

Why We Need
Delta lake

Apache Spark & ACID:-

(2)

A \rightarrow Atomicity (Either All or Nothing)
[possible at job level in ~~append mode~~ but comes at performance impact]
df.write.mode("append").csv("/tmp/test-2")

Means, when using spark writer API, the operation should result in all the data being written (or) nothing written, but the save mode in spark do not utilize any locking and are not atomic.

1) Spark writer API internally relies on a job level commit protocol to ensure some degree of atomicity. (writes to a temp & then to the Outfile)
(temporary)

2) when performing overwrite, the data will be deleted before writing out the new data.

Ex: { for No impact? because of job level Atomicity in Spark?

> spark.range(100).repartition(2).write.
mode("overwrite").
csv("/tmp/test-2")

No Impact { > scala.util.Try(spark.range(100).repartition(1).map{i =>
if(i > 50) {
Thread.sleep(5000)
throw new RuntimeException("Ops!")
}
}.write.mode("append").csv("/tmp/test-2"))

(3)

Spark Job Commit Algorithm

For HDFS :-

FileOutputCommitter version 1 → Default

- ↳ ~~slow~~ slow & Impacts write performance
- ↳ Provide job level Atomicity

FileOutputCommitter version 2 → ^{Alternative Option}

- ↳ Provides ~~not~~ job level Atomicity
- ↳ Writes faster

For EMR :-

EMRFS S3-optimized committer

For S3 :-

Directory committer, Partitioned committer, Magic committer

For Databricks :- DBIO Commit Protocol

∴ Spark Commit Writer API's are not Atomic but mimics atomicity when in "append" mode. However, the Default Implementation comes with performance overhead, especially for Network intensive clusters.

For "Overwrite" spark is not Atomic has spark first Deletes the Data irrespective of job results.

Apache Spark & ACID

→ consistency (Data is always in the valid state)
[Not possible in Spark API]

```
df.write.mode("overwrite").csv("/tmp/test-2")
```

Delete + Write

↓
Time lag,
when Data is not consistent

Ex: Impacted, fails in the middle & Data is lost.

```
Spark.range(100).repartition(2).write.mode("overwrite")  
      .csv("/tmp/test-2")
```

```
scala.util.Try(Spark.range(100).repartition(10).map { i =>  
  if (i > 50) {  
    Thread.sleep(5000)  
    throw new RuntimeException("Oops!")  
  }  
  i  
}.write.mode("overwrite").csv  
("/tmp/test-2"))
```

Due to time lag between Delete & Write operations & prior Deletion of previous Data before Write operation ⑤

Spark Overwrite job writer is not consistent & hence the Data is ~~big~~ Hadoop can only attain eventual consistency.

Apache Spark & ACID :-

⑥

I - Isolation (an operation must be isolated from other concurrent operations)
[Spark doesn't offer isolation even at an API level.]

Spark doesn't offer isolation, which is typically done by transaction level commits in RDBMS.

Spark however offers task level & job level commits.

So, users cannot read a Data-set which a different process is overwriting.

Apache Spark & ACID :-

(7)

D - Durability (once committed data is never lost)

[offers via partition tolerance.]

Durability is offered by the storage layer.

Why ACID is Critical :-

In event of runtime exception data is corrupted or lost.
faulty analytics due to inconsistent data.

Summary:- Spark is a ~~storage engine~~ Processing engine using typically HDFS for storage, YARN/Mesos/Kubernetes for cluster Resource Management & Hive metastore for Database objects creation.

Using Hive metastore is the main reason spark was unable to provide ACID features required for reliable Data processing.

Other problems: —

(8)

Missing schema enforcement.

- Schema on Read Mechanism in Hadoop.
- Spark doesn't support Automatic Schema evolution.
- Data Separation is impossible to identify the corrupted Data.

Small File problems (<128MB)

- File Listing slows down the job performance significantly.
- File Opening/Closing.
- Reduced compression effectiveness due to file overhead.
- Excessive Metadata.

Multish data skipping using partition.

→ using where clause in SQL.

→ Work only for:

- Chronological columns (Ex: Date)
- Low cardinality columns (Ex: Country Code)
- Combination of Chronological & Cardinality columns.

→ lacks indexing