# Visualizing Real Time Data Using ClickHouse and Superset

Srini Kadamati -- Preset
Robert Hodges -- Altinity

1

# Presenter Bios



Robert Hodges - CEO at Altinity

30+ years on DBMS plus virtualization and security. ClickHouse is DBMS #20



Srini Kadamati - Senior Developer Advocate at Preset.io

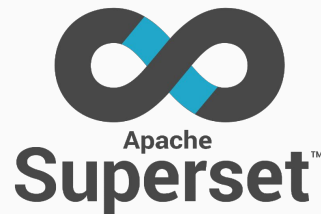6+ years in data science & data science education. Committer to Apache Superset.

Altinity

preset

# Introducing ClickHouse and Superset

## ClickHouse SQL Data Warehouse



- C++ Binary
- Apache 2.0 License
- Active community
- Column storage with compression
- Efficient parallel processing
- Sharding and replication

## Superset Data Visualization and Exploration Platform



- Python + Javascript (Browser)
- Apache 2.0 License
- Active community
- Supports nearly any SQL database
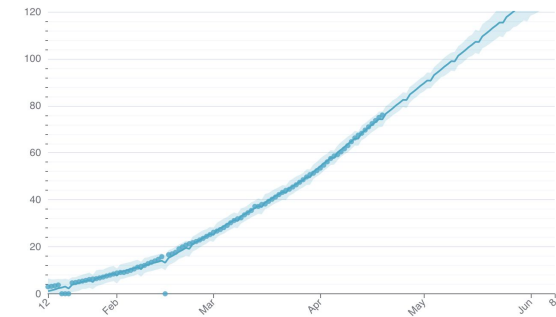- Dozens of chart types
- Preset Cloud managed service

Altinity

preset

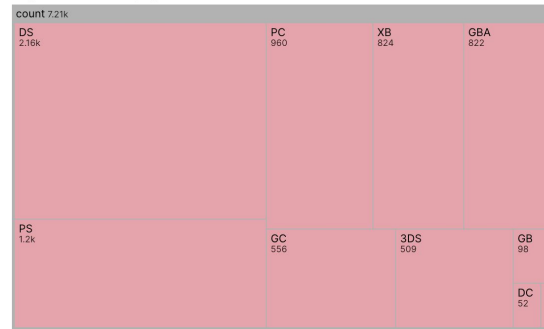# Visualization and data analysis

# The main purpose of data science is insight

# Data + Visualization!

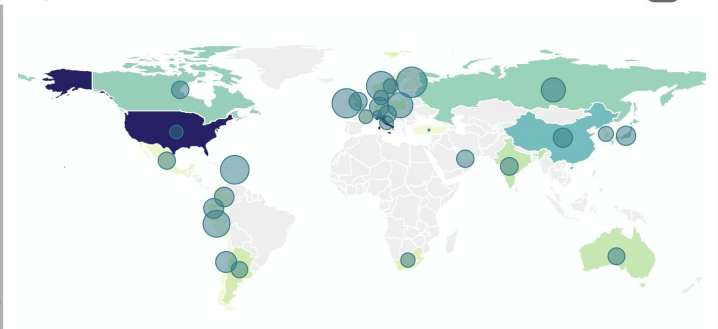Effective data visualization is one of the best ways to get to **insight** from large, complex amounts of data
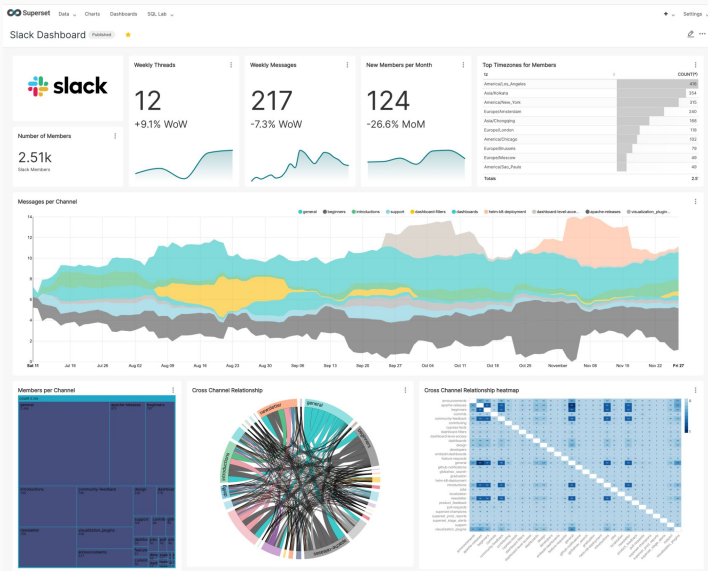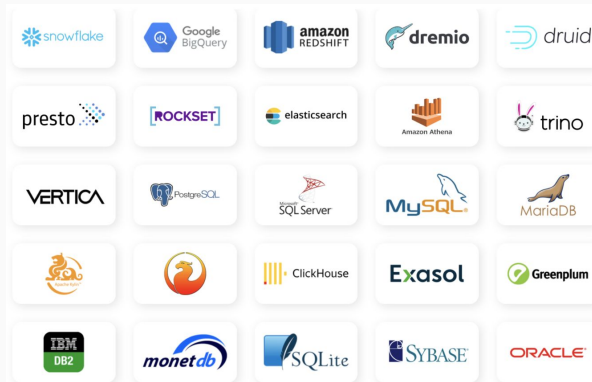
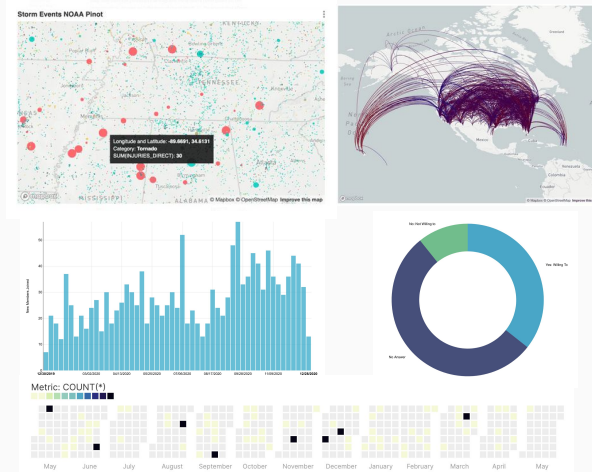# What is Superset & how does it work?

# What is Apache Superset?



Modern open source BI platform

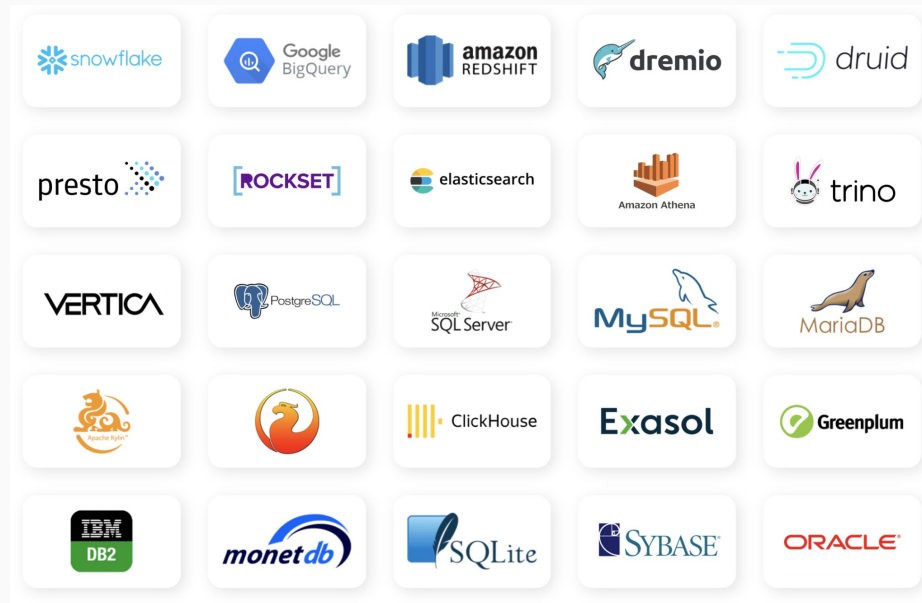Works with nearly any SQL speaking data engine

Large diversity of charts

# Support for SQL Speaking Databases

Superset can query data from nearly any SQL speaking database

Requirements:

- Python DB-API 2.0 driver
  - https://www.python.org/dev/peps/pep-0249/
- SQLAlchemy Dialect
  - Incomplete list here: https://docs.sqlalchemy.org/en/14/dialects/



Blog Post: Building New Database Connectors for Superset

# SQL Lab

Browser based SQL IDE (focused on writing analytical queries)

- Explore data

- Sculpt data for visualization

- Save transformed data as virtual datasets, db VIEWs, or db TABLEs

Altinity

preset

# Superset Semantic Layer

# No Code Chart Builder (Explore)

# Slack Community Dashboard



*Build your own Slack dashboard*

# Introducing ClickHouse and Superset

Altinity

preset

# ClickHouse is a new option for analytic services

An industrial strength, Apache 2.0 SQL data warehouse



| Installs on a laptop in 60 seconds | Runs in cloud, containers, and bare metal | Answers queries in milliseconds |

Altinity

preset

# ClickHouse performance meets or exceeds proprietary SQL data warehouses

Portable C++ binary

Advanced SQL implementation

Column storage with high compression

Distributed, vectorized queries

Built-in sharding and replication

Scales from laptop to 100s of nodes
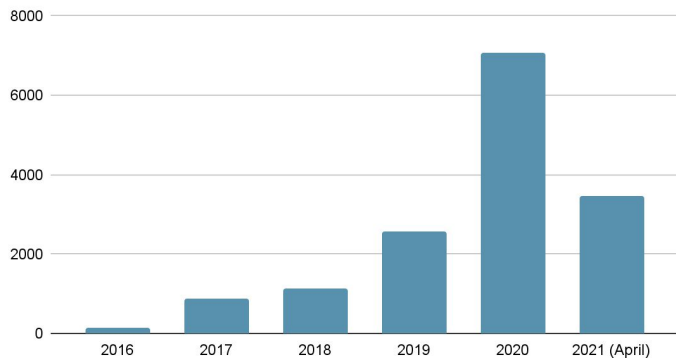
Popular Apache 2.0 licensing

**We use ClickHouse extensively and it's been great.**
John Graham-Cumming, CloudFlare CTO
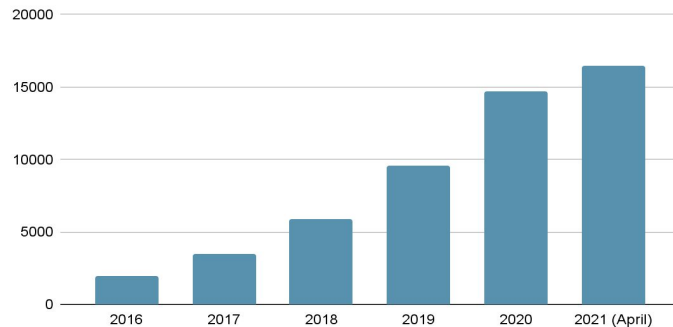Hacker News, 2020

Bloomberg  Flipkart  COMCAST
criteo.  CISCO  INSTANA
CONTENTSQUARE  CLOUDFLARE  LIFESTREET
NuNA  semrush  Spotify  Yandex
SENTRY  Walmart  ebay
Uber  Deutsche Bank

Altinity

preset

# ClickHouse grew from an in-house Yandex project to an international community that rivals ElasticSearch

# Superset connection to ClickHouse

Web Browser



**Superset**

SQLAlchemy

clickhouse-sqlalchemy

clickhouse-driver

**ClickHouse**

(Clear)
:9000

:9440
(TLS)

# ClickHouse SQLAlchemy Drivers (An aside)



sqlalchemy-clickhouse
Apache 2.0
Developed by Marek Vavrusa

- Currently documented in Superset
- Uses ClickHouse HTTP Interface
- No TLS support
- Current pypi.org release:
  0.1.5.post0, Aug 9 2018



clickhouse-sqlalchemy
Apache 2.0
Developed by Konstantin Lebedev

- Supported by Altinity
- Uses ClickHouse Native TCP
- TLS support
- Bug fixes for Superset
- Current pypi.org release:
  0.1.6, Mar 15 2021

Altinity

preset

# Setting up Superset (Ubuntu)

```
# Install Superset
export FLASK_APP=superset
pip install apache-superset
superset db upgrade
superset fab create-admin
superset load_examples
superset init

# Add ClickHouse driver
pip install clickhouse-sqlalchemy

# Start Superset
superset run -p 8088 --with-threads --reload --debugger
```

**Ensure version >= 0.1.6!**

preset™

# ClickHouse connection strings

**SQLAlchemy URL format:**

`clickhouse+native://[user:pw]@host[:port]/database[?options…]`

**ClickHouse on localhost (e.g., your laptop)**

`clickhouse+native://localhost/default`

**ClickHouse public endpoint:**

`clickhouse+native://demo:demo@github.demo.trial.altinity.cloud/default?secure=true`

preset

# Database connection page

# Fun with ClickHouse Data

# Superset dashboard organization

# Creating a time series chart



Create Database → Create Physical Dataset → Create Chart

ClickHouse database

ClickHouse table

ClickHouse query

# Physical dataset creation page



**ClickHouse Connection**

**Database**

**Table**

Add dataset

DATASOURCE

Database: clickhouse clickhouse-public

SCHEMA
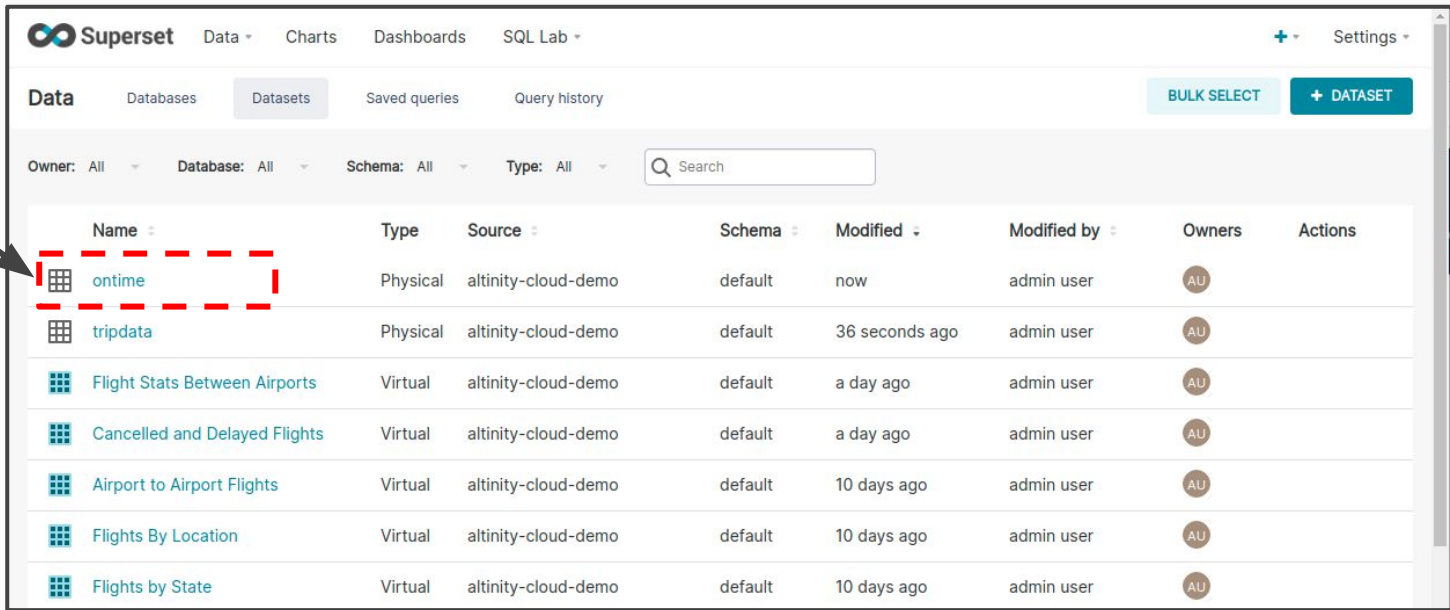
Schema: default

SEE TABLE SCHEMA 5 IN *DEFAULT*

TABLE

ontime

CANCEL     ADD

Altinity

preset

# Starting a chart from a dataset



**Dataset**

Altinity

preset

# Time series chart from physical dataset

# Time series query is generated



```sql
SELECT toStartOfMonth(toDateTime("FlightDate")) AS __timestamp,
       "Carrier" AS "Carrier",
       COUNT(*) AS count
FROM "default".ontime
WHERE "FlightDate" ≥ toDate('2020-01-01')
  AND "FlightDate" < toDate('2021-01-01')
GROUP BY "Carrier",
         toStartOfMonth(toDateTime("FlightDate"))
LIMIT 10000;
```

# Creating a chart on a virtual dataset

**Create Database** → **Create Virtual Dataset** → **Create Chart**

ClickHouse database

ClickHouse query

ClickHouse query

ClickHouse table(s)

ClickHouse subquery
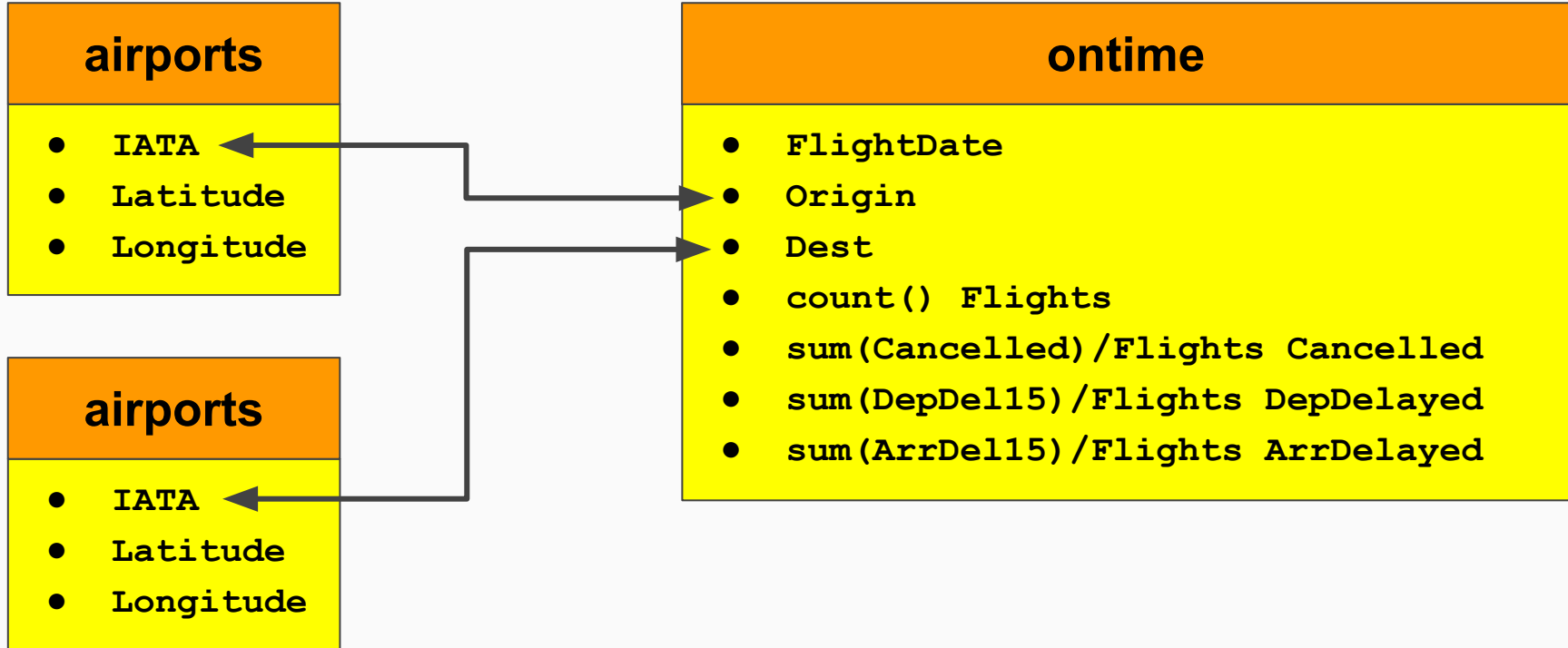
Altinity

preset™

# A query that answers multiple questions

```
SELECT FlightDate, Origin, Dest,
  oa.Latitude as Origin_Latitude, oa.Longitude AS Origin_Longitude,
  od.Latitude as Dest_Latitude, od.Longitude AS Dest_Longitude,
  Flights, Cancelled, DepDelayed, ArrDelayed
FROM (
  SELECT FlightDate, Origin, Dest, count() Flights,
    sum(Cancelled)/Flights Cancelled, sum(DepDel15)/Flights DepDelayed,
    sum(ArrDel15)/Flights ArrDelayed
  FROM ontime
  GROUP BY FlightDate, Origin, Dest ORDER BY FlightDate, Origin, Dest
) AS o
INNER JOIN airports AS oa ON toString(o.Origin) = oa.IATA
INNER JOIN airports AS od ON toString(o.Dest) = od.IATA
```

Altinity

preset

# Using JOIN to add airport LAT/LONG

# Build, run, and save query in SQL Lab

# Use EXPLORE to save as dataset

# Creating a deck.gl Arc chart

**Chart Type***

\* Requires a Mapbox token -- See docs

**Time Dimension**

**Lat/Long**

**Filter**

# How Superset queries virtual datasets

```
SELECT "Origin_Longitude" AS "Origin_Longitude",
       "Dest_Latitude" AS "Dest_Latitude",
       "Origin_Latitude" AS "Origin_Latitude",
       "Dest_Longitude" AS "Dest_Longitude"
FROM
   (
```

**Dataset subquery**

```
   ) AS expr_qry
WHERE "FlightDate" >= toDate('2020-01-01')
  AND "FlightDate" < toDate('2021-01-01')
  AND "Origin" = 'SFO'
  AND "Dest_Latitude" IS NOT NULL AND "Dest_Longitude" IS NOT NULL
  AND "Origin_Latitude" IS NOT NULL AND "Origin_Longitude" IS NOT NULL
LIMIT 5000;
```

**Filters pushed down to base table**

# Creating Word Cloud chart



**Chart Type**

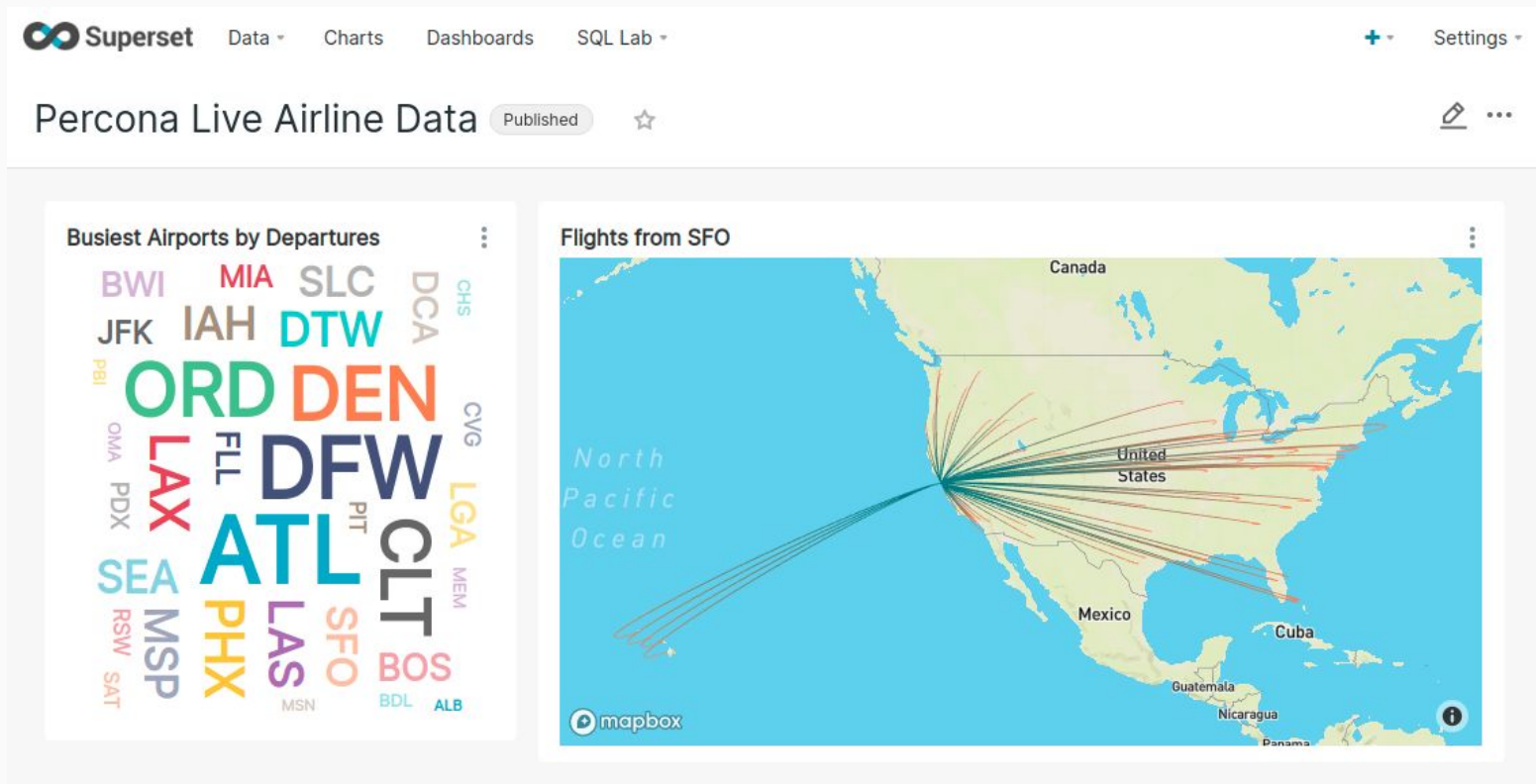**Time Dimension**

**Series**

**Metric**

**Limit & Sort Order**

# Putting charts in a dashboard

# Building from here

- There are *many* more charts available in Superset
- You can build complex dashboards
- Superset caching and clustering can amortize query overhead

# Plans for the future

Altinity

preset

# Roadmap for Superset and ClickHouse

- Improve clickhouse-sqlalchemy driver
  - There are still a few bugs :/
- Build out and document usage examples for users
  - [Blog Post on Clickhouse <> Superset](#)
- Full integration between Preset Cloud and Altinity.Cloud
  - Ensure both services available in common set of AWS regions
  - Automatic setup of AWS PrivateConnect between services

Altinity

preset

# Questions?

Thank you!!

Srini Kadamati
## Preset
https://preset.io

Robert Hodges
## Altinity
https://altinity.com