

Value in the DETAILS

Understanding detailed data through VISUAL EXPLORATION

Richard Brath
Rob Harper



uncharted





I'll be visually exploring all kinds of interesting patterns in tweets about Trump.

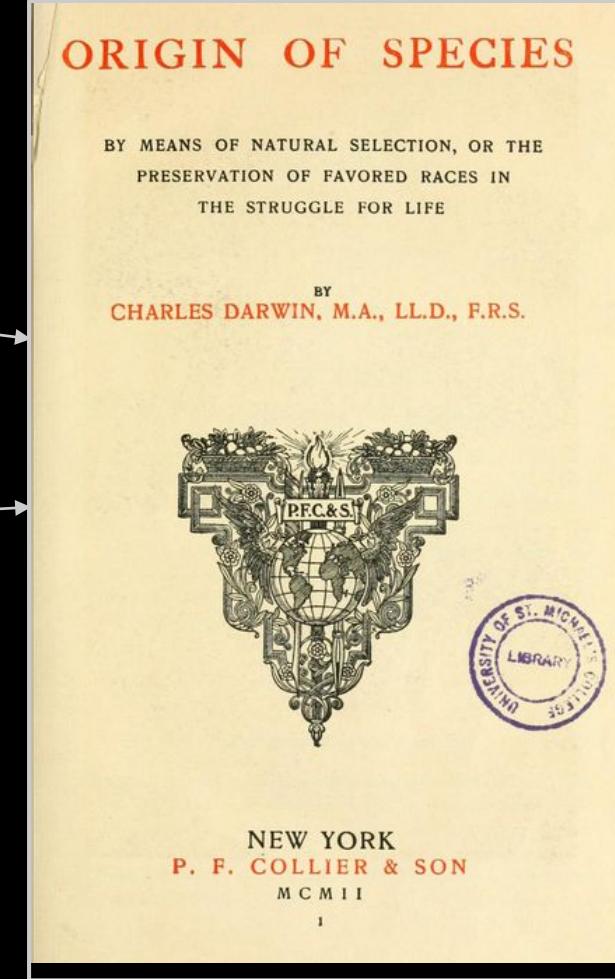
But wait – first what do we mean by visual exploration?



Visual Exploration



Natural
Selection?



Visual Exploration Process

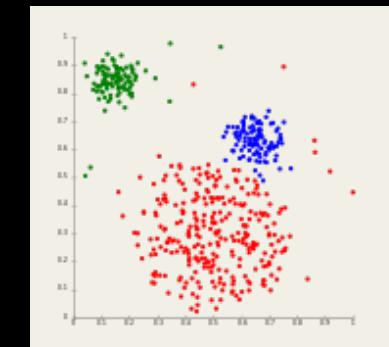
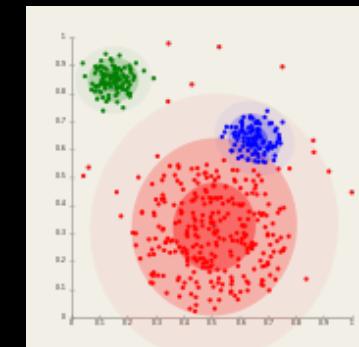
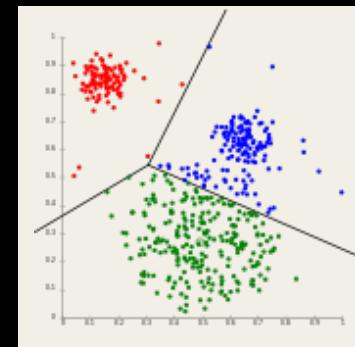
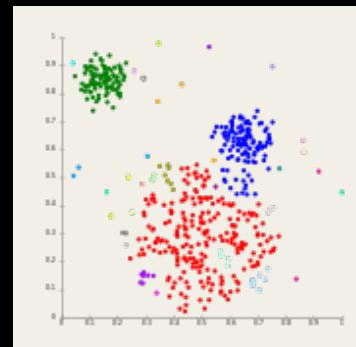
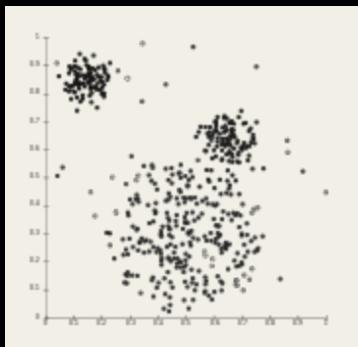
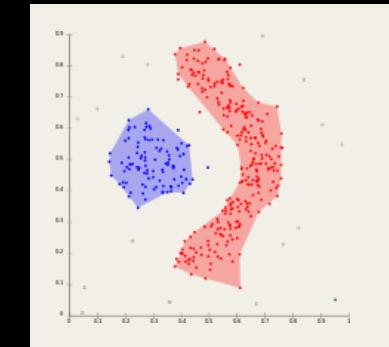
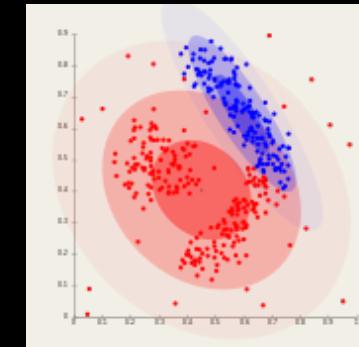
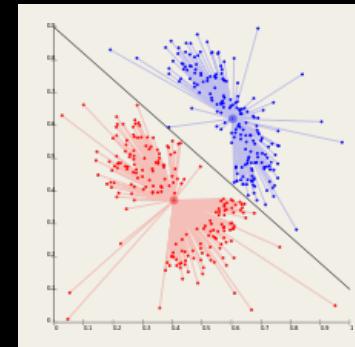
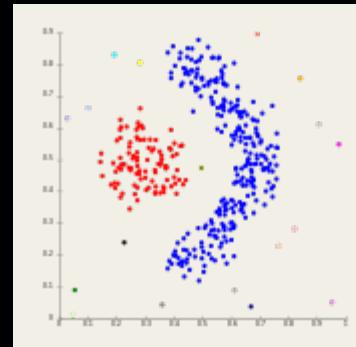
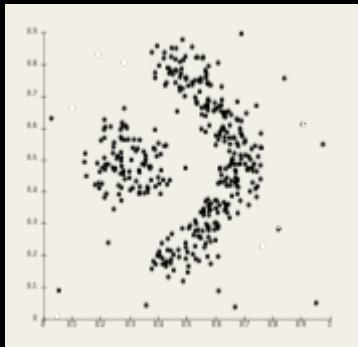
1. Collect a lot of data
2. Observe some interesting patterns
3. Hypothesize about why they exist
4. Refine and build models



Understanding detail data

Can you see the patterns?

How about these algorithms abilities to find the patterns?



No Clustering

Linkage Clustering

K-Means Clustering

Distribution-based Clustering

Density-based Clustering

We see patterns all the time, and quite easily

We tend to group things based on visual cues, such as proximity, alignment and containment.



Seeing detailed patterns

Perception can be whole. Once you see it, you won't un-see it.



First publication of the picture probably in Life Magazine:58;7 1965-02-19, p 120.

Also, check out the movie by Wim van de Grind (http://www.michaelbach.de/ot/cog_dalmatian/index.html)

So what?

Powerful human perception system

- Detects **patterns** in complex data
- Can find **patterns** based on different criteria
- Can find **patterns** at different scales

So what?

Exploratory Data Analysis stems from John Tukey's work in the early 1960s. EDA can be characterized by

- a. understanding "what is going on here?"
- b. graphic representations of data
- c. tentative model building and hypothesis generation
- d. robust measures, re-expression, and subset analysis
- e. skepticism, flexibility.

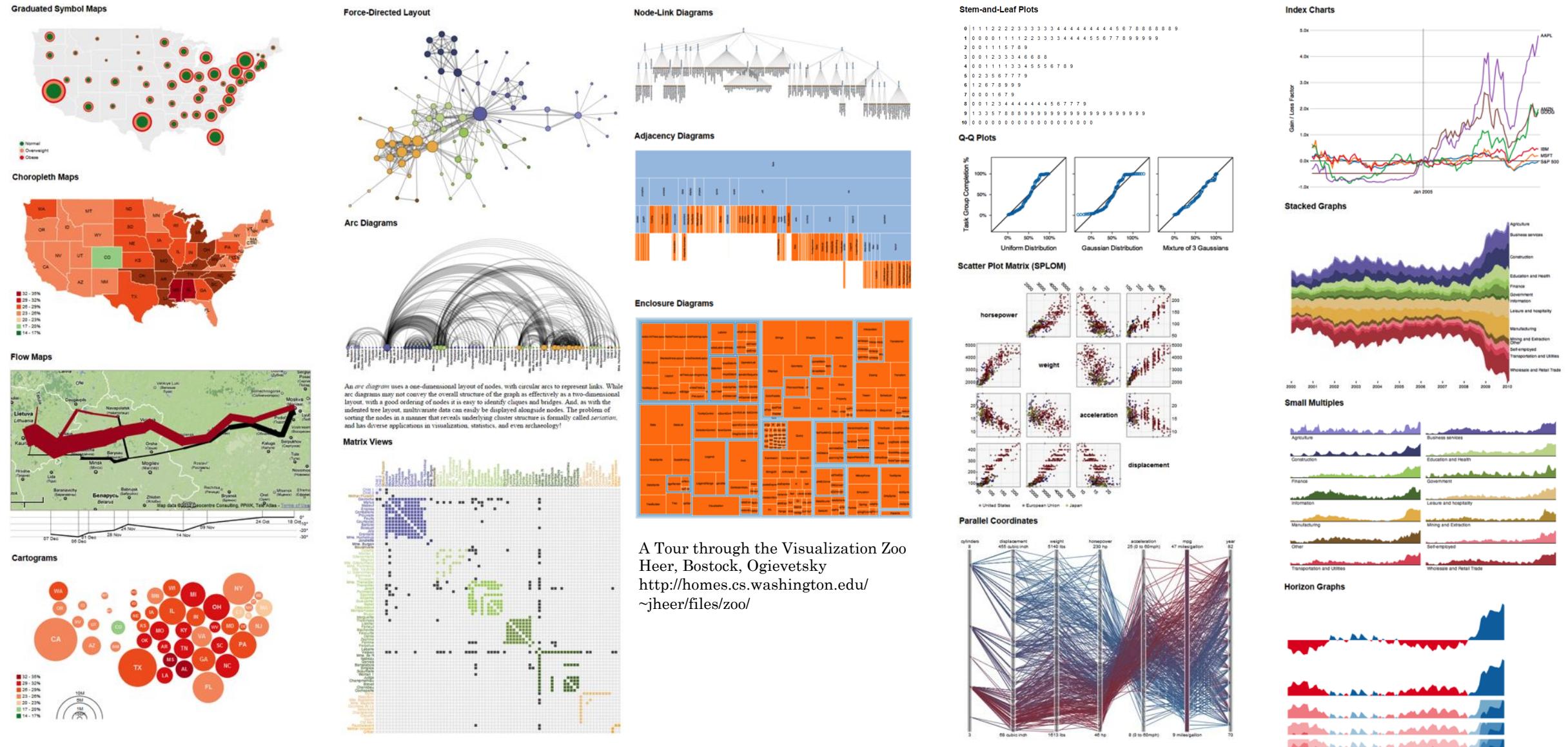
*The goal of Exploratory Data Analysis is to discover **patterns** in data.*



paraphrased from John Behrens, Arizona State University, Principles and Procedures of Exploratory Data Analysis, American Psychological Association, 1997.

So, why are we summarizing
big data into bar charts?

So, why we rolling-up big data into visualizations of 1000 points?



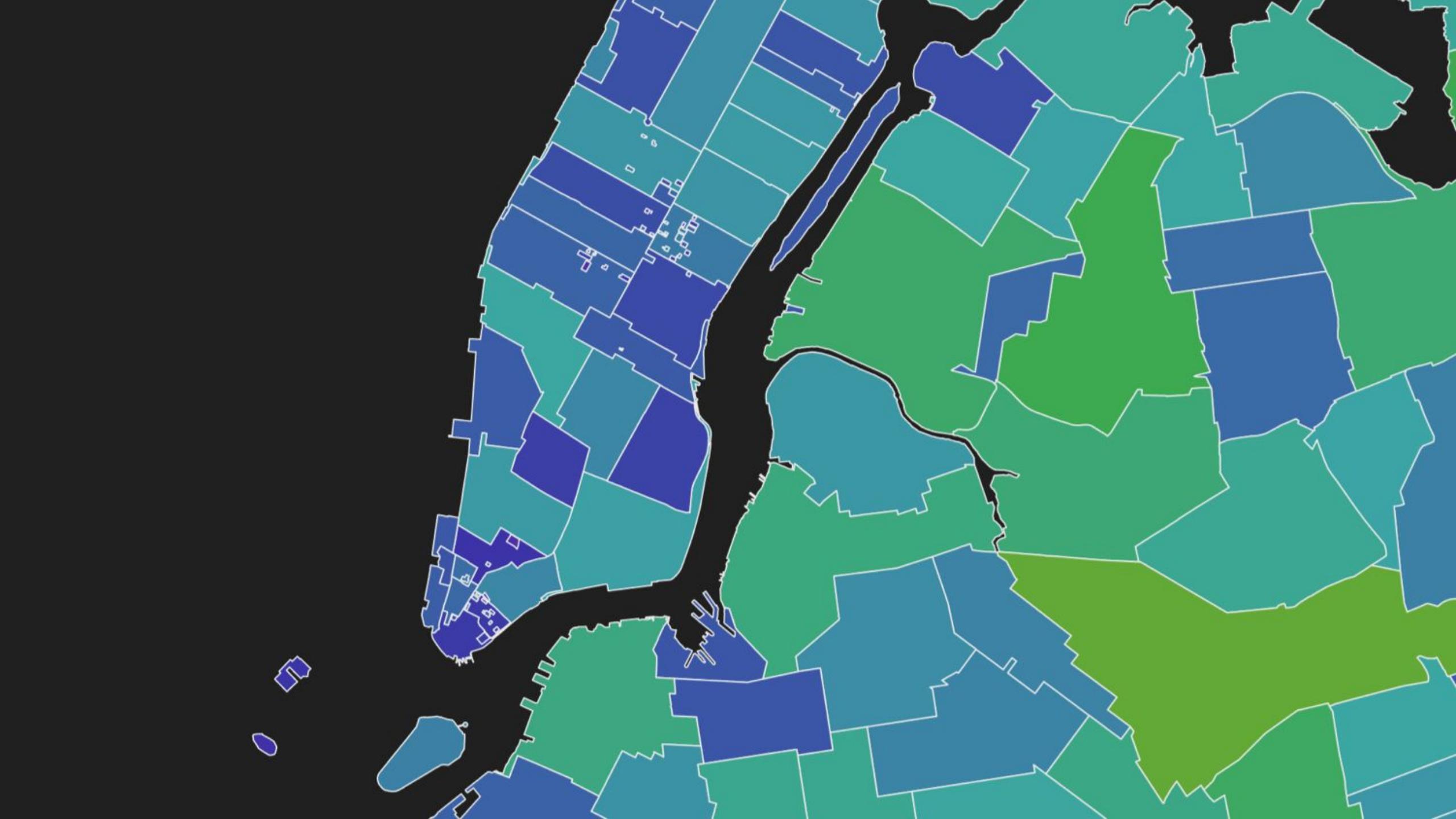
A Tour through the Visualization Zoo
Heer, Bostock, Ogievetsky
<http://homes.cs.washington.edu/~jheer/files/zoo/>

So, how do we visualize 100m data points?

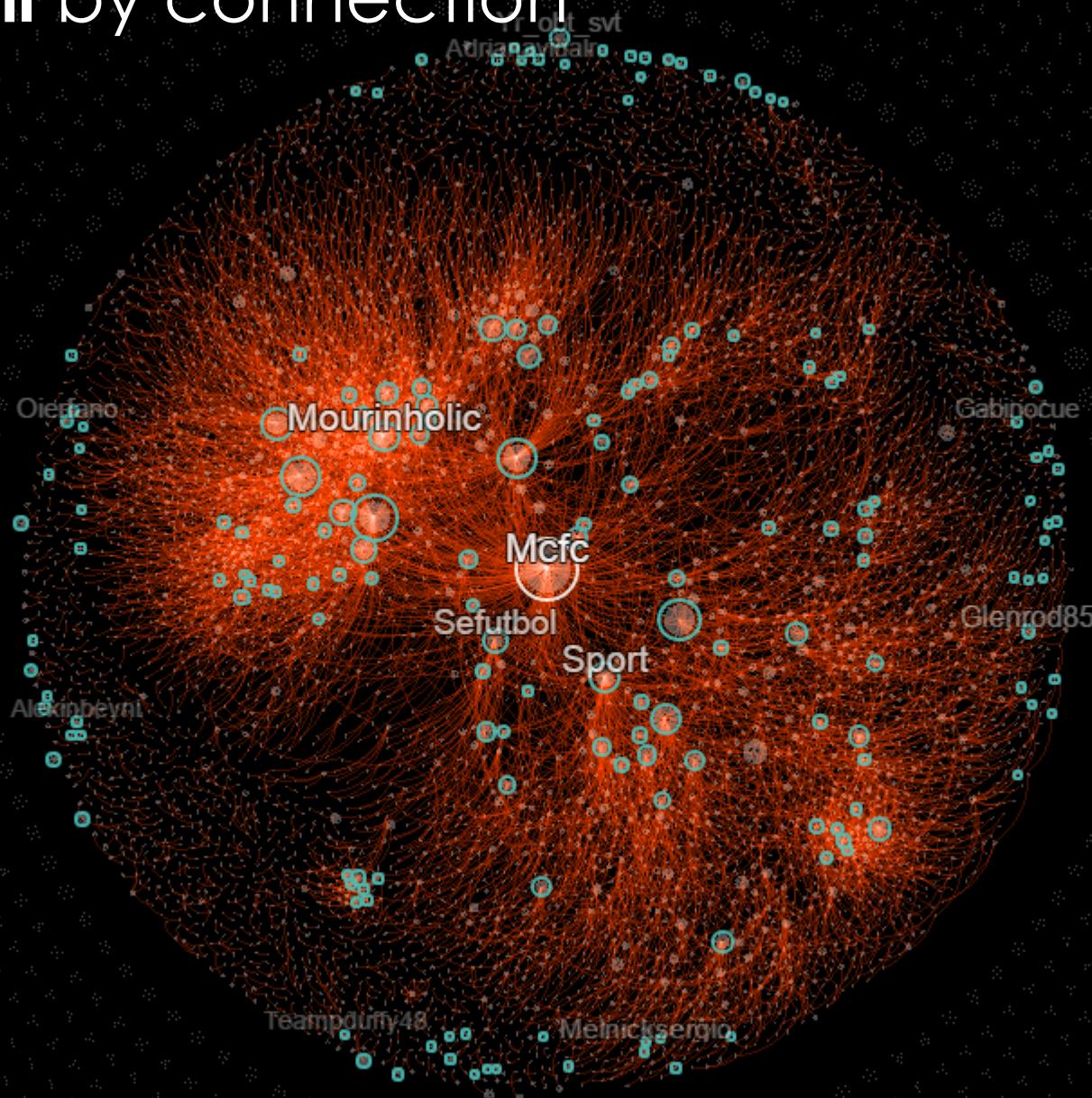


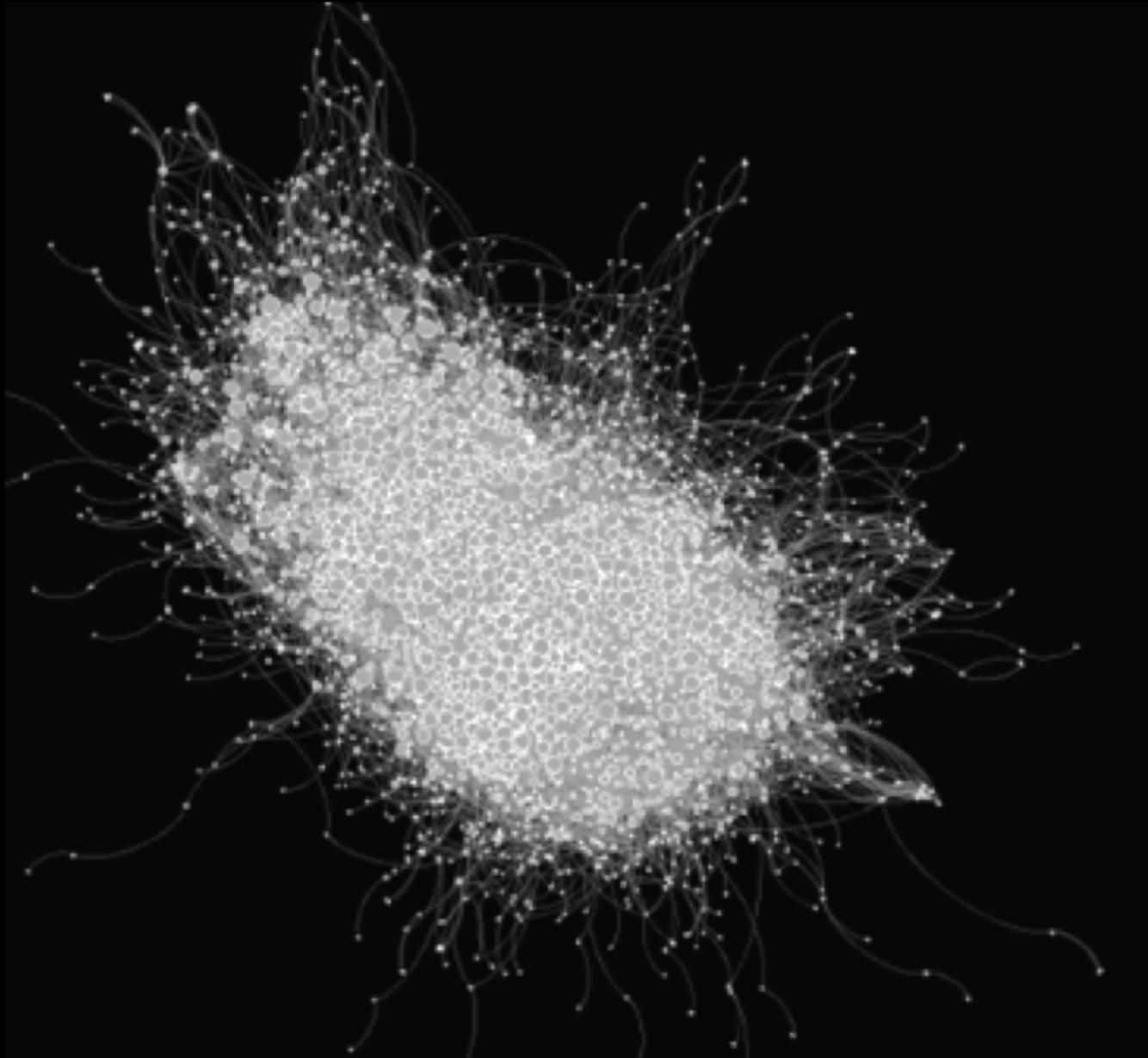
Plot all the detail by map



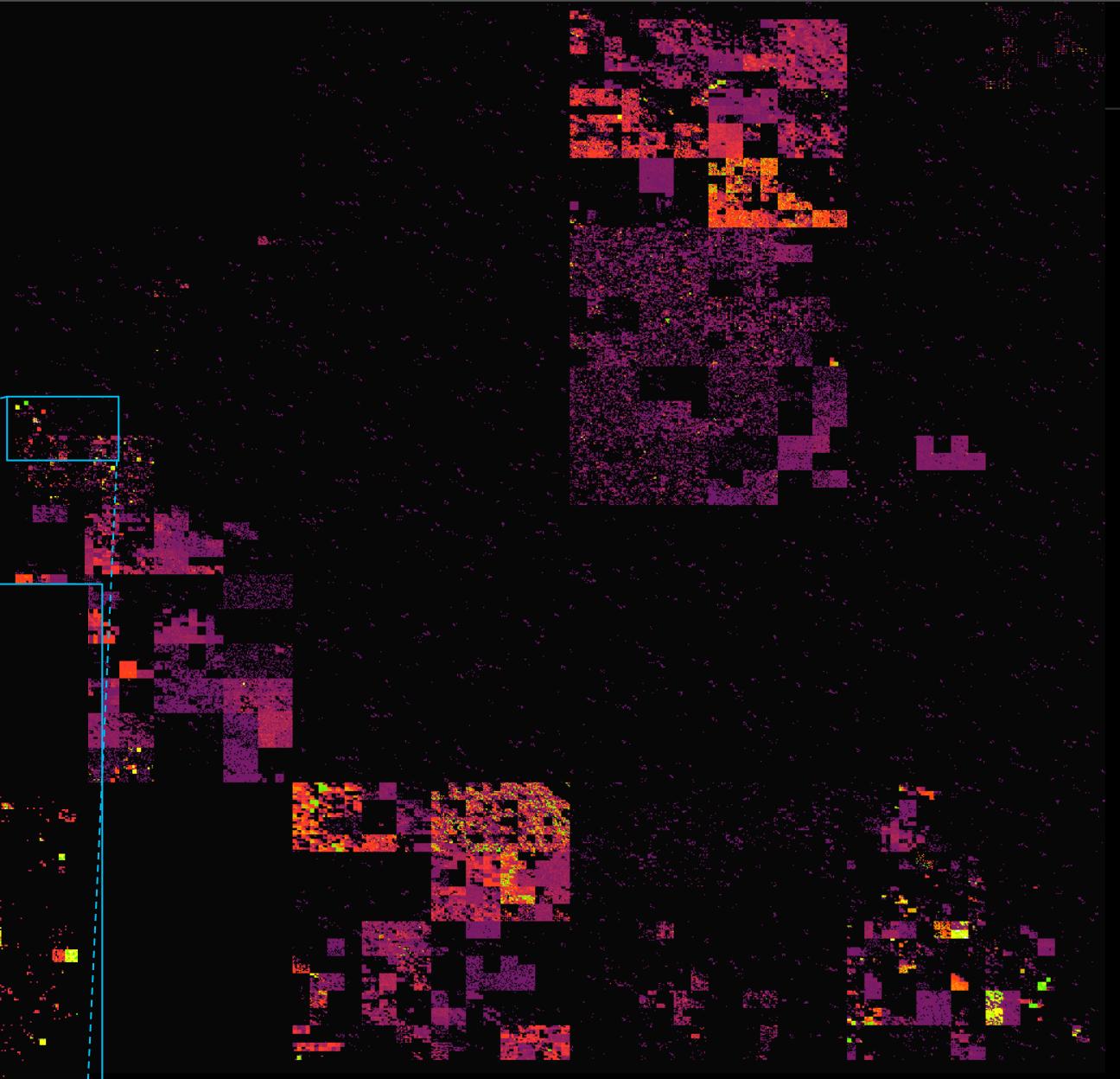
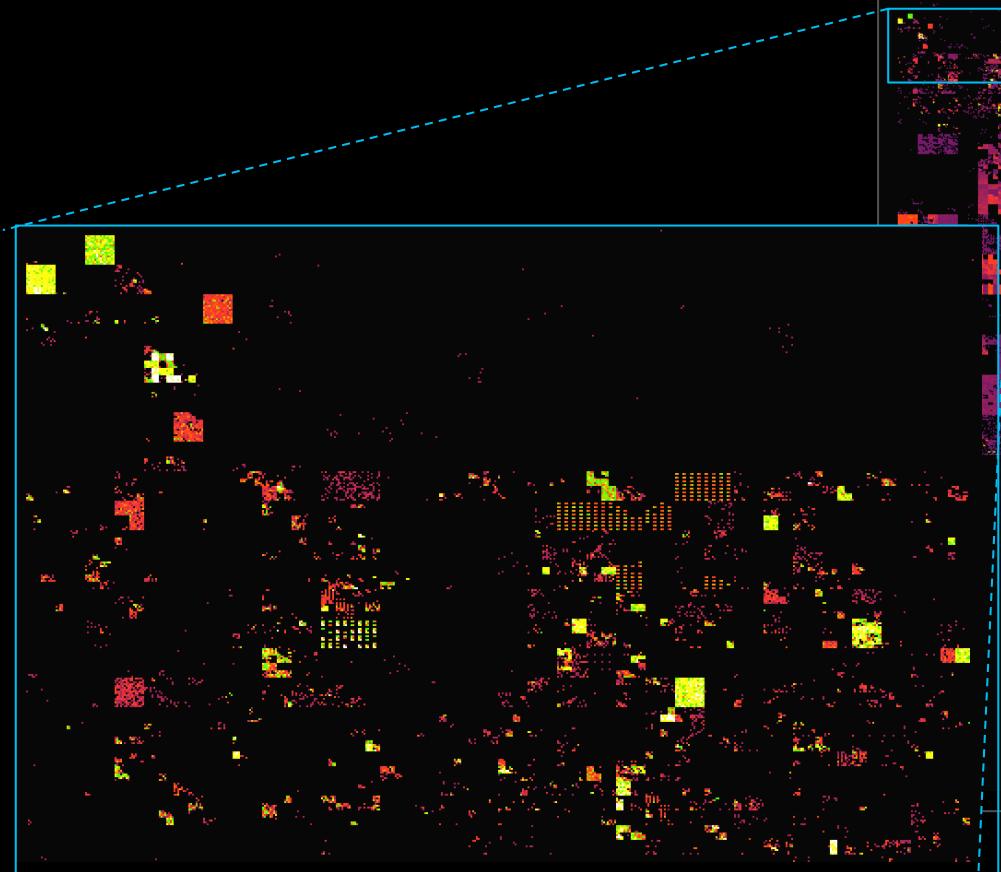


Plot all the detail by connection

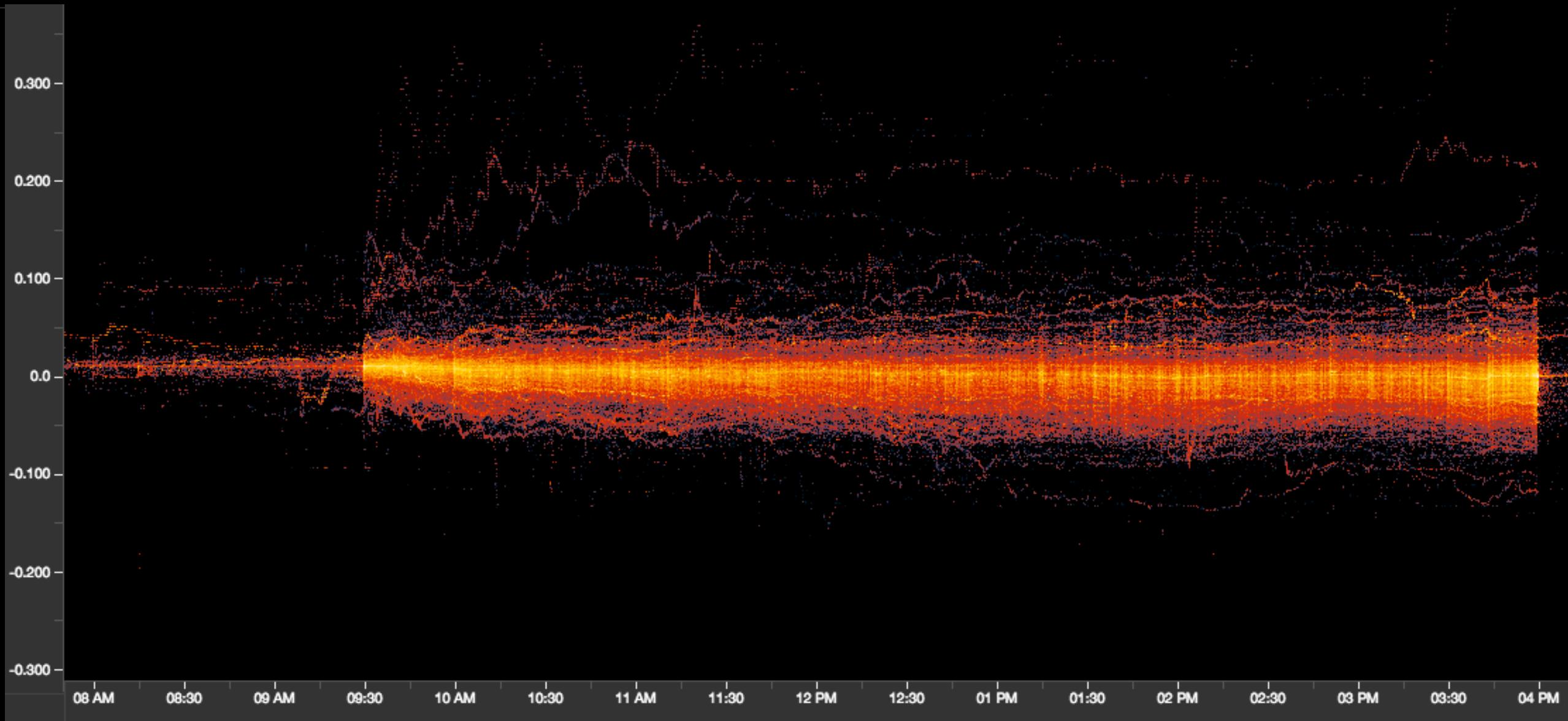




Plot all the detail by order



Plot all the detail by time



Visual Patterns in Time Plots

some other
variable



event

boundary

periodic

anomaly



level (threshold)

trend

time

Bitcoin Transactions



Load & Verify



Bitcoin Address

14sScGvSjGtxNbqFUoStoXN7eXydaN1JmM

bitcoin
Amount:

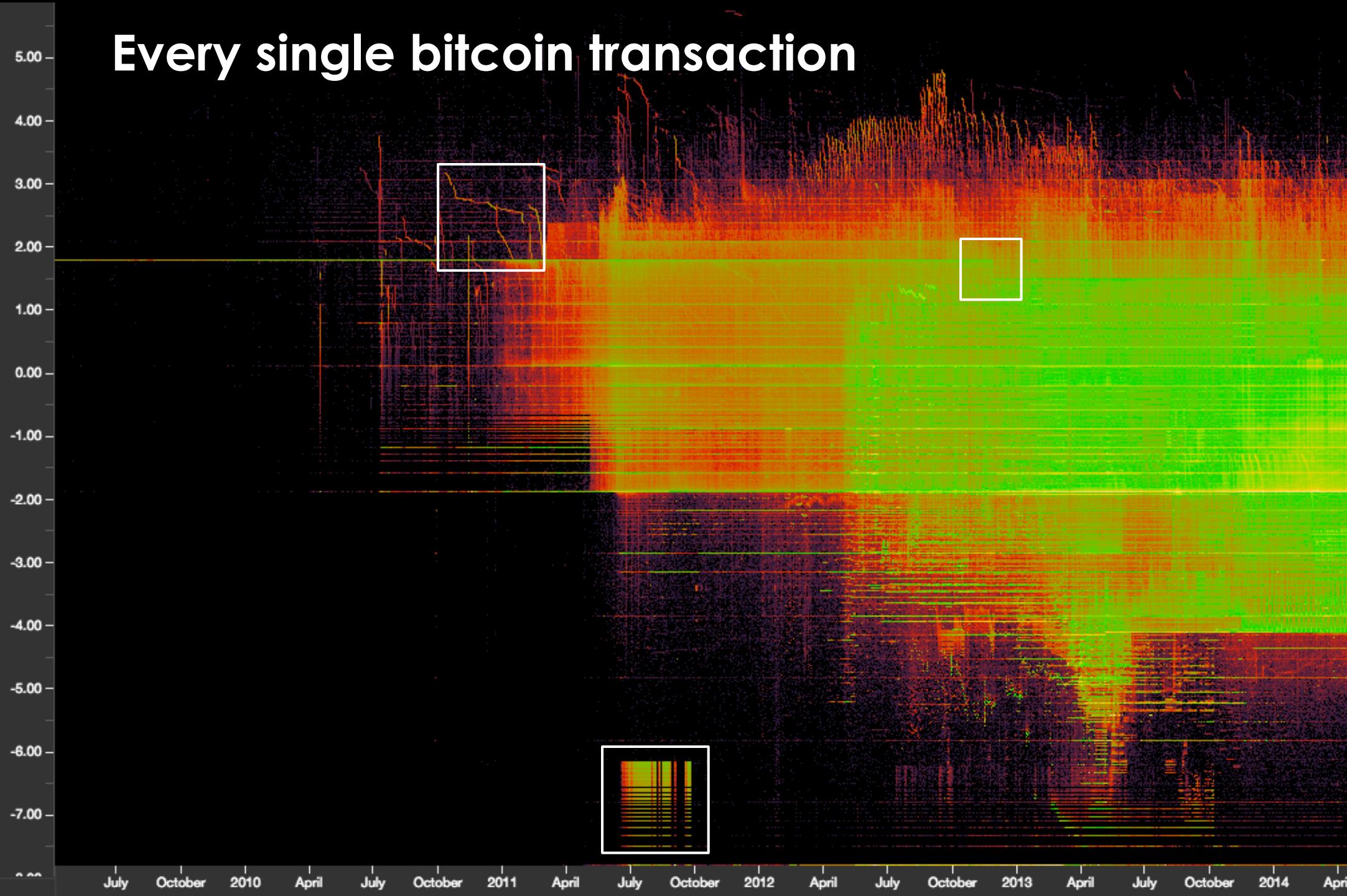
Private Key

HISTORICAL PERSPECTIVE ON THE USE OF SCAFFOLDING



Spend

Every single bitcoin transaction

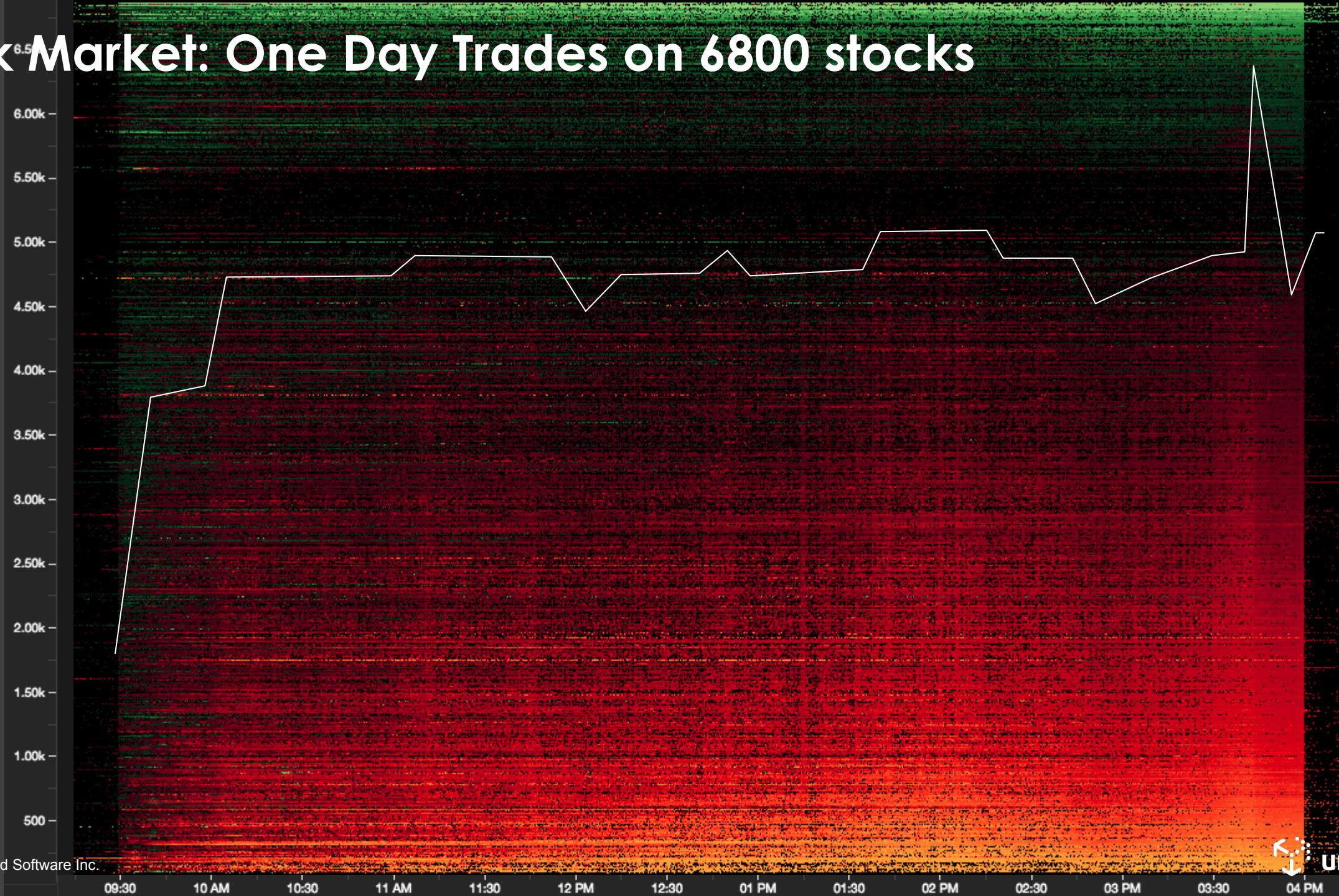
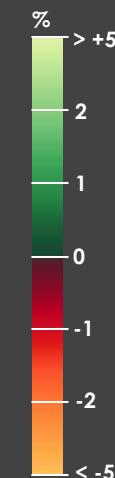


Financial Markets

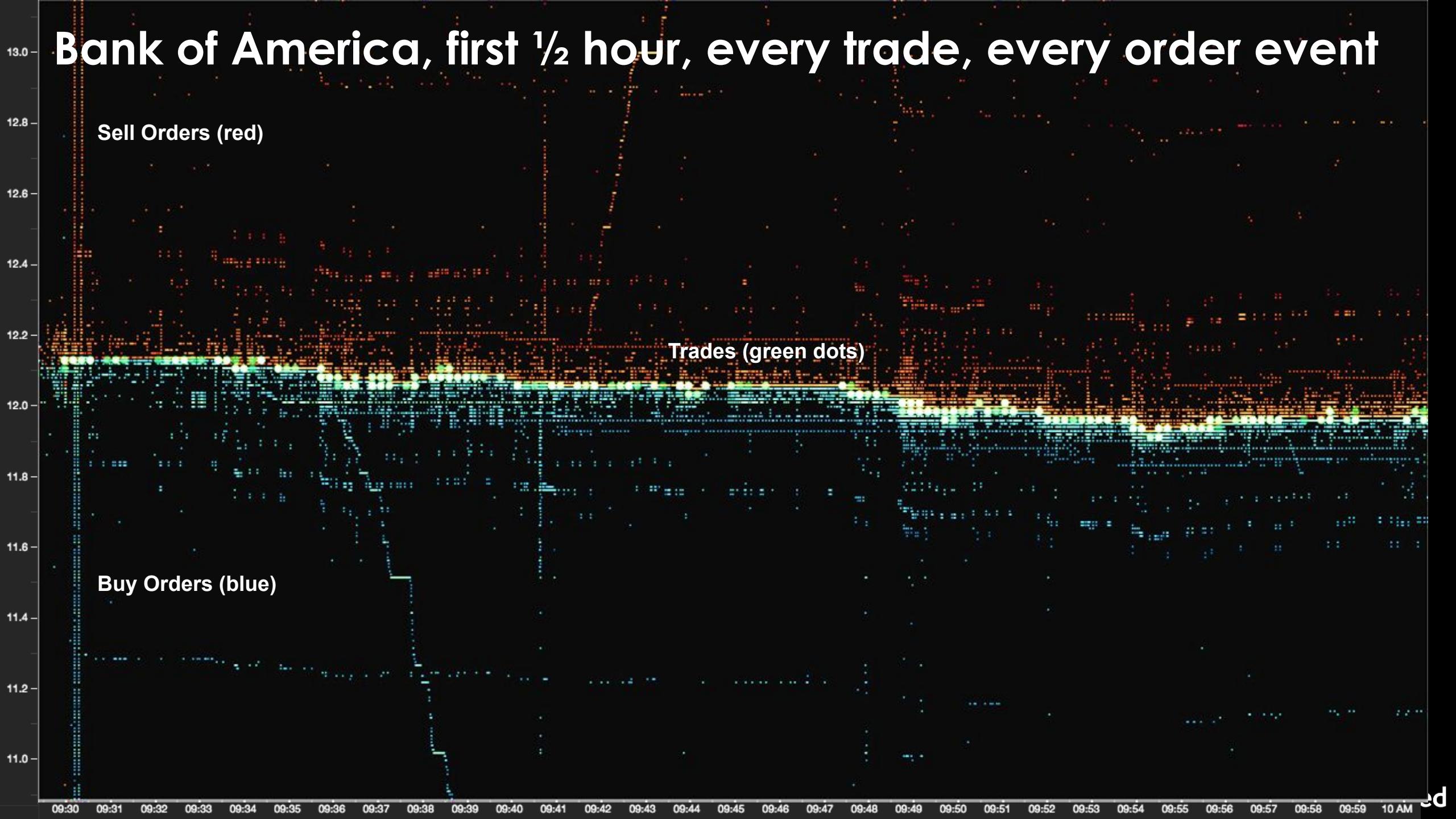


Stock Market: One Day Trades on 6800 stocks

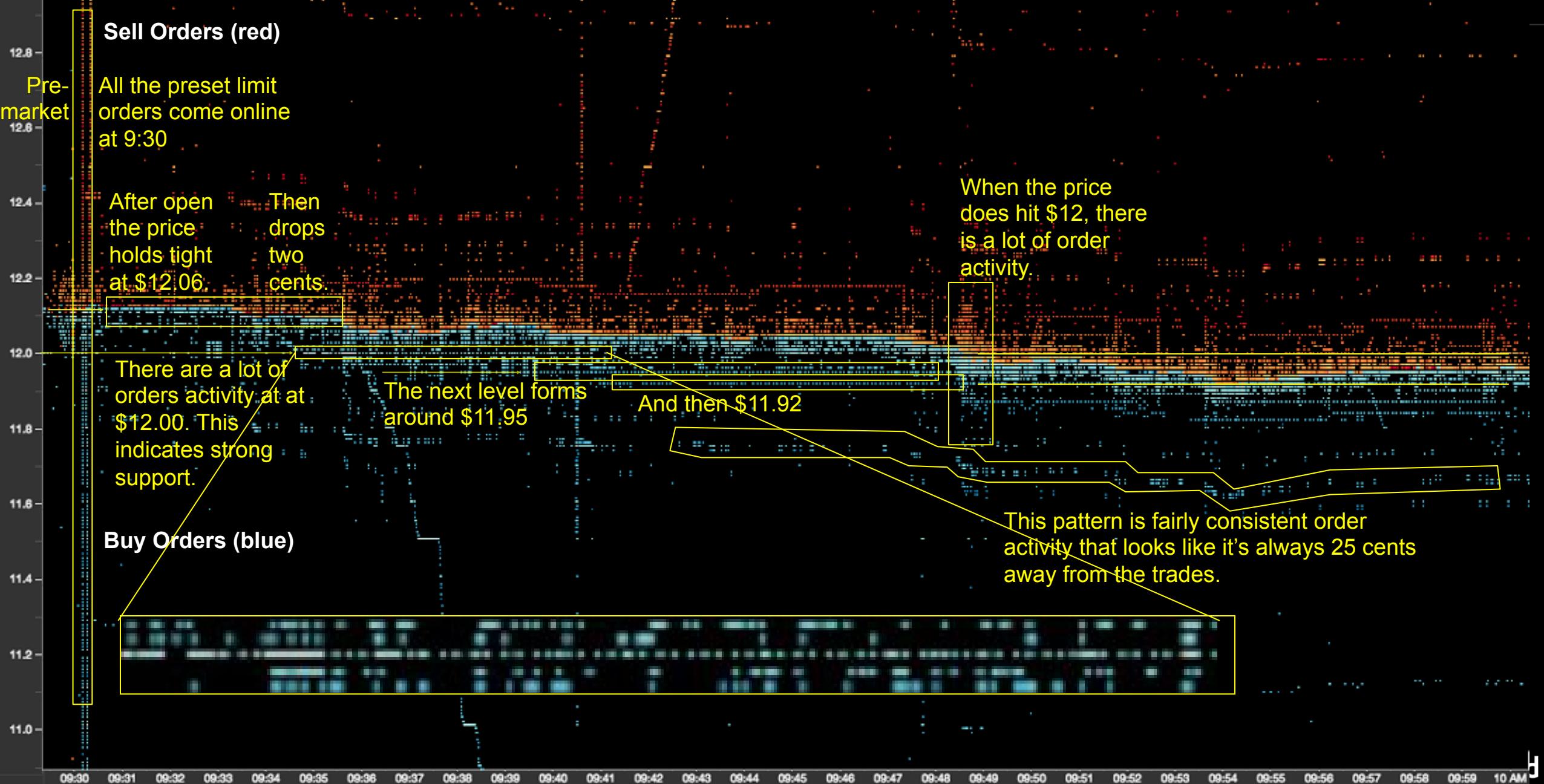
ARCA stocks
ranked by
percent
change on
day



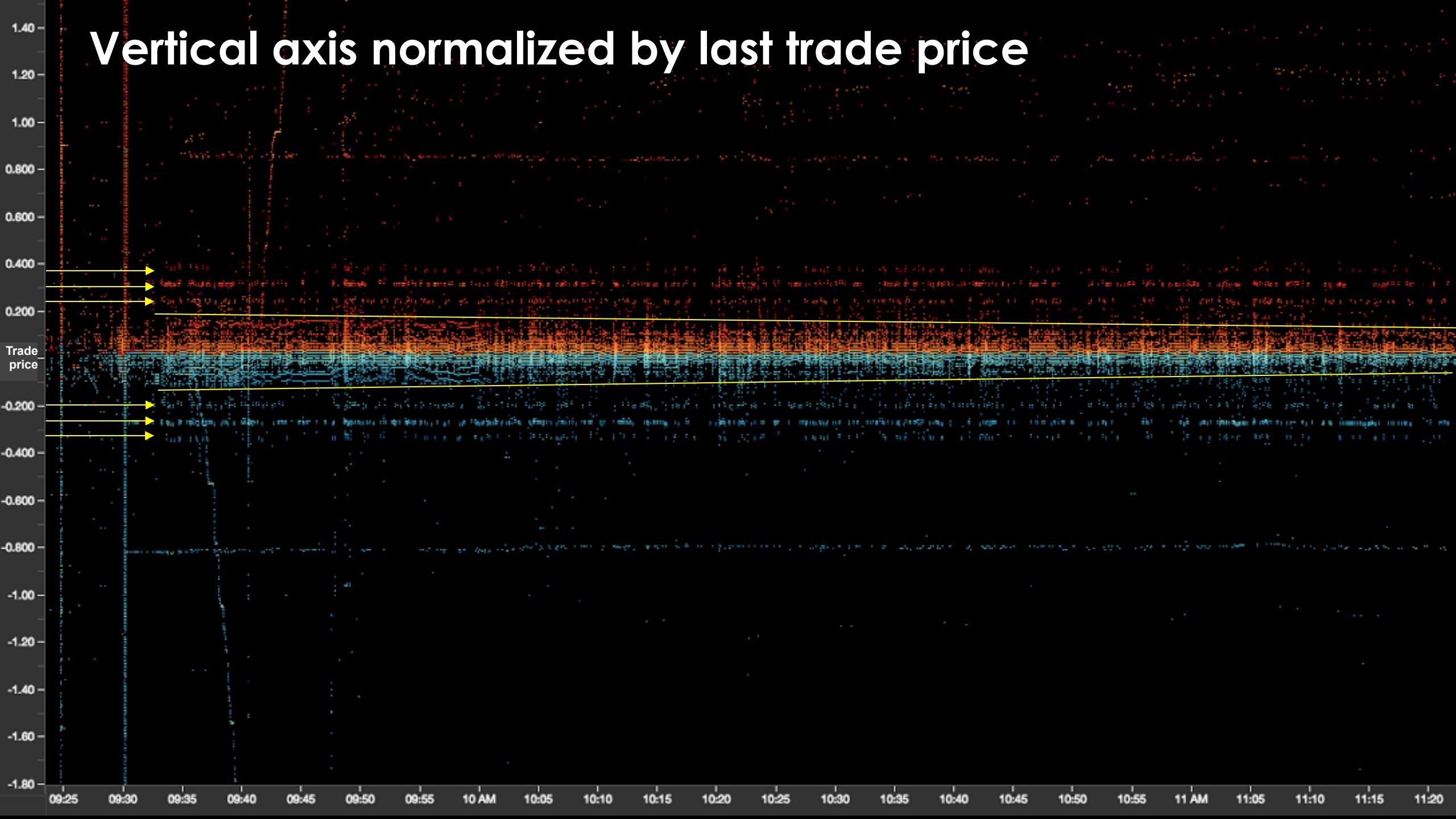
Bank of America, first ½ hour, every trade, every order event



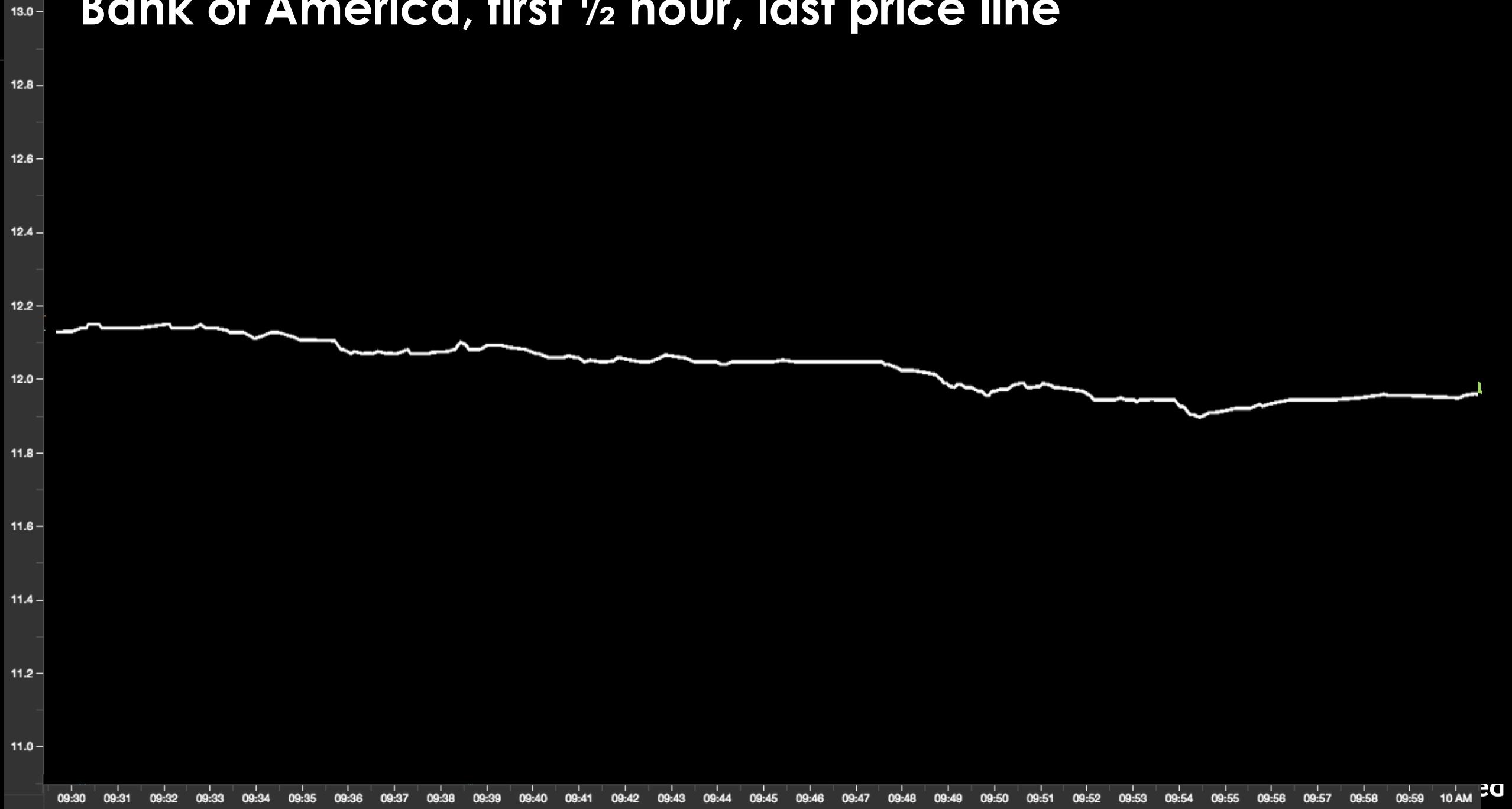
Bank of America, first ½ hour, every order event



Vertical axis normalized by last trade price



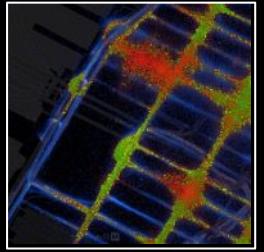
Bank of America, first ½ hour, last price line



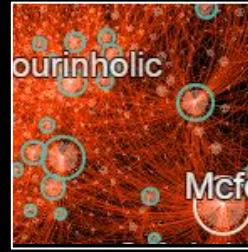
Technical Approach

Exploratory Big Data Analysis

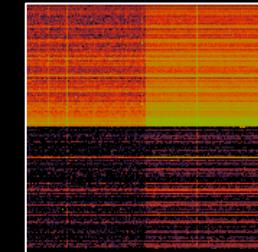
MAP



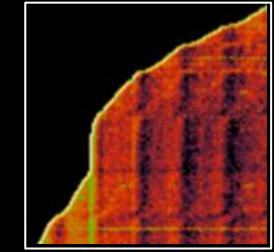
CONNECT

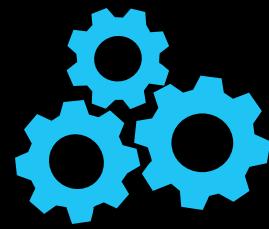
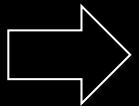
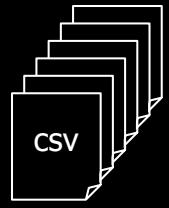


ORDER

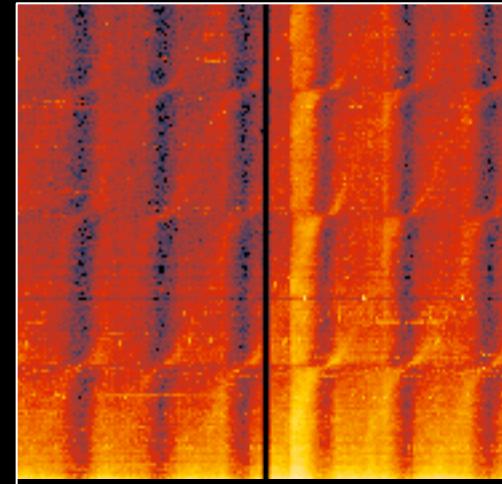
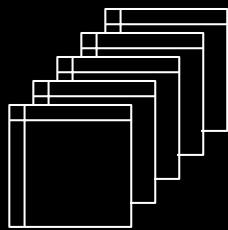
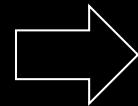


TIME





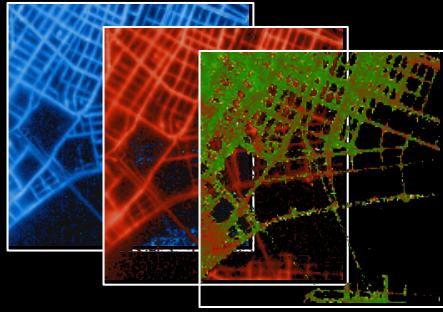
Spark



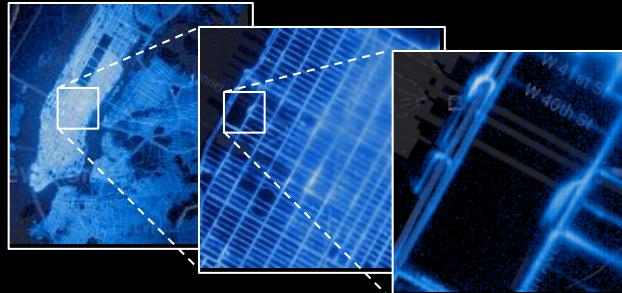
result.png

Exploratory Big Data Analysis – Rich Interactions

LAYERS



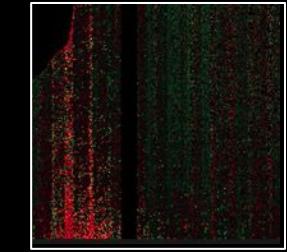
ZOOM



MINE



FILTER





Expressive API

Technical Approach

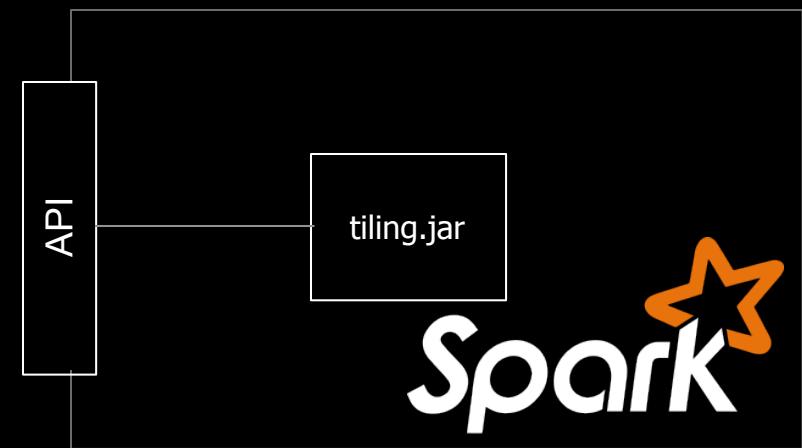
Data Pipeline

- Loading
- Filtering
- Transformation
- Sentiment
- Serialization

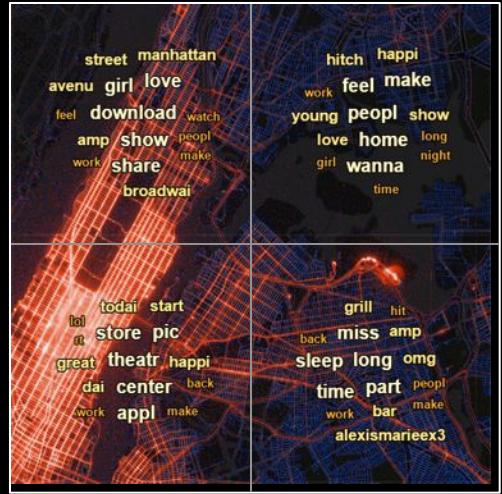
Results Generation

- Input mapping
- Projection
- “Pixel”-level analytics
- Area-level analytics
- Dataset-level analytics

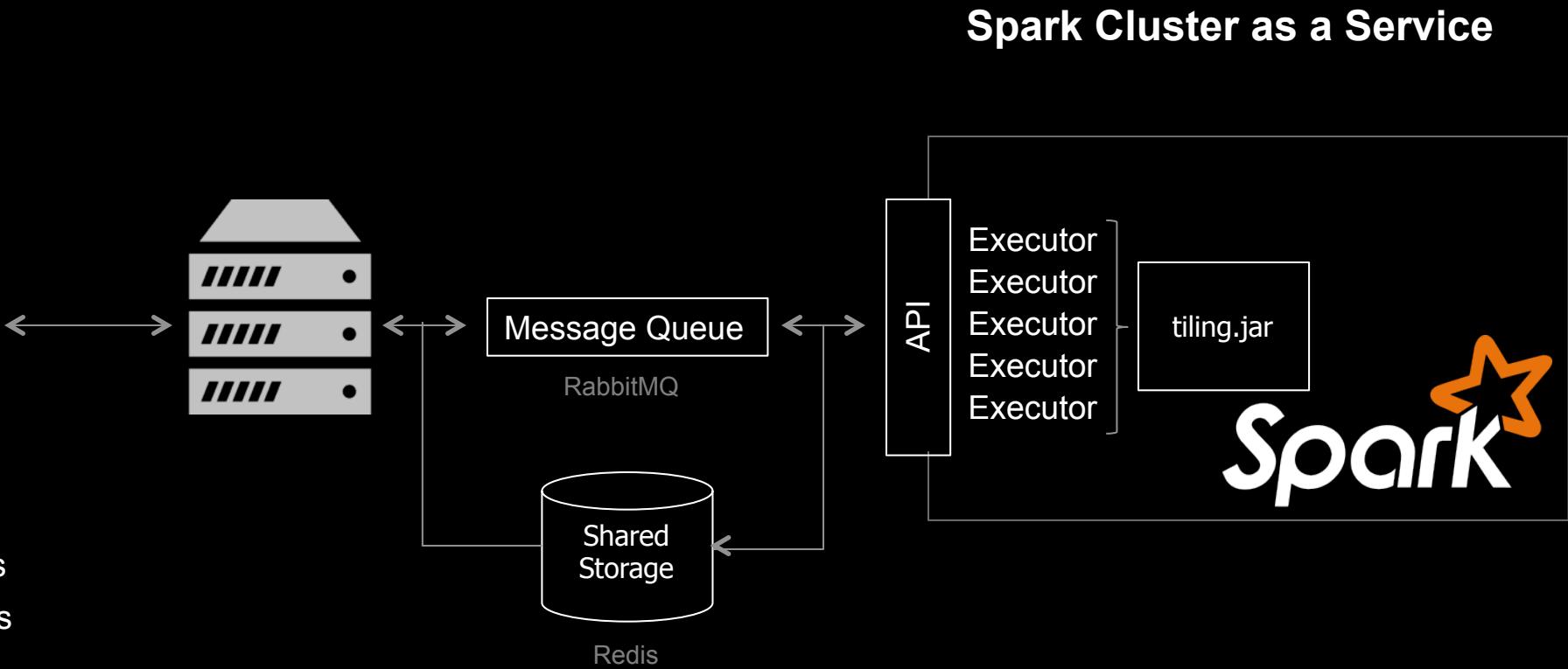
Spark Cluster as a Service



Technical Approach



configuration
density distributions
transformations term frequencies
summary statistics
data drilldown filters



Technical Approach



github.com/unchartedsoftware

Trump Tweets



Donald J. Trump

@realDonaldTrump

Follow

#AskTrump Getting ready to answer your questions.

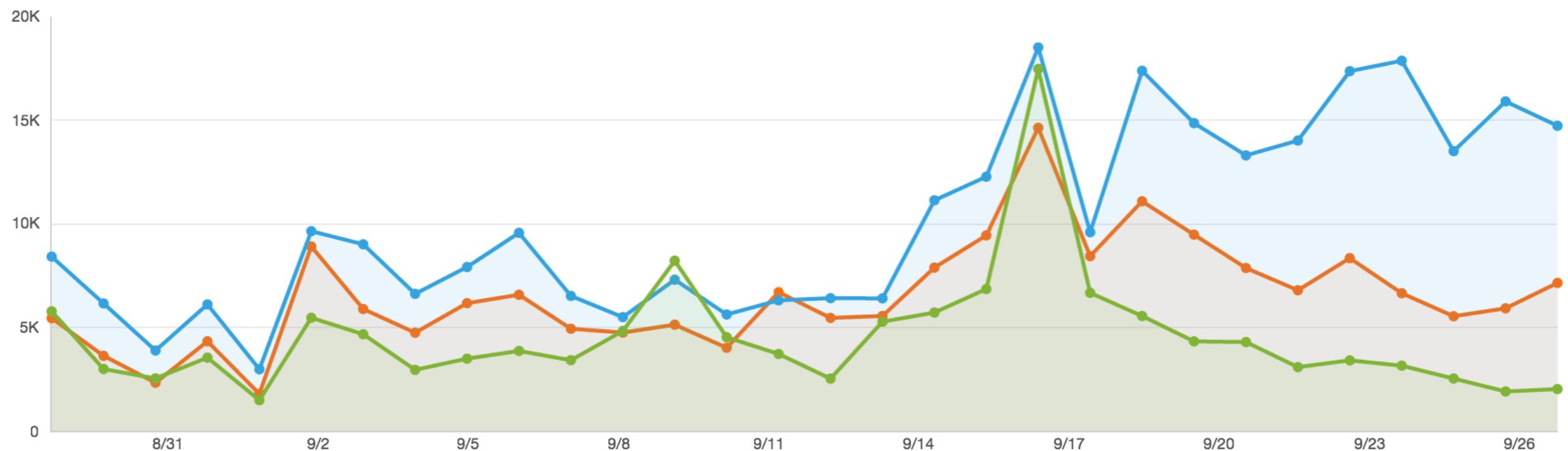
1:07 PM - 21 Sep 2015

4 1,037 ★ 2,429



Tweets per day: #makeamericagreatagain, #trump2016, and #donaldtrump

August 28th — September 27th



#makeamericagreatagain

195,425

#trump2016

304,648

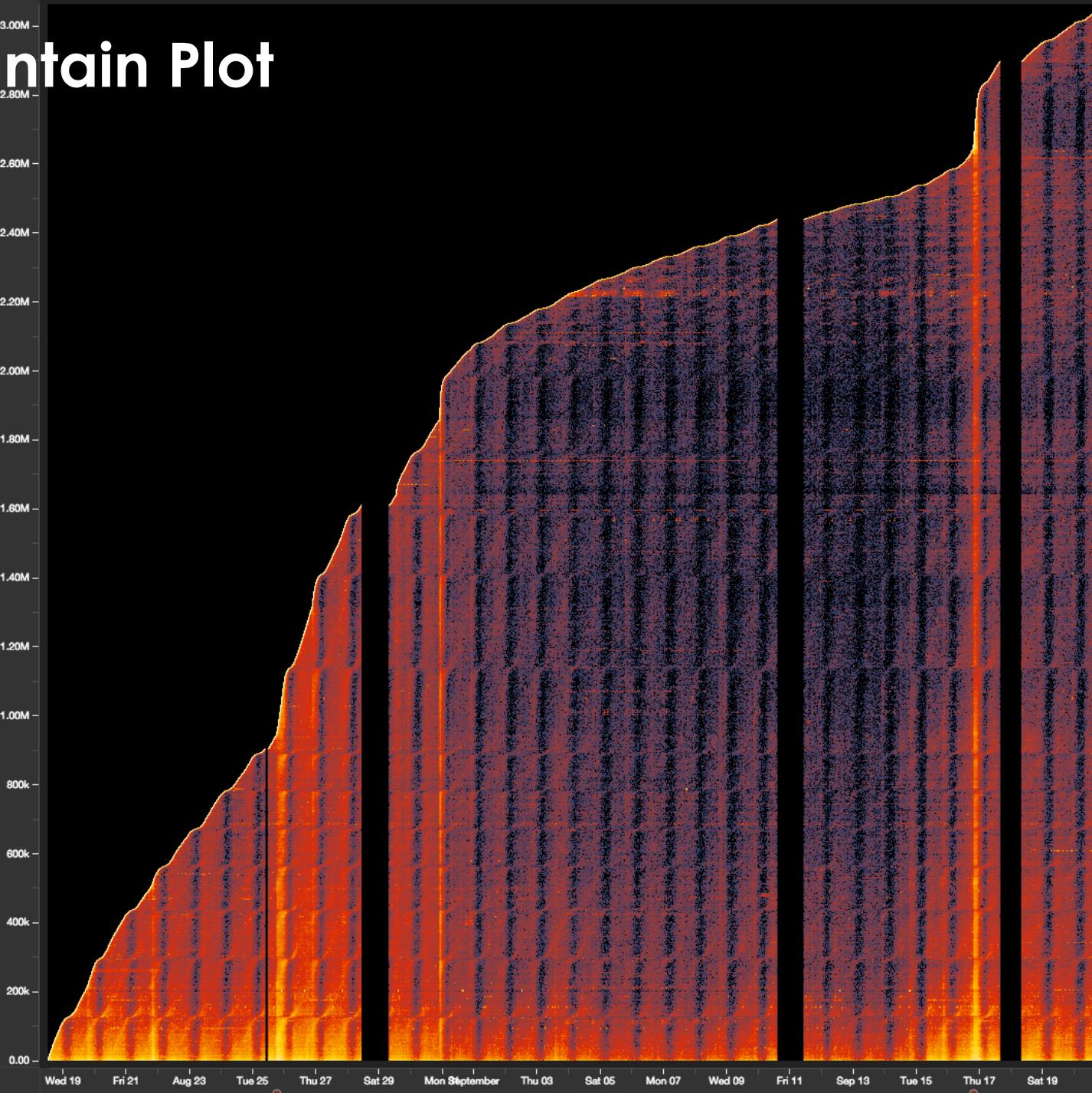
#donaldtrump

136,136

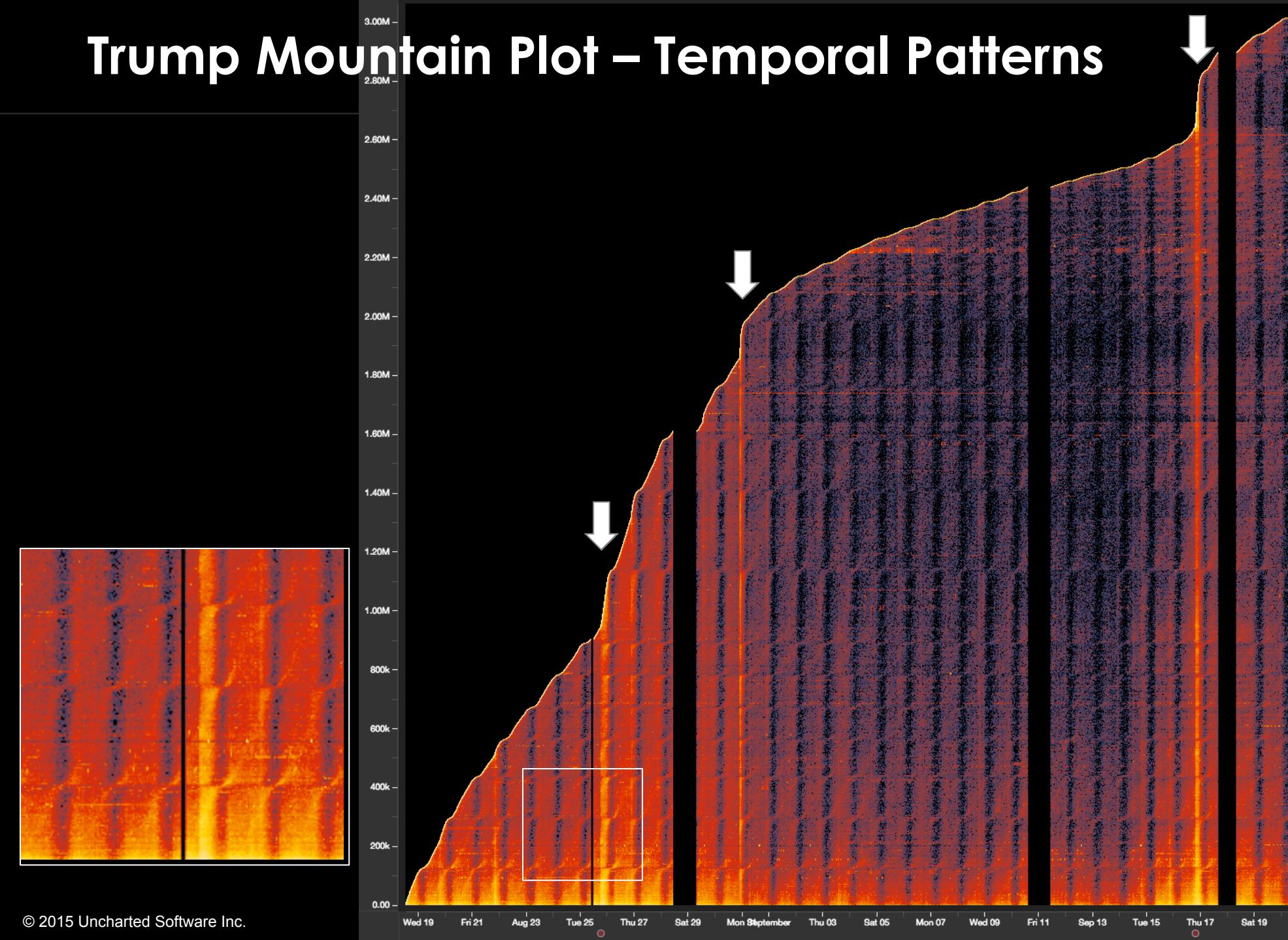
ANALYTICS BY **TOPSY**

[Continue](#)

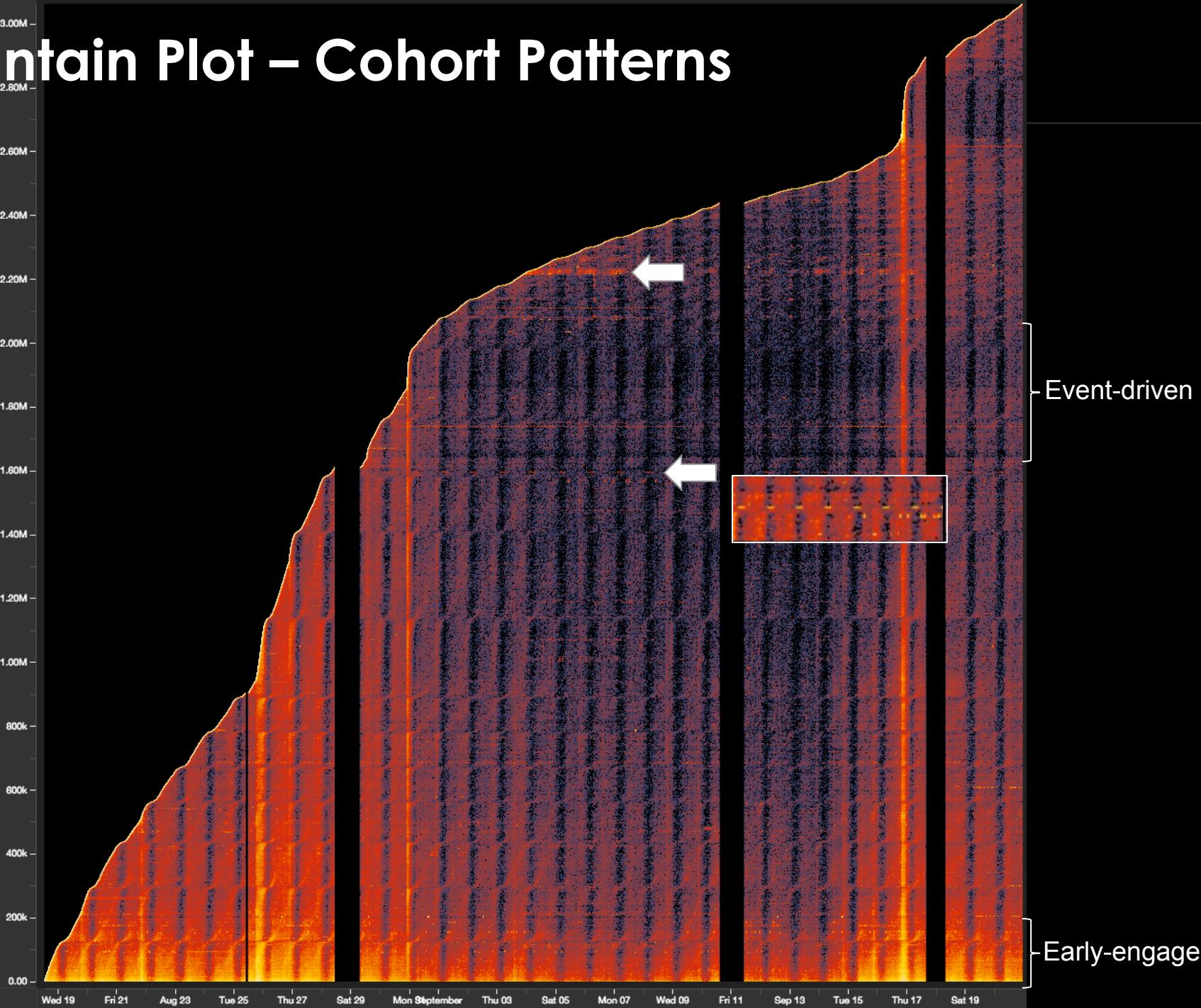
Trump Mountain Plot



Trump Mountain Plot – Temporal Patterns

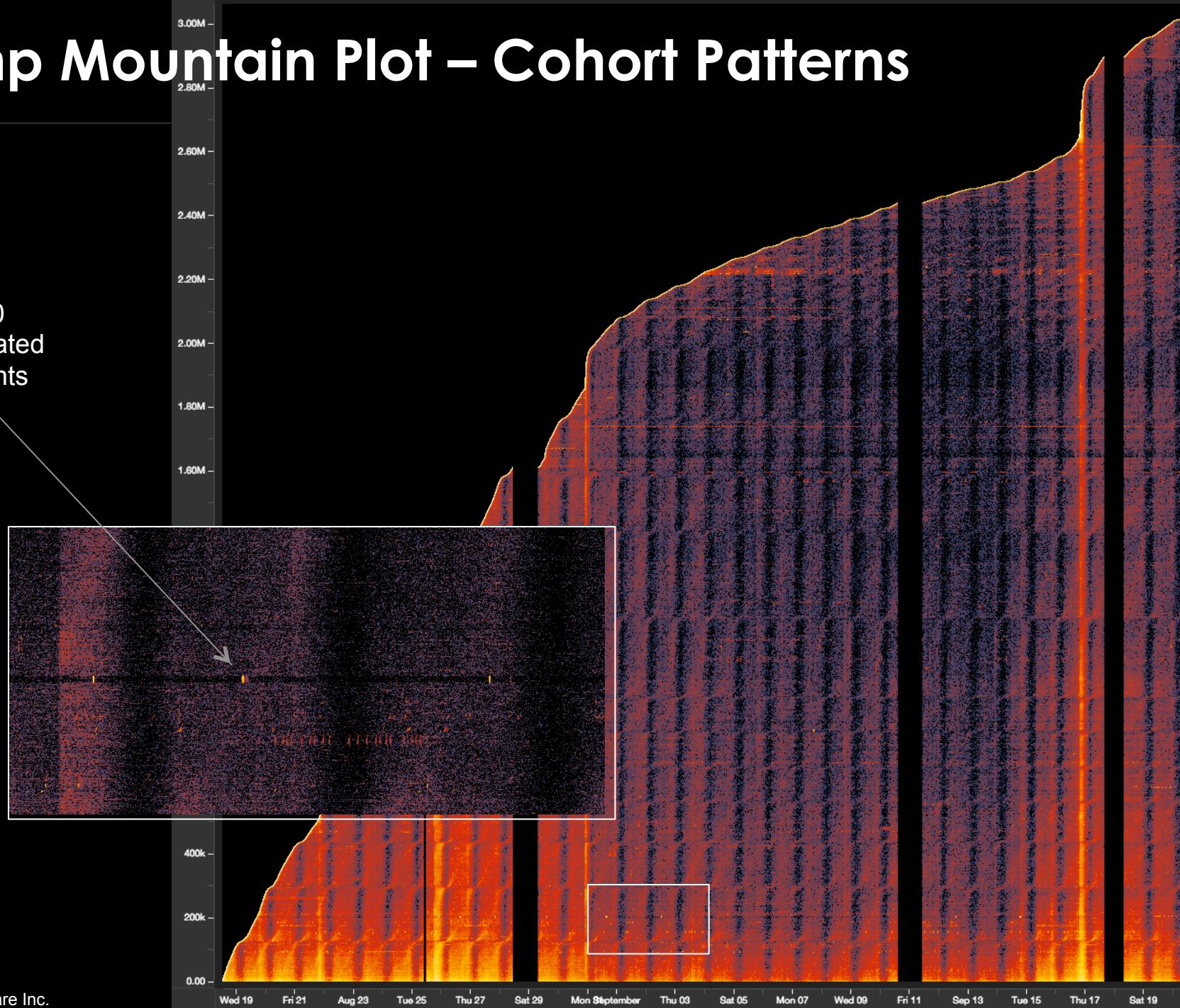


Trump Mountain Plot – Cohort Patterns

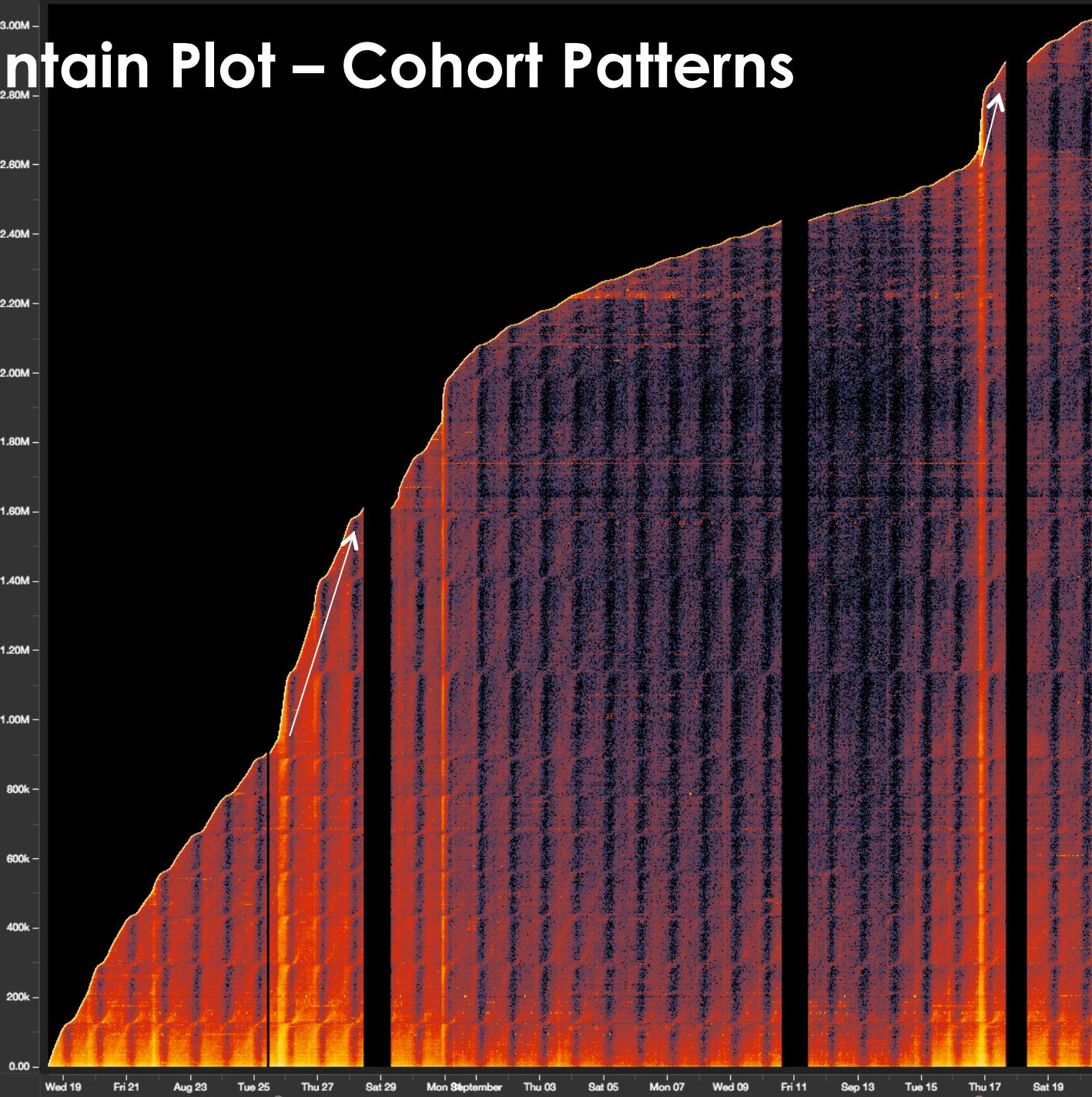


Trump Mountain Plot – Cohort Patterns

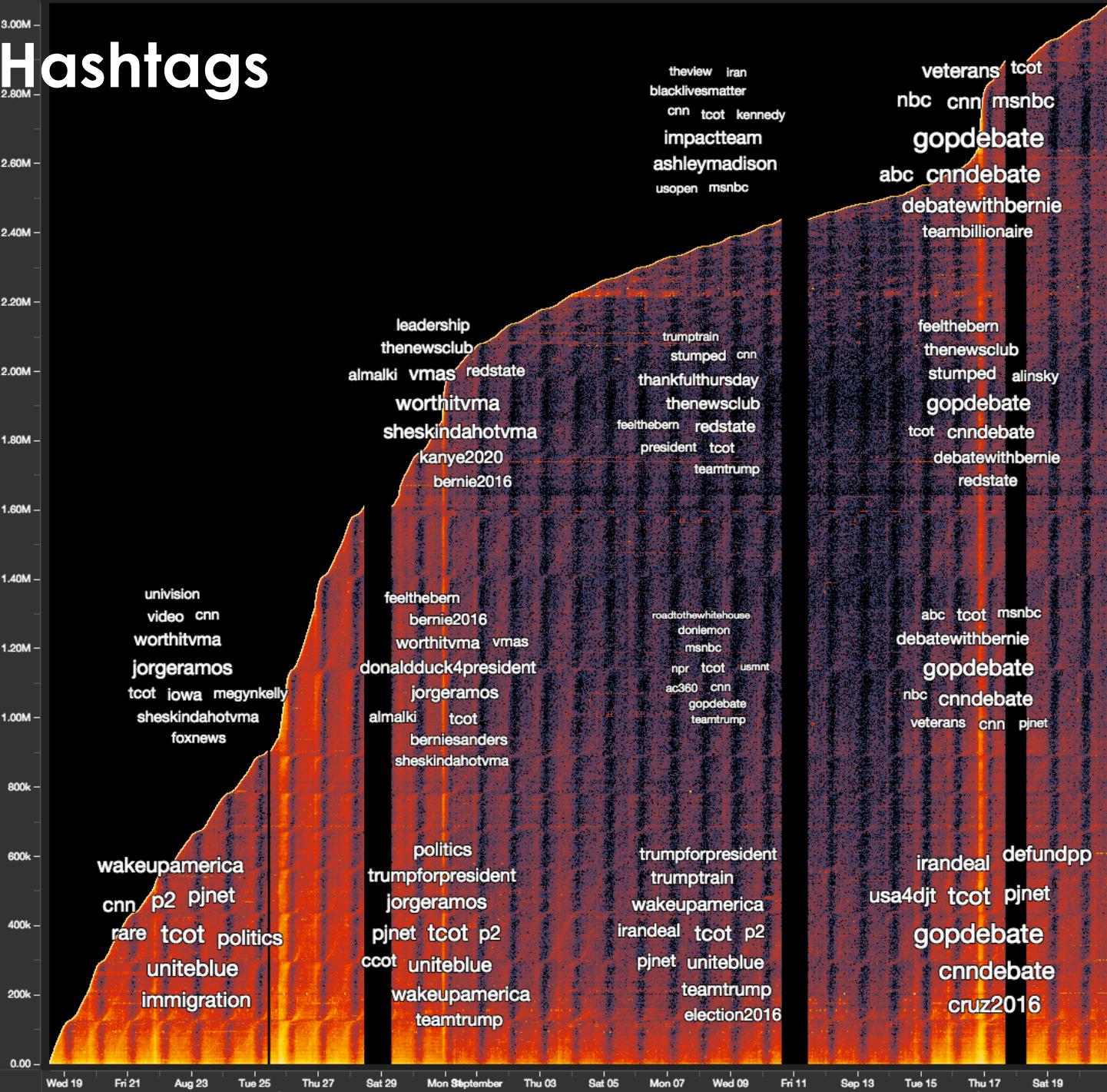
4000
coordinated
accounts



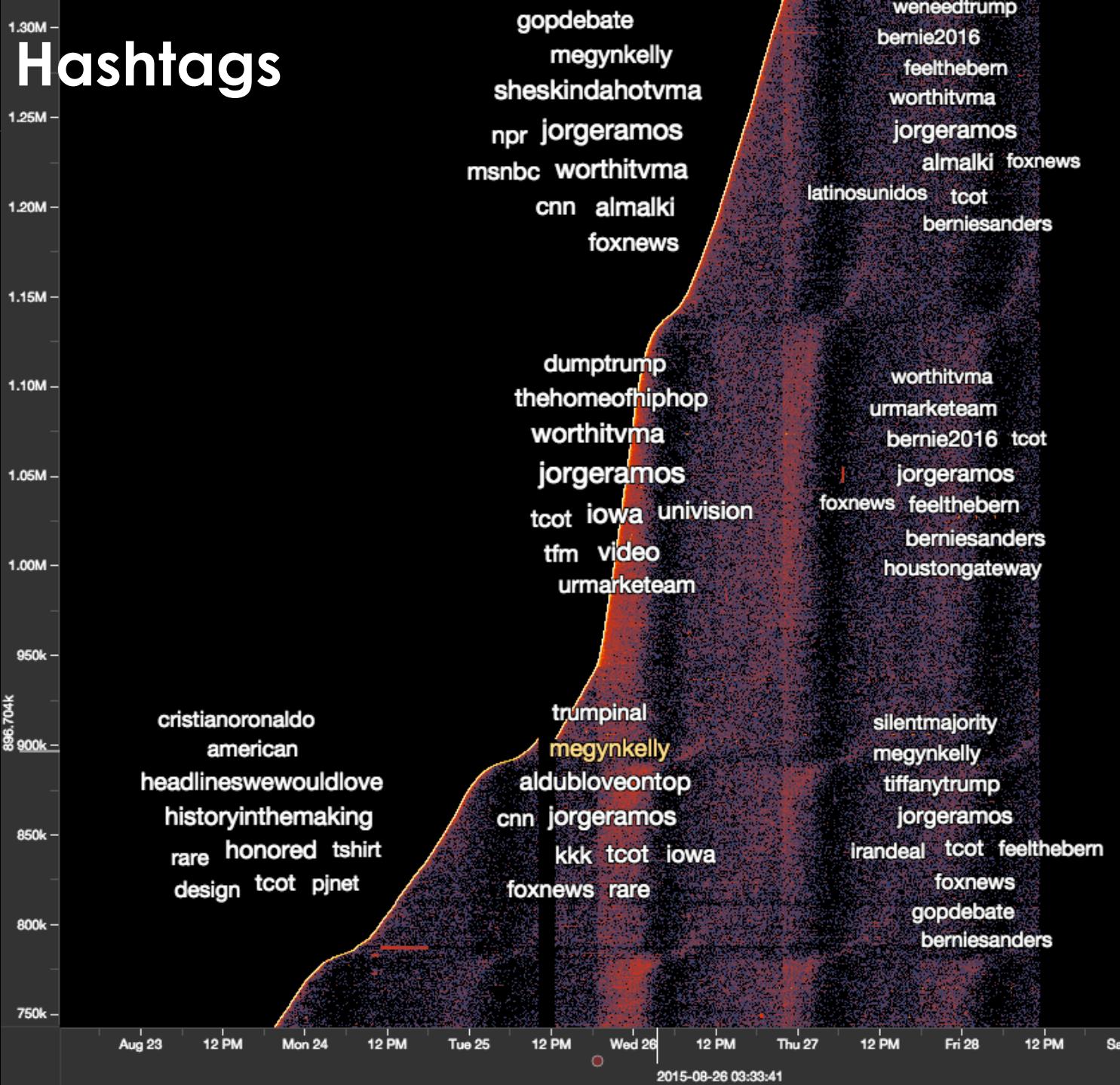
Trump Mountain Plot – Cohort Patterns



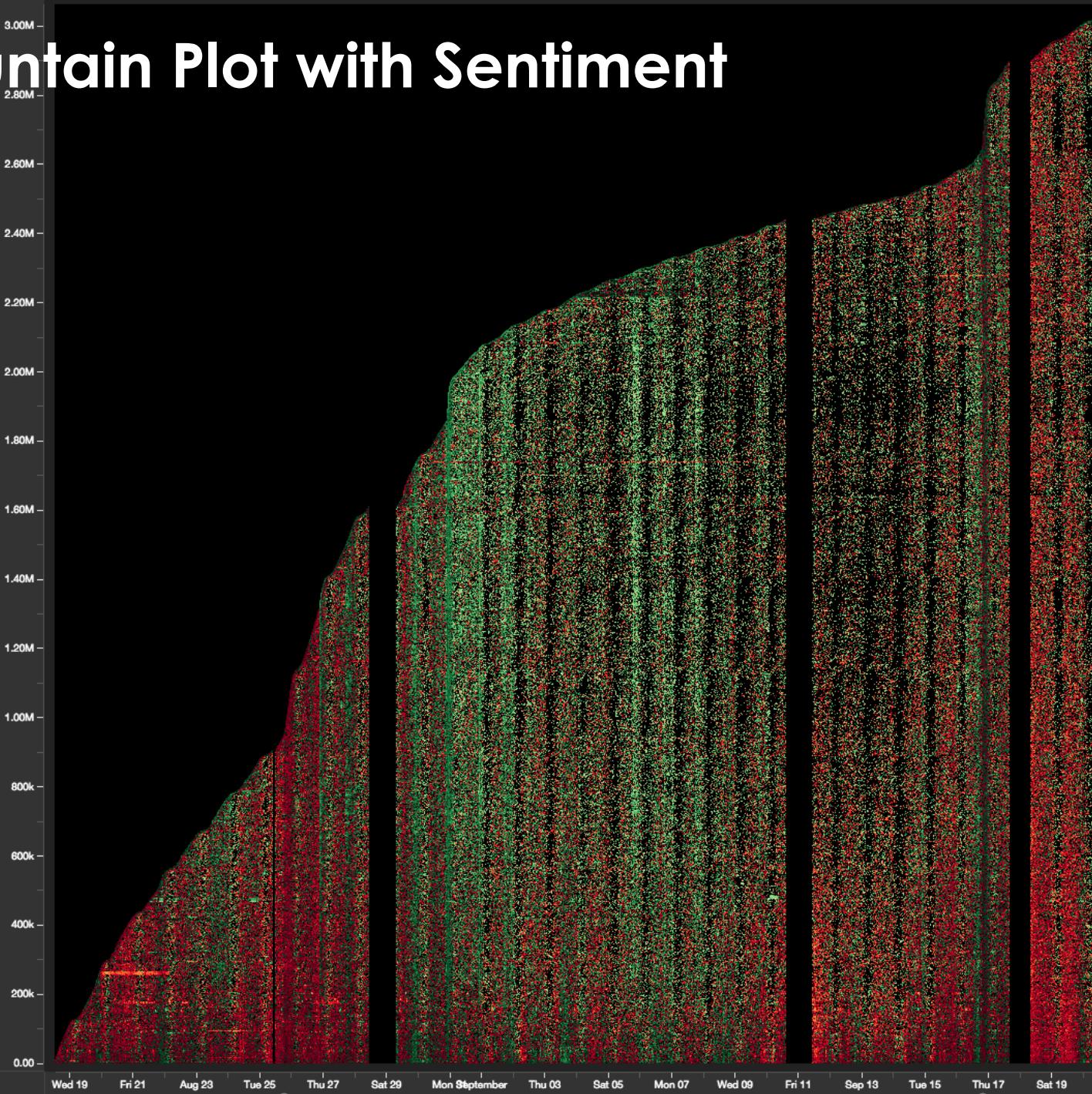
Trump Top Hashtags



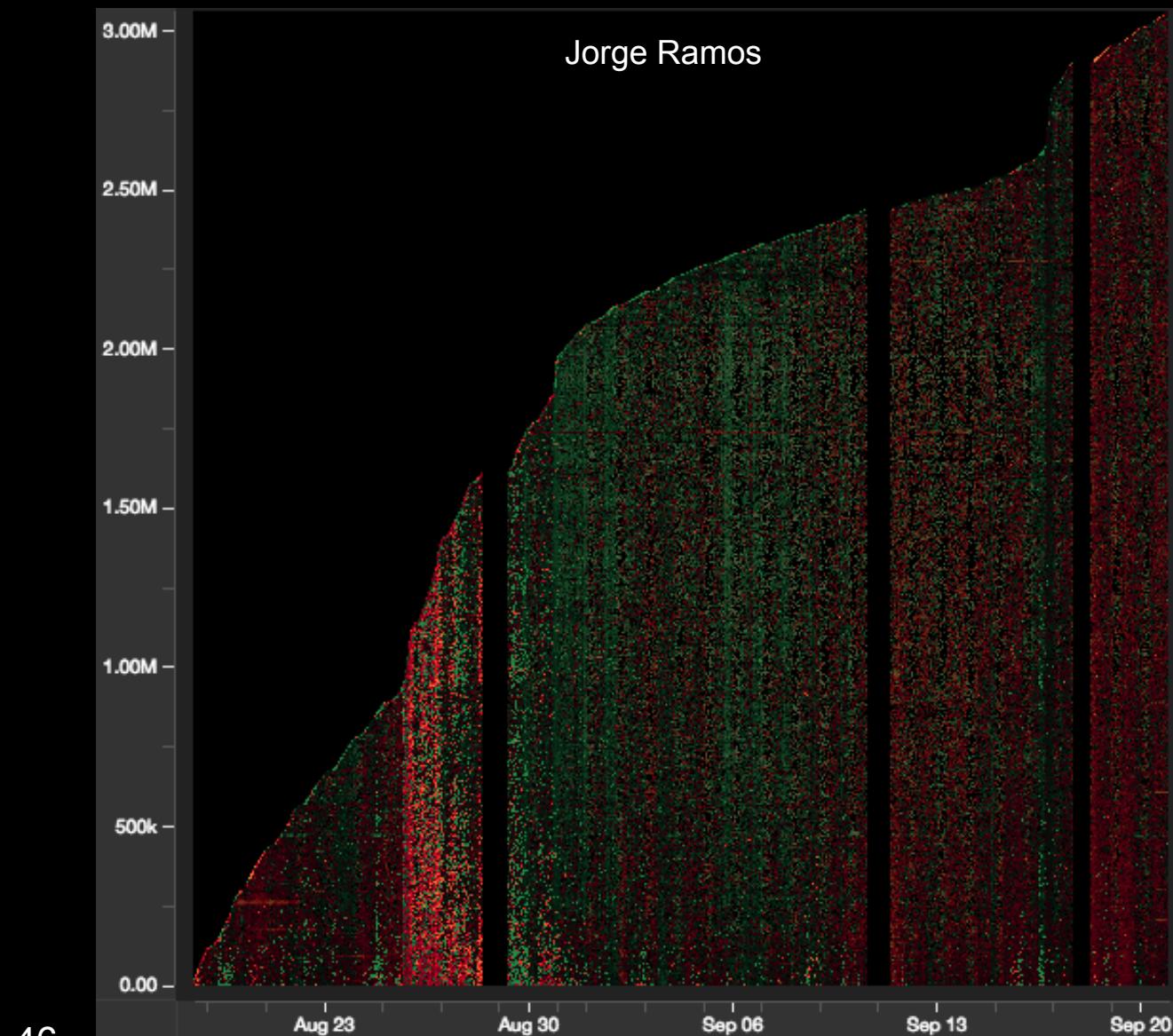
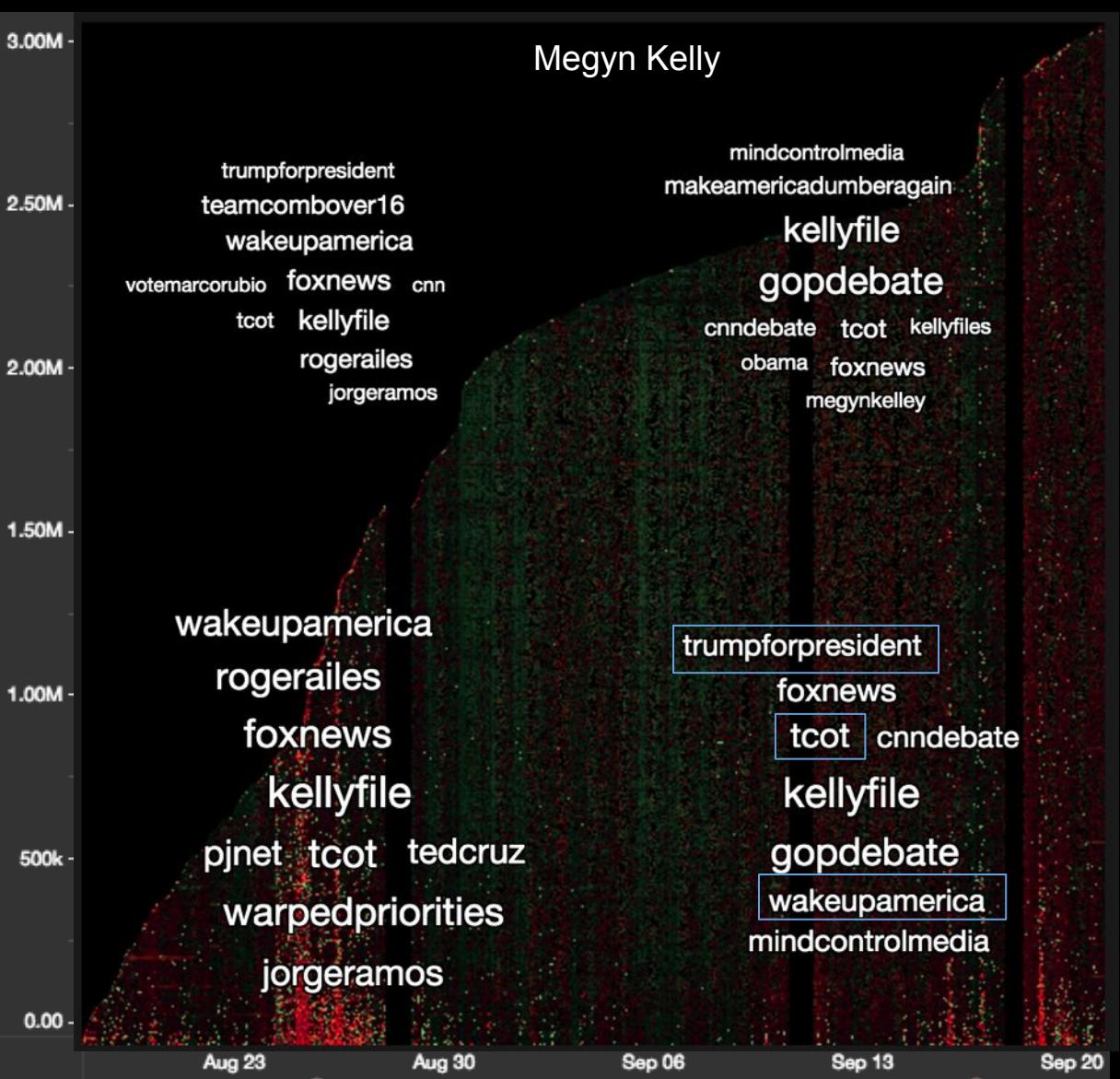
Trump Top Hashtags



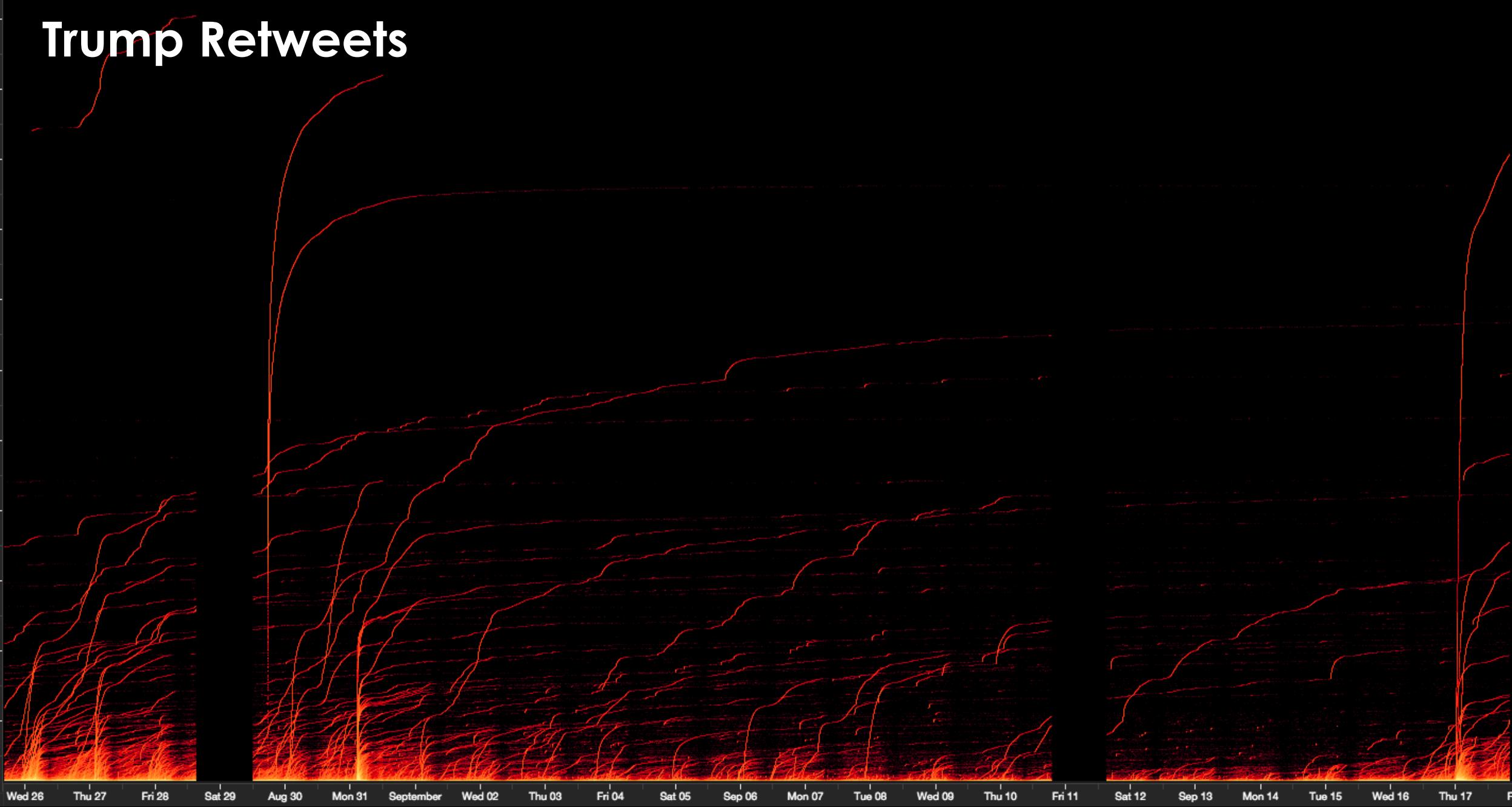
Trump Mountain Plot with Sentiment



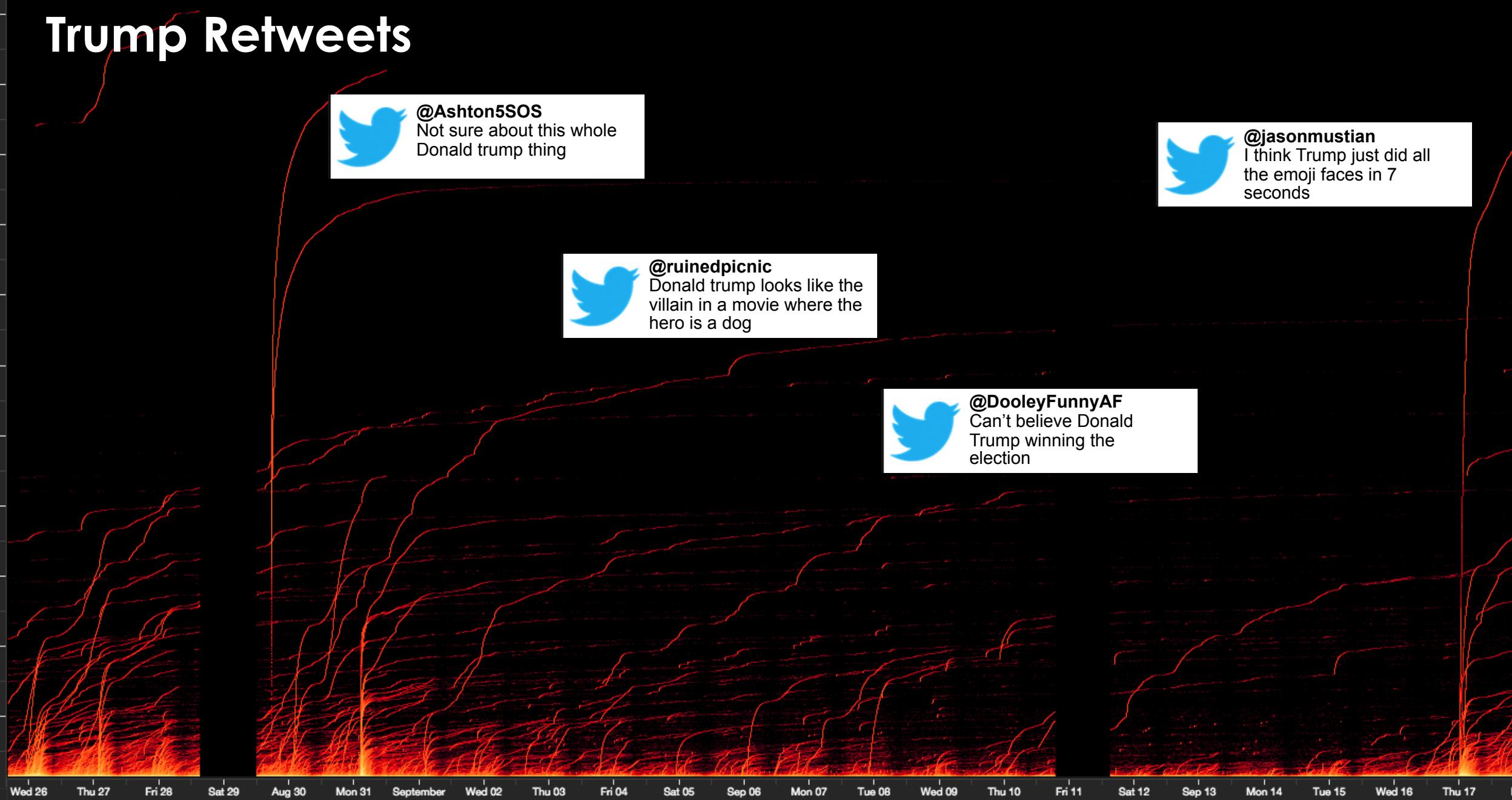
Trump Mountain Plot with Sentiment – Megyn / Jorge



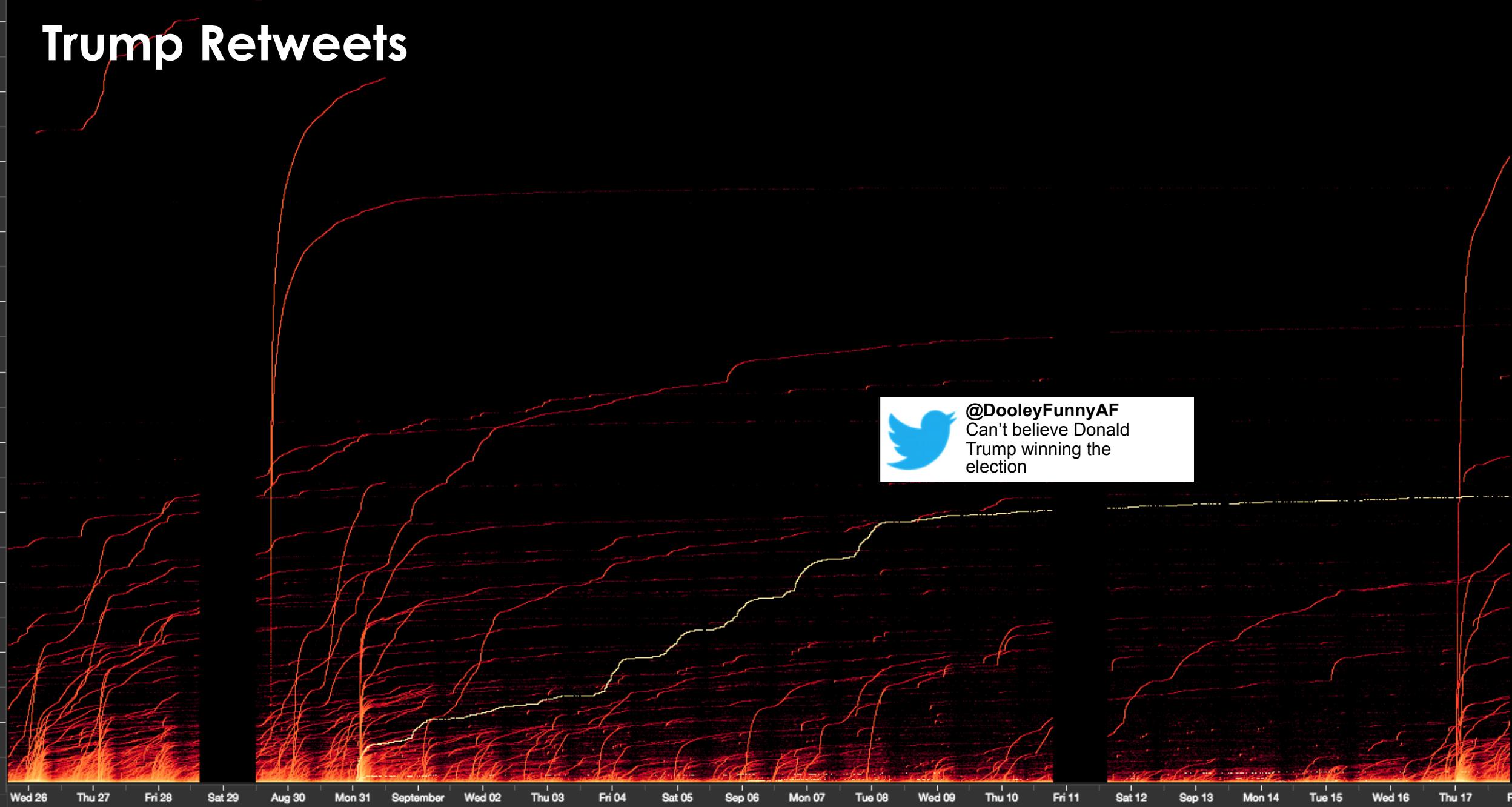
Trump Retweets



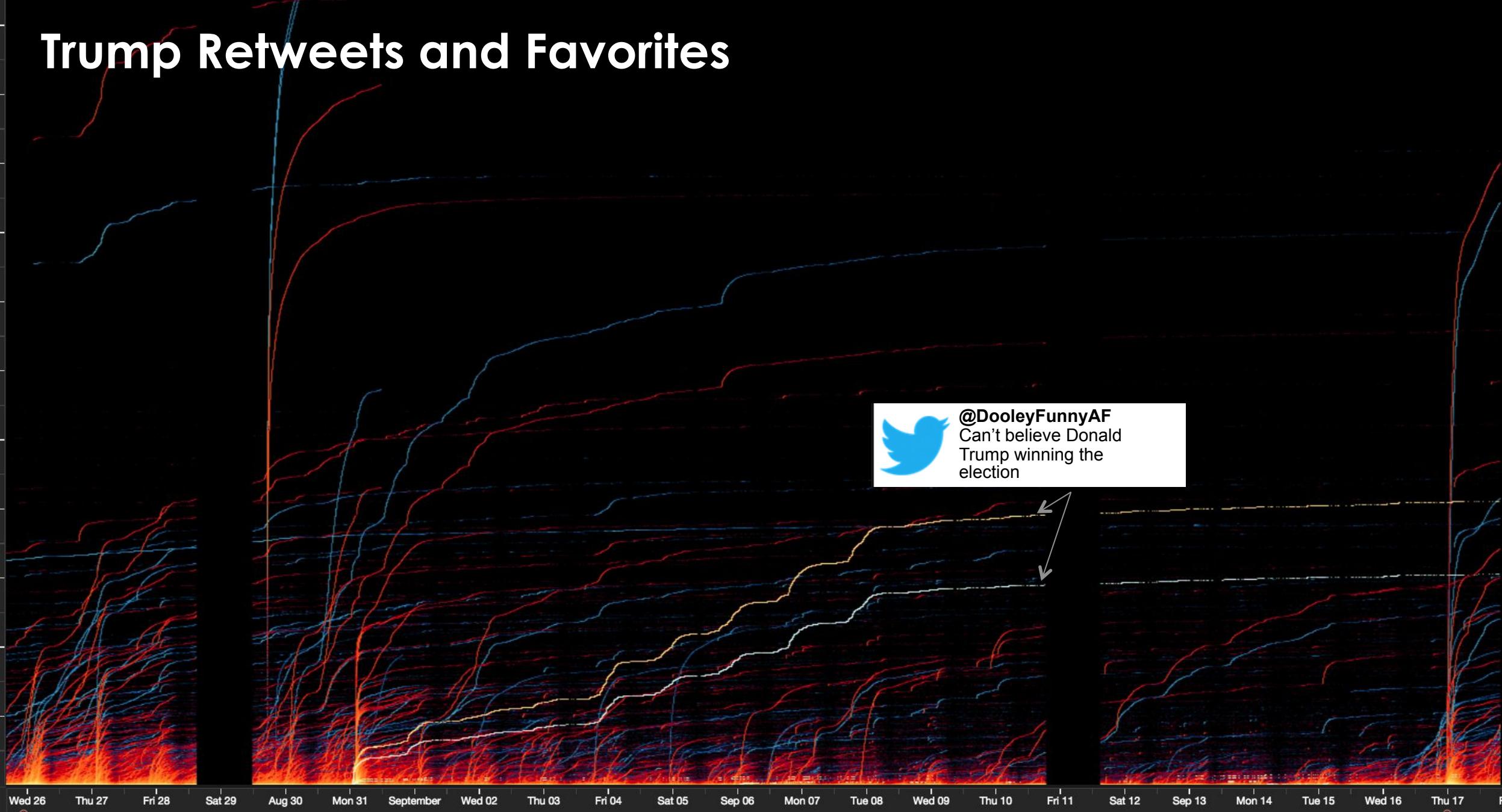
Trump Retweets



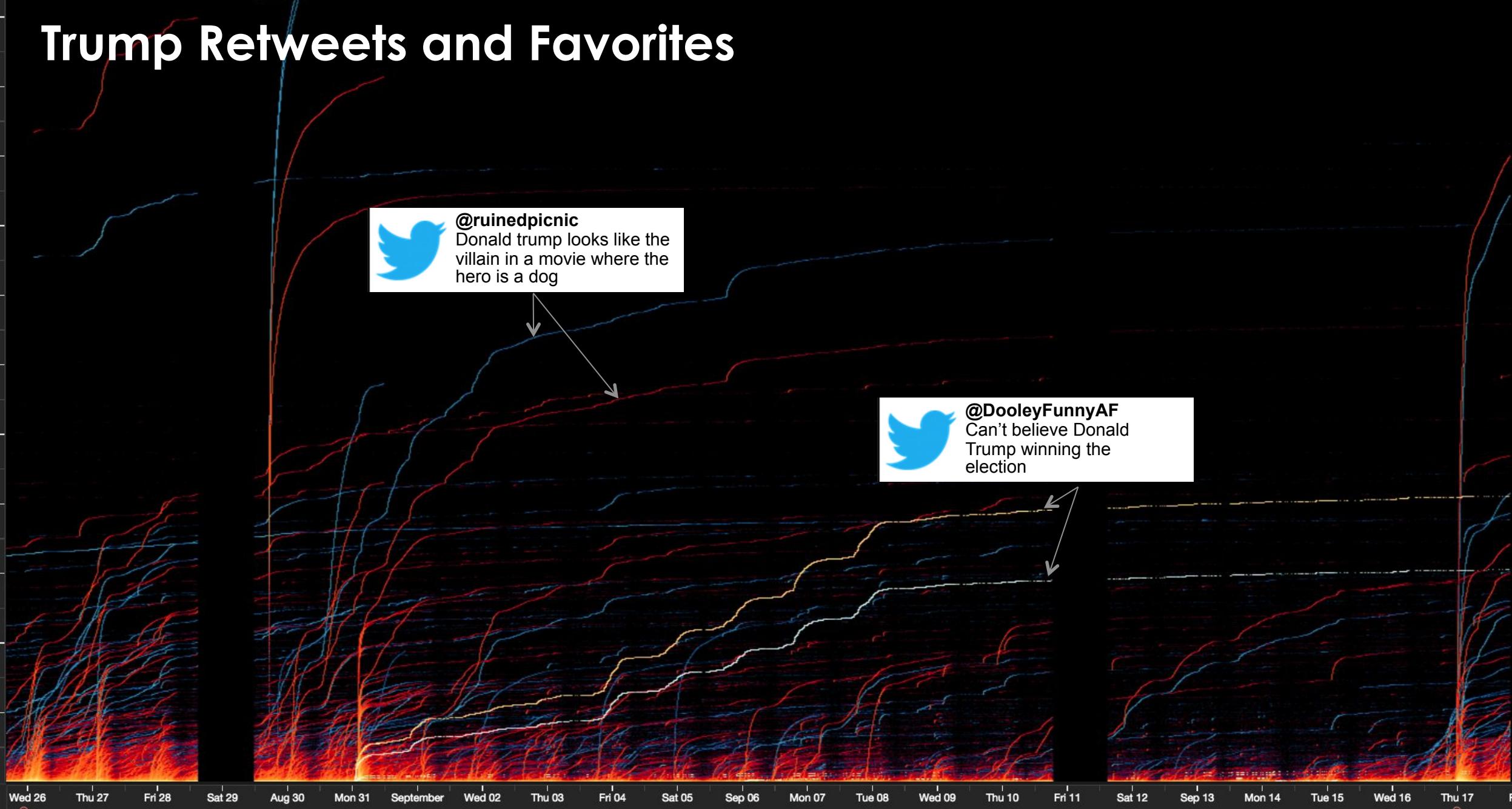
Trump Retweets



Trump Retweets and Favorites



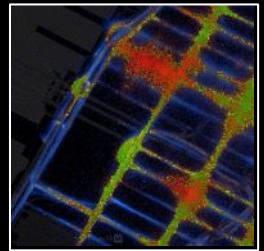
Trump Retweets and Favorites



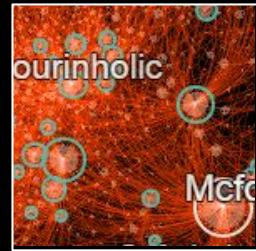
KEY TAKE AWAYS

1. Plot all the data

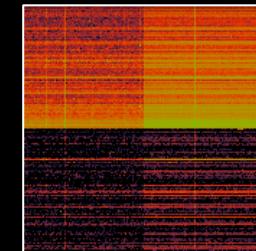
MAP



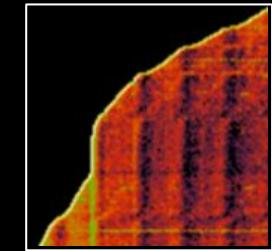
CONNECT



ORDER

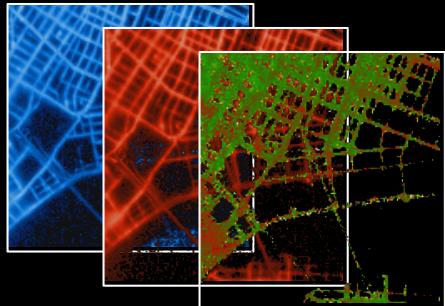


TIME

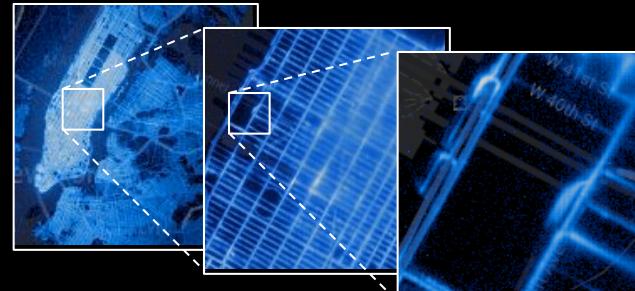


2. Explore it

LAYERS



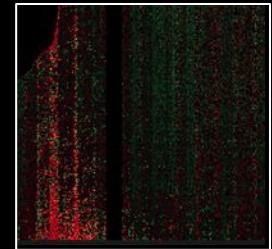
ZOOM



MINE



FILTER



More Info



Richard Brath

rbrath@uncharted.software
416-203-3003 x 242



Robert Harper

rharper@uncharted.software
@rdharper