# Language technology in the service of the humanities

Eetu Mäkelä
Professor (tenure track) in Human Sciences–Computing Interaction
University of Helsinki
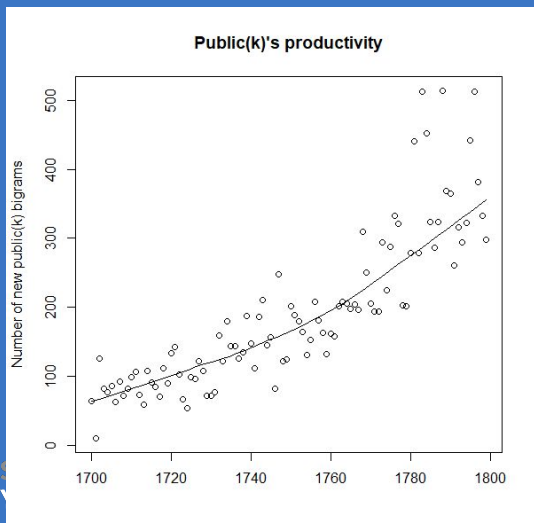
# Analysing public communication in 18th century Britain and 19th century Finland

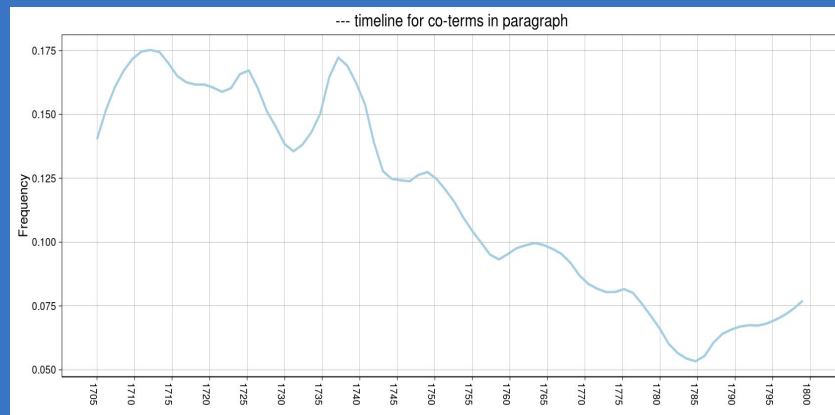with the COMHIS research group, University of Helsinki

# Case public(k)

- If a public sphere formed in Britain in the eighteenth century, can we see this in how uses of the word "public" change? How can we historically interpret the changes that we see?
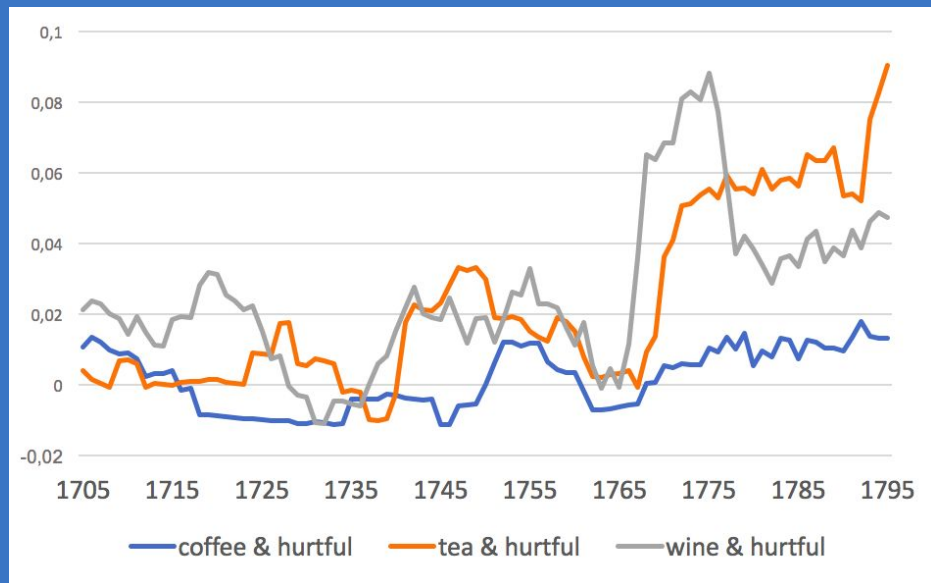


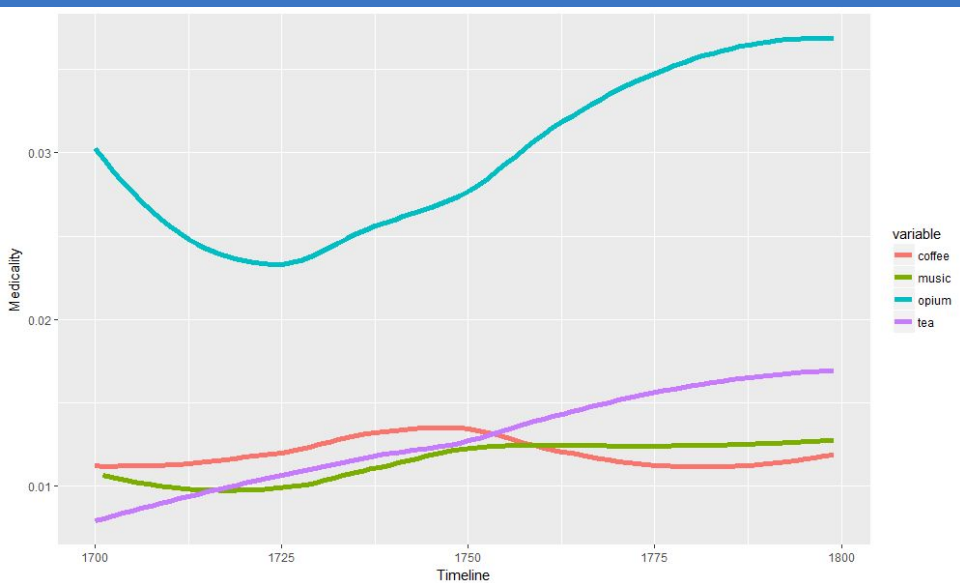**Public(k)'s productivity**



Public + religion

--- timeline for co-terms in paragraph

# Tracking health aspects of consumables through the 18th Century

# Sociohistorical Language Change

with Tanja Säily and Mika Hämäläinen, University of Helsinki

# Questions

- Who creates and adopts new vocabulary in the history of English?
  - Are there differences across social groups?
  - How does new language arise?



| Social rank | # of new words |
|---|---|
| Clergy, lower | 7 |
| Genry, lower | 7 |
| Gentry, upper | 7 |
| Professional | 7 |
| Nobility | 7 |
| Royalty | 6 |
| Clergy, upper | 1 |
| Merchant | 0 |
| Other non-gentry | 0 |

# Etymologies of the new words found

| Etymology type | # of new words |
|---|---|
| Derivative | 19 |
| Borrowing | 17 |
| Compound | 2 |
| Uncertain | 1 |
| Zero derivation | 1 |
| Borrowing/hybrid | 1 |
| Borrowing/zero derivation | 1 |

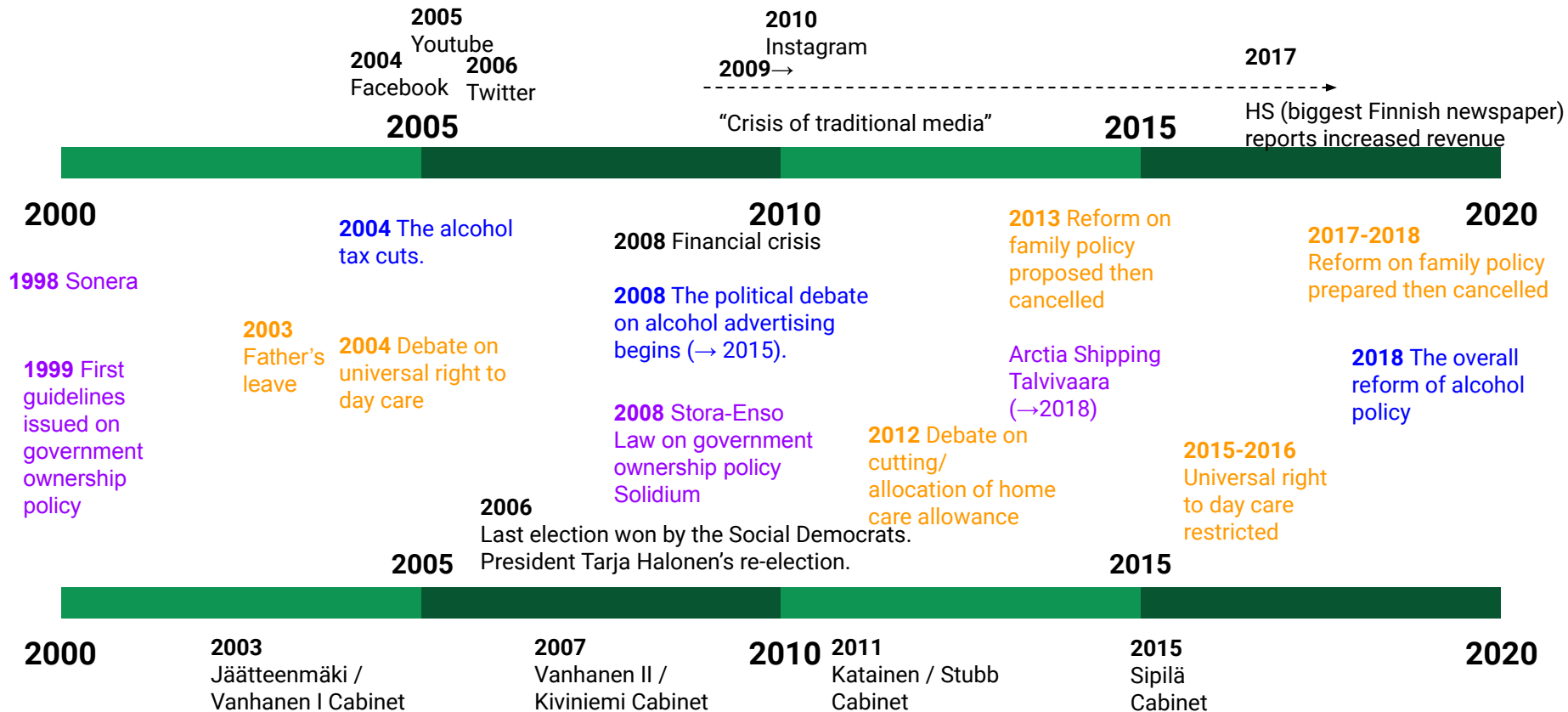| Etymon language | # of new words |
|---|---|
| English | 22 |
| Latin | 11 |
| undetermined | 2 |
| Italian | 2 |
| French | 2 |
| German | 1 |
| Dutch | 1 |

# FLOWS OF POWER

**media as site and agent of politics**

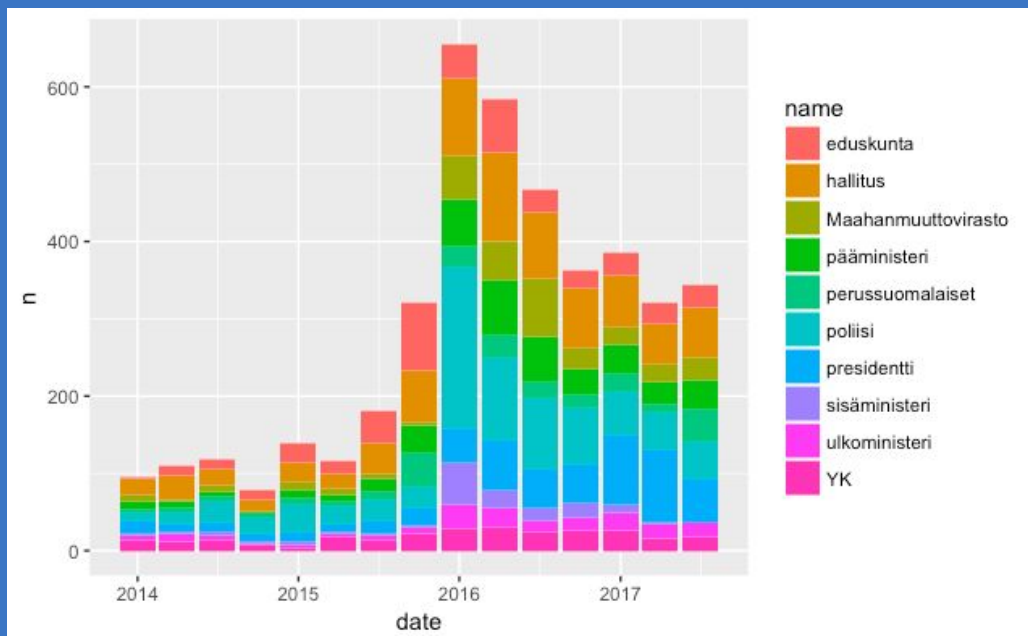Anu Koivunen (Tampere University)
Eetu Mäkelä (University of Helsinki)

Tampere University

HELSINGIN YLIOPISTO
HELSINGFORS UNIVERSITET
UNIVERSITY OF HELSINKI

**2004** Facebook

**2005** Youtube

**2006** Twitter

**2005**

**2010** Instagram

**2009→** "Crisis of traditional media"

**2017** HS (biggest Finnish newspaper) reports increased revenue

**2015**

**2000**

**1998** Sonera

**1999** First guidelines issued on government ownership policy

**2003** Father's leave

**2004** The alcohol tax cuts.

**2004** Debate on universal right to day care

**2005**

**2008** Financial crisis

**2008** The political debate on alcohol advertising begins (→ 2015).

**2008** Stora-Enso Law on government ownership policy Solidium

**2010**

**2006** Last election won by the Social Democrats. President Tarja Halonen's re-election.

**2012** Debate on cutting/ allocation of home care allowance

**2013** Reform on family policy proposed then cancelled

Arctia Shipping Talvivaara (→2018)

**2015-2016** Universal right to day care restricted

**2015**

**2017-2018** Reform on family policy prepared then cancelled

**2018** The overall reform of alcohol policy

**2020**

**2000**

**2003** Jäätteenmäki / Vanhanen I Cabinet

**2005**

**2007** Vanhanen II / Kiviniemi Cabinet

**2010**

**2011** Katainen / Stubb Cabinet
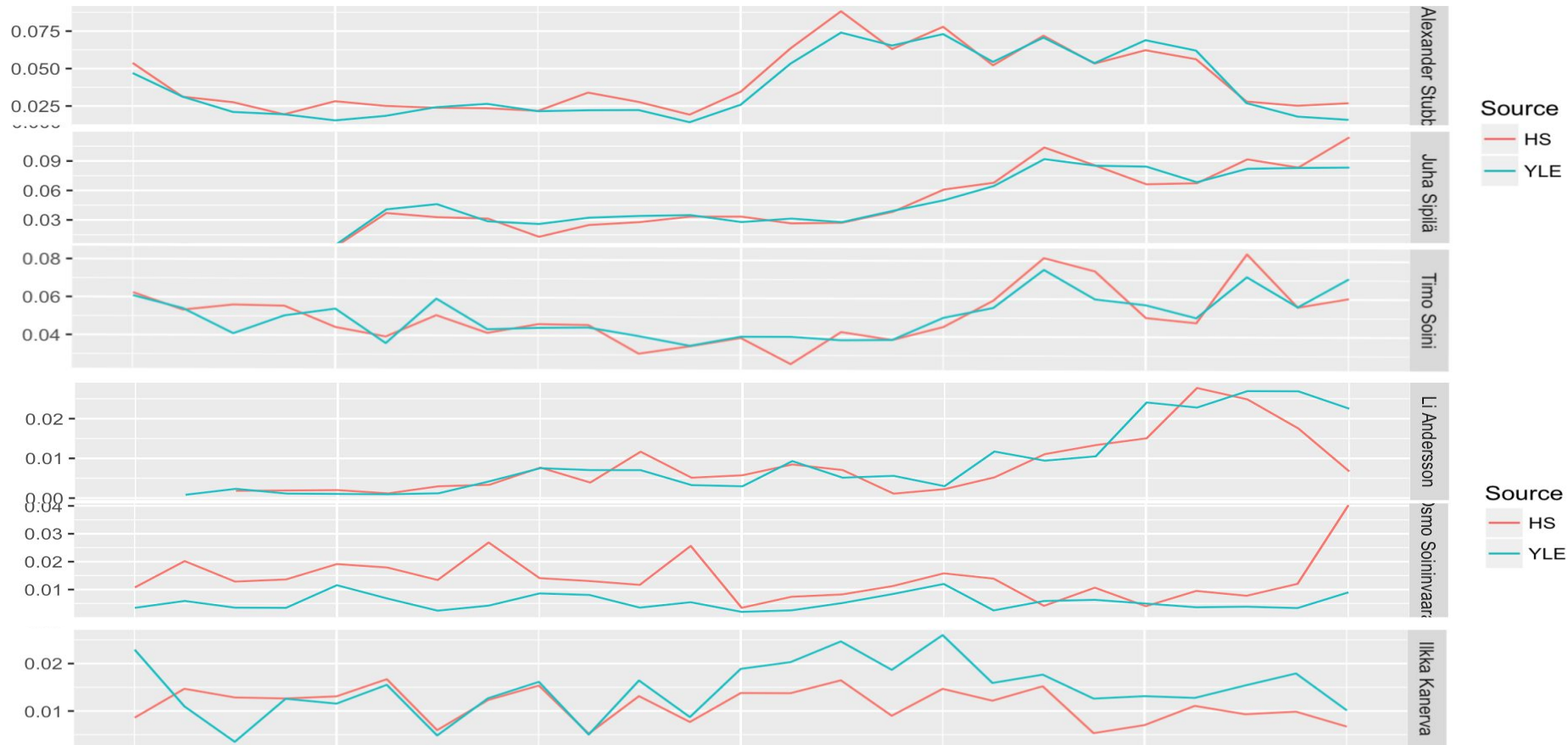
**2015** Sipilä Cabinet

**2015**

**2020**

# Authorities referred to in the debate on immigration in YLE by year

# Media trajectories of prominent Finnish politicians based on proportion of mentions

# Words associated with Ilkka Kanerva in YLE

# Groups of people often appearing together in association with Ilkka Kanerva in YLE

- Matti Louekoski, Matti Vanhala, Sinikka Salo
- Arto Merisalo, Jorma Kalske, Tapani Yli-Saunamäki, Toivo Sukari
- Ari Ruotsalainen, Arto Merisalo, Kyösti Kakkonen, Tapani Yli-Saunamäki, Toivo Sukari

- Anne-Mari Virolainen, Esko Kiviranta, Heli Paasio, Ilkka Kantola, Katja Taimela, Pertti Hemmilä, Petteri Orpo, Stefan Wallin, Ville Niinistö
- Anne-Mari Virolainen, Annika Lapintie, Esko Kiviranta, Heli Paasio, Ilkka Kantola, Katja Taimela, Pertti Hemmilä, Petteri Orpo, Stefan Wallin, Ville Niinistö

# Adjectives associated with Juha Sipilä in YLE

- suuri
- uusi
- vahva
- hurja
- hyvä
- rehti
- avoin

- rauhallinen
- tervetullut
- lupsakka
- uskonnollinen
- raikas poikkeus
- nopea

"Kirjoittaja vertaa Juha Sipilän taivalta heti alussa Julius Caesarin rakettimaiseen voittokulkuun."

# Intertextuality and thematic networks in Finnic regional poetic cultures

with Kati Kallio, Finnish Literature Society

# Research questions

- Tracing multiple levels of commonalities across varying Finnish and Estonian corpora:
    - Characters / individual words
    - Verses
    - Formulae (small repeating bits, e.g. sirkun siivet/linnun siivet/tedre tiivad)
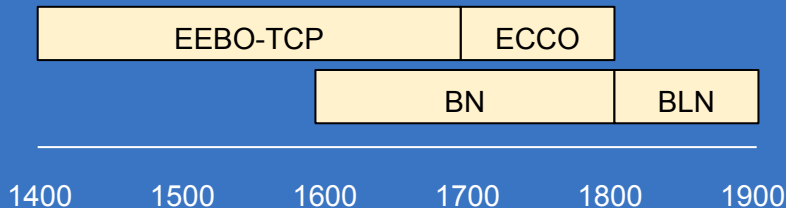    - Themes

# Needs from language technology

# Materials

- Eighteenth Century Collections Online: ~210,000 books (~84% of those known to exist), ~10 billion running words 1700–1800 (ECCO)
- Early English Books Online: ~60,000 books (~53% of those known to exist), ~1.4 billion running words 1400–1700 (EEBO-TCP)
- British Library Newspapers: ~655,000 newspaper issues (uneven selection), ~40 billion running words 1800–1900 (BLN)
- Burney & Nichols Collections: ~230,000 pamphlets and newspapers, ~2.7 billion running words before 1600–1800 (BN)

| EEBO-TCP | ECCO |
|:---:|:---:|

| BN | BLN |
|:---:|:---:|

1400    1500    1600    1700    1800    1900

# Data production pipelines: Early English Books Online


Microfilms


Physical books


Electronic Image Scans


Early English Books Online


OCR

# But what's in there, e.g. genre-wise?



20-30.000 novels published

6.000 novels still extant

3.500 novels available in full text
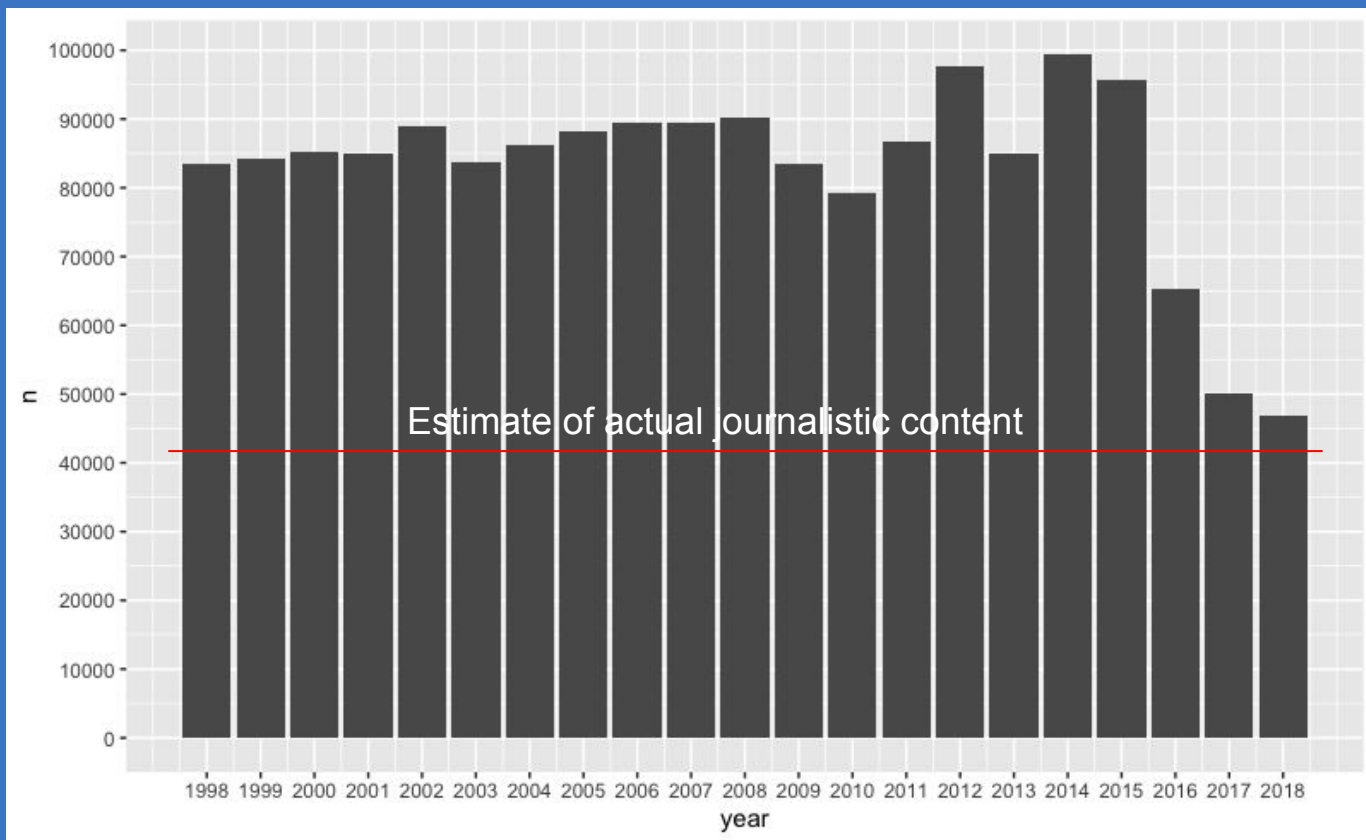
EEBO-TCP and ESTC Text Counts

# Number of articles per year in HS

# Surfacing distributional anomalies

# Relatively good OCR



'Tis not Saint *George* we Sing of here,
Nor *George*, the fatal Duke *Villier*;
Nor *George a Green*, nor Castriot,
Nor *Buchanan* the learned *Scot*;
But 'tis of *George* the Valiant *Monck*,
That made *Van-Trump* in's Blood dead-
And in the Seas his Navy sunck.   (drunk
 *Oh!  this is our brave George!*

lsw-not- Saint George we Sing of here,
Nor George, the fatal Duke Villier ;
Nor George a Green, nor Castriot,
Nor Buchanan the learned S cot q
But us of George the Valiant Monck,
That made Van-Trump in'S Blood deod
and in theseus his Navy snuck. (drunk,
Ok I this is our brave George !

# THE

## Firſte volume of the
## Chronicles of England, Scot-
## lande, and Irelande.

### CONTEYNING,

The deſcription and Chronicles of England, from the firſt inhabiting vnto the conqueſt

The deſcription and Chronicles of Scotland, from the firſt originall of the Scottes nation, till the yeare ofour Lorde. 1571

The deſcription and Chronicles of Yrelande, likewiſe from the firſte originall of that Nation, vntill the yeare. 1547.

Faithfully gathered and ſet forth, by
Raphaell Holinſhed.

## AT LONDON,
Imprinted for Iohn Harriſon.

~ ~k ~

~ l I ~ li ~]J]O DmU~ov O~1i |

~ ~1l ~ ~ -\O~Si~\r<,St~5,o t%,\~t,\~ ~ ~ ~

~' .-bnEIs~l br~; <~5n~1 ~

.~1 1~t ~3mo71~k<~7noostI3o~rsd ~=i~mlm87il fif ~s ~

~' 3,Ilmo~l.6n3~nm/l7~=io\~ ~7g ~i

....~ -,~. ;lIl~1B~ ]8 ~ . ~ ~ ' ~

'.~`~@~ ~ ~`~pA til Sns t' - b~ ~I\U\ `i:~] ~ ~ ~

I I noin~Hodol~o]bsJni~qml '~1 11

1~.1 ~ ~1 11

"' ~ ~ |'? ' ~ 9~ 9~] boO \~

,,.---. ~13 ~ ~ ~

-: ~___ 1

.

# OCR problems

- vectors.most_similar("king") → [('k~ing', 0.88), ('kiing', 0.86), ('"king", 0.84), ('emperor', 0.823), ('icing', 0.83), ('kring', 0.82), ('kcing', 0.82), ('.king', 0.82), ('kting', 0.81), ('iking', 0.81)]

## Configuration

### Search

Max edit distance 1

Required common prefix length 1
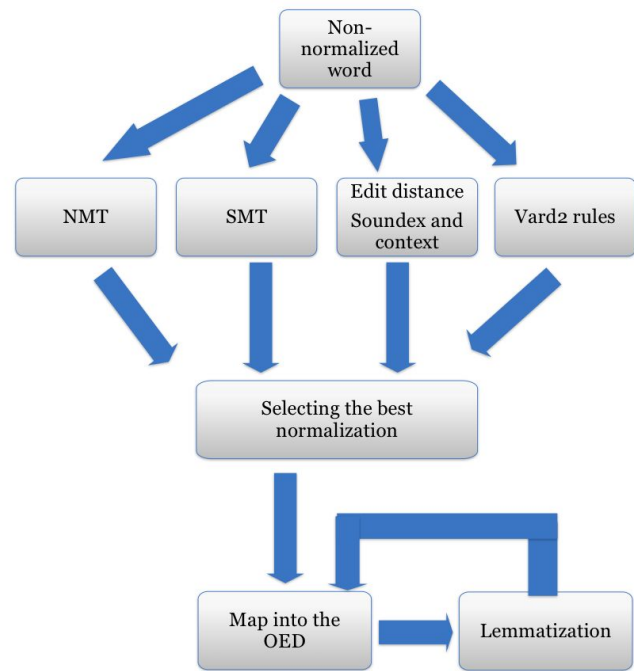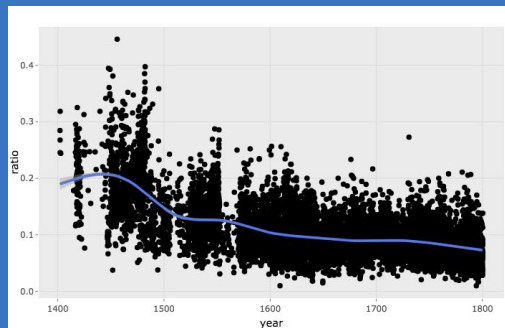
☐ Transposition is a single edit

Query

politeness

## Keywords

☑ politeress (1)  ☑ politeness (19035)  ☑ politzeness (1)
☑ politenehs (4)  ☑ politeess (1)  ☑ politeless (3)
☑ politehess (1)  ☑ politenesa (1)  ☑ politenesj (3)
☑ politeneós (1)  ☑ politen.ess (3)  ☑ politiness (2)
☑ politeneps (3)  ☑ politenees (16)  ☑ politene.s (6)
☑ politeiness (2)  ☑ potiteness (1)  ☑ politeness5 (1)
☑ politenejs (61)  ☑ politene1s (4)  ☑ politene'ss (1)
☑ politenes's (1)  ☑ po'liteness (4)  ☑ politenebs (1)
☑ politenesc (1)  ☑ politenesl (1)  ☑ pcliteness (2)
☑ politness (1)  ☑ polivteness (1)  ☑ poli:eness (1)
☑ politenesss (49)  ☑ politensss (2)  ☑ positeness (2)
☑ politrness (2)  ☑ politenegs (2)  ☑ politeners (585)
☑ politenessa (1)  ☑ politene3s (1)  ☑ politenesr (1)
☑ politenesi (7)  ☑ politcness (4)  ☑ politenes3 (1)
☑ politenezs (1)  ☑ polirteness (1)  ☑ politeneds (7)
☑ politeneos (7)  ☑ polteness (1)  ☑ politeness1 (1)
☑ polileness (1)  ☑ politeiess (1)  ☑ politentess (2)
☑ peliteness (1)  ☑ p.oliteness (6)  ☑ politemess (1)
☑ politeneas (6)  ☑ politencss (3)  ☑ politen'ess (10)
☑ politenels (904)  ☑ politenews (1)  ☑ politenuss (1)
☑ polireness (4)  ☑ políteness (1)  ☑ p'oliteness (9)
☑ politenesk (1)  ☑ politenfss (1)  ☑ potliteness (1)
☑ po.iteness (1)  ☑ politenes.s (2)  ☑ poljteness (1)
☑ politenerss (2)  ☑ politenets (159)  ☑ polite'ness (5)
☑ pol'teness (1)  ☑ polite.ness (1)  ☑ politeneis (311)
☑ politeness3 (1)  ☑ politeniss (1)  ☑ politenese (3)

# Spelling variation

- "... Ships as shall be in <u>readynesse</u> for that service, and this matter <u>requireing</u> the greatest secrecy,... Your very <u>affectionett freind</u>" (Queen Anne, 1704)
- "Right <u>worshipfull</u> and my most <u>entierly beloude moder</u>, in the most <u>louly maner</u> I <u>recommaund</u> me <u>vnto youre gode moderhode</u>," (Elizabeth Poynings, 1459)
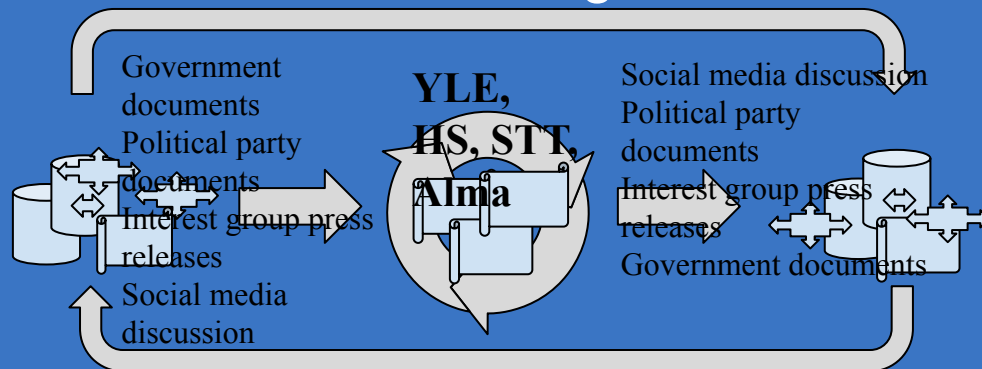
# Linguistic complexity

- "hopia* vaski*"~10 (any word starting with hopia within 10 words of any word starting with vaski)
- (/pi[sk]ku.*/ OR pien*) AND /pi[ij]at.*/
- |infl_fi_~2:Väinämöinen (any [dialectal] variation at max two edit distances from any inflected form of the Finnish name Väinämöinen)

# Focused, versatile, configurable NLP

- Interaction between the computational and the qualitative
- Operationalising indicators to measure the framing of different actors:
  - Tone
  - Standing
  - Polarisation
  - Affect
  - Fact/opinion -balance
  - Narrativity
- Tracking of the flow of vocabulary and agenda-setting power between actors and the media



Government documents
Political party documents
Interest group press releases
Social media discussion

YLE, HS, STT, Alma

Social media discussion
Political party documents
Interest group press releases
Government documents

# **Needs from language technology**

- Quantifying noise and bias (OCR errors, genre balance, variation, …)
  - Quantifying the amount of OCR errors/variation in different parts of a collection
  - Unsupervised genre identification → configurable genre classification
- Methods to clean noise and bias
  - OCR post-correction, better OCR through better language models?
  - Subsampling based on genre information
- Noise/variation-resistant NLP
  - Noise-resistant embeddings
  - Flexible variant normalisation, variation-sensitive querying
- Focused, configurable NLP
  - E.g. from sentiment analysis to the identification of hedging, tone, standing, polarisation, fact/opinion -balance, narrativity, …
  - + everything else you saw before

eetu.makela@helsinki.fi
http://iki.fi/eetu.makela

This presentation:
http://j.mp/lt-hums