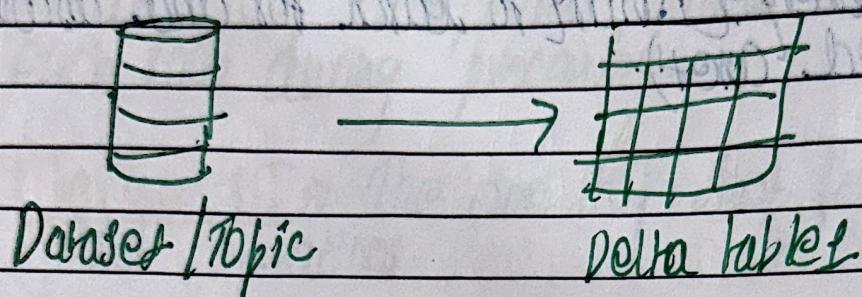


# Databricks Professional DataEngineer Certificate

## Bronze layer Ingestion patterns

- Singleplex → One To One mapping.
- Multiplex → Many to One mapping.

### Singleplex mapping

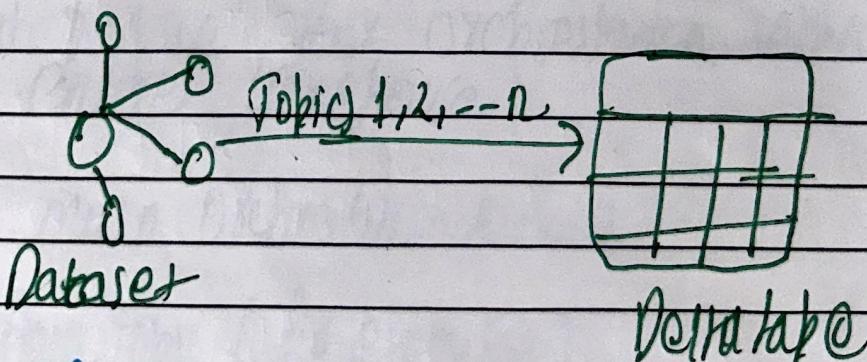


Each ~~per~~ dataset/topic is ingested to a bronze tab called as Singleplex mapping

Note → if we went practically there can be issue of maximum concurrent jobs in workspace

### Multiplex Ingestion Model

~~Topic~~



And after bronze layer there are Segregated in me silver layer.

## Slowly Changing dimension (scd)

Slowly Changing Dimensions (scd) is a data management concept that determines how tables handle data which change over time.

whether you want to overwrite the values or retain their history.

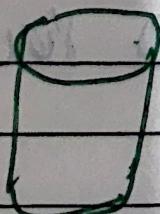
Scd types -

- 4 Type 0 → No changes allowed, these kind of tables are either static tables or Append Only tables.
- 4 Type 1 → Overwrite, No history is retained (use delta time travel if required).
- 4 Type 2 → Add new row of each change and marking old as obsolete, retaining the full history of values.

## Change data Capture (CDC)

It is a process of identifying changes made to data in the source and delivering those changes to the target.

Source



CDC

Target

(inserts, updates, deletes)

There can be row level changes →

- ① New records insertion
- ② Updating existing records
- ③ Deleting existing records.

Merge operation limitation -

- ① merge can not be performed if multiple source rows matched and attempted to modify the same target row in the Delta table
- ② CDC feed with multiple updates for the same key will generate an exception

### Delta lake CDF (Change Data Feed)

- ↳ It automatically generates CDC feeds about delta lake tables.
- ↳ Records row-level changes for all data writes into delta table.

Row data + metadata (whether row was inserted, deleted or updated).

- ↳ It is used in multi-hop architecture to propagate the changes.

## Example of CDF Enabled Table

Delta Table (V1) →

ID	Country	Vaccination Rate
FR	France	0.65
CA	Canada	0.75

After update



Delta Table (V2) →

Country ID	Country	Rate
FR	France	0.65
CA	Canada	0.75

Table changes

ID	Country	Rate	ChangeType	Time	Version
FR	France	0.7	Update_Pre	0:7:00	2
FR	??	0.75	Update_Post	07:00	2
CA	Canada	0.6	Update_Pre	0:9:0	2
CA	??	0.65	Update_Post	0:9:00	2

Preimage → Before update was done.

Post Image → After Update was done.

And if there is insert or delete it will record the delete and insert in same format.

Enabling CDF

4 New Table → Create Table mytable (id int, name string)  
TableProperties (delta-enablechangeflate)

## Existing table

ALTER TABLE my\_table  
SET TBLPROPERTIES ('delta.enable\_change\_data\_feed' = true)

- ↳ And / or enable CDF to all new table
- ↳ spark.databricks.delta.pro-properties.defaults.enabled Chang data feed.

## Use CDF When

- ↳ Table changes include updates and/or delete
- ↳ Small fraction of records updated in each batch (from CDC feed)

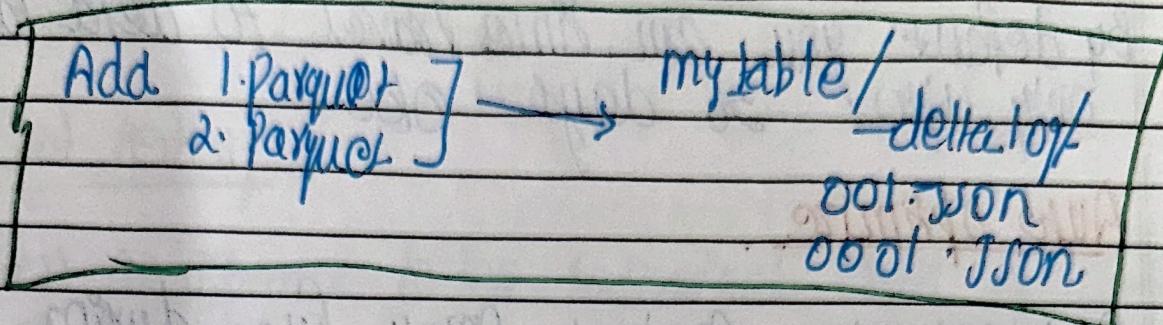
## Stream - static join

- ↳ Streaming tables are over-appending data sources
- ↳ static tables contain data that may be changed or overwritten - break the requirements of an over-appending source for structured streaming

Partitioning, A partition is a subset of rows that share the same value for pre-defined partitioning columns. It is like subfolders being created; it optimizes the query.

## Transaction log

Each Commit to the table is written out as a JSON file



## Transaction log Checkpoints.

- 4 Databricks automatically creates parquet checkpoint files every 10 commits to accelerate the resolution of current table state.
- 4 Then, spark has to perform incremental processing of newly added JSON files.
- 4 Delta lake captures statistics in the transaction log for each added data file.

## Delta lake file statistics.

- 4 These statistics indicate per file:-
  - ① Total number of records:-
  - ② Statistics on the first 32 columns of the table
  - ③ Min value in each col
  - ④ Max " " " "
  - ⑤ Null Value count for each column.

Note:- Running VACUUM doesn't delete the delta log files. And these log files are automatically cleaned up by databricks

By default you can time travel to delta table only upto 30 days old

### Auto optimize.

Automatically compact small files during individual writes to a table

### 4 2 Complementary features:

① Optimized writes: attempts to write 128MB files

② Auto compaction.

After the write completes, it checks if files can further be compacted.

If yes, it runs an OPTIMIZE job forward a file size of 128 MB (9376016B)

### REST API.

Databricks offers many api to create, run, cancel jobs

Data Pipeline testing → Data quality tests : test the quality of the data. It applies check constraints to delta tables

There are standard tests: test the code logic

- 4 Unit testing
- 4 Integration testing
- 4 End-to-End testing

### Unit testing

- 4 It is an approach to test individual units of code, such as functions.
- 4 If you make any change to them in future, you can determine whether they still work as you expect them to.
- 4 This helps you find problems with your code faster and earlier in the development life cycle.
- 4 An assertion is a statement that enables you to test the assumptions you have made in your code.

assert func() == expected\_value

Syntax of assert

### Integration testing

- 4 Approach to test the interaction b/w subystems of an application.
- 4 Software modules are integrated logically and tested as a group.

## End-to-End testing

- 4 Approach to ensure that your application can run properly under real-world scenarios
- 4 Simulate a user experience from start-to-finish

## Certification Overview

- 4 120 mins for exam
- 4 60 questions
- 4 Pass score 70% (42/60)

4 Exam fee = \$200

## Exam questions

- 4 Data processing (18/60)
- 4 Data modelling (12/60)
- 4 Databricks tooling (12/60)
- 4 Security and governance (6/60)
- 4 Testing and deployment (6/60)
- 4 Monitoring and logging (6/60)

## \* Out of scope topics

- 4 DLT (Delta Live Tables)
- 4 Scala
- 4 Orchestration tools (Airflow, NDFD)
- 4 Pub/Sub System Config (Kafka)

- 4 Automation Servers (Jenkin, Azure Devops, etc).
- 4 Managed CI/CD.
- 4 QHollow.
- 4 Cloud-specific Security and integration
- 4 Infrastructure as Code (terraform - -).

## Code Examples

- 4 Code will be in mainly python
- 4 Delta Lake functionality in SQL.

## Question type

- 4 MCQ only questions
- 4 Question types:-

- ① Conceptual Questions
- ② Code-Based Question

## Short - Crash Notes for Databricks Professionals

- ① The minimum permission to view metrics and spark UI of cluster is "Can Attach IO" privilege on the cluster.
- ② For production Databricks job "Job Clusters" are recommended to use.
- ③ Job clusters provide isolated environment for each job (URL).
- ④ "Can Restart" Permission is minimal permission to start and terminate the job.
- ⑤ Whenever an External table is dropped only mandatory table is deleted but the data files are kept on same location.
- ⑥ "Can Manage" is minimum permission to edit cluster.
- ⑦ If you want to print secret keys "REDACTED" is printed.
- ⑧ %sh magic command executes shell code on the local driver machine which leads to significant performance overhead.
- ⑨ file statistics in the delta transaction log is used to leverage the query optimizer.

- (1) Shallow cloned tables doesn't copy the data but refers the Delta transaction logs. If you run VACUUM command on source table you will get error of data not present in the shallow cloned table.
- (2) Syntax for DEP Clone → Create or replace orders\_archive deep clone order.
- (3) Select \* from table@v36 → Table Version 36.
- (4) If data is continuously appended from bronze to silver then we can use structured streaming in batch mode using trigger availableNow option.
- (5) Optimized writer and auto compaction automatically generate smaller data file to reduce the duration of future MERGE operation.
- (6) DBX maintains history of your job run up to 60 days. If you need to preserve job runs.