSI SQL allows administration to grant permissions in their existing data lake
- 2) Define once , Secure everywhere
- 3) Built in Auditing and lineage - user level logs
- 4) data discovery
- 5) system tables



* unity catalog object model:-

* metastore :- top model container for metastore

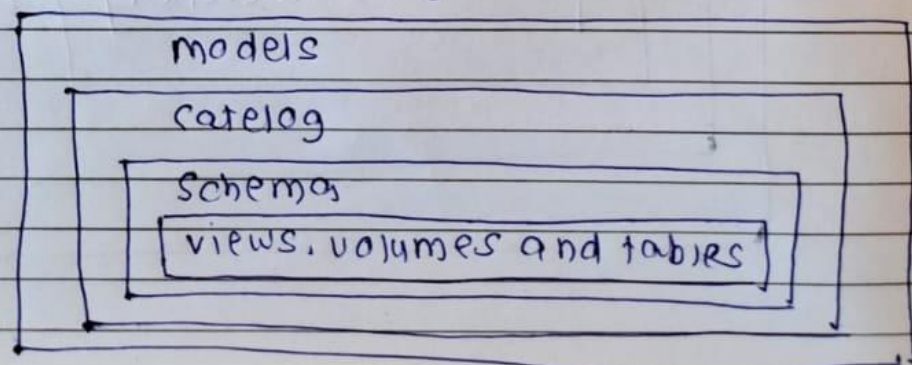
Each metastore exposes a three-level namespace (catalog, schema, table)

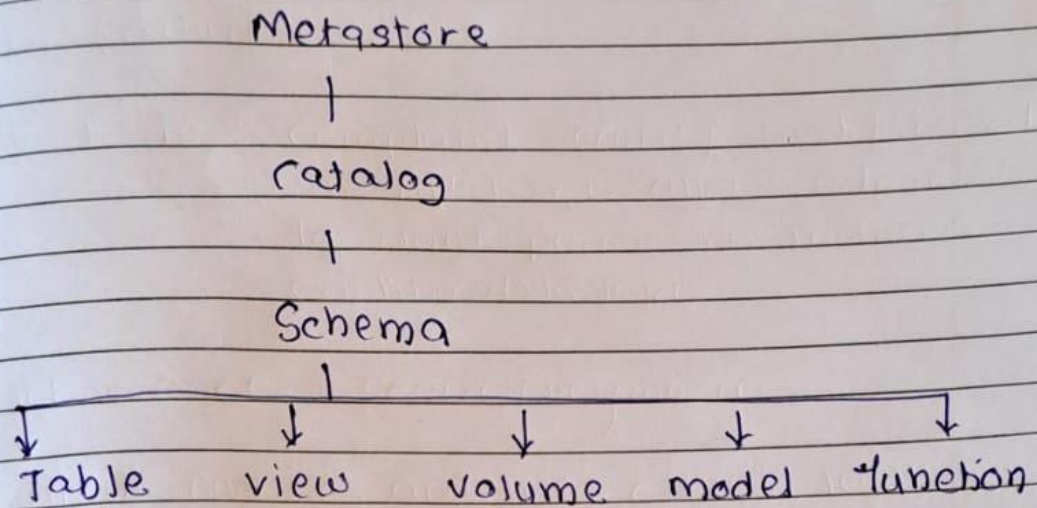
* catalog :- first layer of object hierarchy used to organize data access

* schema :- Databases / schemas are second layer of object hierarchy ; contains table and views

* Table, views and volumes :- At the lowest level are tables, views and volumes
volumes provide governance for non-tabular data

* models : Registered models can be managed in unity catalog (lowest level in hierarchy)





- * Azure databricks account admins
 - create one metastore for each region
 - in which they operate and assign it to DB workspace

Also a metastore can be configured in ADLS Gen2 containers

Externally managed data has file formats

DELTA	CSV	AVRO	PARQUET
ORC	TEXT		

- Unity Catalog

Unity Catalog 17:20:17 Python ☆

File Edit View Run Help Last edit was 1 minute ago New cell UI: ON

Run all Pratik Wani's Cluster Schedule Share

Cell 1

```
from pyspark.sql import SparkSession
```

Cell 2

```
spark=SparkSession.Builder().appName('Unity').getOrCreate()
```

Cell 3

```
spark.sql("CREATE CATALOG IF NOT EXISTS quickstart_catalog")
```

DataFrame[]

Cell 4

```
spark.sql("USE CATALOG quickstart_catalog")
```

DataFrame[]

Unity Catalog 17:20:17 Python ☆

File Edit View Run Help Last edit was 2 minutes ago New cell UI: ON

Run all Pratik Wani's Cluster Schedule

Cell 5

```
display(spark.sql("SHOW CATALOGS"))
```

Table + New result table: OFF

	catalog
1	hive_metastore
2	main
3	quickstart_catalog
4	samples
5	system

5 rows | 0.45 seconds runtime Refreshed 12 minutes ago

Cell 6

```
spark.sql("""
GRANT CREATE, USE CATALOG
ON CATALOG quickstart_catalog
TO `account users`""")
```

DataFrame[]



05:31 PM (1s) Cell 7

```
display(spark.sql("SHOW GRANT ON CATALOG quickstart_catalog"))
```

Table +

New result table: OFF

	Principal	ActionType	ObjectType	ObjectKey
1	account users	CREATE SCHEMA	CATALOG	quickstart_catalog
2	account users	USE CATALOG	CATALOG	quickstart_catalog

2 rows | 0.61 seconds runtime

Refreshed 7 minutes ago

05:31 PM (1s) Cell 8

```
spark.sql("""
CREATE SCHEMA IF NOT EXISTS quickstart_schema
COMMENT 'A new Unity Catalog schema called quickstart_schema'""")
```

DataFrame[]



05:31 PM (1s) Cell 9

```
display(spark.sql("SHOW SCHEMAS"))
```

Table +

New result table: OFF

	databaseName
1	default
2	information_schema
3	quickstart_schema

3 rows | 0.59 seconds runtime

Refreshed 7 minutes ago

05:31 PM (<1s) Cell 10

```
display(spark.sql("DESCRIBE SCHEMA EXTENDED quickstart_schema"))
```

Table +

New result table: OFF

	database_description_item	database_description_value
1	Catalog Name	quickstart_catalog
2	Namespace Name	quickstart_schema
3	Comment	A new Unity Catalog schema called quickstart_schema
4	Location	
5	Owner	pratikwani116@outlook.com
6	Properties	

Unity Catalog 17:20:17Python▼☆

FileEditViewRunHelpLast edit was 3 minutes agoNew cell UI: ON▼

▶ Run allPratik Wani's Cluster▼Schedule

📁🔍👤

▶ 05:32 PM (1s)Cell 11

spark.sql("""
GRANT CREATE TABLE, USE SCHEMA
ON SCHEMA quickstart_schema
TO 'account users' """)

DataFrame[]

▶ 05:32 PM (<1s)Cell 12

spark.sql("USE quickstart_schema")

DataFrame[]

▶ 05:32 PM (<1s)Cell 13

spark.catalog.currentDatabase()

'quickstart_schema'

▶ 05:32 PM (18s)Cell 14

spark.sql("CREATE OR REPLACE TABLE quickstart_table (id STRING)")

▶ (3) Spark Jobs
DataFrame[]

Unity Catalog 17:20:17Python▼☆

FileEditViewRunHelpLast edit was 3 minutes agoNew cell UI: ON▼

▶ Run allPratik Wani's Cluster▼Schedule

📁🔍👤

▶ 05:33 PM (1s)Cell 15

spark.sql("""
GRANT SELECT, MODIFY
ON TABLE quickstart_table
TO 'account users' """)

DataFrame[]

▶ 05:33 PM (<1s)Cell 16

display(spark.sql("SHOW TABLES"))

Table▼ +

New result table: OFF▼

	database	tableName	isTemporary
1	quickstart_schema	quickstart_table	false

1 row | 0.34 seconds runtime

Refreshed 6 minutes ago

▶ 05:33 PM (7s)Cell 17

spark.range(10).selectExpr("id").write.insertInto("quickstart_table")

▶ (6) Spark Jobs

Unity Catalog 17:20:17Python ☆

FileEditViewRunHelpLast edit was 4 minutes agoNew cell UI: ON

Run allPratik Wani's ClusterSchedule

05:33 PM (3s)Cell 18

display(spark.table("quickstart_table"))

(3) Spark Jobs

Table +New result table: OFF

	id
1	2
2	3
3	4
4	7
5	8
6	9
7	0

10 rows | 2.95 seconds runtime

Refreshed 6 minutes ago

4 minutes ago (<1s)Cell 19

display(spark.sql("SHOW TABLES in quickstart_schema"))

Table +New result table: OFF

	database	tableName	isTemporary
1	quickstart_schema	quickstart_table	false

Unity Catalog 17:20:17Python ☆

FileEditViewRunHelpLast edit was 4 minutes agoNew cell UI: ON

Run allPratik Wani's ClusterSchedule

4 minutes ago (1s)Cell 20

import pyspark.pandas as ps
psdf = ps.read_table("quickstart_schema.quickstart_table")

4 minutes ago (1s)Cell 21

display(psdf)

(2) Spark Jobs

Table +New result table: OFF

	id
1	2
2	3
3	4
4	7
5	8
6	9
7	0

10 rows | 0.51 seconds runtime

Refreshed 4 minutes ago

💡1

Unity Catalog 17:20:17

Python

☆

File Edit View Run Help

Last edit was 4 minutes ago

New cell UI: ON

▶ Run all

● Pratik Wani's Cluster

📅 Schedule

📄

📁

🔍

👤

▶ 4 minutes ago (5s)

Cell 22

Python

+

🔗

⋮

📄

```
psdf.to_table("quickstart_schema.quickstart_table_ps", overwriteSchema=True)
```

▶ (6) Spark Jobs

▶ 3 minutes ago (<1s)

Cell 23

spark.sql("""
GRANT SELECT, MODIFY
ON TABLE quickstart_table_ps
TO `account users`""")

DataFrame[]

▶ 3 minutes ago (<1s)

Cell 24

display(spark.sql("SHOW TABLES in quickstart_schema"))

Table

▼

+

New result table: OFF

	database	tableName	isTemporary
1	quickstart_schema	quickstart_table	false
2	quickstart_schema	quickstart_table_ps	false

2 rows | 0.30 seconds runtime

Refreshed 3 minutes ago

Unity Catalog 17:20:17

Python

☆

File Edit View Run Help

Last edit was 5 minutes ago

New cell UI: ON

▶ Run all

● Pratik Wani's Cluster

📅 Schedule

📄

📁

🔍

👤

▶ 3 minutes ago (1s)

Cell 25

Python

+

🔗

⋮

```
display(spark.table("quickstart_table_ps"))
```

▶ (3) Spark Jobs

Table

▼

+

New result table: OFF

	id
1	2
2	3
3	4
4	7
5	8
6	9
7	0

10 rows | 0.97 seconds runtime

Refreshed 3 minutes ago

- Job Scheduling

The screenshot shows the Unity Catalog notebook interface. The top bar displays 'Unity Catalog 17:20:17' and 'Python'. Below the top bar, there are tabs for 'File', 'Edit', 'View', 'Run', and 'Help'. The main area contains a code cell with the following code:

```
from pyspark.sql import SparkSession

spark=SparkSession.Builder().appName('Unity').getOrCreate()

spark.sql("CREATE CATALOG IF NOT EXISTS quickstart_catalog")

DataFrame[]
```

On the right side, there is a 'Job execution history' panel showing four jobs:

- Job - At 05:57 PM - Unity Catalog 172017
Last run: Feb 19 2024, 17:57 PM IST, 1m 3s
- Job - At 05:53 PM - Unity Catalog 172017
Last run: Feb 19 2024, 17:53 PM IST, 1m 31s
- Job - At 05:51 PM - Unity Catalog 172017
Last run: No runs
- Job - At 05:52 PM - Unity Catalog 172017
Last run: Feb 19 2024, 17:52 PM IST, 1m 32s

Each job has a 'Run now' button. At the bottom of the panel, there is an 'Add a schedule' button.

The screenshot shows the 'Jobs' page in the Unity Catalog interface. The top bar displays 'Unity Catalog 172017'. Below the top bar, there are tabs for 'Runs' and 'Tasks'. The main area shows the 'Unity_Catalog_172017' job details, including the task configuration and job execution history.

Task configuration:

- Task name*: Unity_Catalog_172017
- Type*: Notebook
- Source*: Workspace
- Path*: /Users/pratikwani116@outlook.com/Unity Catalog 17:20:17
- Cluster*: Unity_Catalog_172017_cluster 126 GB - 36 Cores - DBR 13.3 LTS - Photon - Spark 3.4...
- Dependent libraries: + Add

Job details:

- Job ID: 49073710216961
- Creator: Pratik Wani
- Run as: Pratik Wani
- Tags: Add tag
- Description: Add description
- Lineage: No lineage information for this job. Learn more

Schedules & Triggers:

- At 05:52 PM (UTC+05:30 — undefined)
- Edit trigger
- Pause
- Delete

Compute:

- Not configured
- Add Git settings

Buttons at the bottom: Cancel, Save task.

