# lyft's Feature Store: Architecture, Optimization, and Evolution

*Rohan Varshney*, Senior Software Engineer

Organized by HOPSWORKS

Lyft's Feature Store, a core infrastructure component in its Data Platform, **optimizes the management & deployment of ML features at scale**. It centralizes feature engineering & ensures uniformity across models & workflows by streamlining feature creation & storage for both offline/online model training & inference, facilitating low-latency access & high-throughput processing. This presentation covers its architecture, practical uses, performance, developer experience, optimization efforts, & evolution over the last 5+ years. We hope to demonstrate its role in empowering Lyft engineers to **develop service components & models more effectively, including for future AI/LLM applications**.

**FEATURE STORE SUMMIT 2025**

Organized by **HOPSWORKS**

| Team / Function | Description | Impact |
|---|---|---|
| Fulfillment | Houses ML models to match drivers to rides & generate + rank rider offers. | Empowers dispatch and offering services. |
| Orchestration | | |
| Pricing | | |
| Integrity | | |
| Growth Platform & Market Expansion | | |

Organized by  HOPSWORKS

| Team / Function | Description | Impact |
|---|---|---|
| Fulfillment | Houses ML models to match drivers to rides & generate + rank rider offers. | Empowers dispatch and offering services. |
| Orchestration | Platform for rider & driver incentive programs. Launch ML models for batch campaigns. | Produce incremental driver hour engagement & growth. |
| Pricing | | |
| Integrity | | |
| Growth Platform & Market Expansion | | |

| Team / Function | Description | Impact |
|---|---|---|
| Fulfillment | Houses ML models to match drivers to rides & generate + rank rider offers. | Empowers dispatch and offering services. |
| Orchestration | Platform for rider & driver incentive programs. Launch ML models for batch campaigns. | Produce incremental driver hour engagement & growth. |
| Pricing | ML models that generate prices, coupons, etc. | 💰💰💰 |
| Integrity | | |
| Growth Platform & Market Expansion | | |

Organized by **HOPSWORKS**

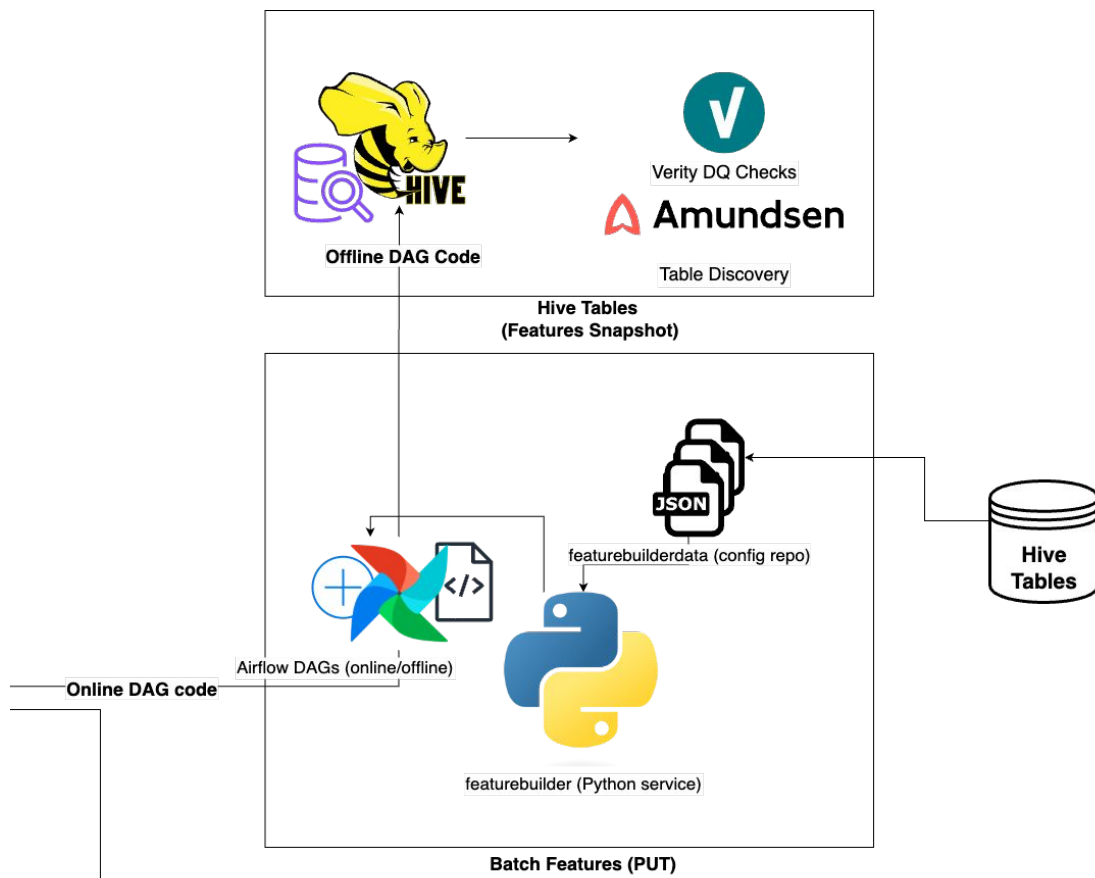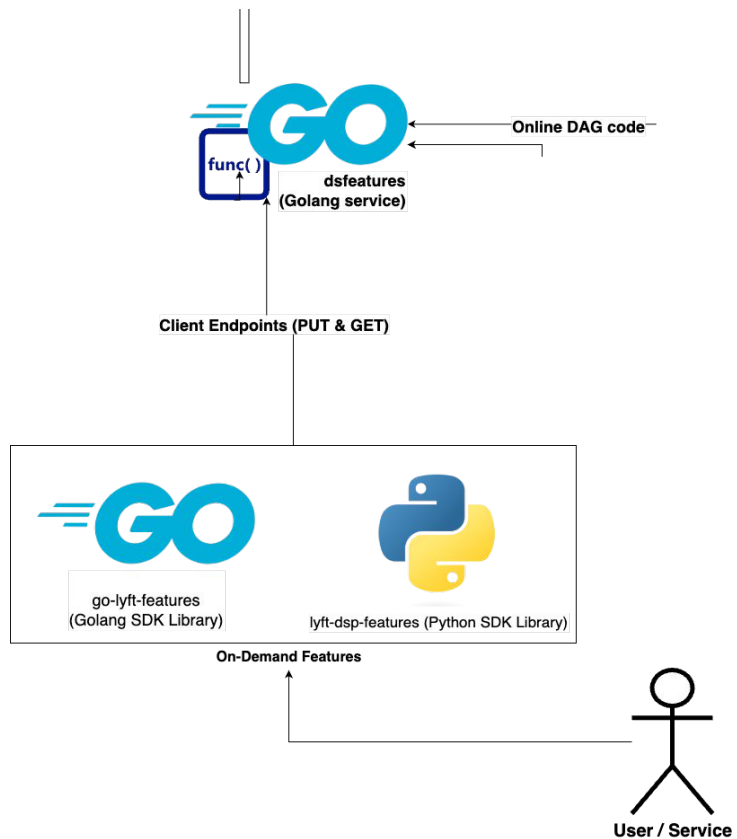| Team / Function | Description | Impact |
|---|---|---|
| Fulfillment | Houses ML models to match drivers to rides & generate + rank rider offers. | Empowers dispatch and offering services. |
| Orchestration | Platform for rider & driver incentive programs. Launch ML models for batch campaigns. | Produce incremental driver hour engagement & growth. |
| Pricing | ML models that generate prices, coupons, etc. | 💰💰💰 |
| Integrity | Fraud detection platform. | Potential blocker of fraudulent activity, real-time and retroactive. |
| Growth Platform & Market Expansion | | |

Organized by **HOPSWORKS**

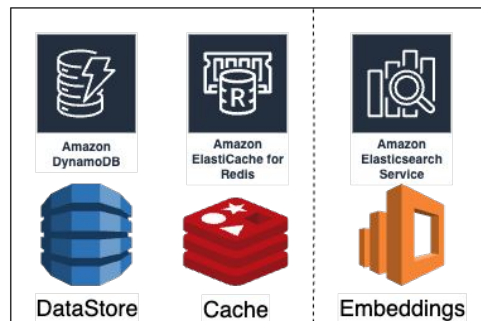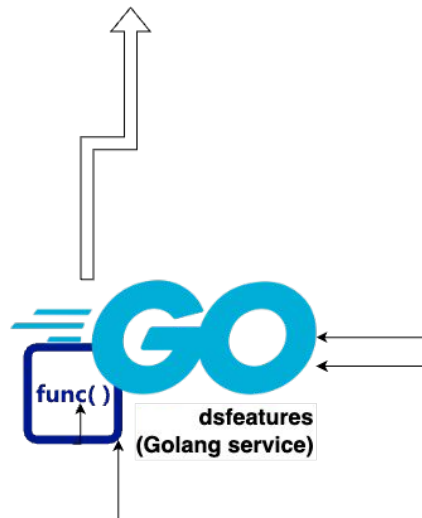| Team / Function | Description | Impact |
|---|---|---|
| Fulfillment | Houses ML models to match drivers to rides & generate + rank rider offers. | Empowers dispatch and offering services. |
| Orchestration | Platform for rider & driver incentive programs. Launch ML models for batch campaigns. | Produce incremental driver hour engagement & growth. |
| Pricing | ML models that generate prices, coupons, etc. | 💰💰💰 |
| Integrity | Fraud detection platform. | Potential blocker of fraudulent activity, real-time and retroactive. |
| Growth Platform & Market Expansion | Comms platform, building "audiences" for targeting marketing messages. | Growth of rider & driver engagement, + incremental ride count & revenue. |

# Architecture

FEATURE STORE SUMMIT 2025

AWS Datastores

Amazon DynamoDB
Amazon ElastiCache for Redis
Amazon Elasticsearch Service

DataStore
Cache
Embeddings

Offline DAG Code

Verity DQ Checks
Amundsen
Table Discovery

Hive Tables (Features Snapshot)

JSON
featurebuilderdata (config repo)

Hive Tables

Airflow DAGs (online/offline)

Online DAG code

dsfeatures (Golang service)

featurebuilder (Python service)

Batch Features (PUT)

Client Endpoints (PUT & GET)

Sink

go-lyft-features (Golang SDK Library)

lyft-dsp-features (Python SDK Library)

On-Demand Features

Apache Beam (spfeatureingest)
Apache Flink

AWS Kafka
Topic
Source

Analytic Events

Streaming Features (PUT)

User / Service

Organized by HOPSWORKS

**Hive Tables (Features Snapshot)**
- Offline DAG Code
- Verity DQ Checks
- Amundsen
- Table Discovery

**Batch Features (PUT)**
- Airflow DAGs (online/offline)
- Online DAG code
- featurebuilderdata (config repo)
- featurebuilder (Python service)
- Hive Tables

Organized by HOPSWORKS

**Online DAG code**

dsfeatures
(Golang service)

func( )

Client Endpoints (PUT & GET)

go-lyft-features
(Golang SDK Library)

lyft-dsp-features (Python SDK Library)

On-Demand Features

User / Service

Organized by **HOPSWORKS**

AWS Datastores

Amazon DynamoDB

Amazon ElastiCache for Redis

Amazon Elasticsearch Service

DataStore

Cache

Embeddings

func( )

dsfeatures (Golang service)

Organized by HOPSWORKS

Developer Experience

**Software Engineer**

FE • BE • Infra

**Data Modeler**

Data Engineer • Data Scientist

**ML Modeler**

ML SWE • Data Scientist

**Researcher**

Data Scientist • Analyst

Organized by **HOPSWORKS**

```
SELECT
    user_id,
    COUNT(requested_at) AS rides_requested_l7,
    COUNT(completed_at) AS rides_completed_l7
FROM core.fact_rides
WHERE ds BETWEEN '{ds}' - INTERVAL '7' DAY AND '{ds}'
GROUP BY 1
```

**SparkSQL
+
?**

Organized by **HOPSWORKS**

```json
{
    "version": 1,
    "entity": "RIDER",
    "feature_group": "rider_stats_l7",
    "query_file": "rider_stats_l7.sql",
    "pagerduty_emails": ["rider@lyft.pagerduty.com"],
    "owner_emails": ["rider-science@lyft.com"],
    "tier": "Tier2",
    "amundsen_tags": [
        "rider"
    ],
    "etl_dependency": [
        {
            "type": "partition",
            "schema": "core",
            "table": "fact_rides"
        }
    ],
    "online_access": true,
    "offline_access": true,
    "features": {
        "rides_requested_l7": {
            "description": "Number of rides requested by the rider in the last 7 days",
            "type": "int",
            "has_personal_data": false
        },
        "rides_completed_l7": {
            "description": "Number of rides completed by the rider in the last 7 days",
    "verity_checks": [
        {
            "type": "check",
            "id": "7-abcdefgh-1234-5678-a1b2-ijklmnopqrst",
            "needs": 
        },
        {
            "type": "check",
```

```sql
SELECT
    user_id,
    COUNT(requested_at) AS rides_requested_l7,
    COUNT(completed_at) AS rides_completed_l7
FROM core.fact_rides
WHERE ds BETWEEN '{ds}' - INTERVAL '7' DAY AND '{ds}'
GROUP BY 1
```

## SparkSQL
## +
## JSON configuration

Organized by **HOPSWORKS**

FEATURE STORE SUMMIT 2025

ML Model Development Life cycle at Lyft

- PROBLEM DEFINITION
- FEATURE & DATA DISCOVERY
- MODEL PROTOTYPING
- MODEL TRAINING
- FEATURE ENGINEERING
- MODEL SERVING
- MODEL MONITORING
- EXPERIMENTATION

Organized by HOPSWORKS

# Evolution

# Batch Features

The biggest family of features at Lyft.

1. Flyte → Self-managed Airflow → Astronomer Ⓐ
2. Deprecate HiveQL & Redshift usage
3. Support for staging (unlocks prototyping & testing)
4. Error Handling (explicit categorization of value type)
5. New Golang library & Offline SDK
6. Data Contracts

# Streaming & Real-Time Features

The growing potential of fresh features.

1. Added OpenSearch integration for use of Embeddings
2. "RealtimeMLPipeline" interface for developing Flink streaming apps with minimal complexity

# Optimization

# Data Generation

How to make a platform of platforms more lightweight?

1. Support for more DAG scheduling windows
2. Monitoring & Dashboards for Failed DAGs (& tasks)
3. Deprecation of underutilized features by R/W activity
4. DAG for offline + DAG for online → **merged**
5. Ownership & Observability Propagation

Organized by **HOPSWORKS**

# Data Retrieval

Can we make the life of our customers significantly better?

1. ElastiCache → ValKey
2. cache payload field reduction
3. EKS pod size reduction
4. Retry + Timeout management
5. Metadata + value cache & datastore TTL
6. Redundant network calls + payload processing steps

Organized by **HOPSWORKS**

# Performance

**8.2ms → 5.5ms**

**-32.9%**

Avg. P95 Read Latency, Oct. 2024 → Oct. 2025

**2,465 → 2,763**

**+12.1%**

Num. Batch Features, Oct. 2024 → Oct. 2025

**62 → 77**
**+24.2%**

Num. Distinct Callers, Oct. 2024 → Oct. 2025

**Q&A**

https://www.linkedin.com/in/rohanvarshney/

FEATURE STORE
SUMMIT
2025