

Deep-Dive into Delta Lake



Gerhard Brueckl Cloud Data Architect @ paigo GmbH

ORGANIZED BY 😂 databricks









gerhard@gbrueckl.at

blog.gbrueckl.at



2

https://github.com/gbrueckl

DatabricksPS

@gbrueckl



Databricks VSCode



PowerBI Connector



www.paiqo.com





Agenda

- What is Delta Lake?
- What is the Delta Log?
- How does it Work?
- File & Storage Management
- Streaming
- Properties
- Conclusion & Lessions Learned

What is Delta Lake?



What is Delta Lake?

https://delta.io

<u>Delta Lake</u> is an open-source storage framework that enables building a <u>Lakehouse architecture</u> with compute engines including Spark, PrestoDB, Flink, Trino, and Hive and APIs for Scala, Java, Rust, Ruby, and Python.

- ACID compliant transactions
 - Optimistic Concurrency Control
- Support for UPDATE / MERGE
- Time-Travel

- Schema enforcement and evolution
- Batch & Streaming
- 100% compatible with Spark



What is Delta Lake?

https://delta.io

- Everything is stored in one folder
 - Meta-data
 - Transaction log / Delta Log
 - Data
- Could basically Copy & Paste whole Delta table
- Supports any storage sub-system
- Consumer only needs location





What is the Delta Log?



What is the Delta Log?

The Transactional Layer

- Contains
 - Table schema + changes
 - References to files
 - Metadata and metrics
- Stored as JSON and Parquet
- One file/version per transaction
- Allows [optimistic] concurrency control
- Used for time-travel, streaming, ...



What is the Delta Log?

DESCRIBE HISTORY

DESCRIBE HISTORY gold.my_big_table

(1) Spark Jobs

Table Data Profile

| | version 🔶 | timestamp 🔺 | userid 🔺 | userName 🔶 | operation 🔶 | operationParameters |
|---|-----------|------------------------------|----------|------------|--------------|---|
| 1 | 1185 | 2022-06-13T16:45:39.000+0000 | | | OPTIMIZE | {"predicate": "[]", "zOrderBy": "[]", "batchId": "0", "auto": "false"} |
| 2 | 1184 | 2022-06-13T16:18:24.000+0000 | | | VACUUM END | {"status": "COMPLETED"} |
| 3 | 1183 | 2022-06-13T16:18:19.000+0000 | | | VACUUM START | {"retentionCheckEnabled": "false", "defaultRetentionMillis": "259200000"} |
| 4 | 1182 | 2022-06-13T13:59:19.000+0000 | | | MERGE | {"predicate": "((target.purchase_id = updates.purchase_id) AND (target.iteration "[{\"predicate\":\"(((NOT (updates.store_fee_rate = target.store_fee_rate)) OR (NO target.store_fee_description))) OR ((NOT (updates.store_fee_absolute = target.sto (updates.net_sales_after_fees = target.net_sales_after_fees))))\",\"actionType\":\"up [{\"actionType\":\"insert\"}]"} |
| 5 | 1181 | 2022-06-13T13:56:44.000+0000 | | | MERGE | {"predicate": "(target.purchase_id = updates.purchase_id)", "matchedPredicates "notMatchedPredicates": "[]"} |



Processing of a simple Query





How does Delta Lake work?



DML Operations – UPDATE

| User | ProductNotebooPCTablet | Price k 900 € 1,500 € 500 € | UPDATE DimProduct SET Price = 1300 WHERE Product = 'PC' | ProductNotebookPCTablet | Price 900 € 1,300 € 500 € |
|------------|--|---|---|---|---|
| _delta_log | 0000.jso "add": { " <mark>pa</mark> | n 1 <mark>rt-01.parquet</mark> ",] | 0001.json "remove": { "path": "add": { "path": " <mark>pa</mark> | " <mark>part-01.parquet</mark> <mark>"t-02.parquet</mark> ", | ", }, } |
| Storage | Parquet part-01 (3 rows) DATA+AI SUMMIT 2022 | | Parquet Parquet part-01 (3 rows) (3 rows) | • | 12 |

DML Operations – DELETE

| | | Product | Price | | Product | Price | |
|---------|---------|--|--------|--|--|-------|--|
| Jser | | Notebook | 900€ | DELETE FROM DimProduct | Notebook | 900€ | |
| | | PC 1,300€ WHERE Product = 'PC' | Tablet | 500€ | | | |
| | | Tablet | 500€ | | | | |
| | | | | | - | | |
| 60 O | c | 0001.json | | 0002.json | | | |
| delta_l | 000.jso | "remove": { "path": " <mark>part-01.parquet</mark> ", }, "add": { "path": " <mark>part-02.parquet</mark> ", } | | arquet ["] , }, "remove": { "path": <mark>Jet</mark> ", } "add": { "path": " <mark>pa</mark> | "remove": { "path": " <mark>part-02.parquet</mark> ", }, "add": { "path": " <mark>part-03.parquet</mark> ", } | | |

| 0001.jsor | n |
|-----------|---|
|-----------|---|

0000.json

DATA+AI

SUMMIT 2022

Storage

0002.json



(3 rows) (3 rows)

part-02

part-01



13

DML Operations – INSERT

DATA+AI

SUMMIT 2022

| | | | | | | _ | |
|------------|------------------------|---|---------------------|--------|--|---------------------------------------|-----------------|
| | | Product | Price | | | Product | Price |
| <u>_</u> | | Notebook | 900€ | INSERT | INTO DimProduct | Notebook | 900€ |
| Jse | | Tablet | 500€ | VALUES | ('Monitor', 200) | Tablet | 500€ |
| | | | | | | Monitor | 200€ |
| | | | | - | | - | |
| _delta_log | 0000.json 0001.json | 0002.json "remove": { "path": "part-02.parquet", }, "add": { "path": "part-03.parquet", } | | | 0003.json "add": { "path": " <mark>pa</mark> i | rt-04.parquet", | . } |
| Storage | | Parquet Parquet Parquet (3 rows) (3 ro | -02 Ws) (2 rows) | | Parquet Parquet part-01 (3 rows) (3 rows) | Parquet part-03 (2 rows) (1 row | uet 04 7) |

DML Operations

- Operations are logged in _delta_log
 - Old files are **logically(!)** removed
 - New files are added
- Most operations create new files! Even a **DELETE** can!

Can create A LOT of files!



File & Storage Management



Data Management – VACUUM

| | | Product | Price | | | Product | Price |
|------------|--|--|------------------------------------|---|---|--|-----------------|
| Jser | | Notebook | 900€ | | | Notebook | 900€ |
| | | Tablet | 500€ | VACUUM DIMPROduct | Tablet | 500€ | |
| | | Monitor | 200€ | | | Monitor | 200€ |
| | | | | | | | |
| _delta_log | 0004.json 5 1 1 1 1 1 1 1 1 | | | 0005 et <mark>", }</mark> {"VA 0006 {"VA | 0005.json {"VACUUM START ", "numFilesToDelete": 2, 0006.json {"VACUUM END ", "numDeletedFiles": 2, } | | |
| Storage | DATA+AI SUMMIT 2022 | Parquet Parc part-O1 part (3 rows) (3 ro | Parquet Part-03 ws) (2 rows) | Parquet part-04 (1 row) | | Parquet part-O3 (2 rows) Parg part- (1 row | uet 04 /) |

Data Management – OPTIMIZE

| <u> </u> | | Product | Price | OPTIMIZE DimProduct | Product | Price |
|------------|--|----------|--------------------------------|-------------------------------|---------------------------------------|--|
| | | Notebook | 900€ | | Notebook | 900€ |
| Use | | Tablet | 500€ | | Tablet | 500€ |
| _ | | Monitor | 200€ | | Monitor | 200€ |
| | | | | | | |
| _delta_log | 0006.json 0007.json "VACUUM END ", "numDeletedFiles": 2, } "remove": { "path": "part-03.parquet", } "remove": { "path": "part-04.parquet", } "add": { "path": "part-05.parquet", } | | | | | <mark>et</mark> ", } <mark>et</mark> ", } } |
| Storage | DATA+AI SUMMIT 2022 | | Parquet part-03 (2 rows) | Parquet part-04 (1 row) | Parquet part-O3 (2 rows) (1 row | Auet -04 w) Parquet part-05 (3 rows) 18 |

VACUUM

- Physically removes unreferenced files older than X days
- Never touches files of latest version of Delta table!

VACUUM events [RETAIN num HOURS] [DRY RUN]

OPTIMIZE

- Collapse small files into bigger files
- Clustering / Ordering
- Improve query performance

```
OPTIMIZE events
[WHERE date = 20200101]
[ZORDER BY (eventType)]
```



VACUUM and OPTIMIZE

- VACUUM DRY RUN
 - Only shows first 1000 files to be deleted
 - Use SCALA to get the actual number of files to be removed!
- Can take a long time!

- OPTIMIZE
 - works per partition level
- Duplicates data!

| _ | | | | | | |
|---|---|--|--|--|--|--|
| | 1 2 | %scala spark.sql("VACUUM gold.my_big_table DRY RUN") | | | | |
| | (12) Spark Jobs | | | | | |
| | 🕨 🥅 res2: org.apache.spark.sql.DataFrame = [path: string] | | | | | |
| | Found 5888 files and directories in a total of 18531 directories t res2: org.apache.spark.sql.DataFrame = [path: string] | | | | | |
| | Com | mand took 1.85 minutes by gbrueckl@paiqo.com at 13/06/2022, 20:46:10 (| | | | |

RESTORE

- Restores a previous state of the Delta table
- At version or timestamp
- Meta-data only operation
- Creates a new version

RESTORE events TO TIMESTAMP AS OF '2022-05-03'

CLONE

- SHALLOW or DEEP
- Forks Delta Log
 - DEEP: copies data files
 - SHALLOW: references data files
- Ideal for testing

CREATE TABLE events_clone SHALLOW CLONE events;

RESTORE and **CLONE**

- You can RESTORE as often as you want
 - To rollback another RESTORE
- RESTORE does not create any new [data] files

• DEEP Clones are incremental and can be used for Backup





Basics

- Delta Tables can be partitioned
 - For ETL performance (usually on Bronze, Silver)
 - For query performance (usually on Gold)
- **Goal:** touch as few partitions as possible/necessary
 - ETL and Query performance can conflict
 - Explicitly specify Partitioning columns

• Partition by Time [and ?]





Advanced

- Avoid over-partitioning!
 - < few 1000s partitions
 - Single partition should be > 1 GB
- Use generated columns
 - EventTimestamp -> partition by CAST(EventTimestamp AS DATE)
 - Delta engine will [try to] push filters on EventTimestamp down to partition
- Used to separate transactions and processing jobs
 - Explicitly specify partitions you touch (e.g. **MERGE** target)!
 - Check Delta Log history for query predicates!



Advanced

 Physical .parquet file does not contain the partitioningcolumns! _delta_log Entry

```
"add": {
    "path": "SalesTerritoryKey=8/SalesDate=20220103/part-....
    "partitionValues": {
        "SalesTerritoryKey": "8",
        "SalesDate": 20220103
    },
    "size": 114365,
    "modificationTime": 1611740902000,
    "dataChange": true,
}
```

• **path** could point anywhere!

• You do not need to specify all partitioning columns sequentially!



Streaming



Streaming

Basics

- Delta Lake can be used as source and target for streaming
- It's technically still [micro-]batches
 - As is Spark Streaming
- Streaming works on a file-level
- Files are processed in order of
 - Version/Transaction number
 - File index (part-XXXXX_...snappy.parquet)



Streaming

Advanced

- Checkpoints
 - Track state of what has already been processed from source
 - One checkpoint per source
 - Could stream from same source multiple times using different checkpoints
- MERGE only with foreachBatch()
- Control the Trigger/Batch size!
- Avoid Trigger.Once
- Can stop/resume stream at any time

Delta Lake Table Properties



Delta Lake Table Properties

- Can be defined on different levels
- Table Properties
 - delta.autoOptimize.optimizeWrite
 - spark.databricks.delta.properties.defaults.optimizeWrite (default for new tables)
- Configured Settings during Execution
 - spark.databricks.delta.optimizeWrite.enabled
- Execution settings have priority over table properties!



Delta Lake Table Properties

Important Table Properties to know

- delta.appendOnly
- delta.autoOptimize.autoCompact
- delta.autoOptimize.optimizeWrite
- delta.deletedFileRetentionDuration
- delta.logRetentionDuration
- delta.dataSkippingNumIndexedCol



Delta Lake Table Properties

- Use defaults for commands
- Define exceptions on table level
- \rightarrow No need to use individual commands per table

• Changing table properties are also a Delta transaction



Conclusion & Lessons Learned



Conclusion

Take Aways

- Delta Lake can solve a lot of problems for you
- File management is crucial
- Data maintenance jobs are mandatory
- Use table properties



Conclusion

References

- The internals of Delta Lake by <u>Jacek Laskowski</u> <u>https://books.japila.pl/delta-lake-internals/</u>
- Delta Transaction Log Protocol <u>https://github.com/delta-io/delta/blob/master/PROTOCOL.md</u>



DATA+AI SUMMIT 2022

Thank you



Gerhard Brueckl

Cloud Data Architect @ paiqo.com