

A apache Hudi:-free/open source

Vinod → ex-employee at es.

Hudi

→
Hadoop Upsert Delete
& incremental it'

- ① Streamline the real time data
- ② Manage SQL Table on your data lake to build the multi-stage incremental pipeline.

Hudi Feature:-

- ① Upsert, Delete with fast pluggable index.
- ② Transaction, Rollback with concurrency control good.
- ③ Automatic file sizing, data durability.
- ④ Built-in Metadata tracking for scalable storage access.
- ⑤ SQL can read/write from file, Spark, Presto
- ⑥ Schema comparison.

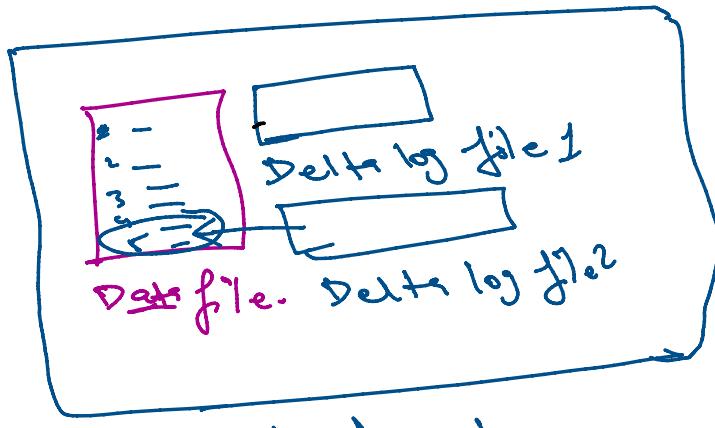
Hudi → MOR → Merge on Read.

Data file / Base file:-→ → Parquet file

Hudi store data in columns Parquet format. It is called

Data file Base file

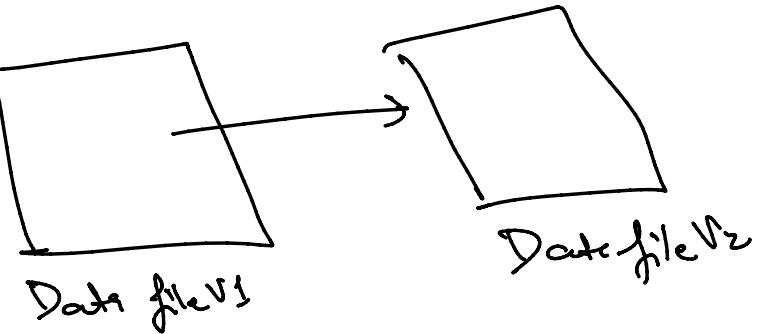
Deltas log file



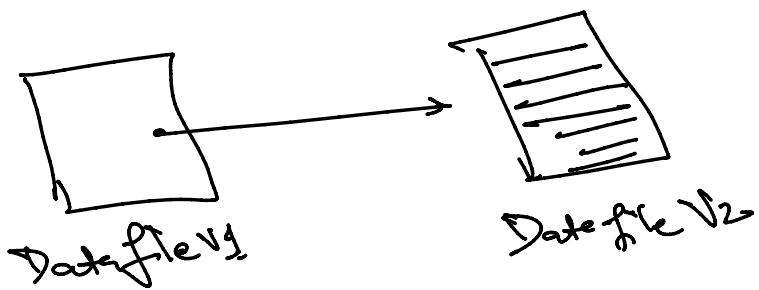
→ MOR Table format, updates are sent to delta log files, which is stored in Avro format. This Deltas log file is always associated with the base / Deltas log files.

file group

File Version :-



COW:- Copy on write

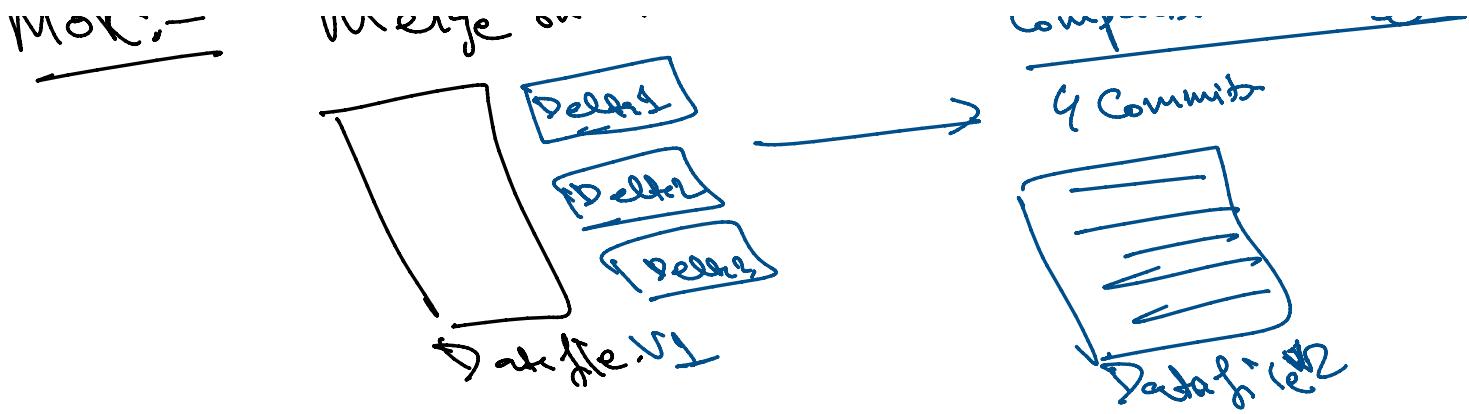


MOR:-

Merge on Read

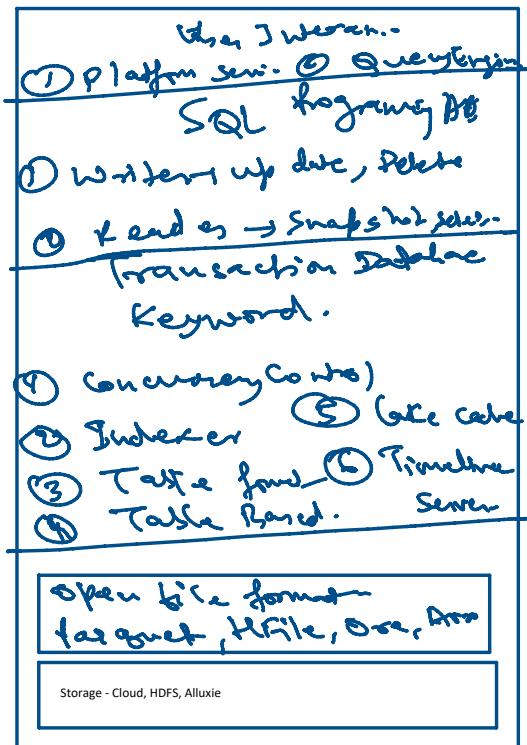
Rollup

Compaction strategy:-
→ 4 commits



- COW
- ① Higher Write latency
 - ② low latency read
 - ③ Cost high
I/O ↑

- MOR
- ① less COW, Synchronous Merge
 - ② high latency reads
 - ③ Cost less
I/O ↓



Lab:-

① S3 Bucket -

Database

1 S3 Bucket -

2 Glue → Spark Query → Database Table.

3- Athena Query the database -

→ it's real time → Commit, it will be available.

→ read optimized view
↳ Query