

① Describe DATABASE EXTENDED db-name ,  
this will show db-name, location where  
it is stored and its properties.



② Describe detail Table\_Name = also prints  
table details

③ Create external table :-

Create table table\_name  
using CSV options (

path 'path of .csv file'

headers "true"

mode "FAILOFAST"

);

④ 3 types of view:- 1) Regular 2) Temporary  
3) Global Temporary

For temporary view specify TEMPORARY  
keyword ~~or after create~~.

e.g Create TEMPORARY View ViewName as  
Select \* from tablename

For global temporary view:- specify  
GLOBAL TEMPORARY keyword

⑤ Temporary views are not accessible in  
other sessions but global temp views  
are accessible in other sessions

But when we restart the cluster  
we won't be able to see ~~a~~ GLOBAL TEMP  
view also

when we want to query the global temp view  
global-temp keyword should be added before  
global-temp-viewname  
select \* from global-temp. global-temp.viewname

⑥ CTE syntax  
with cte-name (

new.colname1,  
new.colname2,  
new.colname3

) AS (

select  
col1, col2, col3  
from tbl-name),

⑦ When we want to query a file data that is stored in parquet.

- Select \* from parquet.'location'

⑧ <sup>TABLE</sup> Describe Extended table-name :-  
to view table metadata

or

Describe detail table-name

- The temp view is not associated with any database. The temp view is only accessible in the current spark session
- Global temp view always registers to global temp database

⑨ To query a single file

Select \* from file format.'path/to/file'

⑩ Create a table which takes external data from .csv file

Create table tbl-name (col\_1 datatype, col\_2 datatype)  
using csv  
options (header = "true", delimiter = "|")  
Location "path to .csv file"

⑪ We can manually refresh the cache of our data by running the REFRESH TABLE command



⑫ Extracting data from SQL databases :-

```
= Create table table-name  
using org.apache.spark.sql.jdbc  
options (  
url "-----",  
dbtable "{dbname}.table",  
user "{jdbcusername}",  
password '{jdbcpassword}'  
)
```

⑬ CTAS statements are useful for external data ingestion from sources with WELL DEFINED SCHEMAS, such as parquet files and tables.

So mostly we can't ingest .csv files with CTAS as there are limitations as there we can't specify options parameters for .csv file

So far .csv file :-

First create temp view and write options, path parameter there

and then use CTAS i.e. create table like as select \* from tempviewname

CTAS statement does not support schema declaration.

⑭ Parquets supports 2 types of constraints :-  
1) NOT NULL  
2) CHECK

## (12) Altering a table and adding CHECK constraint

ALTER TABLE tabl-name ADD CONSTRAINT  
constraint-name CHECK (date > '2020-01-01')  
 LTIMindtree

- 2 addition functions can be used in select

current\_timestamp() 2 metadata  
~~int file input-file-name~~ 3 function

## (13) Cloning Delta Lake tables

- ① Deep clone fully copies data and metadata from a source into a target

- Create or replace table table-name  
DEEP CLONE tabl-name

- ② SHALLOW CLONE just moves the transaction logs but the data doesn't move

## 14 (Complete Overwrites)

### 1) CRAS (Create or replace table)

Statements fully replaces the contents of a table each time they execute

### 2) Insert Overwrite

Insert Overwrite table-name

Select \* from parquet.'pathname'

Insert Overwrite will fail if we try to change our schema

## 15 Appending is much more efficient than Overwriting

- Append rows

Insert into tabl-name

Select \* from parquet.'path'   
 This will append new rows result in ~~duplicate~~ records

## (16) MERGE UPDATES

Merge into target a

Using source b

ON {merge-condition}

when Matched then {matched-action}

when Not Matched then {not-matched}

-action?

(17)

## LOAD INCREMENTALLY

Copy into table-name  
From "path"

FileFormat = PARQUET

for this you can add options  
copy-  
guarantee

copy into my-table,  
from 'path/to/file'  
FILE FORMAT <format>

<format-options>  
<copy-options> ('maxSchemaChangeRate': '1mb')

If for this reparation the data schema  
should be consistent copy into support schema  
~~evolution~~ evolution

(18)

## Cleaning of Data

- count (col) does not count NULL  
values

but count (\*) counts null values  
for all columns

- If count null values, use the  
count-if function or where clause to  
filter for records where the value is  
NULL.

e.g select

count-if(user-id is NULL) as  
missing-id

from tbl

(19) : Syntax to traverse nested data structures

- Select value : device, value : geo : city  
from event strings

(20) UDF

Create or replace function func-name  
(i-p name inp datatype)

RETURNS o/p datatype

RETURN . . . .

(21) EXISTS is a high order function.

Syntax :-

EXISTS ( ~~ed~~ items, i → i.item-name LIKE  
"of-Mattress") AS Mattress

It returns true or false

(22) checkpoint directory - The location for storing metadata about the stream. Checkpoints keep track of streaming progress, while the schema location tracks updates to the fields in the source dataset.

ReadStream consist of formating of data it includes

• format ("cloudfiles") ~~may~~)

• option ("cloudfiles.format", arg)

• option ("cloudfiles.SchemaLocation", arg)

• load (data source)

Write stream includes

Checkpoint location, opt Merge Scheme

table name :

$$r = \frac{1}{2} (1 - \frac{1}{n})$$

$\gamma =$   proves

9: : rec

1

23

→ watermarking and watermarking can be used to add additional functionality to incremental workloads.

- Scanning is not supported on streaming data frames

24) When writing to Delta Lake tables, we typically will only need to worry about setting (Means these options needs to be included in writing of code)

- 1) Checkpointing
- 2) Output modes
- 3) Trigger intervals

### LTM Mindtree

1) Append (complete)  
(default) .outputMode("append")

a) Unspecified .outputMode("append")

b) Fixed interval micro-batches

c) Triggered micro-batches

d) Triggered micro-batches

### Cloudfiles

25) schemaHints is used to specify the data type of column.

26) DLT tables and views will always be preceded by LIVE keyword.

27) The ~~to~~ CONSTRAINT keyword in triggers mostly controls.

ON Violations =

- 1) FAIL Update
- 2) DROP Row
- 3) Omit

These can be

applied to  
streaming or  
DLT only

Trigger Type	Example	By default processing time = 500 ms
1) Unspecified		
2) Fixed interval micro batches	.trigger(processingTime = "2 minutes")	
3) Triggered micro- batch	.trigger(conce = true)	
4) Triggered micro- batches	.trigger(availableNow = true)	

→ Selecting from version

select \* from table-name AS OF VERSION 8

→ How to restore from version

Restore Table table-name to VERSION AS OF 8



LTIMindtree

→ set spark.databricks.delta.retentionDuration  
check.enabled = false

Means without checking the duration of transaction file we can delete it directly if it is set to false

set spark.databricks.delta.vaccum.logging  
- enabled = true

Means if it is set to true the transaction log of vacuum we will be able to see in the describe history command.

→ Create a database with location specified

Create database if not exists db-name  
location

→ If we specify the location of a TABLE while creating it, then it becomes external (unmanaged) table

Create or Replace table-name  
LOCATION ' ... ' AS

Select \* from temp-tables

→ To read from data stored in parquet format

Select \* from parquet.'location'

→ Managed VS External Table

Create or replace table  
table-name  
AS

Select \* from ...

(Managed)

Create or replace table table-name  
LOCATION ' ... '  
AS

Select \* from ...

(External)

→ Query a single file  
Select \* from file-format - 'path-to-file'  
eg select \* from json. 'dataset/raw/events/001.json'  
→ Specifying the options for .csv

Create table tbl-name  
(col-name1 col-name1 datatype, col-name2  
(col-name2 datatype))

USING CSV

OPTIONS (

header = "true",

delimiter = ","

)

LOCATION "path to csv file"

CHECK

→ Adding a table constraint. NOT NULL

ALTER TABLE tbl-name ADD CONSTRAINT constraint-name

CHECK (columnname > '2020-01-01');

→ Deep Clone

Create or Replace table tbl-name

Deep clone tbl-name.

← For Shallow Clone replace Deep clone with Shallow Clone.

→ Insert overwrite - To completely overwrite  
Insert overwrite sales  
Select \* from parquet\_file\_location



→ To insert data into delta table from external table

Insert into deltaTable\_name  
select \* from external\_table\_name

→ Insert with JSON data

select value:device, value: geo:city FROM events\_string.

schema of json

from json - Used to convert string json to json object

from json

schema of json - To convert each value into

proper structure type. So that

we can ~~not~~ flatten the proper result in table.

Select from json (value, schema of json ('  
obj: value')) As json

Now to flatten

Select json.\* from parsed\_events;

Note : Whenever u want to convert a json file to table which has nested attributes then first convert a json file to struct data type using from json and then convert it into table.

- Directly reading from a json file
- select value:device,value:geo:city from events
- After converting or flattening the json file (converting to a struct type)

~~→~~ select geo:city from events.

Hint: [Before flattening use : colon to read and after flattening use . (dot) to read.]

(\*)<sup>5</sup> Design video

→ PIVOT : IN SQL [VERY VERY IMPORTANT]  
(rows to columns).

→ AUTO LOADER is used to incrementally load the data from cloud object storage to Delta table.

→ Syntax to read the data from a delta lake table and write it to streaming temp view

e.g (Spark.readStream

```
    .table("bronze")
    .createOrReplaceTempView("streaming_temp_view"))
```

→ Writing aggregated data to a delta table

```
    .writeStream
        .format("delta")
        .option("checkpointLocation", location/path)
        .outputMode("complete")
        .table("customer_count_by_state")
```

→ If you are directly reading from a table or view then write it as spark.table(tbl-name) and if you are reading from some cloud based storage spark.readStream(<sup>through</sup>)



LTIMindtree

→ Select \* current\_timestamp() receipt\_time, input\_file\_name() source file from recordings - raw-temp

→ It will capture source file name

It will capture timestamp

gold

→ Mostly the output mode of delta table is complete whereas the output mode of silver is append.

→ Delta Live Table

To create a delta live table pipeline

Tools → Delta live table → Create a pipeline

In pipeline configurations

target - type the name of target database

Storage location - It is the location to store logs, tables and other information related to pipeline execution

→ Delta live tables and views will always be preceded by the LIVE keyword

→ Incremental processing via Auto Loader (which uses the same processing model as Structured Streaming), requires the addition of STREAMING keyword in the declaration as seen below.

The cloud\_files() method enables Auto Loader to be natively used with SQL. It takes the following parameters :-

1) source location

2) source data format

3) An arbitrarily sized array of optional read options. In this case we set cloudfiles.inferColumnTypes to true

Syntax - Create OR Refresh Streaming Live Table ~~set-order~~

Comment "The row set orders ingested from databricks - raw

AS Select \* from cloud\_files ("source location/path", "json", map("cloudfiles.inferColumnTypes", "true"))

→ 3 modes of constraint on violations	
# ON Violation	Behaviour
1) FAIL UPDATE	Pipeline failure when constraint is violated LTIMindtree
2) DROP ROW	Skip records that violate constraints
3) Committed, warn	Records violating constraints will be included (but violations will be reported in metrics)

→ References to DLT Tables and Views

References to other DLT tables and views will always include the live. prefix

→ References to Streaming Tables

References to streaming DLT tables use the STREAM( ), supplying the table name as an argument

→ Privileges in Data Explorer

- 1) All Privileges
- 2) Select
- 3) Modify
- 4) Read-Metadata
- 5) Usage
- 6) Create

→ The transfer of ownership of data objects like tables and databases can be carried by the current owner of the object OR Patabricks admin only.

Syntax ALTER SCHEMA Schema-name OWNER TO ↑ user-name

~~Note~~, schema and DATABASE can be used interchangeably.

→ Delta Lake supports generated columns which are special type columns whose values are automatically generated based on a user-specified function over other columns in Delta table

e.g suppose from orderTime which has the timestamp datatype I want to generate order date column which takes date directly from orderTime column

e.g create table orders (orderid INT, orderTime timestamp, orderdate DATE GENERATED ALWAYS AS (cast(orderTime as DATE))

## Data Object Privileges

1) Catalog	Create, Usage
2) Schema	Create, Usage
3) <del>Select</del> Table	Select, Modify (insert, update, alter metadata)
4) View	Select
5) Function	Execute
6) Storage / Credential / External Location	Create table, Read files, Write files

L7 LTIMindtree

→ If a Select privilege on View is given then u don't need access to underlying source tables in views to just read data from a view.

→ In order to create a catalog, you should have metastore admin capability.