**Name: Rohan Vinayak Chaudhari**
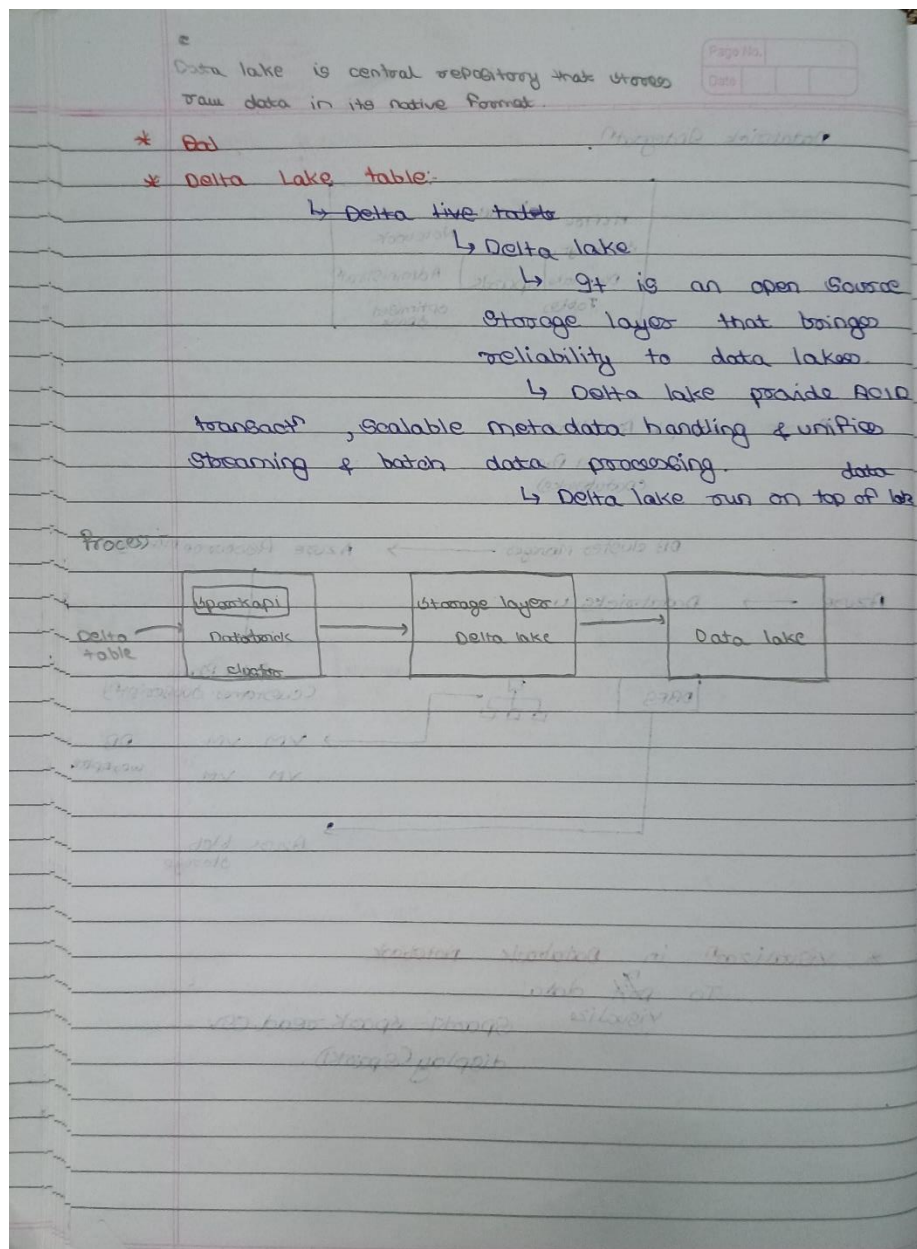
**Batch: Data Engineering**

**Date:14/02/2024**

**Topic: Azure DataBricks**

**Solution:**

# 1.Azure Databricks(Delta lakes):

Data lake is central repository that stores raw data in its native format.

* Ddl
* Delta Lake table:
  ↳ Delta live table
    ↳ Delta lake
      ↳ 'It' is an open source storage layer that brings reliability to data lakes.
        ↳ Delta lake provide ACID transaction, scalable meta data handling & unifies streaming & batch data processing.
          ↳ Delta lake run on top of data lake

Process

| Sparkapi | | Storage layer | | Data lake |
|---|---|---|---|---|
| Databricks cluster | → | Delta lake | → | |

Delta table

## Creating delta lake table with SQL:



## Viewing the data:

## Updating data with overwrite:



## Getting the Updated Data:

## Updating Data without Overwrite:



## Viewing the updated data:

## Deleting Even Data:



## Upsert&Merge operations:

## Viewing Upserted data:



## Getting older data:

# Delta lakes Working with Python:



```python
data = spark.range(0, 5)
data.write.format("delta").save("/tmp/delta-table-withpython")
```

▶ (6) Spark Jobs
▶ 🔲 data: pyspark.sql.dataframe.DataFrame = [id: long]

Command took 2.79 seconds -- by azuser1076_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:39:29 PM on azuser1076_mml.local's Cluster

**Cmd 2**

```python
df = spark.read.format("delta").load("/tmp/delta-table-withpython")
df.show()
```

▶ (3) Spark Jobs
▶ 🔲 df: pyspark.sql.dataframe.DataFrame = [id: long]

```
+---+
| id|
+---+
|  3|
|  4|
|  0|
|  1|
|  2|
+---+
```

Command took 0.75 seconds -- by azuser1076_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:40:08 PM on azuser1076_mml.local's Cluster
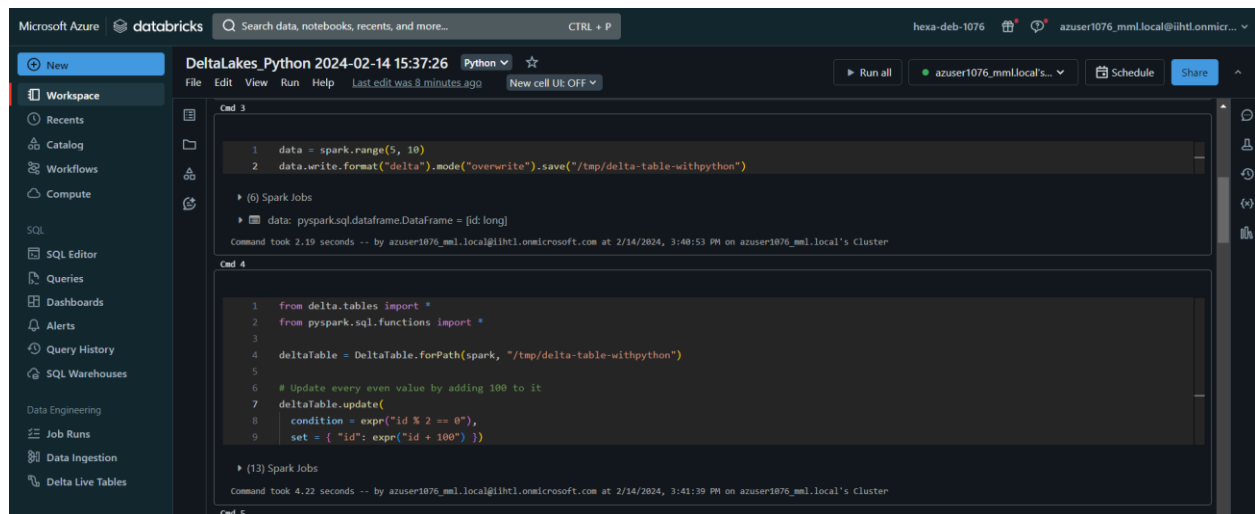
# Overwrite&without overwrite:



**Cmd 3**

```python
data = spark.range(5, 10)
data.write.format("delta").mode("overwrite").save("/tmp/delta-table-withpython")
```

▶ (6) Spark Jobs
▶ 🔲 data: pyspark.sql.dataframe.DataFrame = [id: long]

Command took 2.19 seconds -- by azuser1076_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:40:53 PM on azuser1076_mml.local's Cluster
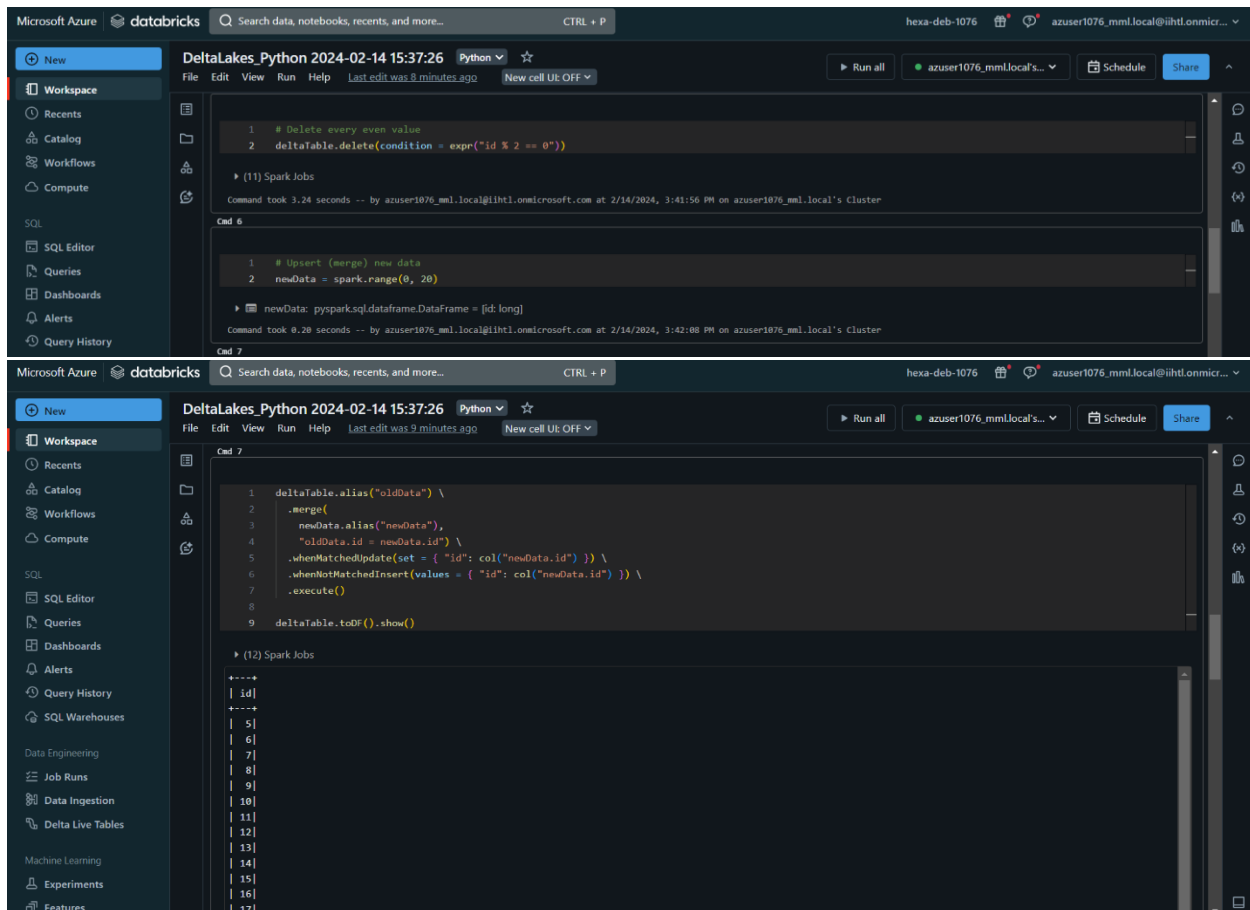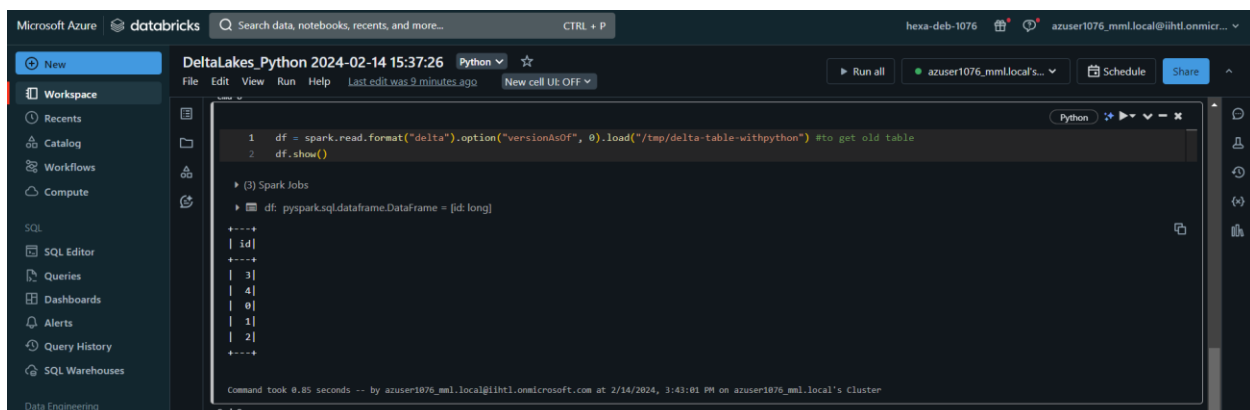
**Cmd 4**

```python
from delta.tables import *
from pyspark.sql.functions import *

deltaTable = DeltaTable.forPath(spark, "/tmp/delta-table-withpython")

# Update every even value by adding 100 to it
deltaTable.update(
  condition = expr("id % 2 == 0"),
  set = { "id": expr("id + 100") })
```

▶ (13) Spark Jobs

Command took 4.22 seconds -- by azuser1076_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:41:39 PM on azuser1076_mml.local's Cluster

**Cmd 5**

# Delete , Upsert&retriving new data in Python:



# Getting older data:

**Streaming Data:**