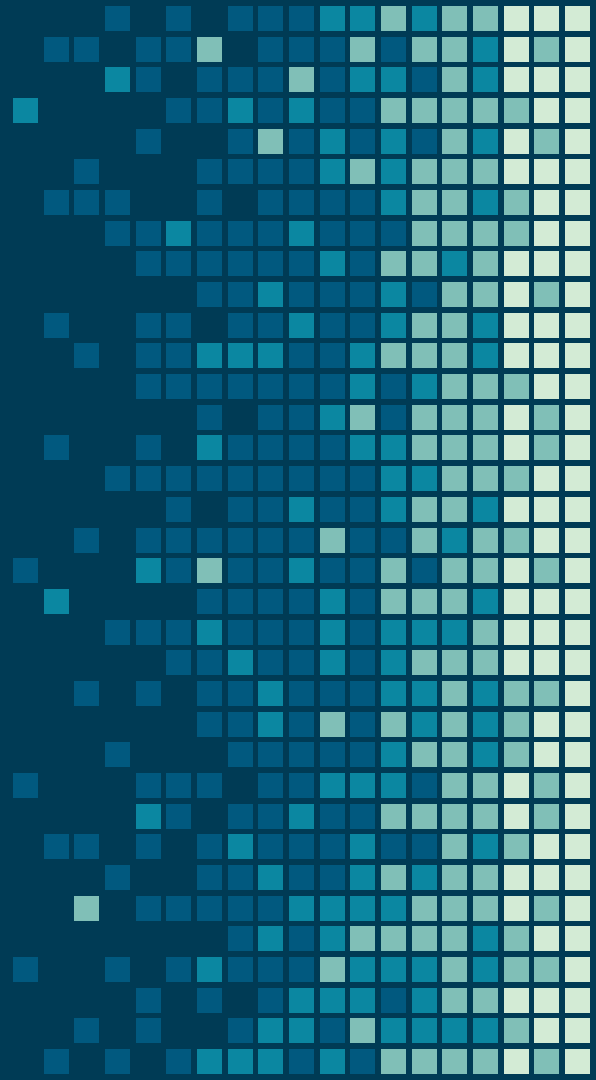# Stylometry with R

Part 1. Textual Forces

Joanna Byszuk, Artjoms Šeļa and Maciej Eder

# What is stylometry?

- A sub-field of computational text analysis that studies **differences** between texts
- Lutosławski 1897: method of "measuring stylistic affinities"
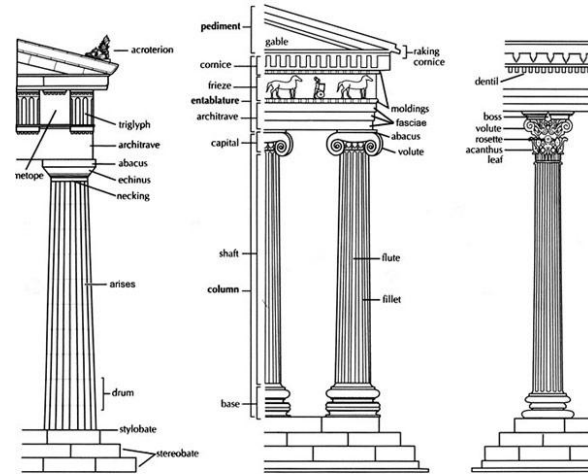
# What is stylometry?

- A sub-field of computational text analysis that studies **differences** between texts
- Lutosławski 1897: method of "measuring stylistic affinities"

Don't mix up with another "stylometry": "the art of measuring columns"!

**Stylometrie', f., stylometry, the art of measuring columns (Säulenmeßkunst).**

# A model of text

- Differences between texts can be expressed in a multitude of ways
- Central question is **how to represent** a text so it could be placed on a quantitative scale? I.e. how to **model** it?
- Short answer: all representations are 'wrong', but some are useful (or more useful than others)
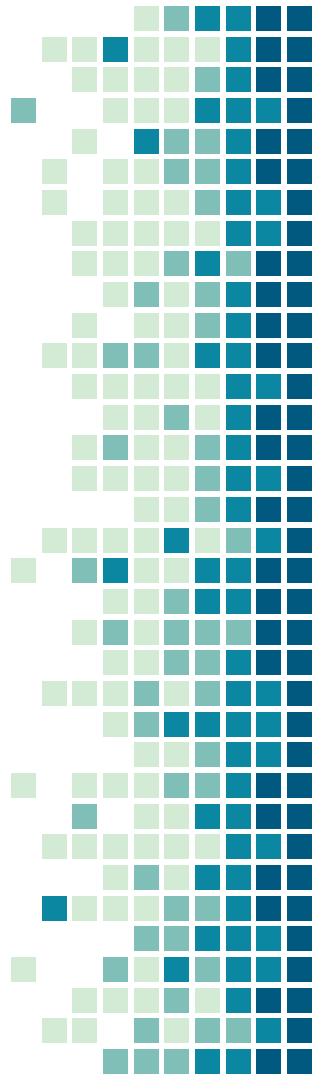
# A model of text

- Differences between texts can be expressed in a multitude of ways
- Central question is **how to represent** a text so it could be placed on a quantitative scale? I.e. how to **model** it?
- Short answer: all representations are 'wrong', but some are useful (or more useful than others)
  - Word frequencies?
  - Algorithmically inferred topics?
  - Part-of-Speech tags?
  - Networks of character connections?
  - Embeddings?
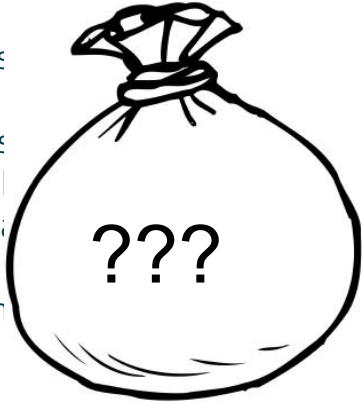  - Sentiment scores?
  - ....

# Silly things: bags of words

Mr. Sherlock Holmes, who was usually very late in the mornings, save upon those not infrequent occasions when he was up all night, was seated at the breakfast table. I stood upon the hearth-rug and picked up the stick which our visitor had left behind him the night before. It was a fine, thick piece of wood, bulbous-headed, of the sort which is known as a "Penang lawyer." Just under the head was a broad silver band nearly an inch across. "To James Mortimer, M.R.C.S., from his friends of the C.C.H.," was engraved upon it, with the date "1884." It was just such a stick as the old-fashioned family practitioner used to carry – dignified, solid, and reassuring.
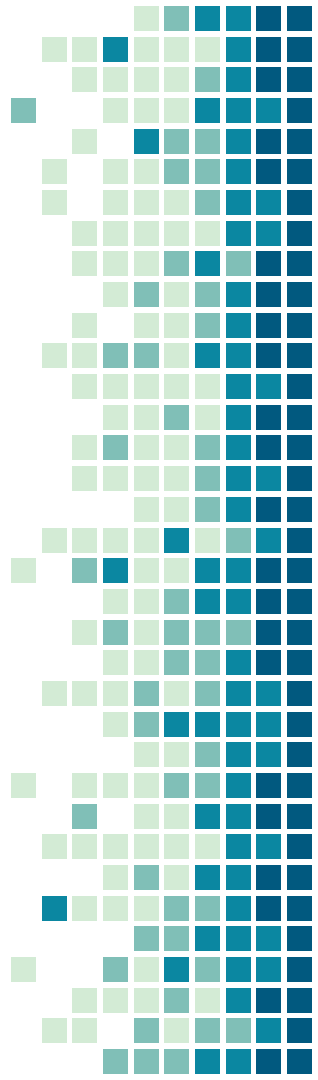
# Silly things: bags of words

Mr. Sherlock Holmes, who was [...] mornings, save upon those not infrequent occasions when he was up all [...] breakfast table. I stood upon the hearth-rug and picked up the s[...] left behind him the night before. It was a fine, thick piece of wood, bul[...] which is known as a "Penang lawyer." Just under the head was a broa[...] ch across. "To James Mortimer, M.R.C.S., from his friends of the C.C.H.," [...] th the date "1884." It was just such a stick as the old-fashioned fam[...] arry – dignified, solid, and reassuring.

???

# Silly things: bags of words

Mr. Sherlock Holmes, who was [...] mornings, save upon those not infrequent occasions when he was up all [...] breakfast table. I stood upon the hearth-rug and picked up the s[...] left behind him the night before. It was a fine, thick piece of wood, bul[...] which is known as a "Penang lawyer." Just under the head was a broa[...] nch across. "To James Mortimer, M.R.C.S., from his friends of the C.C.H.," [...] th the date "1884." It was just such a stick as the old-fashioned fam[...] arry – dignified, solid, and reassuring.

???

the: 10
was: 7
a: 4
it: 3
of: 3
upon: 3
and: 2
as:2

…

# Silly things: bags of words

Mr. Sherlock Holmes, who was ⬚⬚⬚⬚⬚⬚ mornings, save upon those not infrequent occasions when he was up all ⬚⬚⬚⬚⬚⬚ breakfast table. I stood upon the hearth-rug and picked up the s⬚⬚⬚⬚⬚ left behind him the night before. It was a fine, thick piece of wood, bul⬚⬚⬚⬚ which is known as a "Penang lawyer." Just under the head was a broa⬚⬚⬚⬚ch across. "To James Mortimer, M.R.C.S., from his friends of the C.C.H.," ⬚⬚⬚⬚th the date "1884." It was just such a stick as the old-fashioned fam⬚⬚⬚⬚rry – dignified, solid, and reassuring.

**???**

Doyle_Baskerville_p1: (10, 7, 4, 3, 3, 3, 2, 2,…)

# Silly, but useful!

# Silly, but useful!



grey waves mage child cold learned gont port wood house jasper land learned true east found spoke night north rain sleep left hills sail lad run fell till power set dark wise boy red fire south rose evil sun shadow lay spell sea light lost dry ran sky sat low town ogion seas air court day hold days wind boat isle heart master time stood dead staff fear looked door names eyes past serret stone lord white sound voice mountain tower heard raised friend



in the Earthsea Trilogy
by the winner of the
Hugo and Nebula awards
Ursula K. Le Guin
A Wizard
of Earthsea

# Silly, but useful!

# Silly, but useful!

# Silly, but useful!
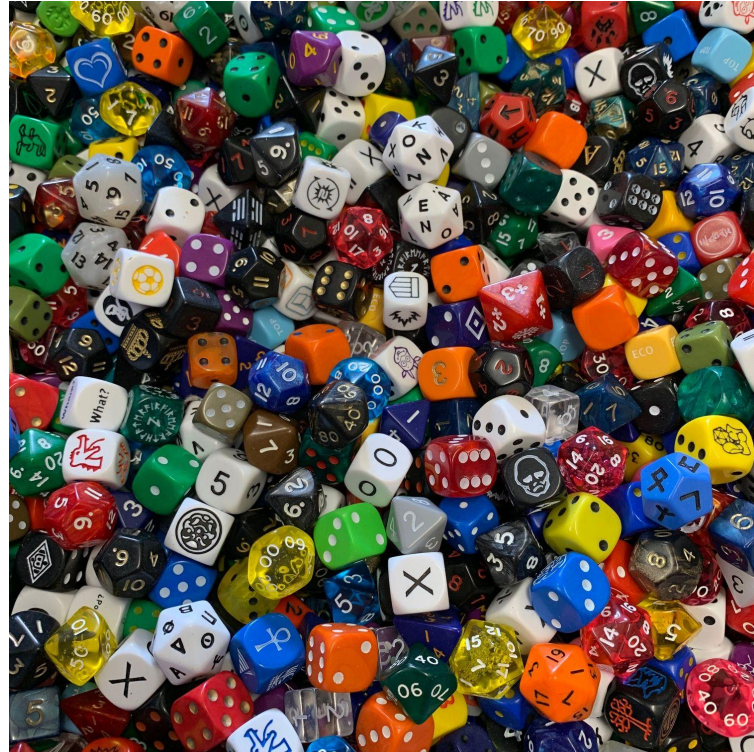
# Silly, but useful!

# Proxies

- Word frequencies may serve as a proxy to **things we care about** in texts

- Word frequencies are the result of word choice -> word choice is a result of of **forces that organize texts**

# Proxies

- Imagine LOTS of colorful dice of different shapes
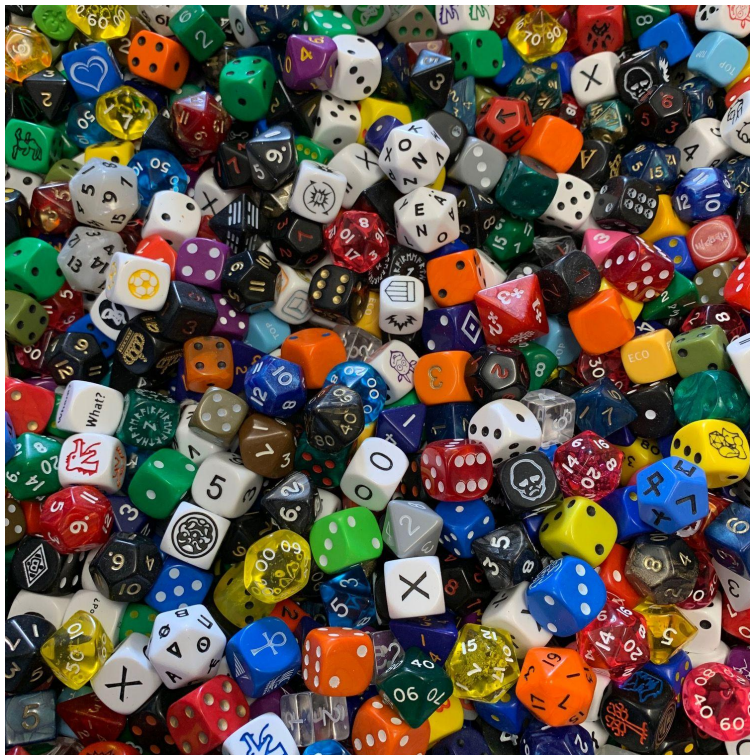- 1 text = 1 complex multidimensional die

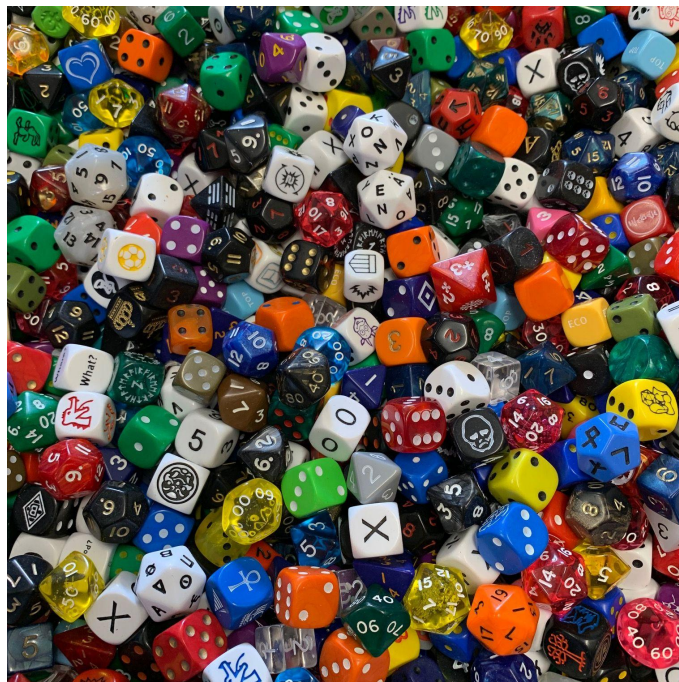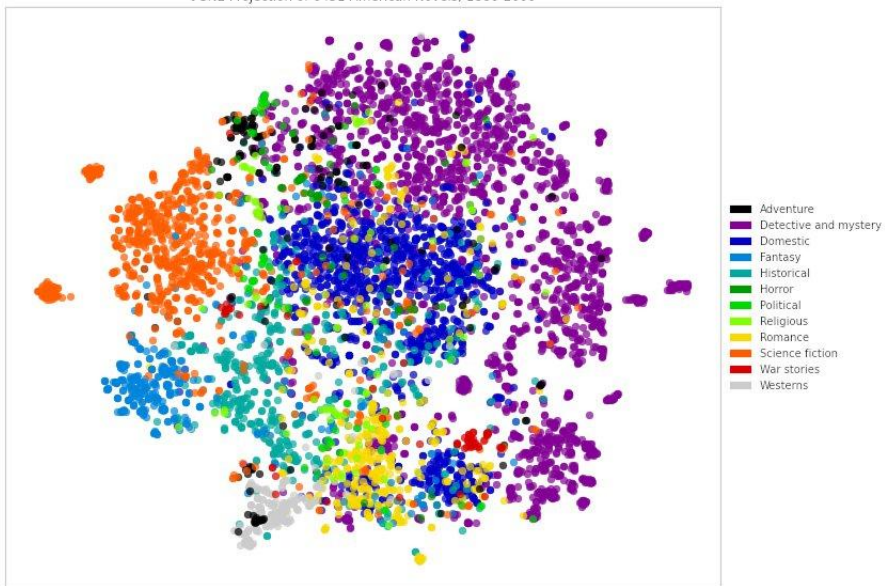Let's toss them into some abstract space!

# Proxies

Normal fair dice:
**roll randomly all over the place, nice!**

Text-dice: **OH NO...**

# Text–dice rolling will be organized by multitude of forces



t-SNE Projection of 6431 American Novels, 1880-2000

- Adventure
- Detective and mystery
- Domestic
- Fantasy
- Historical
- Horror
- Political
- Religious
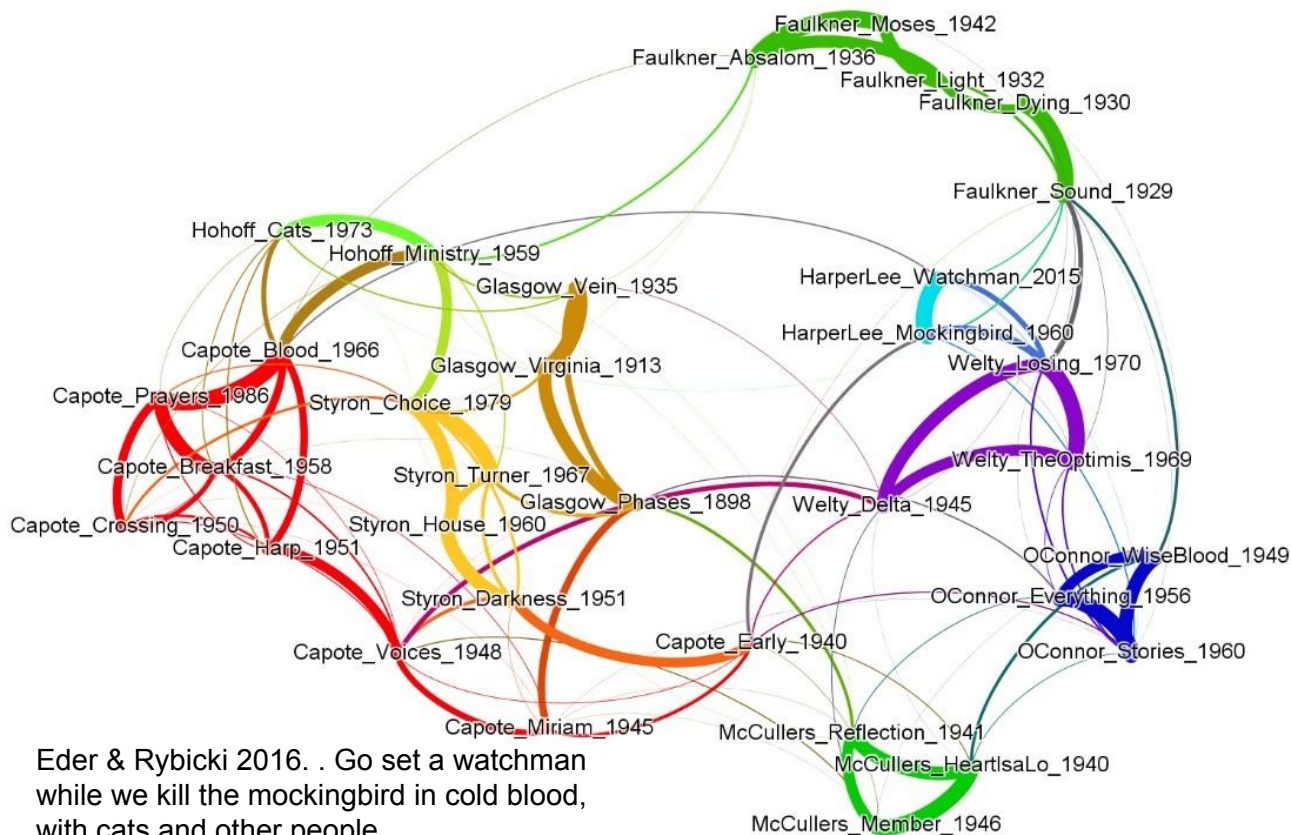- Romance
- Science fiction
- War stories
- Westerns

# Curse & blessing of word frequencies

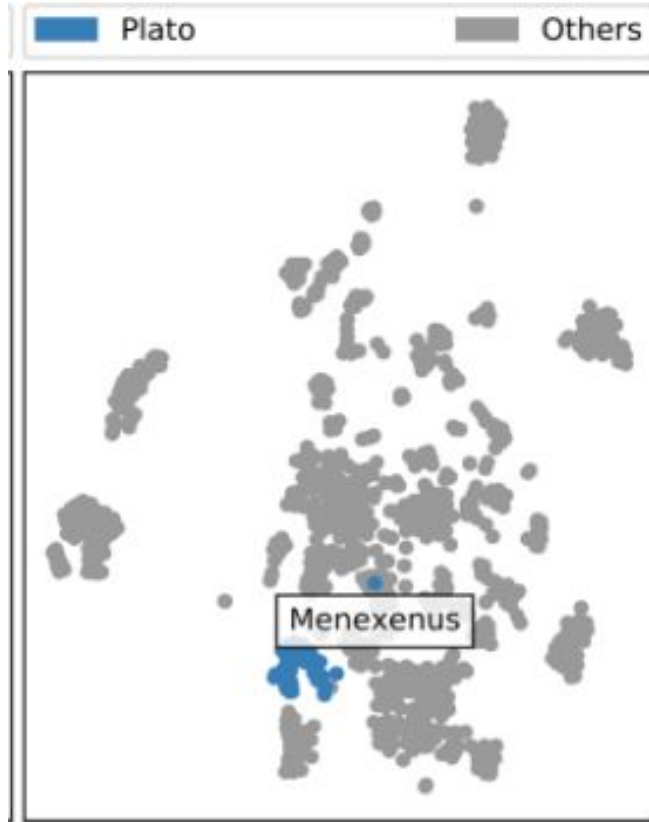Forces that are "naturally" captured by word frequencies in a large corpus:

- **1. Authorship.** Two texts of the same author usually appear the closest to each other than to any outsider text
- **2. Modes of writing.** Fiction and nonfiction grew apart stylistically; Poetry books (esp. regularized verse) will always form a VERY exclusive party with themselves.
- **3. Genre.** Well-formed fiction genres (detective/mystery, sci-fi,
- **4. Chronology.** Global language change: Each generation of writers adopt slightly different version of language than the previous one (also: spelling conventions)
- **5. Gender.** Be aware of a historical gap between women and men writing (socially constructed)
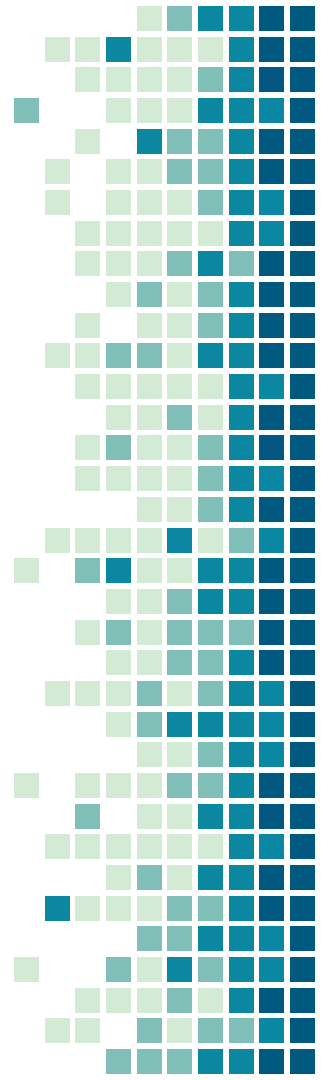- …

# Authorship



Eder & Rybicki 2016. . Go set a watchman while we kill the mockingbird in cold blood, with cats and other people.
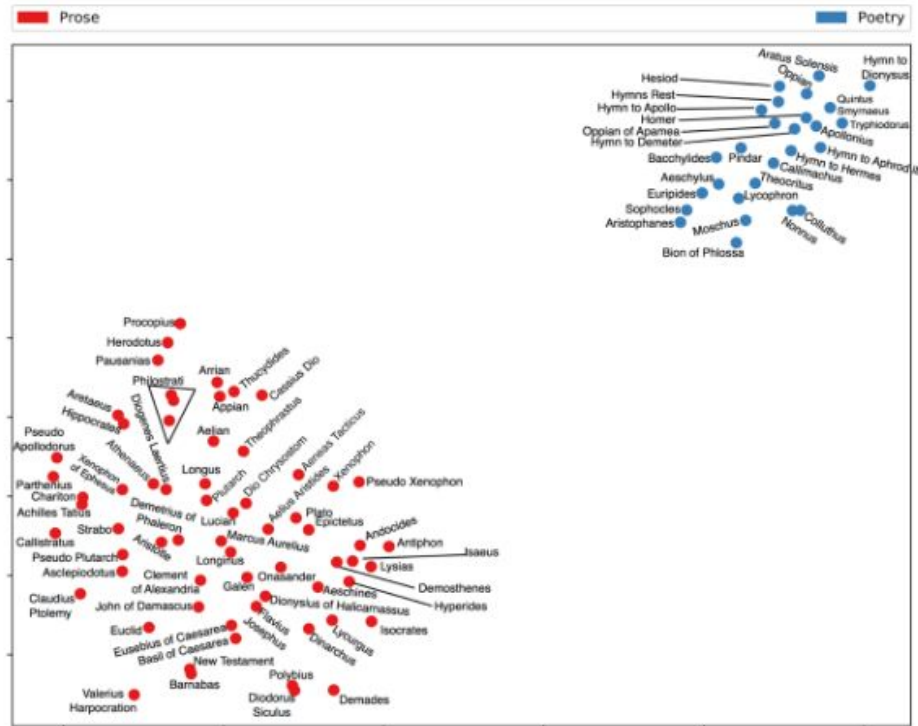
Plato, Others (legend)

Menexenus

- Plato's texts with *Menexenus* as an outlier

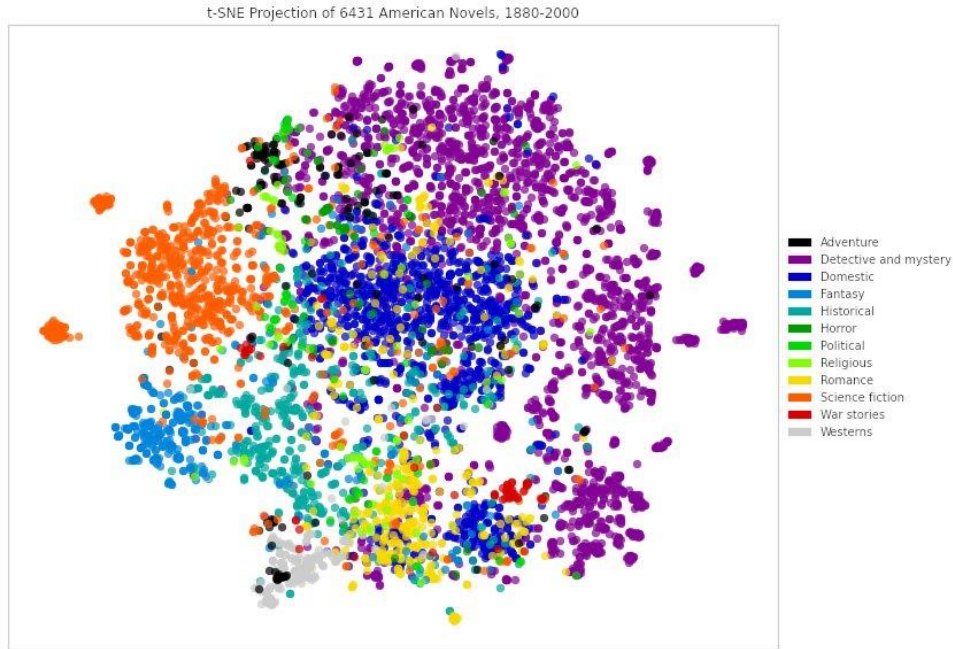Storey & Mimno 2020: Like Two Pis in a Pod: Author Similarity Across Time in the Ancient Greek Corpus
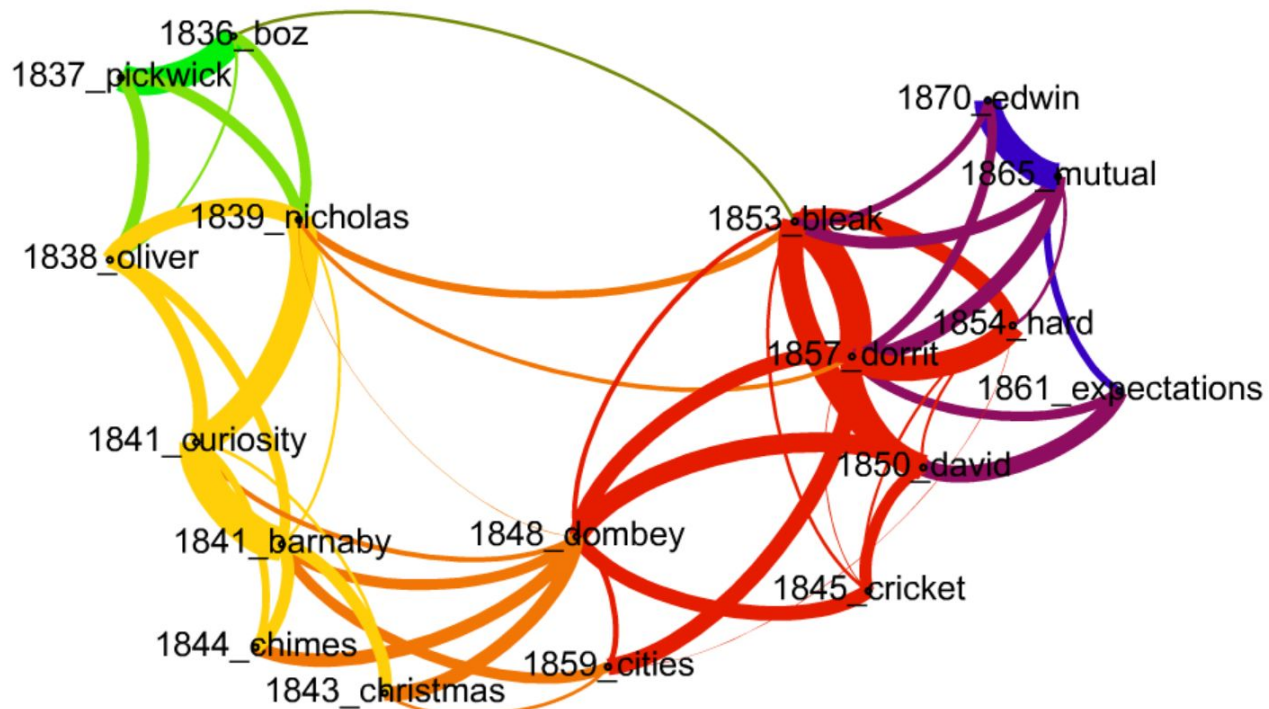
# Poetry parties with poetry



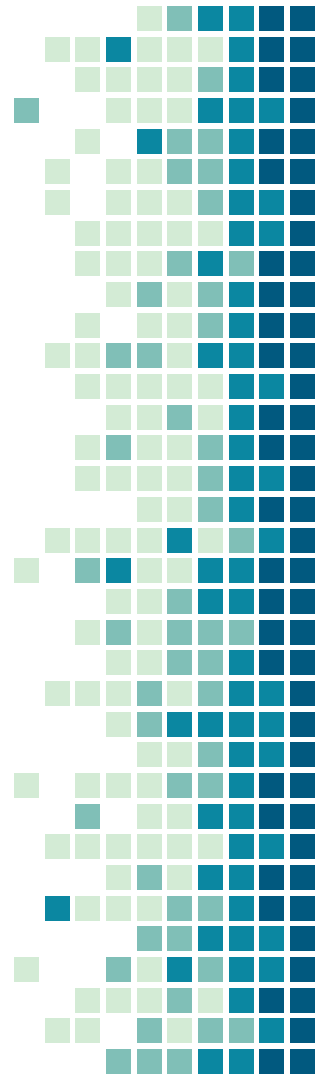Storey & Mimno 2020: Like Two Pis in a Pod: Author Similarity Across Time in the Ancient Greek Corpus
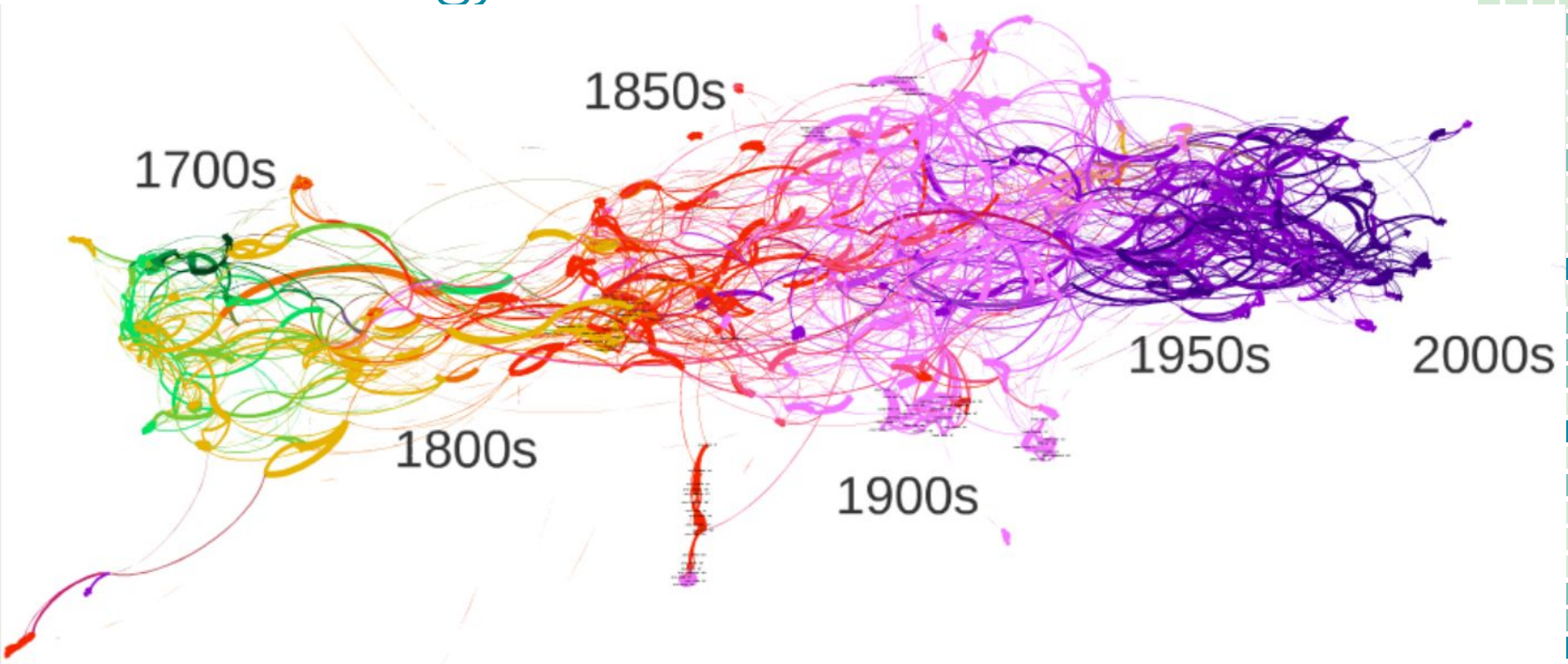
# Genres of fiction



t-SNE Projection of 6431 American Novels, 1880-2000

Adventure
Detective and mystery
Domestic
Fantasy
Historical
Horror
Political
Religious
Romance
Science fiction
War stories
Westerns

Jordan Pruett, *Twitter,* 2020
@pruett_jordan

Rybicky 2016

# Chronology: Global



1700s 1800s 1850s 1900s 1950s 2000s

Rybicky 2016

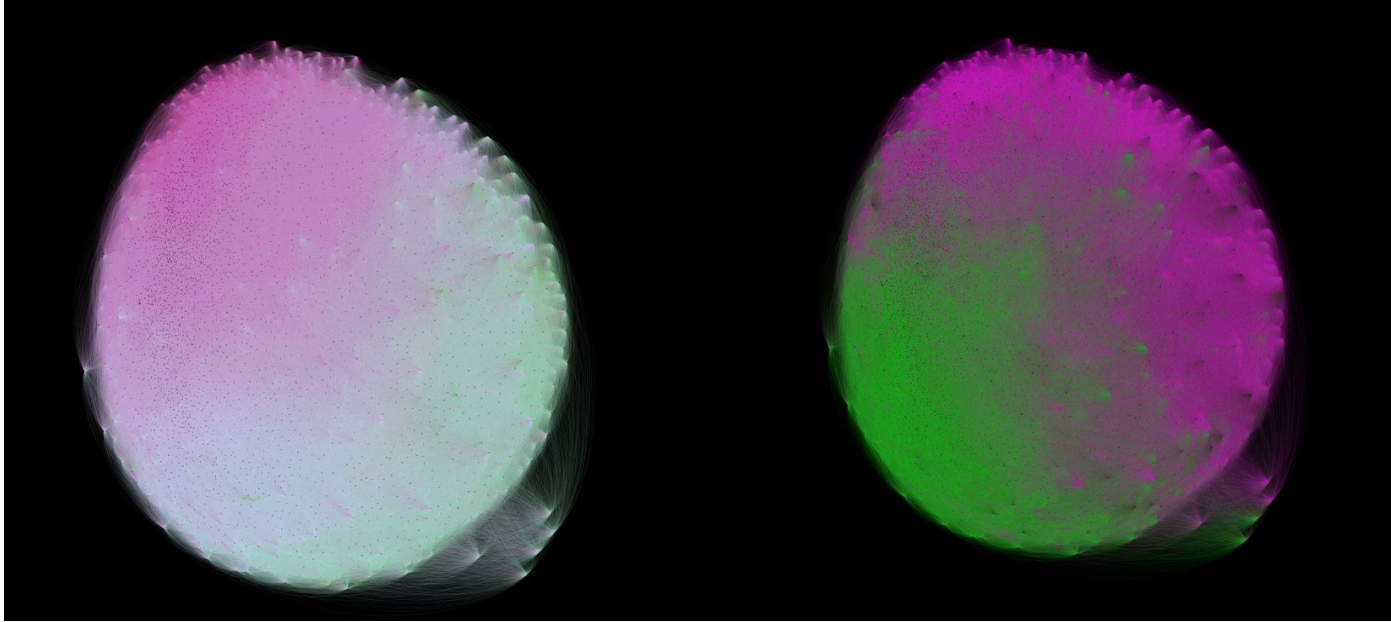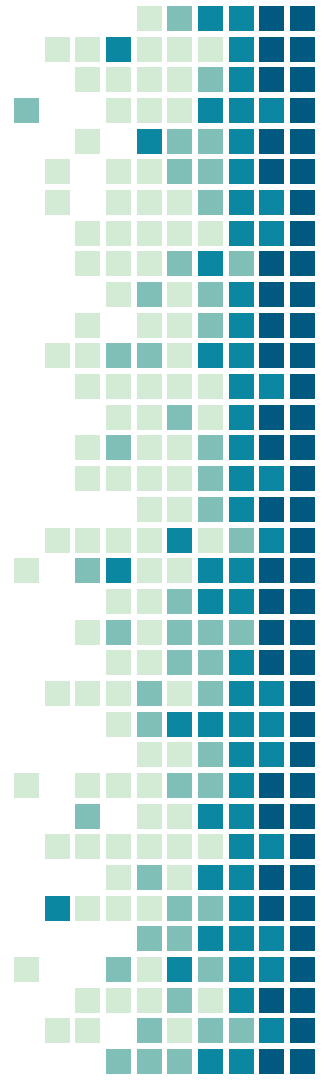# Chronology: Global



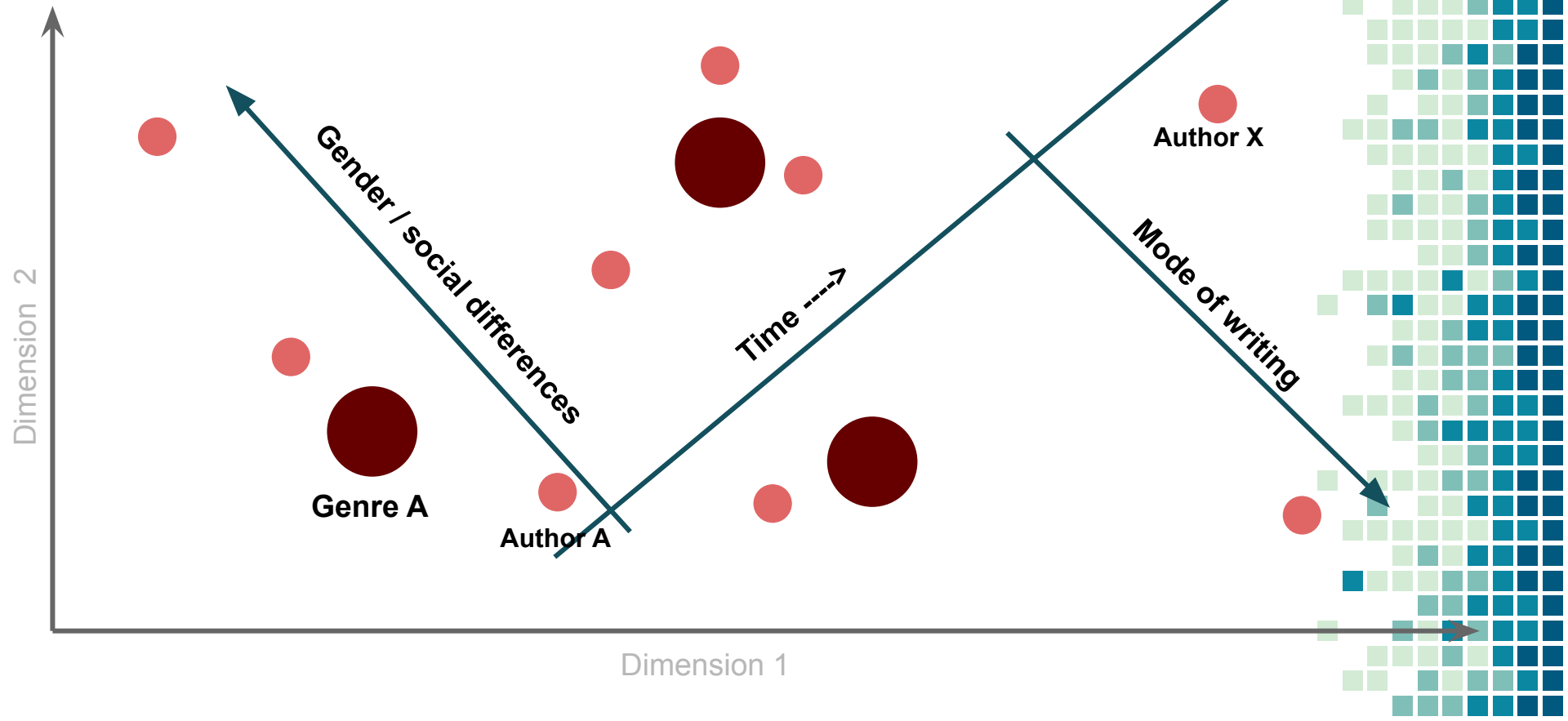Time                                              Gender

Jockers 2013 *Macroanalysis*

# Mapping the force

Thank you!