

What is an open table format? & Why to use one?

with Apache Iceberg

Joseph Machado

Introduction

- Name: Joseph Machado
- Work experience: Data engineer
- Specialty: Data Platform & Pipelines
- Blogs at: [Startdataengineering](#)



Table of Contents

1. Questions to address
2. How do OTFs work?
3. Features
 1. Alter data and partition schema without reprocessing
 2. Time travel
 3. git-like branching of your data
 4. Handle multiple reads/writes concurrently
4. Recap

Questions to address

- What is an open table format(OTF)? & How is it different from file formats like Parquet or ORC?
- What are the pros and cons of using an open table format? & How does it work?
- Is it just a pointer to some metadata files and helps you sift through the data quickly?

By the end of this presentation, you will know what OTFs are, why you use them, and how they work.

TL;DR

Open table formats are wrappers around your data store & use a series of files to

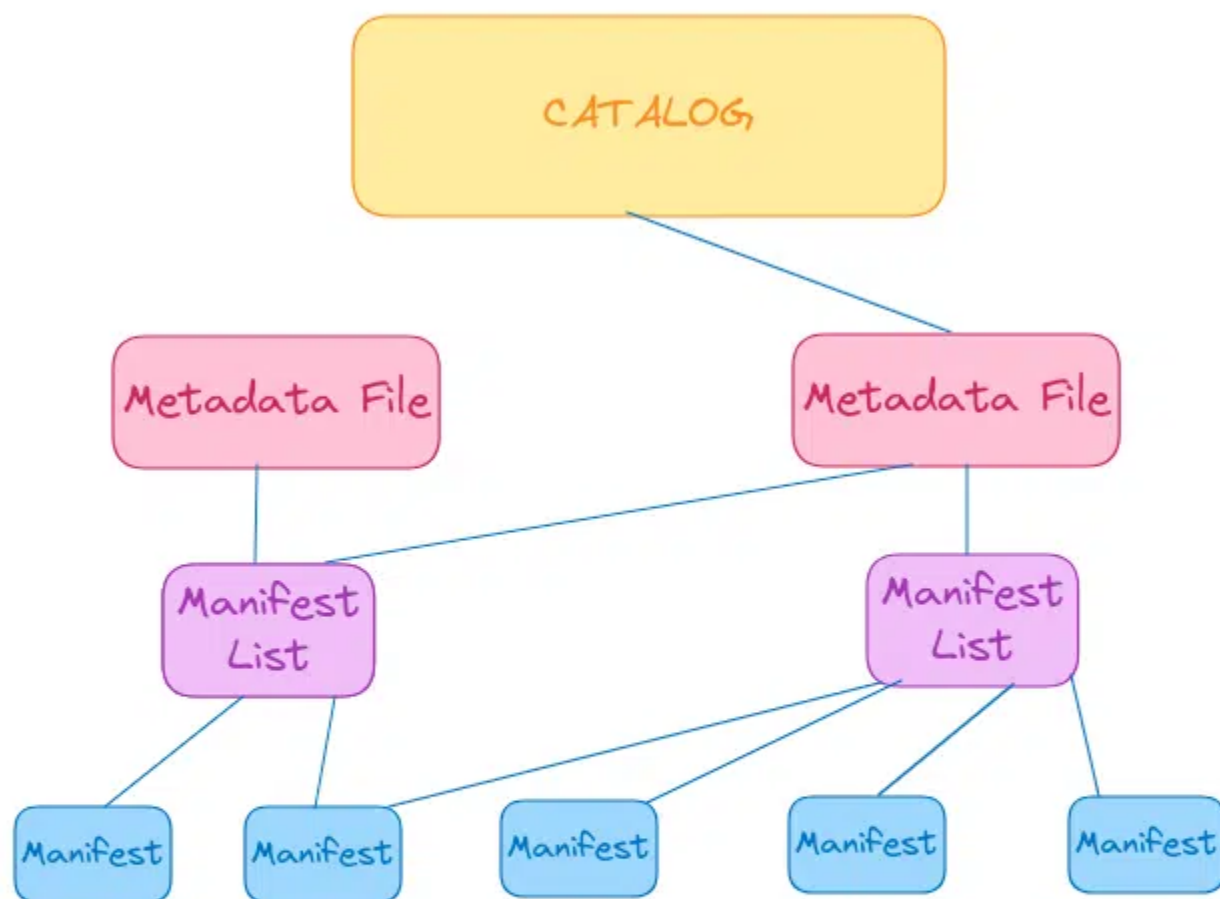
1. Track schema/partition (DDL) changes on your table.
2. Track the table's data files & their column statistics.
3. Track all the Inserts/Updates/Deletes (DML) on your table.

Storing a chronological series of files with all the DDL and DML statements applied to your table & index of the data file locations enables

1. **Schema & Partition Evolution**
2. **Travelling back in time to a previous table state**
3. **Creating table branches & tagging table state (similar to git)**
4. **Handling multiple reads & writes concurrently**

How do OTFs work?

- OTFs are implemented as libraries that handle table management using metadata files. [E.g., Apache Iceberg Java library](#)
- Apache Iceberg uses
 - **Metadata** files to keep track of data/partition schema changes.
 - **Manifest_list** contains partition statistics.
 - **Manifest** files contain column-level statistics and links to the actual data.



1. Check catalog for latest table metadata file

2. Get partitioning, schema and snapshot info.

3. check manifest list to prune based on partition stats

4. check manifests to prune files based on column stats

5. scan files that were not pruned to execute query

```

1 CREATE SCHEMA warehouse;
2 -- Use Apache Iceberg table format
3 DROP TABLE IF EXISTS local.warehouse.orders;
4 CREATE TABLE local.warehouse.orders (
5     order_id string,
6     cust_id INT,
7     order_status string,
8     order_date timestamp
9 ) USING iceberg
10 PARTITIONED BY (date(order_date));
11
12 INSERT INTO local.warehouse.orders VALUES
13 ('e481f51cbdc54678b7cc49136f2d6af7',69,'delivered',CAST('2023-11-01 09:56:3
14 ('e481f51cbdc54678b7cc49136f2d6af7',87,'delivered',CAST('2023-11-01 10:56:3
15 ('e481f51cbdc54678b7cc49136f2d6af7',125,'delivered',CAST('2023-11-01 11:56:
16 ('53cdb2fc8bc7dce0b6741e2150273451',17,'delivered',CAST('2023-11-02 11:56:3
17 ('53cdb2fc8bc7dce0b6741e2150273451',19,'on_route',CAST('2023-11-02 12:56:33
18 ('47770eb9100c2d0c44946d9cf07ec65d',26,'on_route',CAST('2023-11-03 12:56:33
19 ('47770eb9100c2d0c44946d9cf07ec65d',88,'lost',CAST('2023-11-03 12:56:33',AC

```


Files generated as part of our Orders Iceberg table

data/iceberg-warehouse/warehouse/orders

data

- order_date_day=2023-11-01
- order_date_day=2023-11-02
- order_date_day=2023-11-03
- order_date_day=2023-11-04
- order_date_day=2023-11-05
- order_date_day=2023-11-06

metadata

- ddd06a89-b660-4862-a4e5-fbe070b22e75-m0.avro
- snap-3638368400097931127-1-ddd06a89-b660-4862-a4e5-fbe070b22e75.avro
- v1.metadata.json
- v2.metadata.json
- version-hint.text

Manifest

Metadata file

Manifest list

**ALTER data and
partitions without
reprocessing**

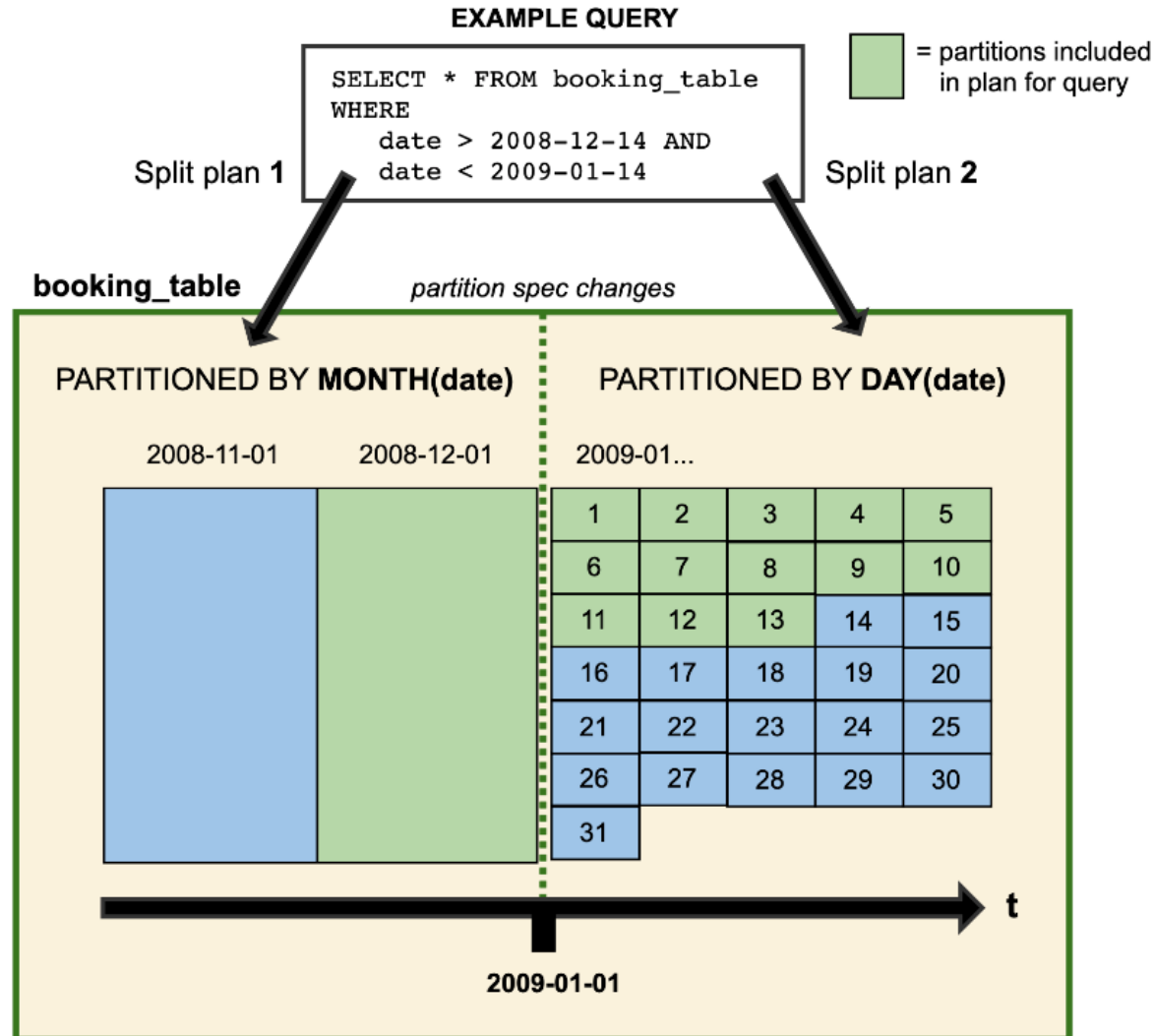
Data Schema & Partition Evolution

- Apache Iceberg enables you to **change your table's data schema or partition schema without reprocessing the existing data.**
- Metadata files track schema & partition changes over time.
- Processing systems use metadata to understand current and past schemas and apply the appropriate schema.

For example:

```
1  -- schema evolution
2  ALTER TABLE local.warehouse.orders ALTER COLUMN cust_id TYPE bigint;
3  ALTER TABLE local.warehouse.orders DROP COLUMN order_status;
4
5  -- parititon evolution
6  ALTER TABLE local.warehouse.orders ADD PARTITION FIELD cust_id;
7  INSERT INTO local.warehouse.orders VALUES
8  (
9      'e481f51cbdc54678b7cc49136f2d6af7',
10     69,
11     CAST('2023-11-14 09:56:33' AS TIMESTAMP)
12 ),
13 (
14     'e481f51cbdc54678b7cc49136f2d6af7',
15     87,
16     CAST('2023-11-14 10:56:33' AS TIMESTAMP)
17 );
18
19 -- check snapshots
```

Reading data which has undergone partition evolution



Hidden partition

- Typically, you would have had to create a separate column to add a partition.
- Apache Iceberg lets you define(& enable the use of) partitions based on column transformations of another column. E.g., you can partition by date on a timestamp column.

```
1  -- The partition is done on a transformation of an existing col
2  CREATE TABLE local.warehouse.orders (
3      order_id string,
4      cust_id INT,
5      order_status string,
6      order_date timestamp
7  ) USING iceberg
8  PARTITIONED BY (date(order_date));
9
10 -- The below query automatically uses the partition to
11 -- prune data files to scan.
12 SELECT cust_id, order_date
13 FROM local.warehouse.orders
14 WHERE order_date
15 BETWEEN '2023-11-01 12:45:33' AND '2023-11-03 12:45:33';
```

The **day** transform function is applied to the **order_date** column for a partition.

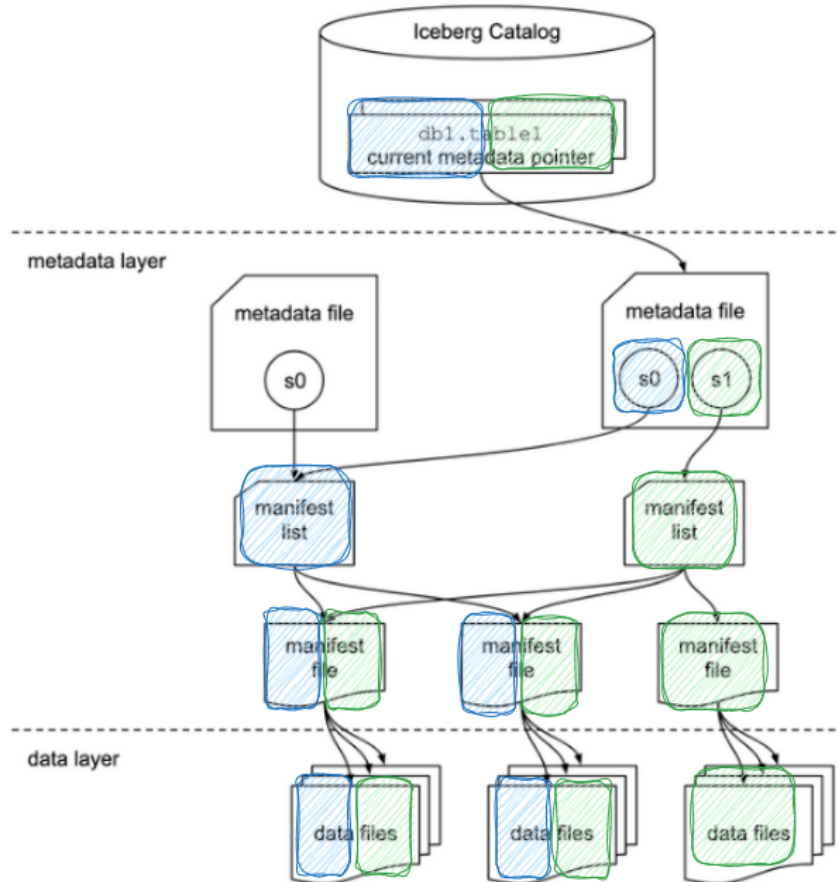
```
"partition-specs" : [ {  
  "spec-id" : 0,  
  "fields" : [ {  
    "name" : "order_date_day",  
    "transform" : "day",  
    "source-id" : 4,  
    "field-id" : 1000  
  } ]  
}
```


See previous point-in-time table state, aka time travel

- We can see our table at a previous point-in-time table state(aka time travel). For example, If you want to see what a table looked like four days ago, you can do this using time travel.
- Time travel is possible because the metadata files keep track of all the DML and DDL operations that ran on the table.
- Apache Iceberg keeps links to the deleted data; this enables the user to see deleted data.

```
1 -- get the time of the first data snapshot
2 select min(committed_at) as min_committed_at
3 from local.warehouse.orders.snapshots;
4 -- e.g. 2023-11-21 12:03:08.833
5
6 -- Query data as of the oldest committed_at (min_committed_at from the above)
7 SELECT *
8 FROM local.warehouse.orders TIMESTAMP AS OF '2023-11-21 12:03:08.833';
9 -- 15 rows
10
11 -- Query without time travel, and you will see all the rows
12 SELECT * FROM local.warehouse.orders;
13 -- 17 rows
```

Note: For Apache Iceberg, as your data size grows, it's recommended to clean up old snapshots.



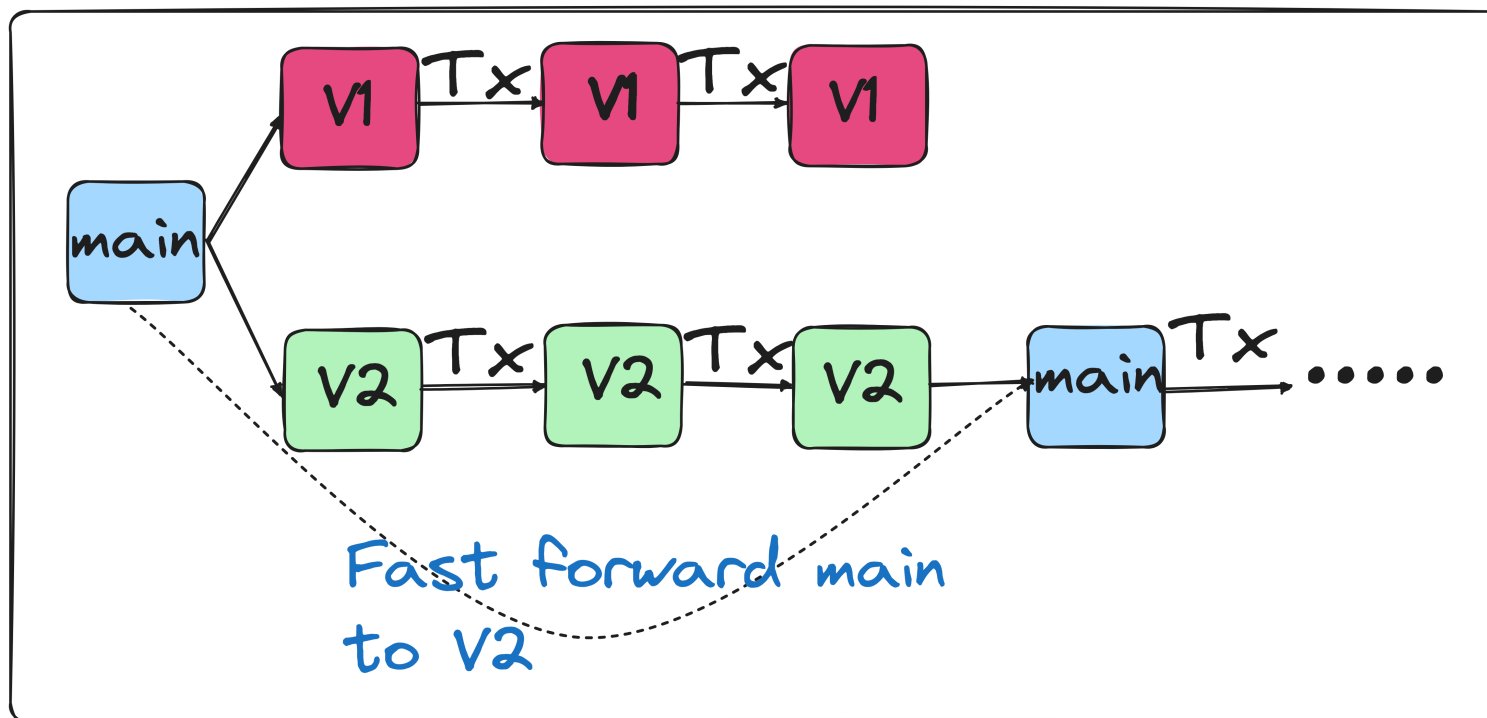
■ Data accessed when we run:
`Select * from table`

■ Data accessed when we run:
`Select * from table as of version s0`

Apache Iceberg Time Travel

Git like branches & tags for your tables

- Apache Iceberg enables the creation of table branches by managing isolated metadata files per branch.
- For example, assume you have to run a pipeline in production for a week to validate that the changes you made are valid. You can do this using branches as shown below:



Create a simple aggregate table that contains aggregated order data.

```
1 DROP TABLE IF EXISTS local.warehouse.orders_agg;
2 CREATE TABLE local.warehouse.orders_agg(
3     order_date date,
4     num_orders int
5 ) USING iceberg;
6
7 INSERT INTO local.warehouse.orders_agg
8 SELECT date(order_date) as order_date,
9 count(order_id) as num_orders
10 from local.warehouse.orders
11 WHERE date(order_date) = '2023-11-02'
12 GROUP BY 1;
```

Create two branches that are both stored for ten days.

```
1 ALTER TABLE local.warehouse.orders_agg
2 CREATE BRANCH `branch-v1` RETAIN 10 DAYS;
3 ALTER TABLE local.warehouse.orders_agg
4 CREATE BRANCH `branch-v2` RETAIN 10 DAYS;
5
6 -- Use different logic for each of the branches
7
8 -- inserting into branch v1
9 INSERT INTO local.warehouse.orders_agg.`branch_branch-v1`
10 SELECT date(order_date) as order_date
11 , count(order_id) as num_orders
12 from local.warehouse.orders
13 WHERE date(order_date) = '2023-11-03'
14 GROUP BY 1;
15
16 INSERT INTO local.warehouse.orders_agg.`branch_branch-v1`
17 SELECT date(order_date) as order_date
18 , count(order_id) as num_orders
19 from local.warehouse.orders
```

We noticed that `branch-v2` has the correct logic, so we fast-forward the main branch to `branch-v2`. The main branch will now have the accurate data for the past two days.

```
1 select *
2 from local.warehouse.orders_agg
3 order by order_date;
4
5 -- Push the main branch to branch v2's state
6 CALL local.system.fast_forward('warehouse.orders_agg', 'main', 'branch-v2')
7
8 select * from local.warehouse.orders_agg order by order_date;
```

Handle multiple reads & writes concurrently

- In traditional OLAP systems (e.g., HIVE) or when creating raw data on cloud storage, if multiple processes read/write to the same table without proper safeguards, there may be inconsistent data reads, or data may get overwritten during writes.
- Apache Iceberg atomically updates its metadata, which forces writers to “commit” their changes one at a time (if multiple writers collide, there will be a [retry for the failed writer](#)).
- When reading data, Apache Iceberg uses the most recent snapshot (using a metadata file) to ensure that no in-process data operations impact reads.

- Since Apache Iceberg is OSS, we can use any system that implements the table format to be able to read and write.
- E.g., We can use DuckDB to read our data:

```
1  INSTALL iceberg;
2  LOAD iceberg;
3
4  -- Count orders by date
5  WITH orders as (
6      SELECT *
7      FROM
8          iceberg_scan(
9              'data/iceberg-warehouse/warehouse/orders',
10             ALLOW_MOVED_PATHS=true
11         )
12     )
13  select strftime(order_date, '%Y-%m-%d') as order_date
14  , count(distinct order_id) as num_orders
15  from orders
16  group by strftime(order_date, '%Y-%m-%d')
17  order by 1 desc;
```

Recap

Open table formats are wrappers around your data store & use a series of files to

1. Track schema/partition (DDL) changes on your table.
2. Track the table's data files & their column statistics.
3. Track all the Inserts/Updates/Deletes (DML) on your table.

Storing a chronological series of files with all the DDL and DML statements applied to your table & index of the data file locations enables

1. **Schema & Partition Evolution**
2. **Travelling back in time to a previous table state**
3. **Creating table branches & tagging table state (similar to git)**
4. **Handling multiple reads & writes concurrently**

Available now

Public preview Apache Iceberg tables on Snowflake !!!