# Delta Lake

# Delta Lake

▶ Open-source storage framework that brings reliability to data lakes

# Delta Lake is/is not

**Is**

▶ Open-source technology
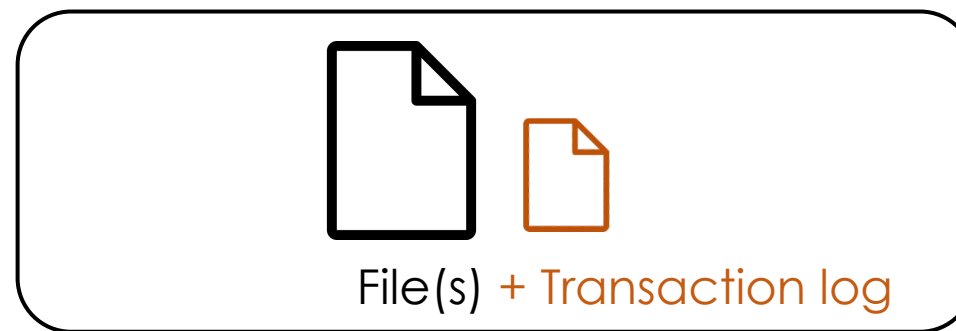
▶ Storage framework/layer

▶ Enabling building Lakehouse

**Is Not**

▶ Proprietary technology

▶ Storage format/medium

▶ Data warehouse/Database service

**Cluster**

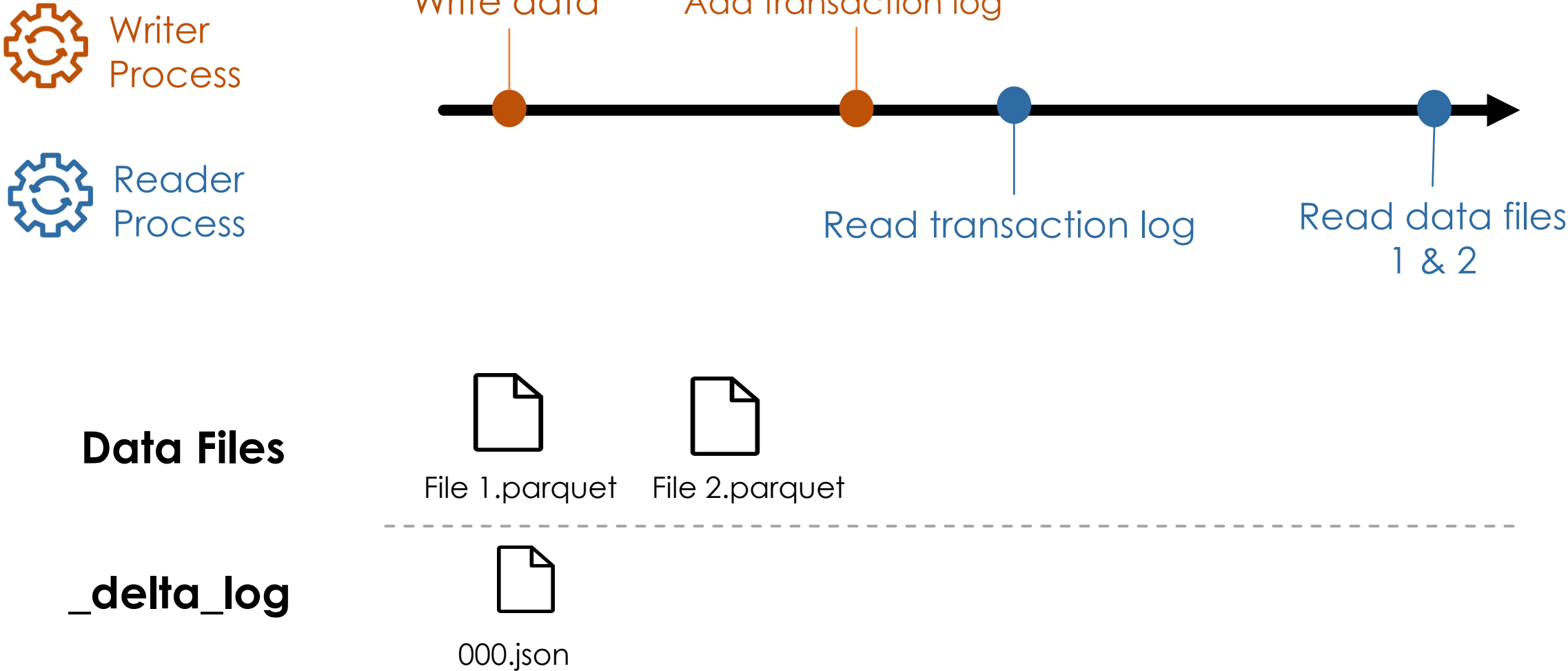Delta Table

DELTA LAKE

**Storage**

File(s) + Transaction log

# Transaction log (Delta log)

▶ Ordered records of every transaction performed on the table

▶ Single Source of Truth

▶ JSON file contains commit information:

   ▶ Operation performed + Predicates used
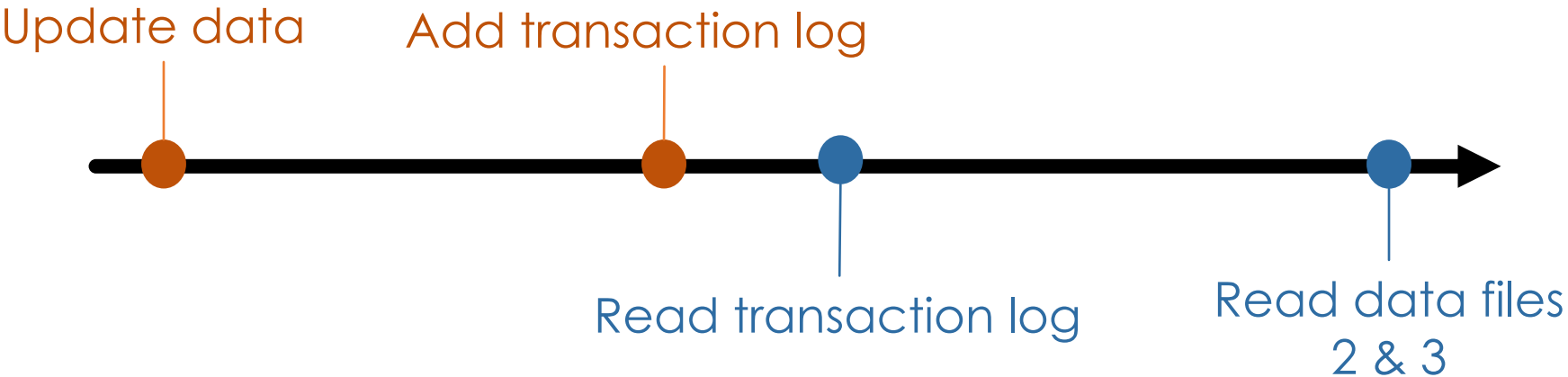
   ▶ data files affected (added/removed)

# Updates

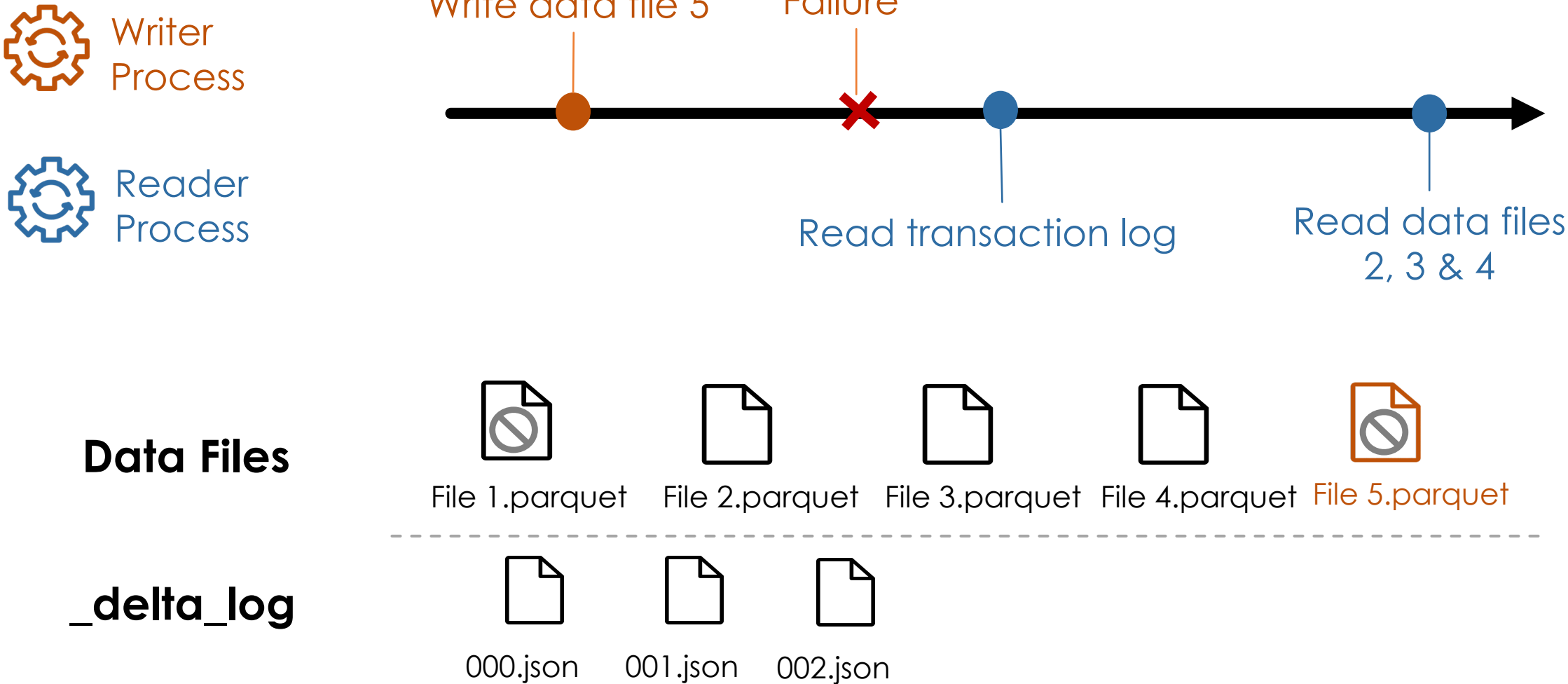# Simultaneous Writes/Reads

# Delta Lake Advantages

▶ Brings ACID transactions to object storage

▶ Handle scalable metadata

▶ Full audit trail of all changes

▶ Builds upon standard data formats: Parquet + Json