# anywaze

## : analytics on Waze data at scale
### (or how to attack google's servers for fun and {hopefully} profit)

Jason Gors
Insight Fellow
Data Engineering 2016

**anywaze.xyz**

# Motivation:

crowdsourced data == controversial:

- Dynamically reroutes users through neighborhoods for quicker driving directions.
- Users report real-time police locations.



CBS NEWS / November 20, 2014, 10:05 AM

**Traffic app facing speed bumps in quiet neighborhoods**

1 Comment / f 20 Shares / Tweet / Stumble / @ Email



CBSNEWS    Video  US  World  Politics  Entertainment  Hea

Sports  Photos  More

AP / January 26, 2015, 10:20 AM

**Police say Waze cop-tracker is threat to officers**

1.7 km
Dodge Ave

In 400 m

Not there



SLASHGEAR    REVIEWS    COLUMNS    FEATURES    HUBS

Trending    Cars    Science    CES    Apple

**RideWith becomes Waze Rider, still in mysterious limited test**

0

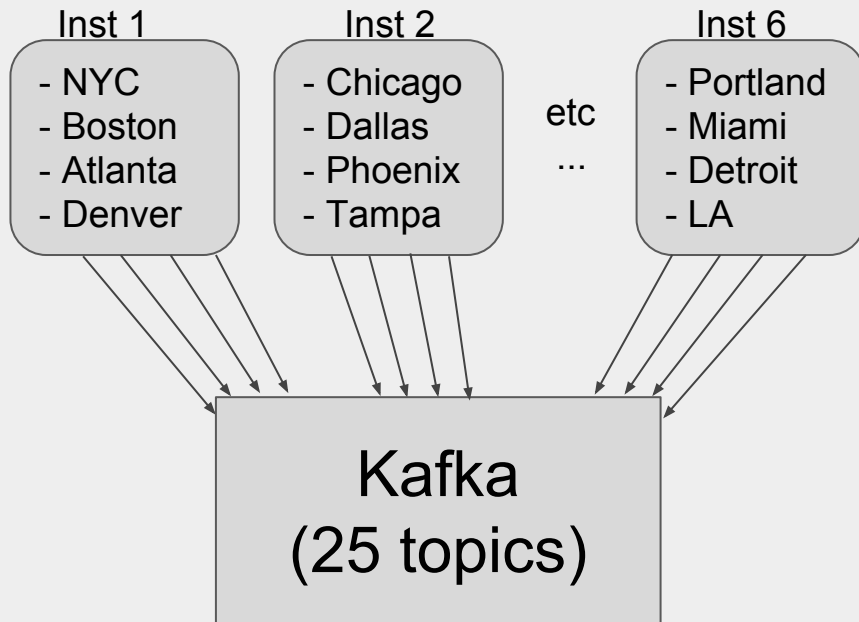JC Torres - Feb 10, 2016

# Data:

- Real-time data from 25 largest U.S. cities:
  - 20 x 20 mile square over city centroid (eg. next slide).
- No official API for Waze:
  - had to be "creative" in getting the data.
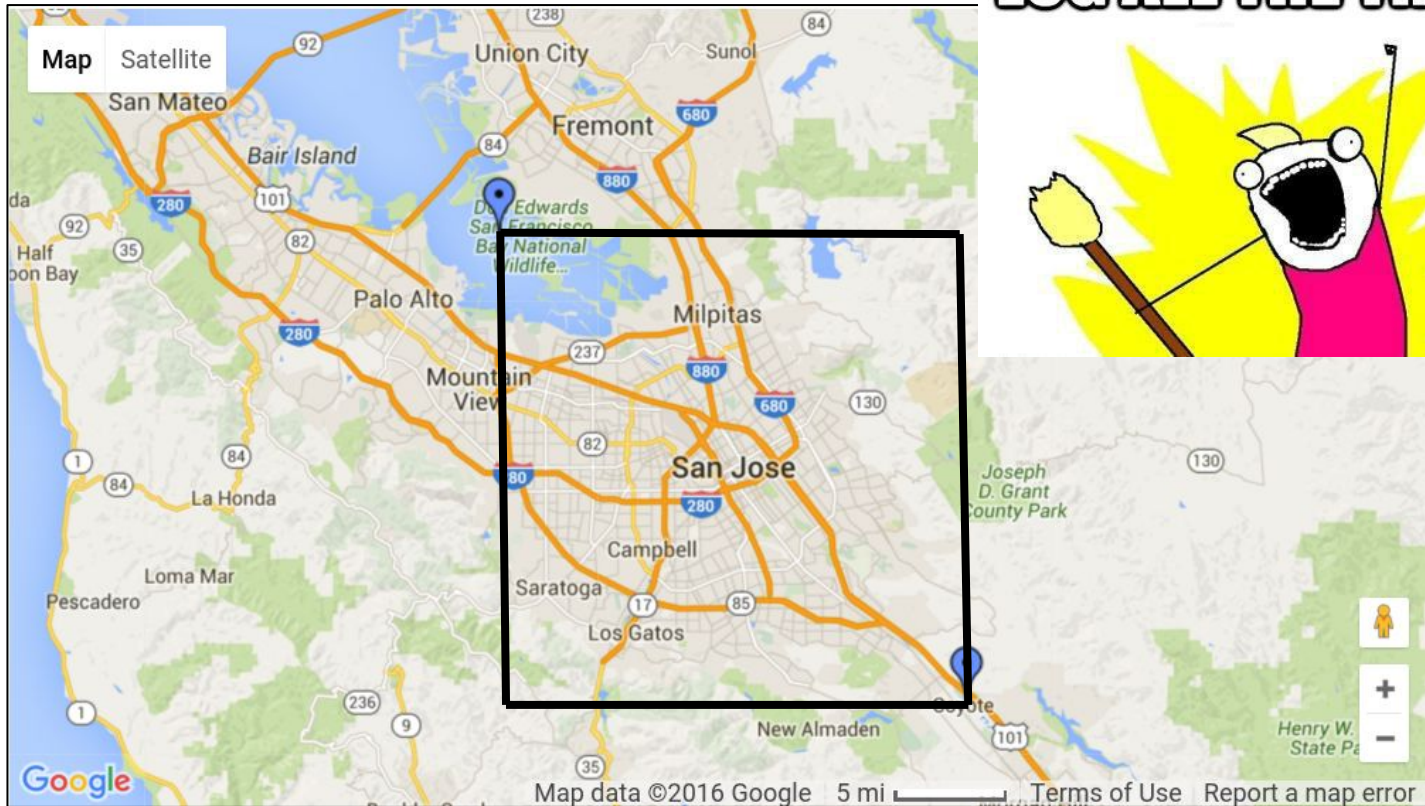  - thus data is very messy!
- Currently ~1.5TB and growing.

# Data:

- Real-time data from 25 largest U.S. cities:
  - 20 x 20 mile square over city centroid (eg. next slide).
- No official API for Waze:
  - had to be "creative" in getting the data.
  - thus data is very messy!
- Currently ~1.5TB and growing.

**Data Collection:**

(read: my **"waze botnet"**)

- 6 AWS instances, each running a web server, each hitting 4 to 5 cities every second:

Inst 1
- NYC
- Boston
- Atlanta
- Denver

Inst 2
- Chicago
- Dallas
- Phoenix
- Tampa

etc
…

Inst 6
- Portland
- Miami
- Detroit
- LA

Kafka
(25 topics)

LOG ALL THE THINGS

```
"time_stamp": 1453696077.374517,
"alerts": [
    {
        "country": "US",
        "latitude": "41.846912",
        "longitude": "-87.643182",
        "numOfThumbsUp": 19,
        "placeNearBy": null,
        "subType": "POLICE_HIDING",
        "type": "POLICE"
    },
    {
        "country": "US",
        "latitude": "41.756379",
        "longitude": "-87.552768",
        "numOfThumbsUp": 5,
        "placeNearBy": null,
        "subType": "",
        "type": "ACCIDENT"
    },
    ...
    {
        "country": "US",
        "latitude": "41.756379",
        "longitude": "-87.552768",
        "numOfThumbsUp": 3,
        "placeNearBy": null,
        "subType": "",
        "type": "JAM"
    },
    {
        "country": "US",
        "latitude": "41.824708",
        "longitude": "-87.719495",
        "numOfThumbsUp": 2,
        "placeNearBy": null,
        "subType": "HAZARD_ON_SHOULDER_CAR_STOPPED",
        "type": "HAZARD"
    }
]
"jams": [
    {
        "country": "US",
        "delayInSec": -1,
        "end": "W Congress Pkwy",
        "endLatitude": "41.875453",
        "endLongitude": "-87.642937",
        "severity": 3,
        "start": null,
```
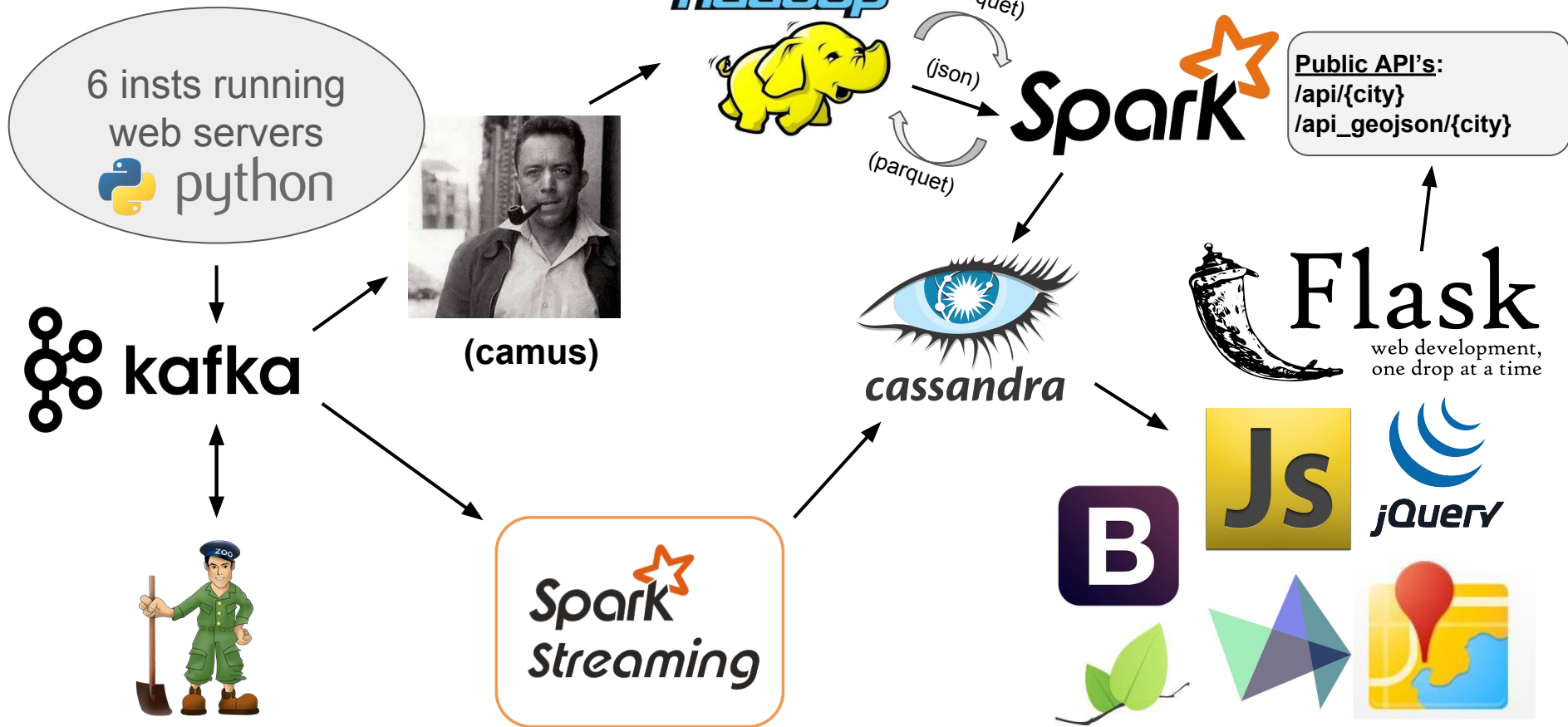
```
{
    "country": "US",
    "latitude": "41.846912",
    "longitude": "-87.643182",
    "numOfThumbsUp": 19,
    "placeNearBy": null,
    "subType": "POLICE_HIDING",
    "type": "POLICE"
},
```
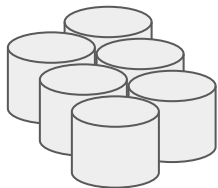
```
{
    "country": "US",
    "latitude": "41.824708",
    "longitude": "-87.719495",
    "numOfThumbsUp": 2,
    "placeNearBy": null,
    "subType": "HAZARD_ON_SHOULDER_CAR_STOPPED",
    "type": "HAZARD"
}
```
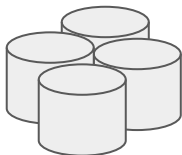
# Pipeline:

# Pipeline (cont'd):

6 instances:
Data collection
web servers

4 instances:
Kafka, Camus,
HDFS, Cassandra

5 instances:
Spark

1 instance:
Flask

c4.xlarge
(network io)

m4.large
(storage)

m4.xlarge
(compute)

m4.large
(general)

# Challenges encountered:

- Data is very messy -- difficult to wrestle into a usable format.
  - Each GET Response from their servers has changing nested json fields.
- Bringing nodes back online after failure.
- Scaling instance storage in real-time (== much pain).
- Adding new nodes into active clusters in real-time.
- How best to represent 25 cities of data in Kafka:
  - 25 topics, 1 Partition each
- Unit & CI testing is hard with these distributed tools.

# About me...

- Previously a Data Scientist
- ...and Software Developer
- Graduate work in Statistics
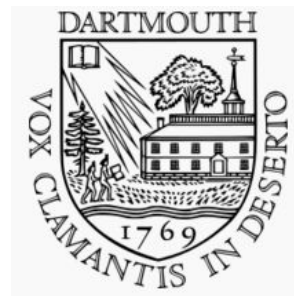- A Bayesian
- Python & Linux evangelist

**Anywaze**, take a look:

**@github**: jgors/anywaze



Questions?