LET'S GO HYBRID

**WORKSHOPS**
APRIL 26th, 2022 ONSITE ONLY

**CONFERENCE**
APRIL 27th, 2022 ONLINE OR ONSITE
APRIL 28th, 2022 ONLINE

www.bigdatatechwarsaw.eu

bigdata
TECHNOLOGY WARSAW SUMMIT

# Scaling Your Data Lake w/ Iceberg

- Victoria Bukta (Shopify)

# Victoria Bukta

- Based in Toronto

- Senior Data Platform Eng

- At Shopify for 4.5 years

  - Toronto & Berlin offices

- Lakehouse (formally Data Acquisition)

- Hobbies

  - Field Hockey

  - Sailing

  - Backcountry Camping

# Agenda

- **Context**
- **Problem**
- **Solution**
  - **What is Apache Iceberg?**
  - **Promise (V2 Spec)**
- **Result**
- **Reflection**
- **Future Challenges**
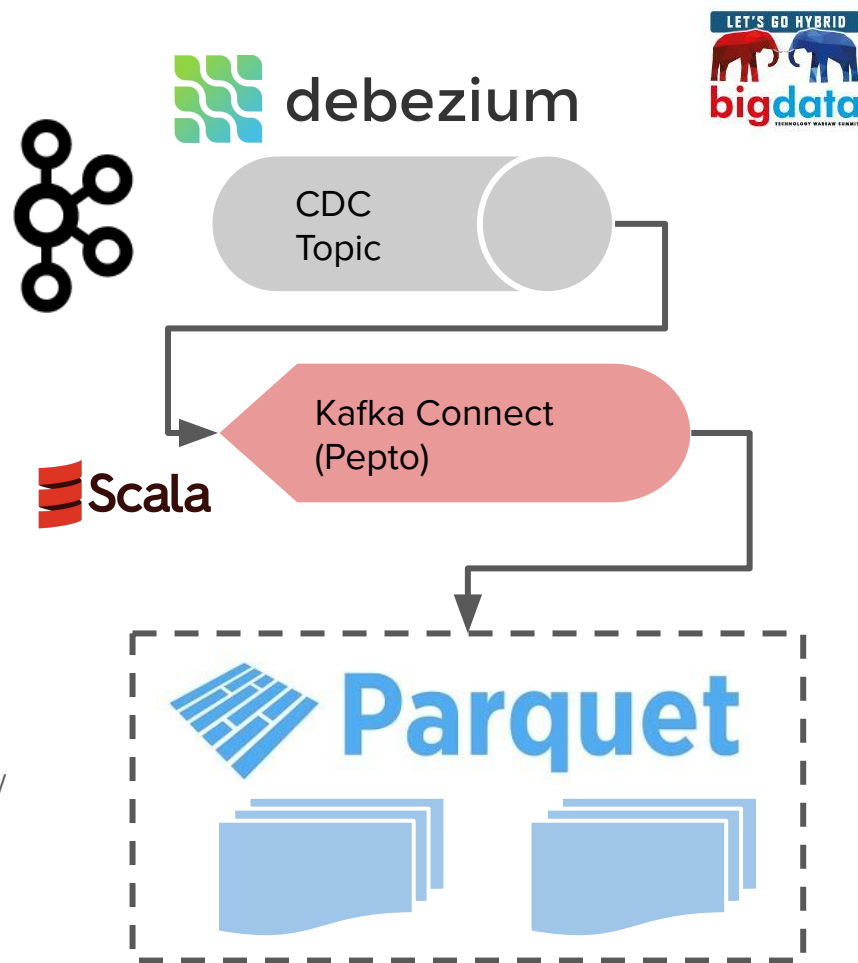
# Context

- New Kafka ingestion tooling is being built to support CDC (change data capture) use case (Pepto)
  - **Streaming ingestion**
  - 15 min SLA
  - **Anticipated huge table, trillions of rows**
  - Columnar schematized datasets
  - Time series data
  - Aggressive schema evolutions
  - **Future use case of supporting Type-1 tables**
    - **new data overwrites the existing data**
- **Kafka Connect application** because of internal support / expertise at scale
- Require read support from Spark, Trino, Flink

# Problem

1. **Transactional Semantics**

2. **Fast upsert to support Type-1 tables**
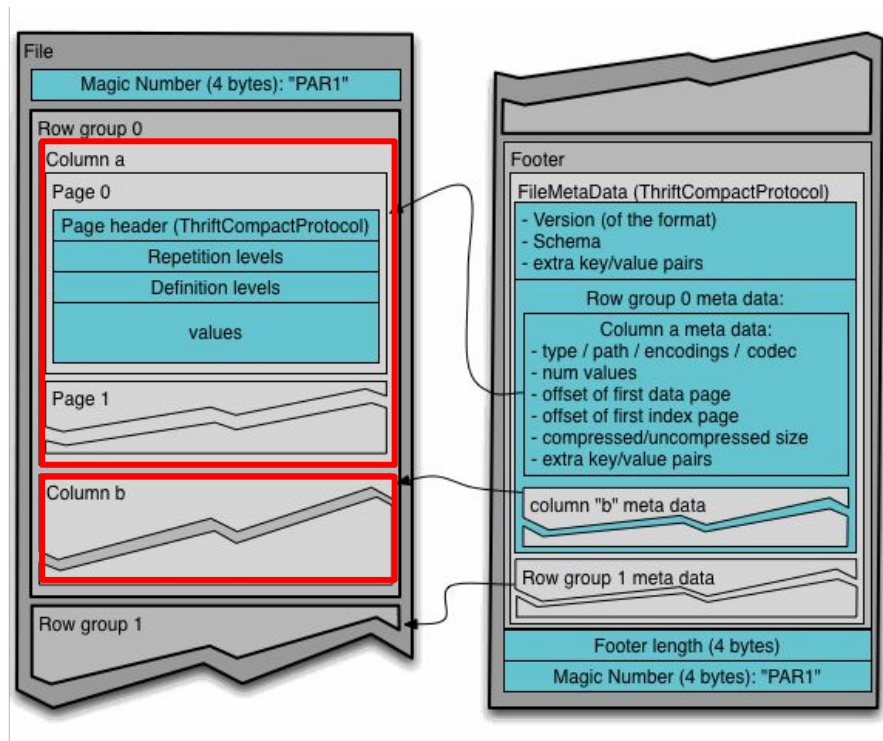
# Problem - Transactional Semantics

- **Modeling tools currently tied to our writing implementation**
  - HDFS vs Object Store **(NOT THE SAME)**
    - FS abstraction is missing
      - Atomic move, rename
  - Timestamp folders on GCS
    - Makes it hard to do maintenance tasks without effecting customers
    - Data scientists refer to datasets by folder location
- **Difficult to innovate** when implementation details are exposed

# Problem - Transactional Semantics

- **Modeling tools currently tied to our writing implementation**
    - HDFS vs Object Store **(NOT THE SAME)**
        - FS abstraction is missing
            - Atomic move, rename
    - Timestamp folders on GCS
        - Makes it hard to do maintenance tasks without effecting customers
        - Data scientists refer to datasets by folder location
- **Difficult to innovate** when implementation details are exposed

```
Buckets > data > shopify > shops
  -> 202204211112/
    -> .metadata
    -> part-0000.parquet
    -> part-0001.parquet
    -> ...
  -> 202204221112/
  -> 202204231112/
  -> ...
```

# Problem - Fast upsert



- **Storing data in columnar format**
  - Efficient compaction of schematized data
  - Optimizing for aggregation analytics over a subset of columns
- **Creating Type-1 dimensions is hard**
  - Columnar files are immutable
  - Rewrite is an expensive operation
  - People want their data **NOW**

# Solution

# Solution - What is Apache Iceberg?

- **Iceberg is a table format**
  - **Just a library**
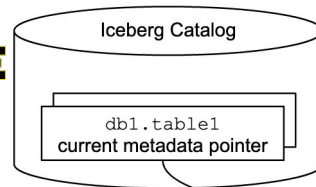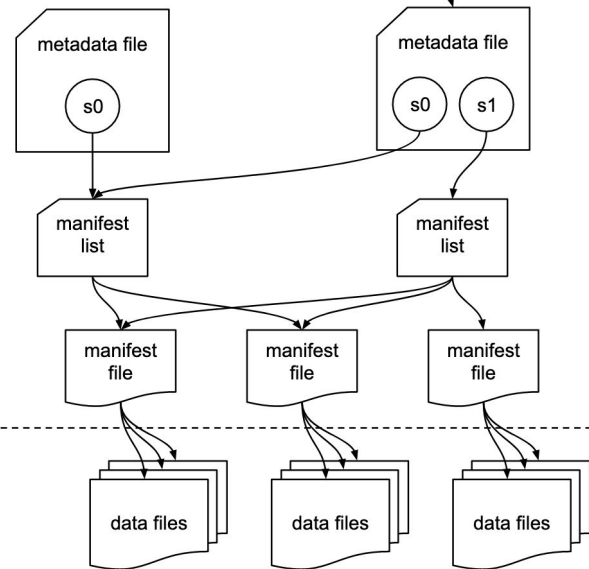  - Contents of a table are identified by traversing through metadata files

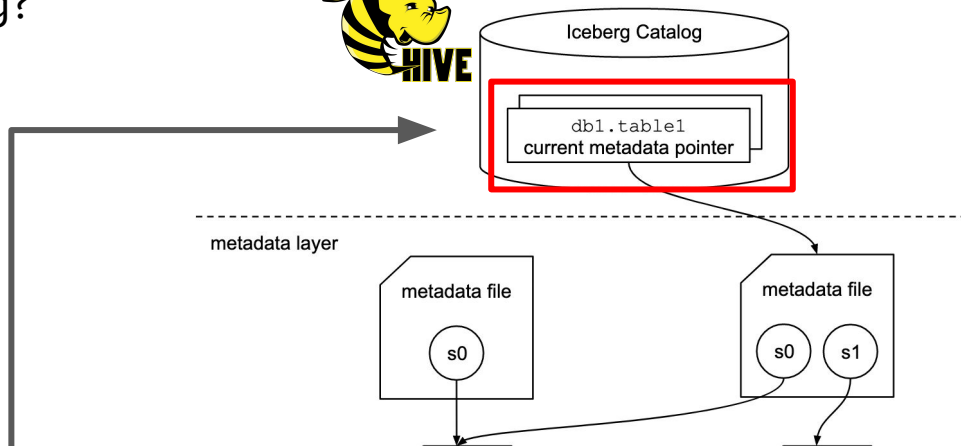**File system**

**Hard Drive**

**Iceberg**

**Object Storage**

# Solution - What is Apache Iceberg?

- **Catalog stores a pointer to a metadata file**

  - This files acts as a ledger

  - Has schema information

  - Has partition information
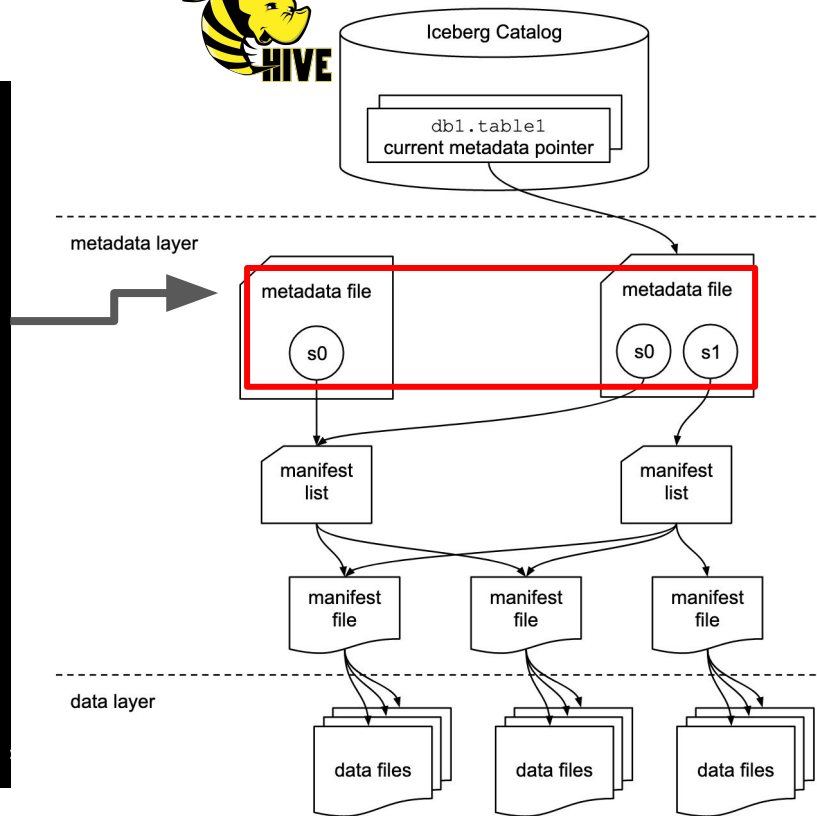
  - Gives us atomic commits



Iceberg Catalog

db1.table1
current metadata pointer

metadata layer

metadata file

s0

metadata file

s0   s1

```
+---------+------------------------+---------------------------------------------------------------+
| TBL_ID  | PARAM_KEY              | PARAM_VALUE                                                   |
+---------+------------------------+---------------------------------------------------------------+
|     292 | EXTERNAL               | TRUE                                                          |
|     292 | metadata_location      | gs://my_bucket/hive-warehouse/table/metadata/00001.metadata.json |
|     292 | numFiles               | 1                                                            |
|     292 | previous_metadata_location | gs://my_bucket/hive-warehouse/table/metadata/00000.metadata.json |
|     292 | table_type             | ICEBERG                                                      |
|     292 | totalSize              | 1624                                                        |
|     292 | transient_lastDdlTime  | 1610742932                                                  |
+---------+------------------------+---------------------------------------------------------------+
7 rows in set (0.02 sec)
```
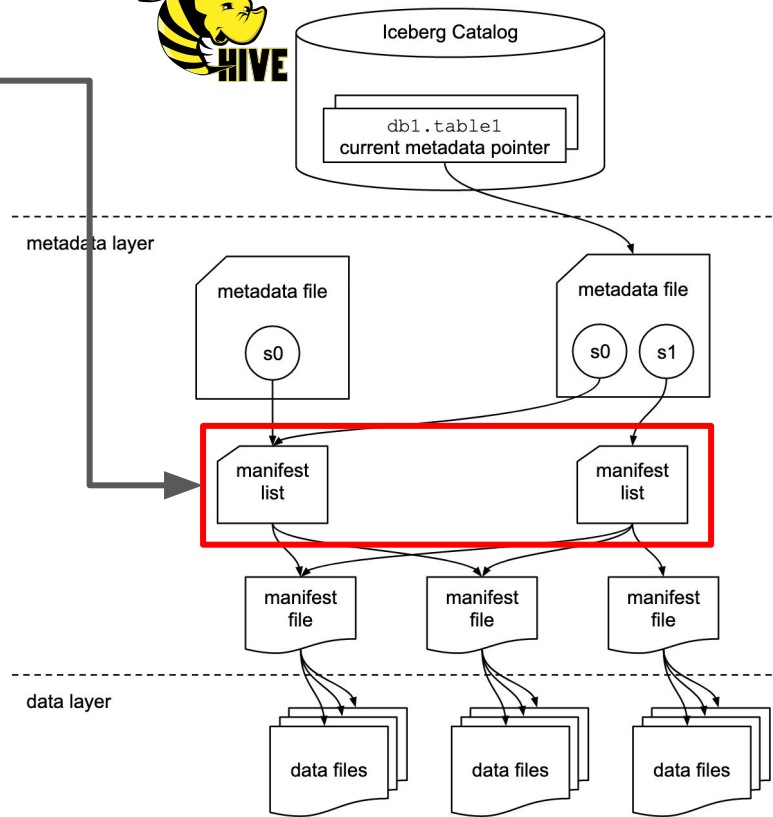
# Solution - What is Apache Iceberg?

```
} ],
"properties" : { },
"current-snapshot-id" : 3327495339618742959,
"snapshots" : [ {
  "snapshot-id" : 3327495339618742959,
  "timestamp-ms" : 1610742940007,
  "summary" : {
    "operation" : "append",
    "added-data-files" : "1",
    "added-records" : "4",
    "added-files-size" : "2524",
    "changed-partition-count" : "1",
    "total-records" : "4",
    "total-data-files" : "1",
    "total-delete-files" : "0",
    "total-position-deletes" : "0",
    "total-equality-deletes" : "0"
  },
  "manifest-list" : "gs://my_bucket/hive-warehouse/table/metadata/snap-00001.avro"
} ],
"snapshot-log" : [ {
  "timestamp-ms" : 1610742940007,
  "snapshot-id" : 3327495339618742959
} ],
"metadata-log" : [ {
  "timestamp-ms" : 1610742928462
  "metadata-file" "gs://my_bucket/hive-warehouse/table/metadata/snap-00001.avro"
} ]
```

# Solution - What is Apache Iceberg?

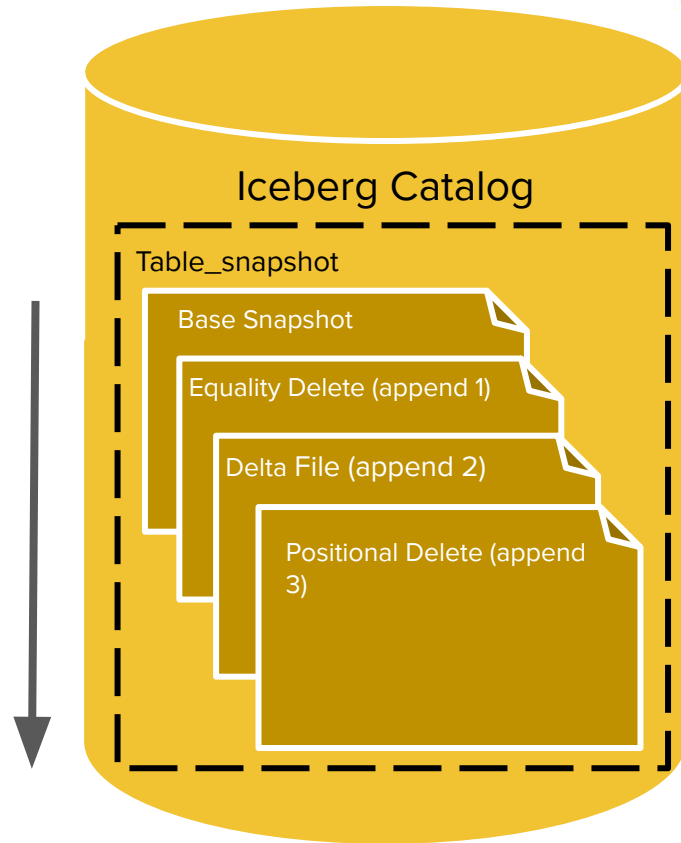# Solution - What is Apache Iceberg?

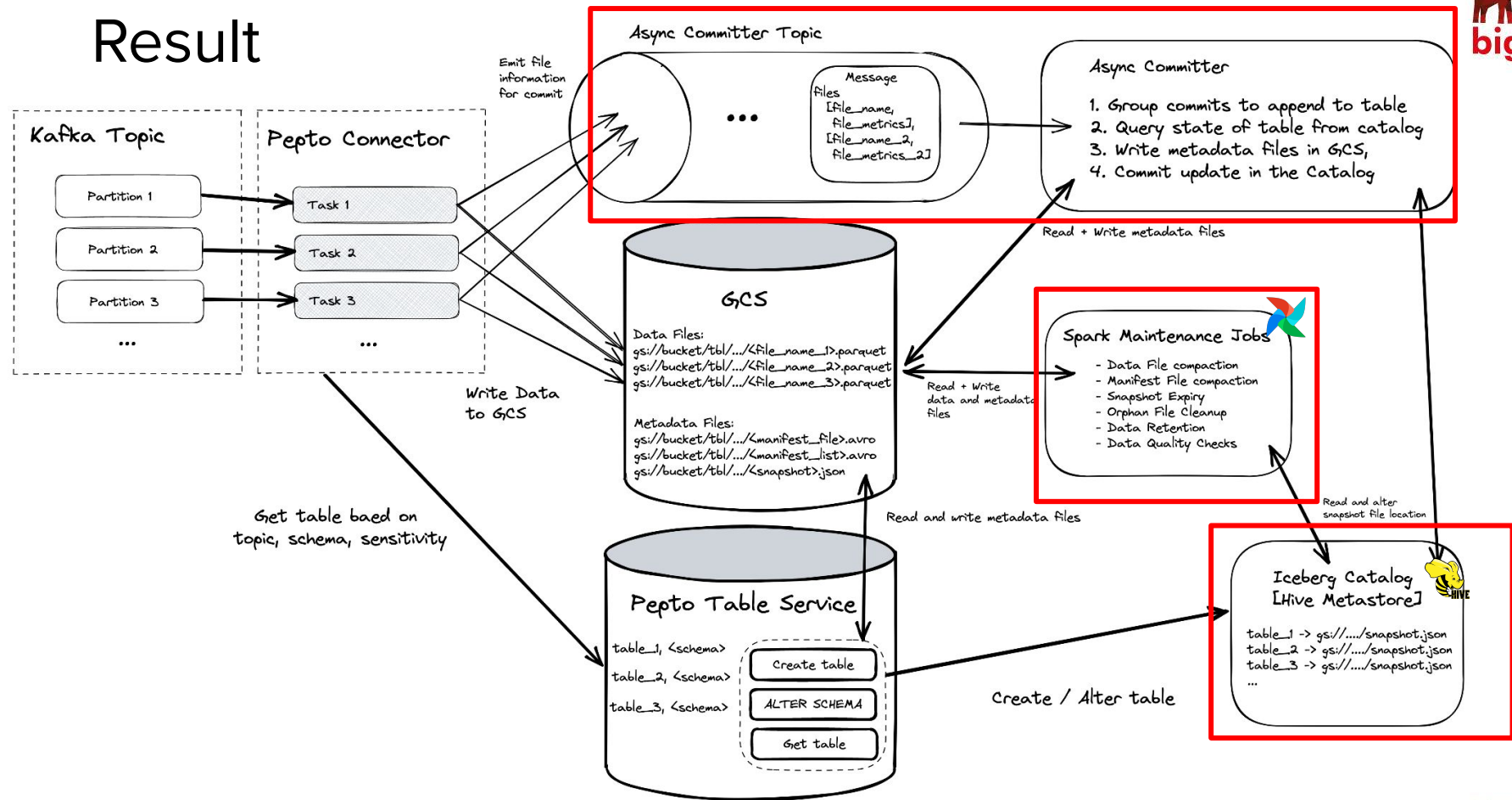# Solution - Promise

- **Merge-on-Read**
  - V2 Spec introduces delete files
    - Positional Delete
    - Equality delete
  - Applied as filters at query time to resolve changes
    - Gets applied to the resultset of your executed query. No full table scan needed!
- **Erik Wright** from Shopify helped write the proposal for Iceberg merge-on-read



Iceberg Catalog

Table_snapshot

Base Snapshot

Equality Delete (append 1)

Delta File (append 2)

Positional Delete (append 3)

# Result

# Result - Zoom Out
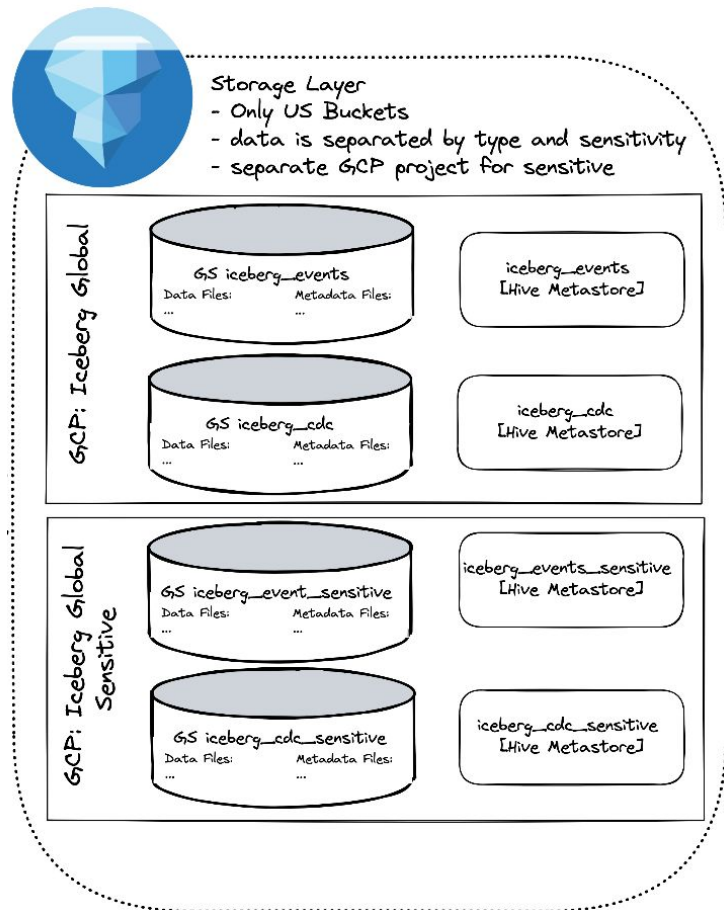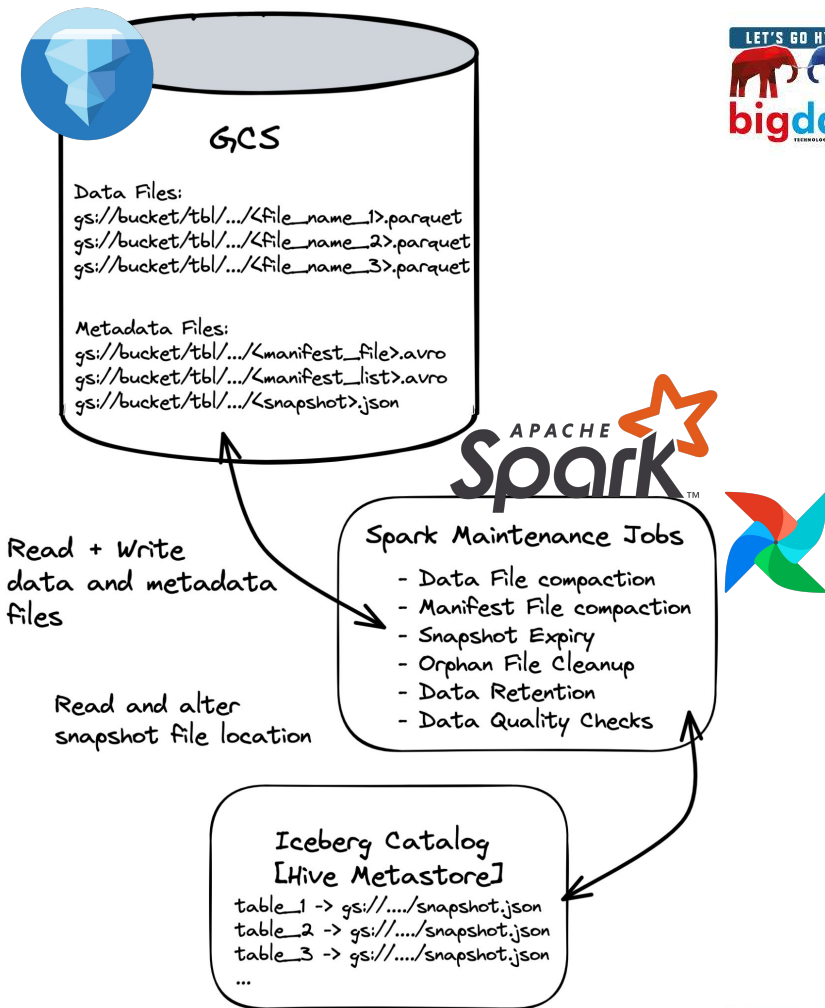
# Result - Storage

- **Bucket per Region + Sensitivity**
  - For us this ended up being a bucket per a catalog
  - **Pushback against** additional infra related work for managing + **spinning up + configuring buckets on the fly**
    - Currently buckets are managed by Terraform
  - Future access restricts could be applied via GCS **prefix IAM restrictions**
- **Catalogs group datasets with the same behaviour**
  - CDC vs events vs raw_type_1 etc.

Storage Layer
- Only US Buckets
- data is separated by type and sensitivity
- separate GCP project for sensitive

GCP: Iceberg Global

GS iceberg_events
Data Files: ...    Metadata Files: ...

iceberg_events
[Hive Metastore]

GS iceberg_cdc
Data Files: ...    Metadata Files: ...

iceberg_cdc
[Hive Metastore]

GCP: Iceberg Global Sensitive

GS iceberg_event_sensitive
Data Files: ...    Metadata Files: ...

iceberg_events_sensitive
[Hive Metastore]

GS iceberg_cdc_sensitive
Data Files: ...    Metadata Files: ...

iceberg_cdc_sensitive
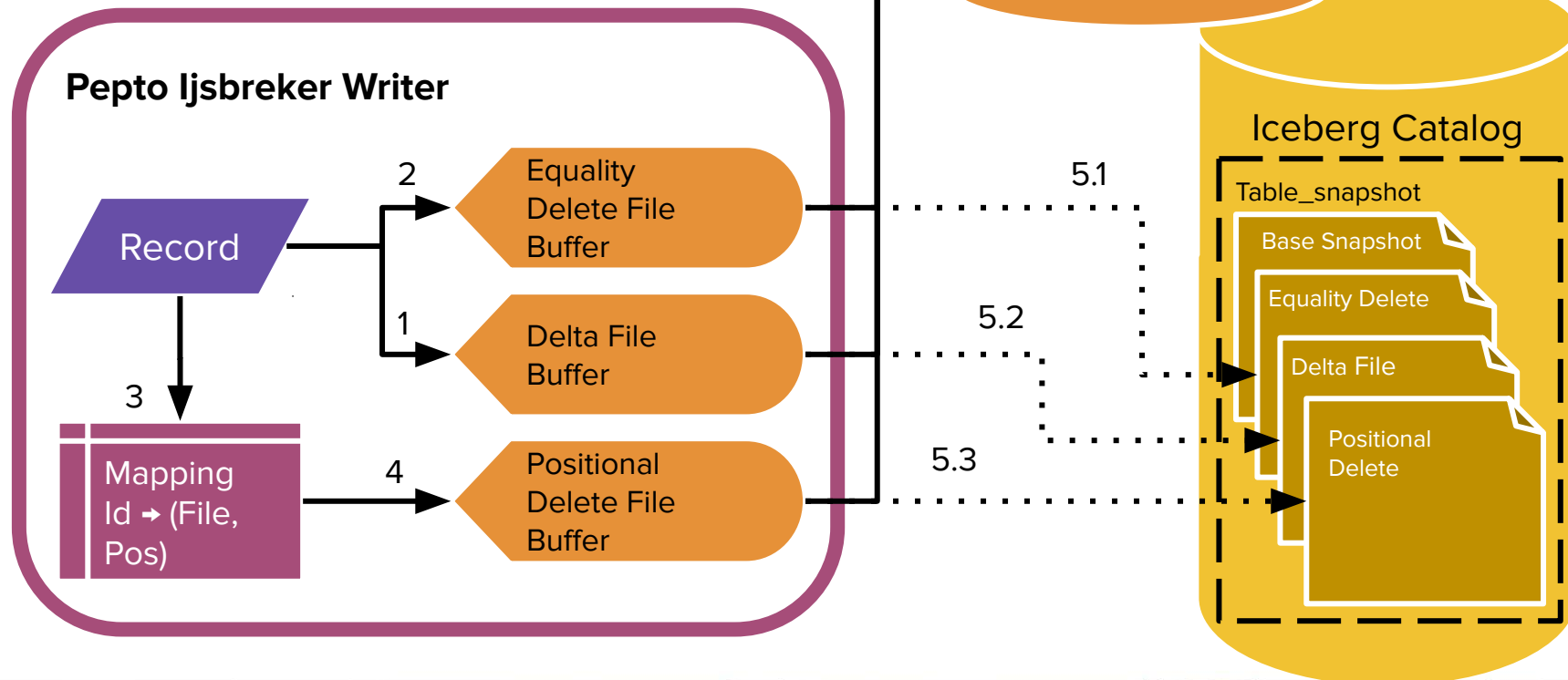[Hive Metastore]
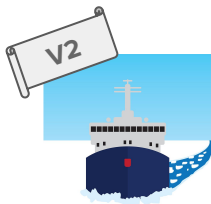
# Result - Maintenance



- **Small data** + **metadata file problem** because of micro-batch processing (streaming)
  - Datafile and manifest compaction solves the small file problem for us
- **Versioning our datasets** by keeping deltas around
- **Privacy**
  - PII is purged after 30 days in GCS
  - Inflight enforcement of data
  - Re-enforcement
- **GCS cleanup**
  - From retention or snapshot expiry, files that are no longer registered to the table are deleted from GCS
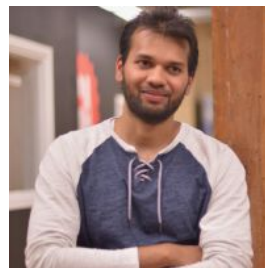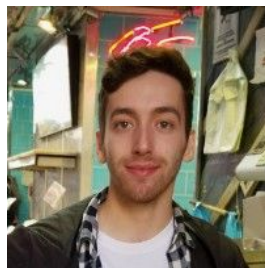
# Reflection

- **I wish we did a bucket per a dataset**

  - A very fine grain of separation

  - Easier to implement specific restrictions / functionality

    - Regulated industry (ex SOX)

  - More upfront work but would have accommodated additional use cases

- **Writing your own engine to utilize Iceberg can be hard** (Kafka connect is not supported)

  - We ended up building our own version of many of the concepts that you see in the Flink Iceberg Connector (async committing to)

  - **Flink connector did not fully exist when we started**

  - Documentation is not great leading to a **divergence between the online docs and their Java API**

# Future Challenges

# Thank you!