# Streaming_practice
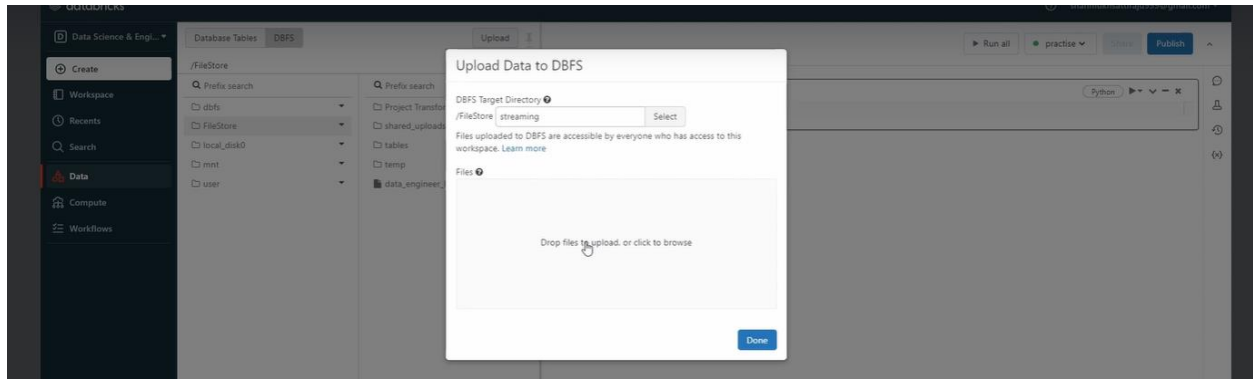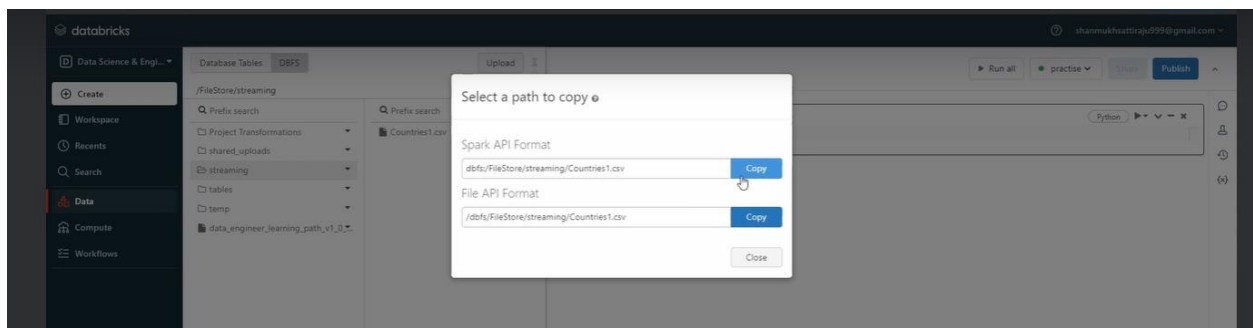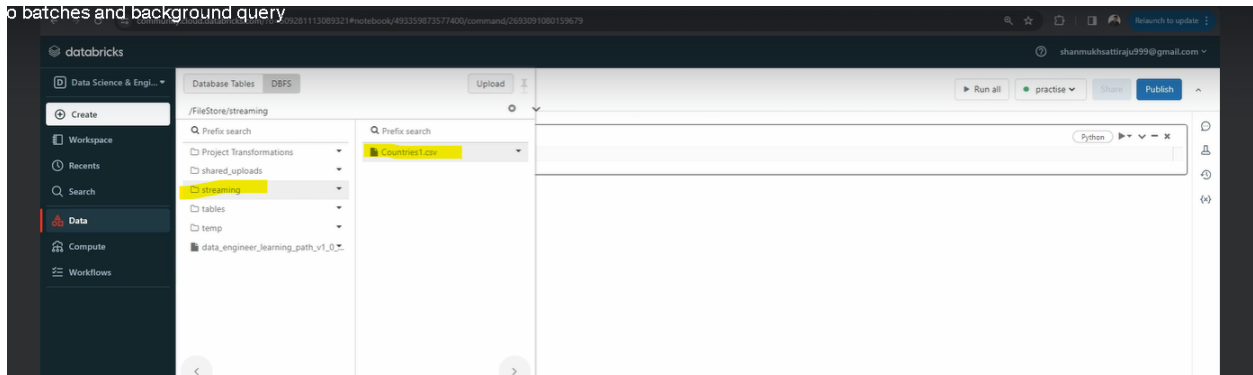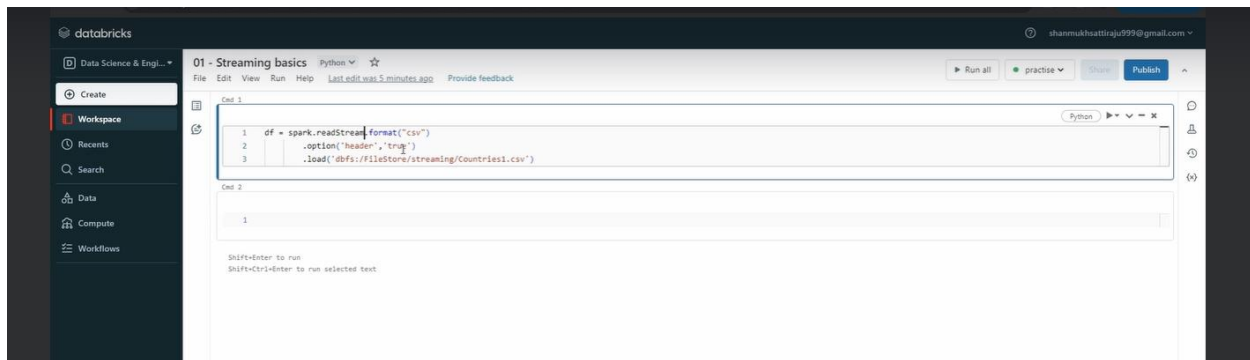


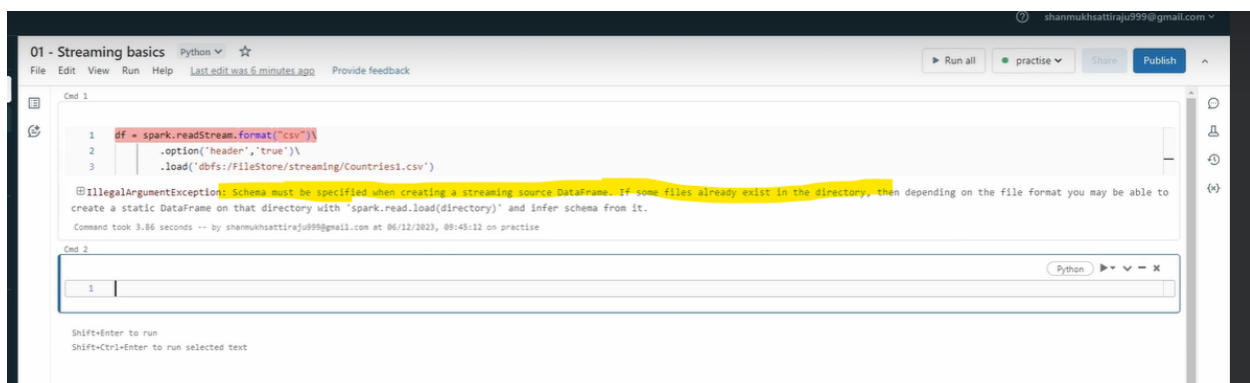Upload a file to dbfs in the workspace creating a folder





Reading the data

Reading a streaming dataset

If we try to read the dataframe directly without specifying or inferring the schema, we will get the below error.
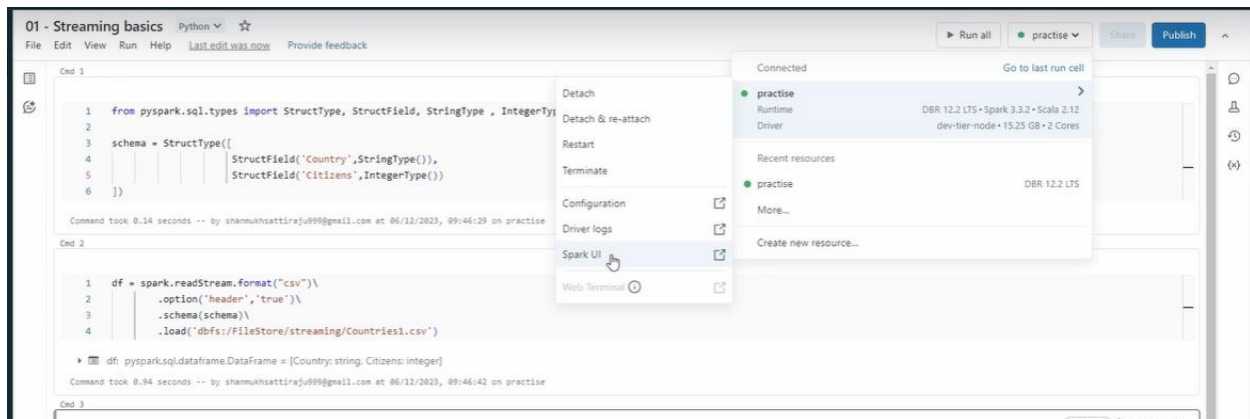


```python
from pyspark.sql.types import StructType, StructField, StringType , IntegerType, FloatType

schema = StructType([
                StructField('Country',StringType()),
                StructField('Citizens',IntegerType())
])
```

Define the schema and re-run the data frame

```python
df = spark.readStream.format("csv")\
        .option('header','true')\
        .schema(schema)\
        .load('dbfs:/FileStore/streaming/Countries1.csv')
```

▸ 🗖 df: pyspark.sql.dataframe.DataFrame = [Country: string, Citizens: integer]

Command took 0.94 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 09:46:42 on practise
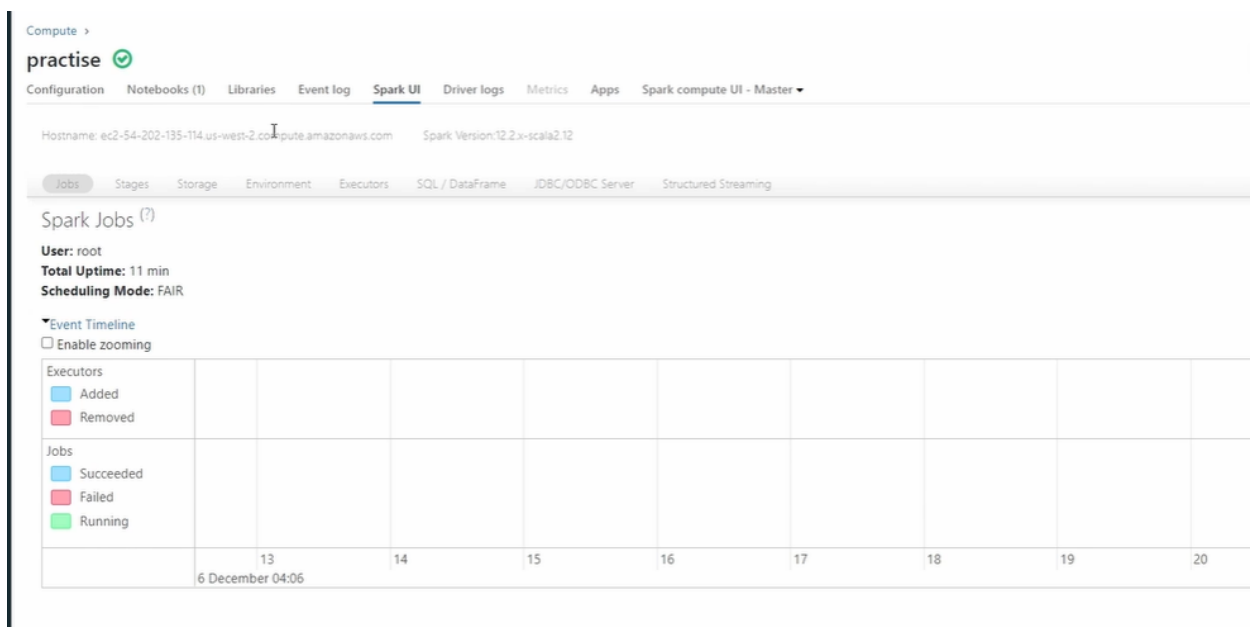
md 3

Rerunning the data frame, specifying the schema will read the stream data

Jobs initiated can be read through the spark UI interface



Initially no jobs were initiated, jobs will be initiated only when trying to fetch some data.



Streaming data can be read through display commad, however basepath must be a directory.

```
1    df = spark.readStream.format("csv")\
2            .option('header','true')\
3            .schema(schema)\
4            .load('dbfs:/FileStore/streaming/')
```
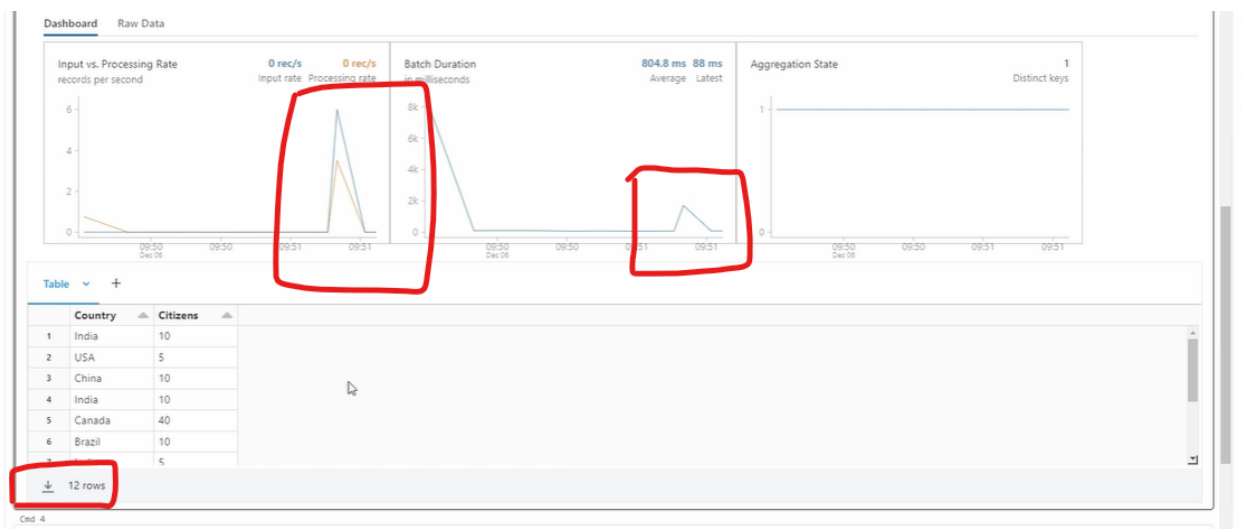
▶ ▦ df: pyspark.sql.dataframe.DataFrame = [Country: string, Citizens: integer]

Command took 0.70 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 09:49:26 on practise

Cmd 3

```
1    display(df)
```

Cancel ▪▪▪

▶ (1) Spark Jobs ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬

⊝ Stream initializing...

Cmd 4

To fix it change the load path to directly to directory level.



To observe the change in the streaming input data, add a new file to the directory. With in few seconds will be noticed in the input and batch streaming graph, without re-executing the display command. Record count is changed to 12 from 6.

Each file addition is read in micro batches.

# Spark Structured Streaming flow



The background streaming query acts as a file watcher, whenever there is a new file added to the directory, a new job will be invoked, the file will be read, transformed and written. The query will be running continuously, until it is stopped explicitly.

Streaming queries can be stopped explicitly by cancelling the query.



## Supported Sources and Sinks

**Sources**

File Source

Kafka Source

Socket Source

Rate Source

Table

**Sinks**

File Sink

Kafka Sink

Foreach Sink

Console Sink

Table

Supported sources and sinks for the Spark Streaming.

# StreamWriter

<StreamingDataframe>.writeStream

        .option('checkpointLocation',<Location>)

        .outputMode('append')

        .toTable('<TableName>')

Checkpoint

To develop fault tolerant and resilient Spark APPS

It maintains intermediate state on fault-tolerant compatible files like DBFS, ADLS and S3 storage systems to recover from the failures.

It is unique to each stream.



```python
(df.writeStream
    .option("checkpointlocation",f"{source_dir}/AppendCheckpoint")
    .outputMode("append")
    .queryName("AppendQuery")
    .toTable("stream.AppendTable"))
```

Out[24]: <pyspark.sql.streaming.query.StreamingQuery at 0x7f404f62c220>

Streaming data frame, data is read, checkpoint location is specified and query name as well specified.

Now after executing the streaming query, a check point location folder is added in the location specified in the above query.



Contains the metadata of the streaming



Assign the streaming query to a variable and to stop the streaming use the command **<var>.stop()**

In the community edition when the cluster is terminated all the schema/database created in the working session will be terminated. Therefore, while using the community edition we should remove the files in the dbfs path, drop the database if exists and re-create it.



When we are trying to run the stream query initially it will verify the check point location, if there are any unprocessed files, only then it will show the updates in the graph

## Output Modes

Deleting the checkpoint location will help to process the older data, without removing the data, just by removing the checkpoint, previous data will be processed.

## Complete Mode



```
Complete

Cmd 10

1    WriteCompleteStream = ( df.writeStream
2        .option('checkpointLocation',f'{source_dir}/CompleteCheckpoint')
3        .outputMode("complete")
4        .queryName('CompleteQuery')
5        .toTable("stream.CompleteTable"))

▶ (4) Spark Jobs

⊞ AnalysisException: Complete output mode not supported when there are no streaming aggregations on streaming DataFrames/Datasets;
FileSource[dbfs:/FileStore/streaming/]

Command took 7.33 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 12:49:33 on newpractise
```

To use the complete output mode aggregation should be used on the data frame.



```
Complete

Cmd 10

1    from pyspark.sql.functions import sum
2    df_complete = df.groupBy("Country").agg(sum('Citizens').alias('Total_Population'))

▶ ⊞ df_complete: pyspark.sql.dataframe.DataFrame = [Country: string, Total_Population: long]

Command took 0.17 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 12:54:45 on newpractise

Cmd 11                                                                          Python  ▶▼ ∨ — ✕

1    WriteCompleteStream = ( df_complete.writeStream
2        .option('checkpointLocation',f'{source_dir}/CompleteCheckpoint')
3        .outputMode("complete")
4        .queryName('CompleteQuery')
5        .toTable("stream.CompleteTable"))

Cancel    Running command...

⊖ Stream initializing...
```

Applied aggregation to execute the **complete output mode**

## Use case of complete mode

When a new file or dataset is added, and we need to apply the aggregation on top it, to fetch the results of the query entire files should be read.

|  | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **Countries1** | | | | | | | | **Result of Countries 1 and 2** | | | |
| 2 | Country | Citizens | | | | | | | | | | |
| 3 | India | 10 | | | | | | | Country | Total Pop | | |
| 4 | USA | 5 | | | | | | | India | 30 | | |
| 5 | China | 10 | | | | | | | USA | 15 | | |
| 6 | India | 10 | | | | | | | China | 15 | | |
| 7 | Canada | 40 | | | | | | | Canada | 50 | | |
| 8 | Brazil | 10 | | | | | | | Brazil | 60 | | |
| 9 | | | | | | | | | | | | |
| 10 | **Countries2** | | | | | | | | | | | |
| 11 | Country | Citizens | | | | | | | | | | |
| 12 | India | 5 | | | | | | | | | | |
| 13 | USA | 10 | | | | | | | | | | |
| 14 | China | 5 | | | | | | | | | | |
| 15 | India | 5 | | | | | | | | | | |
| 16 | Canada | 10 | | | | | | | | | | |
| 17 | Brazil | 50 | | | | | | | | | | |

# Trigger types

## 01. Trigger - default or unspecifed Trigger

Cmd 7

```
1   WriteStream = ( df.writeStream
2       .option('checkpointLocation',f'{source_dir}/AppendCheckpoint')
3       .outputMode("append")
4       .queryName('DefaultTrigger')
5       .toTable("stream.AppendTable"))
```

Cancel

▶ (2) Spark Jobs ━━━━━━━━━━━━━

If no trigger time is specified, by default for every 5sec, streaming query will look for new data. If the new file is added every 5 to 10 min, default trigger is not a good option.

Cmd 9

## 02. Trigger - processingTime

Cmd 10

```
1   WriteStream = ( df.writeStream
2       .option('checkpointLocation',f'{source_dir}/AppendCheckpoint')
3       .outputMode("append")
4       .trigger(processingTime='2 minutes')
5       .queryName('ProcessingTime')
6       .toTable("stream.AppendTable"))
```

Cmd 11

Every 2 min, a micro batch will be processed.

## Streaming Query Statistics

Running batches for **4 minutes 58 seconds** since **2023/12/06 13:08:33** (**4** completed batches)

**Name:** ProcessingTime
**Id:** a279ada4-dcc7-4770-9ceb-72f41d02293f
**RunId:** 6e895f53-48ff-4cfe-9c8c-68896bfcbb36



At 13:08 initial processing occurred, even though data is available only after 2 min next micro batch will be processed.

## 03. Trigger - availablenow

Cmd 12

```
1   WriteStream = ( df.writeStream
2       .option('checkpointLocation',f'{source_dir}/AppendCheckpoint')
3       .outputMode("append")
4       .trigger(availableNow=True)
5       .queryName('AvailableNow')
6       .toTable("stream.AppendTable"))
```

▶ (1) Spark Jobs

▶ ⊖ AvailableNow (id: 9150eaac-ace9-4fd7-a9dc-07c5dc687ac4)   *Last updated: 45 seconds ago*

Is equivalent to incremental batch processing, whenever there is a new data available query will process and stop the streaming

# Triggers

| Triggers | Usage | Description |
|---|---|---|
| Unspecified (default) | | will trigger the microbatch for every 500 ms or half a second |
| processingTime (Fixed Interval) | . trigger(processingTime='2 minutes') | You can set processing time or time interval for each execution . |
| availableNow (OneTime) | .trigger(availableNow = True) | consumes all available records from previous execution as an incremental batch |
| Continuous (experimental) | .trigger(continuous = '1 second') | For ~1ms latency |

Types of streaming triggers

## Autoloader

# Autoloader

- Autoloader is an **optimized data ingestion tool** that incrementally and efficiently processes new data files as they arrive in the cloud storage built into the Databricks Lakehouse.
- Auto Loader incrementally and efficiently processes new data files as they arrive in cloud storage without any additional setup.
- Auto Loader can load data files from Cloud Storages without being vendor specific (AWS S3 , Azure ADLS , Google Cloud Storage, DBFS).
- Auto Loader can ingest JSON, CSV, PARQUET, AVRO, ORC, TEXT, and BINARYFILE file formats
- This Auto loader is beneficial when you are ingesting data into your lakehouse particularly into bronze layer as a streaming query.

```
df_str = (spark.readStream
        .format("cloudFiles")    ## This will tell the spark to use AutoLoader.
        .option("cloudFiles.format","csv")  ## Tells Autoloader to expect csv files
        .option('header','true')
        .schema(schema)
        .load(f'{source_dir}')
)
```

To use the autoloader format("cloudFiles") needs to be specified

## AutoLoader

Cmd 4

```
1    source_dir = 'dbfs:/FileStore/streaming/'
```

Command took 0.14 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 21:41:15 on new

Cmd 5

```
1    df = spark.readStream\
2            .format('cloudFiles')\
3            .option("cloudFiles.format","csv")\
4            .option("cloudFiles.schemaLocation",f'{source_dir}/schemaInfer')\
5            .option("cloudFiles.inferColumnTypes","true")\
6            .option('header','true')\
7            .load(source_dir)
```

▾ ▦ df: pyspark.sql.dataframe.DataFrame
      Country: string
      Citizens: integer
      _rescued_data: string

Command took 0.47 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 21:47:32 on new

To get the correct data types, schema must be inferred.

```
1    %sql
2
3    SELECT *
4    FROM JSON.`dbfs:/FileStore/streaming/schemaInfer/_schemas/0`
```

▸ (2) Spark Jobs

▸ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [_corrupt_record: string, dataSchemaJson: string ... 1 more field]

Table ⌄   +

| | _corrupt_record ▲ | dataSchemaJson | partitionSchemaJson ▲ |
|---|---|---|---|
| 1 | v1 | null | null |
| 2 | null | {"type":"struct","fields":[{"name":"Country","type":"string","nullable":true,"metadata":{}}, {"name":"Citizens","type":"integer","nullable":true,"metadata":{}}]} | {"type":"struct","fields":[]} |

Schema can be referred through the path, rescued data is the additional field created.

### *Schema Hints*

For explicitly defining the data type of a column

## SchemaHints

Cmd 10

```
1   df = spark.readStream\
2         .format('cloudFiles')\
3         .option("cloudFiles.format","csv")\
4         .option("cloudFiles.schemaLocation",f'{source_dir}/schemaInfer')\
5         .option("cloudFiles.inferColumnTypes","true")\
6         .option('cloudFiles.schemaHints',"Citizens LONG")\
7         .option('header','true')\
8         .load(source_dir)
```

▼ ☰ df: pyspark.sql.dataframe.DataFrame
       Country: string
       Citizens: long
       _rescued_data: string

Command took 0.32 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 21:51:01 on new

Cmd 11

For explicitly defining the data type of a column

## Schema Evolution

- Schema evolution is the process of managing changes in data schema as it evolves over time, often due to updates in software or changing business requirements, which can cause schema drift

- Ways to handle schema changes
  - Fail the stream
  - Manually change the existing schema
  - Evolve automatically with change in schema

There are 4 schema evolution types

1. AddNewColunms

   Stream will fail, if new column is added and existing columns do not evolve.

2. FailOnNewColumns

   Stream fails, stream does not restart until schema is updated or offending data is removed.

3. Rescue
   Schema never evolved; stream doesn't change due to schema changes. All new columns are recorded in the rescued data.

4. None

Ignore new columns, (doesn't evolve the schema, new columns are ignored, data is not rescued unless the rescuedDataColumn option is set. Stream doesn't fail due to schema changes)

Using the <mark>rescue type schema evolution</mark>

## Schema Evolution

Cmd 14

```
1   df = spark.readStream\
2         .format('cloudFiles')\
3         .option("cloudFiles.format","csv")\
4         .option("cloudFiles.schemaLocation",f'{source_dir}/schemaInfer')\
5         .option('cloudFiles.schemaEvolutionMode','rescue')\
6         .option('rescuedDataColumn','_rescued_data')\
7         .option("cloudFiles.inferColumnTypes","true")\
8         .option('cloudFiles.schemaHints',"Citizens LONG")\
9         .option('header','true')\
10        .load(source_dir)
```

▼ ☐ df: pyspark.sql.dataframe.DataFrame
        Country: string
        Citizens: long
        _rescued_data: string

Table ∨  +

| | Country | Citizens | _rescued_data |
|---|---|---|---|
| 1 | India | 10 | {"Month":"November","Year":"2023","_file_path":"dbfs:/FileStore/streaming/Countries_newcolumn1.csv"} |
| 2 | USA | 10 | {"Month":"November","Year":"2023","_file_path":"dbfs:/FileStore/streaming/Countries_newcolumn1.csv"} |
| 3 | China | 20 | {"Month":"November","Year":"2023","_file_path":"dbfs:/FileStore/streaming/Countries_newcolumn1.csv"} |
| 4 | Brazil | 10 | {"Month":"November","Year":"2023","_file_path":"dbfs:/FileStore/streaming/Countries_newcolumn1.csv"} |
| 5 | India | 10 | null |
| 6 | USA | 5 | null |
| 7 | China | 10 | null |

Extra columns that are added handled by the rescued_data column

## 02 - addNewColumns - Default

Cmd 19

```
1  df = spark.readStream\
2       .format('cloudFiles')\
3       .option("cloudFiles.format","csv")\
4       .option("cloudFiles.schemaLocation",f'{source_dir}/schemaInfer')\
5       .option("cloudFiles.inferColumnTypes","true")\
6       .option('cloudFiles.schemaHints',"Citizens LONG")\
7       .option('header','true')\
8       .load(source_dir)
```

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [Country: string, Citizens: long ... 1 more field]

Command took 0.26 seconds -- by shanmukhsettiraju999@gmail.com at 06/12/2023, 22:46:17 on new

Cmd 20

Python ▶▾ ∨ ⎯ ✕

```
1  display(df)
```

▶ (1) Spark Jobs

▶ ⊖ display_query_7 (id: 60120d56-4420-4ca8-b356-625185f2b644)    Last updated: 5 seconds ago

⊞ org.apache.spark.sql.catalyst.util.UnknownFieldException: Encountered unknown field(s) during parsing: [Year, Month] in CSV file: dbfs:/FileStore/streaming/Countries_newcolumn1.csv

Command complete

In the addNewColumns(default) type, after adding a file with extra columns, dataframe will be executed, however when we do a display it will fail.

```
1   display(df)
```

Cancel

▶ (1) Spark Jobs ──────────────

▶ ⊖ display_query_8 (id: fb45b5ee-c152-4d90-8227-75c77ce0a14a)    Last updated: 10 seconds ago

Table ∨ ＋

| | Country | Citizens | Year | Month | _rescued_data |
|---|---|---|---|---|---|
| 1 | India | 10 | 2023 | November | null |
| 2 | USA | 10 | 2023 | November | null |
| 3 | China | 20 | 2023 | November | null |
| 4 | Brazil | 10 | 2023 | November | null |
| 5 | India | 10 | null | null | null |
| 6 | USA | 5 | null | null | null |
| 7 | China | 10 | null | null | null |

10 rows

To fix this need to rerun the data frame and display to see the data with evolved schema.

## 03- failOnNewColumns

Cmd 24

```
1    dbutils.fs.ls('dbfs:/FileStore/streaming/schemaInfer/_schemas/')
```

Out[33]: [FileInfo(path='dbfs:/FileStore/streaming/schemaInfer/_schemas/0', name='0', size=274, modificationTime=1701879153000)]

Command took 0.11 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 22:51:58 on new

Cmd 25

```
1    df = spark.readStream\
2            .format('cloudFiles')\
3            .option("cloudFiles.format","csv")\
4            .option("cloudFiles.schemaLocation",f'{source_dir}/schemaInfer')\
5            .option('cloudFiles.schemaEvolutionMode','failOnNewColumns')\
6            .option("cloudFiles.inferColumnTypes","true")\
7            .option('cloudFiles.schemaHints',"Citizens LONG")\
8            .option('header','true')\
9            .load(source_dir)
```

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [Country: string, Citizens: long ... 1 more field]

Command took 0.34 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 22:52:47 on new

Cmd 26

FailOnNewcolumns schema evolution mode, when a new file is added, display will fail, to fix either schema must be redefined or data should be changed.

## None

Cmd 29

```
1    df = spark.readStream\
2            .format('cloudFiles')\
3            .option("cloudFiles.format","csv")\
4            .option("cloudFiles.schemaLocation",f'{source_dir}/schemaInfer')\
5            .option('cloudFiles.schemaEvolutionMode','none')\
6            .option("cloudFiles.inferColumnTypes","true")\
7            .option('cloudFiles.schemaHints',"Citizens LONG")\
8            .option('header','true')\
9            .load(source_dir)
```

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [Country: string, Citizens: long]

Command took 0.27 seconds -- by shanmukhsattiraju999@gmail.com at 06/12/2023, 22:54:44 on new

Cmd 30

```
1    display(df)
```

Cancel •••

▶ (1) Spark Jobs ━━━━━━━━━━━

▶ ⊖ display_query_10 (id: 573f1cbf-356e-4e9e-bf94-23278455c8e7)    Last updated: 10 seconds ago

None type of schema evolution ignores any changes in the schema, add the data to the existing schema.