# WHERE'S MY RIDE?

SANKET PATEL

W205 FINAL PROJECT

# WHY REAL TIME ?

- Helps consumers better plan their commute and save time

- A better consumer experience can potentially result in increased revenue for the transit authority

- Real-time and aggregated data can be used for improving bus scheduling
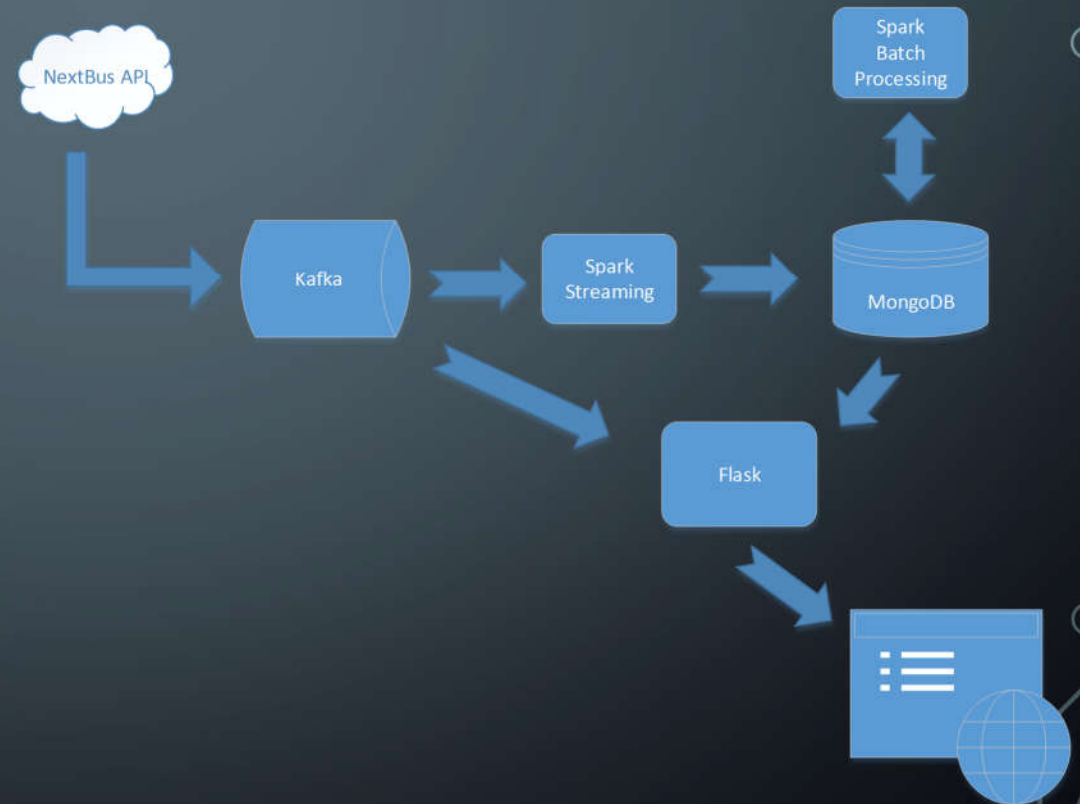
# DATA SOURCES

- Data sources explored

  - NextBus API http://www.nextbus.com/xmlFeedDocs/NextBusXMLFeed.pdf

  - Real Time BART API https://www.bart.gov/schedules/developers

  - OneBusAway http://developer.onebusaway.org/modules/onebusaway-application-modules/1.1.13/api/where/index.html

- Data source used

  - NextBus API http://www.nextbus.com/xmlFeedDocs/NextBusXMLFeed.pdf

# ARCHITECTURE & TECHNOLOGIES

- Lambda Architecture
- All data stored in Mongo DB
- Kafka serves as pipe for all data ingestion
- Streaming layer
  - Spark Streaming job
  - Using Kafka DirectStream and checkpointing to store offsets to achieve exactly-once semantics
- Batch Layer
  - Spark job to aggregate hourly data and store in MongoDB
- Serving Layer
  - Flask
  - SocketIO
  - HTML and Google Maps JavaScript API

# IMPLEMENTATION

- Python script runs as service to query NextBus API every 10 ms and pushes data to kafka

- Spark Stream job cleans the data for storing in MongoDB and also calculates GeoHash for the coordinates

- The streaming data is stored in MongoDB Collection indexed on Geo Coordinates

- A Spark Bath job is kicked off via Cron Job every hour to aggregate data based on GeoHash and stored back into MongoDB collection

- The aggregated data is used for generating HeatMap

- Flask as the serving layer,
  - Implements kafka consumer to get real-time data
  - The real-time data is pushed to WEB UI via WebSockets
  - Serves aggregate data by querying MongoDB aggregate collection

# ALGORITHMS

- GeoHashing
  - Encode(lat,lon) -> "String"
  - Generates same hash for coordinates in a given rectangle
  - The size of the rectangle is determined by the length of the hash
  - Allows easy aggregation based on location

- Other algorithms
  - DBSCAN
    - Works bottom-up by picking a point and looking for more points within a given distance
    - Expands the cluster by repeating this process for new points until the cluster cannot be further expanded
    - Its tuned using 2 parameters:
      - *Epsilon –* This determines how far to search for points near a given point
      - *MinPoints -* This determines how many points should be present in the neighborhood of a given point in order to keep expanding a given cluster

| GeoHash length | Area width x height |
|---|---|
| 1 | 5,009.4km x 4,992.6km |
| 2 | 1,252.3km x 624.1km |
| 3 | 156.5km x 156km |
| 4 | 39.1km x 19.5km |
| 5 | 4.9km x 4.9km |
| 6 | 1.2km x 609.4m |
| 7 | 152.9m x 152.4m |
| 8 | 38.2m x 19m |
| 9 | 4.8m x 4.8m |
| 10 | 1.2m x 59.5cm |
| 11 | 14.9cm x 14.9cm |
| 12 | 3.7cm x 1.9cm |

# TECHNICAL CHALLENGES

- Since the data ingest scripts queries NextBus API every 10ms, it creates bursts of data. Had to make sure streaming job is never overwhelmed by setting *spark.streaming.kafka.maxRatePerPartition*

- Achieving exactly-once semantics

- Serving Real-time locations to Web Page
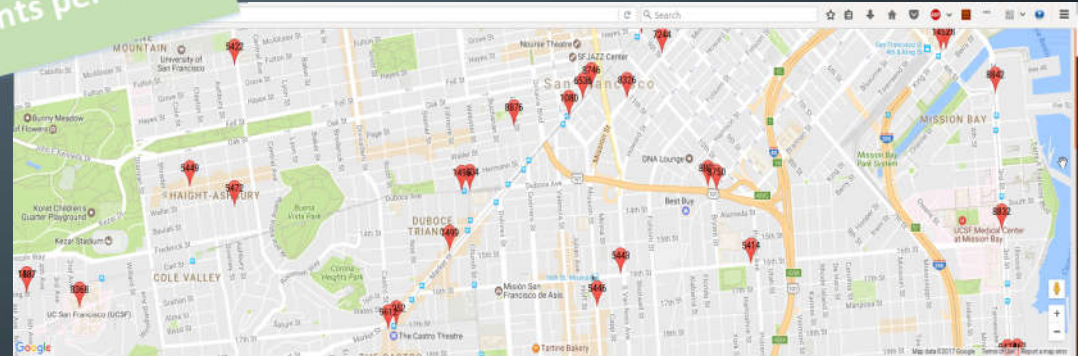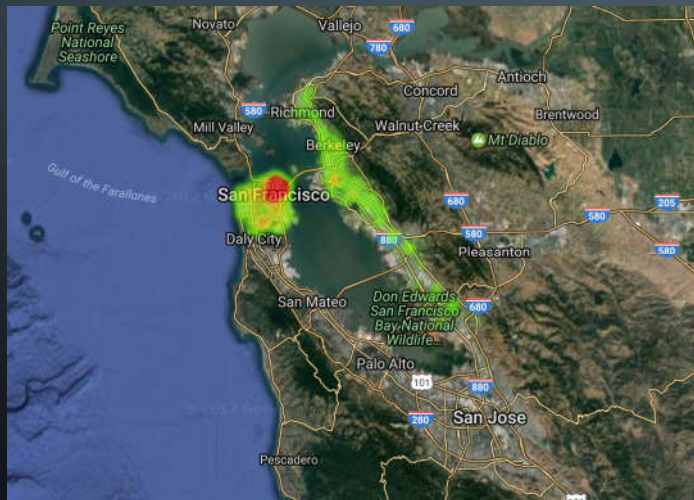
- JavaScript/HTML (not a big fan ☺ )

# RESULTS



40K events per hour

# LIMITATIONS AND EXTENSION

- Limitations
  - Only one client can connect to real-time data. SocketIO can support multiple clients using Rooms and Broadcasting
  - Aggregate job doesn't back fill if system goes down for multiple hours. This can be fixed by having the Cron Job script query MongoDB for missing hours and launch multiple aggregate jobs
  - Only one data source

- Future Extensions
  - Add more data sources (different cities, modes of transport etc)
  - Integrating Uber/Lyft data along with weather data can help consumers make more efficient transport selection