

Apache Iceberg Crash Course

Understanding Apache Iceberg's Partitioning Features



Curriculum

July 11: What is a Data Lakehouse and What is a Table Format?

July 16: The Architecture of Apache Iceberg, Apache Hudi and Delta Lake

July 23: The Read and Write Process for Apache Iceberg Tables

Aug 13: Understanding Apache Iceberg's Partitioning Features

Aug 27: Optimizing Apache Iceberg Tables

Sep 3: Streaming with Apache Iceberg

Sep 17: The Role of Apache Iceberg Catalogs

Oct 1: Versioning with Apache Iceberg

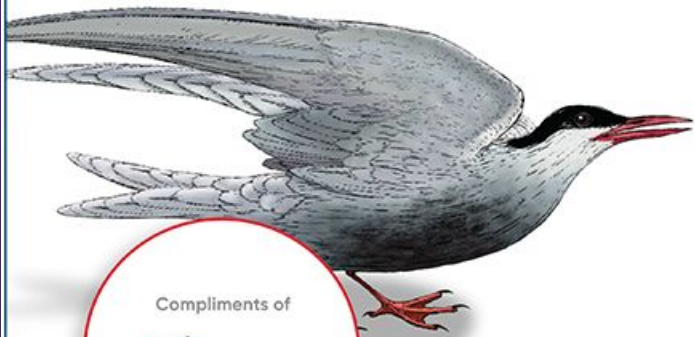
Oct 15: Ingesting Data into Apache Iceberg with Apache Spark

Oct 29: Ingesting Data into Apache Iceberg with Dremio

O'REILLY®

Apache Iceberg The Definitive Guide

Data Lakehouse Functionality, Performance,
and Scalability on the Data Lake



Compliments of



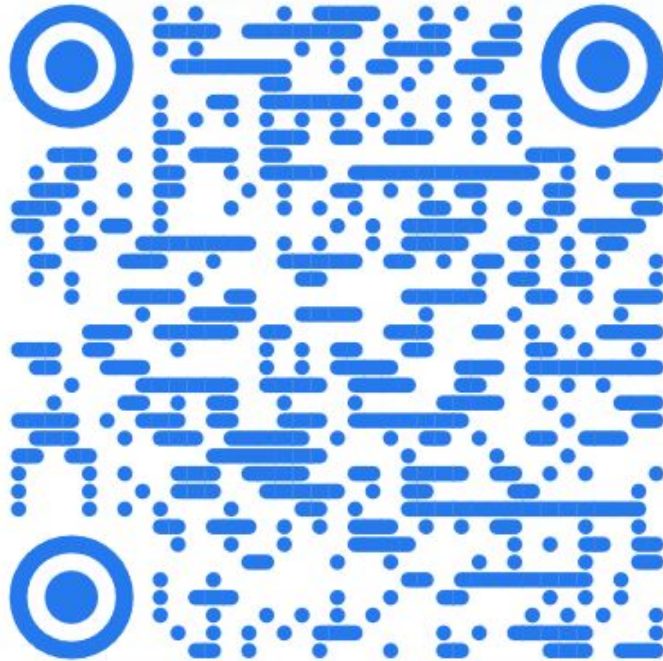
Tomer Shiran,
Jason Hughes &
Alex Merced

Forewords by Gerrit Kazmaier,
Raghu Ramakrishnan & Rick Sears





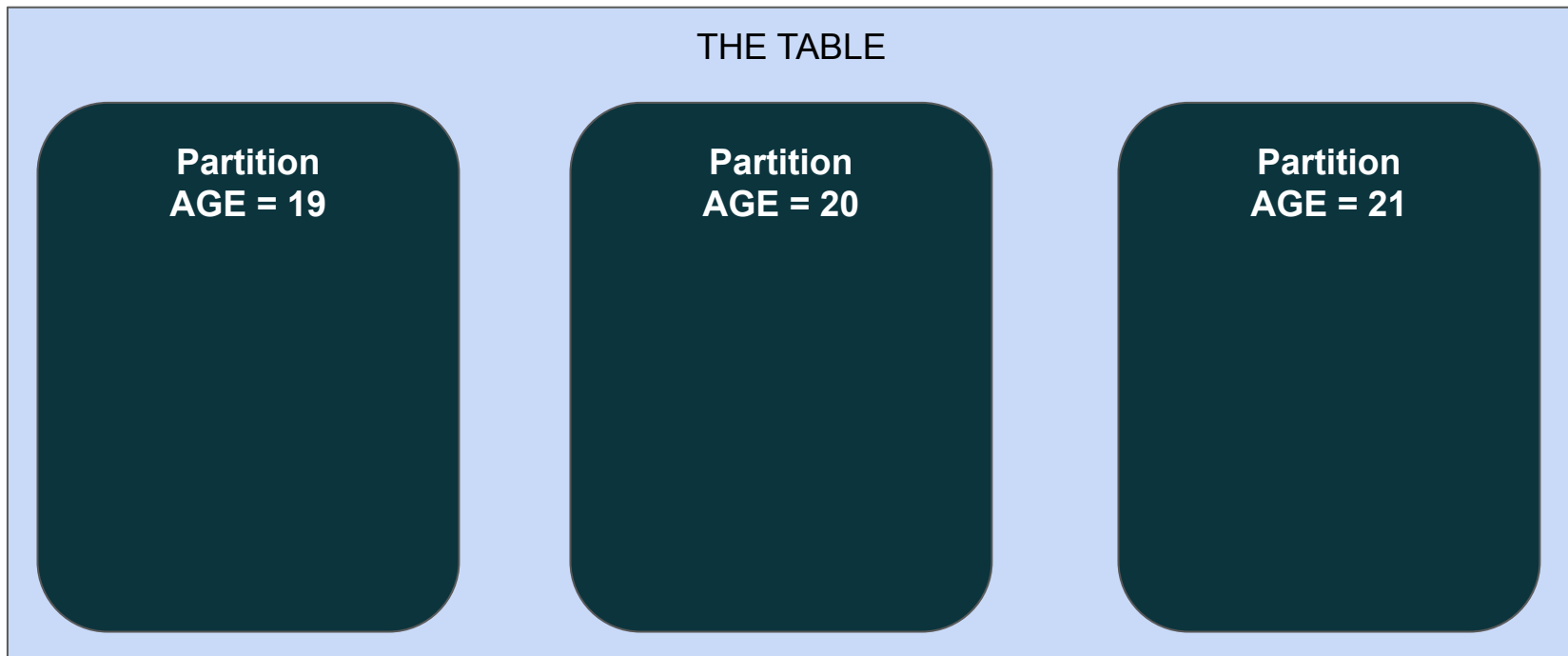
dremio.com/gnarly-data-waves
Youtube | Spotify | iTunes



community.dremio.com
Apache Iceberg Category

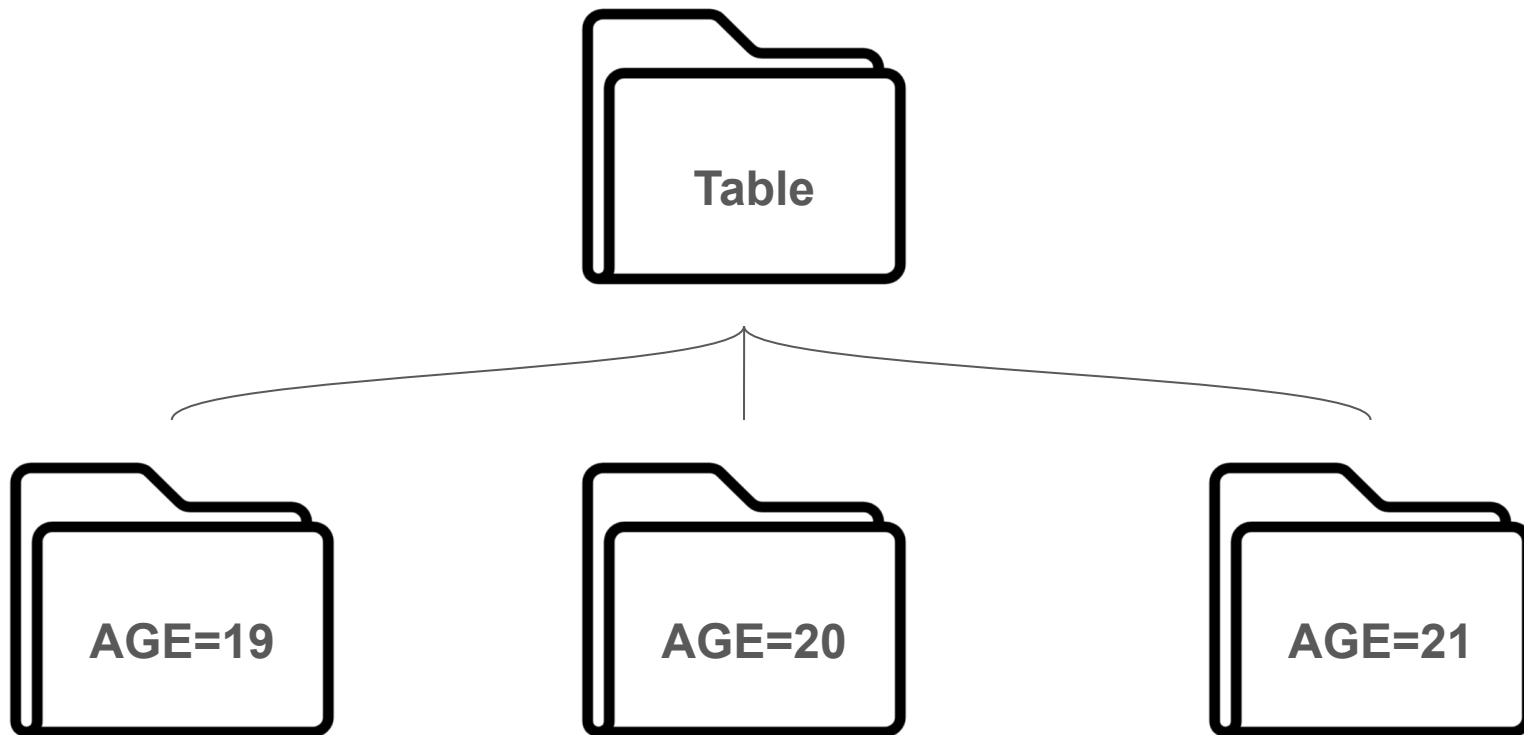
What Is Partitioning?

What is Partitioning?



Partitioning with Hive

Directory Based



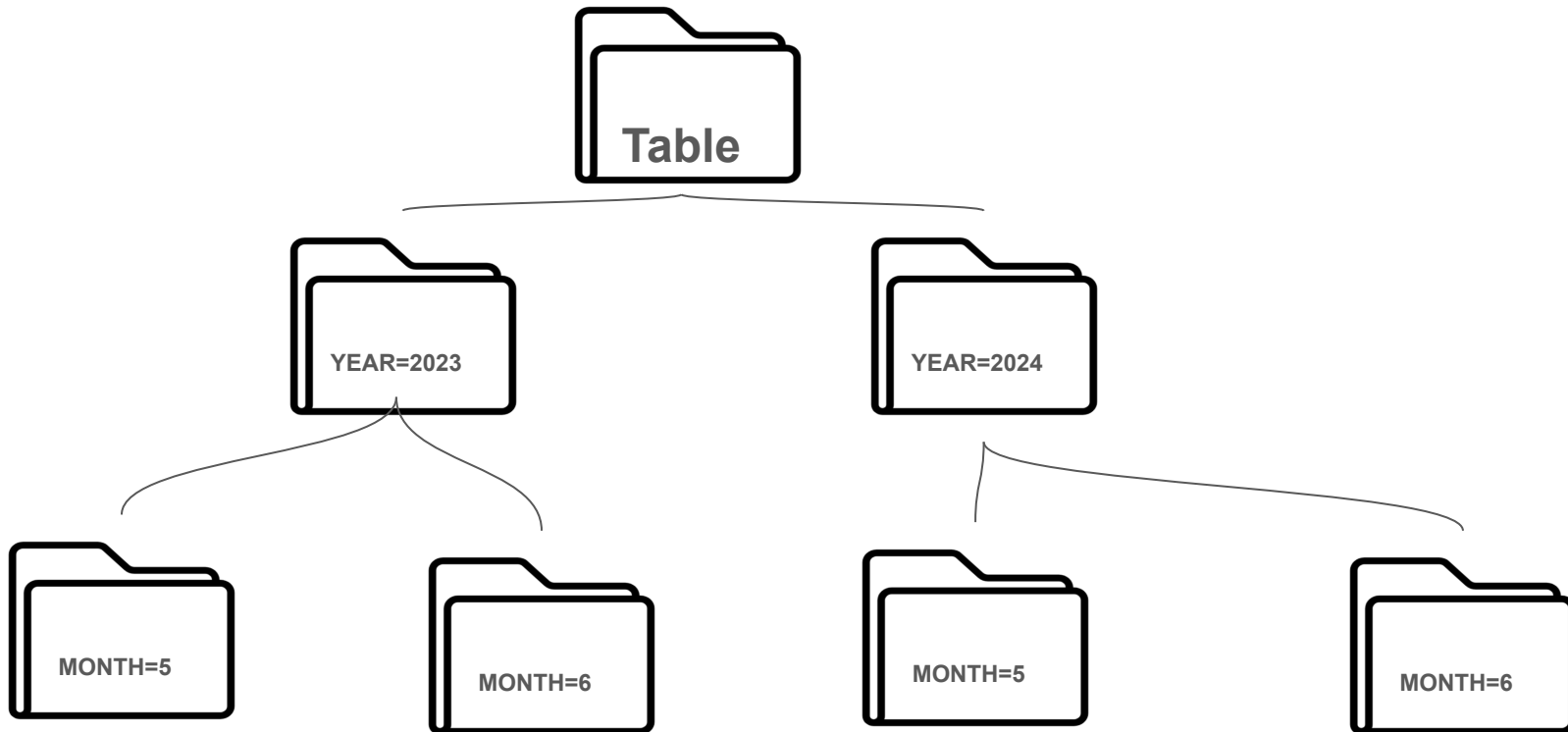
Year/Month/Day Partitioning Difficulties

ID	NAME	AGE	CITY	DATE
1	Alice	28	New York	2024-05-15 10:15:30
2	Bob	34	San Francisco	2024-06-10 11:20:45
3	Carol	23	Chicago	2024-05-22 12:30:10
4	David	45	Los Angeles	2024-07-03 13:45:00
5	Eve	30	Seattle	2024-06-25 14:50:25

Year/Month/Day Partitioning Difficulties

ID	NAME	AGE	CITY	DATE	Month	Year
1	Alice	28	New York	2024-05-15 10:15:30	5	2024
2	Bob	34	San Francisco	2024-06-10 11:20:45	6	2024
3	Carol	23	Chicago	2024-05-22 12:30:10	5	2024
4	David	45	Los Angeles	2024-07-03 13:45:00	7	2024
5	Eve	30	Seattle	2024-06-25 14:50:25	6	2024

Resulting Partitioning



Year/Month/Day Partitioning Difficulties

```

-- FULL TABLE SCAN
SELECT *
FROM people
WHERE date
BETWEEN '2024-06-15 00:00:00' AND '2024-06-30 23:59:59';

-- PARTITION SCAN
SELECT *
FROM people
WHERE (year = 2024 AND month = 6)
      AND date BETWEEN '2024-06-15 00:00:00' AND '2024-06-30
23:59:59';
```

Iceberg's Partitioning

Metadata.json

```
"partition-specs": [  
  {  
    "spec-id": 0,  
    "fields": [  
      {  
        "source-id": 5,  
        "field-id": 1000,  
        "name": "month",  
        "transform": "month"  
      }  
    ]  
  }  
],
```

```
"fields": [  
  {  
    "id": 1,  
    "name": "id",  
    "required": true,  
    "type": "long"  
  },  
  {  
    "id": 2,  
    "name": "name",  
    "required": false,  
    "type": "string"  
  },  
  {  
    "id": 3,  
    "name": "age",  
    "required": false,  
    "type": "integer"  
  },  
  {  
    "id": 4,  
    "name": "city",  
    "required": false,  
    "type": "string"  
  },  
  {  
    "id": 5,  
    "name": "timestamp",  
    "required": true,  
    "type": "timestamp"  
  }  
],
```

Manifest List (Snapshot Level)

Manifest File Path	Manifest Length	Partition Spec ID	Added Files Count	Existing Files Count	Deleted Files Count	Partition Data
s3://bucket/path/manifest1.avro	1024	0	5	3	1	{"month": "2024-06"}
s3://bucket/path/manifest2.avro	2048	1	10	0	0	{"month": "2024-07"}
s3://bucket/path/manifest3.avro	4096	0	2	1	2	{"month": "2024-05"}
s3://bucket/path/manifest4.avro	512	1	6	2	0	{"month": "2024-06"}
s3://bucket/path/manifest5.avro	3072	0	3	3	1	{"month": "2024-07"}

Manifest (Partition/Files Level)

File Path	Partition Data	Lower Bounds	Upper Bounds
s3://bucket/path/data1.parquet	{"month": "2024-06"}	{"id": 1, "name": "A", "age": 20, "city": "A City", "timestamp": "2024-06-01 00:00:00"}	{"id": 100000, "name": "Z", "age": 50, "city": "Z City", "timestamp": "2024-06-30 23:59:59"}
s3://bucket/path/data2.parquet	{"month": "2024-06"}	{"id": 100001, "name": "B", "age": 21, "city": "B City", "timestamp": "2024-06-01 00:00:00"}	{"id": 200000, "name": "Y", "age": 51, "city": "Y City", "timestamp": "2024-06-30 23:59:59"}
s3://bucket/path/data3.parquet	{"month": "2024-06"}	{"id": 200001, "name": "C", "age": 22, "city": "C City", "timestamp": "2024-06-01 00:00:00"}	{"id": 300000, "name": "X", "age": 52, "city": "X City", "timestamp": "2024-06-30 23:59:59"}
s3://bucket/path/data4.parquet	{"month": "2024-06"}	{"id": 300001, "name": "D", "age": 23, "city": "D City", "timestamp": "2024-06-01 00:00:00"}	{"id": 400000, "name": "W", "age": 53, "city": "W City", "timestamp": "2024-06-30 23:59:59"}

Hidden Partitioning

Partition Transforms

- Year
- Month
- Day
- Hour
- Bucket
- Truncate

Result



```
-- Partition Scan thanks to Apache Iceberg  
SELECT *  
FROM your_table_name  
WHERE date  
BETWEEN '2024-06-15 00:00:00' AND '2024-06-30 23:59:59';
```

Partition Evolution

Result

Partitioned by Year
(2020-2023)

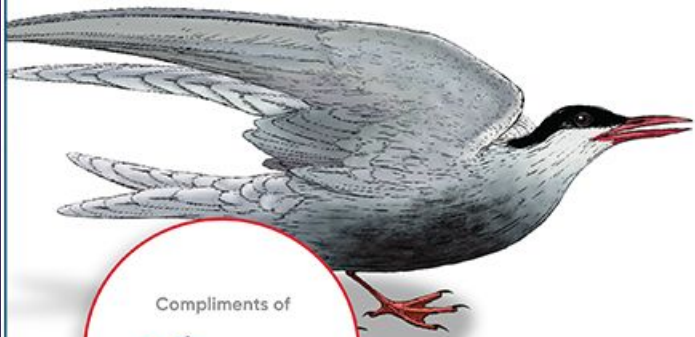
Partitioned by Month
(01/2024-05/2024)

Partitioned by Day
(06/01/2024 -
06/30/2024)

O'REILLY®

Apache Iceberg The Definitive Guide

Data Lakehouse Functionality, Performance,
and Scalability on the Data Lake



Compliments of



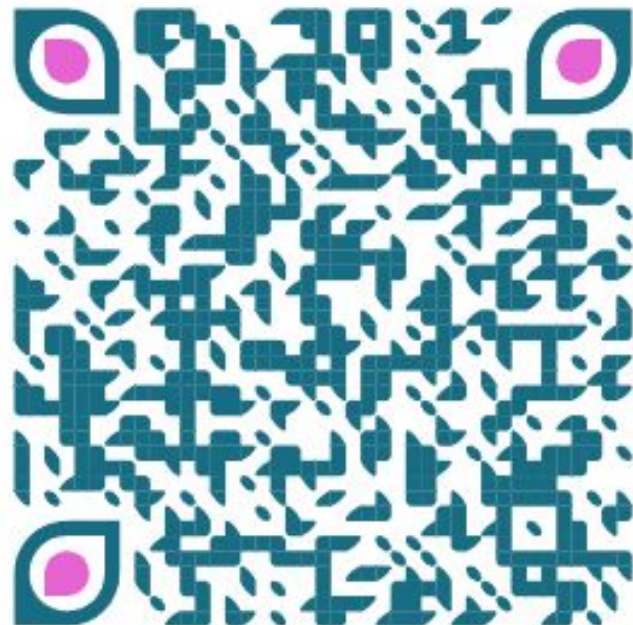
Tomer Shiran,
Jason Hughes &
Alex Merced

Forewords by Gerrit Kazmaier,
Raghu Ramakrishnan & Rick Sears





A Iceberg/Dremio Lakehouse on
your laptop exercise



Deploy Dremio Software or
Dremio Cloud



Postgres -> Iceberg -> Dashboard

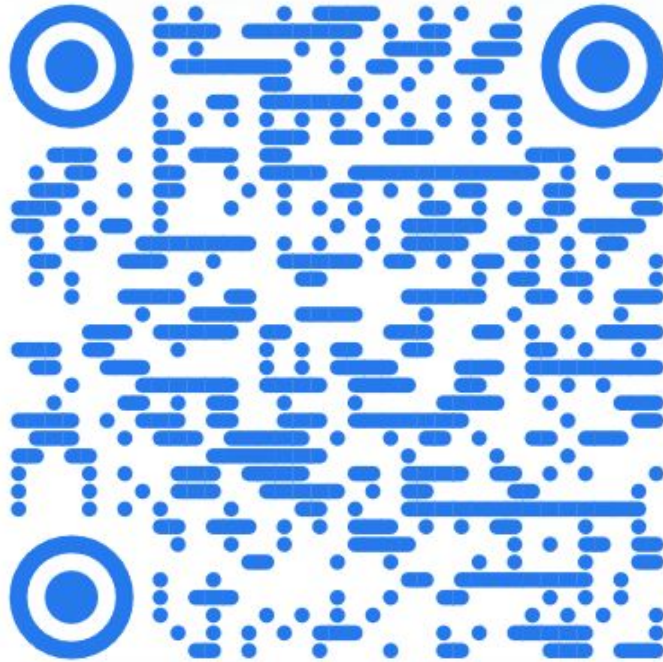


SQLServer -> Iceberg -> Dashboard



MongoDB -> Iceberg -> Dashboard

dremio.com/blog



community.dremio.com
Apache Iceberg Category