

An aerial photograph of a multi-level highway interchange with several ramps and overpasses. The scene is captured from a high angle, showing the complex geometry of the roads. A large white rectangular box is superimposed over the center of the image, containing the main title and authors' names. The Uber logo is positioned in the lower-left corner of this white box. The background image shows asphalt roads, concrete structures, and some sparse vegetation like palm trees.

# Moving the world with Presto

Girish Baliga & Atul Gupte

Uber

# About us



**Girish Baliga**

Engineering Manager, Interactive Analytics



**Atul Gupte**

Product Manager, Product Platform

Uber's mission is to  
**ignite opportunity** by  
setting the world in  
**motion.**

**600+**

Cities

**75M**

Monthly Riders

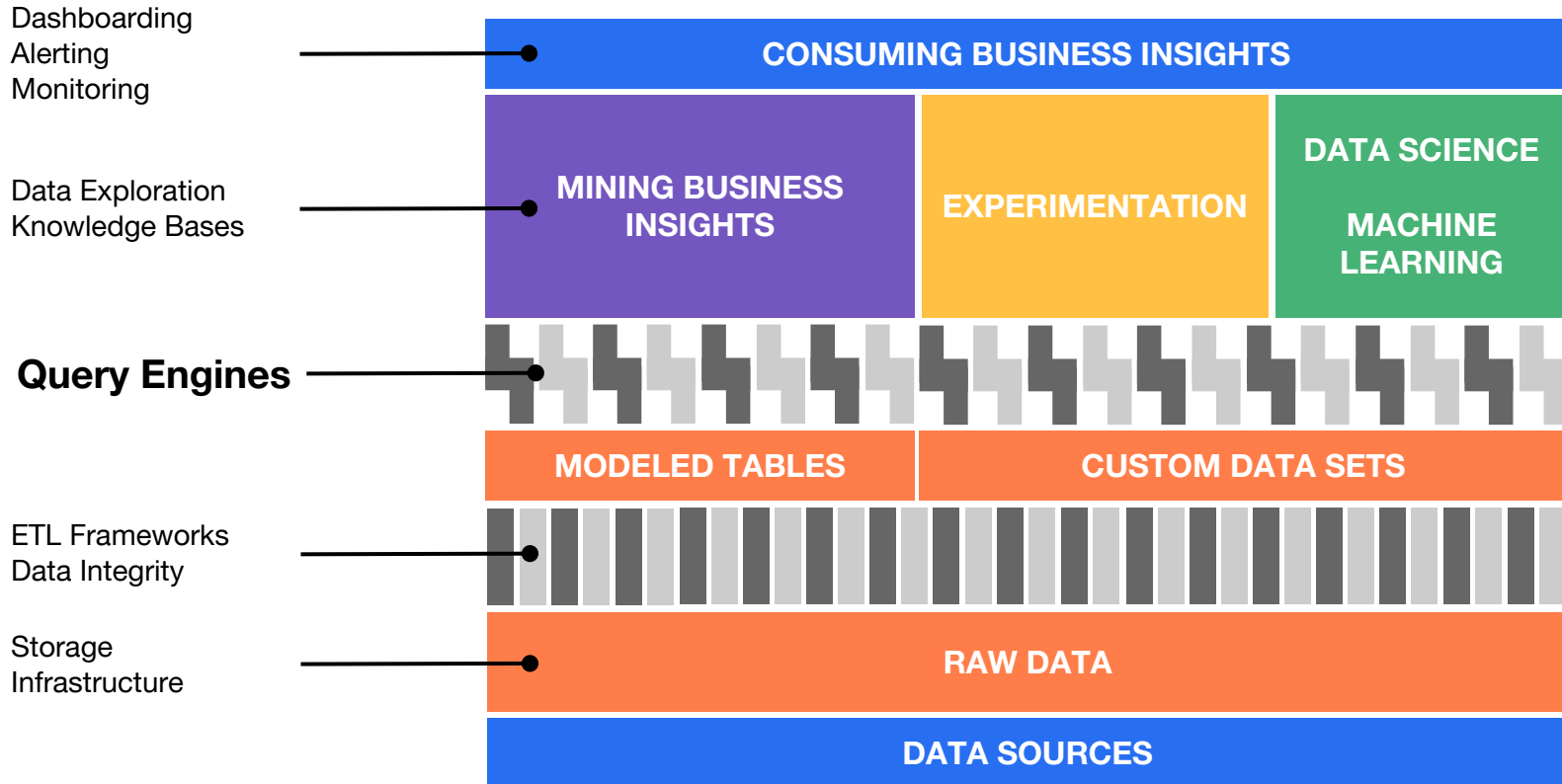
**15M**

Trips/Day

Data informs **every decision** at the company



# Overview of Uber's Data Platform





# Presto @ Uber-scale

**3**

Data Centers

**5K**

Weekly Active Users

**700M**

HDFS files read/day

**2K**

Nodes

**160K**

Queries/day

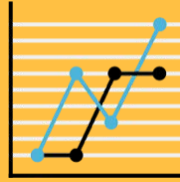
**10PB**

HDFS files  
processed/day

# Presto use cases at Uber



**Marketplace  
Pricing**



**Growth Marketing**



**Data Quality**



**Community  
Operations**



**Data Science**



**Ad-hoc Querying**

# The **people** who rely on us



**Reliant Rebecca**

⚙️ Limited SQL  
Spreadsheets

👥 Marketing Managers  
Entry-level Analysts  
General Managers  
Product Managers



**Monitoring Matt**

⚙️ Intermediate SQL  
Spreadsheets  
Dashboarding

👥 City Operations  
Regional Managers



**Analyst Anna**

⚙️ Advanced SQL  
Spreadsheets  
Limited Statistics  
Limited Python/R

👥 Operations Managers  
Data Analysts  
Product Analysts



**Inventor Ivan**

⚙️ Advanced SQL  
Advanced Statistics  
Scala/Spark, Python/R  
Data Modeling

👥 Data Scientists  
Software Engineers  
ML/AI Researchers

**Technical Skills**



# Exploratory ML & model-training



Engineers

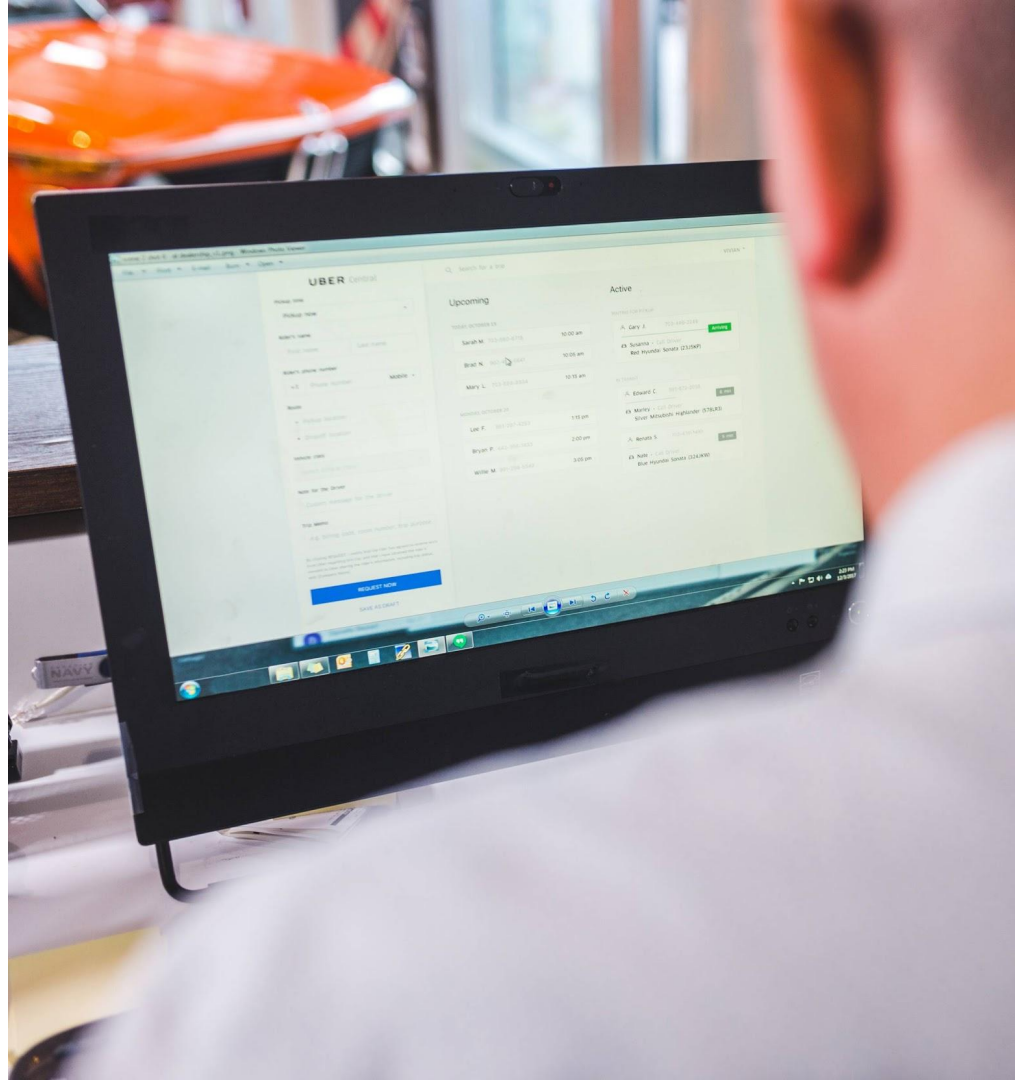


Data Scientists



ML Researchers

Using ML to ensure data security and compliance



# Advanced data science & complex analytics



Data Scientists



Ops Analysts



Support Agents

Surfacing hidden insights to empower restaurants



# Business process automation



Ops Managers



S&P Analysts



Contractors

Using technology to make transportation safer



# Uber Contributions

# Contributions

## Connectors

- Parquet reader
- Elasticsearch connector (release 217)
- Pinot connector enhancements (in-house)

## Optimizations

- Geospatial indexing and operations - 10x or more speedup
- Parquet nested column pushdowns (project, predicate) - 5x speedup

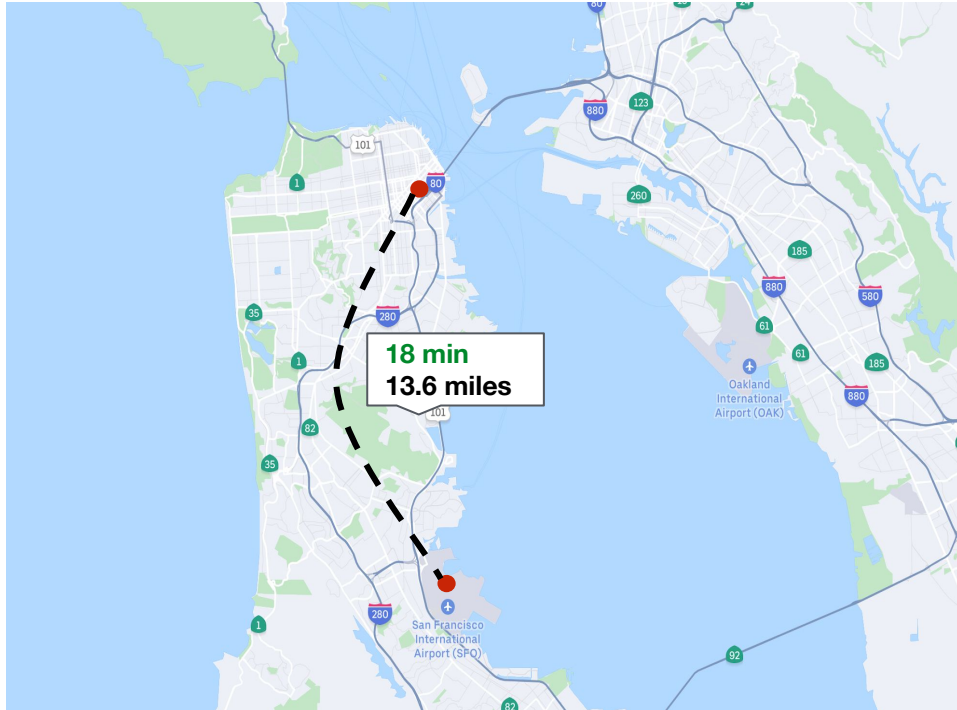
## Production

- Metastore authentication support for Kerberos deployments
- Dispatch Proxy using HTTP redirect for multi-cluster operation

**How many people  
took an Uber to an  
airport last week?**



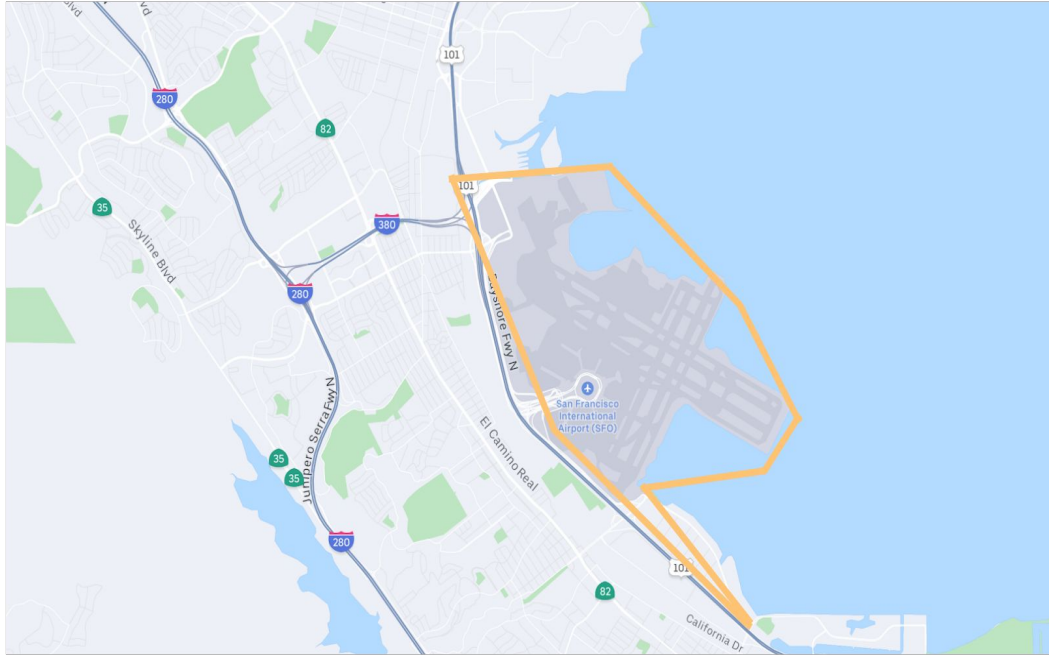
# What does a trip to SFO look like?



## POLYLINE

```
POINT(-73.977540378232 40.756634232442)  
POINT(-77.012723937220 38.984620925742)  
POINT(-112.326839      33.75663)  
POINT(-70.27646950057  -33.63744488957)  
POINT(-103.1284970562  20.78652)  
(5 rows)
```

# Identifying an “airport”



```
presto:default> select airport_code, simplified_shape from dwh.dim_airport where airport_code = 'SFO';  
airport_code | simplified_shape
```

---

```
SFO | POLYGON ((-122.3835454 37.609474654, -122.362637362 37.591782927, -122.454646 37.59  
(1 row)
```



# Geofence (Polygon) - Cross Join

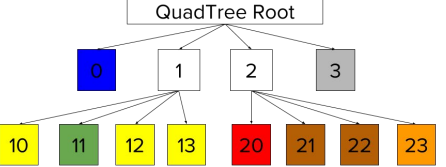
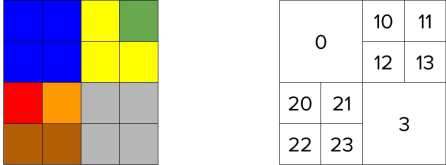
```
SELECT airport_code,  
       COUNT(*)  
FROM rider_events  
JOIN airport  
ON st_contains (geofence,  
               st_point (lng, lat))  
WHERE datestr = '2019-03-11'  
GROUP BY 1
```

event_id	request_lng	request_lat
1	-96.43246	-19.89763
2	-43.9243121	29.084784
3	116.52432433	39.98746
...		
n	-70.3422344	-33.562123

$O(N * M)$

airport	simplified_shape
1. SFO	POLYGON((-122.0..)
2. PEK	POLYGON((3.49...)
...	
m. DMK	POLYGON((-81.7...)

# Geo Index Optimization



event_id	request_lng	request_lat
1	-96.43246	-19.89763
2	-43.9243121	29.084784
3	116.5243243 3	39.98746
...		
n	-70.3422344	-33.562123

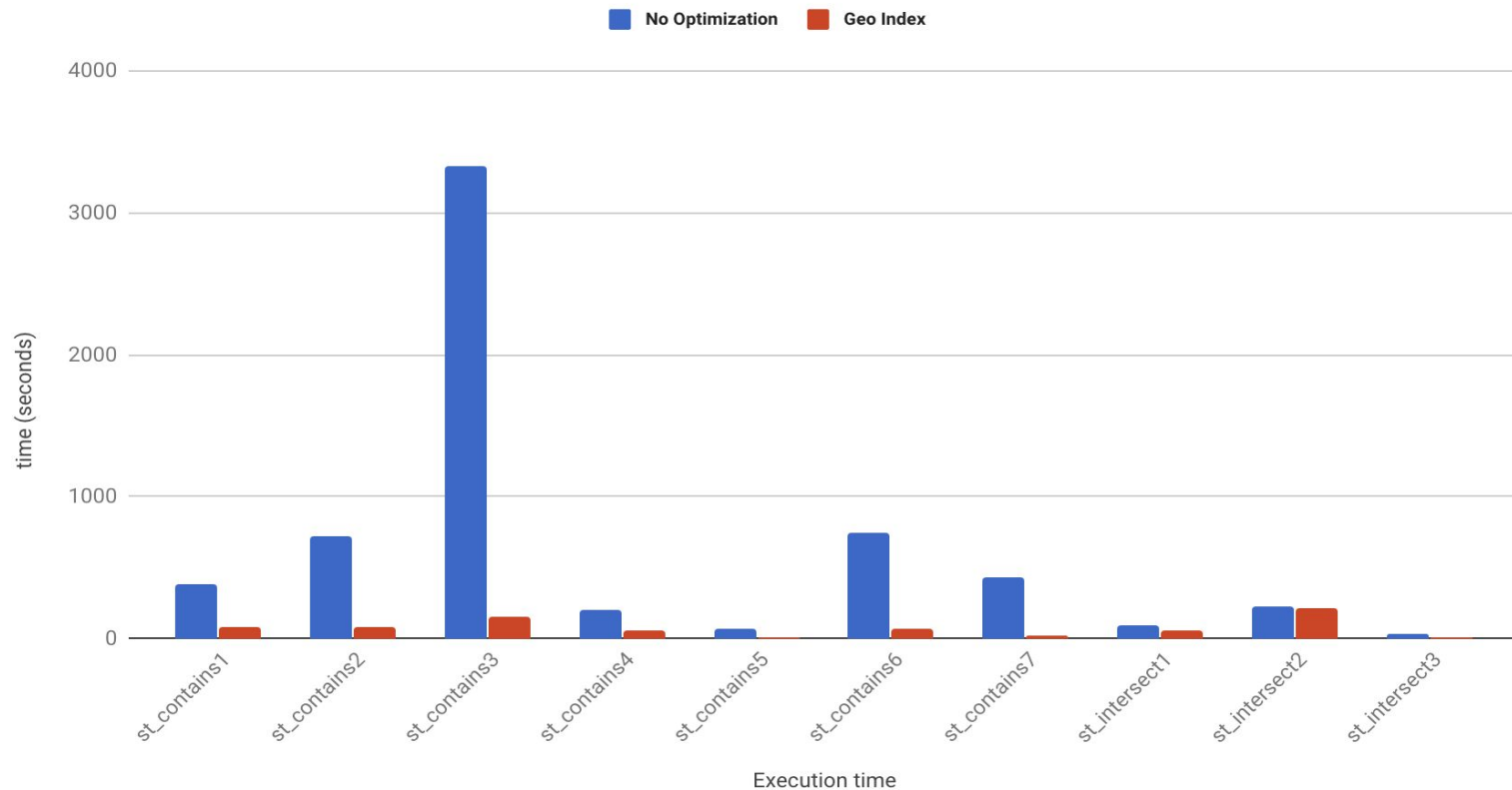
$O(N * \text{Log}(M))$

airport	simplified_shape
1.SFO	POLYGON((( -122.02...))
2.PEK	POLYGON((3.49...))
Fixed number	

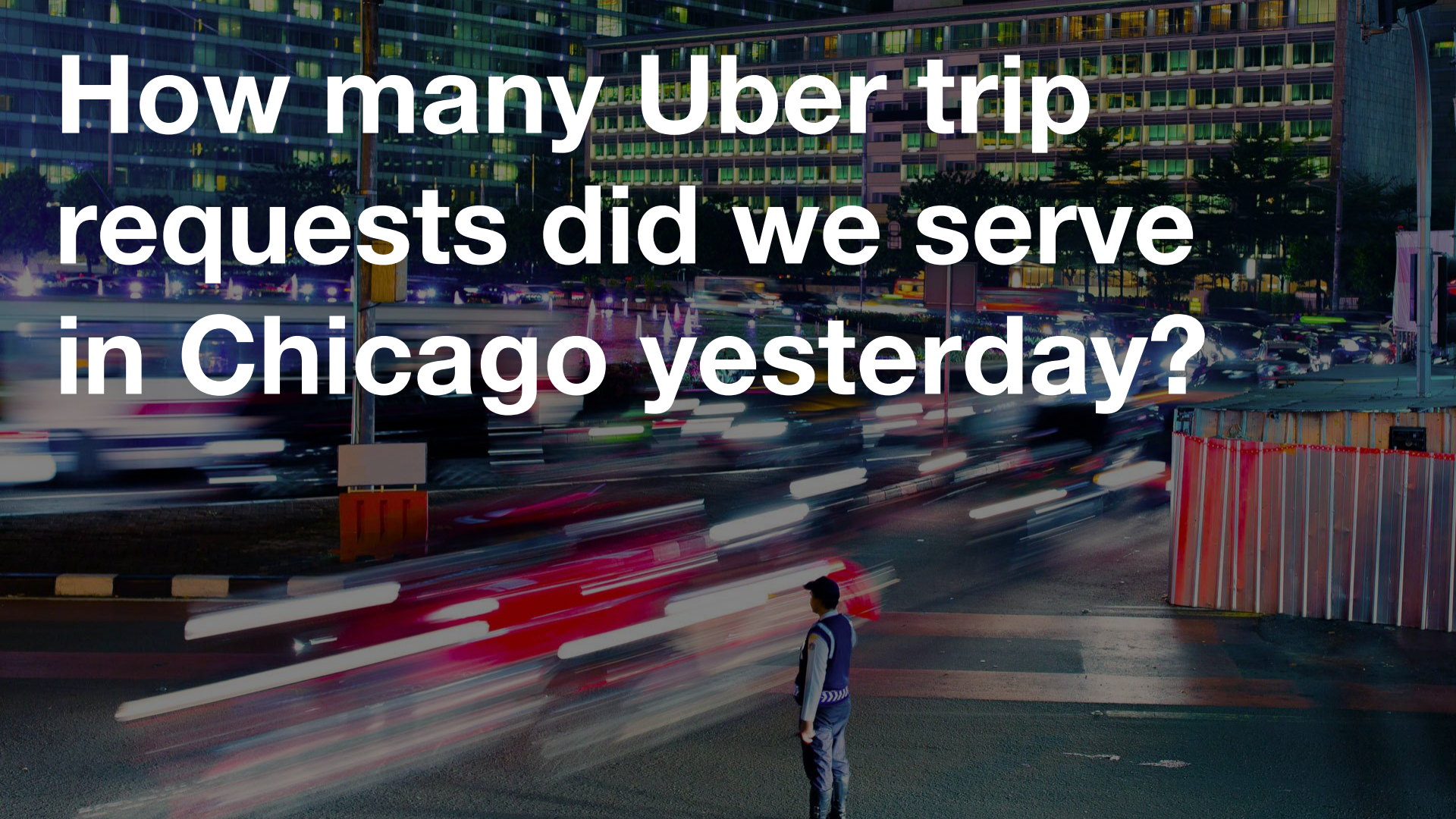
airport	simplified_shape
1. SFO	POLYGON((( -122.02...))
2. PEK	POLYGON((3.49...))
...	
m. DMK	POLYGON((-81.7...))

# Results

## Geo Optimization on Presto



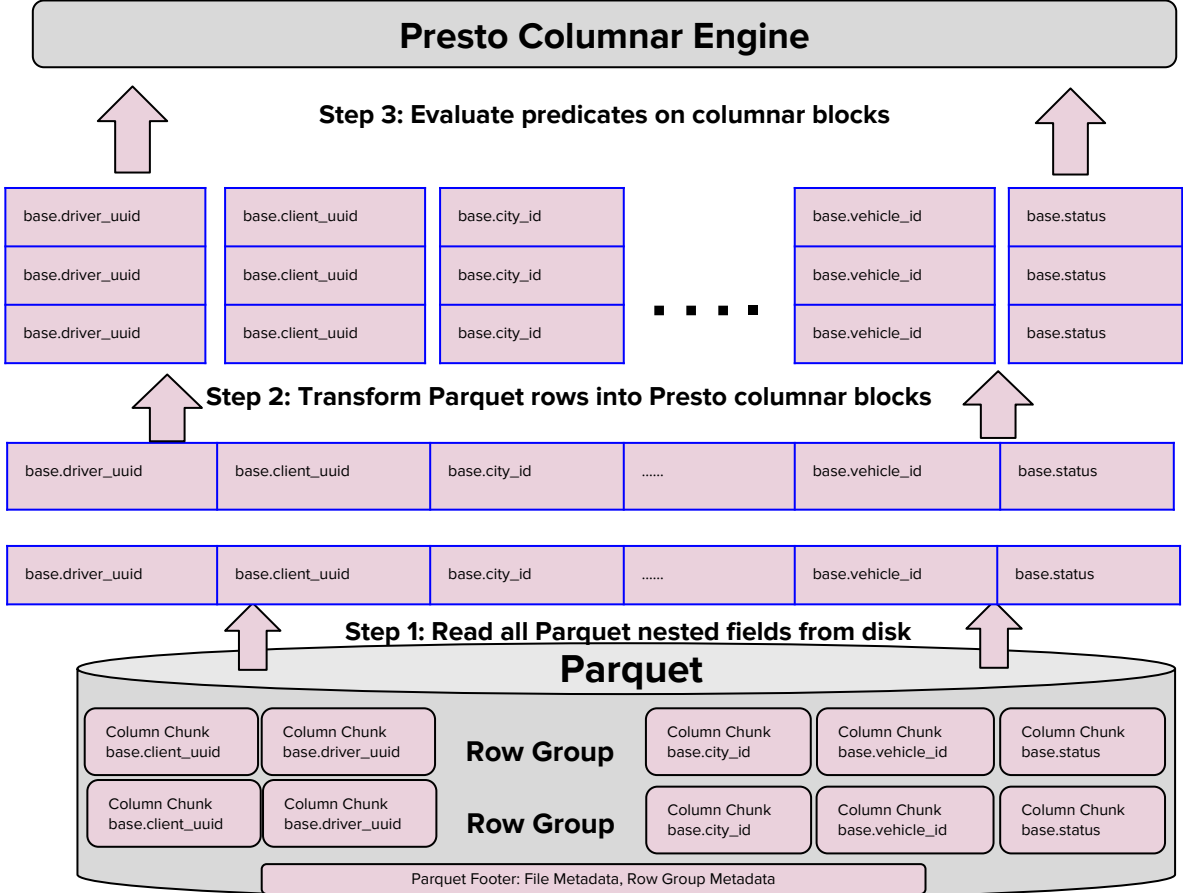
**How many Uber trip requests did we serve in Chicago yesterday?**



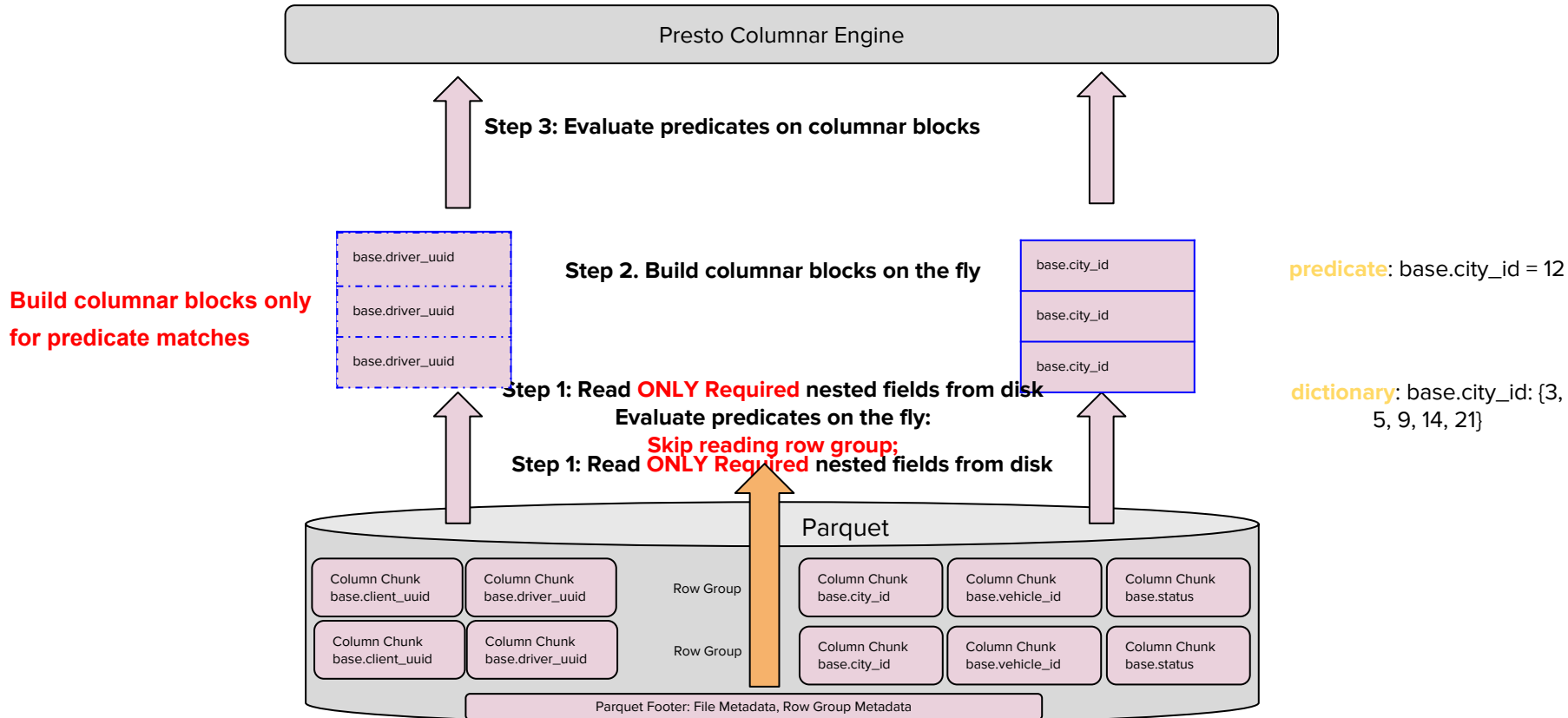
# Fetch daily trip count in seconds

```
SELECT T.base.city_id AS cid,  
       Count(CASE WHEN T.base.status = 'completed' THEN 1 END) AS  
       completed_trips,  
       Count(CASE WHEN T.base.status = 'canceled' THEN 1 END) AS  
       rider_canceled_trips  
  
FROM   trips AS T  
       LEFT JOIN users ON U.id = T.base.id  
  
WHERE  T.datestr = '2019-03-11'  
GROUP BY 1
```

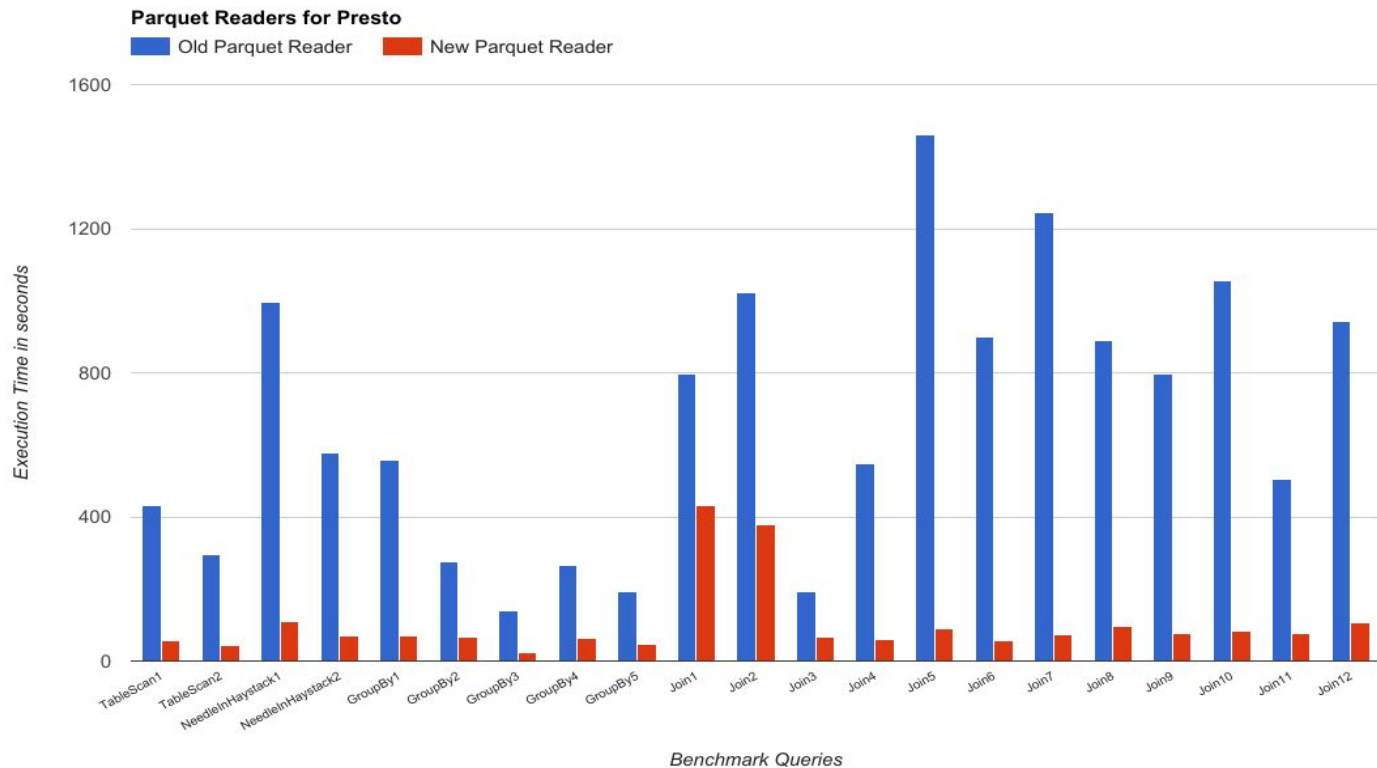
# Default Apache Parquet Reader



# Apache Parquet Reader Optimization



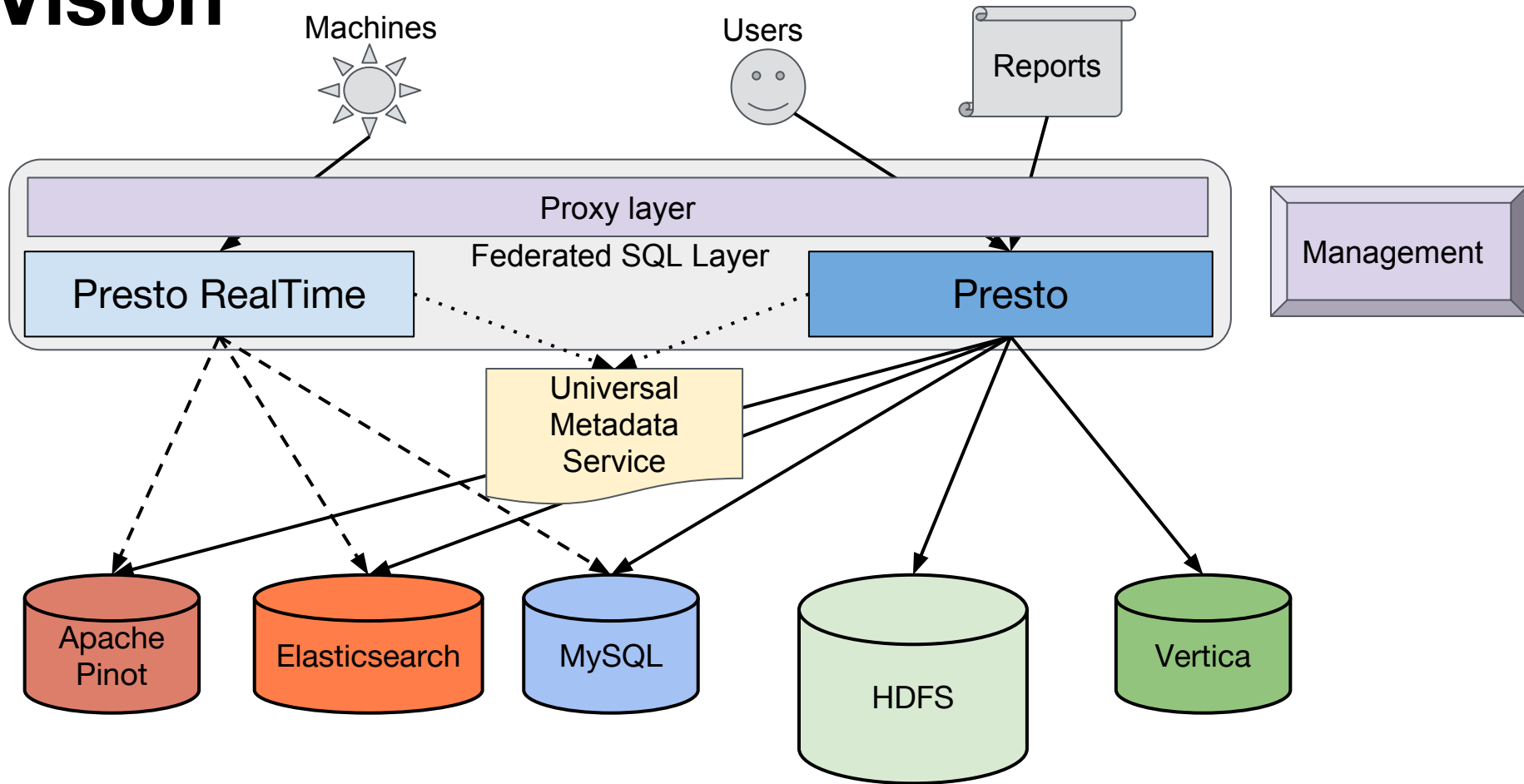
# Results





# Looking forward

# Vision



# Focus areas

## Connectors

- Apache Hive, Apache Pinot, Elasticsearch, Apache Cassandra, Vertica, MySQL, etc
- Aggregation / Join pushdown
- Cross-connector optimizations (hybrid connectors)

## Real-time

- Real-time mode with low latency pass through
- Query plan / result / data cache
- Time-series joins and stitching

## Universal Metadata Service (UMS)

- Logical definitions / physical schemas
- Column stitching and joins
- Table and partition caching

# What's important?

## Reliability

- Production ecosystem
- Deployment and testing @ scale

## Engineering

- Leadership - long term perspective
- Velocity - stay close to head
- Integration hooks - logging, monitoring, security, etc

## Community

- Continuity - multi-year plans
- Collaboration - focus groups & meetups

# Thank you

Proprietary © 2018 Uber Technologies, Inc. All rights reserved. No part of this document may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage or retrieval systems, without permission in writing from Uber. This document is intended only for the use of the individual or entity to whom it is addressed. All recipients of this document are notified that the information contained herein includes proprietary information of Uber, and recipient may not make use of, disseminate, or in any way disclose this document or any of the enclosed information to any person other than employees of addressee to the extent necessary for consultations with authorized personnel of Uber.

