

Databricks Autoloader Cookbook



Rahul Singha

Follow

10 min read · Mar 16, 2023

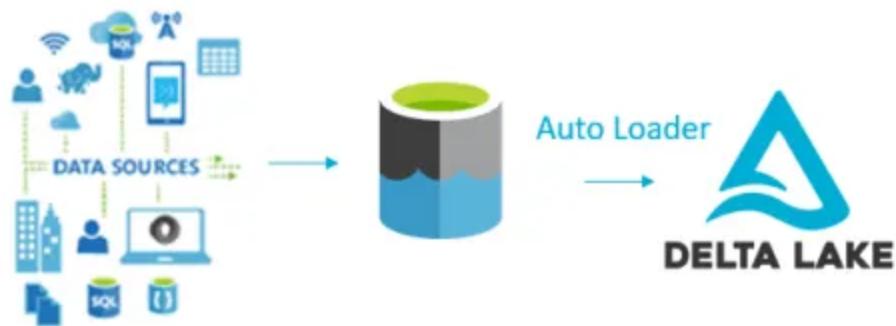
90

1

+

▶

↑



In this article, we are going to discuss the following topics:

1. How Autoloader handles empty files and file names starting with an underscore
2. When to use the compression codec in Autoloader and what are the best practices for compressed files and various file formats
3. modifiedAfter and modifiedBefore in Autoloader

4. partitionColumns

5. overWrites

6. ignoreMissingFiles

7. pathGlobFilter

8. Moving Autoloader Job from one place to another Workspace to Another

Open in app ↗

Sign up

Sign in

Medium



Search



Write



In Databricks, when data is streamed using an autoloader, it should be made sure that the file names must not begin with an underscore '_', Otherwise, files will be ignored by the autoloader.

This can be explained with an example. Initially, three CSV files are kept in a directory in Azure Data Lake Storage, a non-empty CSV (sample1.csv), an empty CSV (empty_file.csv) and a non-empty CSV file whose name starts with an underscore '_' (_sample2.csv). This can be verified using the dbutils.fs.ls.

```
1 display(dbutils.fs.ls("abfss://<container>@<storage_acc>.dfs.core.windows.net/directory"))
```

▶ (2) Spark Jobs

Table +

	path	name	size	modificationTime
1	abfss://<container>@<storage_account>.dfs.core.windows.net/directory/_sample2.csv	_sample2.csv	3240604	1674124682000
2	abfss://<container>@<storage_account>.dfs.core.windows.net/directory/empty_file.csv	empty_file.csv	0	1674124680000
3	abfss://<container>@<storage_account>.dfs.core.windows.net/directory/sample1.csv	sample1.csv	3240604	1674124972000

↓ Showing all 3 rows. | 1.10 seconds runtime

Autoloader is then used to load the files from the directory.

```
spark.readStream  
  .format("cloudFiles")\  
  .option("cloudFiles.format", "csv")\  
  .option("cloudFiles.schemaLocation", "<schema_location>")\  
  .option("cloudFiles.useIncrementalListing", "auto")\  
  .load("abfss://<container>@<storage_account>.dfs.core.windows.net/directory")  
  .writeStream  
  .option("checkpointLocation", "<path_to_checkpoint_location>")\  
  .trigger(availableNow=True) \  
  .table("table_name")
```

The cloud_files_state function of Databricks, which keeps track of the file-level state of an autoloader cloud-file source, confirmed that the autoloader processed only two files, non-empty CSV (sample1.csv) and empty CSV (empty_file.csv). The autoloader ignored the file ‘_sample2.csv’, whose name starts with an ‘_’.

```
1 %sql  
2 select * from cloud_files_state("dbfs:/user/<path>/<to>/<checkpoint>");
```

▶ (1) Spark Jobs

▶ _sqldf: pyspark.sql.DataFrame = [path: string, size: long ... 4 more fields]

Table	+				
path	size	create_time	discovery_time	commit_time	archive_time
1 <container>@<storage_account>/directory/empty_file.csv	0	2023-01-19T10:38:00.000+0000	null	null	null
2 <container>@<storage_account>/directory/sample1.csv	3240604	2023-01-19T10:42:52.000+0000	null	null	null

2. When to use the compression codec in Autoloader and what are the best practices for compressed files and various file formats

This article provides a detailed comparison of the processing time of Autoloader for a variety of file types in Databricks.

The codec used is “ .option(“io.compression.codecs”,
‘nl.basjes.hadoop.io.compress.SplittableGzipCodec’) ”.

In order to use the codec in the Databricks notebook, the appropriate library must be installed in the cluster. Go to cluster details > Libraries tab > Install new > Maven > Search packages > Enter the name of the package that is “splittablegzip” and install.

The cluster configuration used for this comparison is given below.

DBR version 11.3 LTS

Driver and Worker Type: i3.2xlarge 61GB Memory, 8 cores

Number of workers 8

Autoscaling disabled

Used photon acceleration

For convention, the same data is converted to various file formats and compression and is stored in 2 files only.

The size of the file varies significantly for different file formats when the same data is stored.

Case 1: CSV and CSV GZIP

CSV: 15.1 GB x 2

CSV Gzip: 6.86 GB x 2

When processed, the following observations are noted.

File Type	Using Schema	Using Codec	Processing time*
CSV	✗	NA	18.96 mins
CSV	✓	NA	1.35 mins
CSV Gzip file	✗	✗	21.26 mins
CSV Gzip file	✓	✗	9.88 mins
CSV Gzip file	✗	✓	21.87 mins
CSV Gzip file	✓	✓	8.68 mins

Case 2: JSON and JSON GZIP

JSON: 20.88 GB x 2

JSON Gzip: 8.25 GB x 2

When processed, the following observations are note

File Type	Using Schema	Using Codec	Processing time
JSON file	x	NA	3.79 mins
JSON file	✓	NA	1.85 mins
JSON Gzip file	x	x	41.47 mins
JSON Gzip file	✓	x	17.96 mins
JSON Gzip file	x	✓	25.91 mins
JSON Gzip file	✓	✓	4.46 mins

Case 3: TSV and TSV Gzip

TSV: 2.74 GB x 12 files

TSV Gzip: 1.21 GB x 12 files

When processed, the following observations are noted.

File Type	Using Schema	Using Codec	Processing time
TSV file	x	NA	4.53 mins
TSV file	✓	NA	2.31 mins
TSV Gzip file	x	x	4.47 mins
TSV Gzip file	✓	x	2.47 mins
TSV Gzip file	x	✓	4.39 mins
TSV Gzip file	✓	✓	3.85 mins

Case 4: Parquet and Snappy

Parquet: 6.22 GB x 2 files

Snappy: 6.16 GB x 2 files

When processed, the following observations are noted.

File Type	Using Schema	Processing time
Parquet file	✗	4.38 mins
Parquet file	✓	3.87 mins
Parquet snappy file	✗	53 secs
Parquet snappy file	✓	27 secs

Conclusion:

1. Parquet is faster and the size of Parquet compression is less as compared to other file formats when the same data is stored.
 2. Assigning schema manually also will improve the performance because it does not do the schema inference with the huge set of files. Thus, performance will be enhanced. But at the same time, the schema evolution, which is an advantage of Autoloader, will be lost.
 3. The codec works with JSON Gzip compression and helps to process the data faster in Databricks.
 4. The codec does not work with CSV and TSV Gzip compression.
 5. It is recommended to use binary format like parquet for faster processing.
- 3. modifiedAfter and modifiedBefore in Autoloader**

modifiedBefore and modifiedAfter are options that can be applied together or separately in order to achieve greater granularity over which files may load during a Spark batch query.

- modifiedBefore: an optional timestamp to only include files with modification times occurring before the specified time.
- modifiedAfter: an optional timestamp to only include files with modification times occurring after the specified time.

The provided timestamp must be in the following format: YYYY-MM-DDTHH:mm:ss (e.g. 2020-06-01T13:00:00). When a timezone option is not provided, the timestamps will be interpreted according to the Spark session timezone (spark.sql.session.timeZone). The default value is None.

This modification time can be converted from unix to utc using the following command which is the required timestamp format of modifiedAfter and modefiedBefore.

```
import datetime
def unix_to_utc(unix_timestamp):
    return datetime.datetime.utcfromtimestamp(unix_timestamp).strftime('%Y-%m-%d %H:%M:%S')
```

The following line of code can be used to load files using autoloader with modifiedAfter option and write the data to a new delta table

```
spark.readStream
    .format("cloudFiles")
```

```

.option("cloudFiles.format", "csv")
.option("modifiedAfter", str(unix_to_utc(<unix_time>/1000)))
.option("cloudFiles.schemaLocation", "<schema_location>")
.option("cloudFiles.useIncrementalListing", "auto")
.load("<location_of_files>")
.writeStream
.option("checkpointLocation", "<checkpoint_location>")
.trigger(availableNow=True)
.table("table_name")

```

This can be easily understood with the help of an example.

A directory contains 4 files with different modification time stamps which can be seen using **dbutils.fs.ls** command.

```
display(dbutils.fs.ls("directory_path"))
```

▶ (4) Spark Jobs

Table +

	path	name	size	modificationTime
1	file1_path/file1.csv	file1.csv	23	1675069906000
2	file2_path/file2.csv	file2.csv	23	1675069985000
3	file3_path/file3.csv	file3.csv	23	1675070015000
4	file4_path/file4.csv	file4.csv	23	1675070033000

↓ Showing all 4 rows. | 1.08 seconds runtime

When modifiedAfter is used with a timestamp equal to the time stamp of file2, data of only those files where the modification time is strictly after the modifiedAfter timestamps are read. This can be observed by displaying the metadata column.

```

df_mA = spark.readStream\
    .format("cloudFiles")\
    .option("cloudFiles.format", "csv")\
    .option("modifiedAfter", str(unix_to_utc(1675069985000/1000)))\
    .option("cloudFiles.schemaLocation", "<schema_location>")\
    .option("cloudFiles.useIncrementalListing", "auto")\
    .load("<files_location>")\
    .select("*", "_metadata").display()

```

▶ (2) Spark Jobs

↪ ⓧ display_query_5 (id: 3e6d4a10-5633-4aea-926a-73d57cccd86e) Last updated: 19 minutes ago

Table					
	Name	Data	_rescued_data	_metadata	
1	person_4	264	null	‣ {"file_path": '<file_path_of_file_4>', 23, "file_modification_time": "2023-01-30T09:13:53.000+0000"} /file4.csv", "file_name": "file4.csv", "file_size":	
2	person_3	636	null	‣ {"file_path": '<file_path_of_file_3>', 23, "file_modification_time": "2023-01-30T09:13:35.000+0000"} /file3.csv", "file_name": "file3.csv", "file_size":	

Showing all 2 rows.

These options can be used together to get data of files with modification time within a certain interval of time bounded by the modifiedAfter and modifiedBefore. Here, modifiedAfter is used with a timestamp equal to the timestamp of file1 and modifiedBefore is used with a timestamp equal to the timestamp of file4. Thus, only file2 and file3 are loaded.

```

1 df = spark.readStream\
2     .format("cloudFiles")\
3     .option("cloudFiles.format", "csv")\
4     .option("modifiedAfter", str(unix_to_utc(1675069906000/1000)))\
5     .option("modifiedBefore", str(unix_to_utc(1675070033000/1000)))\
6     .option("cloudFiles.schemaLocation", "<schema_location>")\
7     .option("cloudFiles.useIncrementalListing", "auto")\
8     .load("<location_of_files>")\
9     .select("*", "_metadata").display()

```

▶ (1) Spark Jobs

↪ ⓧ display_query_2 (id: 72aeba36-835e-4068-a68f-dd32e8ac63ea) Last updated: 3 minutes ago

Table					
	Name	Data	_rescued_data	_metadata	
1	person_2	720	null	‣ {"file_path": '<file_path_of_file2>', 23, "file_modification_time": "2023-01-30T09:13:05.000+0000"} /file2.csv", "file_name": "file2.csv", "file_size":	
2	person_3	636	null	‣ {"file_path": '<file_path_of_file3>', 23, "file_modification_time": "2023-01-30T09:13:35.000+0000"} /file3.csv", "file_name": "file3.csv", "file_size":	

Showing all 2 rows.

4. partitionColumns inference

Get Rahul Singha's stories in your inbox

Join Medium for free to get updates from this writer.

Enter your email

Subscribe

Partition columns are columns that are inferred from the directory structure of the files.

They are key-value pairs combined by an equality sign such as
<base_path>/a=x/b=y/c=z/file.format.

For example, if directories are created in this order:

root_path>year=2022>month=january>day=21>file1.csv

root_path>year=2023>month=july>day=20>file2.csv

root_path>year=2023>month=october>day=12>file6.csv

The partition columns in the above example are year, month, and date. By default, these columns will be automatically added to the schema if schema inference is used and the <root_path> is provided to load data from. In case the partition column does not appear, the following code can be used to specify the partition columns.

```
df = spark.readStream\
    .format("cloudFiles")\
    .option("cloudFiles.format", "csv")\
    .option("cloudFiles.partitionColumns", "year,month,day")\
```

```
.option("cloudFiles.schemaLocation", "<schema_location>")\
.option("cloudFiles.useIncrementalListing", "auto")\
.load("<root_path>")
```

```
df.display()
```

- ▶ (1) Spark Jobs
- ▶ ⚡ display_query_1 (id: 7094ca6d-9676-43cd-b888-d91c8755e45f) *Last updated: 2 hours ago*

Table ▾ +

	Name	Number	year	month	day	_rescued_data
1	pr1	1234	2022	january	21	null
2	pr2	647	2022	january	21	null
3	pr9	2734	2023	july	20	null
4	pr10	2744	2023	july	20	null
5	pr11	748	2023	october	12	null
6	pr12	7572	2023	october	12	null

If these columns are not required as part of the schema, “” can be specified to hide these columns.

```
df = spark.readStream\
    .format("cloudFiles")\
    .option("cloudFiles.format", "csv")\
    .option("cloudFiles.partitionColumns", "")\
    .option("cloudFiles.schemaLocation", "<schema_location>")\
    .option("cloudFiles.useIncrementalListing", "auto")\
    .load("<root_path>")
```

```
1 df =spark.readStream\  
2     .format("cloudFiles")\  
3     .option("cloudFiles.format", "csv")\  
4     .option("cloudFiles.partitionColumns", "")\  
5     .option("cloudFiles.schemaLocation", "<schema_location>")\  
6     .option("cloudFiles.useIncrementalListing", "auto")\  
7     .load("<location_of_root_folder>")
```

▶ 📄 df: pyspark.sql.dataframe.DataFrame = [Name: string, Number: string ... 1 more field]

Command took 0.35 seconds -- by rahul.singha@databricks.com at 10/02/2023, 11:10:29 on test-rahul

Cmd 8

```
1 df.display()
```

▶ (1) Spark Jobs

▶ ⏱ display_query_4 (id: 3dedaed1-9cb2-4b25-a20f-d99c546a29e0) Last updated: 1 minute ago

Table			
	Name	Number	_rescued_data
1	pr5	3743	null
2	pr6	2542	null
3	pr3	374	null
4	pr4	5763	null
5	pr1	1234	null
6	pr2	647	null

Showing all 12 rows.

5. cloudFiles.allowOverwrites

In Databricks, autoloader by default does not process a file if it is processed once even if the file is modified.

In order to solve this problem, cloudFiles.allowOverwrites option is used to process the files again if the files have been modified.

This option is available in Databricks Runtime 7.6 and above. The default value is false.

This may result in another problem. If files are modified and processed again by the autoloader, duplicate records will be generated.

To tackle this problem, merge insert can be used in each micro batch using Upsert class.

```
sql_query = """
MERGE INTO table_upsert_update a
USING stream_updates b
ON a.Name=b.Name
WHEN MATCHED AND a.Number!=b.Number THEN UPDATE SET *
WHEN NOT MATCHED THEN INSERT *
"""

class Upsert:
    def __init__(self, sql_query, update_temp="stream_updates"):
        self.sql_query = sql_query
        self.update_temp = update_temp

    def upsert_to_delta(self, microBatchDF, batch):
        microBatchDF.createOrReplaceTempView(self.update_temp)
        microBatchDF._jdf.sparkSession().sql(self.sql_query)

#passing the sql_query to Upsert class
streaming_merge = Upsert(sql_query)
```

Since SQL is being used to write to the Delta table, it is required to create the table before using it.

```
%sql
CREATE TABLE IF NOT EXISTS table_upsert_update
```

```
(Name String, Number String)  
USING DELTA
```

Now, if the streaming operation is performed with the Autoloader keeping the allowOverwrite as True, records can be obtained from the overwritten file without allowing duplicate records.

```
query = (spark.readStream\  
    .format("cloudFiles")\  
    .option("cloudFiles.format", "csv")\  
    .option("cloudFiles.allowOverwrites", True)\  
    .option("cloudFiles.schemaLocation", "<schema_location>")\  
    .option("cloudFiles.useIncrementalListing", "auto")\  
    .load("<file_location>")\  
    .writeStream  
    .foreachBatch(streaming_merge.upsert_to_delta)  
    .outputMode("update")  
    .option("checkpointLocation", "<checkpoint_location>")\  
    .start())
```

The records will be stored in table_upsert_update table.

However, Files are processed exactly once unless cloudFiles.allowOverwrites is enabled. If a file is appended to or overwritten, Databricks does not guarantee which version of the file is processed. It is, therefore, recommended to use Auto Loader to ingest only immutable files.

6. ignoreMissingFiles

Autoloader loads files in two steps. In the first step, Autoloader records all the files in the checkpoint directory. This is followed by the actual processing of files from the list of files in the checkpoint directory..

If the file is removed or deleted before the Autoloader processes the file, an error or exception is raised and the streaming stops.

In order to avoid this exception and continue the streaming uninterrupted, the option ignoreMissingFiles is present in Autoloader.

If this option is set to true, the Spark jobs will continue to run when encountering missing files and the contents that have been read will still be returned. Available in Databricks Runtime 11.0 and above.

The default value is False (True for COPY INTO).

```
spark.readStream\  
  .format("cloudFiles")\  
  .option("cloudFiles.format", "csv")\  
  .option("ignoreMissingFiles", True)\  
  .option("cloudFiles.schemaLocation", "<schema_location>")\  
  .option("cloudFiles.useIncrementalListing", "auto")\  
  .load("<files_location>")\  
  .writeStream\  
  .option("checkpointLocation", "<checkpoint_directory>")\  
  .table("table_name")
```

7. pathGlobFilter

pathGlobFilter is used to only include or exclude files with file names matching the pattern. The syntax follows org.apache.hadoop.fs.GlobFilter. It does not change the behavior of partition discovery. **Note: This is validated only for file listing mode.**

```
df=spark.readStream\  
    .format("cloudFiles")\  
    .option("cloudFiles.format", "csv")\  
    .option("cloudFiles.schemaLocation", "<schema_location>")\  
    .option("cloudFiles.useIncrementalListing", "auto")\  
    .option("pathGlobfilter","*pattern*")\  
    .load("<file_location>")
```

This will process files with the following pattern in the file name.

In pathGlobfilter,

- ? matches only a single character
- * matches zero or more than one character.
- ^ Matches a single character that is not from the character set or range {a}. Note that the ^ character must occur immediately to the right of the opening bracket.

Pattern	Description
pattern	Process files if the given pattern is present in the file name
*.csv	Processes only CSV files
{pattern1,pattern2} *{p1,p2,p3,p4}* *{p1,p2,p3,p4}*{p5,p6,p7,p8}*{p9,p10,p11,p12}*{p13,p14,p15,p16}*{p17,p18,p19,p20}*{p21,p22,p23,p24}*{p25,p26,p27,p28}*{p29,p30,p31,p32}*{p33,p34,p35,p36}*{p37,p38,p39,p40}*{p41,p42,p43,p44}*{p45,p46,p47,p48}*{p49,p50,p51,p52}*{p53,p54,p55,p56}*{p57,p58,p59,p60}*{p61,p62,p63,p64}*{p65,p66,p67,p68}*{p69,p70,p71,p72}*{p73,p74,p75,p76}*{p77,p78,p79,p80}*{p81,p82,p83,p84}*{p85,p86,p87,p88}*{p89,p90,p91,p92}*{p93,p94,p95,p96}*{p97,p98,p99,p100}*{p101,p102,p103,p104}*{p105,p106,p107,p108}*{p109,p110,p111,p112}*{p113,p114,p115,p116}*{p117,p118,p119,p120}*{p121,p122,p123,p124}*{p126,p127,p128,p129}*{p131,p132,p133,p134}*{p136,p137,p138,p139}*{p141,p142,p143,p144}*{p147,p148,p149,p150}*{p153,p154,p155,p156}*{p159,p160,p161,p162}*{p165,p166,p167,p168}*{p171,p172,p173,p174}*{p177,p178,p179,p180}*{p183,p184,p185,p186}*{p189,p190,p191,p192}*{p195,p196,p197,p198}*{p199,p200,p201,p202}*{p205,p206,p207,p208}*{p211,p212,p213,p214}*{p217,p218,p219,p220}*{p223,p224,p225,p226}*{p229,p230,p231,p232}*{p235,p236,p237,p238}*{p241,p242,p243,p244}*{p247,p248,p249,p250}*{p253,p254,p255,p256}*{p259,p260,p261,p262}*{p265,p266,p267,p268}*{p271,p272,p273,p274}*{p277,p278,p279,p280}*{p283,p284,p285,p286}*{p289,p290,p291,p292}*{p295,p296,p297,p298}*{p299,p300,p301,p302}*{p305,p306,p307,p308}*{p311,p312,p313,p314}*{p317,p318,p319,p320}*{p323,p324,p325,p326}*{p329,p330,p331,p332}*{p335,p336,p337,p338}*{p341,p342,p343,p344}*{p347,p348,p349,p350}*{p353,p354,p355,p356}*{p359,p360,p361,p362}*{p365,p366,p367,p368}*{p371,p372,p373,p374}*{p377,p378,p379,p380}*{p383,p384,p385,p386}*{p389,p390,p391,p392}*{p395,p396,p397,p398}*{p399,p400,p401,p402}*{p405,p406,p407,p408}*{p411,p412,p413,p414}*{p417,p418,p419,p420}*{p423,p424,p425,p426}*{p429,p430,p431,p432}*{p435,p436,p437,p438}*{p441,p442,p443,p444}*{p447,p448,p449,p450}*{p453,p454,p455,p456}*{p459,p460,p461,p462}*{p465,p466,p467,p468}*{p471,p472,p473,p474}*{p477,p478,p479,p480}*{p483,p484,p485,p486}*{p489,p490,p491,p492}*{p495,p496,p497,p498}*{p499,p500,p501,p502}*{p505,p506,p507,p508}*{p511,p512,p513,p514}*{p517,p518,p519,p520}*{p523,p524,p525,p526}*{p529,p530,p531,p532}*{p535,p536,p537,p538}*{p541,p542,p543,p544}*{p547,p548,p549,p550}*{p553,p554,p555,p556}*{p559,p560,p561,p562}*{p565,p566,p567,p568}*{p571,p572,p573,p574}*{p577,p578,p579,p580}*{p583,p584,p585,p586}*{p589,p590,p591,p592}*{p595,p596,p597,p598}*{p599,p600,p601,p602}*{p605,p606,p607,p608}*{p611,p612,p613,p614}*{p617,p618,p619,p620}*{p623,p624,p625,p626}*{p629,p630,p631,p632}*{p635,p636,p637,p638}*{p641,p642,p643,p644}*{p647,p648,p649,p650}*{p653,p654,p655,p656}*{p659,p660,p661,p662}*{p665,p666,p667,p668}*{p671,p672,p673,p674}*{p677,p678,p679,p680}*{p683,p684,p685,p686}*{p689,p690,p691,p692}*{p695,p696,p697,p698}*{p699,p700,p701,p702}*{p705,p706,p707,p708}*{p711,p712,p713,p714}*{p717,p718,p719,p720}*{p723,p724,p725,p726}*{p729,p730,p731,p732}*{p735,p736,p737,p738}*{p741,p742,p743,p744}*{p747,p748,p749,p750}*{p753,p754,p755,p756}*{p759,p760,p761,p762}*{p765,p766,p767,p768}*{p771,p772,p773,p774}*{p777,p778,p779,p780}*{p783,p784,p785,p786}*{p789,p790,p791,p792}*{p795,p796,p797,p798}*{p799,p800,p801,p802}*{p805,p806,p807,p808}*{p811,p812,p813,p814}*{p817,p818,p819,p820}*{p823,p824,p825,p826}*{p829,p830,p831,p832}*{p835,p836,p837,p838}*{p841,p842,p843,p844}*{p847,p848,p849,p850}*{p853,p854,p855,p856}*{p859,p860,p861,p862}*{p865,p866,p867,p868}*{p871,p872,p873,p874}*{p877,p878,p879,p880}*{p883,p884,p885,p886}*{p889,p890,p891,p892}*{p895,p896,p897,p898}*{p899,p900,p901,p902}*{p905,p906,p907,p908}*{p911,p912,p913,p914}*{p917,p918,p919,p920}*{p923,p924,p925,p926}*{p929,p930,p931,p932}*{p935,p936,p937,p938}*{p941,p942,p943,p944}*{p947,p948,p949,p950}*{p953,p954,p955,p956}*{p959,p960,p961,p962}*{p965,p966,p967,p968}*{p971,p972,p973,p974}*{p977,p978,p979,p980}*{p983,p984,p985,p986}*{p989,p990,p991,p992}*{p995,p996,p997,p998}*{p999,p1000,p1001,p1002}*{p1005,p1006,p1007,p1008}*{p1011,p1012,p1013,p1014}*{p1017,p1018,p1019,p1020}*{p1023,p1024,p1025,p1026}*{p1029,p1030,p1031,p1032}*{p1035,p1036,p1037,p1038}*{p1041,p1042,p1043,p1044}*{p1047,p1048,p1049,p1050}*{p1053,p1054,p1055,p1056}*{p1059,p1060,p1061,p1062}*{p1065,p1066,p1067,p1068}*{p1071,p1072,p1073,p1074}*{p1077,p1078,p1079,p1080}*{p1083,p1084,p1085,p1086}*{p1089,p1090,p1091,p1092}*{p1095,p1096,p1097,p1098}*{p1099,p1100,p1101,p1102}*{p1105,p1106,p1107,p1108}*{p1111,p1112,p1113,p1114}*{p1117,p1118,p1119,p1120}*{p1123,p1124,p1125,p1126}*{p1129,p1130,p1131,p1132}*{p1135,p1136,p1137,p1138}*{p1141,p1142,p1143,p1144}*{p1147,p1148,p1149,p1150}*{p1153,p1154,p1155,p1156}*{p1159,p1160,p1161,p1162}*{p1165,p1166,p1167,p1168}*{p1171,p1172,p1173,p1174}*{p1177,p1178,p1179,p1180}*{p1183,p1184,p1185,p1186}*{p1189,p1190,p1191,p1192}*{p1195,p1196,p1197,p1198}*{p1199,p1200,p1201,p1202}*{p1205,p1206,p1207,p1208}*{p1211,p1212,p1213,p1214}*{p1217,p1218,p1219,p1220}*{p1223,p1224,p1225,p1226}*{p1229,p1230,p1231,p1232}*{p1235,p1236,p1237,p1238}*{p1241,p1242,p1243,p1244}*{p1247,p1248,p1249,p1250}*{p1253,p1254,p1255,p1256}*{p1259,p1260,p1261,p1262}*{p1265,p1266,p1267,p1268}*{p1271,p1272,p1273,p1274}*{p1277,p1278,p1279,p1280}*{p1283,p1284,p1285,p1286}*{p1289,p1290,p1291,p1292}*{p1295,p1296,p1297,p1298}*{p1299,p1300,p1301,p1302}*{p1305,p1306,p1307,p1308}*{p1311,p1312,p1313,p1314}*{p1317,p1318,p1319,p1320}*{p1323,p1324,p1325,p1326}*{p1329,p1330,p1331,p1332}*{p1335,p1336,p1337,p1338}*{p1341,p1342,p1343,p1344}*{p1347,p1348,p1349,p1350}*{p1353,p1354,p1355,p1356}*{p1359,p1360,p1361,p1362}*{p1365,p1366,p1367,p1368}*{p1371,p1372,p1373,p1374}*{p1377,p1378,p1379,p1380}*{p1383,p1384,p1385,p1386}*{p1389,p1390,p1391,p1392}*{p1395,p1396,p1397,p1398}*{p1399,p1400,p1401,p1402}*{p1405,p1406,p1407,p1408}*{p1411,p1412,p1413,p1414}*{p1417,p1418,p1419,p1420}*{p1423,p1424,p1425,p1426}*{p1429,p1430,p1431,p1432}*{p1435,p1436,p1437,p1438}*{p1441,p1442,p1443,p1444}*{p1447,p1448,p1449,p1450}*{p1453,p1454,p1455,p1456}*{p1459,p1460,p1461,p1462}*{p1465,p1466,p1467,p1468}*{p1471,p1472,p1473,p1474}*{p1477,p1478,p1479,p1480}*{p1483,p1484,p1485,p1486}*{p1489,p1490,p1491,p1492}*{p1495,p1496,p1497,p1498}*{p1499,p1500,p1501,p1502}*{p1505,p1506,p1507,p1508}*{p1511,p1512,p1513,p1514}*{p1517,p1518,p1519,p1520}*{p1523,p1524,p1525,p1526}*{p1529,p1530,p1531,p1532}*{p1535,p1536,p1537,p1538}*{p1541,p1542,p1543,p1544}*{p1547,p1548,p1549,p1550}*{p1553,p1554,p1555,p1556}*{p1559,p1560,p1561,p1562}*{p1565,p1566,p1567,p1568}*{p1571,p1572,p1573,p1574}*{p1577,p1578,p1579,p1580}*{p1583,p1584,p1585,p1586}*{p1589,p1590,p1591,p1592}*{p1595,p1596,p1597,p1598}*{p1599,p1600,p1601,p1602}*{p1605,p1606,p1607,p1608}*{p1611,p1612,p1613,p1614}*{p1617,p1618,p1619,p1620}*{p1623,p1624,p1625,p1626}*{p1629,p1630,p1631,p1632}*{p1635,p1636,p1637,p1638}*{p1641,p1642,p1643,p1644}*{p1647,p1648,p1649,p1650}*{p1653,p1654,p1655,p1656}*{p1659,p1660,p1661,p1662}*{p1665,p1666,p1667,p1668}*{p1671,p1672,p1673,p1674}*{p1677,p1678,p1679,p1680}*{p1683,p1684,p1685,p1686}*{p1689,p1690,p1691,p1692}*{p1695,p1696,p1697,p1698}*{p1699,p1700,p1701,p1702}*{p1705,p1706,p1707,p1708}*{p1711,p1712,p1713,p1714}*{p1717,p1718,p1719,p1720}*{p1723,p1724,p1725,p1726}*{p1729,p1730,p1731,p1732}*{p1735,p1736,p1737,p1738}*{p1741,p1742,p1743,p1744}*{p1747,p1748,p1749,p1750}*{p1753,p1754,p1755,p1756}*{p1759,p1760,p1761,p1762}*{p1765,p1766,p1767,p1768}*{p1771,p1772,p1773,p1774}*{p1777,p1778,p1779,p1780}*{p1783,p1784,p1785,p1786}*{p1789,p1790,p1791,p1792}*{p1795,p1796,p1797,p1798}*{p1799,p1800,p1801,p1802}*{p1805,p1806,p1807,p1808}*{p1811,p1812,p1813,p1814}*{p1817,p1818,p1819,p1820}*{p1823,p1824,p1825,p1826}*{p1829,p1830,p1831,p1832}*{p1835,p1836,p1837,p1838}*{p1841,p1842,p1843,p1844}*{p1847,p1848,p1849,p1850}*{p1853,p1854,p1855,p1856}*{p1859,p1860,p1861,p1862}*{p1865,p1866,p1867,p1868}*{p1871,p1872,p1873,p1874}*{p1877,p1878,p1879,p1880}*{p1883,p1884,p1885,p1886}*{p1889,p1890,p1891,p1892}*{p1895,p1896,p1897,p1898}*{p1899,p1900,p1901,p1902}*{p1905,p1906,p1907,p1908}*{p1911,p1912,p1913,p1914}*{p1917,p1918,p1919,p1920}*{p1923,p1924,p1925,p1926}*{p1929,p1930,p1931,p1932}*{p1935,p1936,p1937,p1938}*{p1941,p1942,p1943,p1944}*{p1947,p1948,p1949,p1950}*{p1953,p1954,p1955,p1956}*{p1959,p1960,p1961,p1962}*{p1965,p1966,p1967,p1968}*{p1971,p1972,p1973,p1974}*{p1977,p1978,p1979,p1980}*{p1983,p1984,p1985,p1986}*{p1989,p1990,p1991,p1992}*{p1995,p1996,p1997,p1998}*{p1999,p2000,p2001,p2002}*{p2005,p2006,p2007,p2008}*{p2011,p2012,p2013,p2014}*{p2017,p2018,p2019,p2020}*{p2023,p2024,p2025,p2026}*{p2029,p2030,p2031,p2032}*{p2035,p2036,p2037,p2038}*{p2041,p2042,p2043,p2044}*{p2047,p2048,p2049,p2050}*{p2053,p2054,p2055,p2056}*{p2059,p2060,p2061,p2062}*{p2065,p2066,p2067,p2068}*{p2071,p2072,p2073,p2074}*{p2077,p2078,p2079,p2080}*{p2083,p2084,p2085,p2086}*{p2089,p2090,p2091,p2092}*{p2095,p2096,p2097,p2098}*{p2099,p2100,p2101,p2102}*{p2105,p2106,p2107,p2108}*{p2111,p2112,p2113,p2114}*{p2117,p2118,p2119,p2120}*{p2123,p2124,p2125,p2126}*{p2129,p2130,p2131,p2132}*{p2135,p2136,p2137,p2138}*{p2141,p2142,p2143,p2144}*{p2147,p2148,p2149,p2150}*{p2153,p2154,p2155,p2156}*{p2159,p2160,p2161,p2162}*{p2165,p2166,p2167,p2168}*{p2171,p2172,p2173,p2174}*{p2177,p2178,p2179,p2180}*{p2183,p2184,p2185,p2186}*{p2189,p2190,p2191,p2192}*{p2195,p2196,p2197,p2198}*{p2199,p2200,p2201,p2202}*{p2205,p2206,p2207,p2208}*{p2211,p2212,p2213,p2214}*{p2217,p2218,p2219,p2220}*{p2223,p2224,p2225,p2226}*{p2229,p2230,p2231,p2232}*{p2235,p2236,p2237,p2238}*{p2241,p2242,p2243,p2244}*{p2247,p2248,p2249,p2250}*{p2253,p2254,p2255,p2256}*{p2259,p2260,p2261,p2262}*{p2265,p2266,p2267,p2268}*{p2271,p2272,p2273,p2274}*{p2277,p2278,p2279,p2280}*{p2283,p2284,p2285,p2286}*{p2289,p2290,p2291,p2292}*{p2295,p2296,p2297,p2298}*{p2299,p2300,p2301,p2302}*{p2305,p2306,p2307,p2308}*{p2311,p2312,p2313,p2314}*{p2317,p2318,p2319,p2320}*{p2323,p2324,p2325,p2326}*{p2329,p2330,p2331,p2332}*{p2335,p2336,p2337,p2338}*{p2341,p2342,p2343,p2344}*{p2347,p2348,p2349,p2350}*{p2353,p2354,p2355,p2356}*{p2359,p2360,p2361,p2362}*{p2365,p2366,p2367,p2368}*{p2371,p2372,p2373,p2374}*{p2377,p2378,p2379,p2380}*{p2383,p2384,p2385,p2386}*{p2389,p2390,p2391,p2392}*{p2395,p2396,p2397,p2398}*{p2399,p2400,p2401,p2402}*{p2405,p2406,p2407,p2408}*{p2411,p2412,p2413,p2414}*{p2417,p2418,p2419,p2420}*{p2423,p2424,p2425,p2426}*{p2429,p2430,p2431,p2432}*{p2435,p2436,p2437,p2438}*{p2441,p2442,p2443,p2444}*{p2447,p2448,p2449,p2450}*{p2453,p2454,p2455,p2456}*{p2459,p2460,p2461,p2462}*{p2465,p2466,p2467,p2468}*{p2471,p2472,p2473,p2474}*{p2477,p2478,p2479,p2480}*{p2483,p2484,p2485,p2486}*{p2489,p2490,p2491,p2492}*{p2495,p2496,p2497,p2498}*{p2499,p2500,p2501,p2502}*{p2505,p2506,p2507,p2508}*{p2511,p2512,p2513,p2514}*{p2517,p2518,p2519,p2520}*{p2523,p2524,p2525,p2526}*{p2529,p2530,p2531,p2532}*{p2535,p2536,p2537,p2538}*{p2541,p2542,p2543,p2544}*{p2547,p2548,p2549,p2550}*{p2553,p2554,p2555,p2556}*{p2559,p2560,p2561,p2562}*{p2565,p2566,p2567,p2568}*{p2571,p2572,p2573,p2574}*{p2577,p2578,p2579,p2580}*{p2583,p2584,p2585,p2586}*{p2589,p2590,p2591,p2592}*{p2595,p2596,p2597,p2598}*{p2599,p2600,p2601,p2602}*{p2605,p2606,p2607,p2608}*{p2611,p2612,p2613,p2614}*{p2617,p2618,p2619,p2620}*{p2623,p2624,p2625,p2626}*{p2629,p2630,p2631,p2632}*{p2635,p2636,p2637,p2638}*{p2641,p2642,p2643,p2644}*{p2647,p2648,p2649,p2650}*{p2653,p2654,p2655,p2656}*{p2659,p2660,p2661,p2662}*{p2665,p2666,p2667,p2668}*{p2671,p2672,p2673,p2674}*{p2677,p2678,p2679,p2680}*{p2683,p2684,p2685,p2686}*{p2689,p2690,p2691,p2692}*{p2695,p2696,p2697,p2698}*{p2699,p2700,p2701,p2702}*{p2705,p2706,p2707,p2708}*{p2711,p2712,p2713,p2714}*{p2717,p2718,p2719,p2720}*{p2723,p2724,p2725,p2726}*{p2729,p2730,p2731,p2732}*{p2735,p2736,p2737,p2738}*{p2741,p2742,p2743,p2744}*{p2747,p2748,p2749,p2750}*{p2753,p2754,p2755,p2756}*{p2759,p2760,p2761,p2762}*{p2765,p2766,p2767,p2768}*{p2771,p2772,p2773,p2774}*{p2777,p2778,p2779,p2780}*{p2783,p2784,p2785,p2786}*{p2789,p2790,p2791,p2792}*{p2795,p2796,p2797,p2798}*{p2799,p2800,p2801,p2802}*{p2805,p2806,p2807,p2808}*{p2811,p2812,p2813,p2814}*{p2817,p2818,p2819,p2820}*{p2823,p2824,p2825,p2826}*{p2829,p2830,p2831,p2832}*{p2835,p2836,p2837,p2838}*{p2841,p2842,p2843,p2844}*{p2847,p2848,p2849,p2850}*{p2853,p2854,p2855,p2856}*{p2859,p2860,p2861,p2862}*{p2865,p2866,p2867,p2868}*{p2871,p2872,p2873,p2874}*{p2877,p2878,p2879,p2880}*{p2883,p2884,p2885,p2886}*{p2889,p2890,p2891,p2892}*{p2895,p2896,p2897,p2898}*{p2899,p2900,p2901,p2902}*{p2905,p2906,p2907,p2908}*{p2911,p2912,p2913,p2914}*{p2917,p2918,p2919,p2920}*{p2923,p2924,p2925,p2926}*{p2929,p2930,p2931,p2932}*{p2935,p2936,p2937,p2938}*{p2941,p2942,p2943,p2944}*{p2947,p2948,p2949,p2950}*{p2953,p2954,p2955,p2956}*{p2959,p2960,p2961,p2962}*{p2965,p2966,p2967,p2968}*{p2971,p2972,p2973,p2974}*{p2977,p2978,p2979,p2980}*{p2983,p2984,p2985,p2986}*{p2989,p2990,p2991,p2992}*{p2995,p2996,p2997,p2998}*{p2999,p3000,p3001,p3002}*{p3005,p3006,p3007,p3008}*{p3011,p3012,p3013,p301	

pathGlobfilter can only work with the file names and not with the file path and directories.

pathGlobfilter can also be used in the path e.g., `.load('/source-path/2022/06/13/05*/')`, or as a combination of both.

In order to attain filtering based on a particular path or directories, metadata column can be used to filter post-landing of data in the sink. For example

```
df.select("*","_metadata").select("*","_metadata.file_path").filter("file_path not REGEXP '/folder_name/' ").display()
```

This will exclude all files from the given folder_name. However, this is just a simple filter (using SQL REGEXP) applied to the whole data frame and excludes data post-landing in Autoloader.

This can be easily understood with the help of an example.

Let's assume that there are some files in a directory containing sales data for shop1, shop2, and online sales.

```
display(dbutils.fs.ls("<file_location>"))
```

▶ (4) Spark Jobs

Table +

	path		name	size	modificationTime
1	<file_location>	/online_01012023.csv	online_01012023.csv	54	1678992892000
2	<file_location>	!/shop1_01012023.csv	shop1_01012023.csv	52	1678992906000
3	<file_location>	!/shop1_02012023.csv	shop1_02012023.csv	52	1678992915000
4	<file_location>	!/shop2_01012023.csv	shop2_01012023.csv	54	1678992924000
5	<file_location>	!/shop2_02012023.csv	shop2_02012023.csv	56	1678992196000

↓ 5 rows | 1.04 seconds runtime

Now if it is required to process files of only shop1, the following code can be used.

```
df_shop1=spark.readStream\  
    .format("cloudFiles")\  
    .option("cloudFiles.format", "csv")\  
    .option("cloudFiles.schemaLocation", "<schema_location>")\  
    .option("cloudFiles.useIncrementalListing", "auto")\  
    .option("pathGlobfilter","*{shop1}*")\  
    .load("<file_location>")
```

This will load data only from files that have the pattern “shop1” in the file name which can be verified by displaying the metadata column of the dataframe.

df_shop1.select("*", "_metadata").display()					
▶ (1) Spark Jobs ▶ ⓧ display_query_4 (id: 183915f8-9c28-4bc0-88ca-bd8695231e14) Last updated: 10 minutes ago					
Table ▾ +					
ProductID	Quantity	CostPerItem	_rescued_data	_metadata	
1	pr2	9	100	null	▶ {"file_path": "<-----file_location----->/shop1_02012023.csv", "file_name": "shop1_02012023.csv", "file_size": 52, "file_modification_time": "2023-03-16T18:55:15.000+0000"}
2	pr4	1	300	null	▶ {"file_path": "<-----file_location----->/shop1_02012023.csv", "file_name": "shop1_02012023.csv", "file_size": 52, "file_modification_time": "2023-03-16T18:55:15.000+0000"}
3	pr5	1	100	null	▶ {"file_path": "<-----file_location----->/shop1_01012023.csv", "file_name": "shop1_01012023.csv", "file_size": 52, "file_modification_time": "2023-03-16T18:55:06.000+0000"}
4	pr6	5	500	null	▶ {"file_path": "<-----file_location----->/shop1_01012023.csv", "file_name": "shop1_01012023.csv", "file_size": 52, "file_modification_time": "2023-03-16T18:55:06.000+0000"}
↓	4 rows				

Similarly, if to load data of shop1 and online sales,

“*{shop1,online}*” can be used in pathGlobfilter.

Now, if it is required to load data from all files except certain files, let's say, sales data from shop2, pathGlobfilter cannot be used. In such cases, SQL

REGEXP can be used.

```
df_notshop2 =spark.readStream\  
    .format("cloudFiles")\  
    .option("cloudFiles.format", "csv")\  
    .option("cloudFiles.schemaLocation", "<schema_location>")\  
    .option("cloudFiles.useIncrementalListing", "auto")\  
    .load("<file_location>")\  
    .select("*","_metadata").select("*","_metadata.file_path")\  
    .filter("file_path not REGEXP 'shop2' ")
```

8. Moving Autoloader Job from one Workspace to another to Another

Autoloader jobs can be easily moved from one workspace to another by just replicating the same code in the target workspace and directing the checkpoint to the same checkpoint used before.

If the checkpoint point location is also changed, the Autoloader will start from the beginning and there may be duplicate records in the sink.

An example will make it more clear.

Autoloader is used to load files from a source with the code as shown below:

```
spark.readStream\  
    .format("cloudFiles")\  
    .option("cloudFiles.format", "csv")\
```

```

.option("cloudFiles.schemaLocation", "<schema_location>")\
.option("cloudFiles.useIncrementalListing", "auto")\
.load("<location_of_files>")\
.withColumn("Processing_timestamp", current_timestamp())\
.writeStream\
.option("checkpointLocation", "<initial_checkpoint_location>")\
.trigger(availableNow=True)\
.table("table_move")

```

In this sample of code, a column “Processing_timestamp” is created which will store the timestamp when the record is processed by the autoloader.

```
%sql
select * from table_move
```

- ▶ (2) Spark Jobs
- ▶ _sqldf: pyspark.sql.dataframe.DataFrame = [Name: string, Number: string ... 2 more fields]

	Name	Number	_rescued_data	Processing_timestamp
1	pr1	1234	null	2023-03-16T21:46:11.043+0000
2	pr2	364	null	2023-03-16T21:46:11.043+0000
3	pr3	731	null	2023-03-16T21:46:11.043+0000

↓ 3 rows | 0.38 seconds runtime

New files are then loaded in the source location.

Now, the autoloader code is moved to another workspace with a new checkpoint location and the autoloader is fired. It is observed that duplicate records are generated in the sink, table_move. The Processing_timestamp column signifies that the old files are processed again along with the new files.

```
%sql
select * from table_move
```

▶ (3) Spark Jobs

▶ _sqldf: pyspark.sql.dataframe.DataFrame = [Name: string, Number: string ... 2 more fields]

Table +

	Name	Number	_rescued_data	Processing_timestamp	
1	pr2	364	null	2023-03-16T21:46:11.043+0000	
2	pr3	731	null	2023-03-16T21:46:11.043+0000	
3	pr1	1234	null	2023-03-16T21:46:11.043+0000	
4	pr4	873	null	2023-03-16T21:53:03.743+0000	
5	pr1	1234	null	2023-03-16T21:53:03.743+0000	
6	pr3	731	null	2023-03-16T21:53:03.743+0000	
7	pr2	364	null	2023-03-16T21:53:03.743+0000	

↓ 7 rows | 0.52 seconds runtime

Duplicate Records

If it is required to resume the Autoloader from where it stopped last to avoid duplicate records in the sink, two approaches can be followed:

1. the autoloader should point to the old checkpoint location
2. the files of the old checkpoint location should be copied to the new checkpoint location.

However, after changing the checkpoint location, the autoloader code can be modified to perform only stateless transformations (like filter). If stateful transformations (like groupBy) are added to the code, an exception will be thrown.

```
spark.readStream
  .format("cloudFiles")
  .option("cloudFiles.format", "csv")
  .option("cloudFiles.schemaLocation", "<schema_location>")
  .option("cloudFiles.useIncrementalListing", "auto")
  .load("xfile_location")
  .withColumn("curr_timestamp", current_timestamp())
  .groupBy("Name").agg(sum("Number").alias("Sum"))
  .writeStream
  .option("checkpointLocation", "<new_checkpointlocation>")
  .trigger(availableNow=True)
  .table("table_move")
```

AnalysisException: Append output mode not supported when there are streaming aggregations on streaming DataFrames/DataSets without watermark;
Aggregate [Name#13119, [Name#13119, sum(cast(Number#13120 as double)) AS Sum#13136]
+= Project [Name#13119, Number#13120, _rescued_data#13121, current_timestamp#13125]
+= StreamingRelationDataSource(org.apache.spark.sql.SparkSession@3042d2fc,cloudfiles.List(),None,List(),None,Map(cloudFiles.format -> csv, cloudFiles.schemaLocation -> abfss://move@rahultestcodec.dfs.core.windows.net/autoloaderData/schema, cloudFiles.useIncrementalListing -> auto, path -> abfss://move@rahultestcodec.dfs.core.windows.net/data),None), cloudFiles, [Name#13119, Number#13120, _rescued_data#13121]

Command took 1.78 seconds -- by rahul.singha@databricks.com at 17/03/2023, 03:43:28 on test-rahul

[Autoloader](#)[Spark Streaming](#)[Databricks](#)[Data Engineering](#)[Data Engineering 101](#)**Rahul Singha**

24 followers · 4 following

[Follow](#)

Responses (1)

[Write a response](#)

What are your thoughts?

**Brett M**

May 6, 2023

...

Rahul,

This Databricks article was excellent about Autoloader with great examples. I hope you write more Databricks articles!

Thanks,

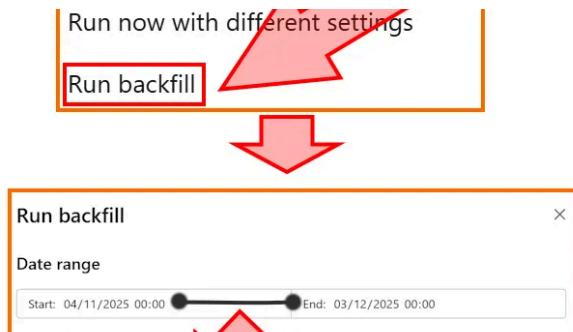
Brett



5

[Reply](#)

Recommended from Medium



Hubert Dudek

Databricks Lakeflow Jobs Workflow Backfill

In modern data pipelines, incremental processing is the standard approach for...

6d ago 9



Pallavi Sinha

Powering SQL with AI Functions in Databricks

Building Modern Data Platforms with Delta Lake



From Data Warehouses & Data Lakes to the Unified Lakehouse

In EndToEndData by Prem Vishnoi(cloudvala)

The Complete Guide To Lakehouse Architecture

1.2 Data Lakes: Scale Without Trust

Dec 30, 2025 75



Namaste Databricks

Why Databricks Lakehouse?

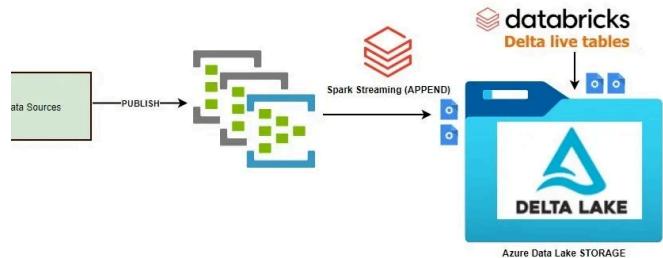
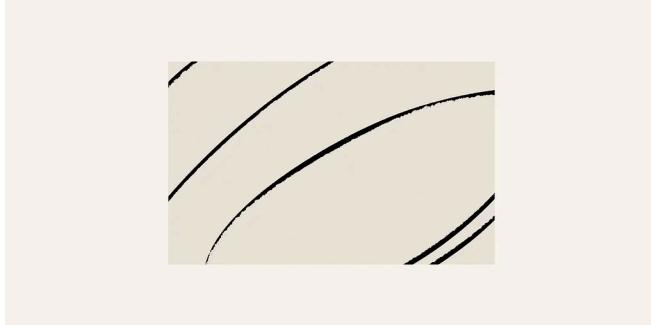
Storage is the foundation of all data systems, providing the means to store, organise and...

Data engineers and analysts spend much of their time in SQL—cleaning data,...

Oct 11, 2025



Jul 29, 2025 24 1



Umesh Pawar

Delta Lake Change Data Feed (CDF) Reduces ETL Costs and...

GitHub repository: umeshpawar2188/Delta-change-data-feed

Jul 25, 2025



Tharaniesh

Building a Real Time Lakehouse Data Platform (LDP) with Delta Li...

From Raw Streams to Business Insights: My Hands On Journey with Medallion...

Nov 18, 2025 1



[See more recommendations](#)