

Understanding the Internals of Delta Tables

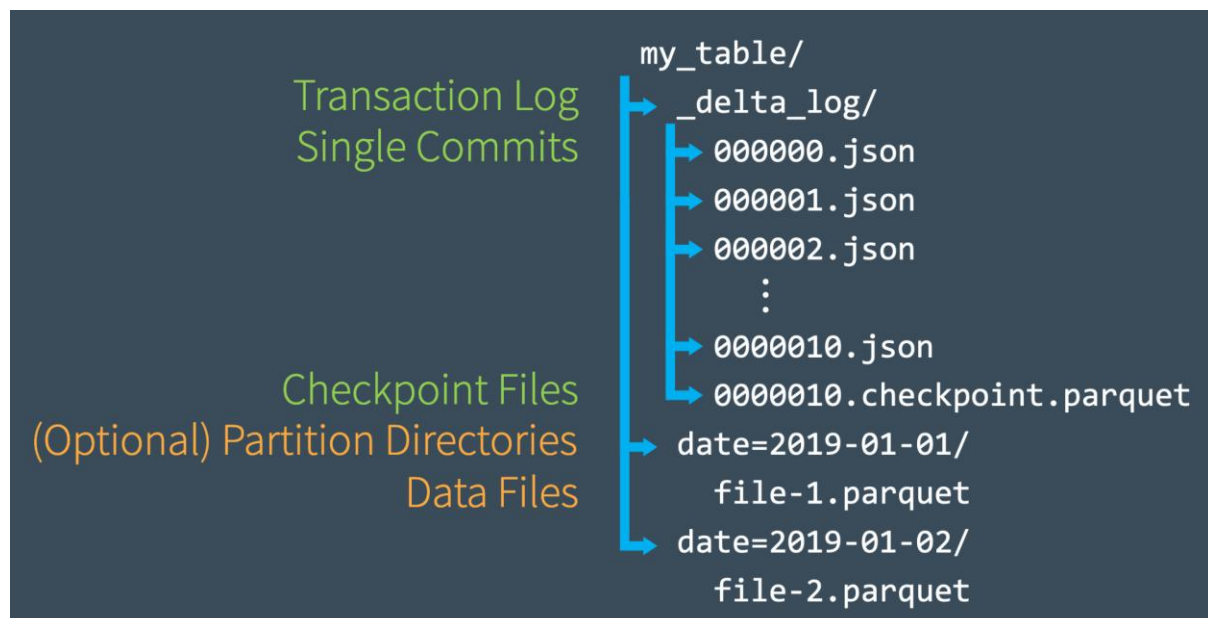
1. Introduction

Delta Lake is an open-source storage layer that brings **ACID transactions** to Apache Spark and big data workloads. It enables reliable data lakes by combining the scalability of data lakes with the reliability of databases. Delta Tables are the core storage units in Delta Lake, offering features such as **time travel**, **schema enforcement**, **upserts**, and **optimized query performance**.

2. Core Components of Delta Tables

2.1 Transaction Log (Delta Log)

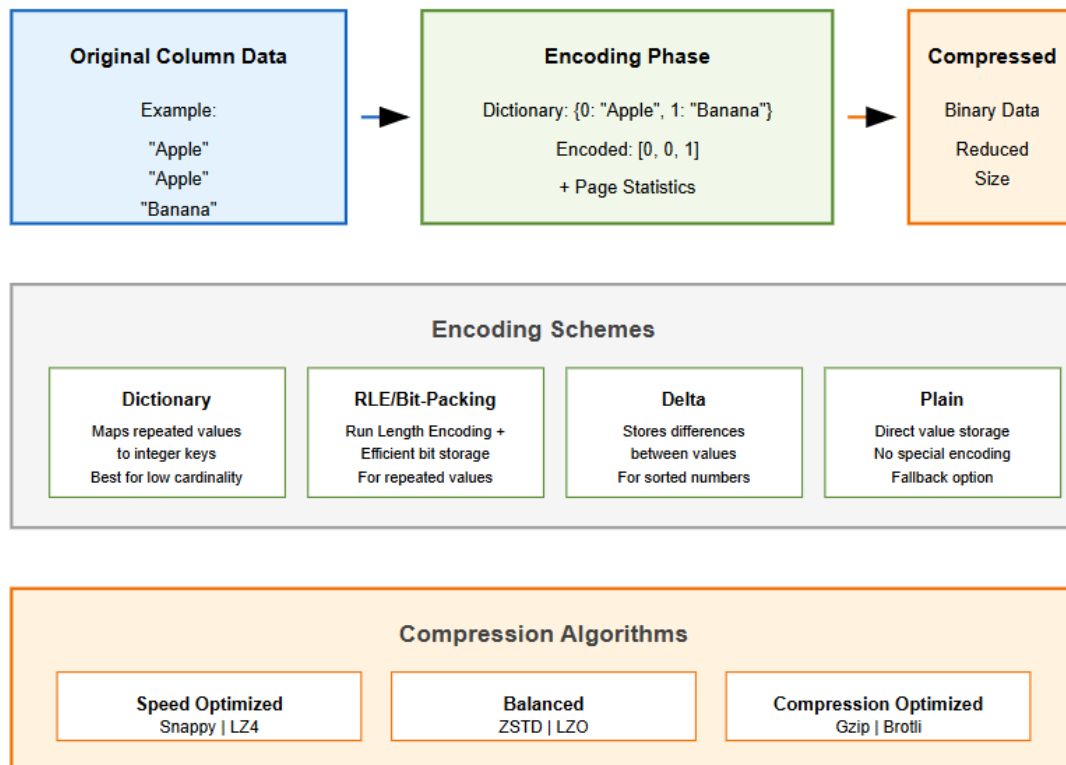
- The **transaction log** is the heart of a Delta Table.
- Stored in the **_delta_log** directory within the table's storage path.
- Maintains an **ordered record of all transactions** on the table.
- Each transaction creates a new JSON file (e.g., 000000.json, 000001.json).
- Ensures **atomicity** and **isolation** by recording:
 - Files added
 - Files removed
 - Schema changes
 - Operation details



2.2 File Organization

- Data is stored in **Parquet format**.
- **Immutable files**: Once written, files are never modified directly.
- Updates/deletes result in new files replacing old ones (**copy-on-write mechanism**).
- Older files are retained for **time travel** and historical queries.

Parquet Data Processing Pipeline



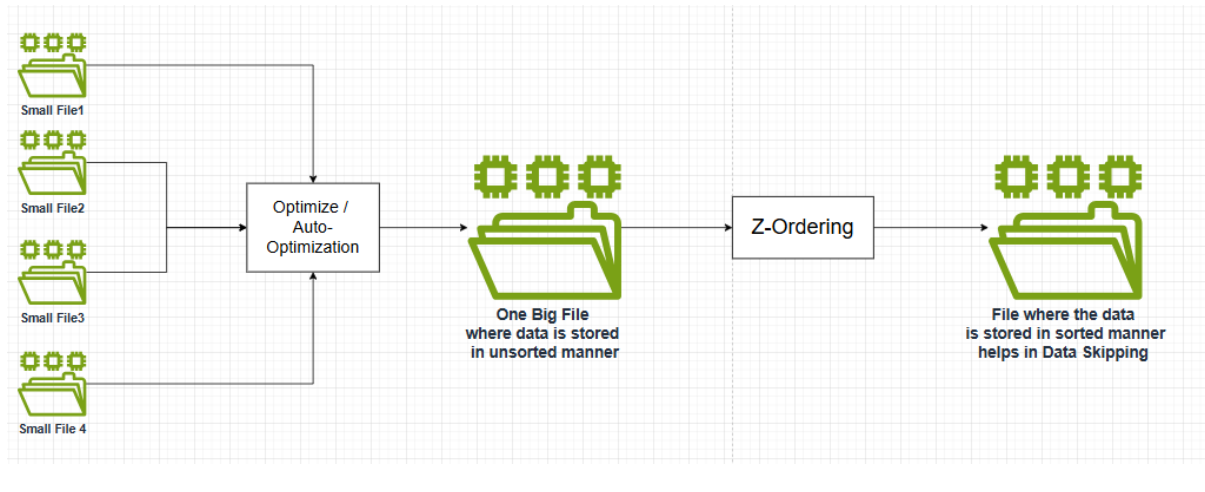
2.3 Metadata Management

- Schema information is stored in the transaction log.
 - Supports **schema evolution**:
 - Adding columns
 - Changing data types (with constraints)
 - Enforces **schema validation** to prevent invalid writes.
-

2.4 Optimizations

Delta Tables support various optimizations for performance:

- **Z-ordering:** Clustering data to improve data skipping during queries.
- **Compaction:** Merging many small files into larger ones for better read performance.
- **Data clustering:** Organizing data files based on query patterns.

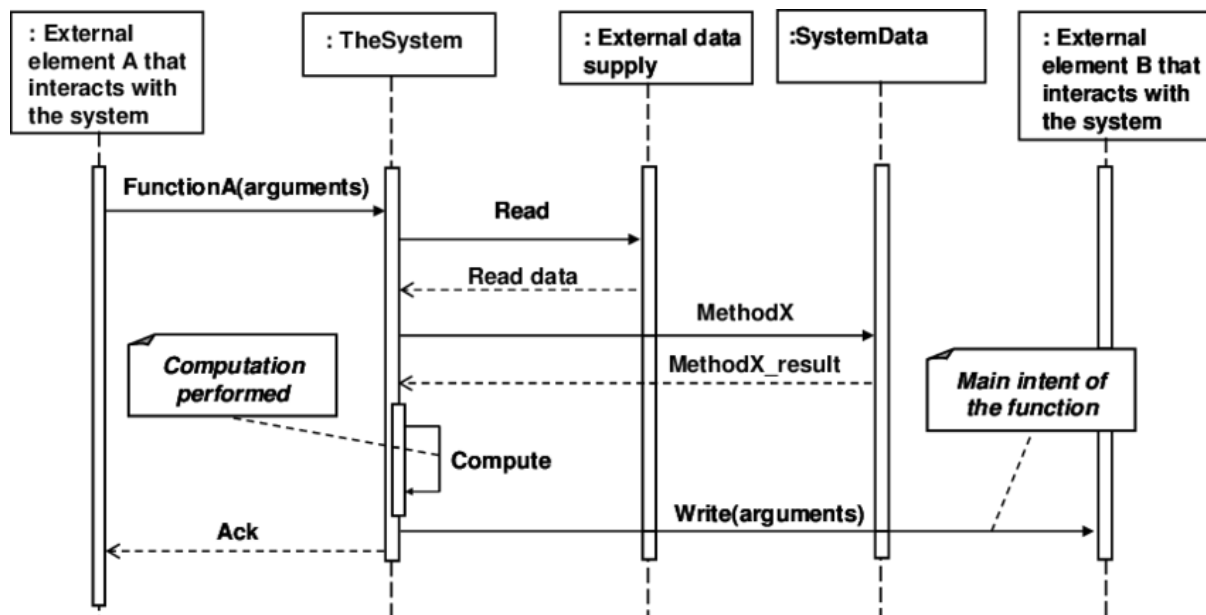


3. How Delta Tables Work

3.1 Write Operations

When data is written:

1. New **Parquet files** are created.
2. A new **transaction log entry** is added with:
 - Files added
 - Files removed
 - Schema changes
 - Operation metrics
3. For updates/deletes:
 - Files containing affected rows are identified.
 - They are rewritten with changes (**copy-on-write**).
 - Changes are recorded in the log.



3.2 Read Operations

When reading data:

1. The system checks the latest **transaction log**.
2. Determines the set of Parquet files for the current version.
3. Reads only relevant files (using statistics for data skipping).
4. Can access older versions using **time travel**.

4. Key Features Enabled by This Architecture

- **ACID Transactions:** Guaranteed by the transaction log.
 - **Time Travel:** Access to previous versions via log history.
 - **Upserts (MERGE):** Modify existing records or insert new ones.
 - **Scalable Metadata:** Handles billions of files efficiently.
 - **Audit History:** Complete record of all changes for compliance.
-

5. File Structure Example

```
/path/to/delta_table/  
  _delta_log/  
    000000.json  
    000001.json  
    000002.json  
    ...  
  part-00000-...-c000.snappy.parquet  
  part-00001-...-c000.snappy.parquet  
  ...
```

6. Conclusion

Delta Tables combine the **scalability** of data lakes with the **reliability** of databases. Their architecture, built around the transaction log, enables features such as **ACID guarantees**, **time travel**, and **optimized queries**, making them ideal for modern big data workloads.

7. Recommended Diagrams for the Assignment

1. **Delta Table Architecture Overview** – Showing Transaction Log, Parquet files, and Metadata.
2. **Write Operation Flow** – How updates create new files and update the log.
3. **Read Operation Flow** – How the latest log version determines the read set.
4. **Time Travel Mechanism** – How older versions are accessed.