# Create a git-like experience for Data Lake analytics

October 2021, Itai Admi





**Itai Admi** R&D Team Lead Treeverse

• Tel Aviv, Israel

#### Data Lake Architecture Reference





#### The Data Lake Advantage

- 1. Scalability and cost effectiveness
- 2. Accessibility and ease of use
- 3. High throughput
- 4. Rich application ecosystem



#### The Data Lake Challenge

- 1. In-ability to experiment, compare and reproduce Example: add new metric for BI
- 2. Difficult to enforce data best practices Example: schema, format enforcement
- **3. Hard to ensure high quality data** Example: validate statistical properties of the data



In a perfect world

we would manage data from dev to production the way we manage code





Atomic versioned data lake on top of object storage



#### Manageability & Reliability Layer





#### How it works

## **API Compatibility** with object stores

Versioning Engine

Transactional, atomic, isolated

#### **Git Terminology**

Branch, commit, merge





Azure

 $\begin{array}{c} \Box \\ \Box \\ \Box \\ \Box \end{array} \rightarrow \begin{array}{c} \Box \\ \Box \\ \Box \end{array}$ 





#### How it works

### s3://data-bucket/collections/foo

#### s3://data-bucket/<u>main</u>/collections/foo



#### Integrates with your existing tools



lakeFS can be used in any stage of your data development lifecycle



#### In Development

- **Experiment -** try new tools, upgrade versions, and evaluate code changes in isolation.
- **Debug -** checkout specific commits in a repository's commit history to materialize consistent, historical versions of your data.
- **Collaborate** leverage isolated branches managed by metadata (not copies of files) to work in parallel.



#### Experiment Safely

.





#### Experimenting with Spark

lakectl branch create \
 lakefs://example-repo@testing-spark-3 \
 --source lakefs://example-repo@main

# output:

# created branch 'testing-spark-3.0', pointing to commit ID: 'd1e9adc71c10a'

val dfExperiment1 = sc.read.parquet("s3a://example-repo/experiment-1/events/by-date")
val dfExperiment2 = sc.read.parquet("s3a://example-repo/experiment-2/events/by-date")

```
dfExperiment1.groupBy(" ... ").count()
dfExperiment2.groupBy(" ... ").count() // now we can compare the properties of the data itself
```



#### Experimenting with Presto

```
lakectl branch create \
    lakefs://example-repo@testing-spark-3 \
    --source lakefs://example-repo@main
# output:
# output:
```

# created branch 'testing-spark-3.0', pointing to commit ID: 'd1e9adc71c10a'

```
CREATE TABLE master.request_logs (
  request_time timestamp,
  url varchar,
  ip varchar,
  user_agent varchar
)
WITH (
  format = 'TEXTFILE',
  external_location = 's3a://example/main/data/logs/'
);
```



#### **During Deployment**

- **Data update safety** Ingest new data onto an isolated branch, perform data validations, then add to production through a merge operation.
- **Test -** define pre-merge and pre-commit hooks to run tests that enforce schema and validate properties of the data to catch issues before they reach production.



#### **Enforce Best Practices**





•

### CI with PyArrow

import lakefs import pyarrow.parquet as pq # Setup a PyArrow FileSystem that we can use to query data in the source ref fs = lakefs.get\_filesystem(repo, source\_ref) for change in lakefs.diff(repo, source\_ref, target\_branch, prefix='public/'): if not change.path.endswith('.parquet'): continue # Read Parquet column metadata schema = pq.read\_schema(fs.open\_input\_file(change.path)) if filter(lambda column: column.name == 'user\_id', schema): raise ValidationError('user\_id column not allowed')



#### Streaming Data with Kafka Connect





#### In Production

- **Roll Back -** recover from errors by instantly reverting data to a former, consistent snapshot of the data lake.
- **Troubleshoot -** investigate production errors by starting with a snapshot of the inputs to the failed process.
- **Cross-collection Consistency** provide consumers multiple synchronized collections of data in one atomic, revertable action.



#### Troubleshoot - Reproduce a bug in production





# Simplify workflows at each step of the data lifecycle



### lakeFS Architecture





	Code	Data Lake
Amount of files tracked	Thousands- <u>3.5 million</u>	Millions-billions
File format	Mutable text files	Immutable binary data files, unstructured data
Change velocity	Hundreds of changes/day	Potentially millions/day
Change resolution	Lines in files	files
Mode of operation	Clone locally, push changes	<u>Unless you're rich</u> , you work remotely



#### Architecture





#### Data Format









#### Additional Resources





Check out the docs



Join the lakeFS Slack Channel



Contribute and star the <u>repo</u>



# Thanks!









einat.orr@treeverse.io