

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

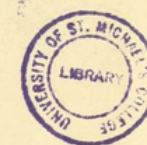
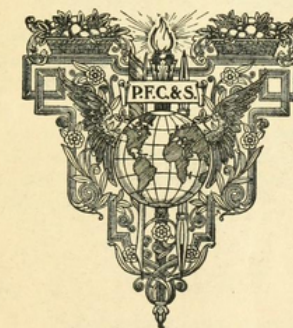


Natural
Selection?

ORIGIN OF SPECIES

BY MEANS OF NATURAL SELECTION, OR THE
PRESERVATION OF FAVORED RACES IN
THE STRUGGLE FOR LIFE

BY
CHARLES DARWIN, M.A., LL.D., F.R.S.



NEW YORK
P. F. COLLIER & SON
MCMII

1

Visual Exploration Process

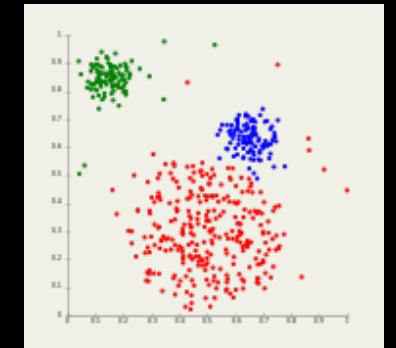
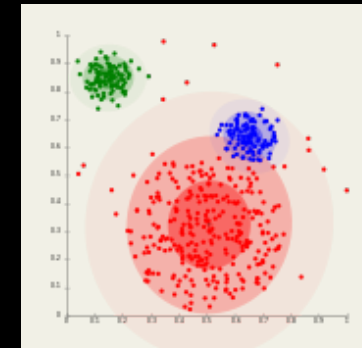
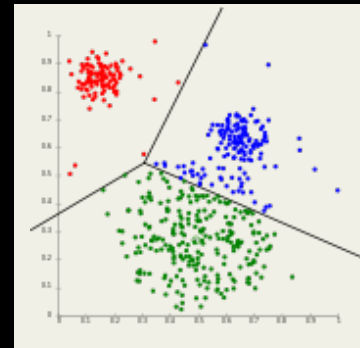
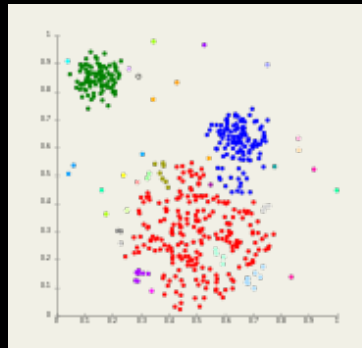
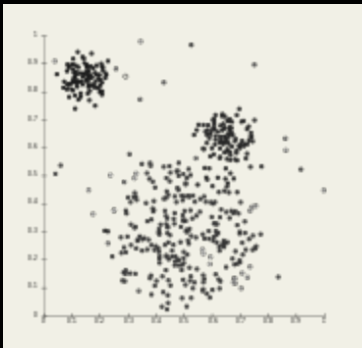
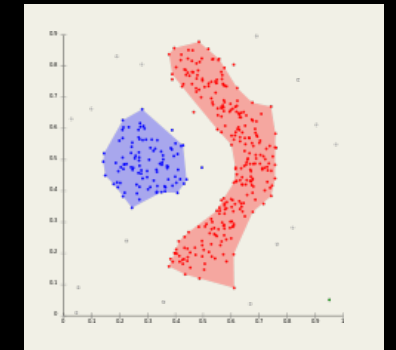
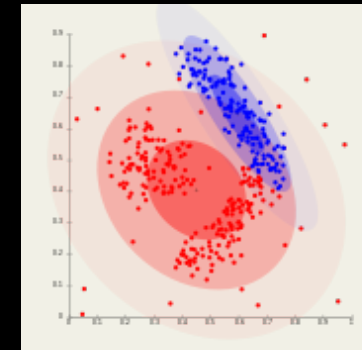
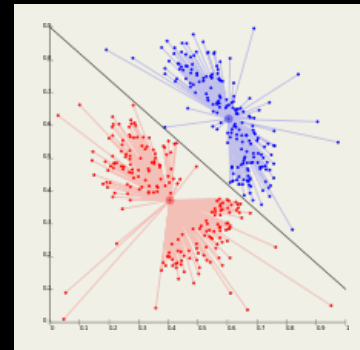
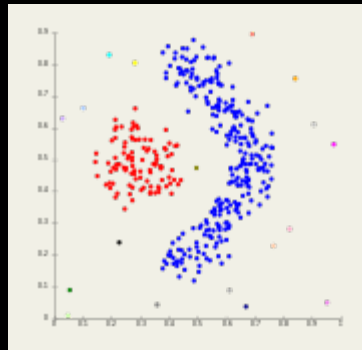
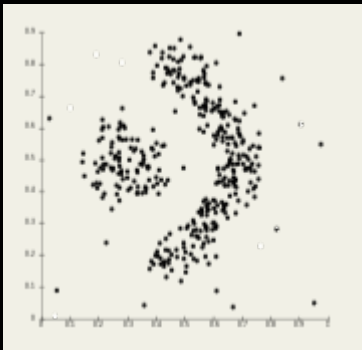
1. Collect a lot of data
2. Observe some interesting patterns
3. Hypothesize about why they exist
4. Refine and build models



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.*

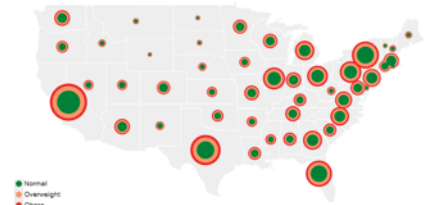


paraphrased from John Behrens, Arizona State University, Principles and Procedures of Exploratory Data Analysis, American Psychological Association, 1997.

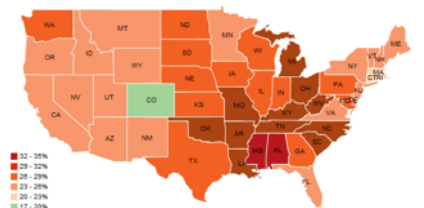
So, why are we summarizing
big data into bar charts?

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

Graduated Symbol Maps



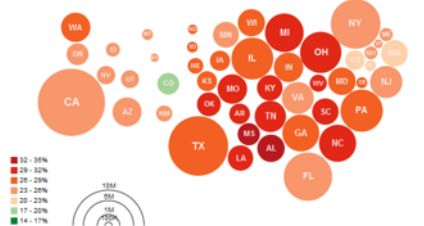
Choropleth Maps



Flow Maps



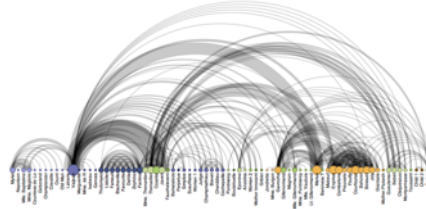
Cartograms



Force-Directed Layout

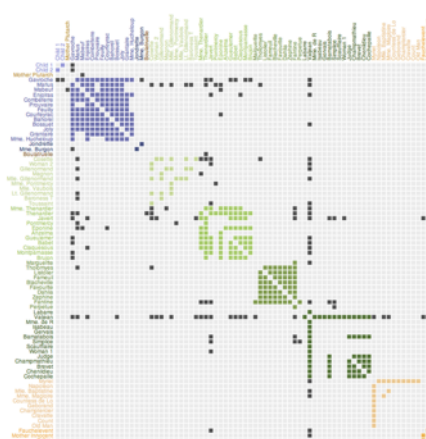


Arc Diagrams

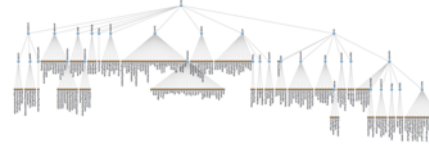


An arc diagram uses a one-dimensional layout of nodes, with circular arcs to represent links. While arc diagrams may not convey the overall structure of the graph as effectively as a two-dimensional layout, with a good ordering of nodes it is easy to identify cliques and bridges. And, as with the indented tree layout, multivariate data can easily be displayed alongside nodes. The problem of sorting the nodes in a manner that reveals underlying cluster structure is formally called *seriation*, and has diverse applications in visualization, statistics, and even archaeology!

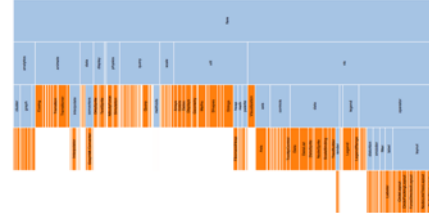
Matrix Views



Node-Link Diagrams



Adjacency Diagrams

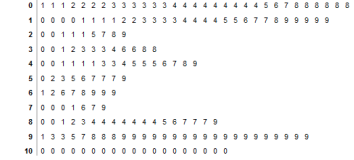


Enclosure Diagrams

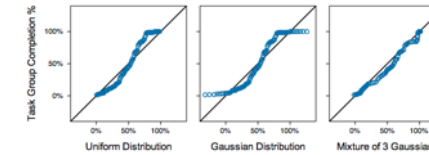


A Tour through the Visualization Zoo
Heer, Bostock, Ogievetsky
<http://homes.cs.washington.edu/~jheer/files/zoo/>

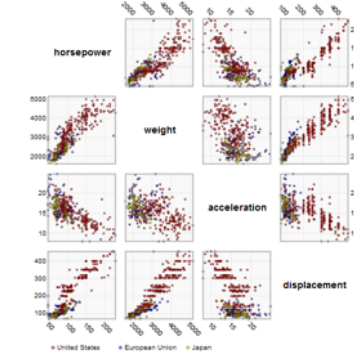
Stem-and-Leaf Plots



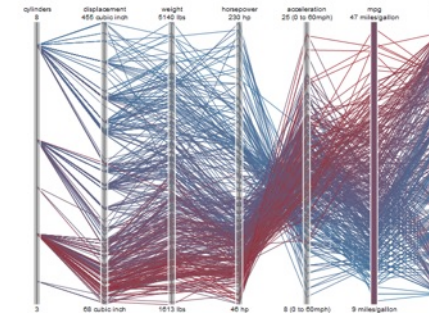
Q-Q Plots



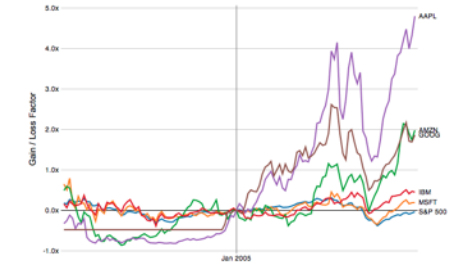
Scatter Plot Matrix (SPLOM)



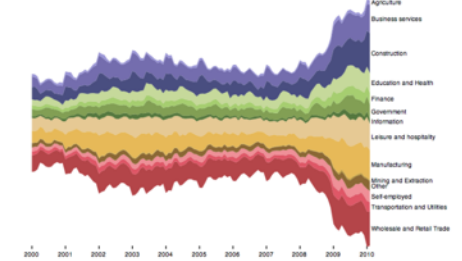
Parallel Coordinates



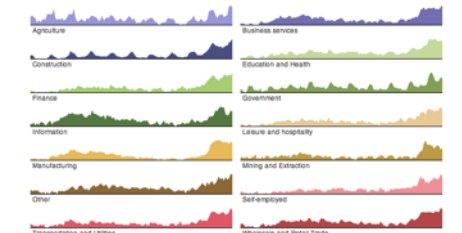
Index Charts



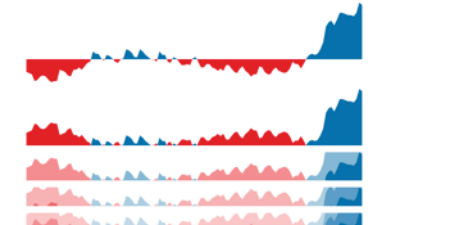
Stacked Graphs



Small Multiples

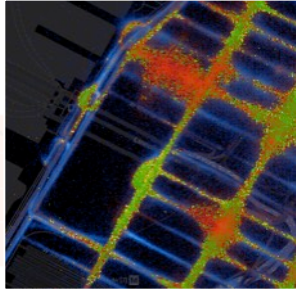


Horizon Graphs

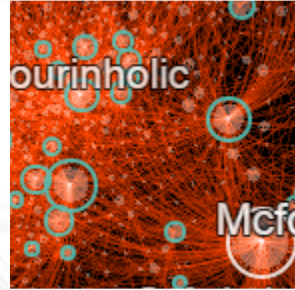


So, how do we visualize 100m data points?

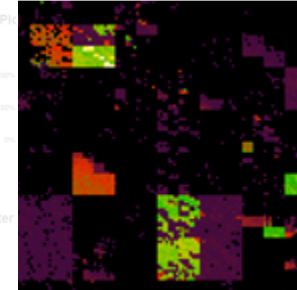
MAP



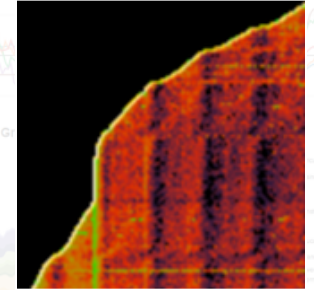
CONNECT



ORDER



TIME



Strata NYC 2014

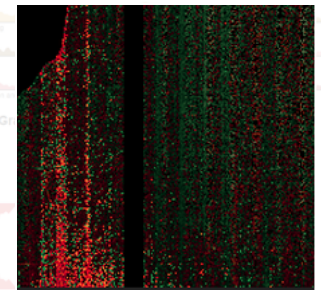
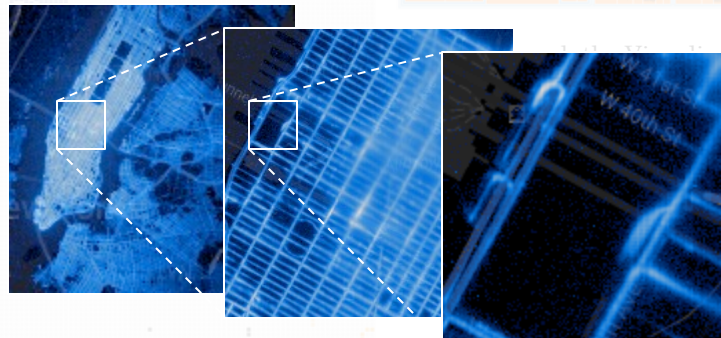
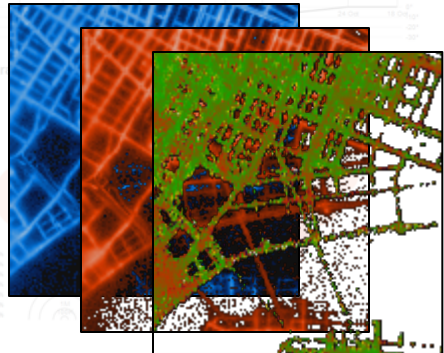
Strata London 2015

LAYERS

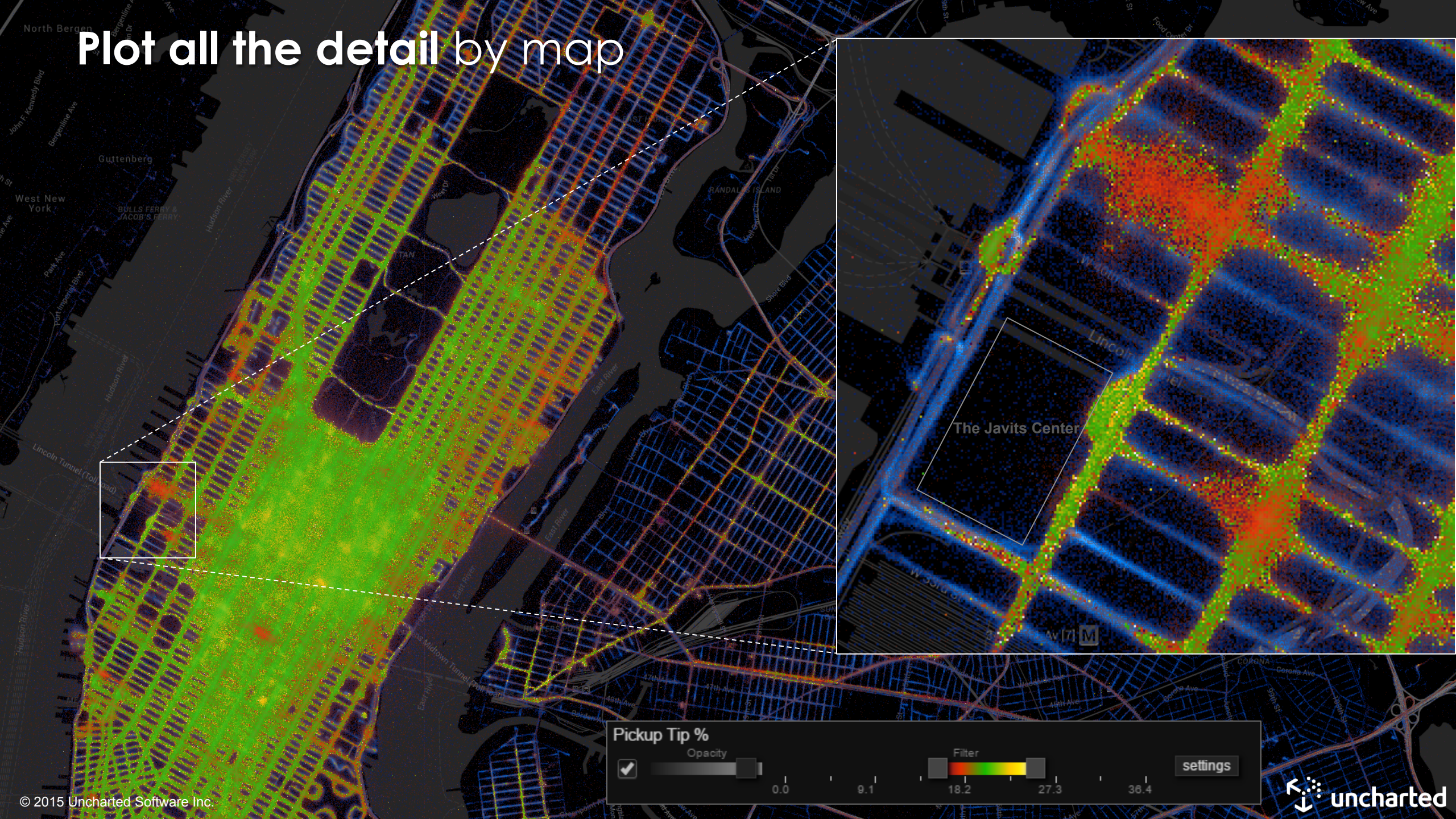
ZOOM

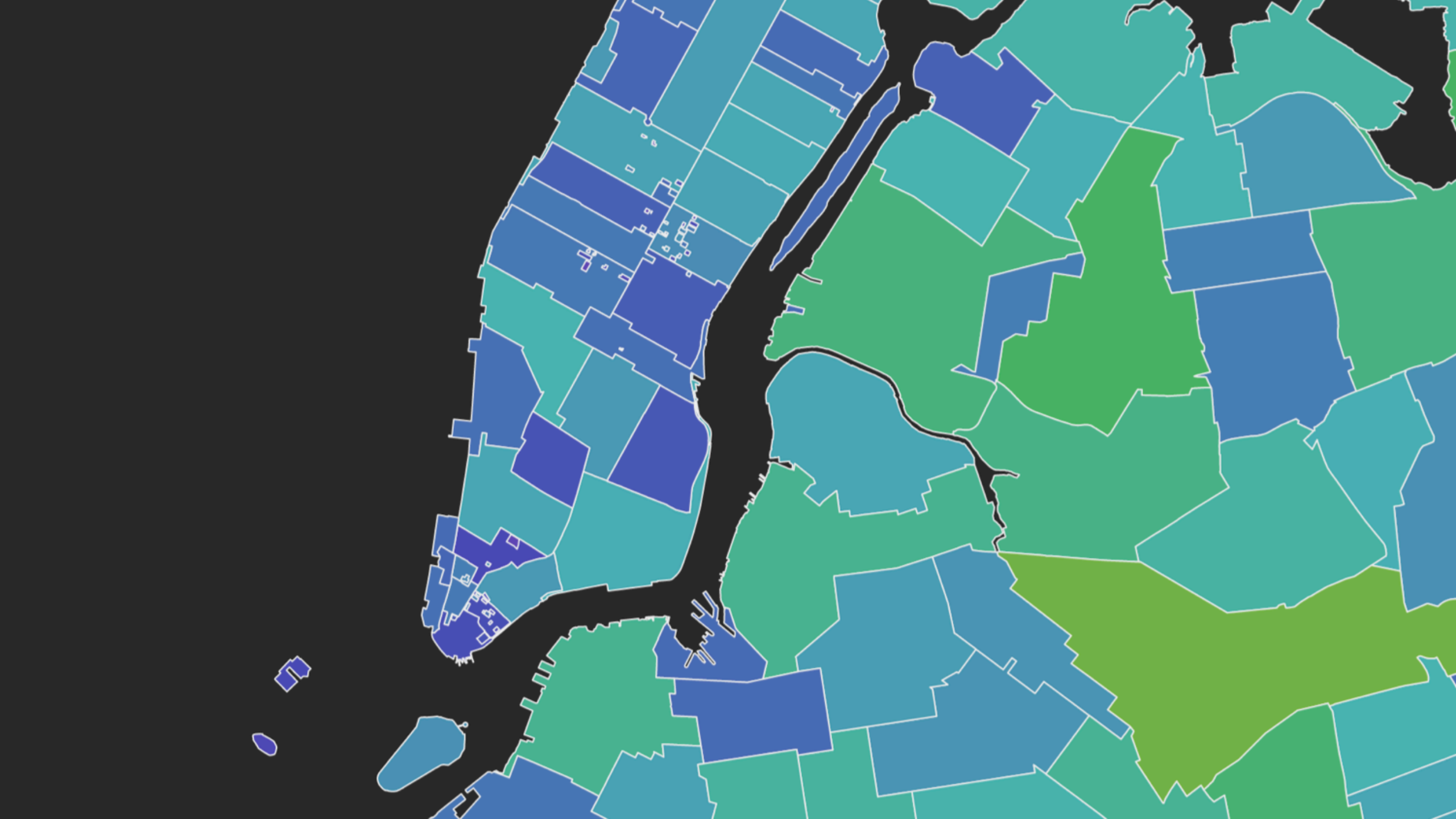
MINE

FILTER

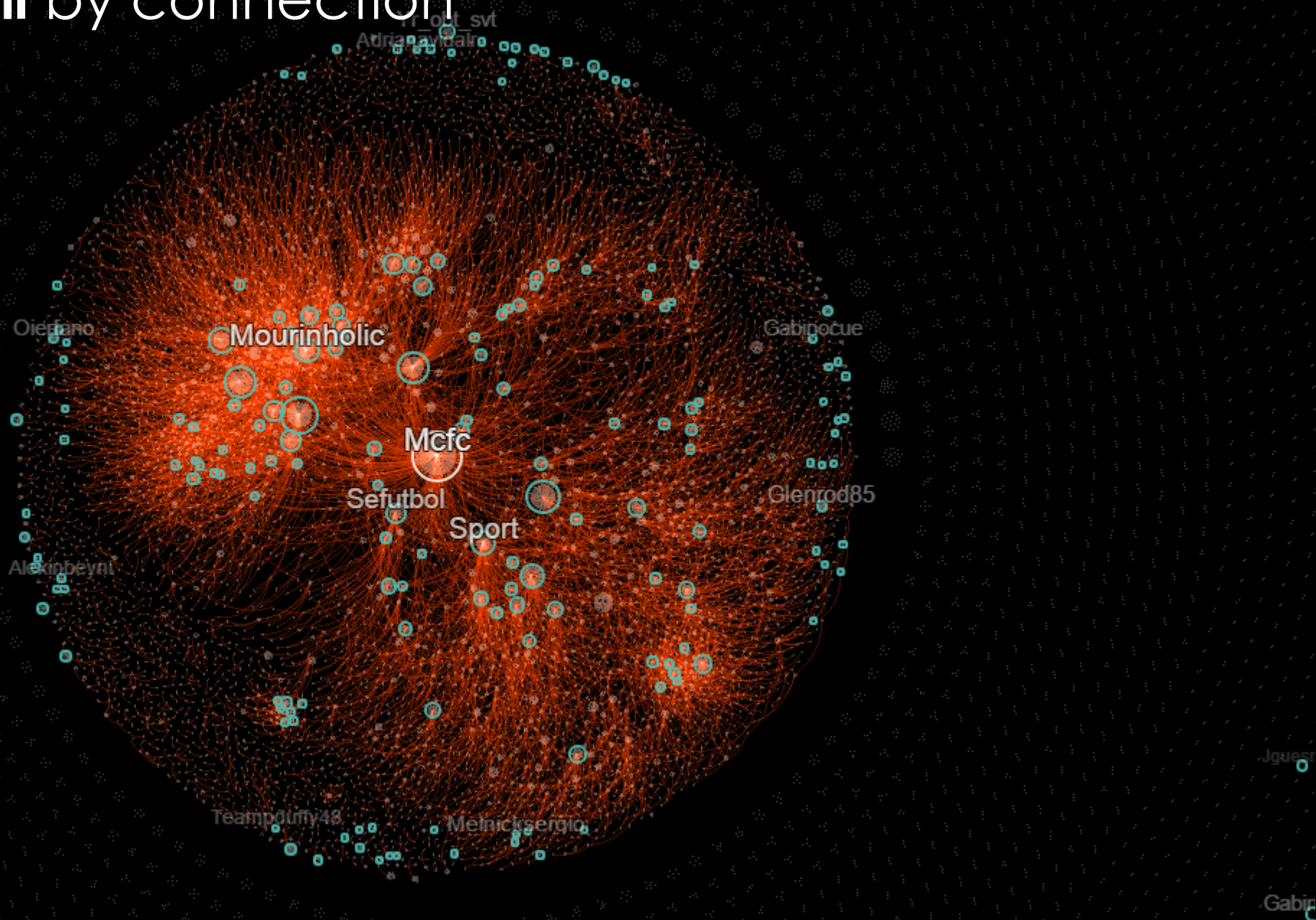


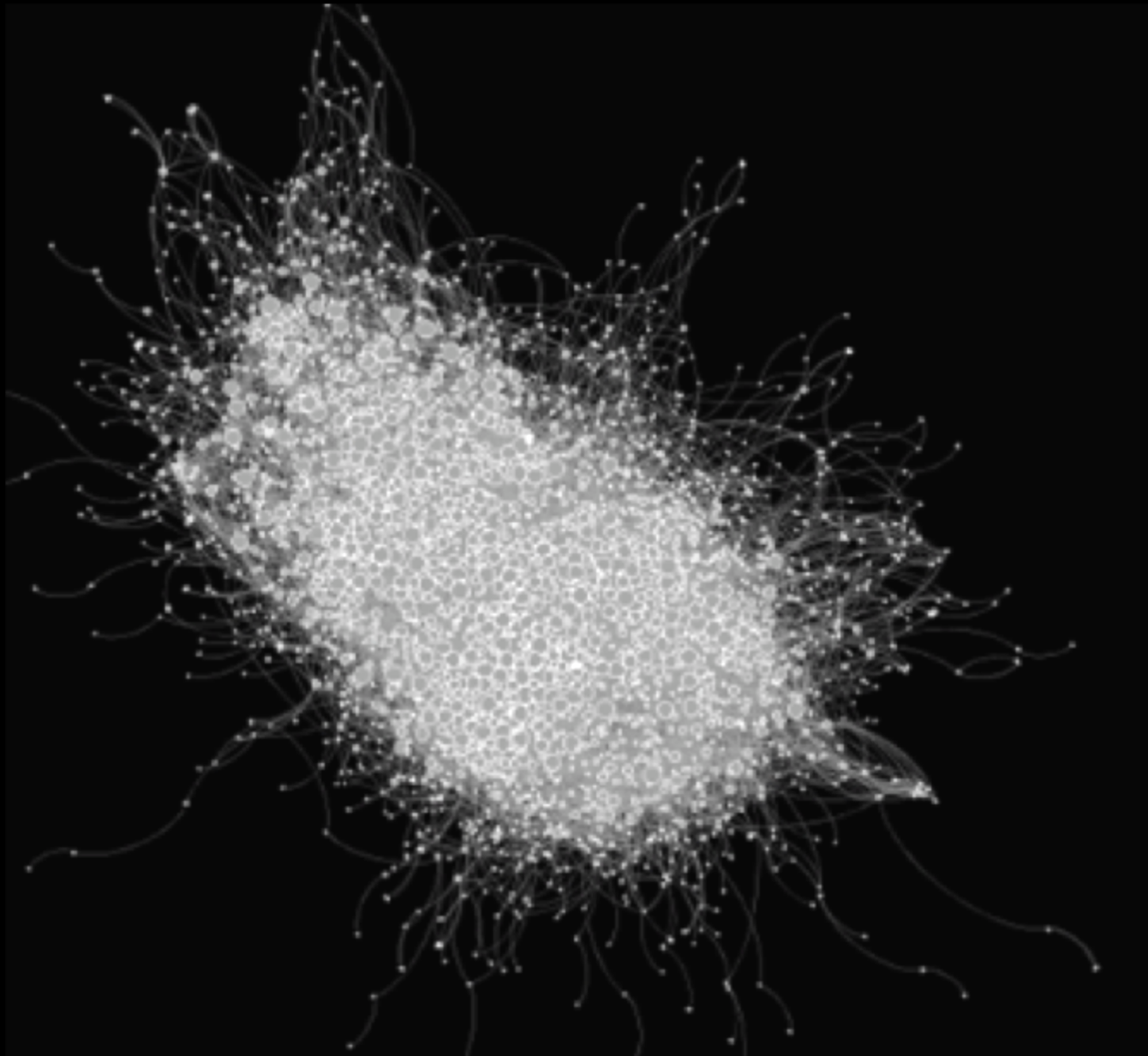
Plot all the detail by map



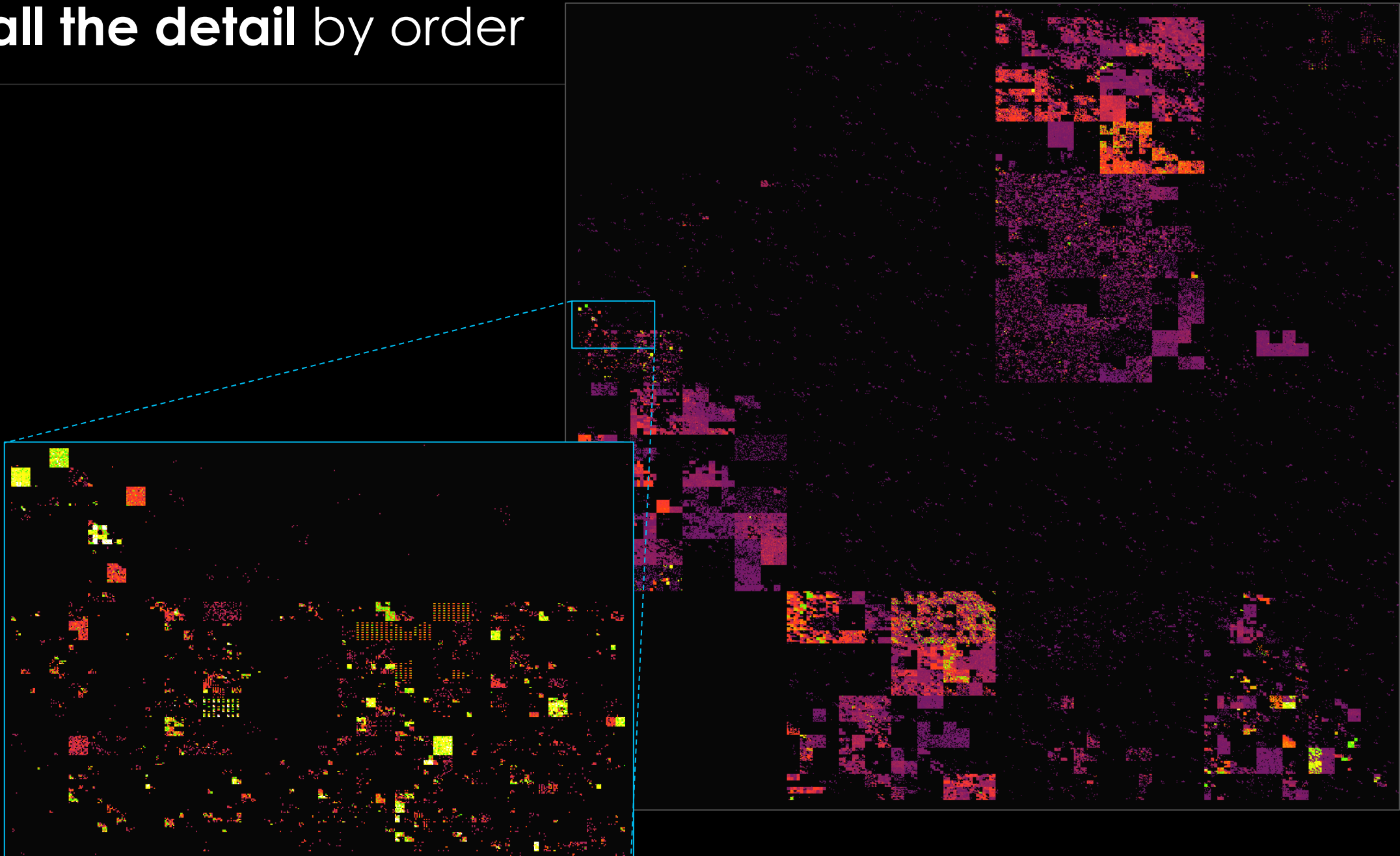


Plot all the detail by connection

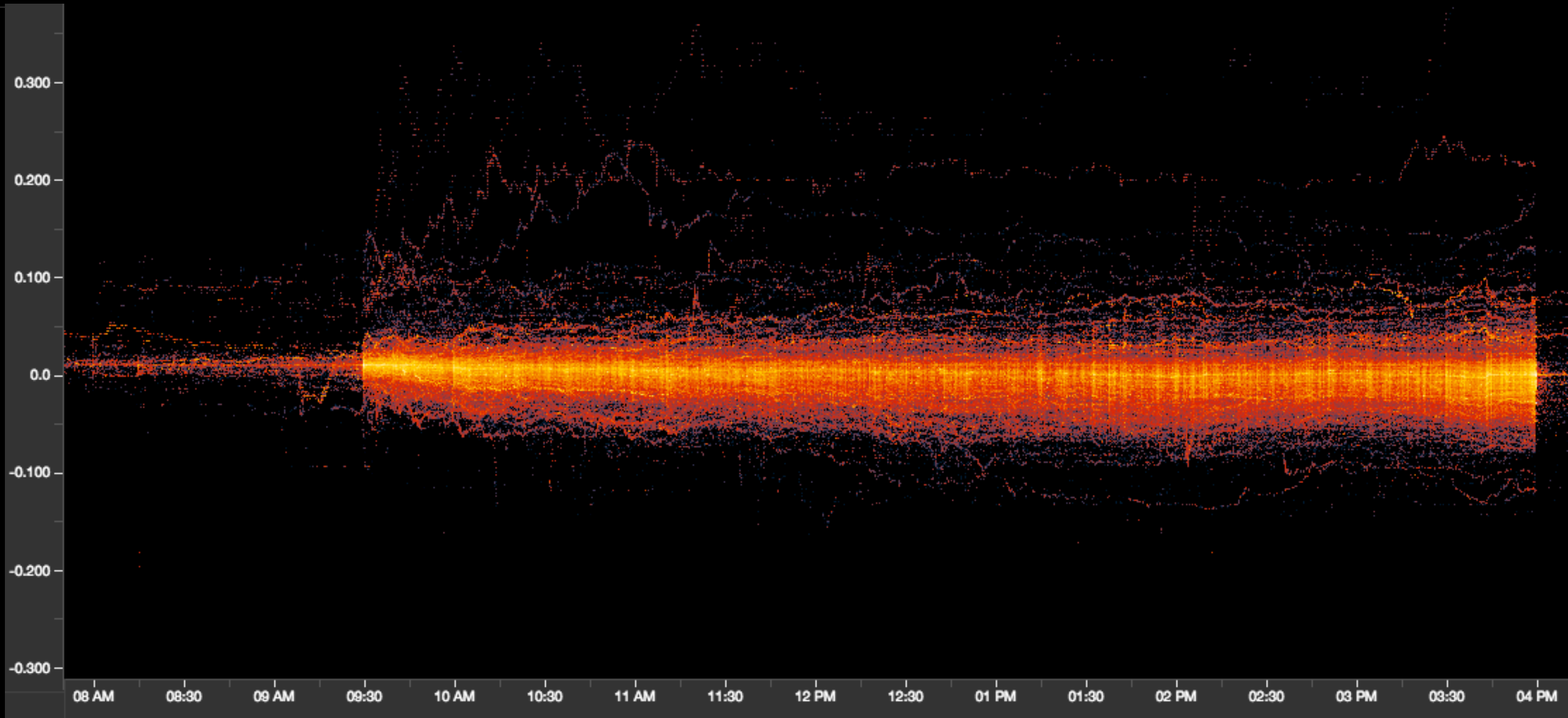




Plot all the detail by order



Plot all the detail by time



Visual Patterns in Time Plots

some other
variable

event

boundary

periodic

level (threshold)

anomaly

trend

time

Bitcoin Transactions



Load & Verify

Bitcoin Address

14sScGvSjGtxNbqFUoStoXN7eXydaN1JMM

bitcoin
Amount:

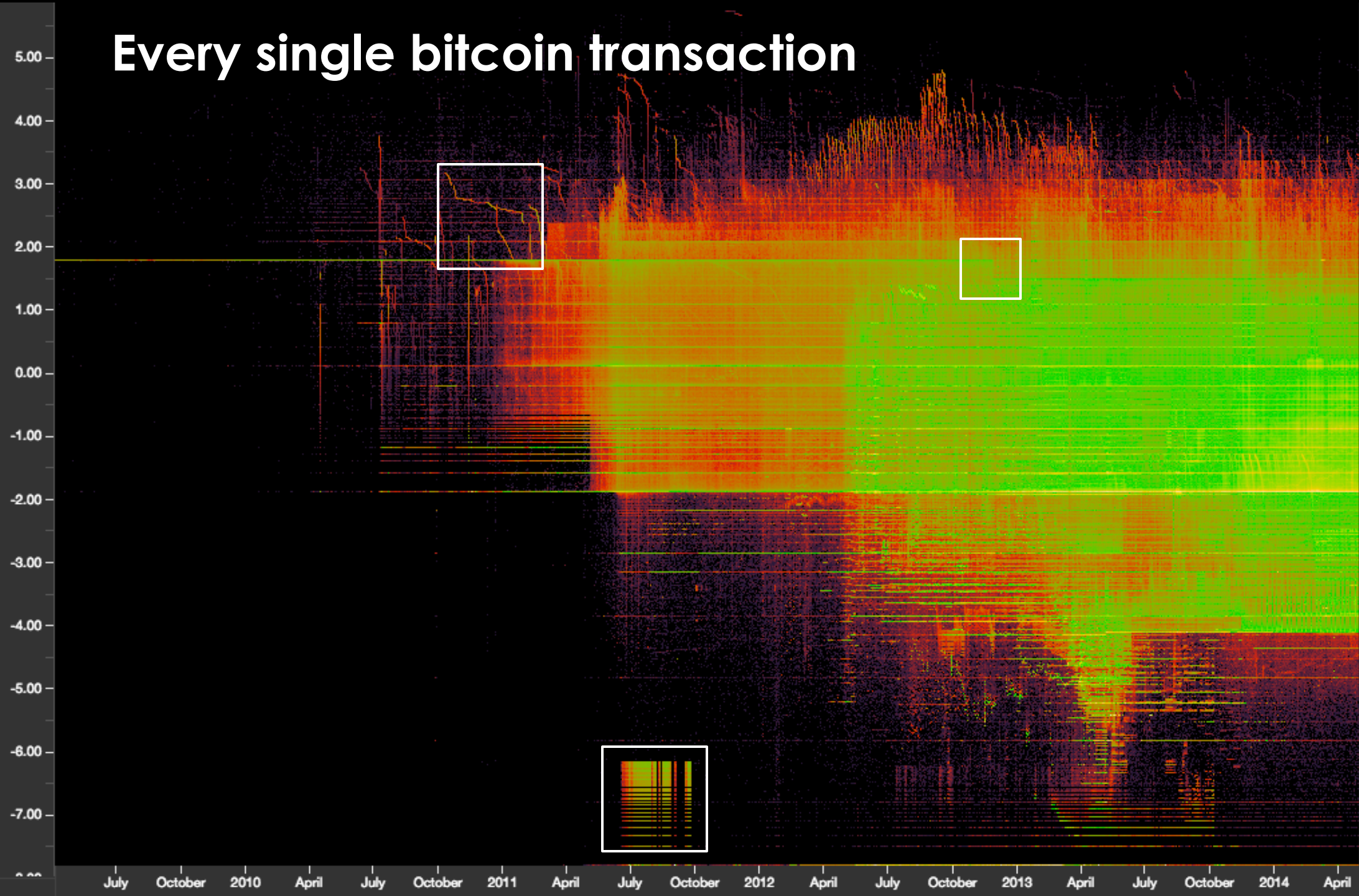
Private Key

5JGnozK YRtHhpap6DeUrHRmccTHCh8f7wcszpHGepCxS998Hfqqh



Spend

Every single bitcoin transaction

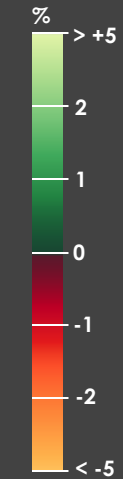


Financial Markets

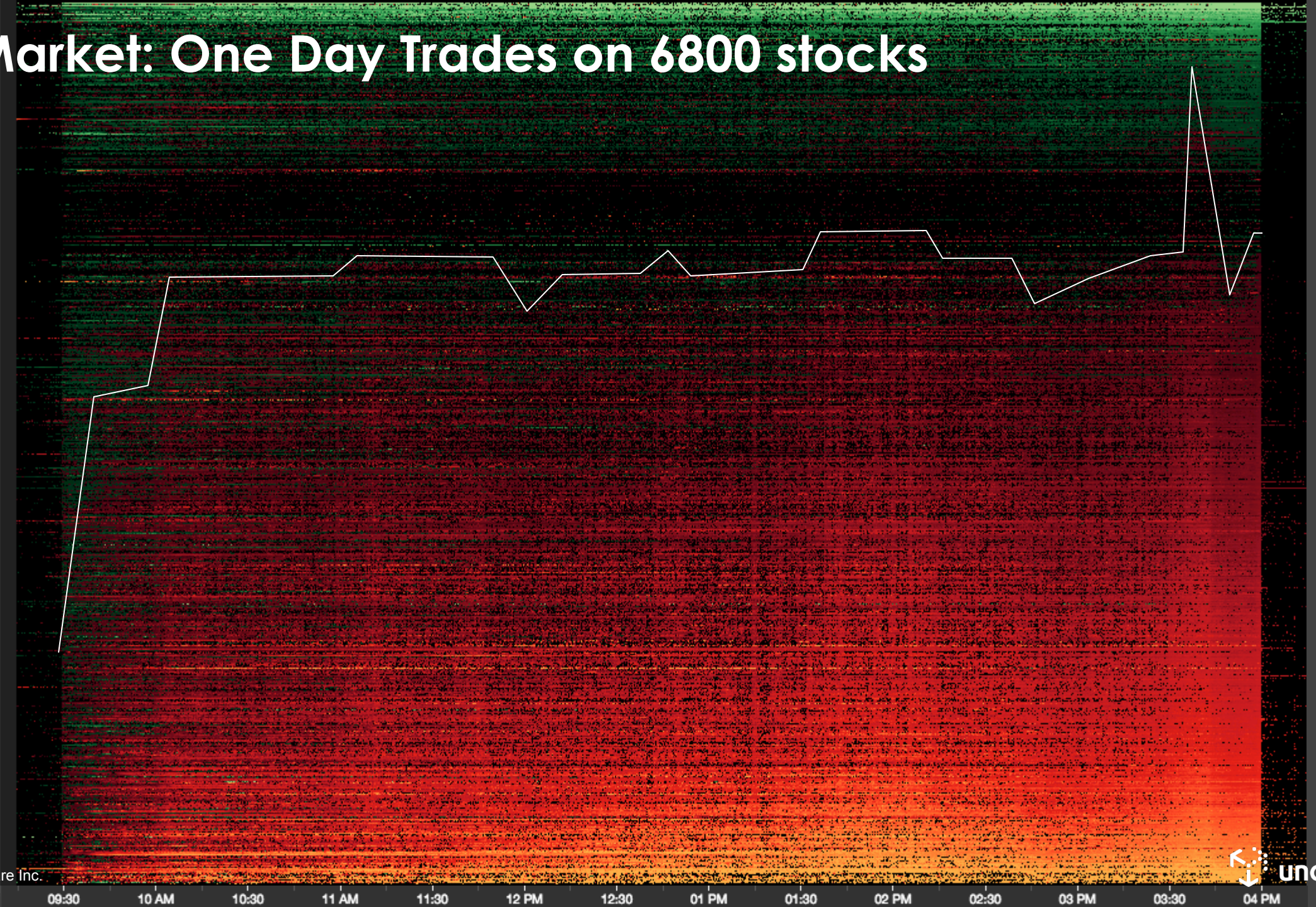


Stock Market: One Day Trades on 6800 stocks

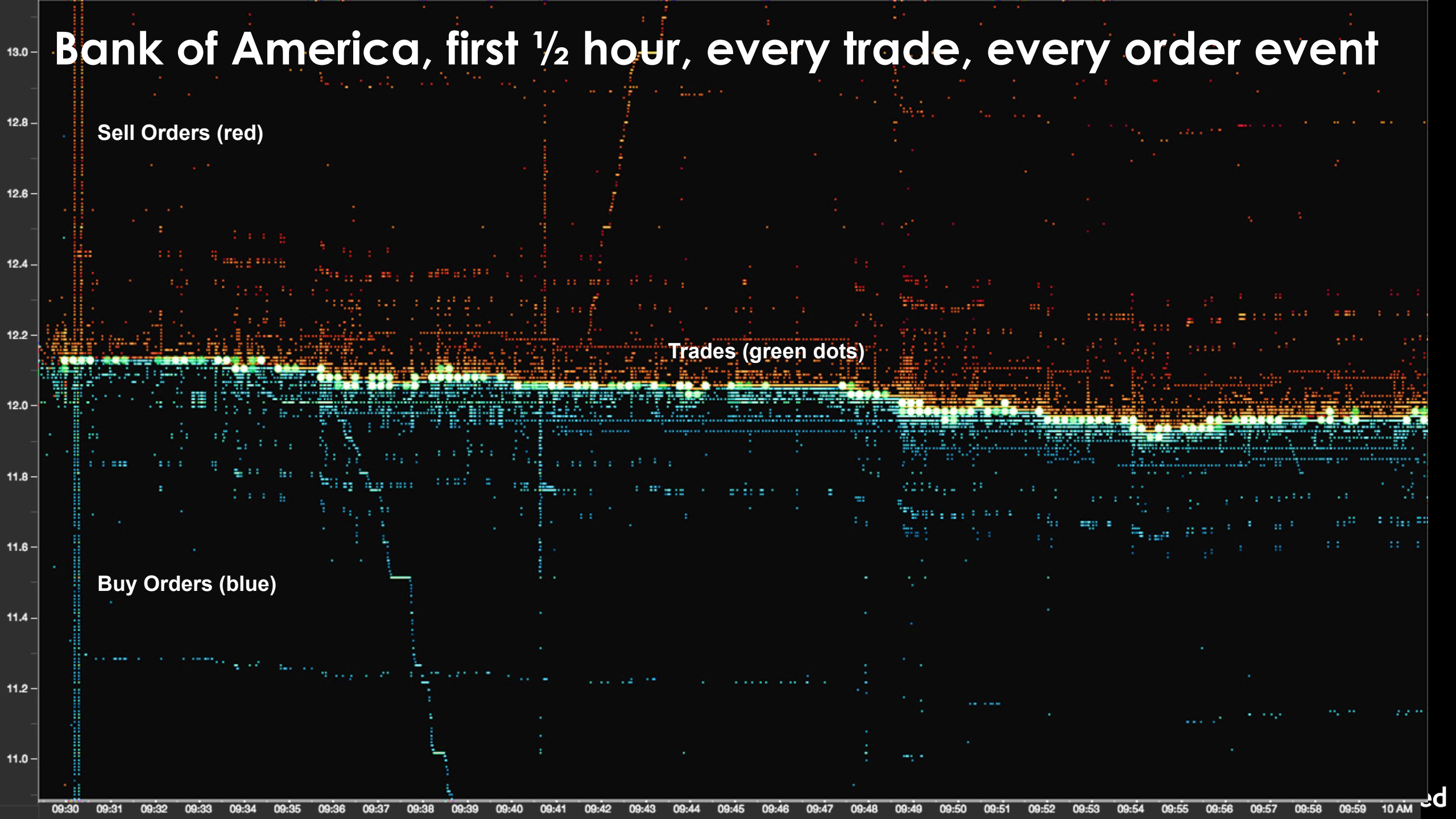
ARCA stocks
ranked by
percent
change on
day



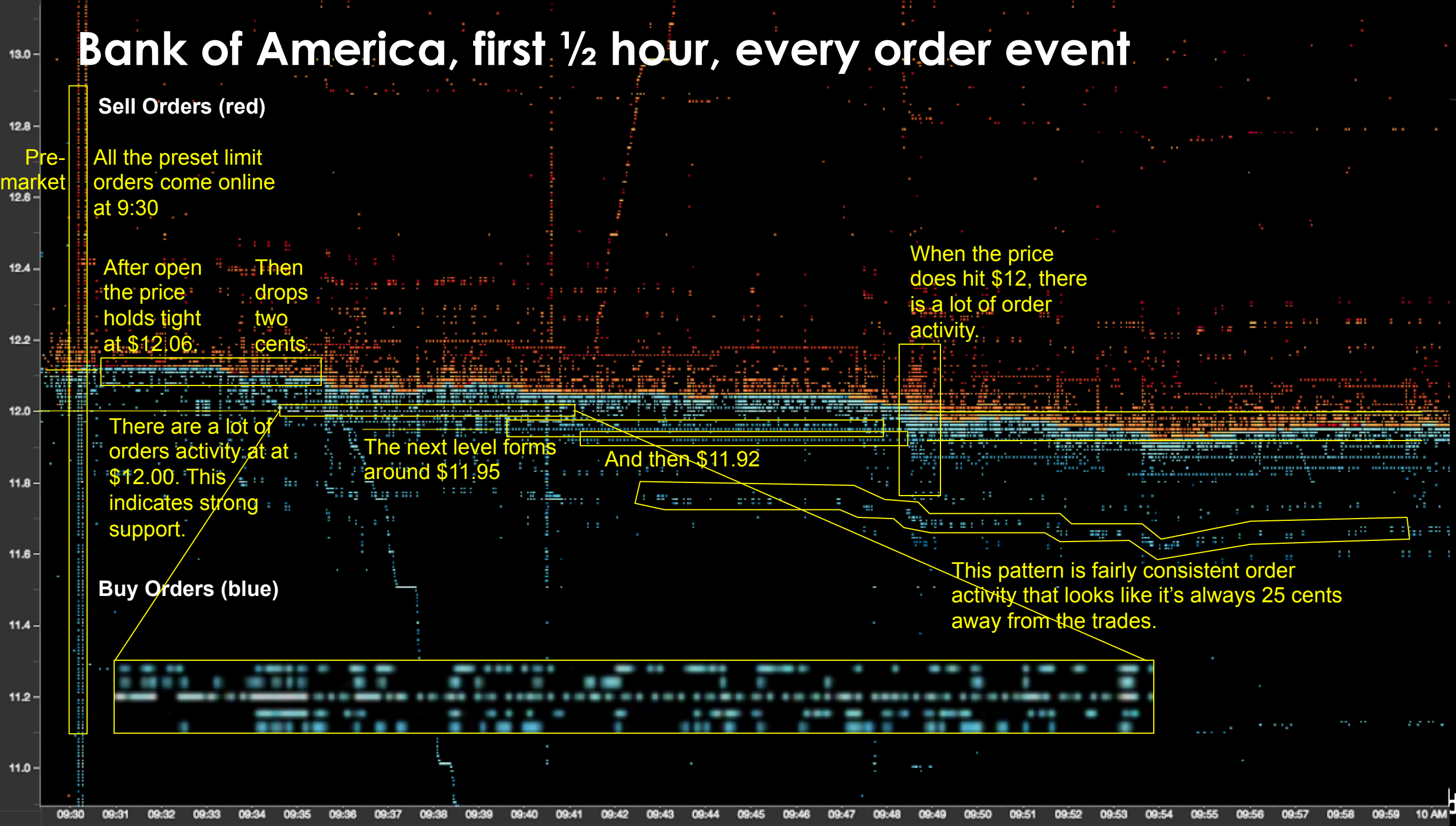
6.5k
6.00k
5.50k
5.00k
4.50k
4.00k
3.50k
3.00k
2.50k
2.00k
1.50k
1.00k
500



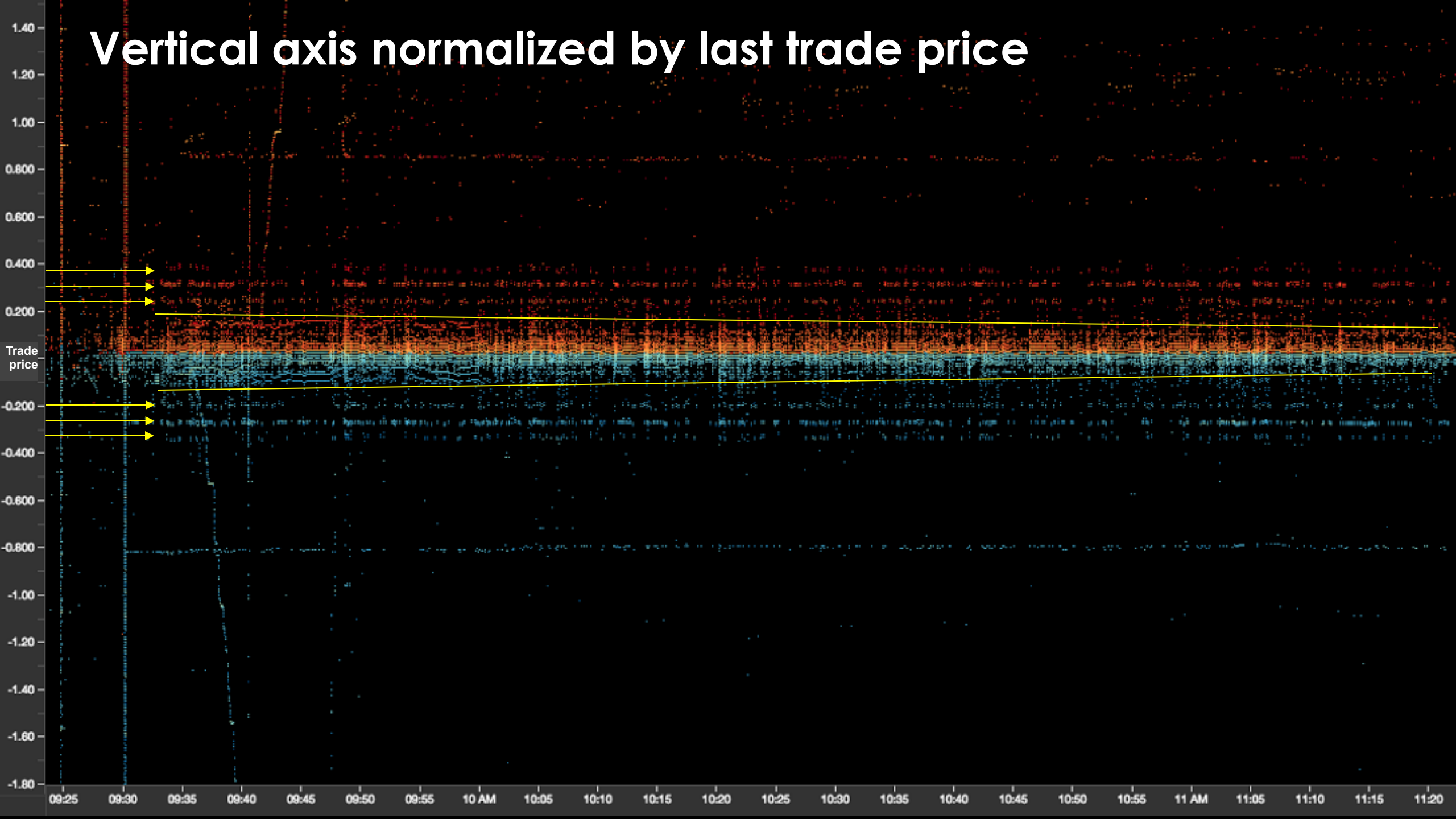
Bank of America, first ½ hour, every trade, every order event



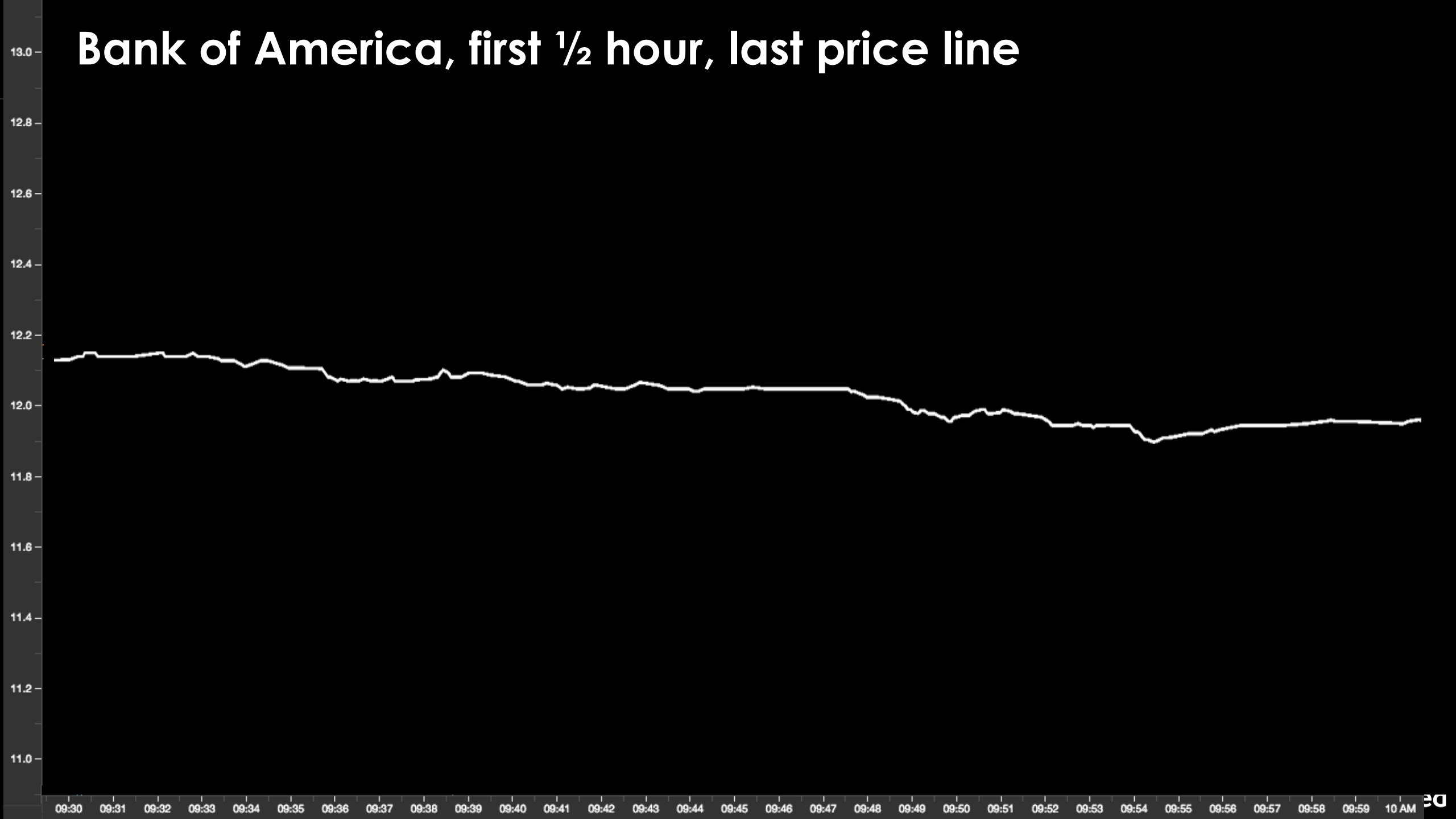
Bank of America, first ½ hour, every order event



Vertical axis normalized by last trade price



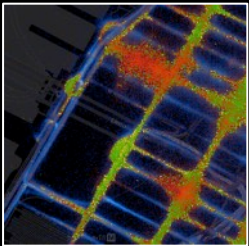
Bank of America, first 1/2 hour, last price line



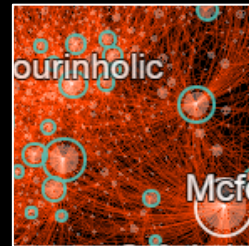
Technical Approach

Exploratory Big Data Analysis

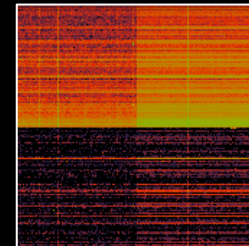
MAP



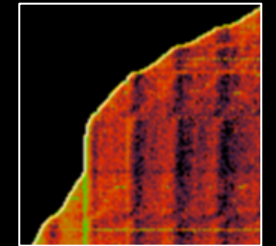
CONNECT

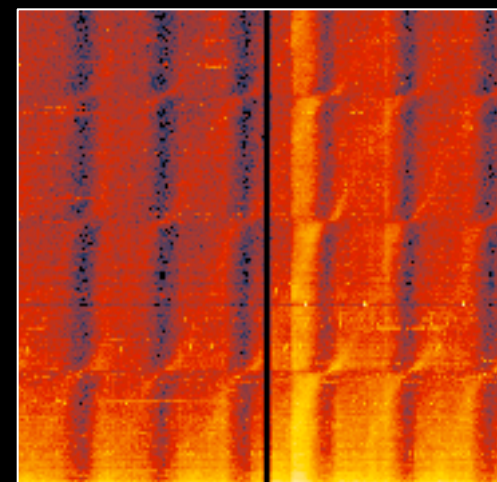
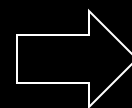
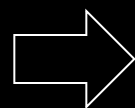
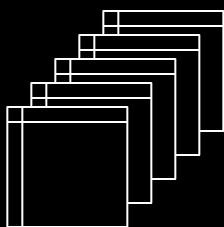
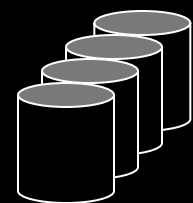
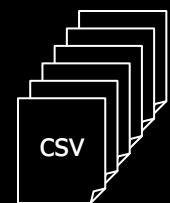


ORDER



TIME

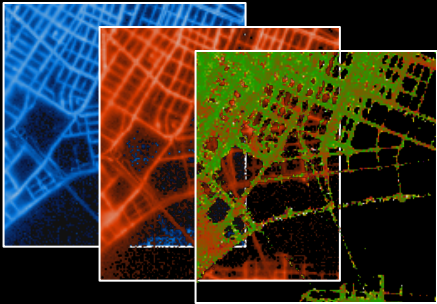




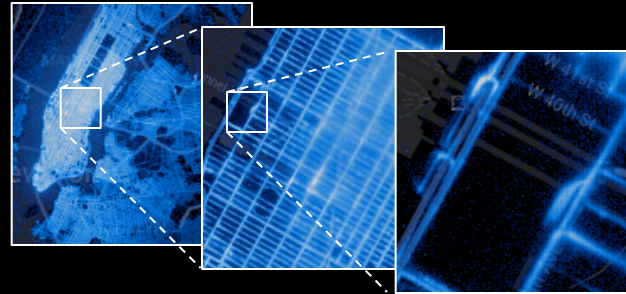
result.png

Exploratory Big Data Analysis – Rich Interactions

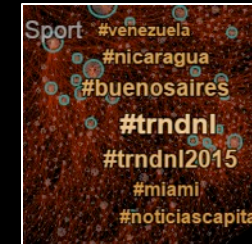
LAYERS



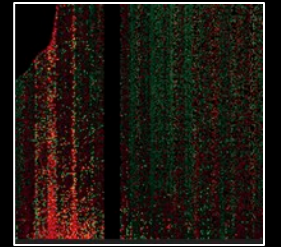
ZOOM



MINE



FILTER



Spark as a Service

Expressive API

Technical Approach

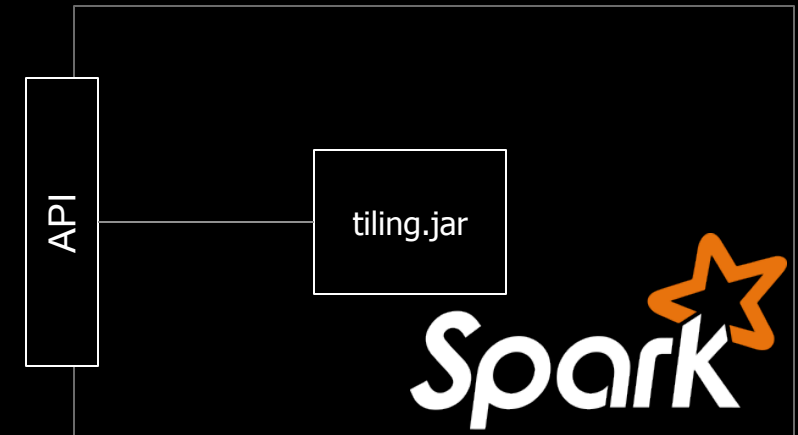
Data Pipeline

- Loading
- Filtering
- Transformation
- Sentiment
- Serialization

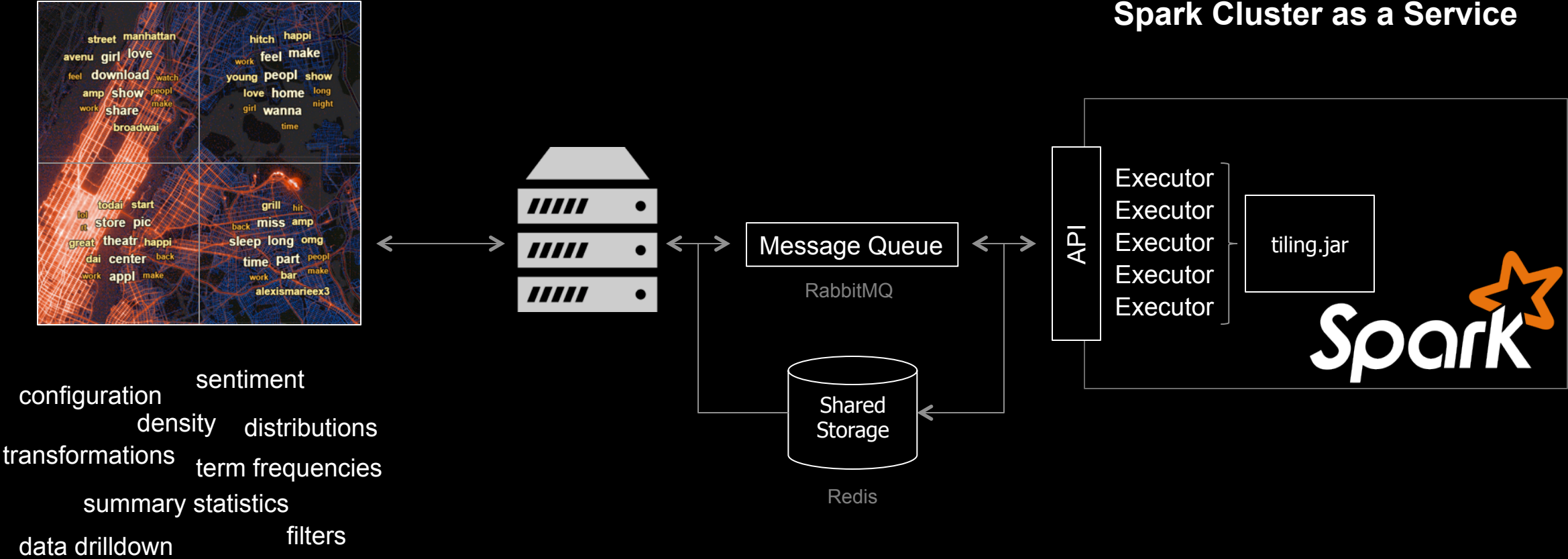
Results Generation

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

Spark Cluster as a Service



Technical Approach



Technical Approach



github.com/unchartedsoftware

Trump Tweets



Donald J. Trump ✓
@realDonaldTrump

Follow

[#AskTrump](#) Getting ready to answer your questions.

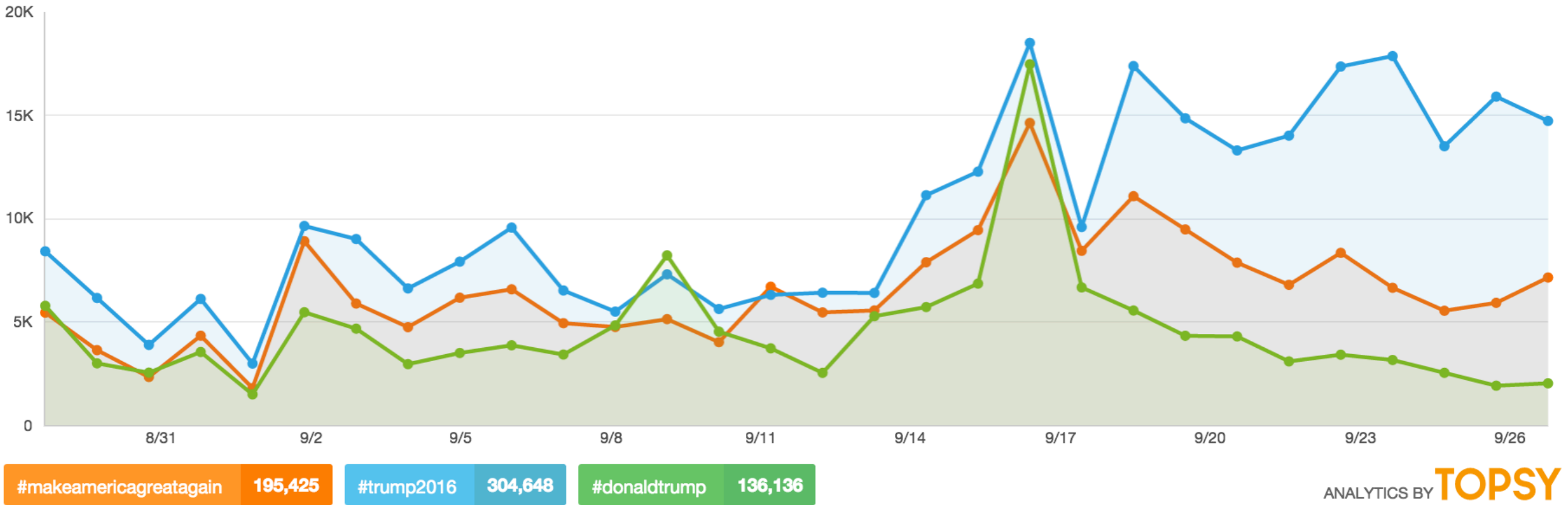
1:07 PM - 21 Sep 2015

1,037 2,429



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

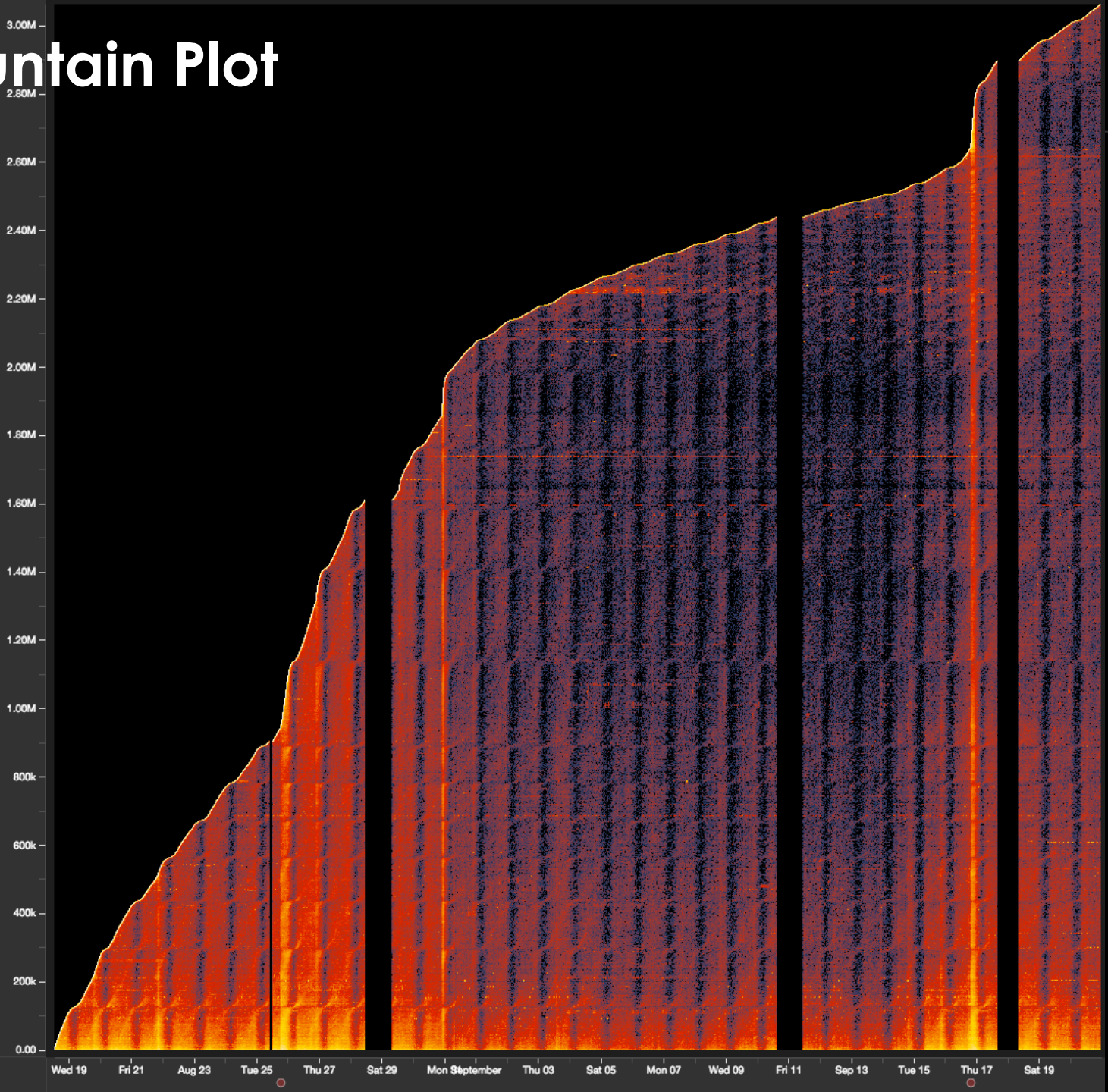
August 28th — September 27th



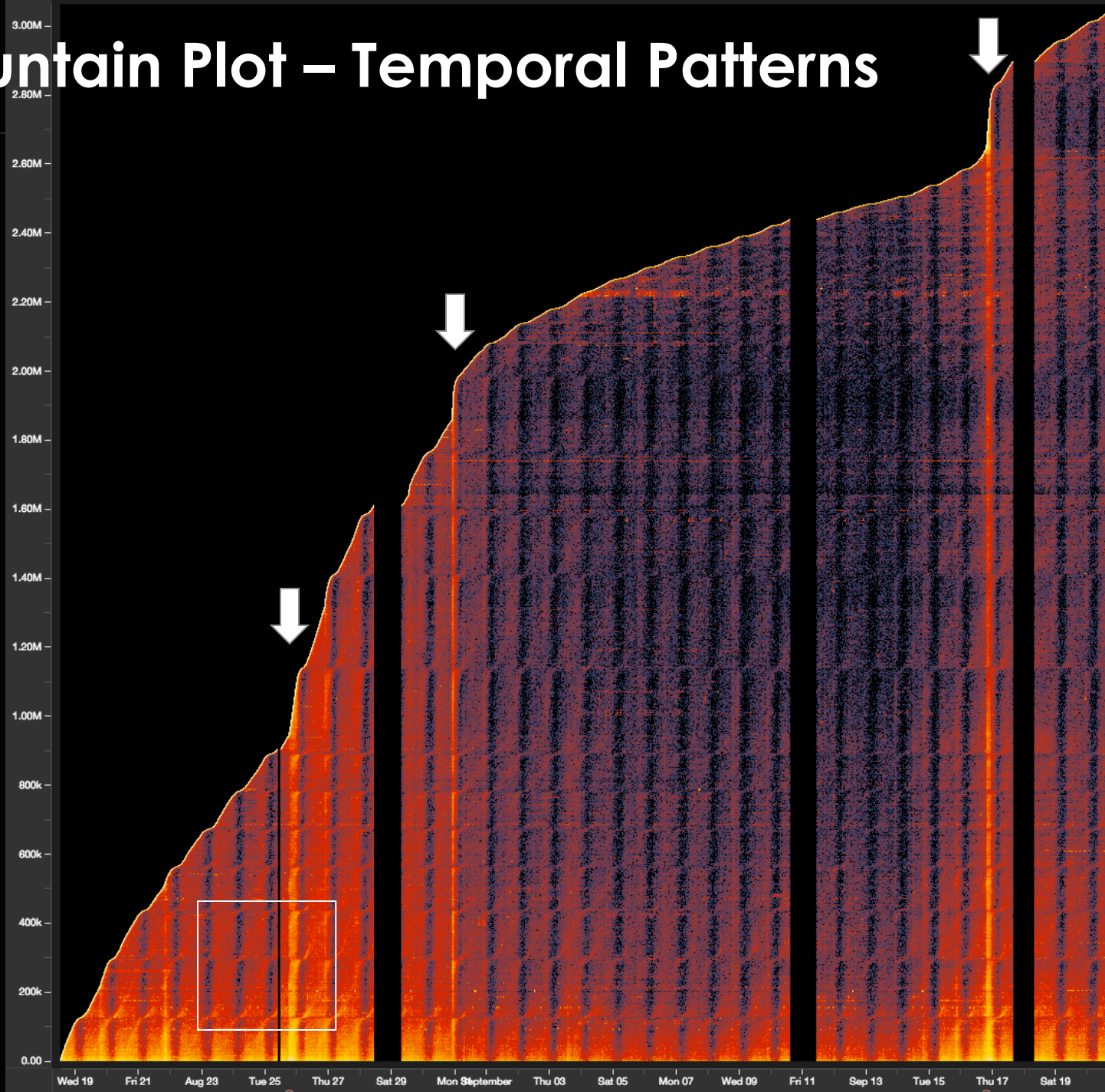
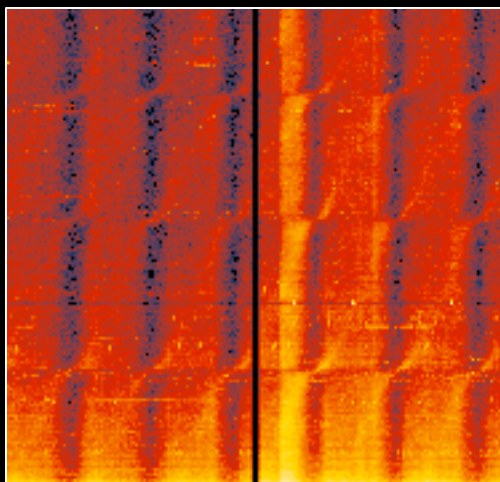
ANALYTICS BY **TOPSY**

[Continue](#)

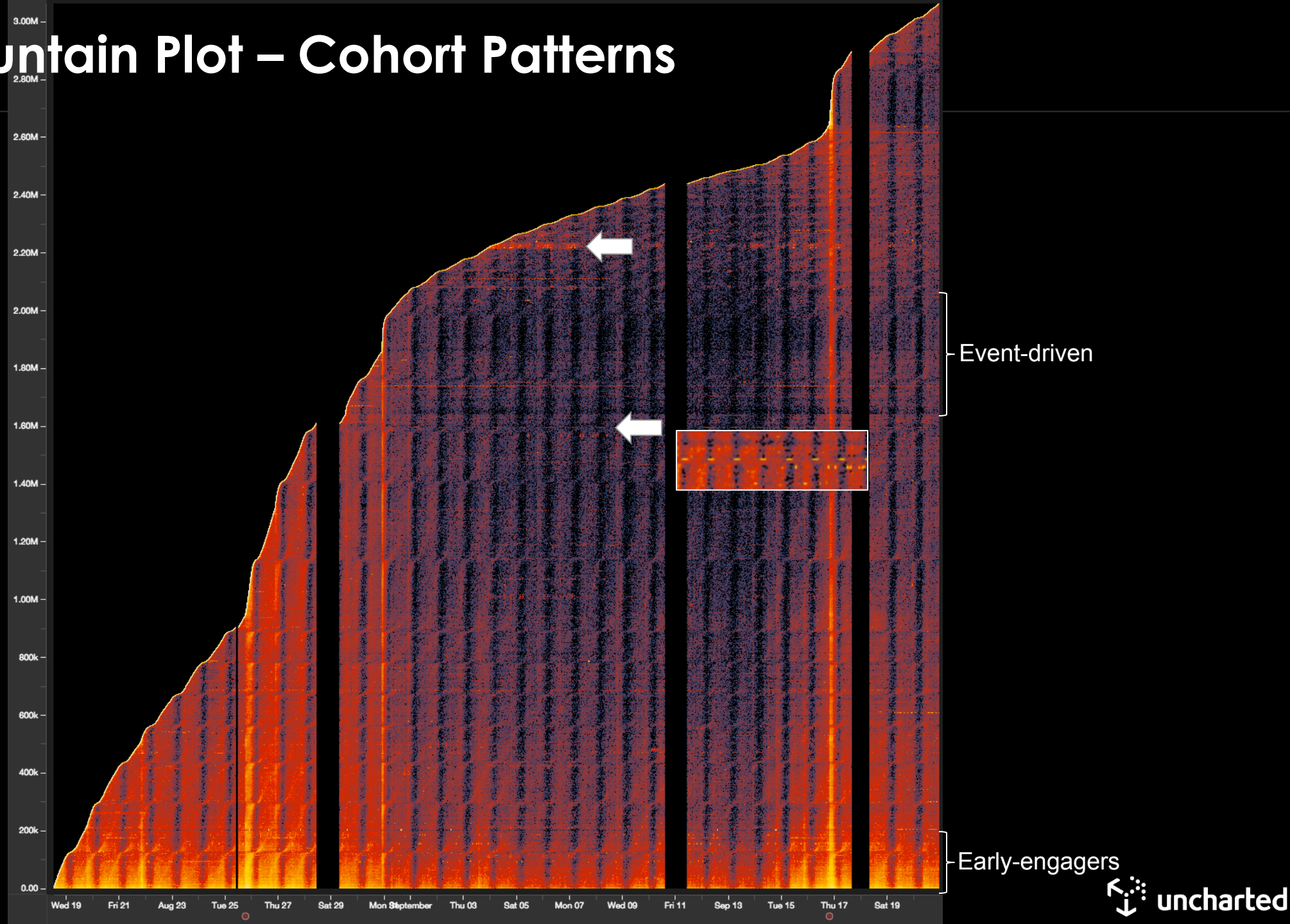
Trump Mountain Plot



Trump Mountain Plot – Temporal Patterns

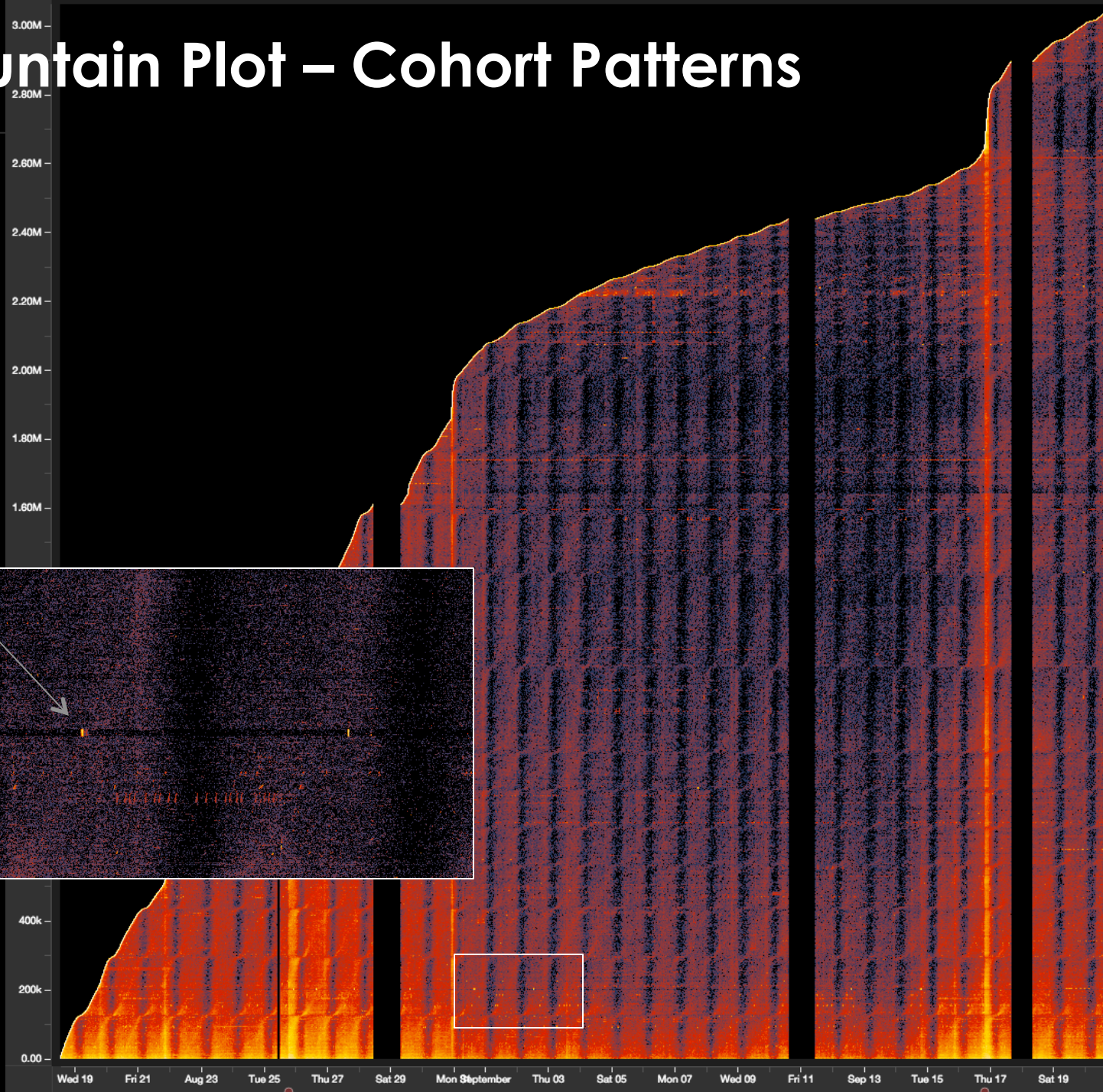
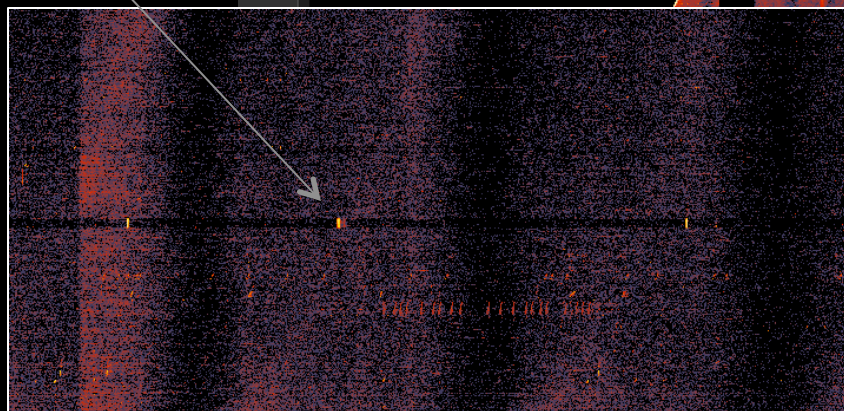


Trump Mountain Plot – Cohort Patterns

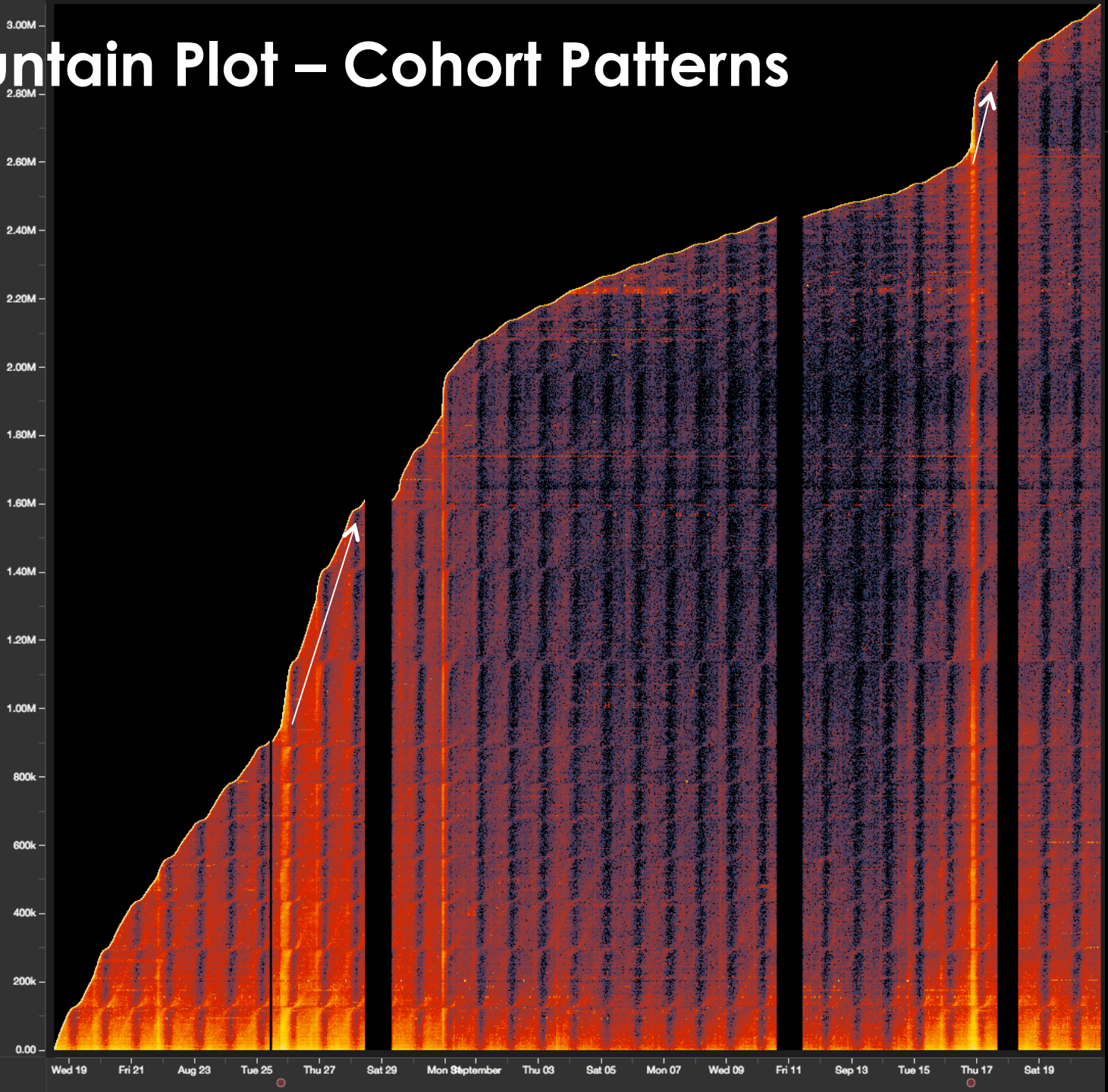


Trump Mountain Plot – Cohort Patterns

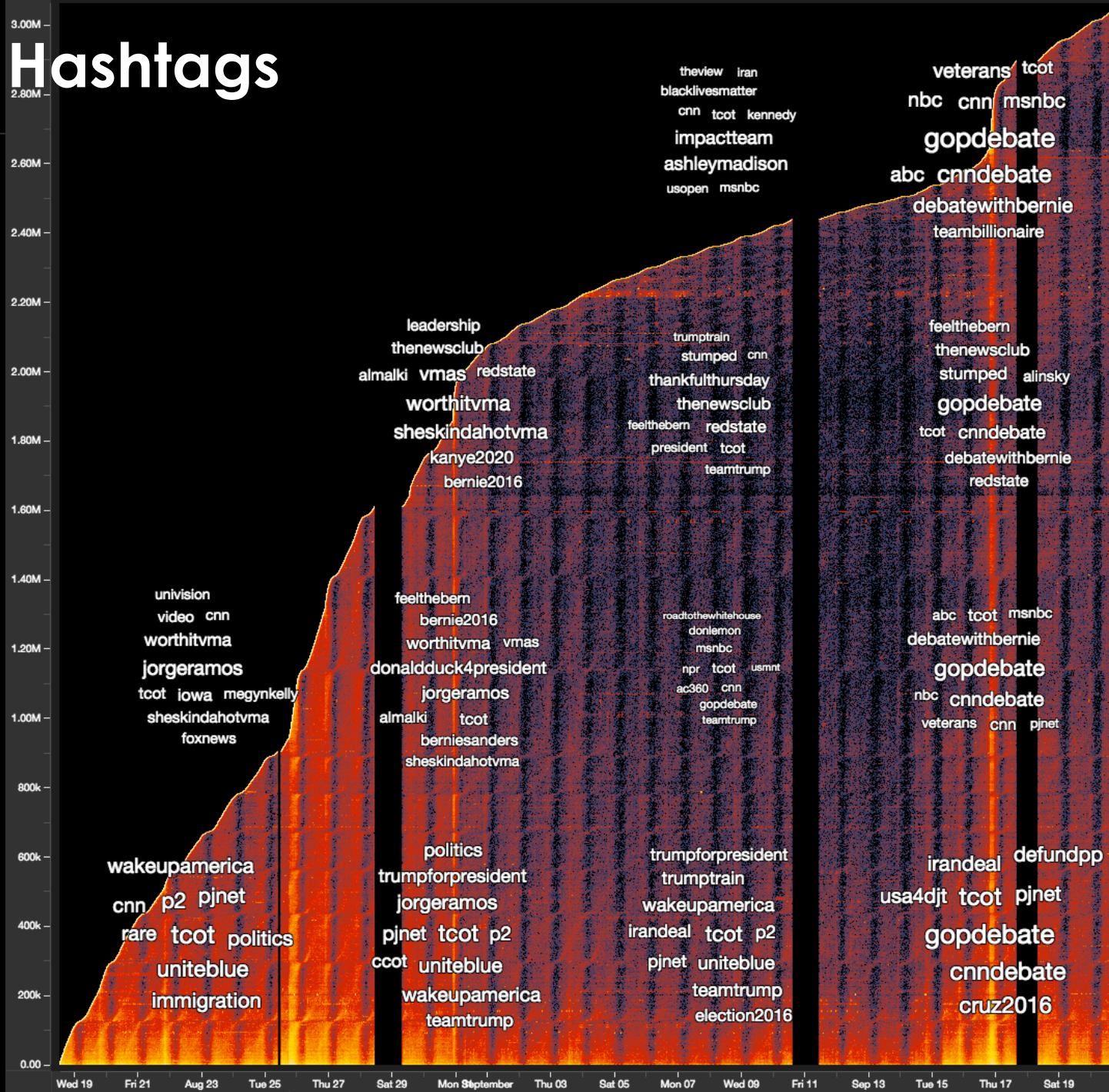
4000
coordinated
accounts



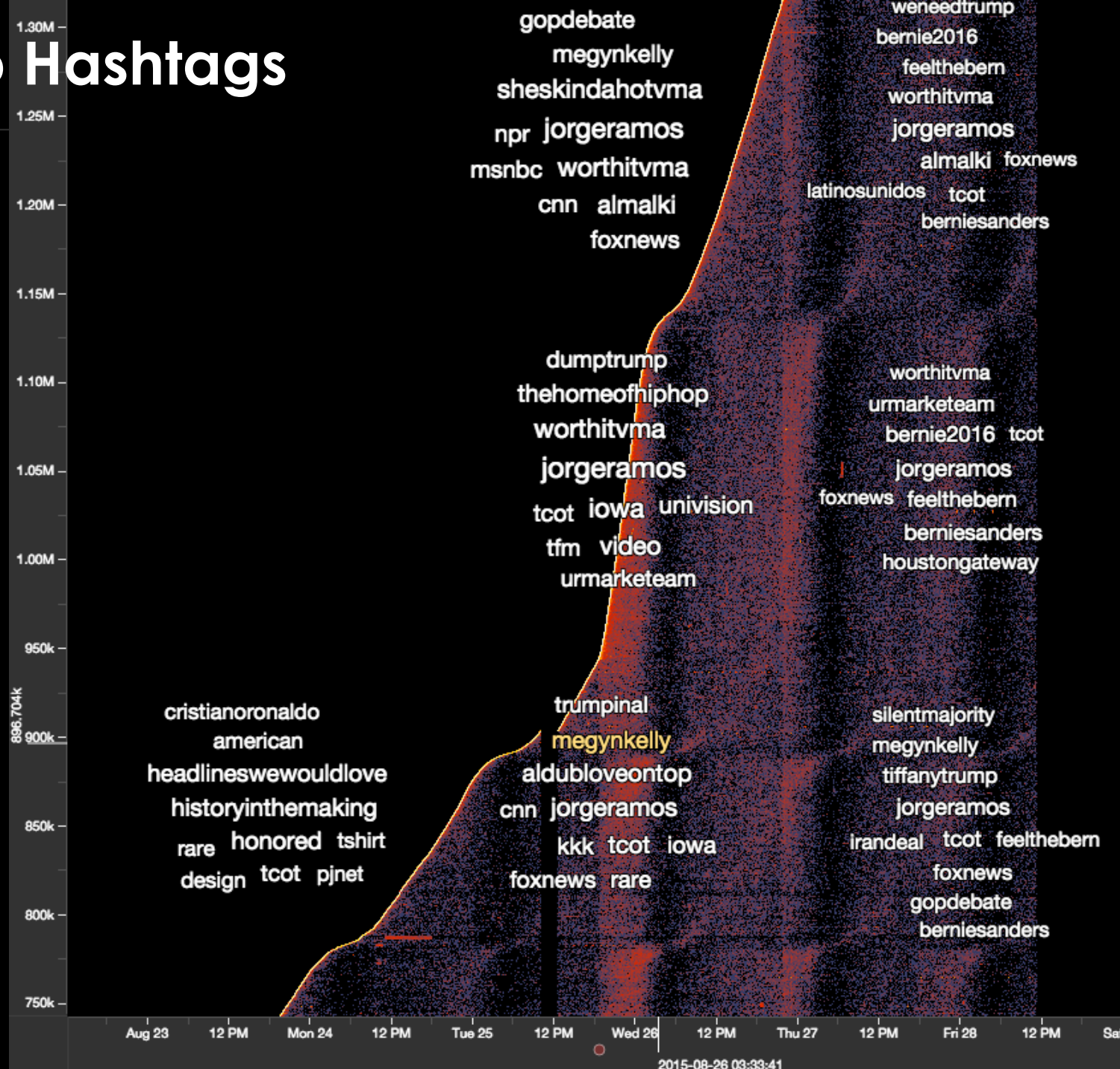
Trump Mountain Plot – Cohort Patterns



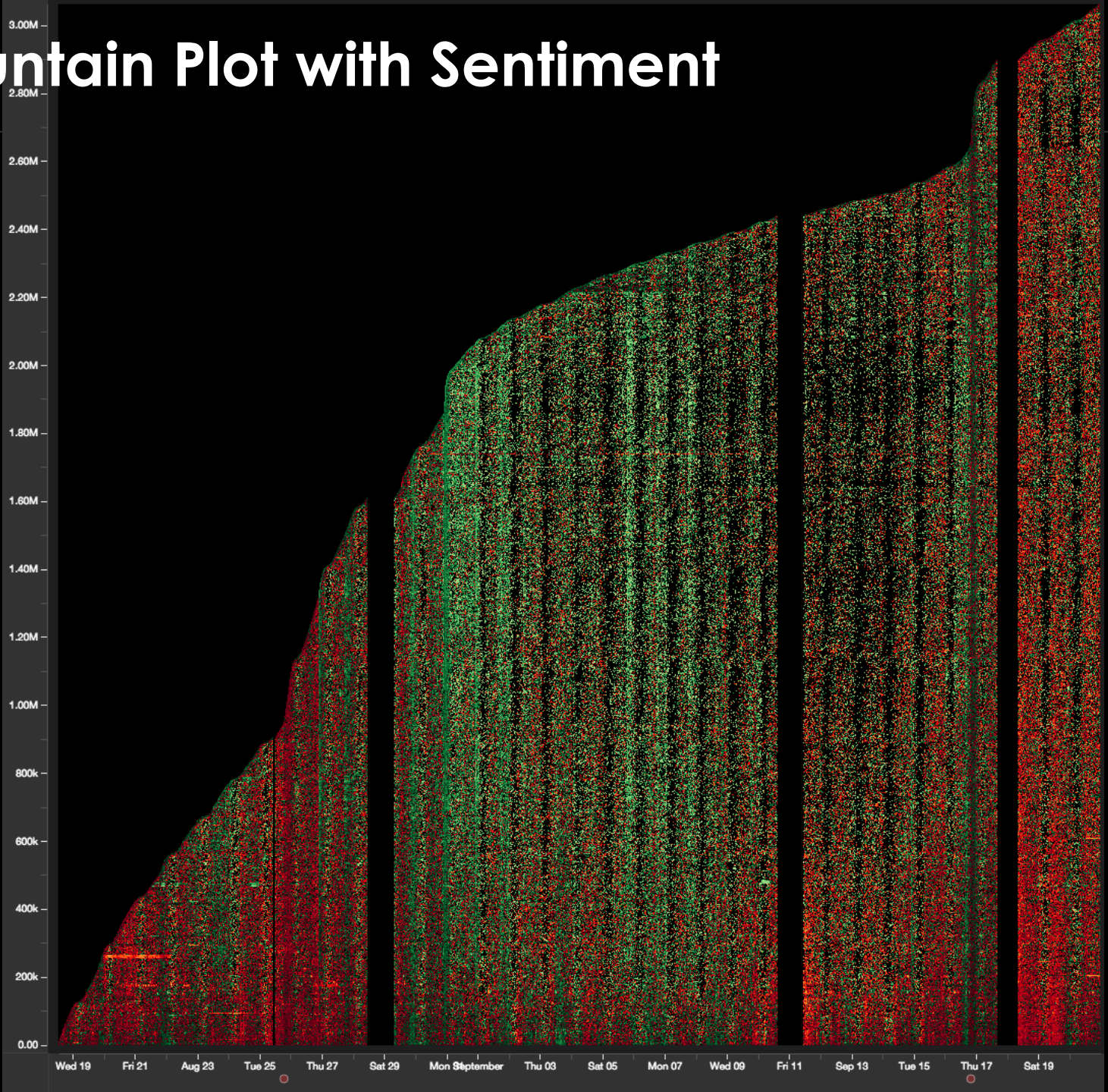
Trump Top Hashtags



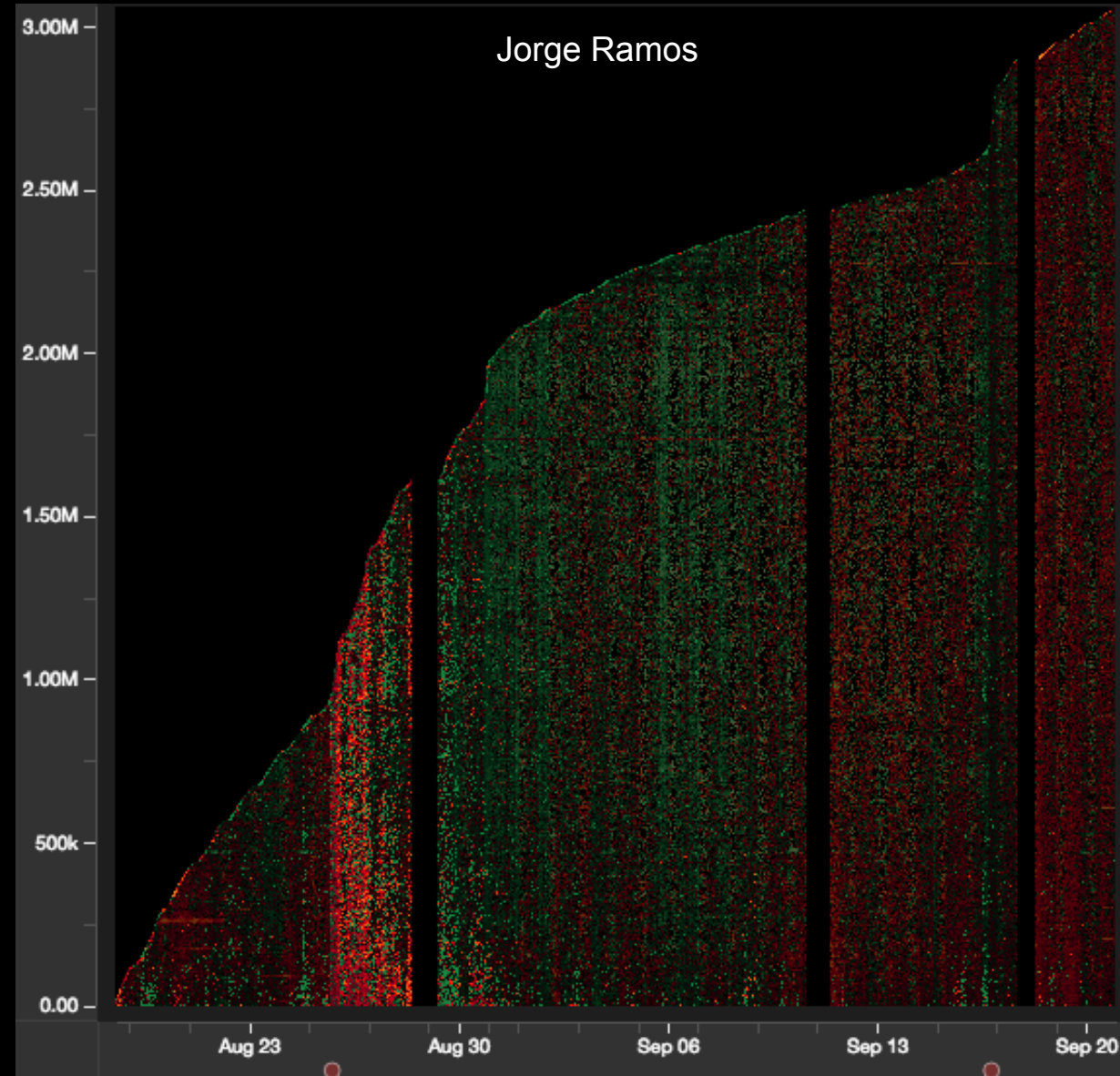
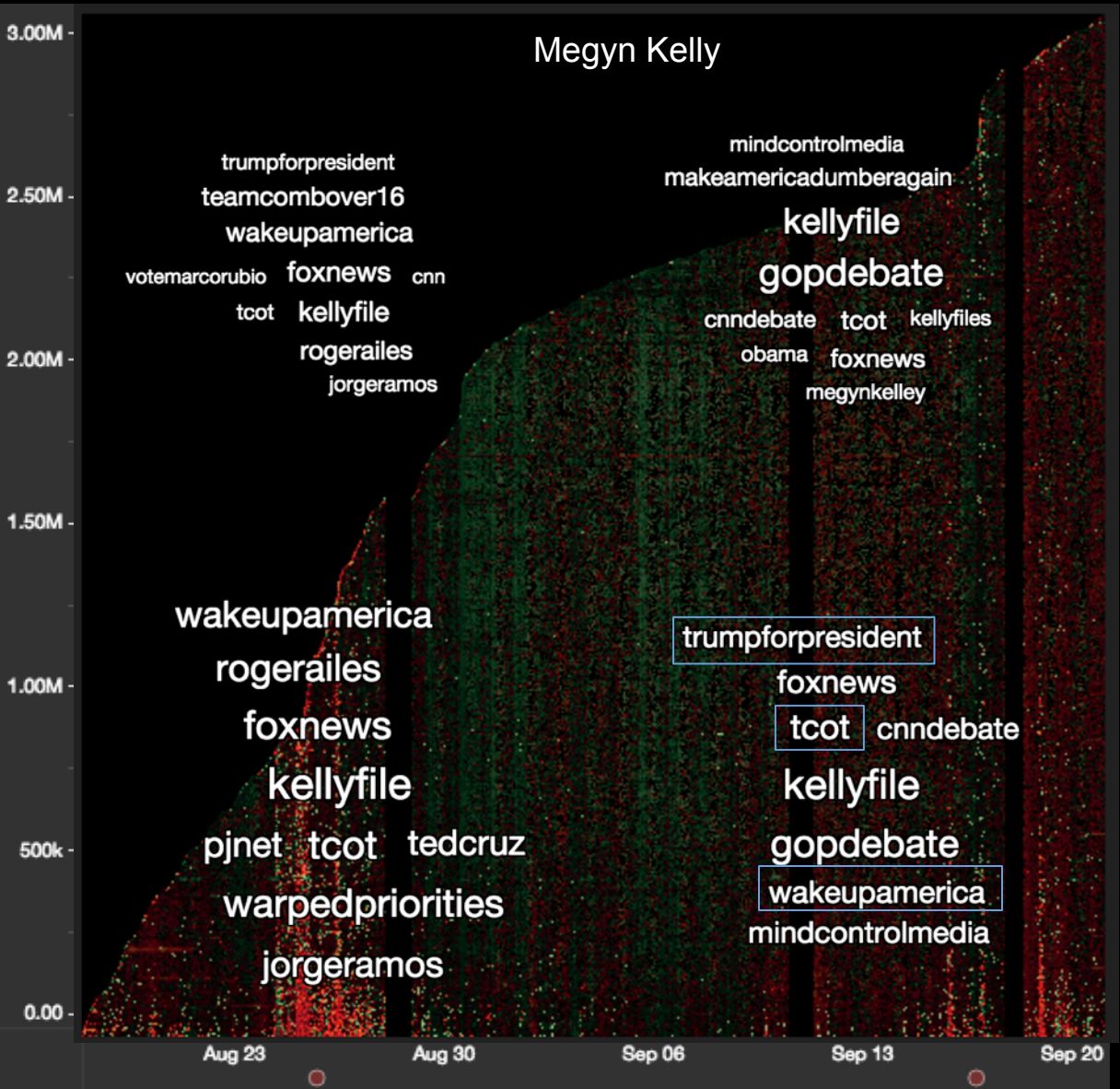
Trump Top Hashtags



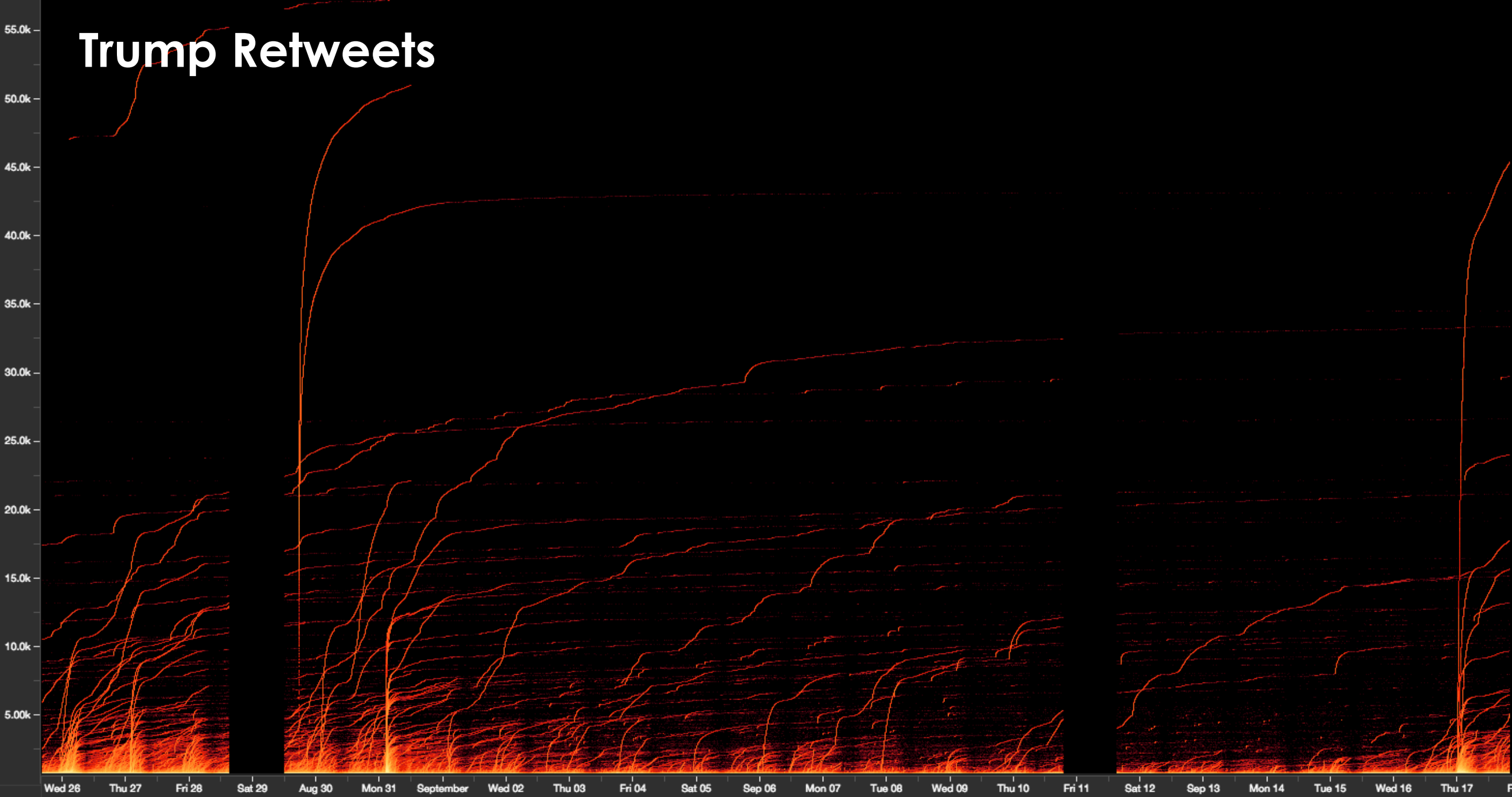
Trump Mountain Plot with Sentiment



Trump Mountain Plot with Sentiment – Megyn / Jorge



Trump Retweets



Trump Retweets



@Ashton5SOS

Not sure about this whole
Donald trump thing



@jasonmustian

I think Trump just did all
the emoji faces in 7
seconds



@ruinedpicnic

Donald trump looks like the
villain in a movie where the
hero is a dog

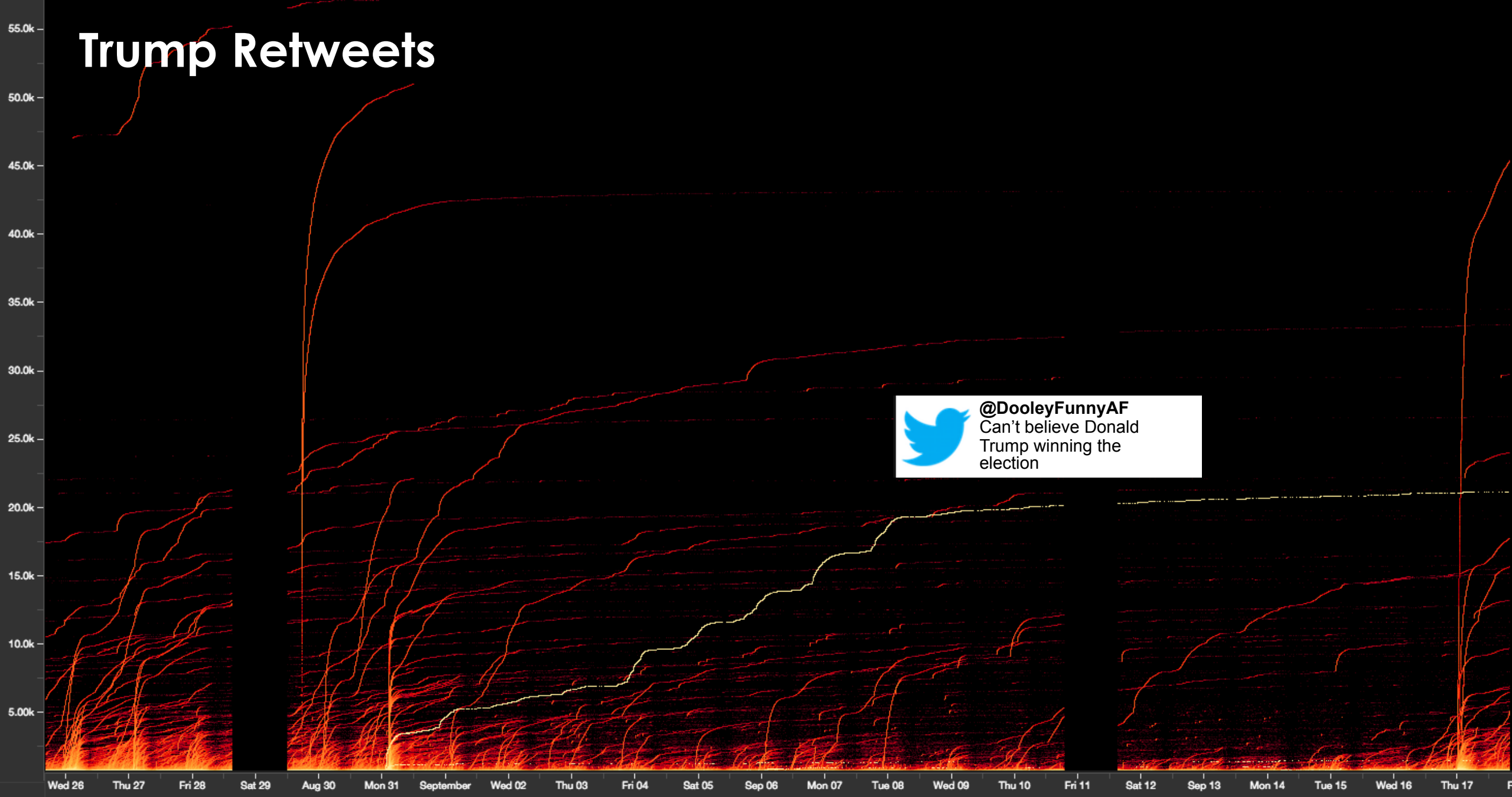



@DooleyFunnyAF

Can't believe Donald
Trump winning the
election

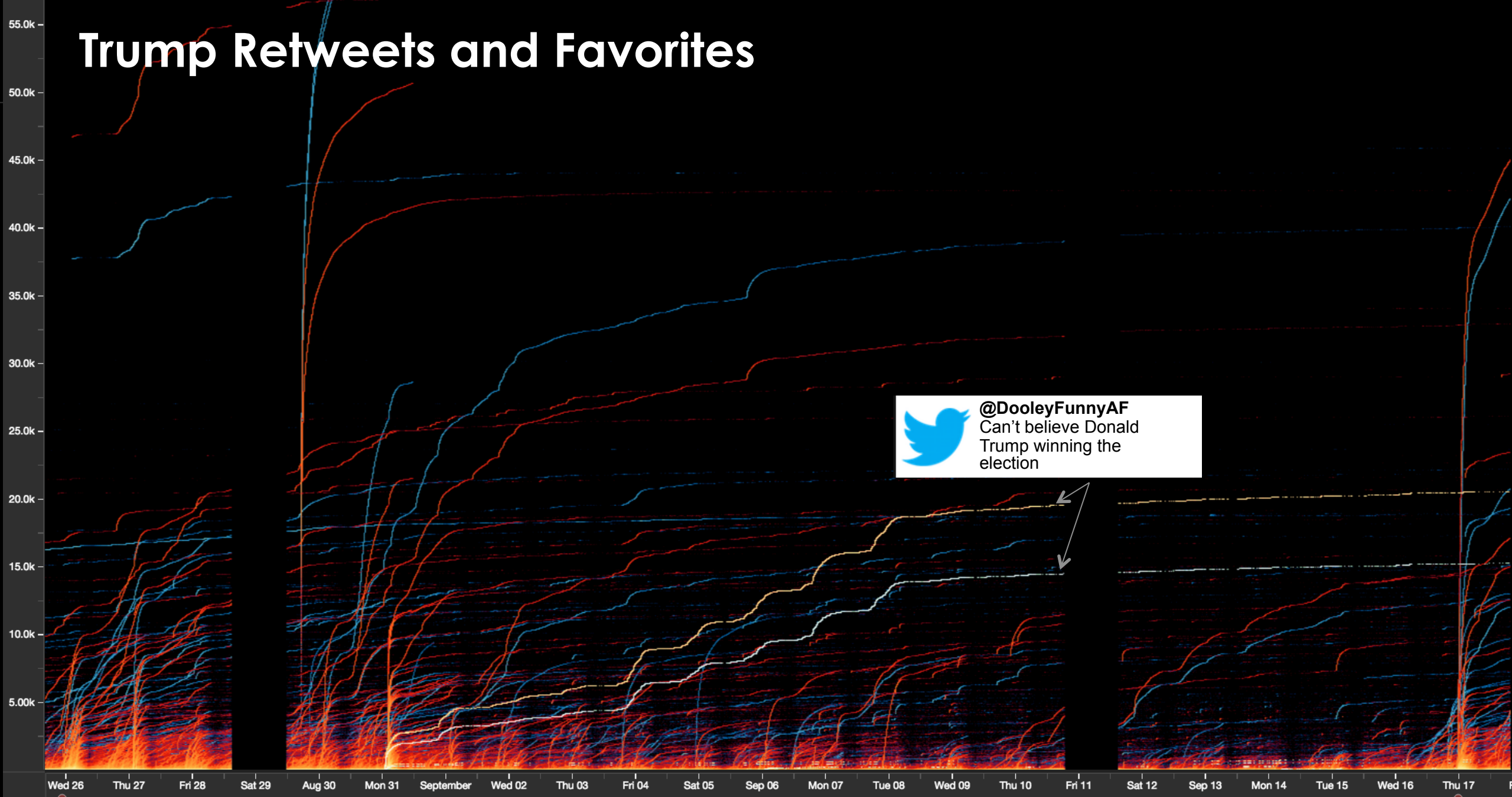
Wed 26 Thu 27 Fri 28 Sat 29 Aug 30 Mon 31 September Wed 02 Thu 03 Fri 04 Sat 05 Sep 06 Mon 07 Tue 08 Wed 09 Thu 10 Fri 11 Sat 12 Sep 13 Mon 14 Tue 15 Wed 16 Thu 17 F

Trump Retweets

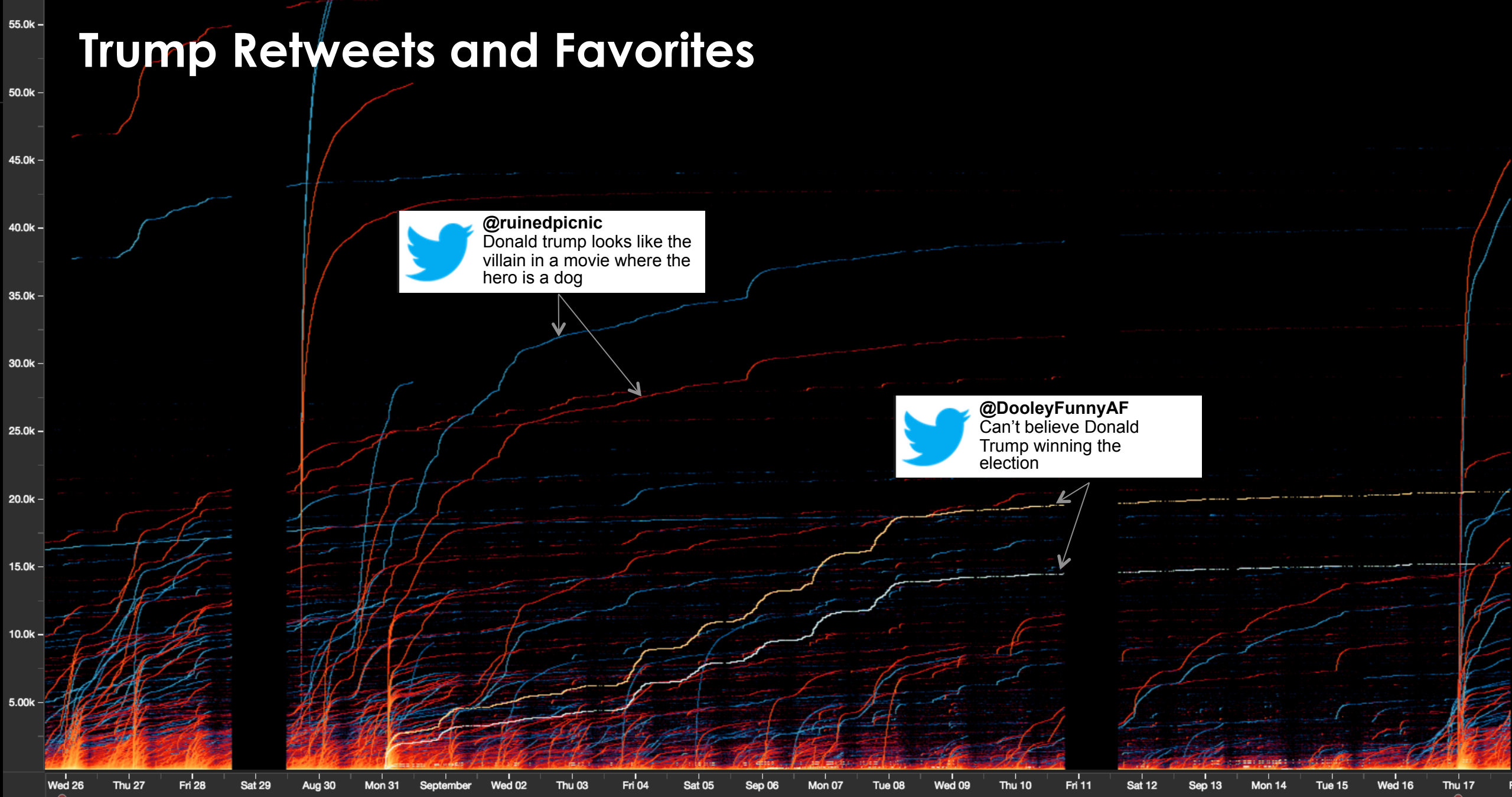


 **@DooleyFunnyAF**
Can't believe Donald
Trump winning the
election

Trump Retweets and Favorites



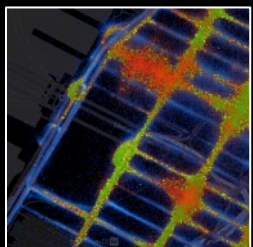
Trump Retweets and Favorites



KEY TAKE AWAYS

1. Plot all the data

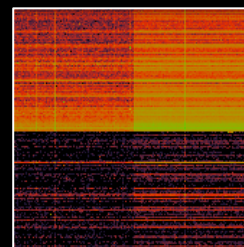
MAP



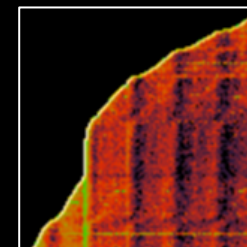
CONNECT



ORDER

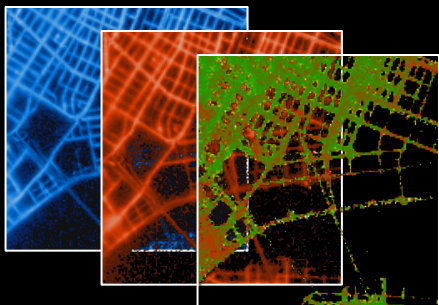


TIME

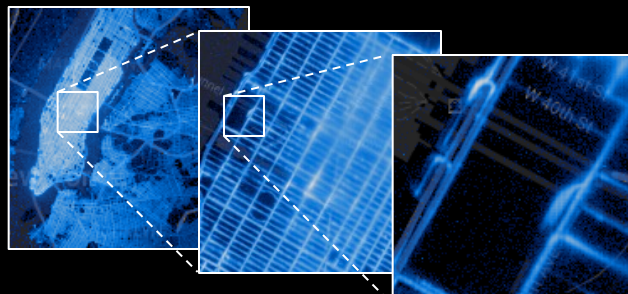


2. Explore it

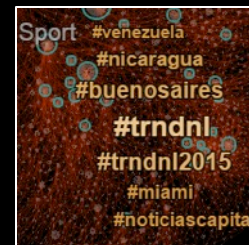
LAYERS



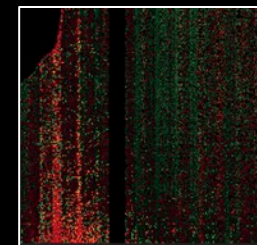
ZOOM



MINE



FILTER



More Info



Richard Brath

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



Robert Harper

rharper@uncharted.software
@rdharper