

Visualization-Guided Text Analysis in Python

Jason Kessler

Presentation for Indie Thinkers

June 22, 2021

Overview

- If you'd like to follow along, please install Scattertext and related packages
 - <https://github.com/jasonkessler/IndithinkersTutorial>
- Visualize variation in language
 - Between groups (often by using proxies for membership)
 - Left vs. right
 - Recipients of the Pfizer or Moderna vaccines
 - By gender
 - Over time
- We will first discuss some work in visualizing categorized text
- Then discuss how to use Scattertext to analyze social media chatter around Vaccines
- Discuss work in visualizing language change over time
- Demo a new feature of Scattertext which allows you to see how topics evolve over time

OKCupid: an online dating site

Which words and phrases statistically distinguish ethnic groups and genders?

hobos

almond butter



Our source:
526,000 OkCupid
profiles

The image shows three side-by-side OKCupid profile pages. Each page has a header with tabs: 'About', 'Photos', 'The Rest of Us', and 'Personality'. Below the header is a 'My Self-Summary' section, followed by a 'What I'm doing with my life' section, and finally a 'I'm really good at' section. The 'My Self-Summary' sections all mention 'hobos' (highlighted in yellow). The 'What I'm doing with my life' sections mention 'almond butter' (highlighted in yellow) and 'Bikram yoga' (highlighted in yellow). The 'I'm really good at' sections mention 'judging' (highlighted in yellow) and 'Japanese' (highlighted in yellow). A blue arrow points from the word 'hobos' to the first profile's self-summary. Another blue arrow points from the word 'almond butter' to the second profile's self-summary. A blue arrow points from the word 'Bikram yoga' to the third profile's self-summary.

Profile 1 (Left):

- My Self-Summary:** In my free time I train **hobos** for the circus.
<http://www.youtube.com/watch?v=3V7MBULRoo>
- What I'm doing with my life:** study languages, drink gin and pickle brine **Martini**, bake, eat **almond butter** with a spoon, draw things with a pen, write things with a different pen.
- I'm really good at:** **judging**.
- The first things people usually notice about me:** if you're **japanese**, it's that I look like hermione granger.

Profile 2 (Middle):

- My Self-Summary:** I'm a Colorado native who loves to ski but gave up the mountains of Colorado for the city five years ago. I feel like a **Rocky Mountain girl** living a Big City life - not really firmly belonging in either place. I love it here though, especially when I discover something I didn't know about before, recapturing that sense of adventure even though this is the place I live. I like **adventures!**
- What I'm doing with my life:** applying to graduate programs, fighting the urge to say fuck it and go to **pastry school**.
- I'm really good at:** **judging**.

Profile 3 (Right):

- My Self-Summary:** I'm a California girl, currently in a New York state of mind. I moved to NYC **1st** for **grad school** and stayed on for work. I'm enjoying my time here exploring Manhattan, but ultimately planning on moving back to CA. I'm pretty **easy-going** and I like to try new things.
- What I'm doing with my life:** I am composed, thoughtful, and sarcastic.
- My favorite books, movies, music, and food:**
 - Books:** Anything by Jhumpa Lahiri, **100 Years of Solitude**, The Great Gatsby
 - Movies:** Jurassic Park, Coming to America, Women on the Verge of a Nervous Breakdown
 - Music:** Tupac, 5-45, Jay-Z, Mariah Carey, Tamia, Justin Timberlake, Drake
 - Food:** Japanese, Thai, New American, California Cuisine

100 Years of Solitude

Christian Rudder. "The Real 'Stuff White People Like'" 2010. <https://www.gwern.net/docs/psychology/okcupid/thererealstuffwhitepeoplelike.html>
Rudder was OKCupid's founding data scientist. This blog post has since been removed from okcupid.com.



OKCupid: Words and phrases that distinguish white men.

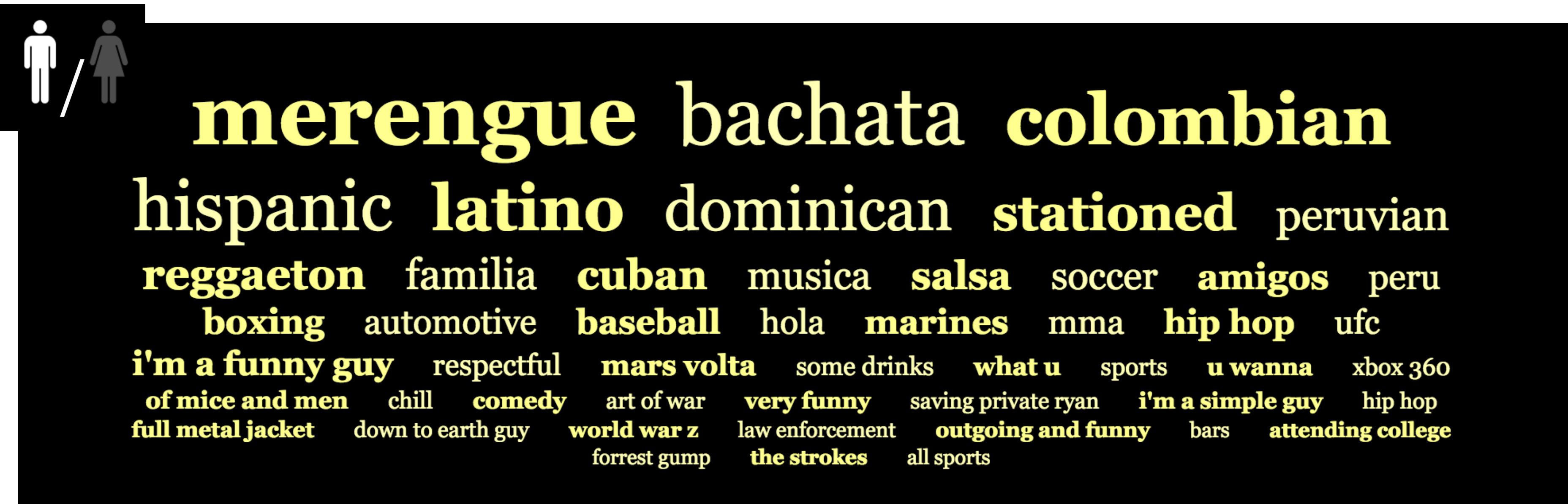


tom clancy van halen **golfing**
harley davidson **ghostbusters** phish
the big lebowski soundgarden **brew** boating **nofx**
groundhog day **hockey** jeep **blazing saddles** the red sox
the dropkick murphys megadeth **grilling** ccr
robert heinlein boats **skiing** zappa **nascar** motorcycles
software dark tower **the hitchhiker's guide to the galaxy** breaking bad
band of brothers burn notice **coen brothers** michael crichton **bad religion** tenacious d
mostly rock i'm a country boy **building things** queens of the stone age **mountain biking**
i can fix anything **the offspring** a few beers **apocalypse now** lock, stock, and two smoking barrels
hunting and fishing most sports **world war z** guitar

In general, I won't comment too much on these lists, because the whole point of this piece is to let the groups speak for themselves, but I have to say that the mind of the white man is the world's greatest sausagefest. Unless you're counting **Queens of the Stone Age**, there is not even one vaguely feminine thing on his list, and as far as broad categories go we have: sweaty guitar rock, bro-on-bro comedies, things with engines, and dystopias.



OKCupid: Words and phrases that distinguish Latin men.



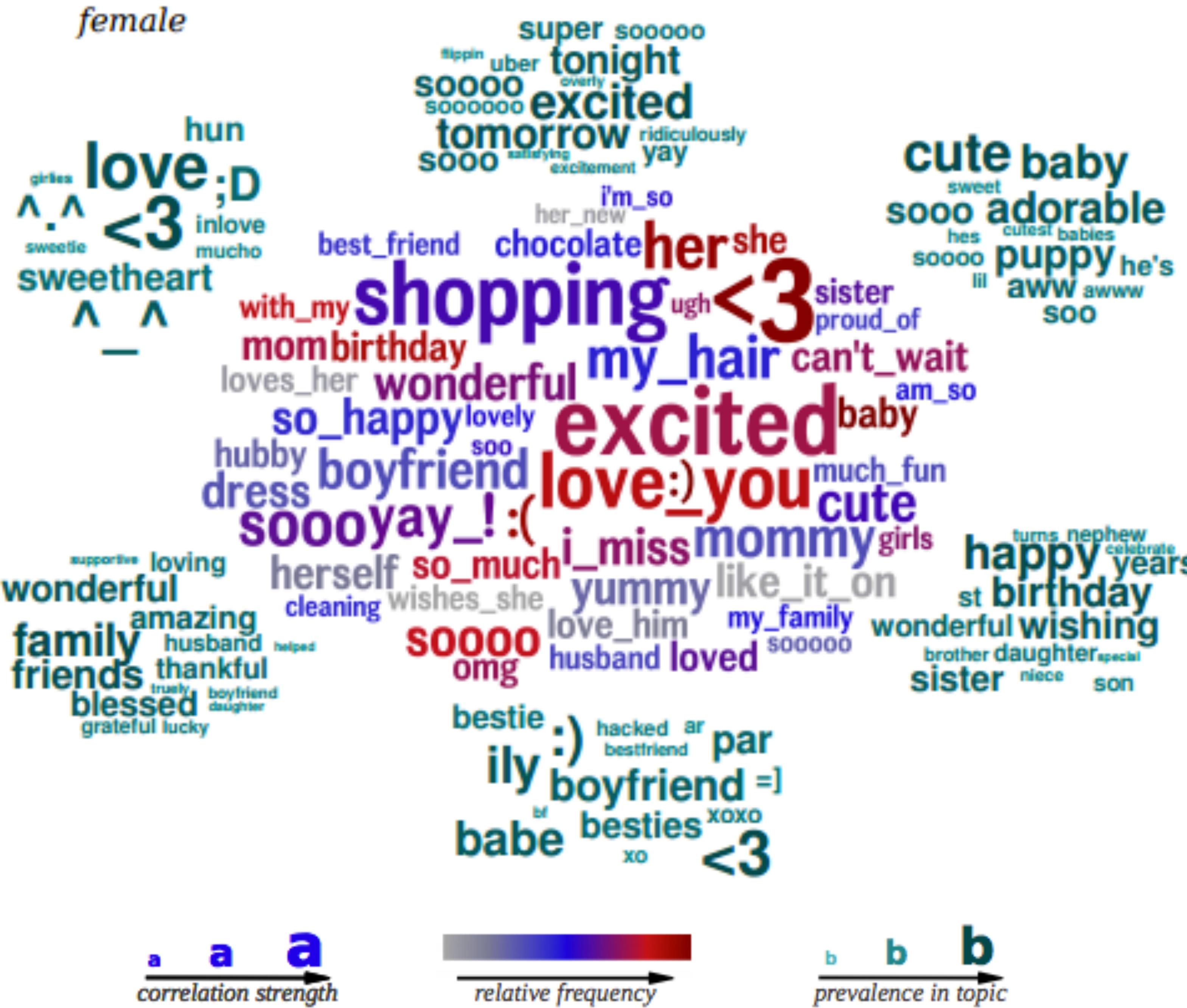
Explanation

Music and dancing—**merengue**, **bachata**, **reggaeton**, **salsa**—are obviously very important to Latinos of both genders. The men have two other fascinating things going on: an interest in telling you about their sense of humor (**i'm a funny guy**, **very funny**, **outgoing and funny**, etc.) and an interest in industrial strength ass-kicking (**mma**, **ufc**, **boxing**, **marines**, etc.) Basically, if a Latin dude tells you a joke, you should laugh.



Christian Rudder. "The Real 'Stuff White People Like'" 2010. <https://www.gwern.net/docs/psychology/okcupid/theralstuffwhitepeoplelike.html>
Rudder was OKCupid's founding data scientist. This blog post has since been removed from okcupid.com.

Word Use and Gender in Facebook Statuses



The good:

- Word clouds force you to hunt for the most impactful terms
- You end up examining the long tail in the process
- Compactly represent a lot of phrases and topics

Lyle Ungar
2013 AAAI
Tutorial



Word Use and Gender in Facebook Statuses



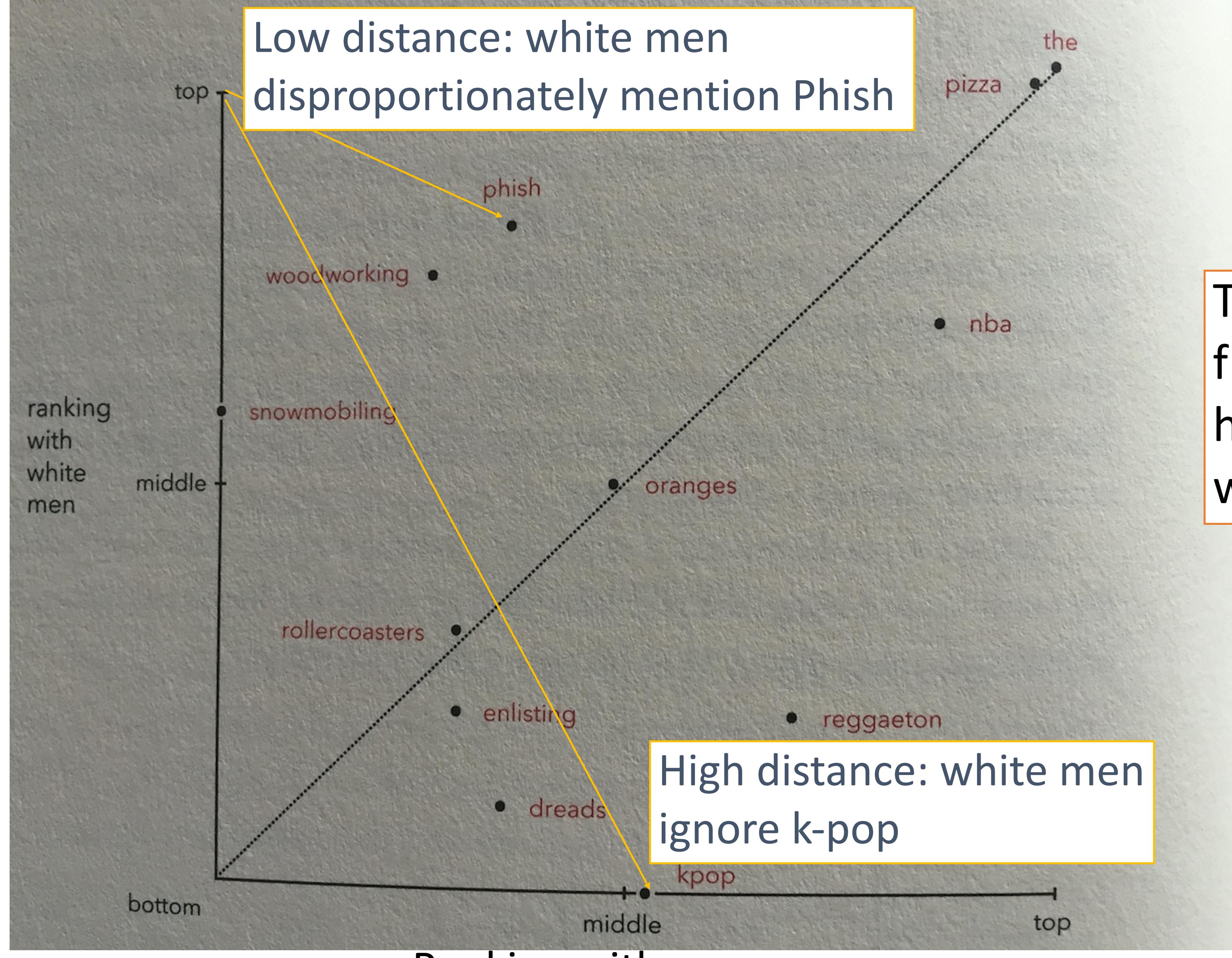
The bad:

- “Mulletts of the Internet”
 - Jeffrey Zeldman, 2005
 - Longer phrases are more prominent.
 - Ranking is unclear
 - Does size indicate higher frequency?

Lyle Ungar

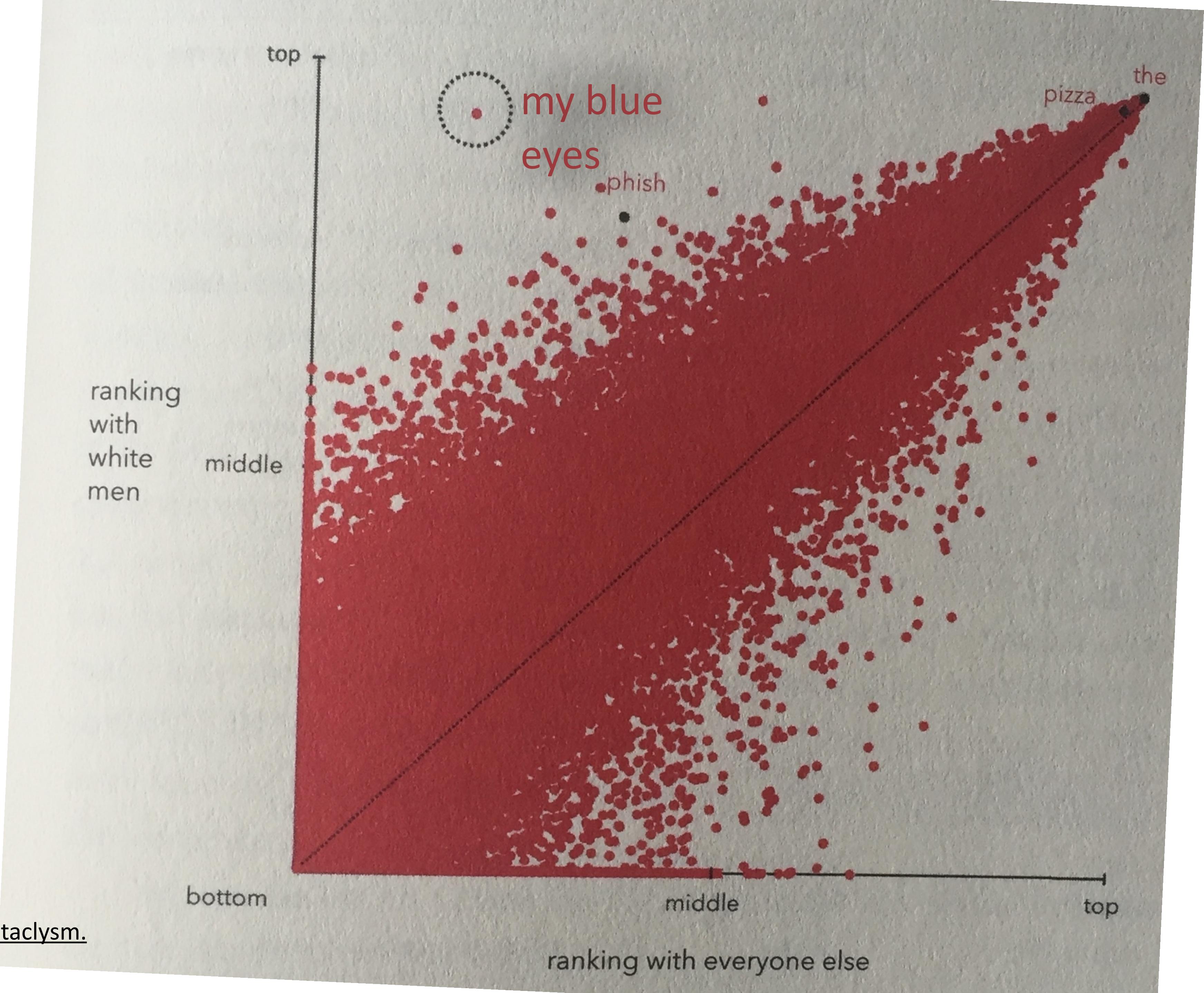
2013 AAAI Tutorial





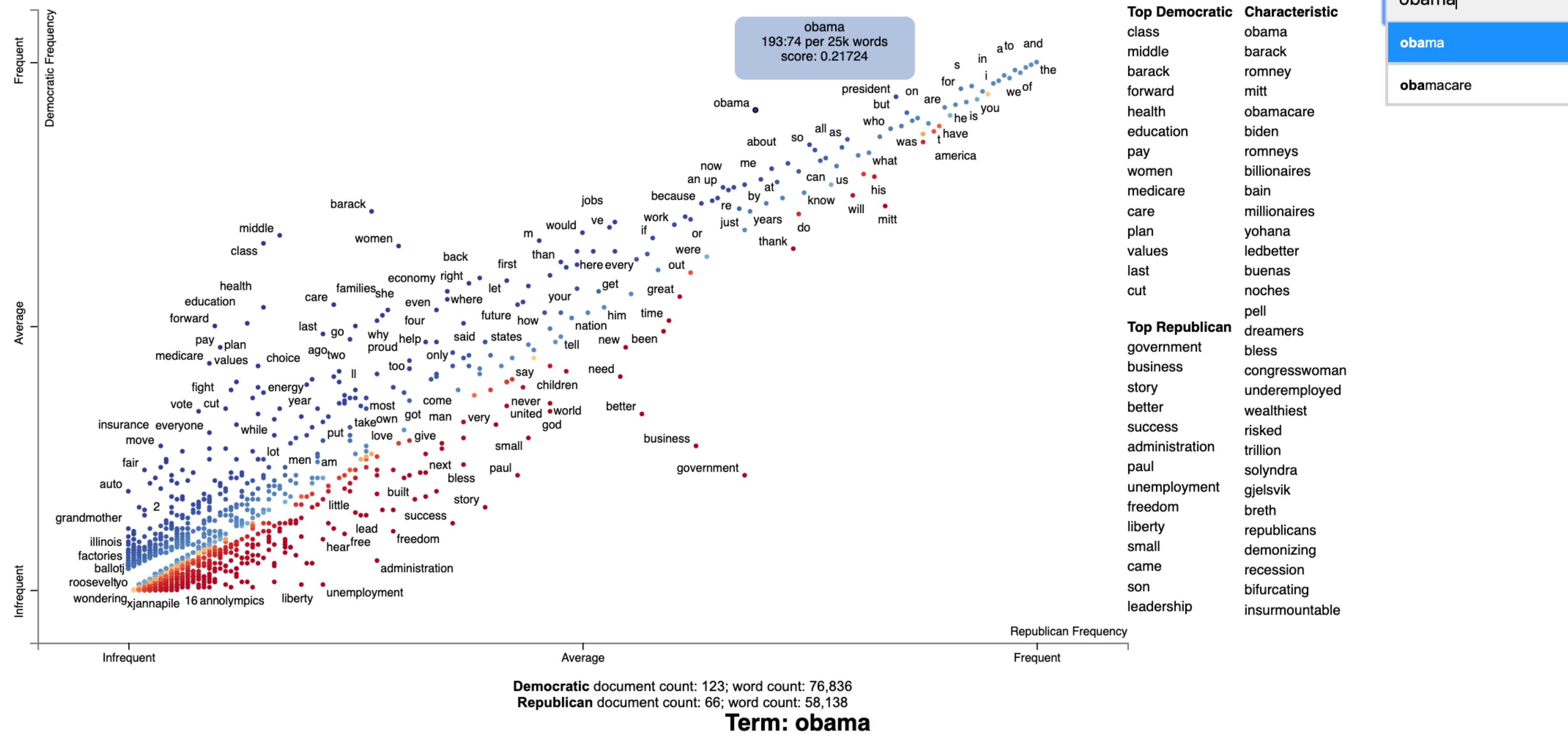
The smaller the distance from the top left, the higher the association with white men.





Source: Christian Rudder. Dataclysm.
2014.





Democratic frequency:

193 per 25,000 terms

919 per 1,000 docs

Some of the 535 mentions:

RICHARD DURBIN

It was a cold, cold January afternoon when Barack **Obama** lifted his hand from Abraham Lincoln's Bible and looked out on an America facing an economic collapse.

Well, President **Obama** and millions of American families think it's a great idea for America.

And with President **Obama** and Vice President Joe Biden in the White House, we will.

Republican frequency:

74 per 25,000 terms

606 per 1,000 docs

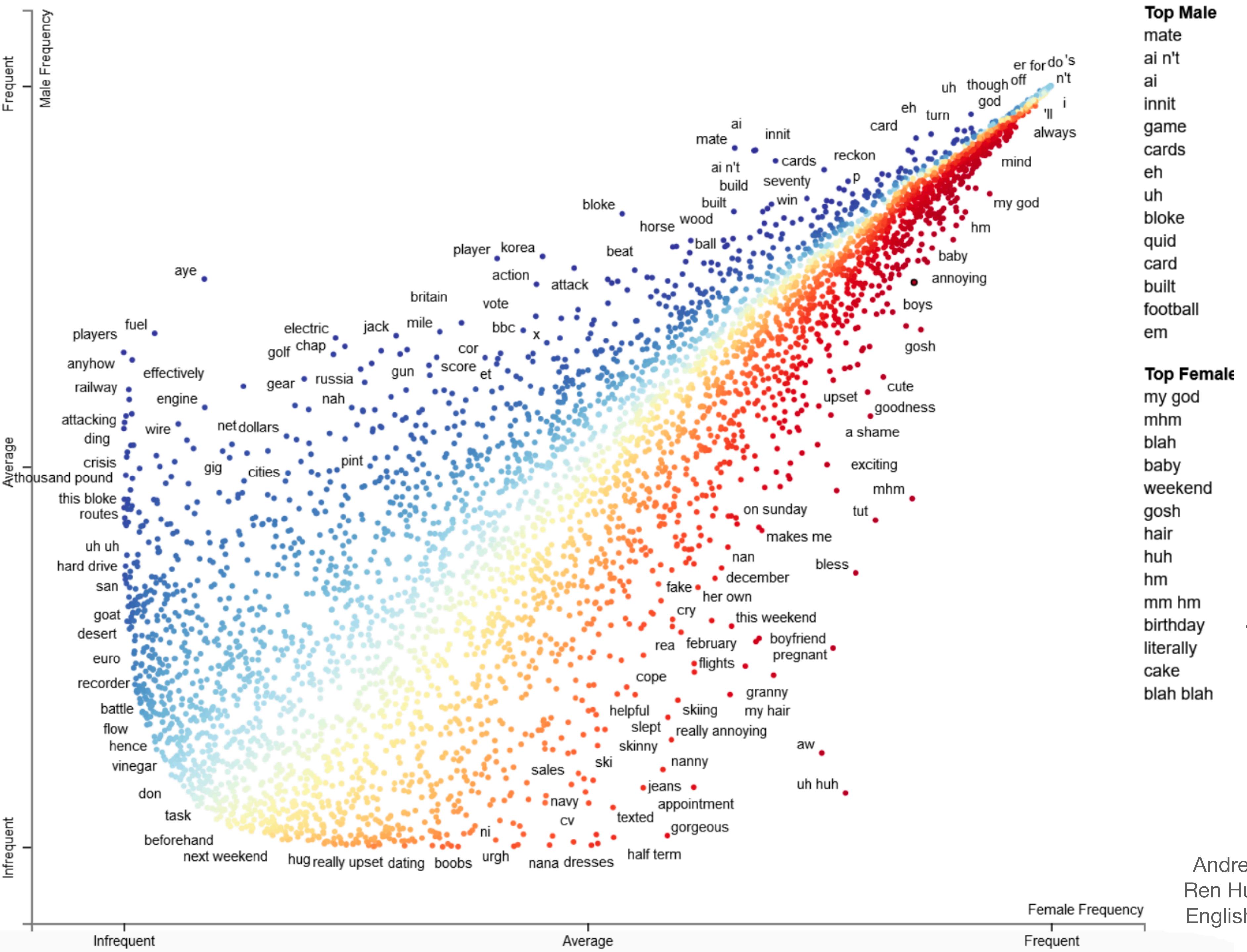
Some of the 167 mentions:

MITT ROMNEY

I wish President **Obama** had succeeded because I want America to succeed.

If you felt that excitement when you voted for Barack **Obama**, shouldn't you feel that way now that he's President Obama?

Some of the companies we helped start are names you — you know and you've heard from tonight: an office company called Staples where I'm pleased to see the **Obama** campaign's been shopping — — today Steel

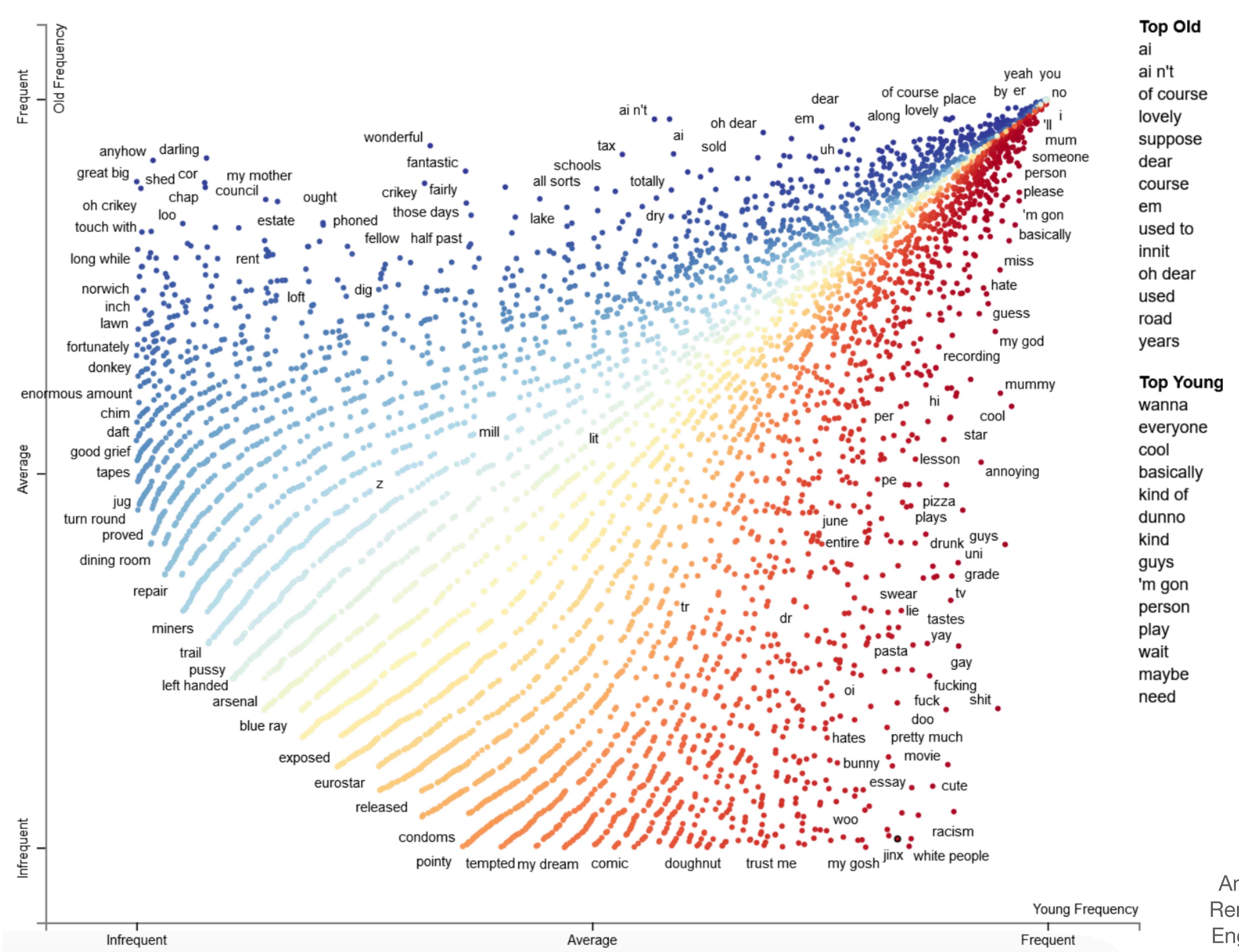


Scattertext in the wild:

- Identifying differences between gender and age in SpokenBNC2014

- Differences in “minimal particles” are particularly interesting
 - M: Eh, um, em,
 - F: mhm, hum, hm, mm hm

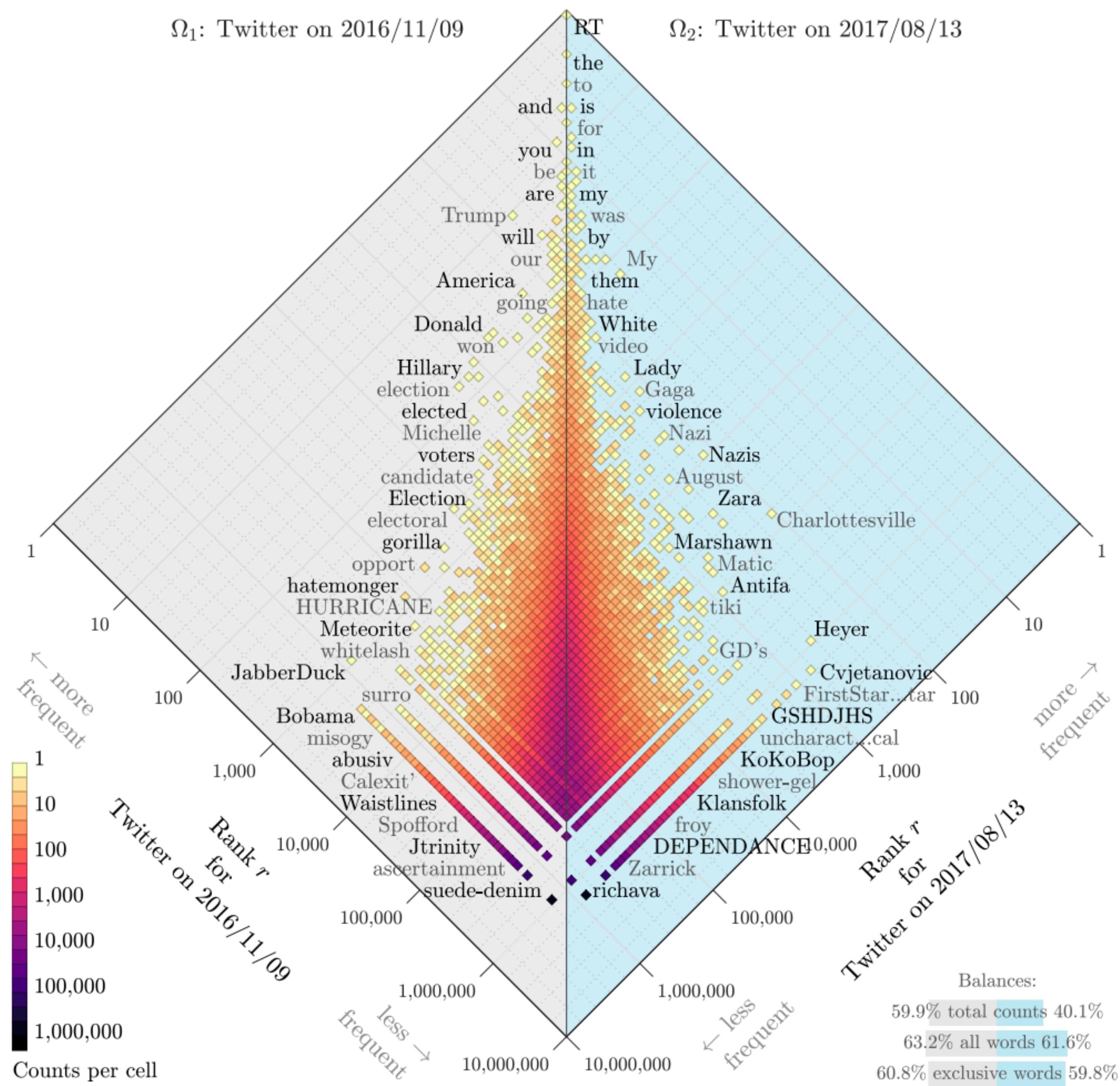
Plotting roughly on frequency rank; ties are broken alphabetically



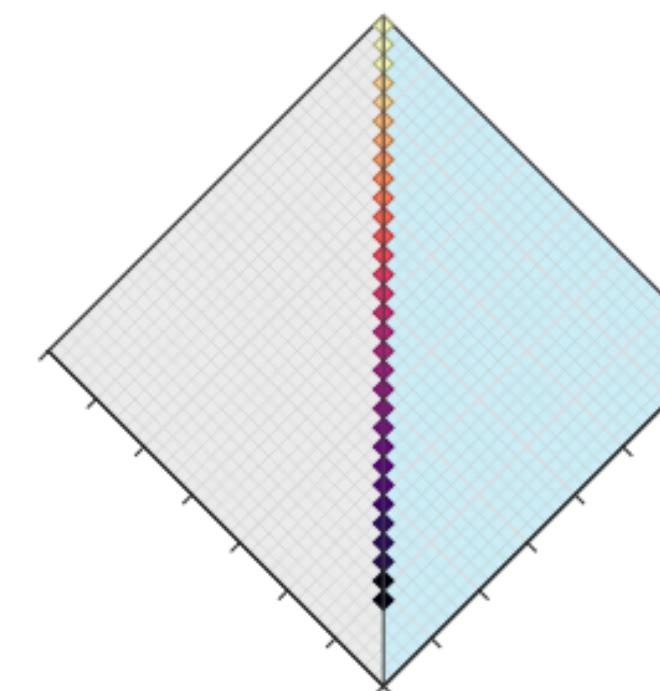
Scattertext in the wild: -Old (>=70) vs. young (<=19)

Author's note (Oct 2020): statement on the use of social categories in this study This work involves the labelling of participants for social categories related to gender and age. We caution against the use of this heuristic due to the risk of promoting a biased view on the topics. We would like to encourage those interested in the computational modelling of social categories to join the discussion on these concerns and consider participating in efforts to build more inclusive resources for the study of the topics.

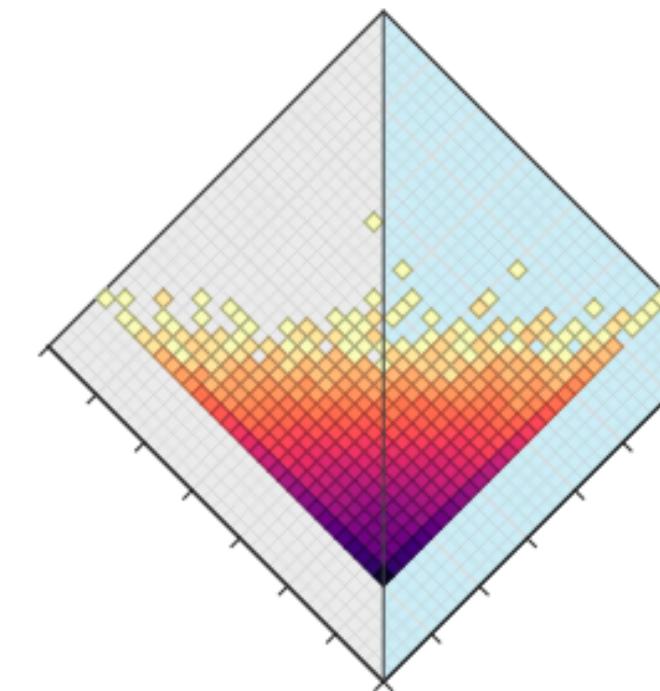
A. Rank-turbulence histogram:



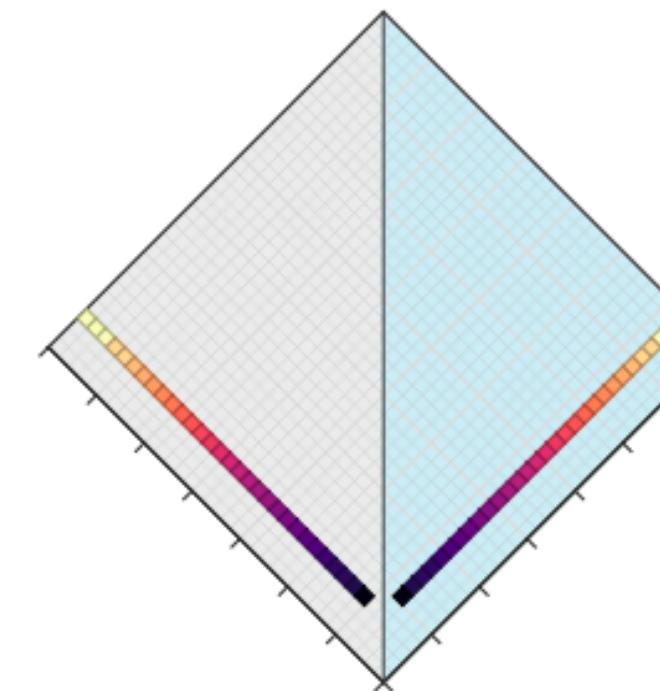
B. Identical systems:



C. Randomized systems



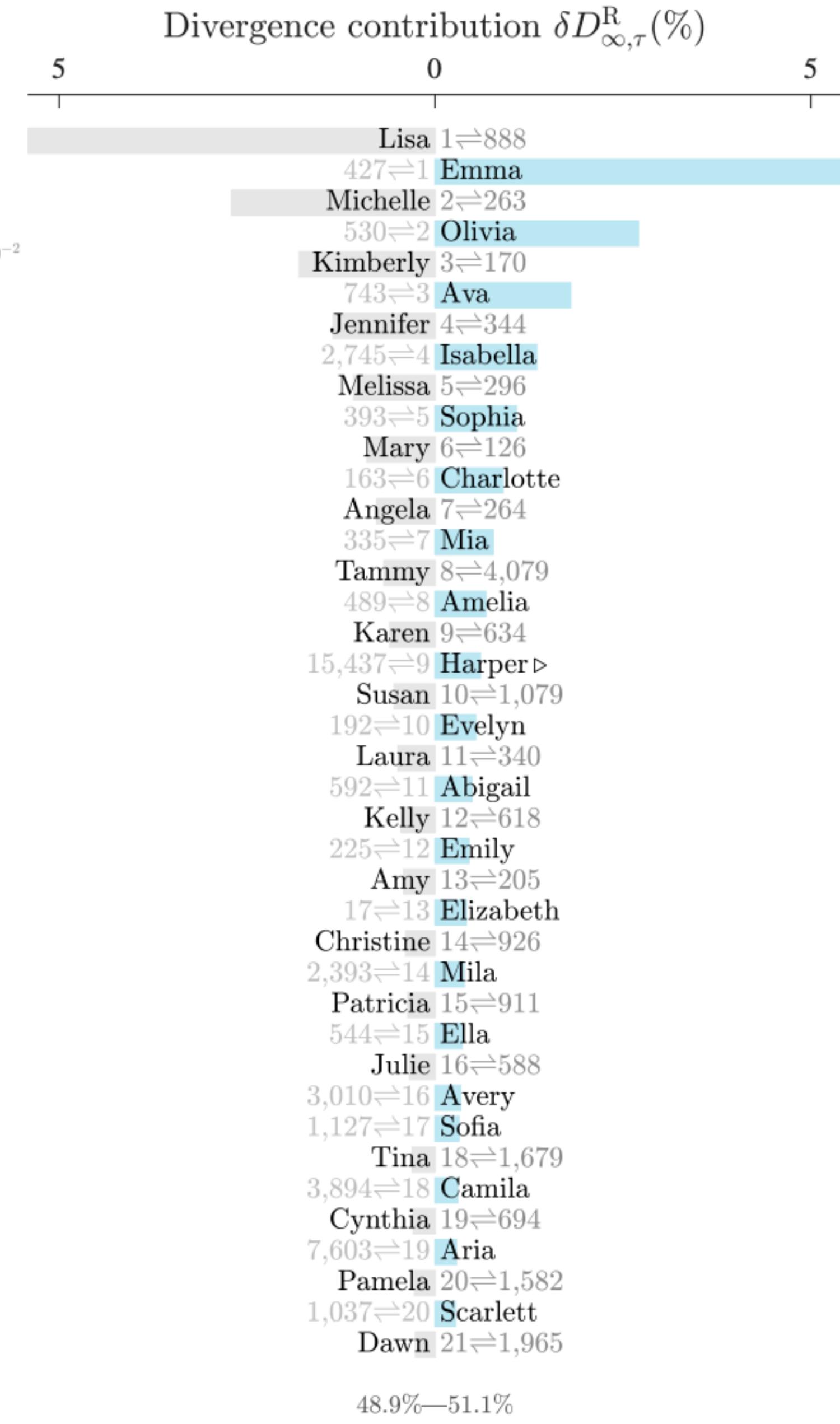
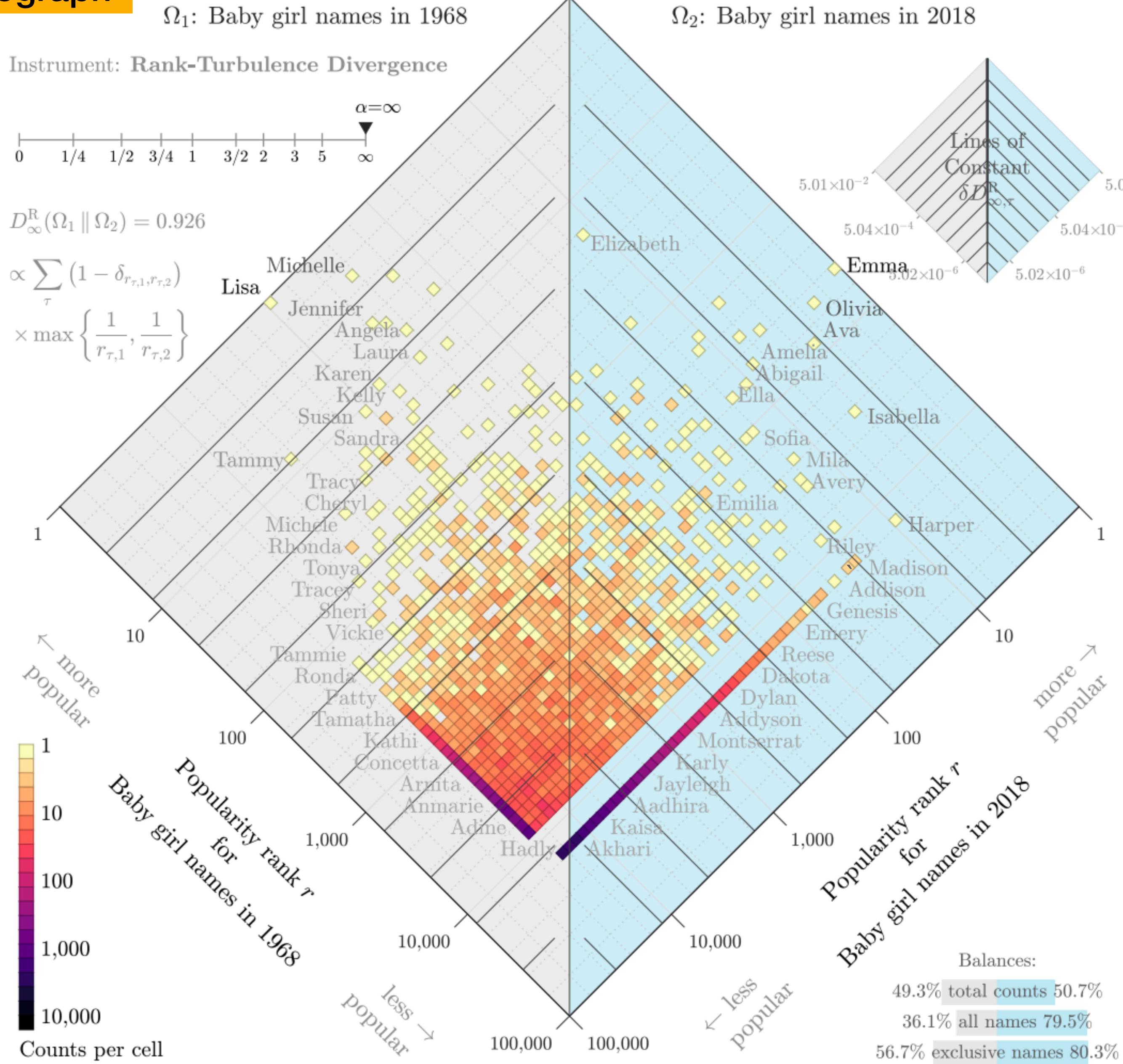
D. Disjoint systems:



Balances:

59.9% total counts	40.1%
63.2% all words	61.6%
60.8% exclusive words	59.8%

Allotaxonograph



Using Scattertext Part 1

- What can we learn from how people tweet about Moderna and Pfizer?
- Jupyter notebook demo
 - <https://nbviewer.jupyter.org/github/JasonKessler/IndieThinkersTutorial/blob/main/Moderna%20vs.%20Pfizer%20Analysis.ipynb>
- This includes enough code for you to easily replicate the study
 - And explore your own corpora

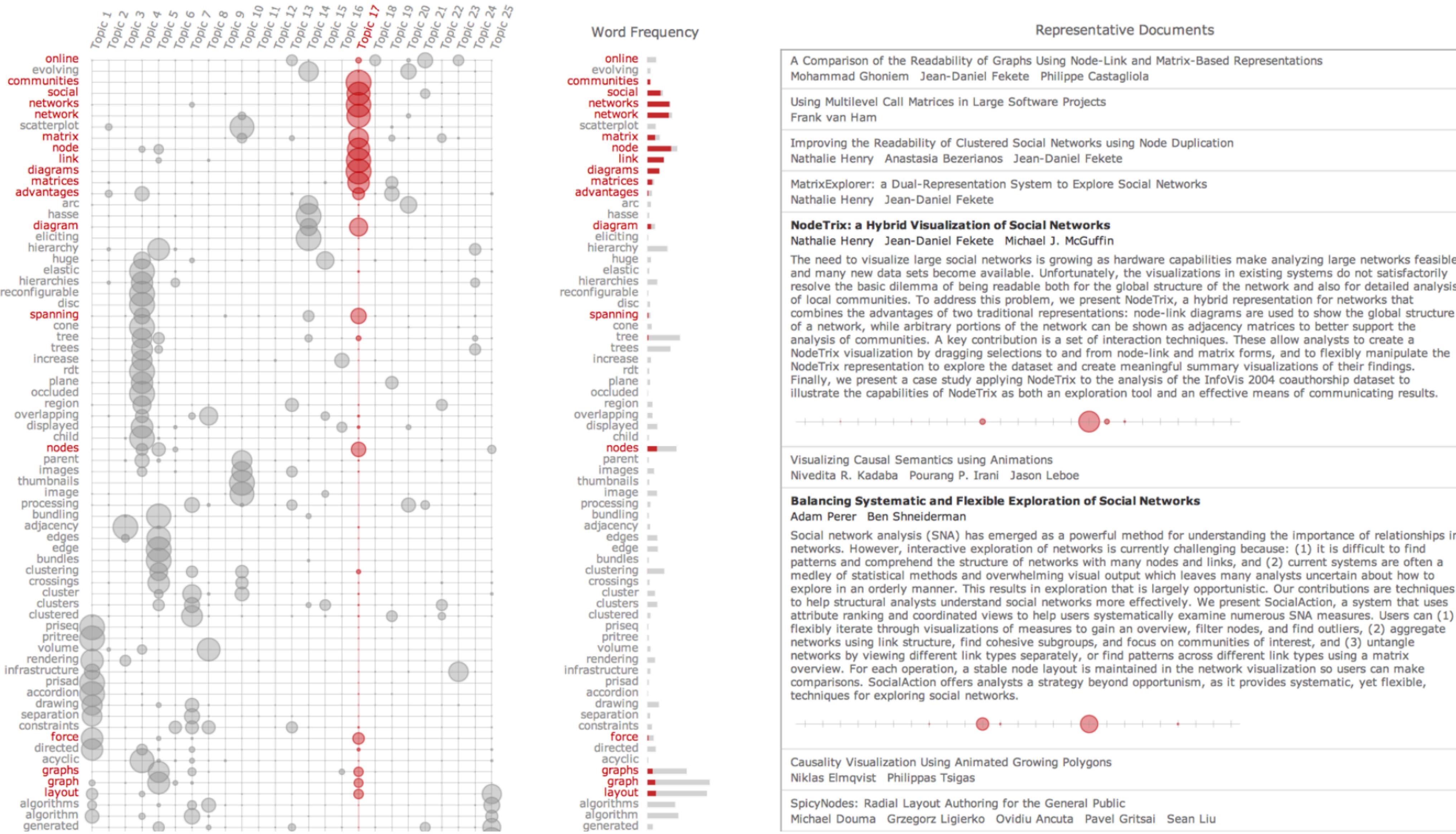
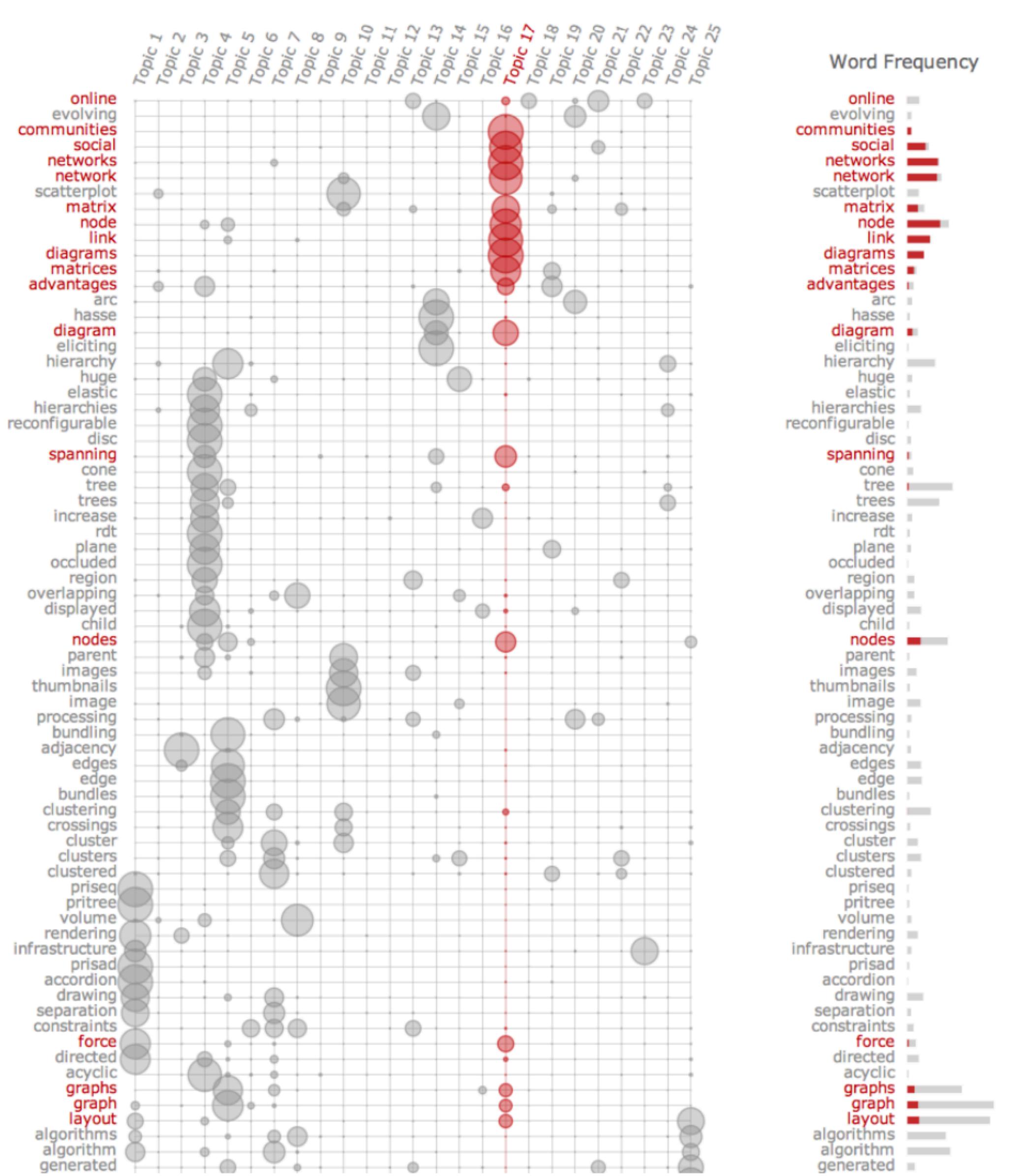


Figure 3: The Termite system. When a topic is selected in the term-topic matrix (left), the systems visualizes the word frequency distribution relative to the full corpus (middle) and shows the most representative documents (right).



Words are included based on “saliency”

- 10-250 terms shown; monitor dependent
- Distinctiveness
 - Sum $KL\text{-divergence}(P(\text{topic}|\text{word})||P(\text{topic}))$

$$distinctiveness(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

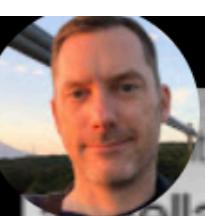
$$saliency(w) = P(w) \times distinctiveness(w)$$

KLD is used for circle sizes

Ordered using “seriation”

- If word B closely follows word A in the list, “A B” is likely to be a phrase
- Neighboring words tend to contain phrases in order
- Online, evolving, communities, social, networks, network
- G^2 for collocation; Bond Energy Algorithm for ordering

Figure 3: The Termite system. When a topic is selected its word frequency distribution relative to the full corpus (middle)



Parallel Tag Clouds to Explore and Analyze Faceted Text Corpora



Visual Encoding

27

- One column per *Circuit*
- Top N words only (scrollable)
- Text Size → significance score
- Order → alphabetic



Significance score:

Based on G^2 :

- Similar to a chi-square contingency test

Size indicates significance score, top K most significant terms are included

Order is alphabetical, which keeps the same tags at similar horizontal levels

- Position changes could falsely imply difference in significance

Difference in position is more noticeable than difference in size

Can't directly indicate change in freq.

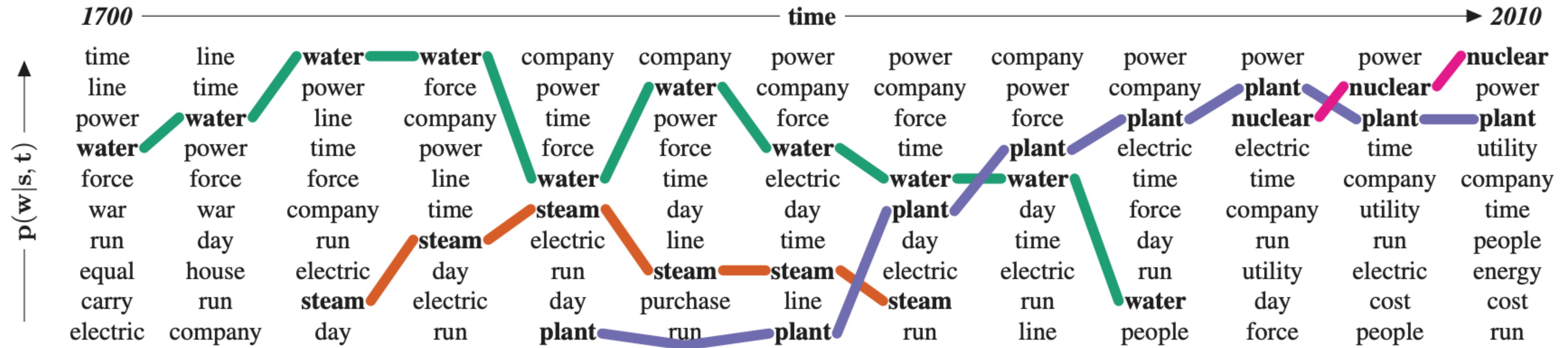


Figure 5: Sense-internal temporal dynamics for the *energy* sense of the word *power* (violet/6 in Figure 4). Columns show the ten most highly associated words for each time interval for the period between 1700 and 2010 (ordered by decreasing probability). We highlight how four terms characteristic of the sense develop over time (see {water, steam, plant, nuclear} in the figure).



Figure 1: Two-dimensional visualization of semantic change in English using SGNS vectors.² **a**, The word *gay* shifted from meaning “cheerful” or “frolicsome” to referring to homosexuality. **b**, In the early 20th century *broadcast* referred to “casting out seeds”; with the rise of television and radio its meaning shifted to “transmitting signals”. **c**, *Awful* underwent a process of pejoration, as it shifted from meaning “full of awe” to meaning “terrible or appalling” (Simpson et al., 1989).

Using Scattertext Part 2

- What can we learn about the evolution of the terms “woke” and “BIPOC” on Twitter?
 - Early demo of Scattertext’s newest feature: diachronic tables
 - Inspired by Parallel Tag Clouds
 - Columns are sorted by l1-regression coefficients, sizes are proportional to frequency
 - Also includes residual dispersion scatterplot. Novel measure of dispersion.
- Jupyter notebook demo
 - [https://nbviewer.jupyter.org/github/JasonKessler/IndieThinkersTutorial/blob/main/
Evolutions%20of%20Woke%20and%20BIPOC.ipynb](https://nbviewer.jupyter.org/github/JasonKessler/IndieThinkersTutorial/blob/main/Evolutions%20of%20Woke%20and%20BIPOC.ipynb)



SparkClouds

Users studies show this outperforms Parallel Tag Clouds

Terms not used in epoch being examined are blurred. Size is relative to frequency at epoch.

Issues: number of terms is limited; most frequent terms aren't always the most informative

New York Times Word Usage Frequency (1970-2018)



<https://twitter.com/DavidRozado/status/1134041329292460032>