# Agenda

- What is Data Discovery?

- Challenges in Data Discovery

- Introducing Amundsen

- Amundsen Architecture

- Impact and Future Work

# What is Data Discovery?

# Data is used to make informed decisions

| Analysts | Data Scientists | Product Managers | General Managers | Engineers | Experimenters |

*Make data the heart of every decision*

Data-driven decision making process:

1. Search & find data
2. Understand the data
3. Perform an analysis/visualization
4. Share insights and/or make a decision

# What is Data Discovery?

Consider a data-driven decision making process:

1. Search & find data    ***Data Discovery***
2. Understand the data
3. Perform an analysis/create a visualization
4. Share insights and/or make a decision

# Challenges in Data Discovery

# Hi! I'm a new Analyst!

- My first project is predict the attendance for IDEAS conference

- Goal: Help the office team make a decision on number of chairs to provide?

- Idea: Let's take a look into attendance from previous conferences... but where do I look?

# Step 1: Search & find data

- Ask a friend/manager/coworker

- Ask in a wider Slack channel

- Search in the Github repos

We end up finding tables: `hosted_events`

that seems to be the right one

# Step 2: Understand the data

- You find several columns that might be what you're looking for:

  - `booked`, `registered`, and `attendance`

- But you still have many questions such as:

  - Does `attendance` include staff?

  - What's the difference between `booked` and `registered`?

  - How accurate are these figures?

# Step 2: Understand the data

- Look for further documentation on these columns

    - Where does this documentation live?

- Ask an expert who knows this table

    - Who is an expert?

- Run some queries to try to figure it out at the risk of being wrong

```
SELECT * FROM schema.host_events

LIMIT 100;
```

# Nearly 1/3 of Data Scientist time is spent in Data Discovery

- **Data discovery is a problem** because of the lack of understanding of what data exists, where, who owns it, & how to use it.
- Data Discovery provides little to no intrinsic value
- **Impactful work happens in Analysis**

**Data Scientists Time Spent**

# Introducing Amundsen

# What is Amundsen?

- Built at Lyft, official launch in late 2018

- Inspired by Google Search, Airbnb Data Portal, and

  Apache Gobblin

- Named after Norwegian explorer Roald Amundsen

  - Led the first expedition to the South Pole

  - Led the first expedition through the Northwest Passage
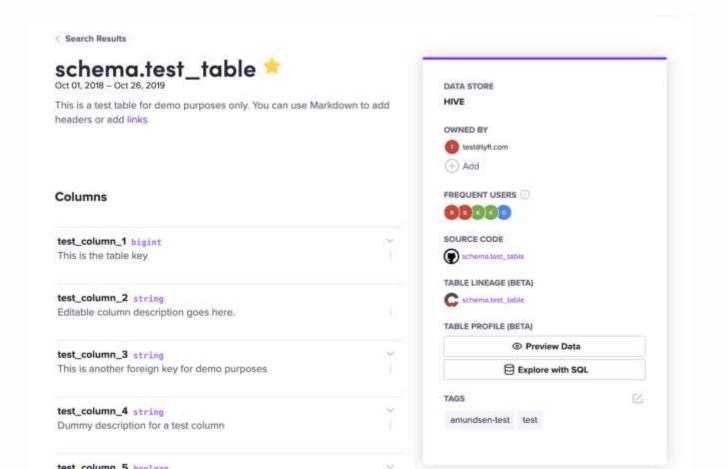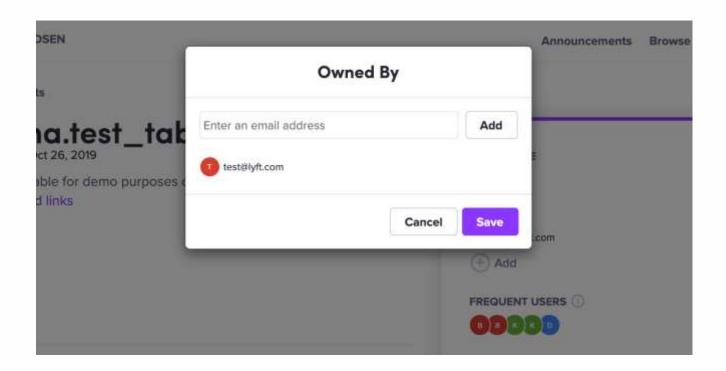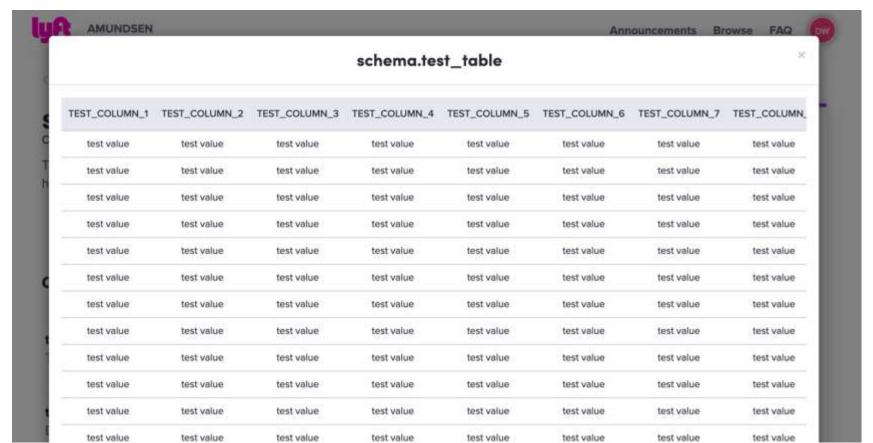
# Home Page

# Search

# Resource Metadata

# Resource Ownership

# Data Preview

# Computed Column Statistics

**test_column_4** `string`

Dummy description for a test column

| COUNT | COUNT_NULL | COUNT_DISTINCT | LEN_MAX |
|---|---|---|---|
| 1234567 | 0 | 22 | 3 |

| LEN_MIN | LEN_AVG | LEN_SUM |
|---|---|---|
| 3 | 3 | 1234567 |

*Stats reflect data collected on Oct 26, 2019 only. (daily partition)*

**test_column_5** `boolean`

*Disclaimer: these stats are arbitrary.*

# Requesting Descriptions



**Amundsen Resource Request** ✕

From

test-user@lyft.com

To

table-owner@lyft.com

Request Type

☐ Table Description

☑ Column Descriptions

Additional Details

Description requested for column: test_column_1
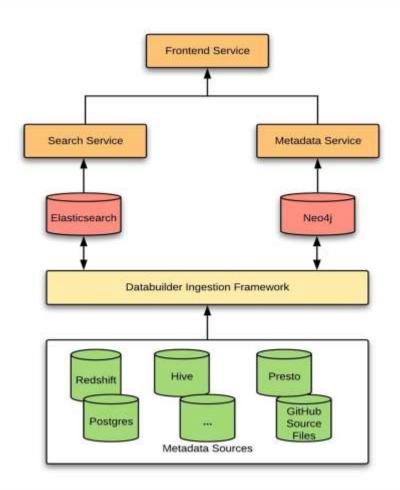
**Send Request**

# User Profile

# In-Application User Feedback

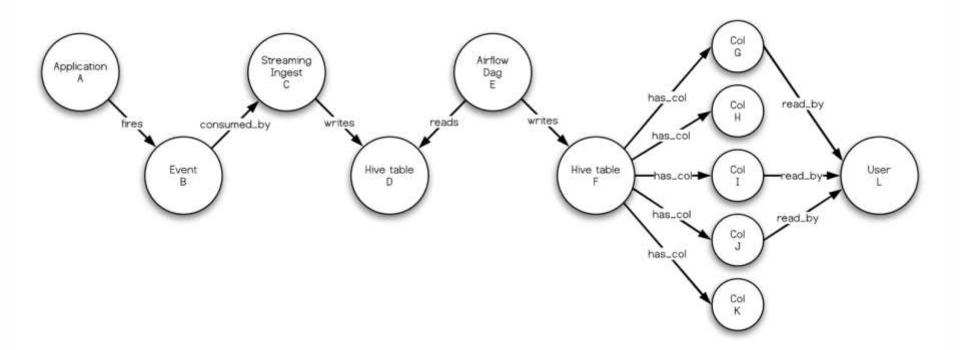# Amundsen Architecture

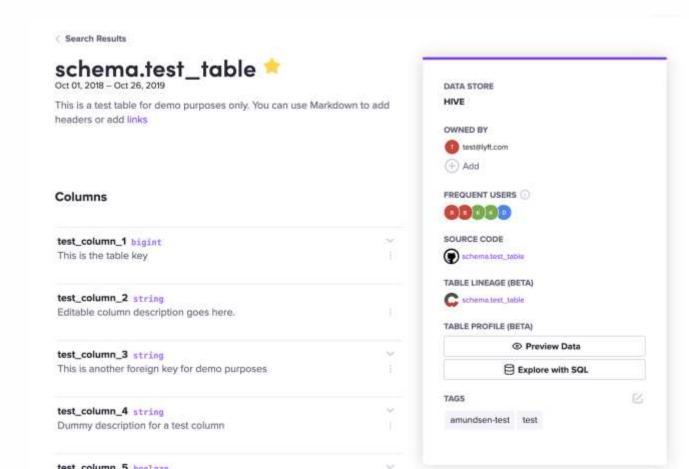# Amundsen Architecture
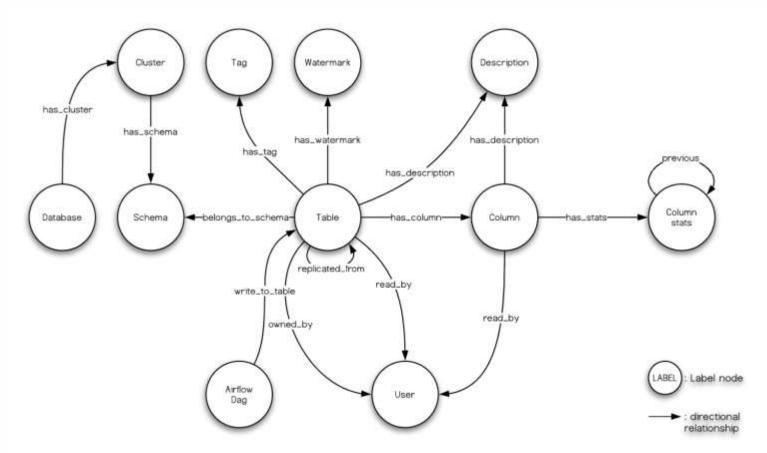
# **Why** choose a graph database?
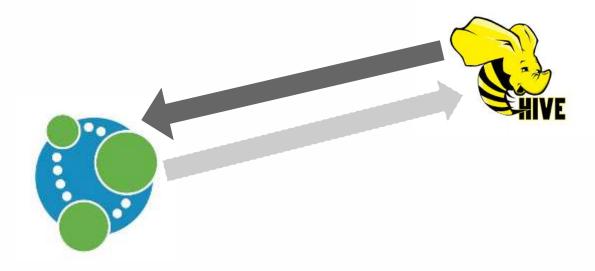
# Why Graph database? (1/2)
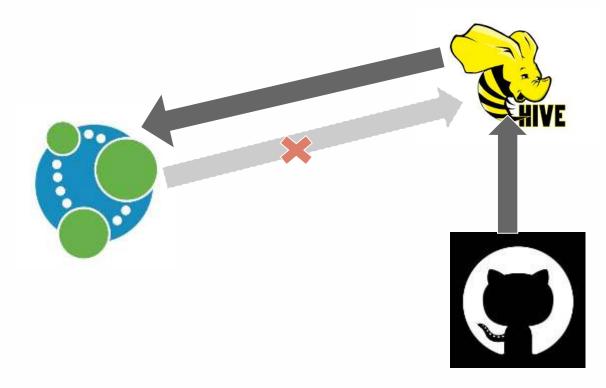
# View Resource Metadata

# Why Graph database? (2/2)

# Neo4j is the source of truth for editable metadata

# Why not propagate the editabled metadata back to source

# Why not propagate the editabled metadata back to source

# Why not propagate the editabled metadata back to source

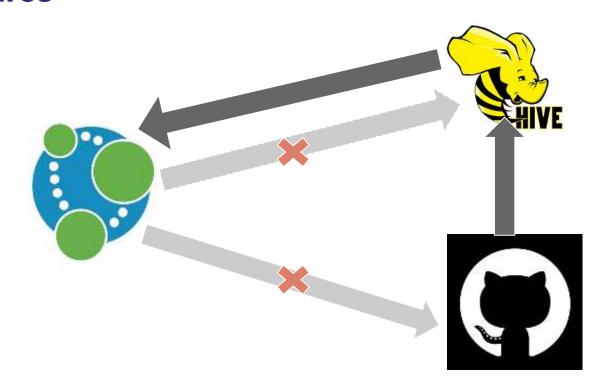# Why not propagate the editabled metadata back to source
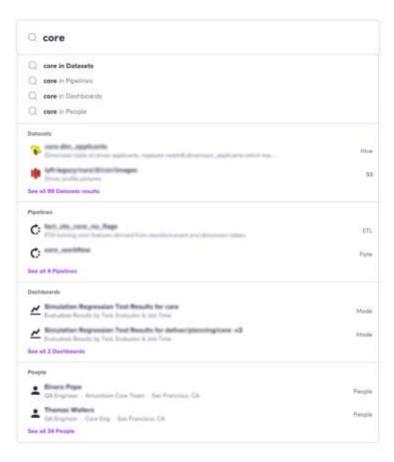
# Impact at Lyft

# Amundsen's Impact at Lyft

- Deployed at Lyft for over 1 year

- Over 700 Weekly Active Users

- 90% penetration among Data Scientists

- Reduced mean time to discovery by 75%

- Also used by Data Eng, Software Eng, PMs, Ops, Marketing Managers, and more
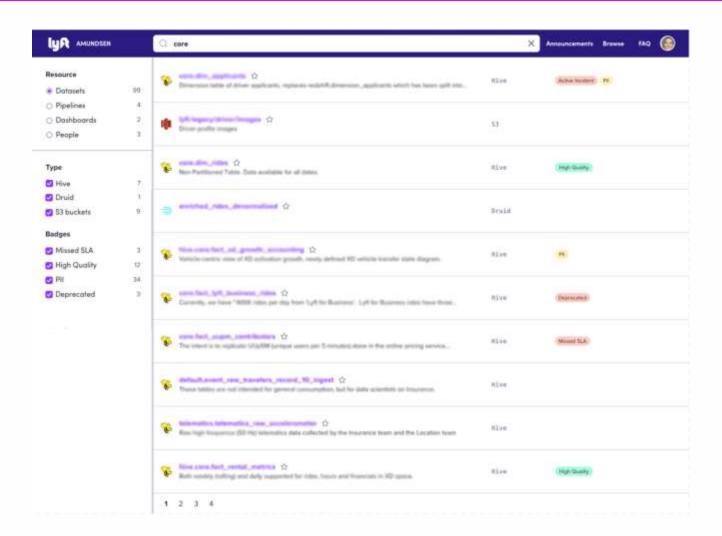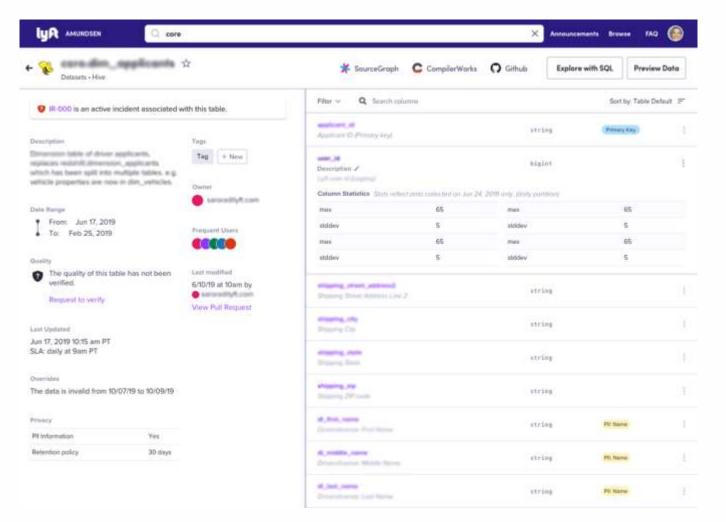
# Future Work

# Search Preview

# Advanced Search

## More Metadata

# We're Open Source

# Amundsen is Open Source!

- [github.com/lyft/amundsen](github.com/lyft/amundsen)

- 200+ github stars, 10+ companies contributing back

- Slack channel 250+ people from 30+ companies

- Presented at conferences in San Francisco, Barcelona, Vilnius, Moscow, LA, NYC by Lyft employees and community

# Community Overview

# Thank You

**Alagappan Sethuraman | /in/alagappanut**
**Daniel Won | /in/danwon**

**Project Code @** github.com/lyft/amundsen