Custom Learnings
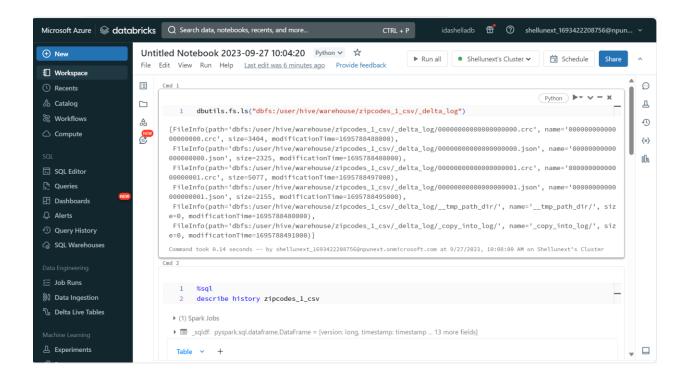
Day 20

**Databricks**

Structured Streaming: Modes:

1. Append: Optional when not using aggregation.
2. Complete: Mandatory in case of aggregation.
3. Update: On the go will check for the updating if required it will do.

This involves more cost than ADF.

Parquet is columnar format which goes well with big data. But it doesn't support deletion operation.

Delta Table: Whenever we create a table in databricks, delta table is created.

We can do time-travelling in delta table. We can either mimic the table for that time period or rollback to that particular version. Inside the delta.log folder, the json file will store these versions. Vacuum command is used to clean the versions in periodic intervals.

Microsoft Azure | databricks | Search data, notebooks, recents, and more... CTRL + P | idashelladb | shellunext_1693422208756@npun... 

New
Workspace
Recents
Catalog
Workflows
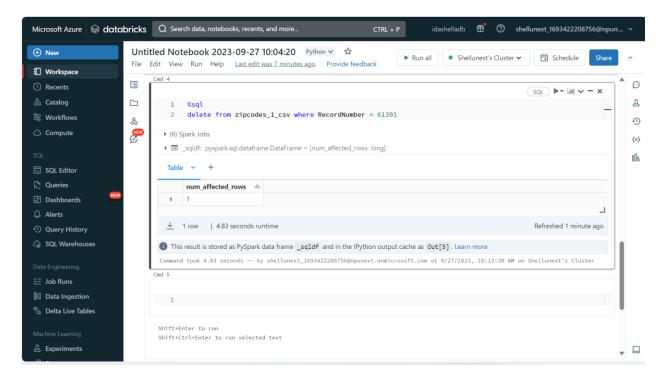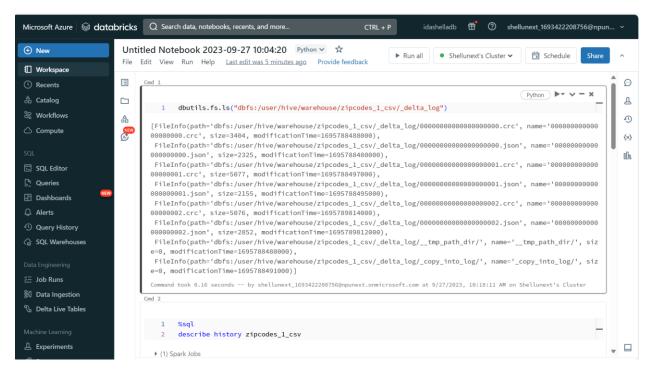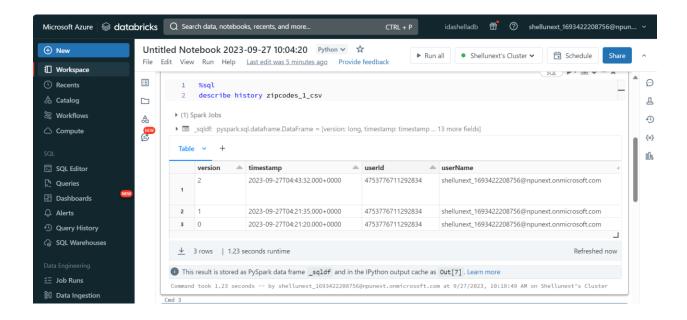Compute

SQL
SQL Editor
Queries
Dashboards
Alerts
Query History
SQL Warehouses

Data Engineering
Job Runs

Untitled Notebook 2023-09-27 10:04:20 | Python | Last edit was 7 minutes ago | Provide feedback | Run all | Shellunext's Cluster | Schedule | Share

Cmd 2

```sql
%sql
describe history zipcodes_1_csv
```

▶ (1) Spark Jobs
▶ _sqldf: pyspark.sql.dataframe.DataFrame = [version: long, timestamp: timestamp ... 13 more fields]

Table +

| | version | timestamp | userId | userName |
|---|---|---|---|---|
| 1 | 1 | 2023-09-27T04:21:35.000+0000 | 4753776711292834 | shellunext_1693422208756@npunext.onmicrosoft.com |
| 2 | 0 | 2023-09-27T04:21:20.000+0000 | 4753776711292834 | shellunext_1693422208756@npunext.onmicrosoft.com |

2 rows | 1.91 seconds runtime | Refreshed 5 minutes ago

ⓘ This result is stored as PySpark data frame _sqldf and in the IPython output cache as Out[3] . Learn more

Command took 1.91 seconds -- by shellunext_1693422208756@npunext.onmicrosoft.com at 9/27/2023, 10:09:18 AM on Shellunext's Cluster

Cmd 3

---



Microsoft Azure | databricks | Search data, notebooks, recents, and more... CTRL + P | idashelladb | shellunext_1693422208756@npun... 

New
Workspace
Recents
Catalog
Workflows
Compute

SQL
SQL Editor
Queries
Dashboards
Alerts
Query History
SQL Warehouses

Data Engineering
Job Runs
Data Ingestion
Delta Live Tables

Machine Learning
Experiments

Untitled Notebook 2023-09-27 10:04:20 | Python | Last edit was 7 minutes ago | Provide feedback | Run all | Shellunext's Cluster | Schedule | Share

Cmd 3

```sql
%sql
select * from zipcodes_1_csv
```

▶ (1) Spark Jobs
▶ _sqldf: pyspark.sql.dataframe.DataFrame = [RecordNumber: integer, Zipcode: integer ... 18 more fields]

Table +

| | RecordNumber | Zipcode | ZipCodeType | City | State | LocationType | Lat |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 704 | STANDARD | PARC PARQUE | PR | NOT ACCEPTABLE | 17.96 |
| 2 | 2 | 704 | STANDARD | PASEO COSTA DEL SUR | PR | NOT ACCEPTABLE | 17.96 |
| 3 | 10 | 709 | STANDARD | BDA SAN LUIS | PR | NOT ACCEPTABLE | 18.14 |
| 4 | 61391 | 76166 | UNIQUE | CINGULAR WIRELESS | TX | NOT ACCEPTABLE | 32.72 |
| 5 | 61392 | 76177 | STANDARD | FORT WORTH | TX | PRIMARY | 32.75 |
| 6 | 61393 | 76177 | STANDARD | FT WORTH | TX | ACCEPTABLE | 32.75 |
| 7 | | 704 | STANDARD | URB EUGENE RICE | PR | NOT ACCEPTABLE | 17.96 |

21 rows | 1.77 seconds runtime | Refreshed 3 minutes ago

ⓘ This result is stored as PySpark data frame _sqldf and in the IPython output cache as Out[4] . Learn more

Command took 1.77 seconds -- by shellunext_1693422208756@npunext.onmicrosoft.com at 9/27/2023, 10:11:49 AM on Shellunext's Cluster

Cmd 4

```
%sql
delete from zipcodes_1_csv where RecordNumber = 61391
```

(6) Spark Jobs

_sqldf: pyspark.sql.dataframe.DataFrame = [num_affected_rows: long]

| num_affected_rows |
|---|
| 1 | 1 |

1 row | 4.83 seconds runtime                    Refreshed 1 minute ago

This result is stored as PySpark data frame `_sqldf` and in the IPython output cache as `Out[5]`. Learn more

Command took 4.83 seconds -- by shellunext_1693422208756@npunext.onmicrosoft.com at 9/27/2023, 10:13:30 AM on Shellunext's Cluster

Cmd 5

Shift+Enter to run
Shift+Ctrl+Enter to run selected text

After deleting:



```
dbutils.fs.ls("dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log")
```

```
[FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/00000000000000000000.crc', name='00000000000
00000000.crc', size=3404, modificationTime=1695788488000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/00000000000000000000.json', name='00000000000
000000000.json', size=2325, modificationTime=1695788480000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/00000000000000000001.crc', name='00000000000
00000001.crc', size=5077, modificationTime=1695788497000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/00000000000000000001.json', name='00000000000
000000001.json', size=2155, modificationTime=1695788495000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/00000000000000000002.crc', name='00000000000
00000002.crc', size=5076, modificationTime=1695789814000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/00000000000000000002.json', name='00000000000
000000002.json', size=2852, modificationTime=1695789812000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/__tmp_path_dir/', name='__tmp_path_dir/', siz
e=0, modificationTime=1695788480000),
 FileInfo(path='dbfs:/user/hive/warehouse/zipcodes_1_csv/_delta_log/_copy_into_log/', name='_copy_into_log/', siz
e=0, modificationTime=1695788491000)]
```

Command took 0.16 seconds -- by shellunext_1693422208756@npunext.onmicrosoft.com at 9/27/2023, 10:18:11 AM on Shellunext's Cluster

Cmd 2

```
%sql
describe history zipcodes_1_csv
```

(1) Spark Jobs

**Dataframe Activity:**

We optimize the partitions so that each partitioned file consists of the same no. of records.

By Partition creating table:

Restore previous version:
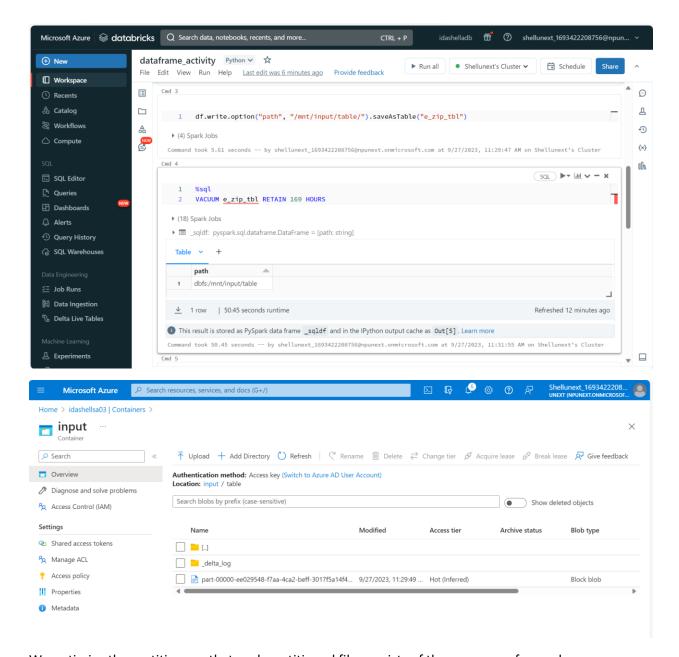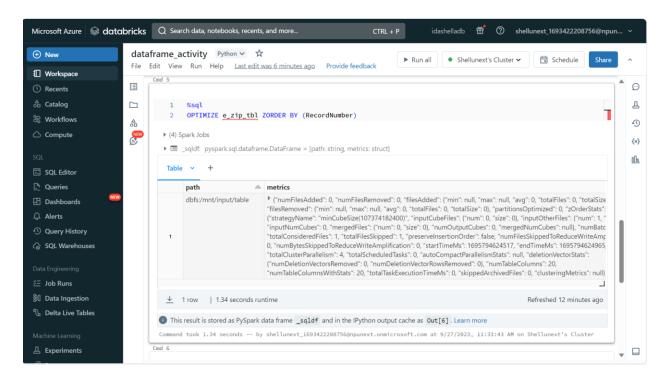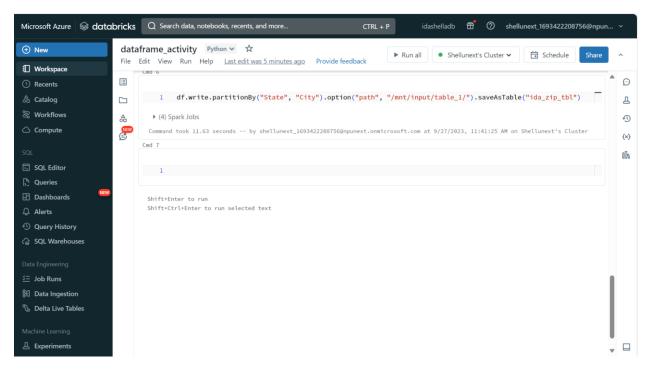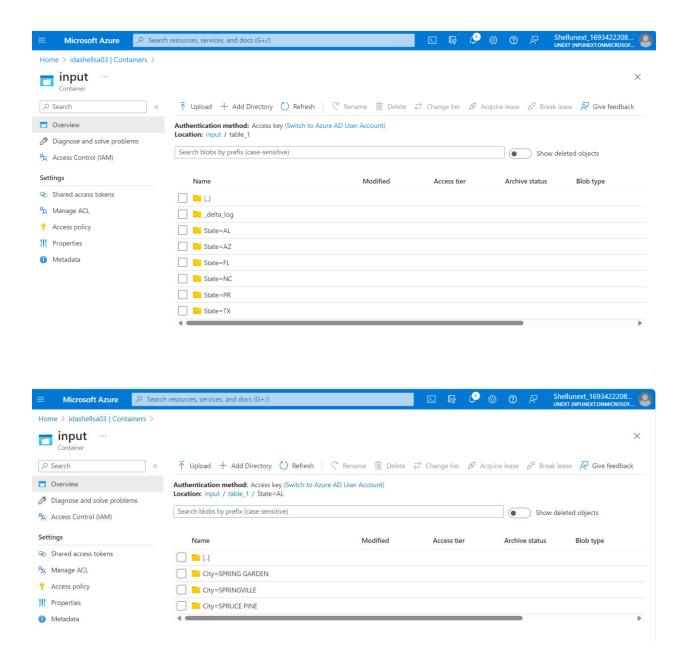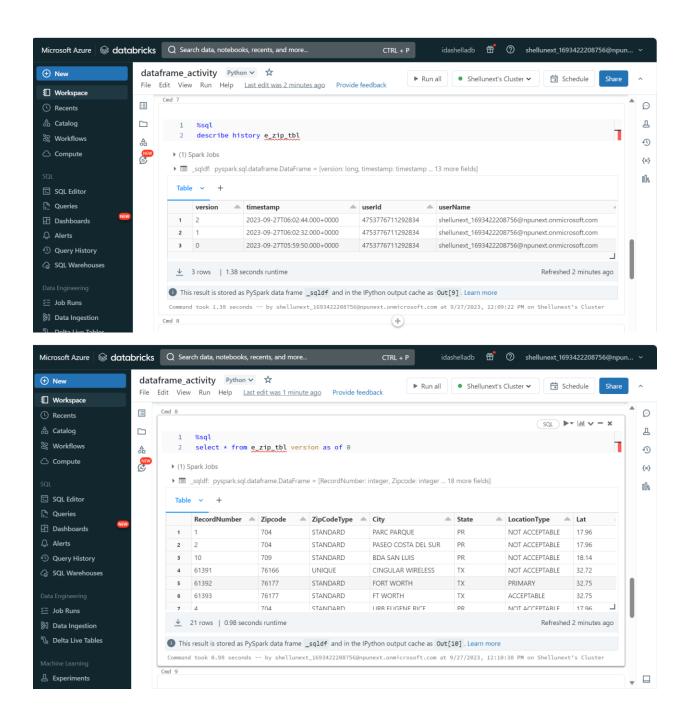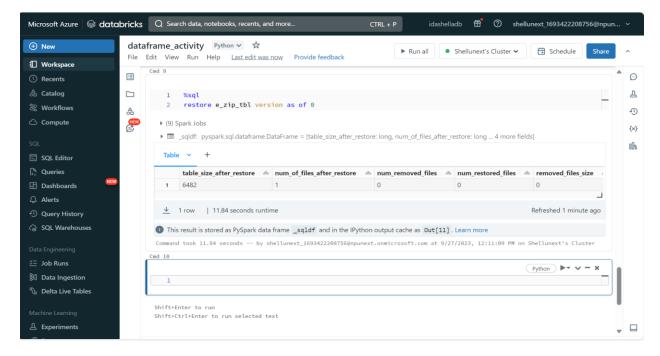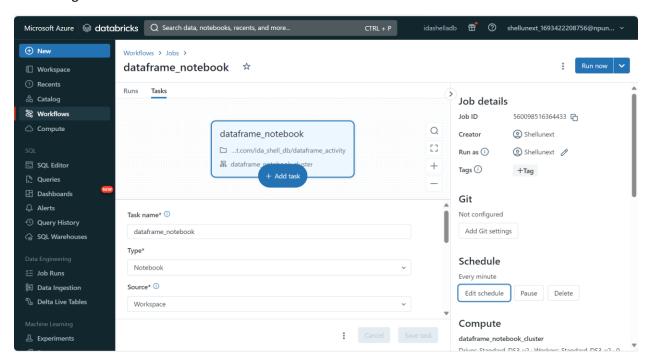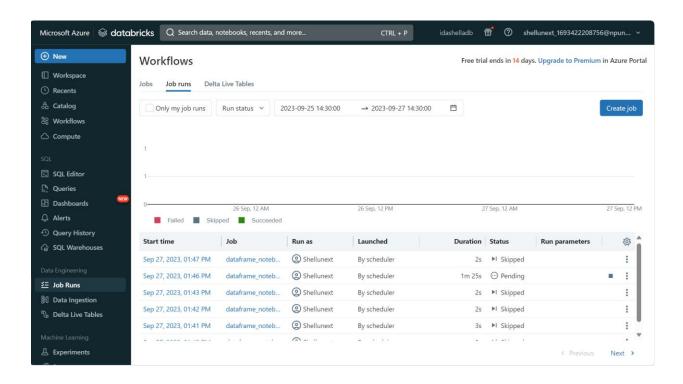
**dataframe_activity** Python ⌄ ☆

File Edit View Run Help · Last edit was 2 minutes ago · Provide feedback

▶ Run all · ● Shellunext's Cluster ⌄ · 📅 Schedule · Share ⌃

Cmd 7

```
1  %sql
2  describe history e_zip_tbl
```

▶ (1) Spark Jobs

▶ 🖻 _sqldf: pyspark.sql.dataframe.DataFrame = [version: long, timestamp: timestamp ... 13 more fields]

Table ⌄ +

| | version | timestamp | userId | userName |
|---|---|---|---|---|
| 1 | 2 | 2023-09-27T06:02:44.000+0000 | 4753776711292834 | shellunext_1693422208756@npunext.onmicrosoft.com |
| 2 | 1 | 2023-09-27T06:02:32.000+0000 | 4753776711292834 | shellunext_1693422208756@npunext.onmicrosoft.com |
| 3 | 0 | 2023-09-27T05:59:50.000+0000 | 4753776711292834 | shellunext_1693422208756@npunext.onmicrosoft.com |

⬇ 3 rows | 1.38 seconds runtime · Refreshed 2 minutes ago

ⓘ This result is stored as PySpark data frame `_sqldf` and in the IPython output cache as `Out[9]` . Learn more

Command took 1.38 seconds -- by shellunext_1693422208756@npunext.onmicrosoft.com at 9/27/2023, 12:09:22 PM on Shellunext's Cluster

Cmd 8

---

Cmd 8

SQL ▶ ⌄ 📊 ⌄ − ✕

```
1  %sql
2  select * from e_zip_tbl version as of 0
```

▶ (1) Spark Jobs

▶ 🖻 _sqldf: pyspark.sql.dataframe.DataFrame = [RecordNumber: integer, Zipcode: integer ... 18 more fields]

Table ⌄ +

| | RecordNumber | Zipcode | ZipCodeType | City | State | LocationType | Lat |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 704 | STANDARD | PARC PARQUE | PR | NOT ACCEPTABLE | 17.96 |
| 2 | 2 | 704 | STANDARD | PASEO COSTA DEL SUR | PR | NOT ACCEPTABLE | 17.96 |
| 3 | 10 | 709 | STANDARD | BDA SAN LUIS | PR | NOT ACCEPTABLE | 18.14 |
| 4 | 61391 | 76166 | UNIQUE | CINGULAR WIRELESS | TX | NOT ACCEPTABLE | 32.72 |
| 5 | 61392 | 76177 | STANDARD | FORT WORTH | TX | PRIMARY | 32.75 |
| 6 | 61393 | 76177 | STANDARD | FT WORTH | TX | ACCEPTABLE | 32.75 |
| 7 | 4 | 704 | STANDARD | URB EUGENE RICE | PR | NOT ACCEPTABLE | 17.96 |

⬇ 21 rows | 0.98 seconds runtime · Refreshed 2 minutes ago

ⓘ This result is stored as PySpark data frame `_sqldf` and in the IPython output cache as `Out[10]` . Learn more

Command took 0.98 seconds -- by shellunext_1693422208756@npunext.onmicrosoft.com at 9/27/2023, 12:10:38 PM on Shellunext's Cluster

Cmd 9

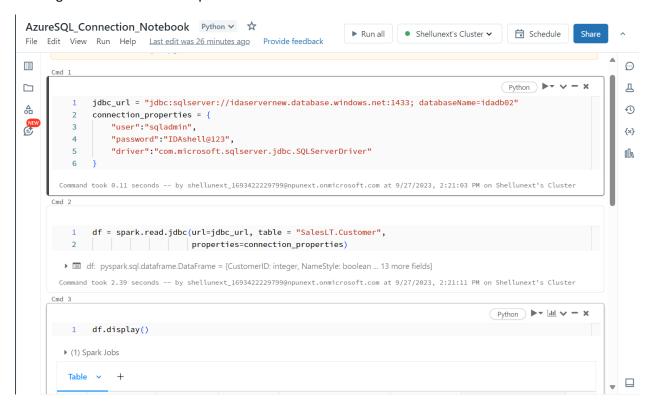## Medallion/ Multi Hop Architecture

We follow three folder architecture:

1. Raw/Bronze: All raw data from all the sources are kept in this folder.
2. Silver: Cleansing and formatting the data files into one particular file format.
3. Gold: Performing some aggregations and joins and keeping the data files needed according to the requirements.

Scheduling the task:

Reading the data from azure sql:



```python
jdbc_url = "jdbc:sqlserver://idaservernew.database.windows.net:1433; databaseName=idadb02"
connection_properties = {
    "user":"sqladmin",
    "password":"IDAshell@123",
    "driver":"com.microsoft.sqlserver.jdbc.SQLServerDriver"
}
```

Command took 0.11 seconds -- by shellunext_1693422229799@npunext.onmicrosoft.com at 9/27/2023, 2:21:03 PM on Shellunext's Cluster

Cmd 2

```python
df = spark.read.jdbc(url=jdbc_url, table = "SalesLT.Customer",
                properties=connection_properties)
```

▶ ▤ df: pyspark.sql.dataframe.DataFrame = [CustomerID: integer, NameStyle: boolean ... 13 more fields]

Command took 2.39 seconds -- by shellunext_1693422229799@npunext.onmicrosoft.com at 9/27/2023, 2:21:11 PM on Shellunext's Cluster

Cmd 3

```python
df.display()
```

▶ (1) Spark Jobs

Table ∨    +

- For the same workspace, there can be multiple users that work on different notebooks.
- For development, production, UAT env there are different workspaces.
- To share notebooks present under different users in different workspaces, the admin must provide access separately to every user for any file present in different workspaces.
- This is a very inefficient and time-consuming process.
- With unity catalog, workspace sharing can be enabled.
- This provides data access without needing admin permission.
- Azure Access Connector is used to provide connection to different workspace files.