



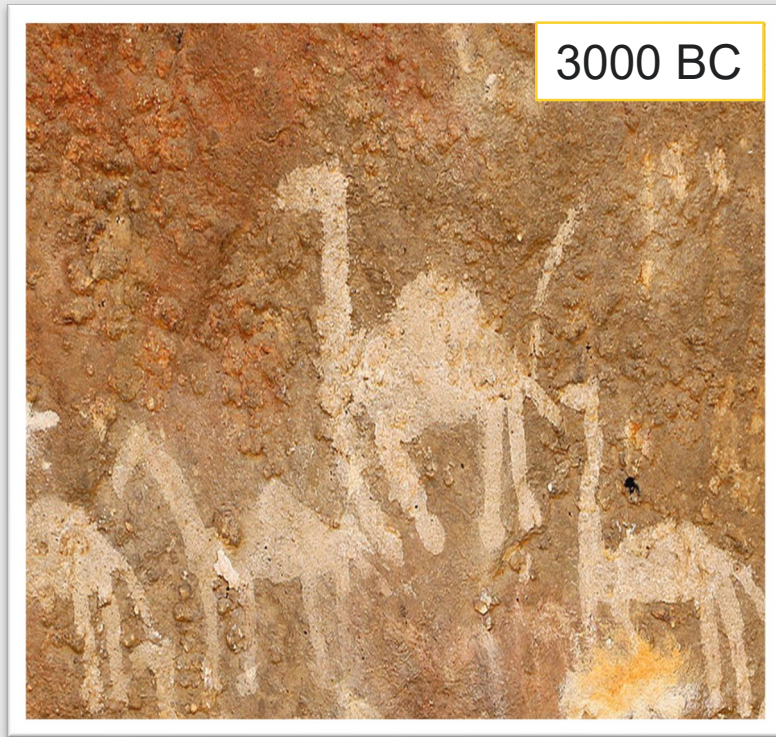
Visualize my Corpus

Dr Mahmoud El-Haj

Lecturer

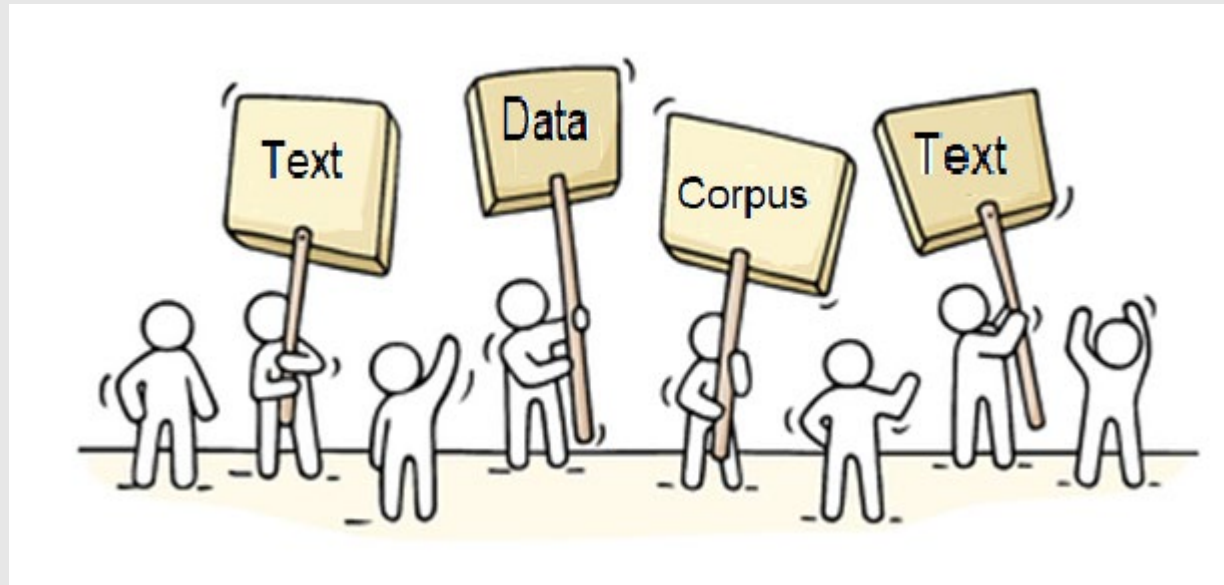
SCC, Lancaster University

What is visualisation?



communicate a message

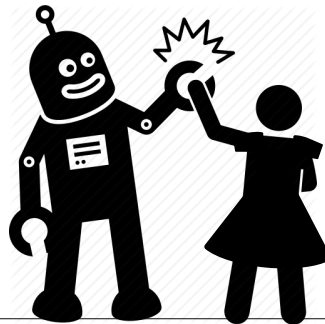
What to visualise?





WHY VISUALISE?

Machine: I can read it!



Human: I can read it too!

Isn't data enough as it is?

...keite sich ueber Computer virenz. Zwei keener die mit karpn burger (Autor des empty"Bl...
...en war die niederlaendische Gruppe hackz-TIC die eine Maschine (ca. 1500\$) vurfuehrte, die t...
...es CCC) arrangiert wurde um die Unsecurityz des Plastikgeldes zu demonstrieren. Waehrend...
...eichen CCC attackze auf die BTX/HASPA Konten war) sagte, das dies unmoeglich waere, beka...
...apie seine Karte. Allerdings wurde bei dem Versuch Geld aus einem Automaten zu beko...
...Workshop ueber 'Phreaking'. Versuche und Methoden wie man 'so weit, mit so vielen Phr...
...beschrieben. Genauso wie Tricks mit der 130-Nummer der Post(sowie Verbindungen zu den...
...kuriert wurden, um preiswertes Telefonieren zu koennen. Konferenzen und Voice Mailbox...
...l von einem US phreak zu hoeren, dass d...
...bericht war der Kommunikation zwisch...
...rd. Verwendung von Mailboxenv...
...e Sitzung beschaeftigte sich mit der...
...en Hilfe von Computerhaendlern, d...
...e das Projekt mit der Wiedervereinig...
...udent (Tommi) seine Diplomarbeit: vo...
...p Workshop Workshop Workshop Hier...
...und wie man sich anschliessen kann. Recht interessant, um einen kleinen Einblick in den Hint...
...nd Irc wurden den Anwesenden vorgebracht. Besonders hervorzuheben ist der Beitrag von...
...en der News rueberbrachte und auch eine Diskussion ueber die Zensur von Newsgruppen ir...
...fundiert war, die Meinungen jedoch nicht zusammenzubringen waren. Empty Workshoperr...
...ty"Datenschutz"empty" wurde hauptsaechlich ueber die Arbeit und die Aufgaben des Da...
...wieder die implizite Unsecurityz der heutigen Betriebssysteme Unix und DOS (PC) angefue...
...h gegen die Ueberschrift empty"Techno-terrorz"empty" gewandt hatte und erstmal eine Be...
...e"empty" sein? Die diversen zu diskutierenden Begriffe? Eingeschraenkt auf das Thema "empt...
...Pilotprojekt empty"DAWIN"empty" vorgestellt Funknetze: C-Netz, D-Netz, Cityruf etc. - Wie...
...e abgehört werden? Die von engagierten MailBoxen geschaffenen unabhängigen Bürger...
...edens- und Menschenrechtsgruppen ermöglichen. Seit vielen Monaten erhalten wir über die...
...slawien. Auf dem Congress wird eine Direktverbindung zu verschiedenen MailBoxen dort h...
...as sind Computervirenz; - Wie arebietne virenz? - Gibt es Schutz vor virenz? - Umgang mit...
... "hackz-Jäger" Protokolle: IP, TCP, UDP, ICMP Routing, Adressen Funktionsweisen: DNS, Eth...
...mputer Club treffen sich zum jährlichen Austausch organisatorischer Informationen. Insbesor

Text
Mining



software
source code
change
processor machine
server
at attack internet
signal

talk use right software
system challenge video picture platform gateway step year eine work law design seit device enclave light communication key
bug chip processor machine
computer victim framework test instrumentation process domain model files system communication key
borow server
technologies at attack internet
encryption technique paradigm programmer discussion password
leak car module team
project program calendar hardware variety information cloud component data bereits implementation
math initiative experience protection surveillance binaries einer climate performance research stellen deployment researchers hacker
encryption attention
leak car module team
project program calendar hardware variety information cloud component data bereits implementation
math initiative experience protection surveillance binaries einer climate performance research stellen deployment researchers hacker

environment infrastructure
software
source code
change
processor machine
server
at attack internet
signal

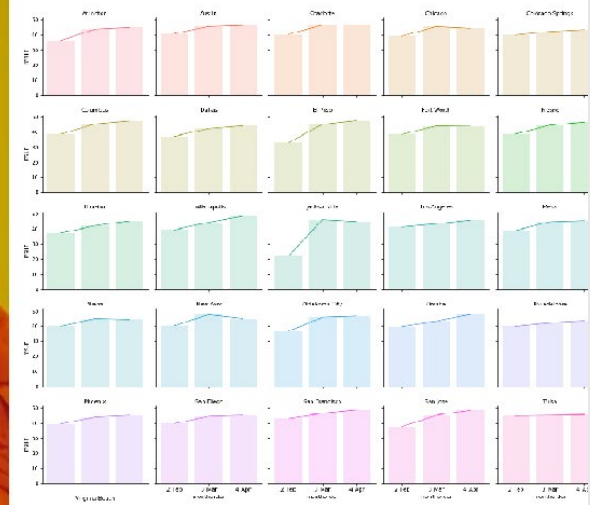


Data in plain text



Data in a plot

...eine Sitzung ... DDR zuhelfen, indem ein ziviles Computernetz (DDRNET) eingeführt werden sollte. Trotz der grossen Hilfe von Computerhändlern, die spontan PC's, Software und modems spendeten, und trotz lokaler Interessen, Kosten und Organisationsproblemen musste das Projekt mit der Wiedervereinigung gestoppt werden. Früheren Diskussionen des folgenden, ueber soziologische Aspekte des hackzens, beschreibt ein Student (Tommi) seine Diplomarbeit: Die Verwandheit von Computer und Psychologie. Workshop Workshop Workshop Diskussion Workshop Workshop Workshop Workshop Workshop Workshop Workshop Workshop Workshop Workshop Workshop Hier wurden die privaten Netze IN, SubNet, Maus und andere beschrieben, was fuer Moeglichkeiten die Netze haben und wie man sich anschliessen kann. Recht interessant, um einen kleinen Einblick in den Hintergrund der Netzwerke zu bekommen. Mail, News, FTP, Remote Login und Irc wurden den Anwesenden nahegebracht. Besonders hervorzuheben ist der Beitrag von Princess ueber die News, der engagiert und gut verstaendlich das Wesen der News rueberbrachte und auch eine Diskussion ueber die Zensur von Newsgruppen in Gang setzte, die der Zielgruppe des Vortrags (Netz-Laien) recht fundiert war, die Meinungen jedoch nicht zusammenzubringen waren. empty"Workshopempty" Cornflakes-Pfeifen (BlueBoxing)" In der Podiumsdiskussion empty"Datenschutz"empty" wurde hauptsaechlich ueber die Arbeit und die Aufgaben des Datenschutzbeauftragten geredet sowie speziell von einer Seite immer wieder die implizite Unsecurityz der heutigen Betriebssysteme Unix und DOS (PC) angefuehrt." Es fing dann an, dass Wau zur Vorgeschichte sagte, dass er sich gegen die Ueberschrift empty"Techno-terrorz"empty" gewandt hatte und erstmal eine Begriffsdefinition wuenschte. Das sollte also die empty"Definitionsfrage"empty" sein? Die diversen zu diskutierenden Begriffe? Eingeschraenkt auf das Thema "empty"Technoterrorz"empty"? Hier wurde ausfuehrlich das Muensteraner Pilotprojekt empty"DAWIN"empty" vorgestellt



What can we do with text!



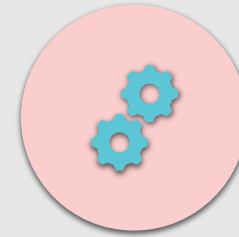
TOKENIZATION



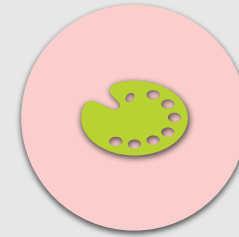
ANNOTATION



TOPIC
MODELLING



MACHINE
LEARNING



VISUALISATION

Or as in the language of memes



tokenization



annotation



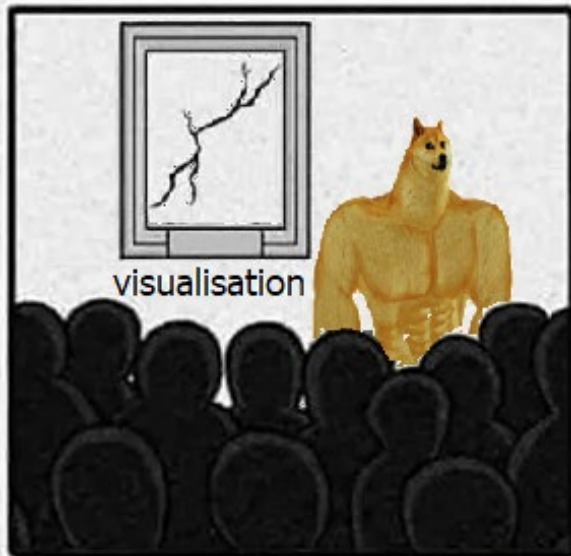
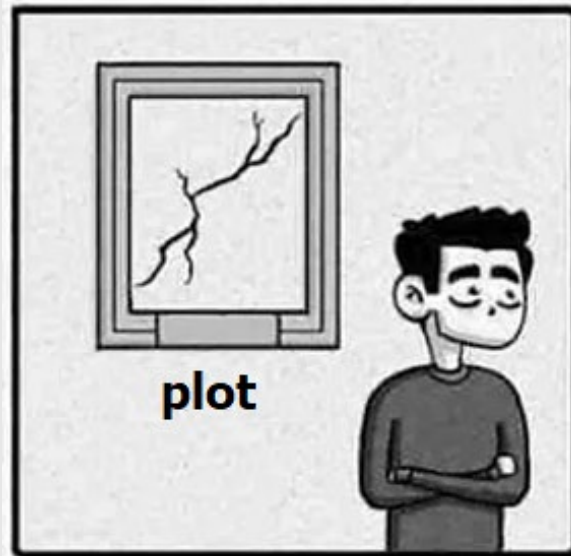
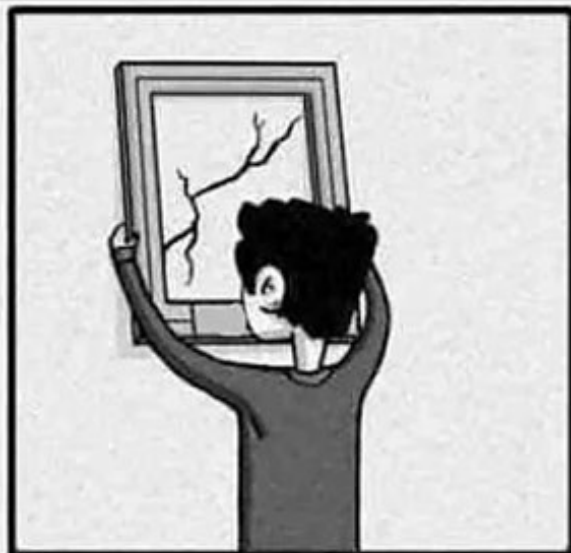
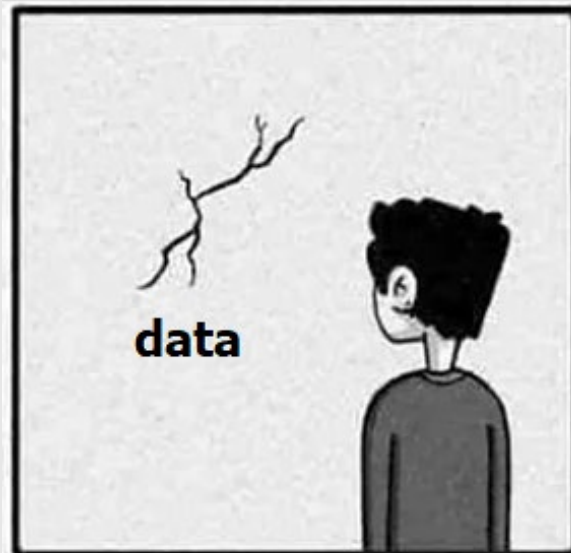
topic
modelling



machine
learning



VISUALISATION



Plain text sample

`<s>Visualise my corpus.</s>`

`<s>Text visualisation is awesome!</s>`

`<s>Visually representing the content of a text document is important.</s>`

`<s>Today is March 18th and Mahmoud is showing us how to visualising text at Lancaster University, well online!</s>`

Added tags (annotations)

<s>Visualise [VERB] my [DET] corpus [NOUN].</s>

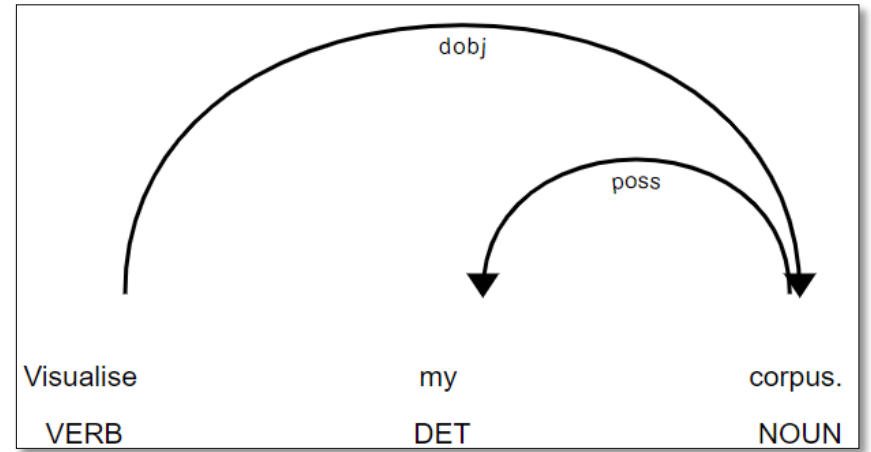
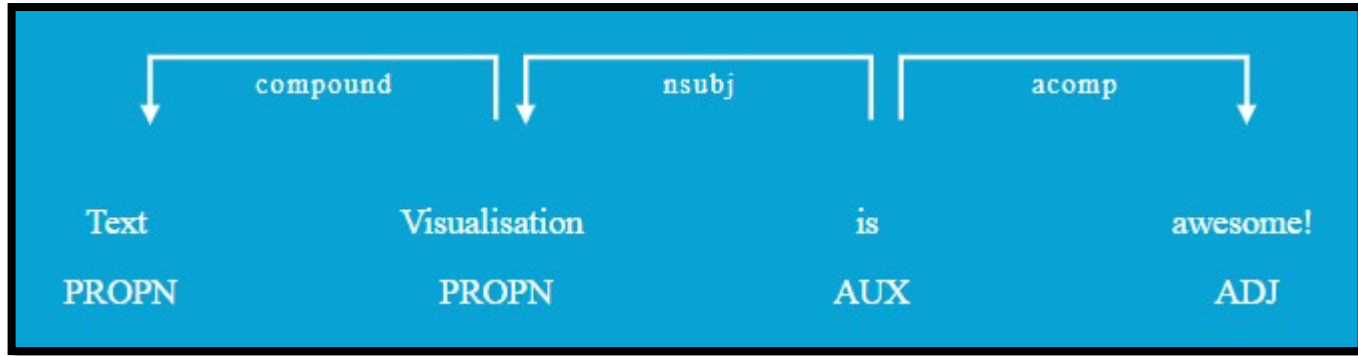
<s>Text [PROPN] visualisation [PROPN] is [AUX] awesome! [ADJ]</s>

<s>Visually [ADV] representing [VERB] the [DET] content [NOUN] of [ADP] a [DET] text [NOUN] document [NOUN] is [AUX] important [ADJ].</s>

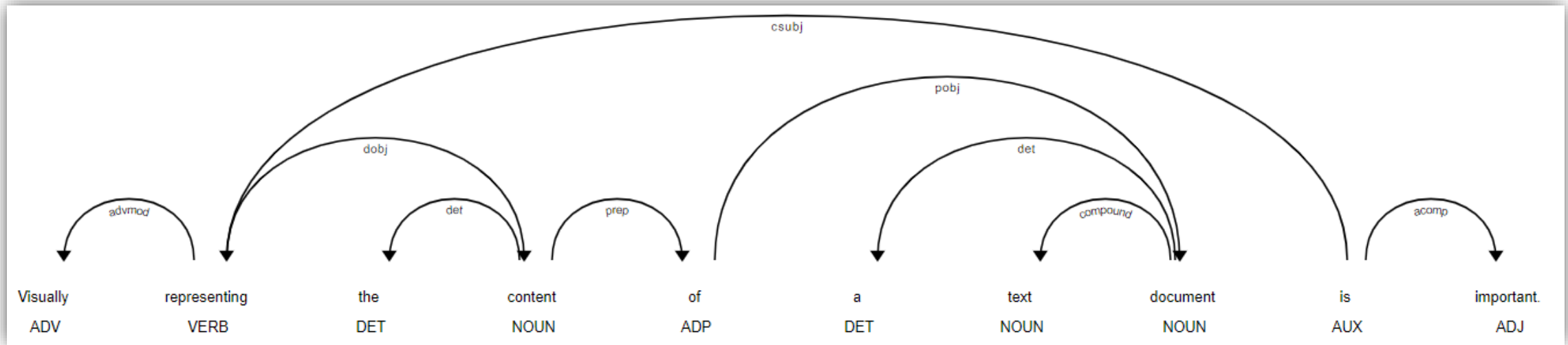
<s>{Today is March 18th}[DATE] and {Mahmoud}[PERSON] is showing us how to visualising text at {Lancaster University}[ORG].</s>

POST

NER



Today is March 18th **DATE** and **Mahmoud PERSON** is showing us how to visualise text at **Lancaster University ORG**





What can visualisation tell us?

Things we don't directly see!

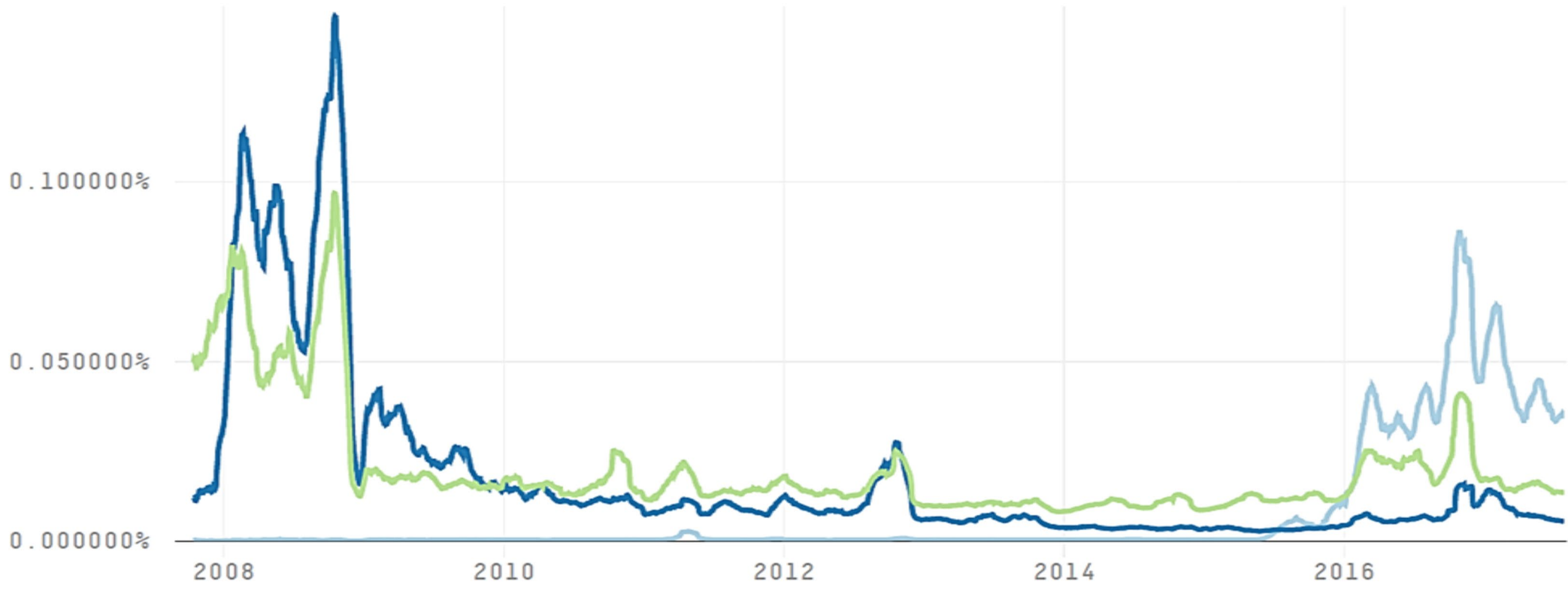
Reddit.com



a social news aggregation, web content rating, and discussion website

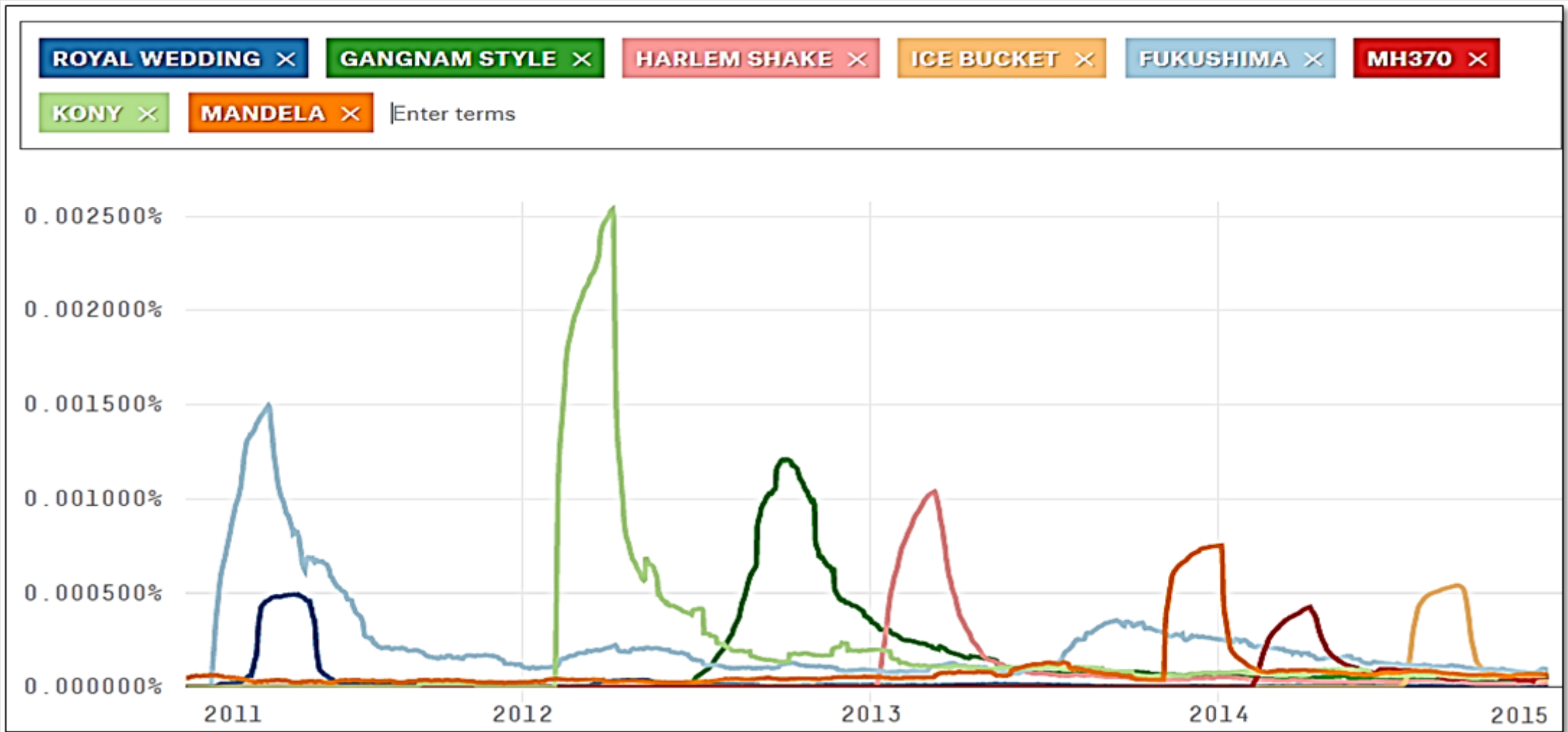
TRUMP vs OBAMA vs VOTE

TRUMP × OBAMA × VOTE × Enter terms



<https://projects.fivethirtyeight.com/reddit-ngram/?keyword=trump.obama.vote&start=20071015&end=20170731&smoothing=22>

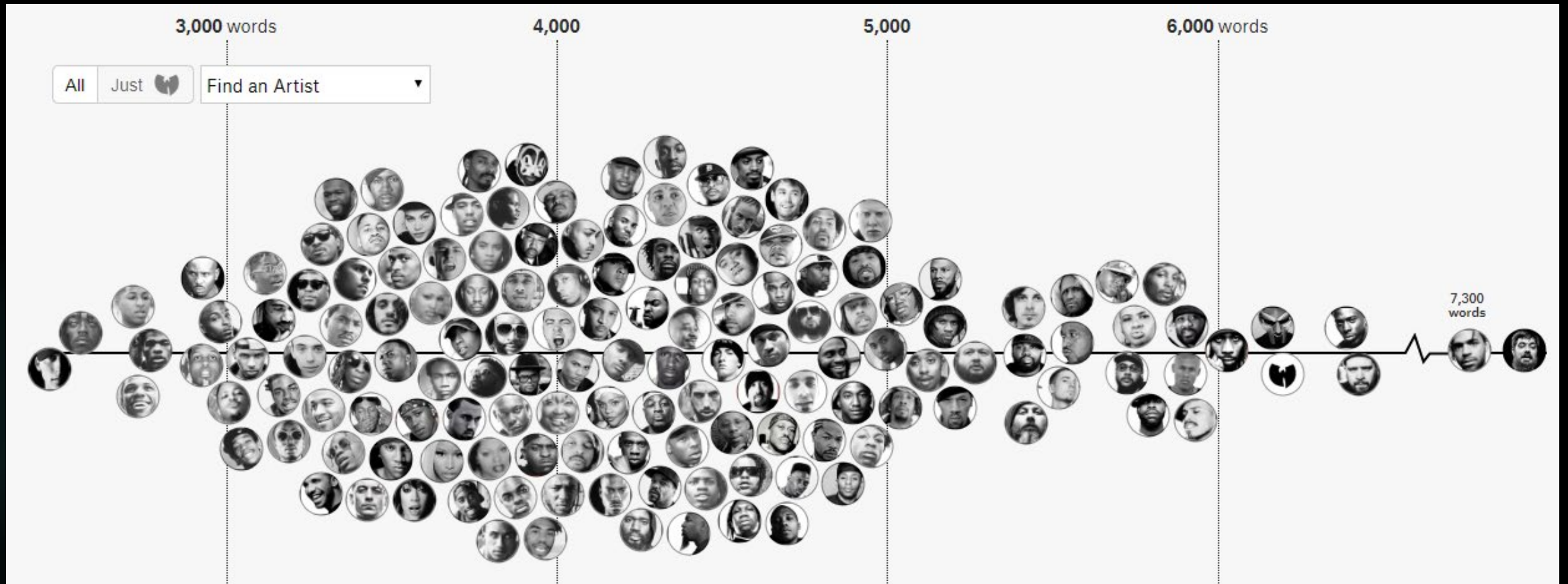
How the Internet Talks!



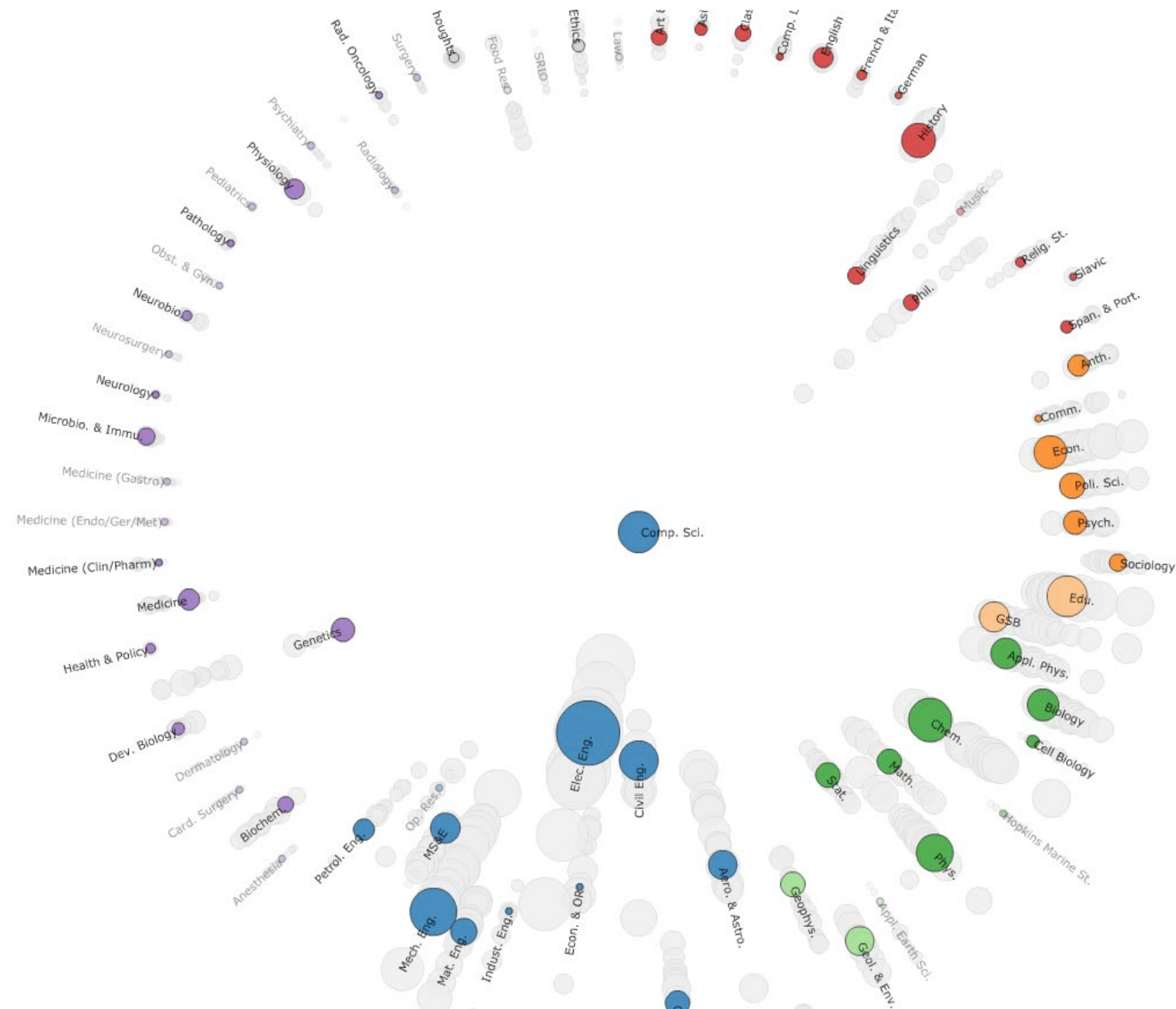
Major events/trends between 2011 - 2015

The Largest Vocabulary In Hip Hop

Rappers, ranked by the number of unique words used in their lyrics



In May 2014, a study by Matt Daniels found that **Aesop Rock**'s vocabulary in his music surpassed 85 other major hip-hop and rap artists, as well as Shakespeare's works and Herman Melville's *Moby Dick*, being named the largest vocabulary in Hip Hop ^{[43][44]} [Wikipedia](#)



Stanford PhD Theses

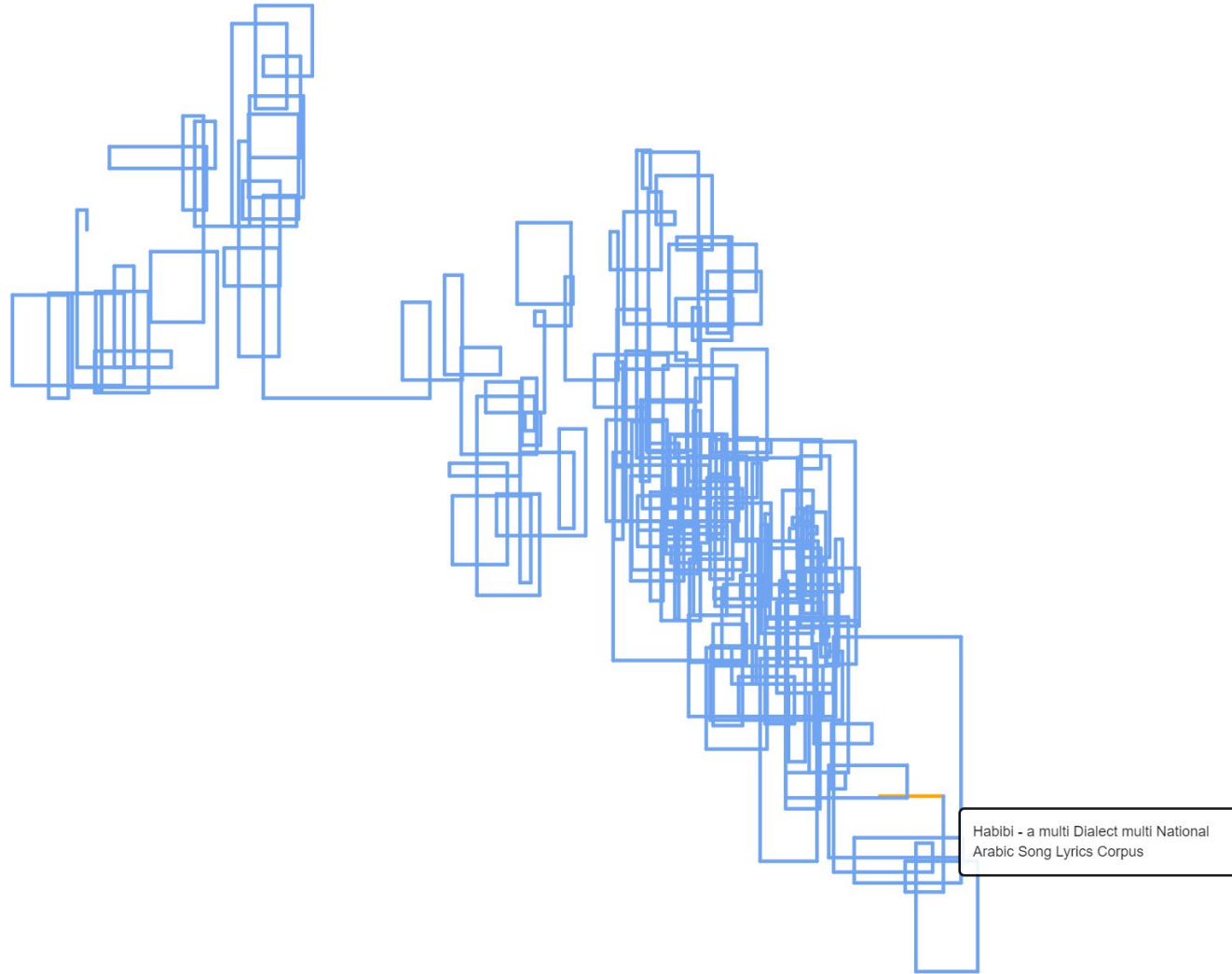
<https://nlp.stanford.edu/projects/dissertations/browser.html>

Stanford's PhD dissertation abstracts from 1993-2008



HABIBI PAPER

Habibi paper ▾



Drawing one of my papers*

Try it out:

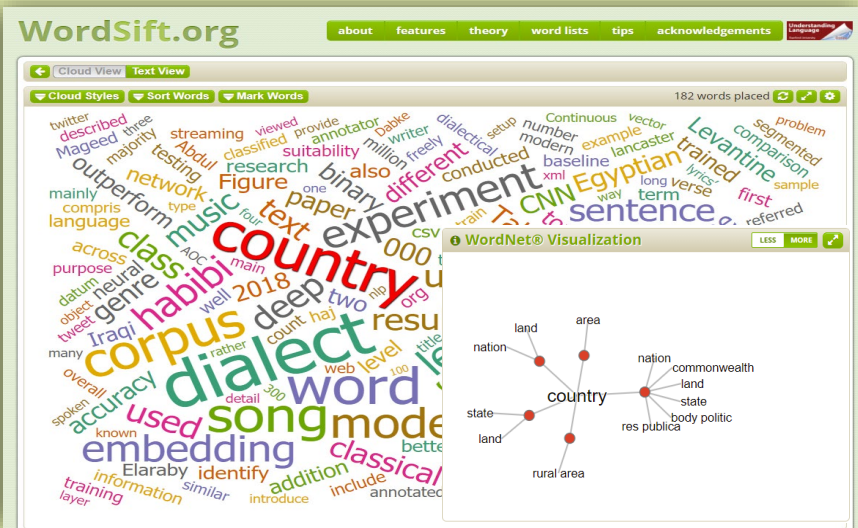
<https://www.lancaster.ac.uk/staff/elhaj/draw/index.html#habibi>

[*https://www.lancaster.ac.uk/staff/elhaj/docs/habibi.pdf](https://www.lancaster.ac.uk/staff/elhaj/docs/habibi.pdf)

[GitHub](#)

[Tyler Rinker Blog](#)

Stefanie Posavec
Understand more about text through Art!
Unlock hidden data!
<https://youtu.be/y1wkGMLEktQ>



<https://wordsift.org>



<http://wordwanderer.org/>

More Text Visualisation Browser

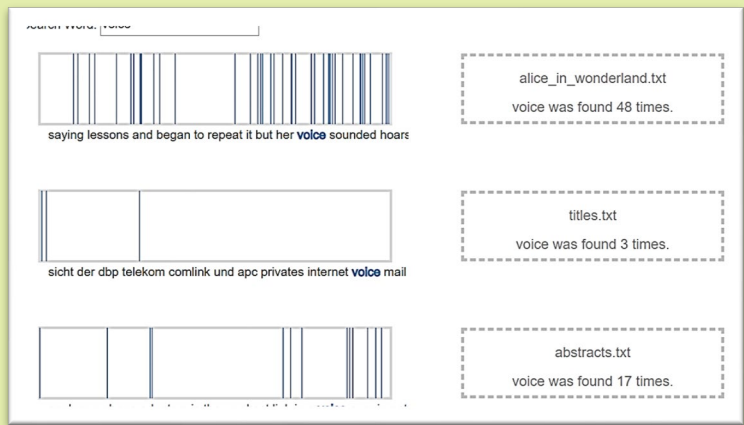
<https://textvis.lnu.se/>

Data Driven Documents

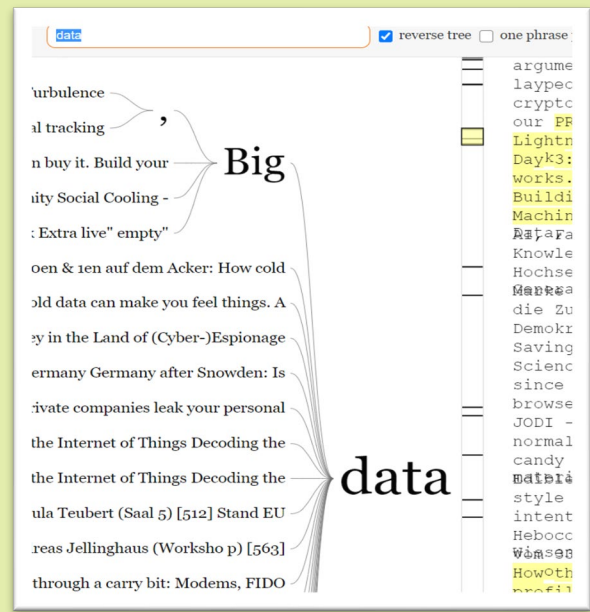
<https://d3js.org/>

Jim Vallandingham

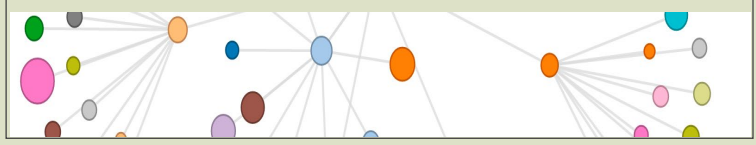
<http://vallandingham.me>



http://vallandingham.me/concordance_plot/



<https://www.jasondavies.com/wordtree>





What are we visualising today?

Chaos Communication Congress talks
(<https://www.ccc.de/en/>)

Tutorial Material:

<https://github.com/drelhaj/NLP-ML-Visualization-Tutorial>

Conference Data:

<https://gitlab.com/maxigas/cc-congresstalks/>

Visualisation Roadmap

