GNARLY
Data_Waves

PRESENTED BY dremio

EPISODE 8

# Managing Your Data as Code with Dremio Arctic

Jeremiah Morrow
Director, Product Marketing for Dremio Arctic
jeremiah.morrow@dremio.com
www.linkedin.com/in/jeremiahmorrow

# Agenda

1. The need for data as code
2. Dremio Arctic
3. Use Cases for Data as Code
4. Demonstration

# The need for data as code

# The Data Explosion



**More data**

175 ZB of data worldwide by 2025, with as much residing in the cloud as on-prem



**More consumers**

Growing number of data requests from technical & non-technical users



**More complexity**

New sources, users and use cases mean we need to change the way we work with data.

# GitHub/GitLab is Central to How We Build Products



- ✓ Code storage
- ✓ Atomic changes
- ✓ Documentation
- ✓ Access old versions
- ✓ Recover from mistakes
- ✓ Trace history
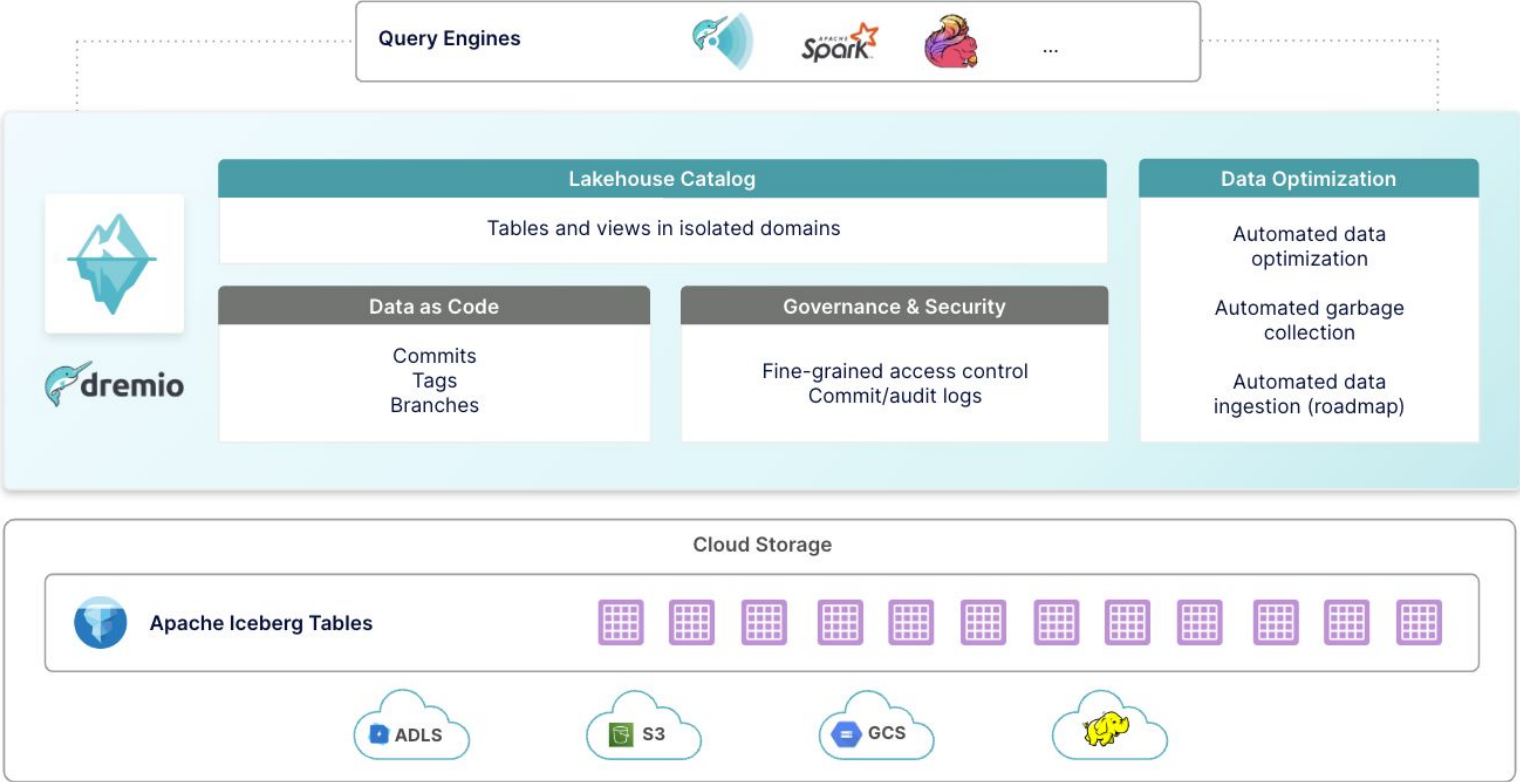- ✓ Isolated development
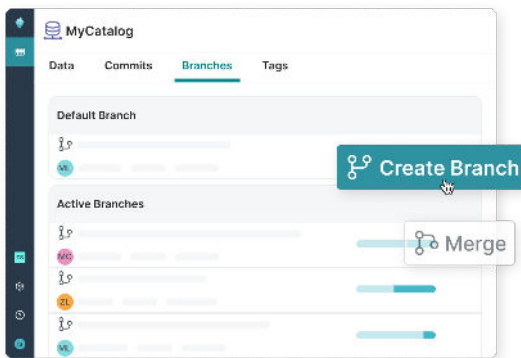- ✓ CI/CD
- ✓ Collaboration

# Dremio Arctic

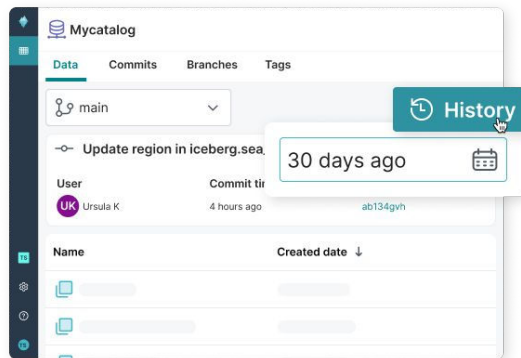# Dremio Arctic is a Data Lakehouse Management Service

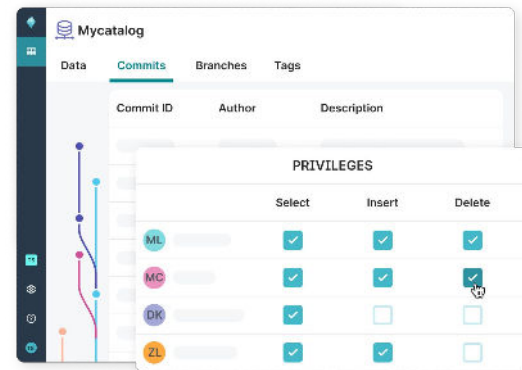# Dremio Arctic Enables Data as Code Management



## ISOLATION

- Experiment with data without impacting other users

- Ingest, transform and test data before exposing it to other users in an atomic merge

## VERSION CONTROL

- Reproduce models and dashboards from historical data based on time or tags

- Recover from any mistake by instantly undoing accidental data or metadata changes
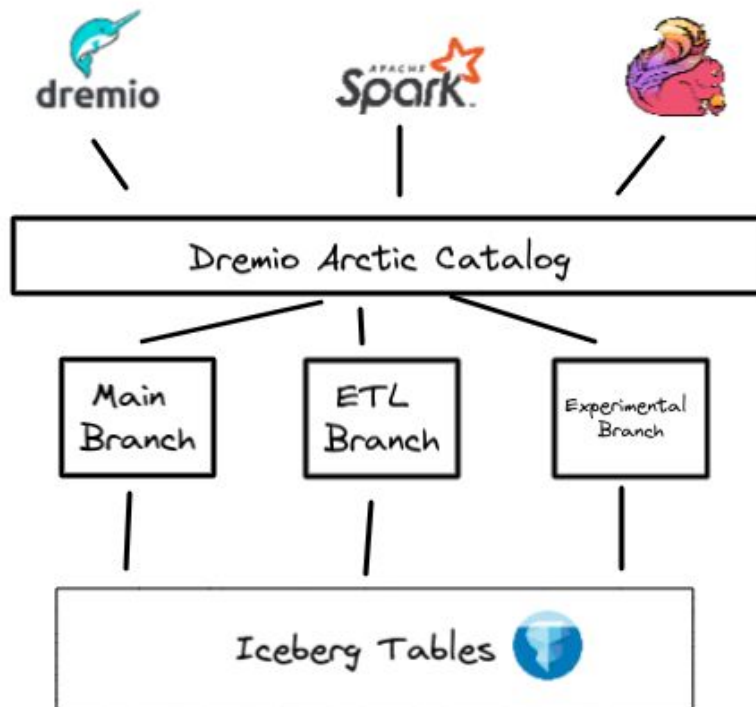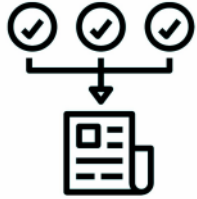
## GOVERNANCE

- All changes to the data and metadata are tracked: who accessed what data and when

- Fine-grained privileges to control access to the data at the table, column and row level

# Branches, Tags, and Commits: A No-Copy Solution for Data Management
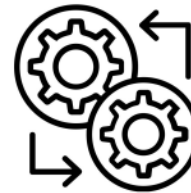
# Benefits of Data as Code

**Consistency**

**Experimentation**

**Collaboration**

**Governance**

**Reproducibility**

# 5 Use Cases for Data as Code

# Ensure Data Quality with ETL Branches

Create an ETL branch and ingest the data with COPY INTO, CTAS or Spark:

```
CREATE BRANCH events_etl_9_28_22
USE BRANCH events_etl_9_28_22
COPY INTO web.events …
```
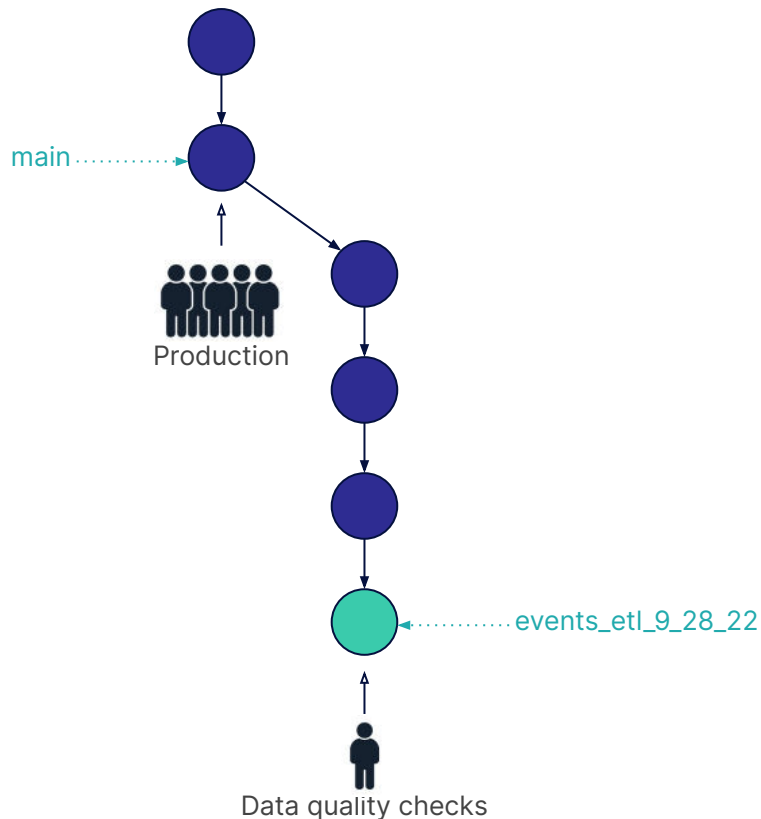
Run queries to test data quality:

```
SELECT COUNT(*) FROM web.events WHERE
length(ip_address) >= 7
```

Test the dashboard to see that it looks ok:



Fix the problems and merge into main:

```
DELETE FROM web.events WHERE length(ip_address) >= 7
USE BRANCH main
MERGE BRANCH events_etl_9_28_22
```



main

Production

events_etl_9_28_22

Data quality checks

# Experiment with Data in Transient Branches

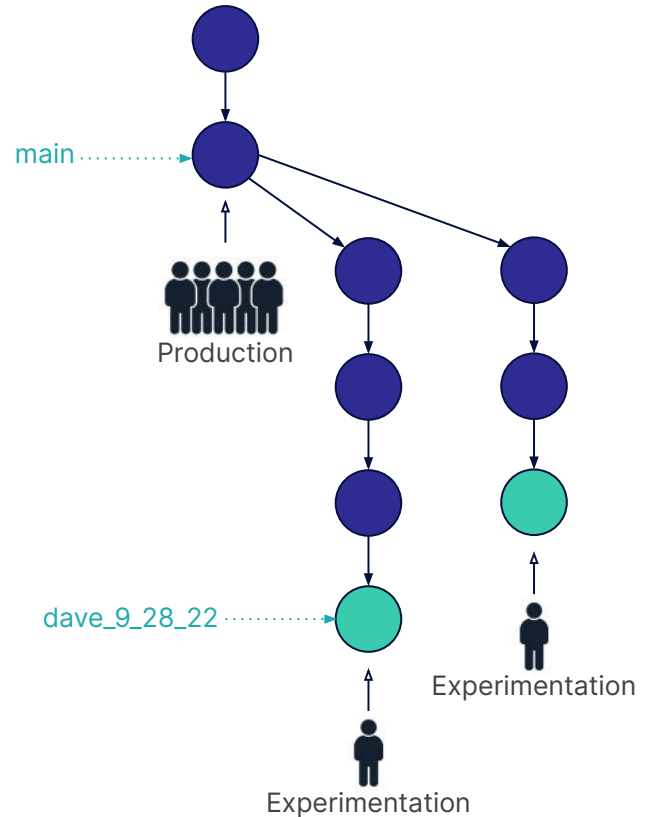Create a transient branch and perform data explorations and transformations in it:

```
CREATE BRANCH dave_9_28_22
USE BRANCH dave_9_28_22
CREATE TABLE t AS SELECT ...
UPDATE t ... SET ...
```

Create ad-hoc visualizations on the branch via a Notebook:



Delete the branch or merge it when experimentation is complete:

```
DROP BRANCH dave_9_28_22
```



main

Production

dave_9_28_22

Experimentation

Experimentation

# Reproduce Models

Change context to a named tag:

```
spark.sql("USE REFERENCE modelA in arctic;")
```

Create ML model based on historic data:

```
val trainingData = spark.read.table("arctic.t")
val lr = new LogisticRegression()
// configure logistic regression...
val paramMap = ParamMap(...)
val model = lr.fit(trainingData, paramMap)
```

Select a tag, commit or branch to query in SQL Runner:

# Recover From Mistakes

Move the branch head to a historical commit:

```
ALTER BRANCH main ASSIGN COMMIT ...f724
```

...a233

...b84c

...9bc8

...f724 ← main'
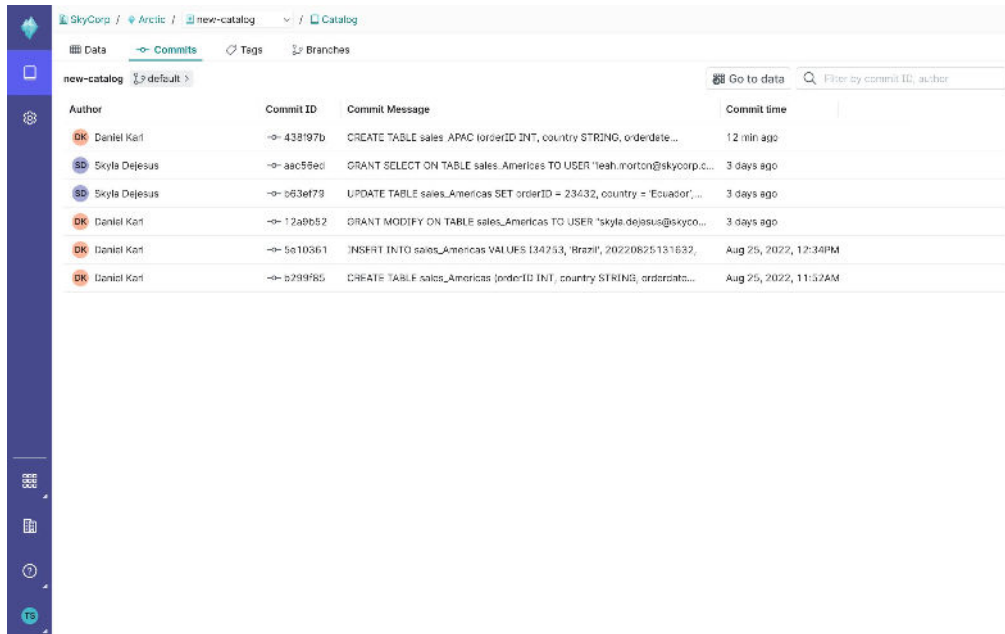
...2563

...4231 ← main

# Troubleshoot

Get the commit history for a branch:

```
SHOW LOGS AT REFERENCE etl;
```

Get the commit history for a specific table:

```
curl -X GET -H 'Authorization: Bearer
<PAT>' <Catalog API
Endpoint>/trees/tree/<reference
name>/log\?filter="operations.exists(op,op.
key=='<table name>')"
```
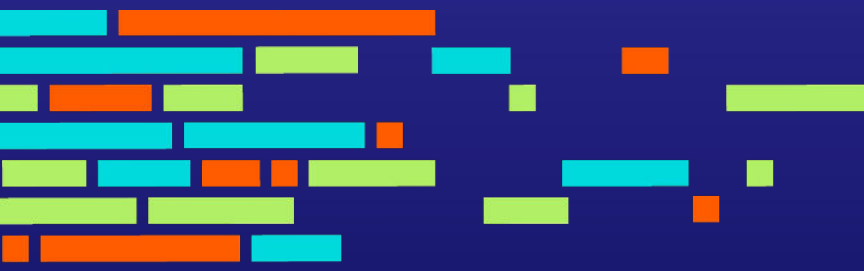
# Dremio Arctic Demo