

Highly Scalable Tile-Based Visualization for Exploratory Data Analysis

Strata NY: Hadoop and Beyond, 10/17/2014
David Jonker, Rob Harper



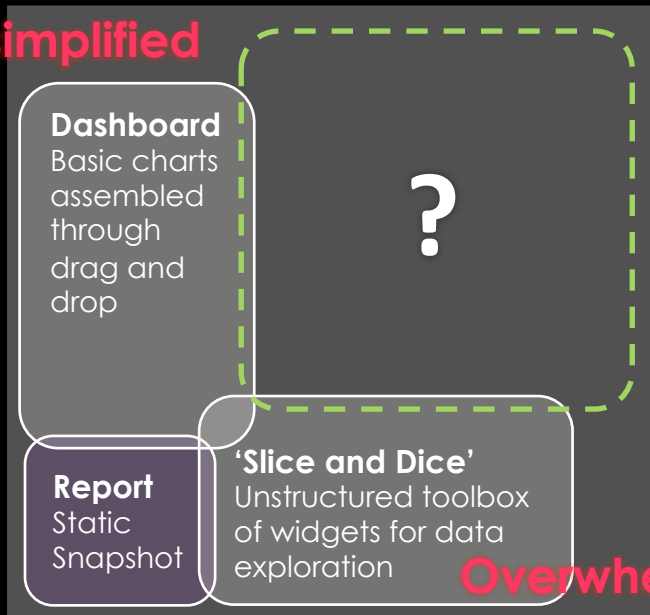
Making Sense of Big Data

“Big Data for us is in complexity.
We have complex data.
We have big complexity.”
– Analyst

“No one wants more data,
everyone wants better data.”
– Analyst

**Less Time
Required**
to acquire
insights

Oversimplified



More Information Available
To make an effective decision

Challenges of Effective Visual Analytics

ESSENTIAL

- Richly **informative** and true to reality.
- Answers easily arrived at, **easily understood**.



ALSO, every business problem is unique, but cannot afford an entirely unique solution for every problem.

- Need relatively universal, **repeatable** technical approaches.

Tile-Based Visualization for Big Data

WEB MAPPING APPS have an established model for intuitively navigating and understanding big data, based on zoomable multiresolution **tiles** and **layers**.

BUT geospatial data is relatively static. How can a tiled, layered approach be applied with **dynamic** data, and non-geospatial problems, at scale?



Aperture Tiles

TILE-BASED VISUAL ANALYTICS

- Hierarchical data tiling using cluster computing.
- Interactive on-demand image tile generation.
- Layers of raw data and derivative analytics.

OPEN SOURCE

- Oculus research product.
- Built on Apache Spark and Hadoop.
- Preparing for version one product release.



Example Applications



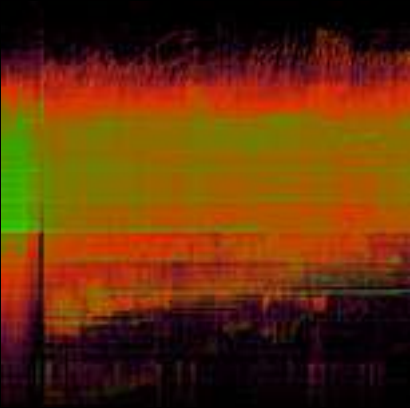
Geospatial



Social Media



Biomedical



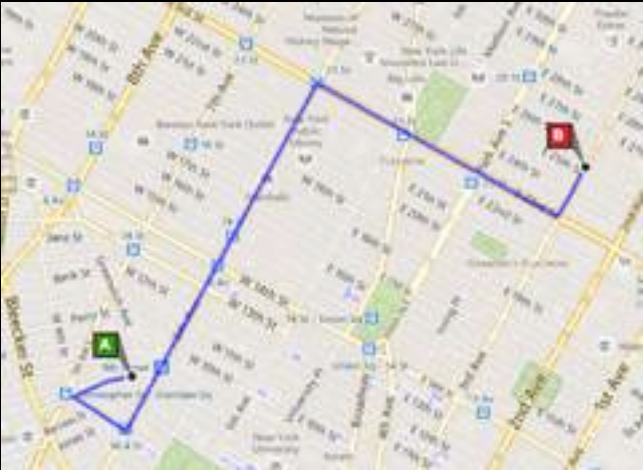
Financial

187,000,000 NYC Taxi Trips



NYC Taxi Data

FOIL request by Chris Whong, March 2014



Taxi ID	Hack, Medallion, Vendor
Origin	Location, Time
Dest	Location, Time
Trip	Duration, Distance, # Passengers
Fare \$	Base, Tip, Tolls, Taxes, Payment method



2013 NYC Taxi Trips



enf

Visualizing NYC Taxi Data

Eric Fischer, MapBox



www.mapbox.com/blog/nyc-taxi/



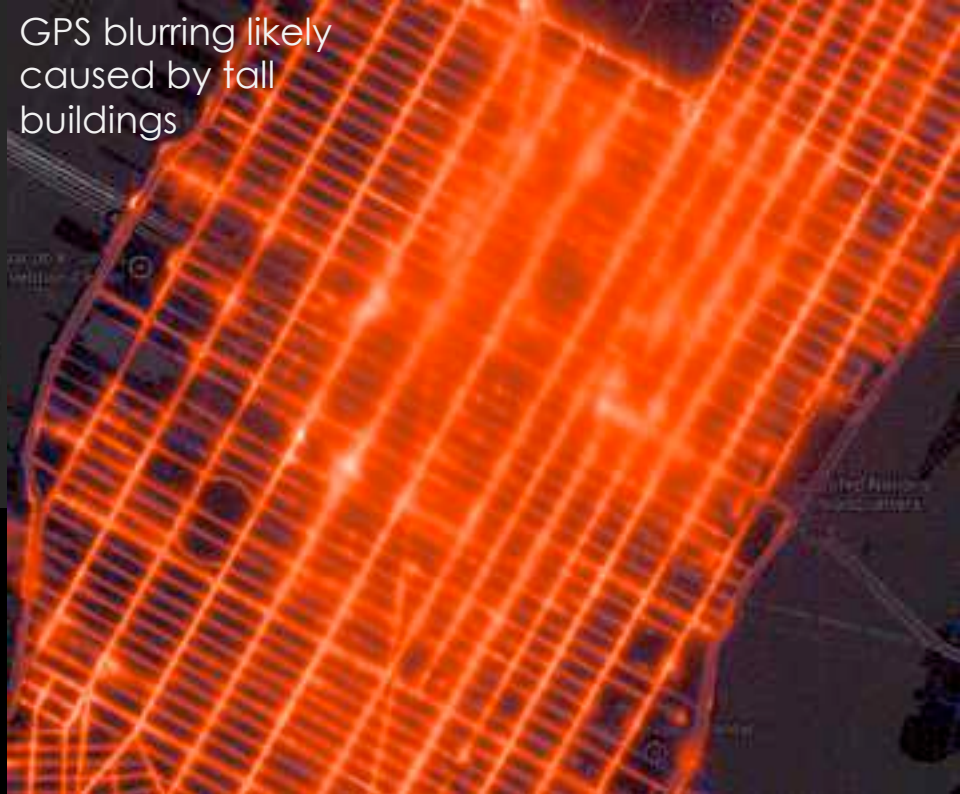
Demo

[Continue presentation](#)

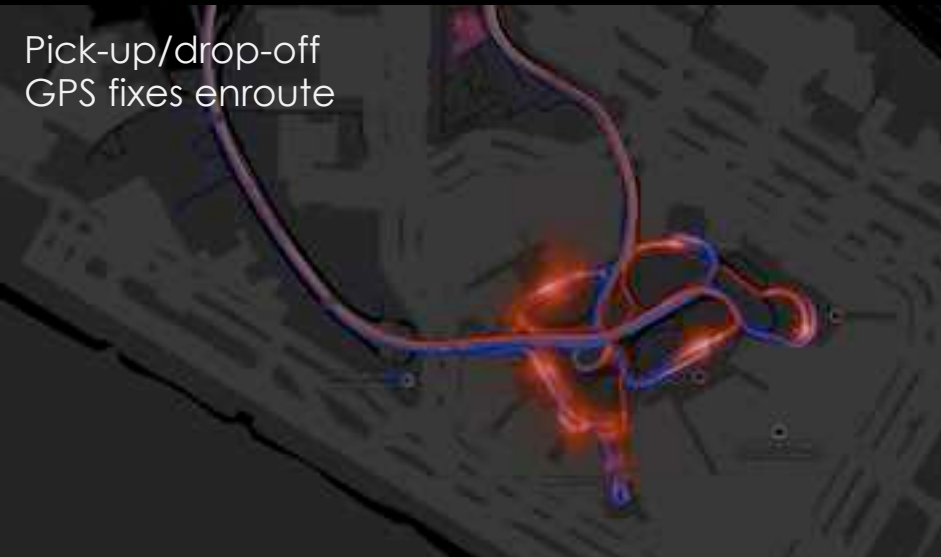
Location data errors



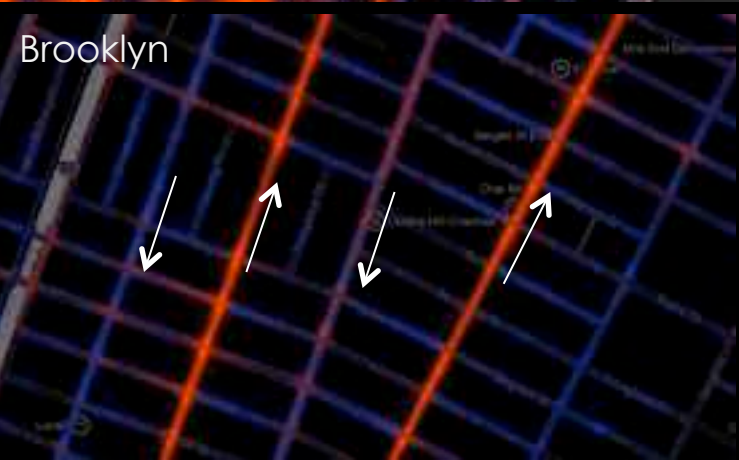
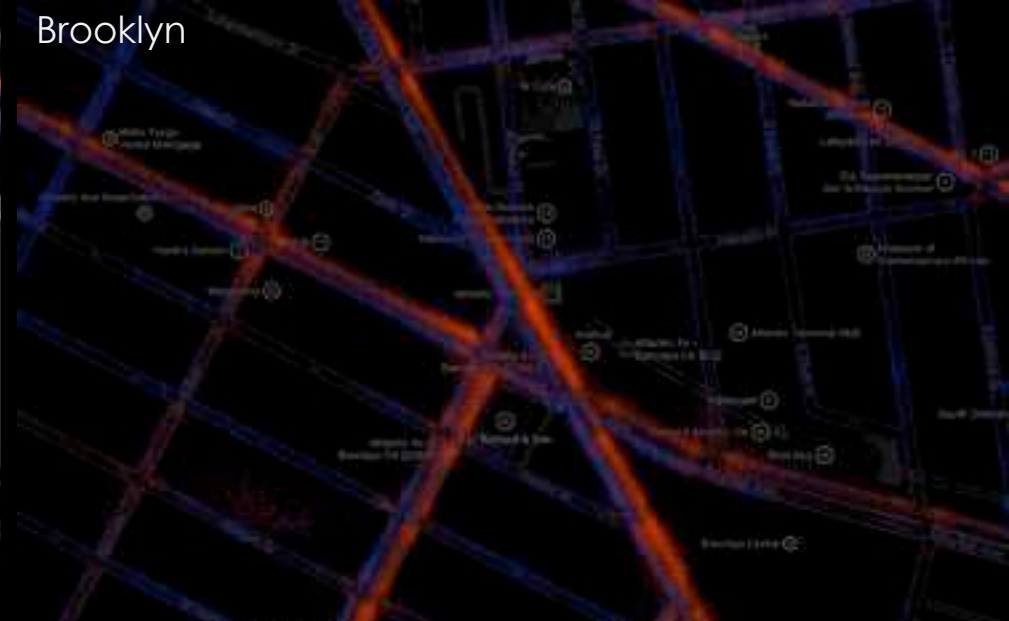
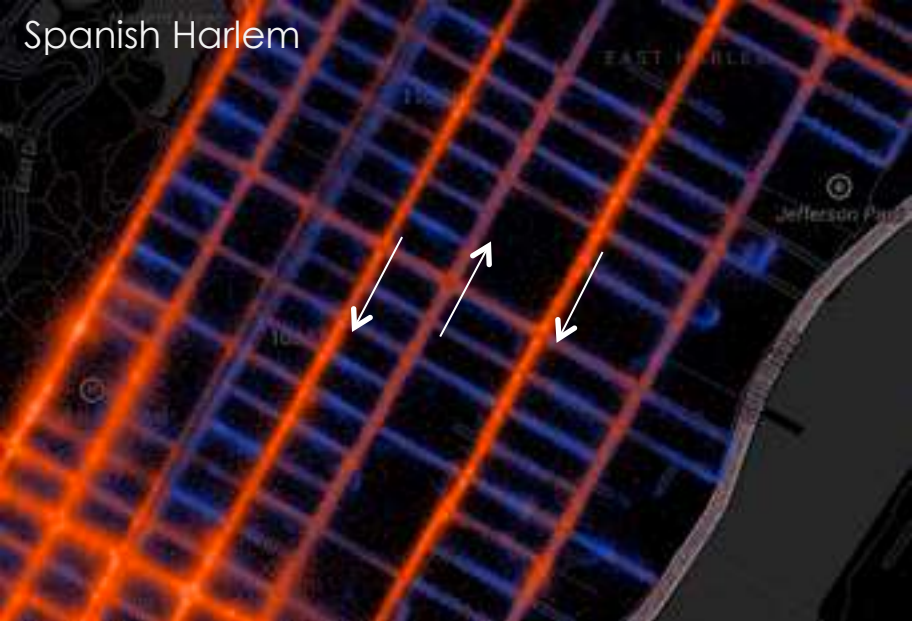
GPS blurring likely caused by tall buildings



Pick-up/drop-off GPS fixes enroute



Visualizing all of the data including “bad data” helps identify source of errors, understand how errors manifest, and how they may affect analysis



Comparing Pick-up and Drop-off Locations

- Red – Pick-up Locations
- Blue – Drop-off Locations

Multi-scale visualization allows exploration of macro and micro trends.



Average Tip % by Pick-up Location*

- Red ~15%
- Green ~20%
- Yellow ~25%

LaGuardia passengers tip more than those from JFK?

Downtown bankers tip less than mid-town shoppers?

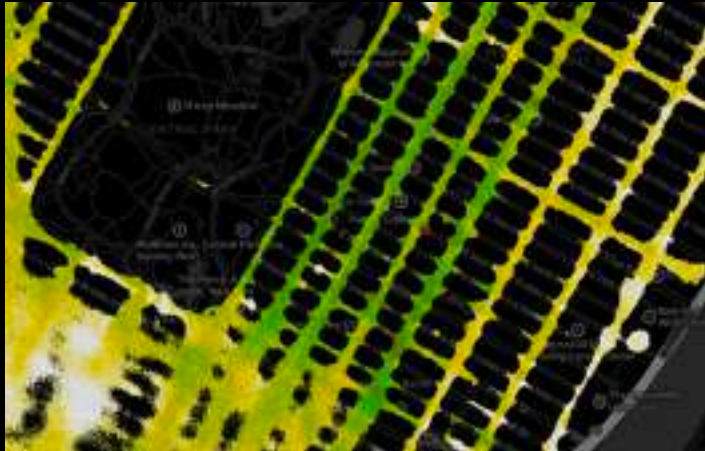
*CC payments only

Using a spectral color ramp helps us compare differences across the data (luminance better for multi-series).

Tiles uses aggregate data to generate rasters, not pre-computed images

Easily switching between fully zoomable multi-scale layers makes the task of understanding the data easier.

Very short ride anomaly in upper east side exposed by range filtering – hospital sending patients on a short ride to clinic?



Distance
Red = 1m
Yellow = 18m



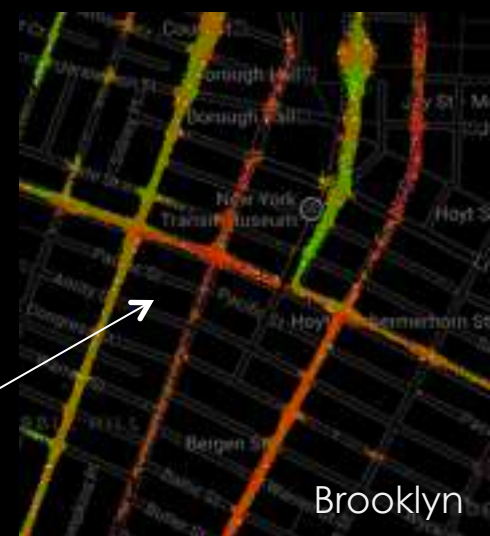
Average \$/hr by Pick-up Location

Red \$30/hr
Yellow \$75/hr

Calculated using:
 $(\text{fare} + \text{tip}) / (\text{duration of trip} + \text{time until next fare})$

Hot zones at ends of
tunnels due to GPS fix
lag observed earlier?

Profitability correlates
with pickup likelihood
on one-way streets



Brooklyn



Up until now have focused on pick-up location, what about drop-off?

Average \$/hr by Drop-off Location

Red	\$30/hr
Yellow	\$60/hr

Calculated using:
 $(\text{fare} + \text{tip}) / (\text{duration of trip} + \text{time until next fare})$

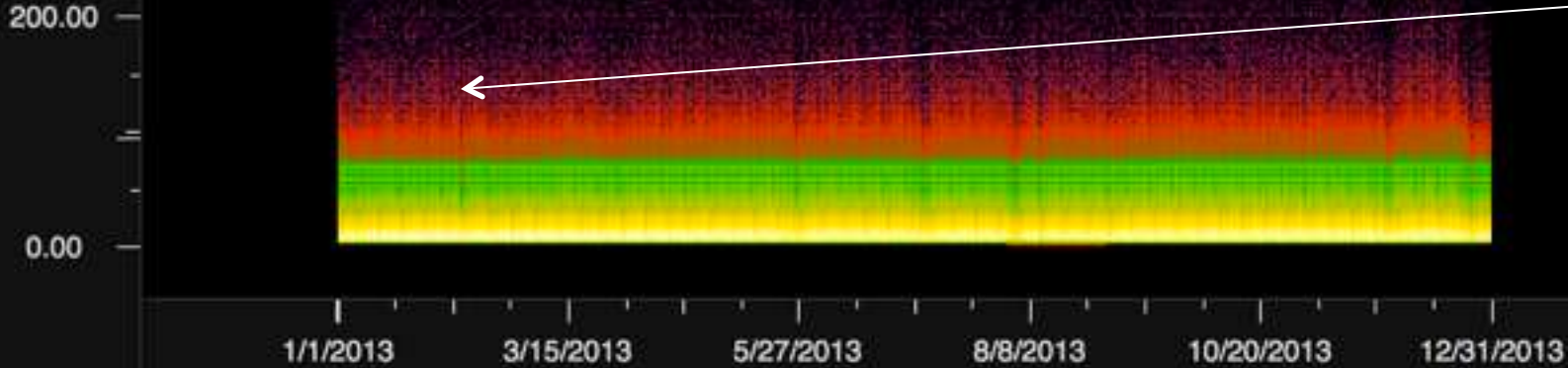


Average annual taxi
income by pick-up
location

Red \$80k
Yellow \$100k

Fare Total vs Time

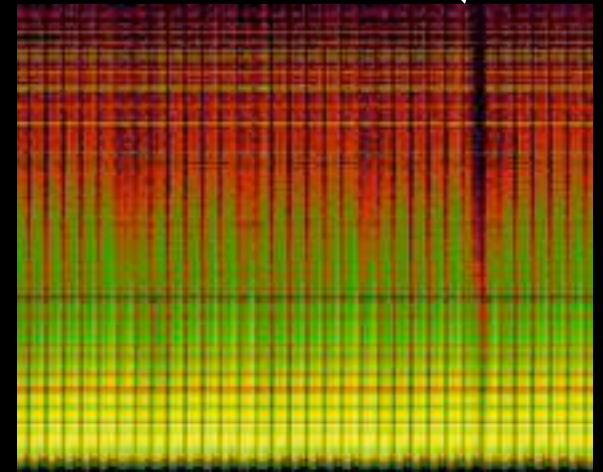
Feb 2013 "Nor'easter storm"
Expensive fares drop off by
low fares unaffected



Tiles approach isn't limited to geographic plots

X-Y crossplots can show 3 arbitrary dimensions and
expose macro and micro patterns.

e.g. Time + Fare Total + Frequency

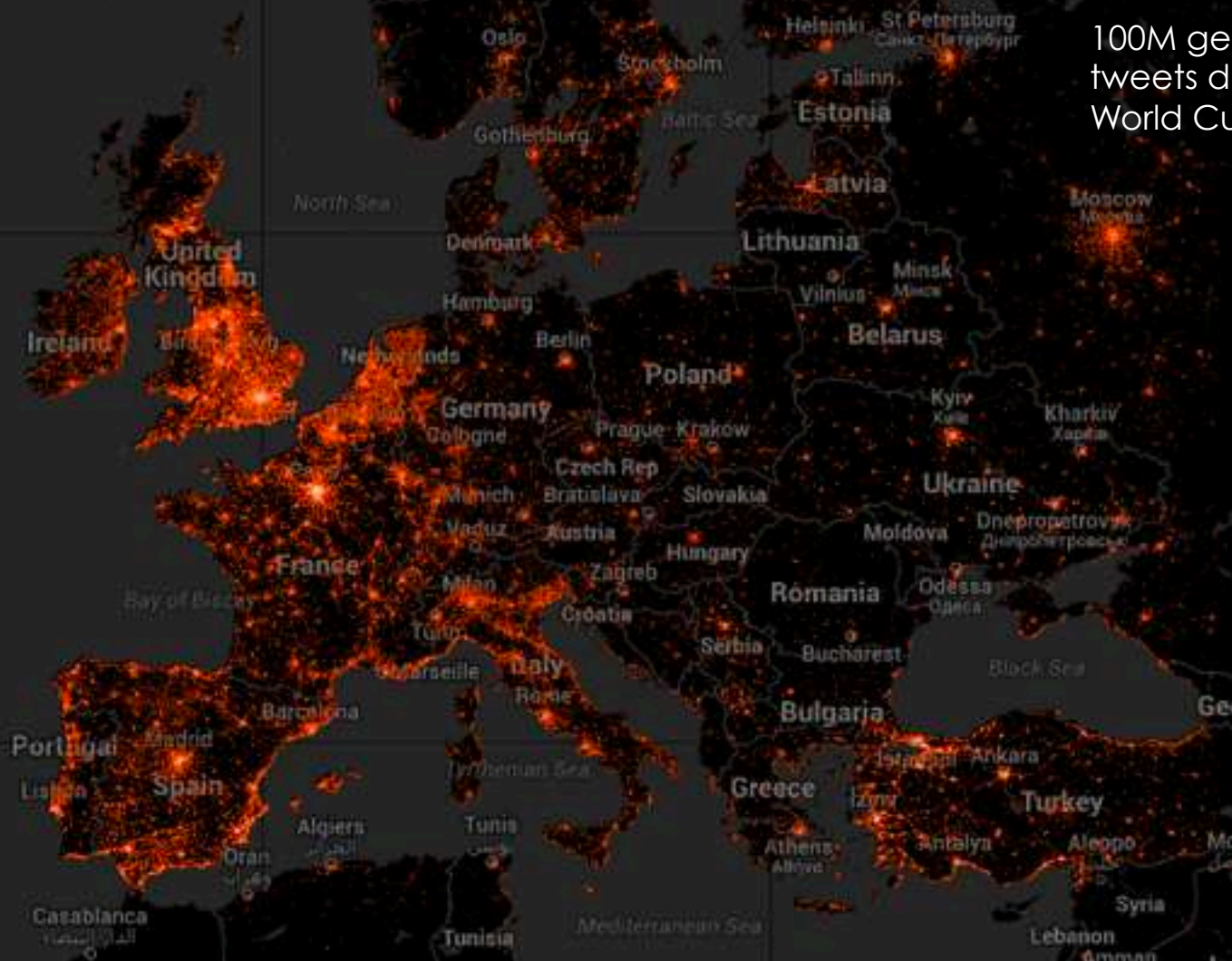


World Cup Tweets

player @mentions



Photo by [paulisson miura](#), CC



100M geolocated tweets during the World Cup



Analytic overlays showing @mention summaries

Antofagasta

Spain
Portugal
Chile
Argentina
Brazil
Columbia

Chile

Columbia
Brazil
Spain

Chile
Argentina
Portugal

Mendoza

Santiago

Argentina
France
Portugal

Argentina

Brazil
Spain
Columbia

Rosario

Uruguay

Buenos Aires

Buenos Aires

1128

Paraguay

Asuncion

Brazil
Portugal

Argentina

France
Spain
Columbia

Cordoba

GRUSSU
DO SUL

STATE OF
SAO PAULO

STATE OF
PARANA

Columbia
Spain

Brazil
Portugal

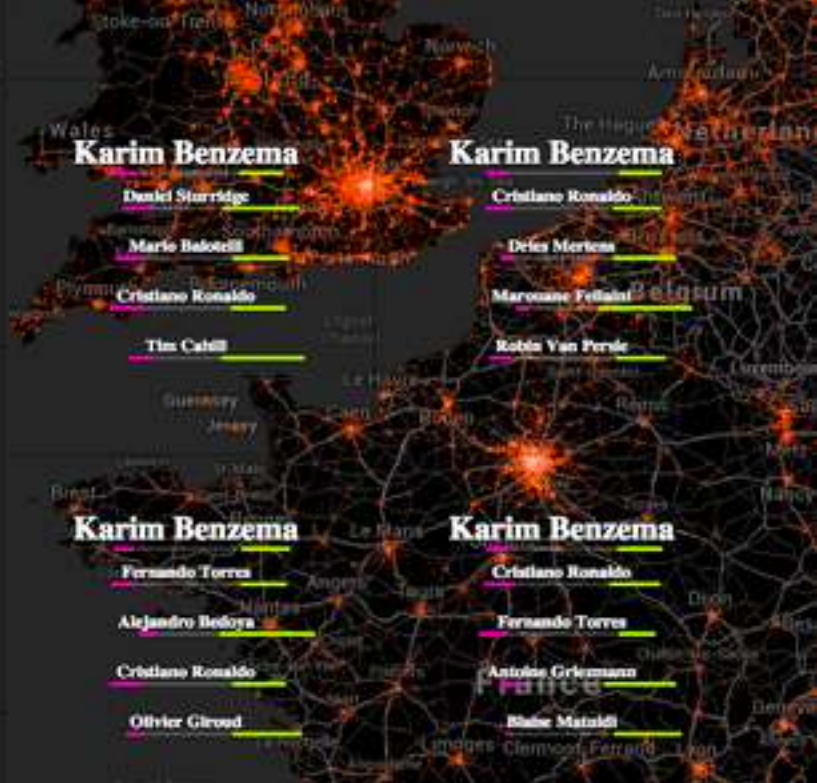
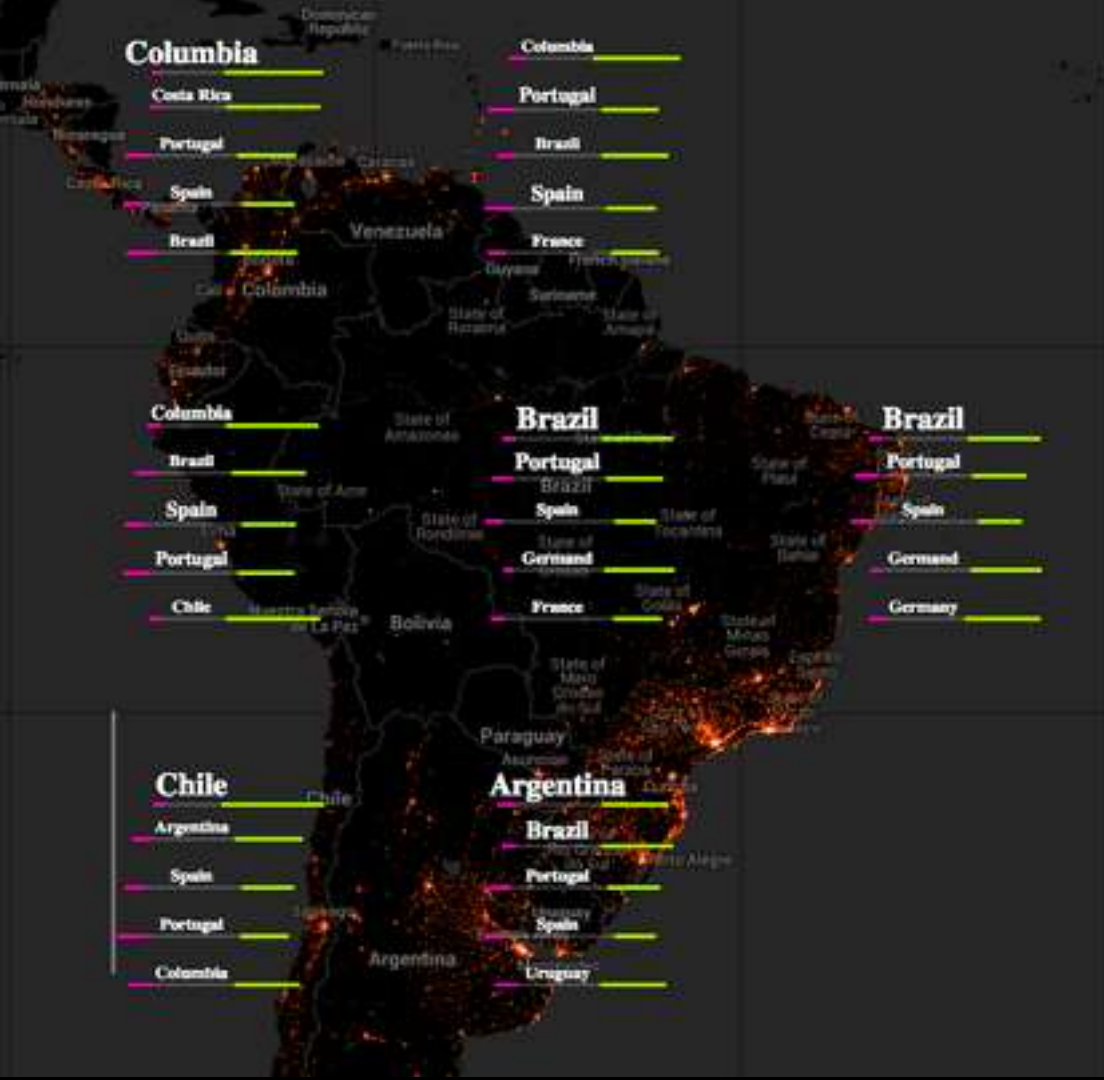
STATE OF
RIO
DO SUL

Germany
France

Porto Alegre

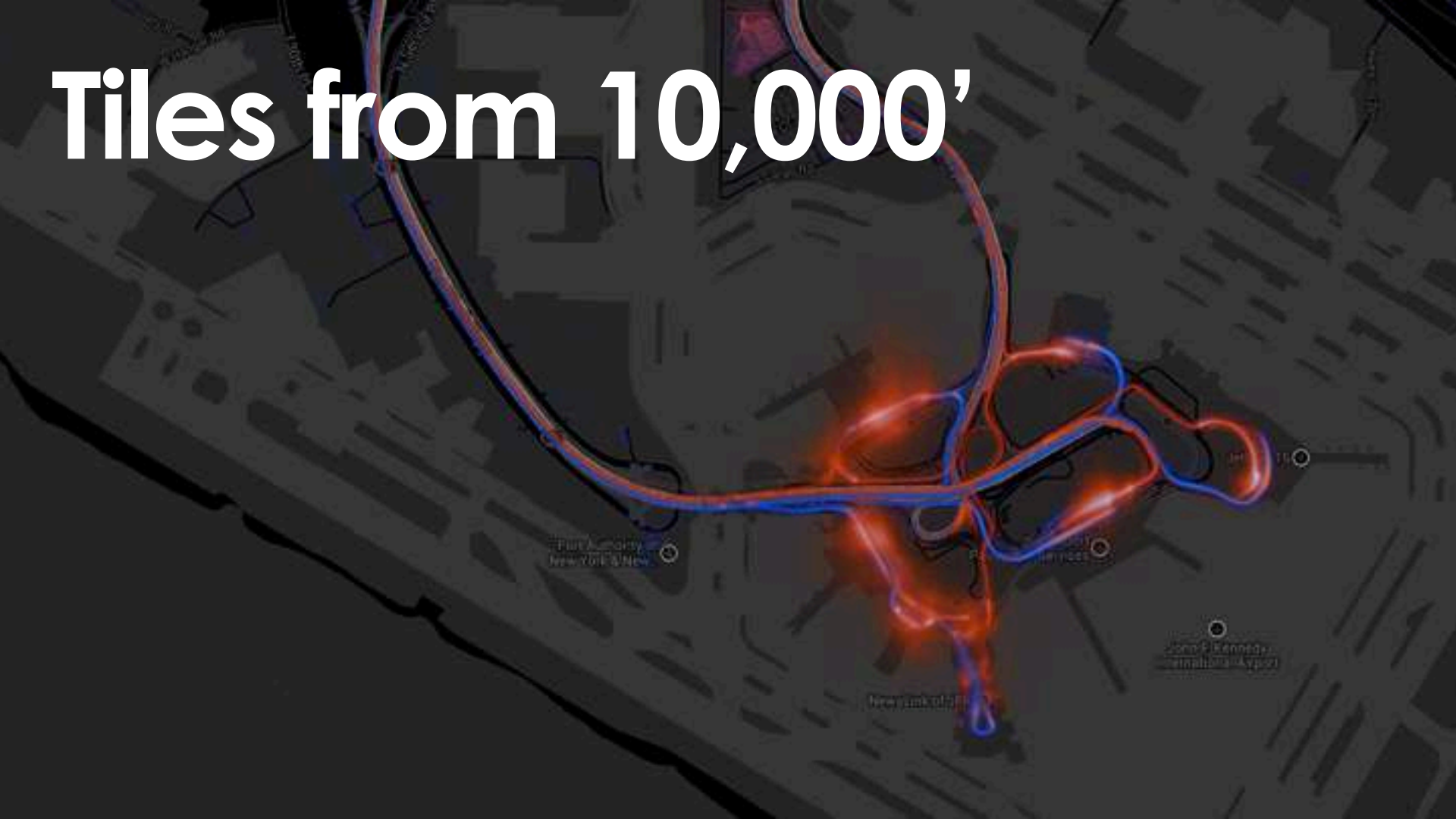
Columbia
Italy

Uruguay
Spain
Portugal
Argentina



Team and Player @mention sentiment
Interactive overlays

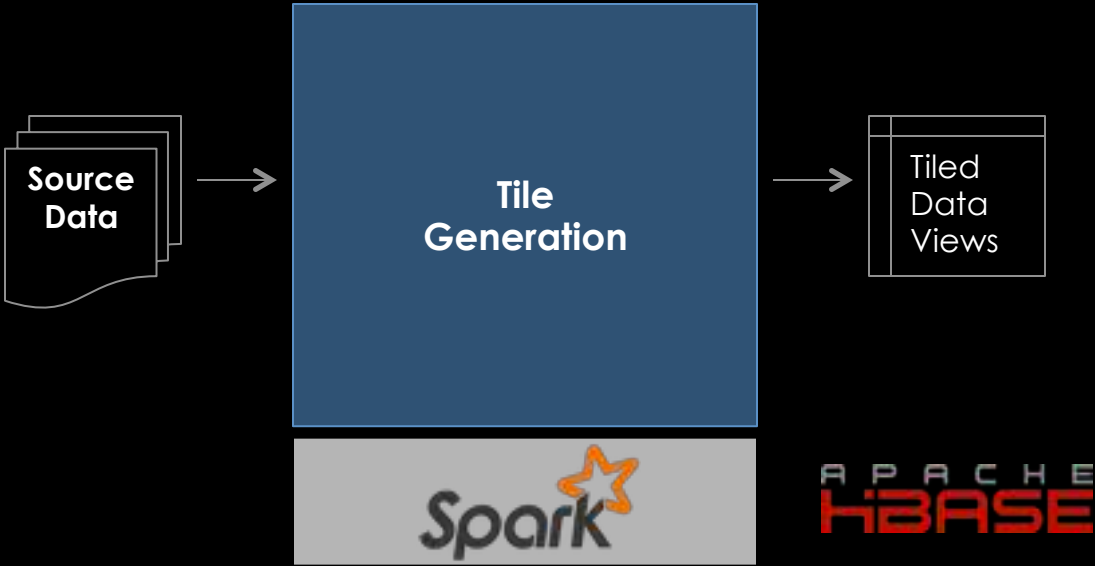
Tiles from 10,000'



Three Tiers



Tile Generation



Tiled Data Views

Aggregate views of source data designed
for answering **analytic questions**

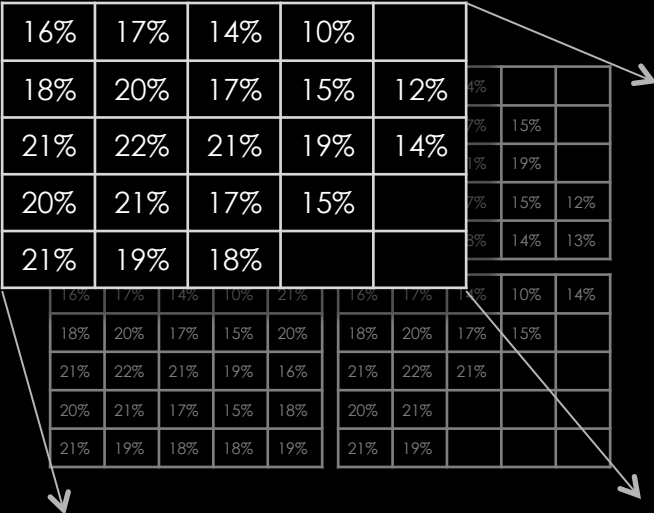
Optimized for query speed

Indexed by tile key - one view/row per tile

AVRO / Thrift

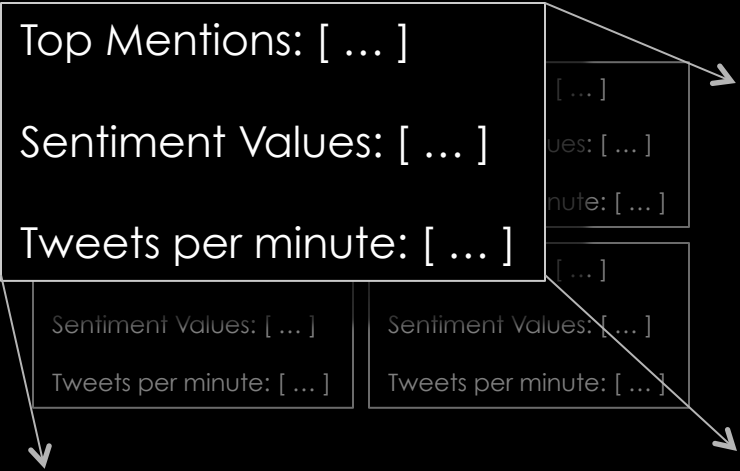
Tiled Data Views

Taxi Tip %



Tiled Data Views

World Cup Sentiment



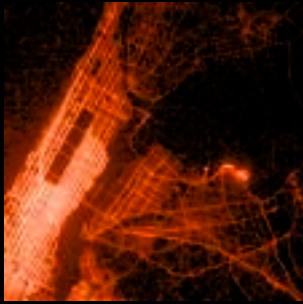
Tile Generation

16%	17%	14%	10%	
18%	20%	17%	15%	12%
21%	22%	21%	19%	14%
20%	21%	17%	15%	
21%	19%	18%		

APACHE
HBASE



← query



Tile Client

Tile Map Service (TMS)

`http://{rootURL}/{layer}/{zoom}/{x}/{y}.png`



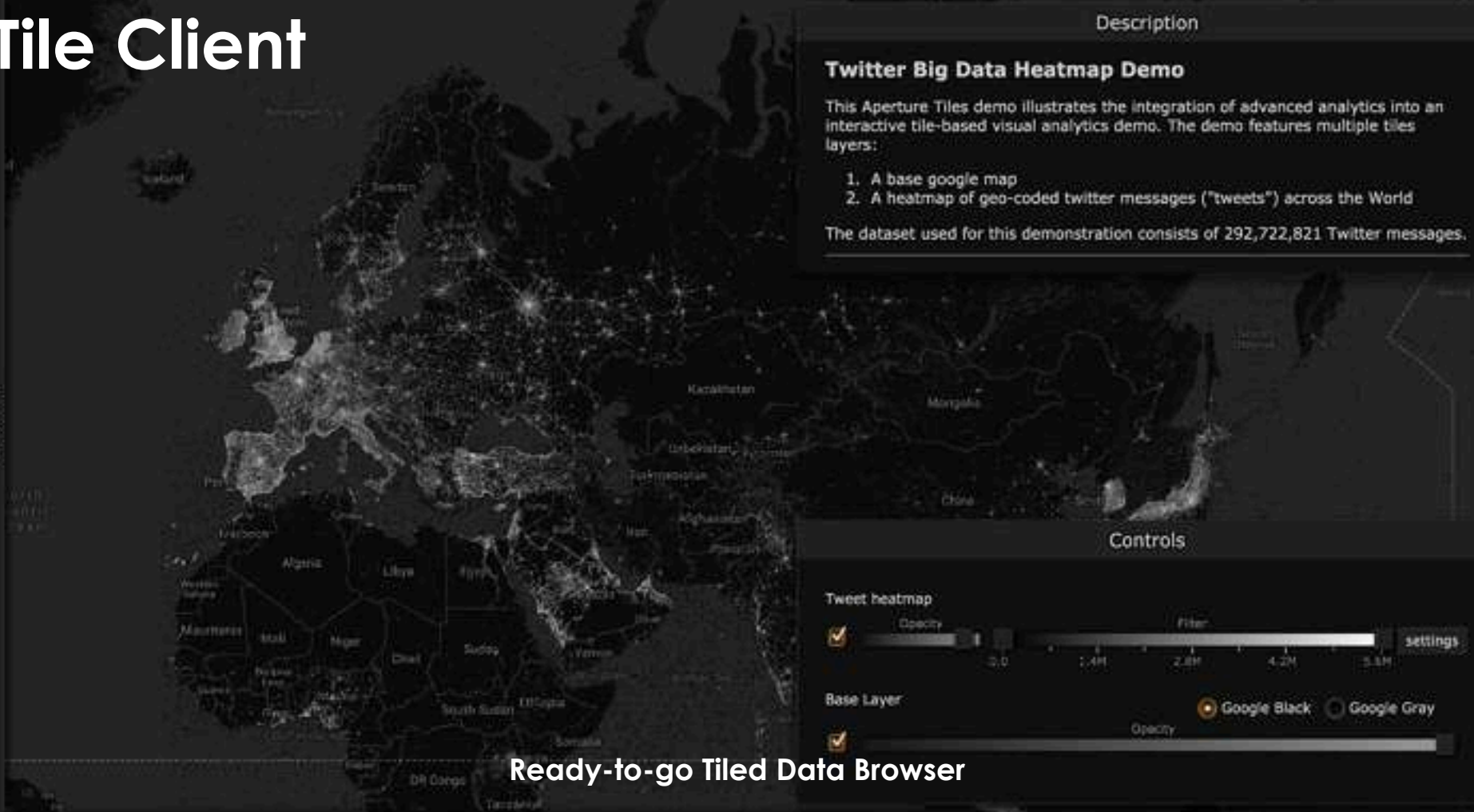
OpenLayers



ESRI

Tile Client

Latitude



Description

Twitter Big Data Heatmap Demo

This Aperture Tiles demo illustrates the integration of advanced analytics into an interactive tile-based visual analytics demo. The demo features multiple tiles layers:

1. A base google map
2. A heatmap of geo-coded twitter messages ("tweets") across the World

The dataset used for this demonstration consists of 292,722,821 Twitter messages.

Controls

Tweet heatmap

Opacity Filter

0.0 1.4M 2.8M 4.2M 5.6M

Base Layer

Opacity Google Black Google Gray

Ready-to-go Tiled Data Browser

Longitude



<http://tiles.oculusinfo.com>



<https://github.com/oculusinfo/aperture-tiles>

Key Points of Value

PLOTTING ALL of the data for richer, truer insights.

- Reveals **truths** in the data that summaries cannot.
- Provides important **context** and **evidence** for analytics.

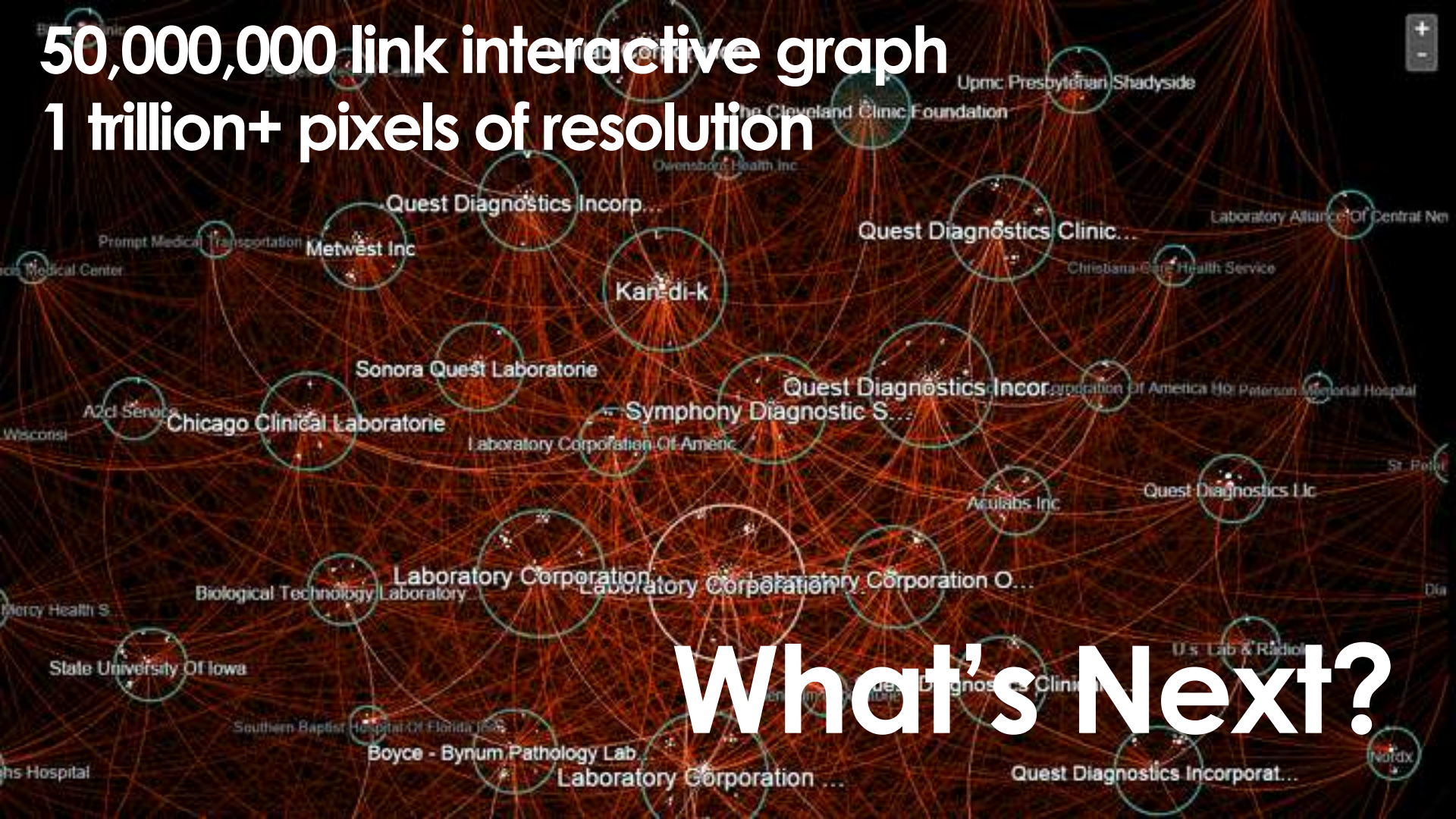
TILED, LAYERED user interface for natural, scalable exploration.

- **Intuitive** navigation and understanding.
- Analytic layers at **levels of detail**, down to raw provenance.

INTERACTIVE analysis with data tiling and analytics.

- Rapidly repeatable **tailored** solutions, with filtering + selection.

50,000,000 link interactive graph
1 trillion+ pixels of resolution



What's Next?

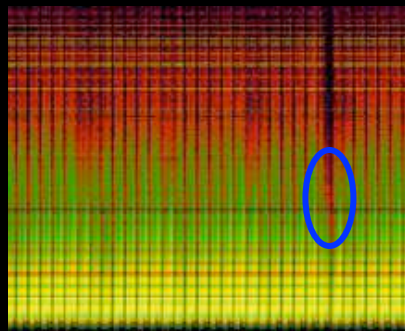
What's Next?



Interactive
Tiled Graphs



Streaming Data
and live query



More advanced
filtering and analytics

Questions?

Oculus gratefully acknowledges the support of DARPA in research and development of this work.