Building an ML Experimentation Platform for Easy Reproducibility





About Me

- I am Vino Duraisamy.
- SWE -> Data/ML Engineer -> Developer Advocate
- Open-source contributor @lakeFS

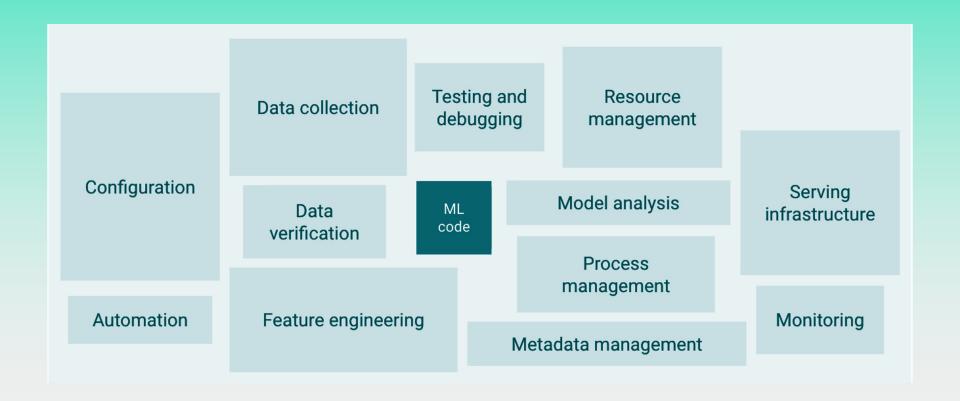




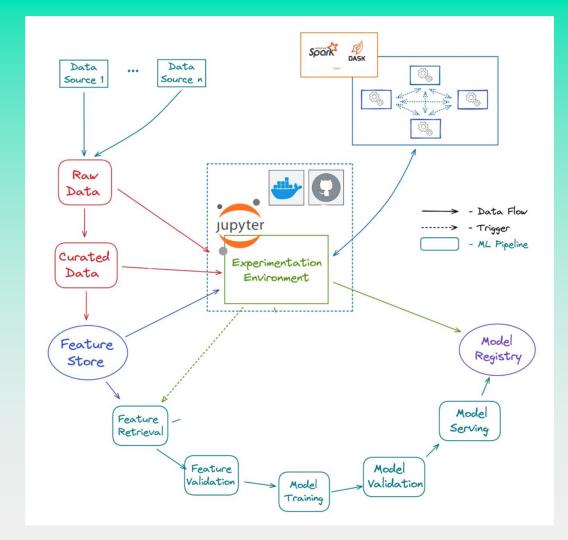




MLOps in a Complex Data World



MLOps in a Complex Data World



Uncovering the Gaps in MLOps

- ML experimentation infrastructure
 - Automation of experiments
 - Crisis of Reproducibility
 - Explainability of ML models
- Data/Feature ownership and Lineage
- Collaboration
 - MLEs working with same training data/feature set
- Version Control all ML assets atomically
 - Data, Model Artifacts, Code, Config, Metrics

MLOps landscape today





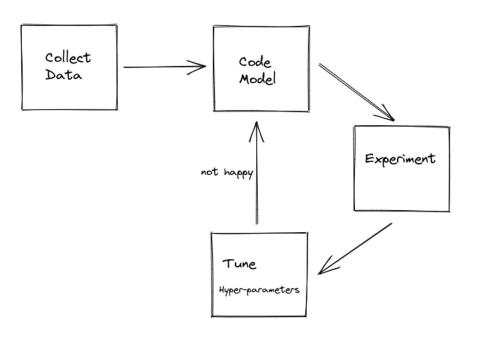
1. ML Experimentation



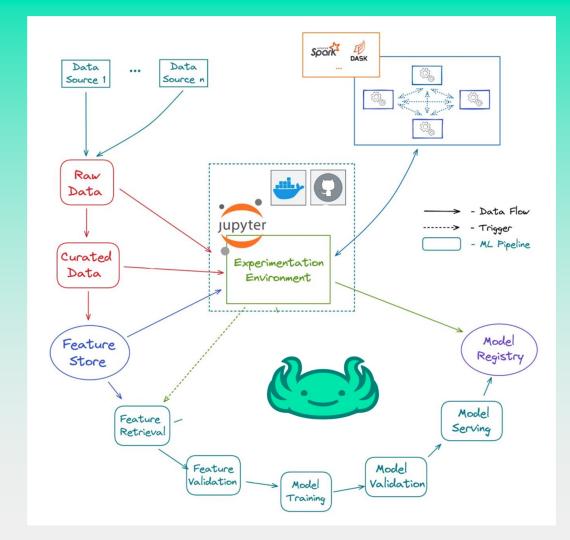


ML Experimentation: Challenges

- Iterative development
- Ensure data consistency
- Trust your ML pipeline
- Compliance to Regulations
- Reproducibility of ML training
- Explainable ML models

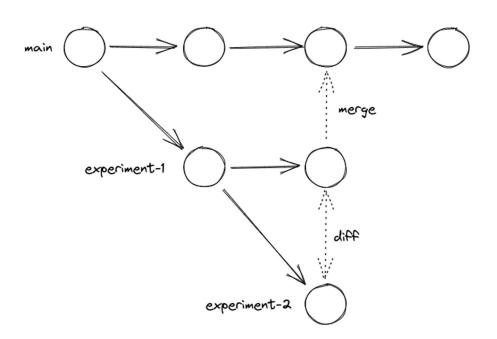


MLOps with lakeFS

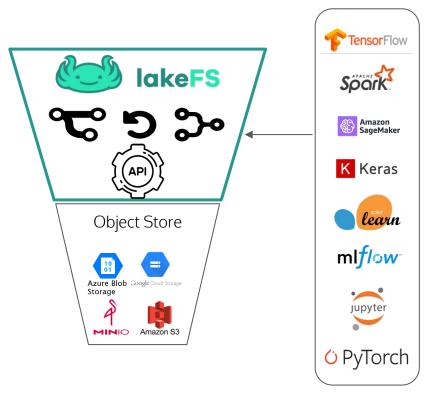


ML Experimentation Infrastructure

Git-like branching and experimenting with lakeFS



Git for Data - lakeFS

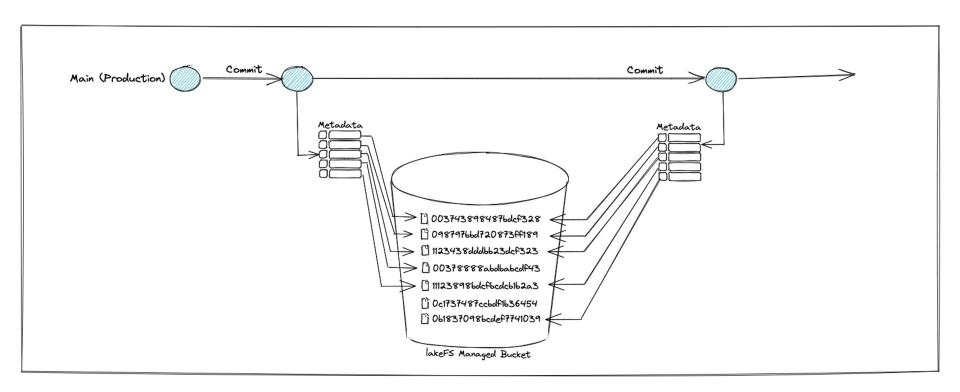


s3://repo/collections/foo

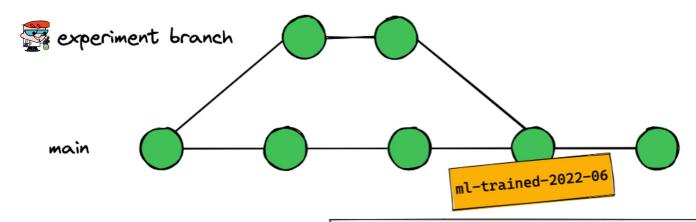
lakefs://repo/main/collections/foo



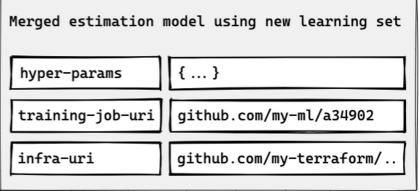
How does lakeFS work?



lakeFS for ML experimentation



```
df = spark.read.parquet(
'lakefs://my-repo/ml-trained-2022-06/inputs/vists/')
```

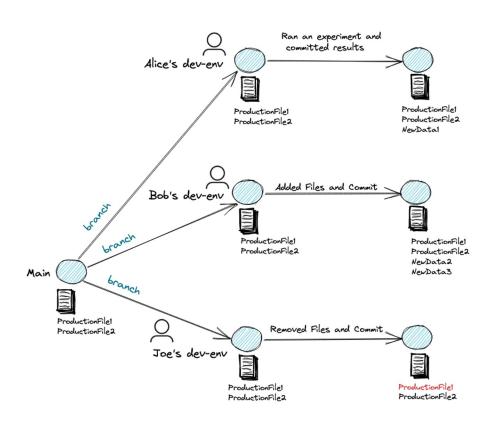


2. Unlock Collaboration





Modularize & Collaborate at each step



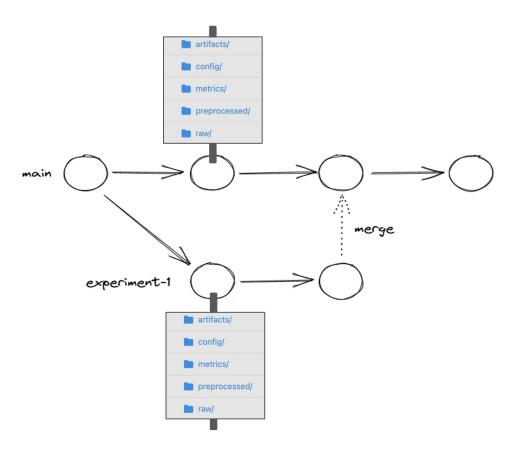
Enable ML engineers to work on same training data and features without duplication using lakeFS

3. Version Control all ML assets atomically

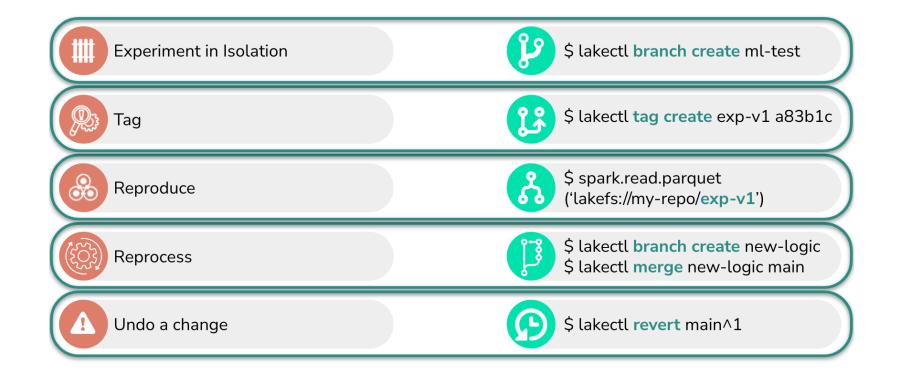




Version control all ML assets atomically



Using lakeFS



Demo Time!





lakeFS Community











BEDROCK







//w pollinate

























lakefs.io/slack

