



# Text Analysis and Political Rhetoric

---

Martijn Schoonvelde

05 August 2019

SICSS – Bamberg

# **Text Analysis and Political Rhetoric**

---

**Opportunities**

**Challenges**

**New Directions**

# **Text Analysis and Political Rhetoric**

---

**Opportunities**

**Challenges**

**New Directions**

(Also: plea for increased collaboration across disciplines)

# Who am I?

---

- Political scientist at  
University College Dublin

# Who am I?

---

- Political scientist at  
University College Dublin
- Text as data and rhetoric  
of political leaders

# Who am I?

---

- Political scientist at  
University College Dublin
- Text as data and **rhetoric  
of political leaders**
- Previously at Stony  
Brook, Exeter, EUI, VU

# Who am I?

- Political scientist at University College Dublin
- Text as data and **rhetoric of political leaders**
- Previously at Stony Brook, Exeter, EUI, VU
- #SICSS2018



Politicians talk...



... quite often



... really quite often



... really, really quite often



# Possibilities

---

# Benefits of Analyzing Political Speech

- Learn about all the politics that happens **in between elections**
  - **Testing ground** for new policy ideas
  - Agenda-setting / framing / priming: influences what the **media** writes about and what the public **thinks**
- Unobtrusive, accessible and retroactive



## Benefits of Analyzing Political Speech

---

We learned a lot of cool new  
things we didn't know before!

# Benefits of Analyzing Political Speech

We learned a lot of cool new things we didn't know before!

To name just a few ...

- Institutional power dynamics in the European Union (Cross & Hermansson, 2017)
  - Edit distances between legislative proposals and legislative outcomes as a measure legislative power in the EU.

Article

## Legislative amendments and informal politics in the European Union: A text reuse approach

**James P Cross**

School of Politics & International Relations, University College Dublin, Dublin, Ireland

**Henrik Hermansson**

Centre for European Politics, Department of Political Science, University of Copenhagen, Copenhagen, Denmark



European Union Politics

2017, Vol. 18(4) S81–602

© The Author(s) 2017

Reprints and permissions:

sagepub.co.uk/journalsPermissions.nav

DOI: 10.1177/1465116517717071

journals.sagepub.com/home/eup



# Benefits of Analyzing Political Speech

We learned a lot of **cool new things** we didn't know before!

To name just a few ...

- Blame shifting by heads of government in response to economic downturns  
(Traber, Schoonvelde & Schumacher, 2019)
  - **LDA topic models + sentiment analysis**



*“As a Social Democrat, I attribute exceptional importance to solidarity. (But) you also have obligations. You cannot spend all the money on drinks and women and then ask for help.”*

# Benefits of Analyzing Political Speech

We learned a lot of **cool new things** we didn't know before!

To name just a few ...

- Political leaders follow rather than lead discussion of public issues (Barberá et al. 2019)
  - Analysis of **topical content** in tweets of American legislators
  - Politicians more responsive to their supporters than to the general public

# Benefits of Analyzing Political Speech

---

Plus we are getting access to more and more data:

**Machine translation across languages**

**Automated speech recognition**

# Machine Translation and Bag of Words models

PA

## No Longer Lost in Translation: Evidence that Google Translate Works for Comparative Bag-of-Words Text Applications

Erik de Vries<sup>1</sup>, Martijn Schoonvelde<sup>2</sup> and Gijs Schumacher<sup>3</sup>

<sup>1</sup> Department of Media and Social Sciences, University of Stavanger, Stavanger, Norway. Email: [erik.devries@uis.no](mailto:erik.devries@uis.no)

<sup>2</sup> Department of Political Science and Public Administration, Vrije Universiteit, Amsterdam, The Netherlands.  
Email: [h.j.m.schoonvelde@vu.nl](mailto:h.j.m.schoonvelde@vu.nl)

<sup>3</sup> Department of Political Science, University of Amsterdam, Amsterdam, The Netherlands. Email: [g.schumacher@uva.nl](mailto:g.schumacher@uva.nl)

### Abstract

Automated text analysis allows researchers to analyze large quantities of text. Yet, comparative researchers are presented with a big challenge: across countries people speak different languages. To address this issue, some analysts have suggested using Google Translate to convert all texts into English before starting the analysis (Lucas *et al.* 2015). But in doing so, do we get lost in translation? This paper evaluates the usefulness of machine translation for bag-of-words models—such as topic models. We use the *europarl* dataset and compare term-document matrices (TDMs) as well as topic model results from gold standard translated text

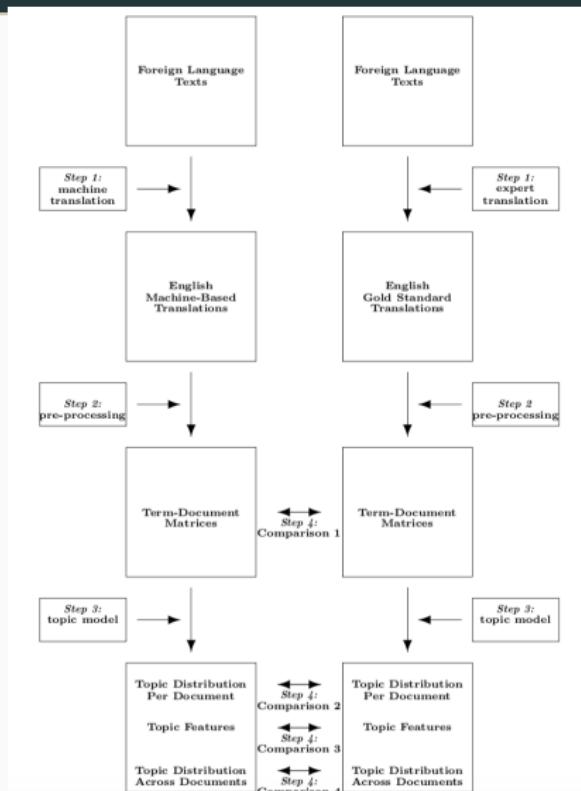
Question: is machine-translated text of good enough quality so that we can use it for our comparative research questions?

# Data

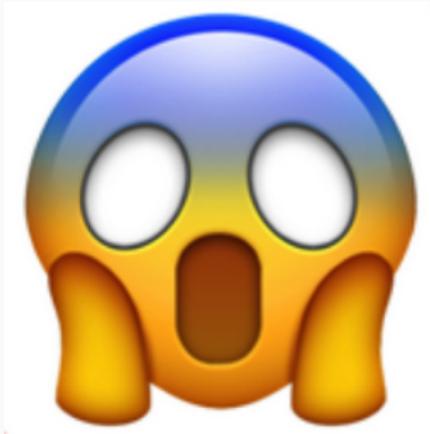
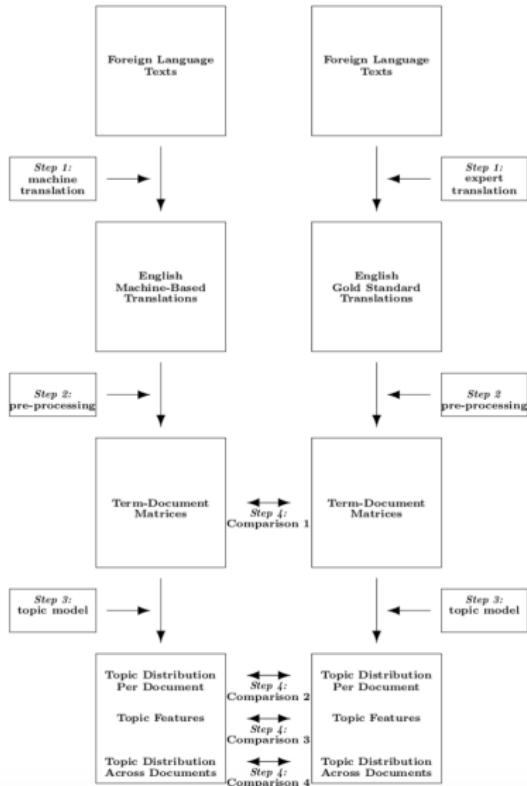
---

- Europarl dataset (Koehn 2005) → transcriptions of European Parliament debates by official translators
  - Danish, German, Spanish, French and Polish for the period of January 2007 to November 2011
  - Lots of cleaning to match individual chapters (i.e., debates, questioning, etc) per language pair
  - Turned 11,469 individual chapters we could match into html files – translated using Google Translate Web Plugin
- Take source language contributions and compare machine- and human-translations

# Research design



# Research design



# Similarity at document level

Figure 3: Distribution of cosine similarity per language pair

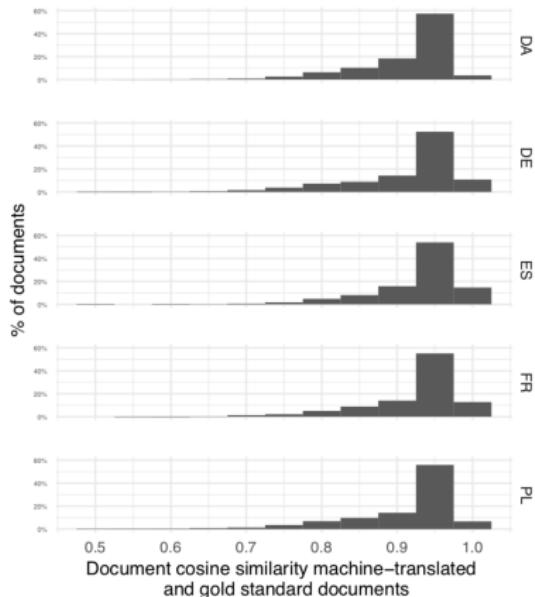
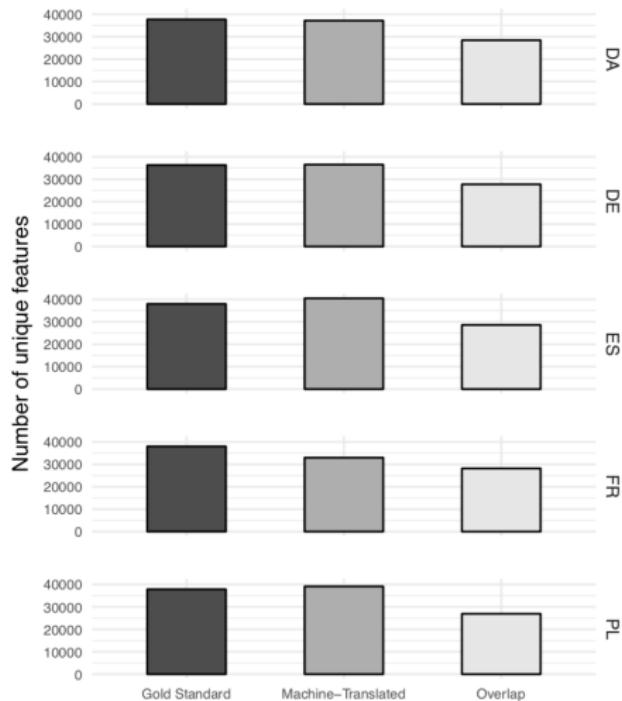


Table 2: Cosine similarity distribution per language

Language	N	Mean	St. Dev.	Min	Max
Danish	2,301	0.915	0.063	0.549	0.992
German	2,148	0.915	0.074	0.488	0.991
Spanish	2,335	0.929	0.059	0.483	0.991
French	2,347	0.925	0.064	0.564	0.989
Polish	2,338	0.913	0.073	0.475	0.989
Total:	11,469	0.919	0.066	0.475	0.992

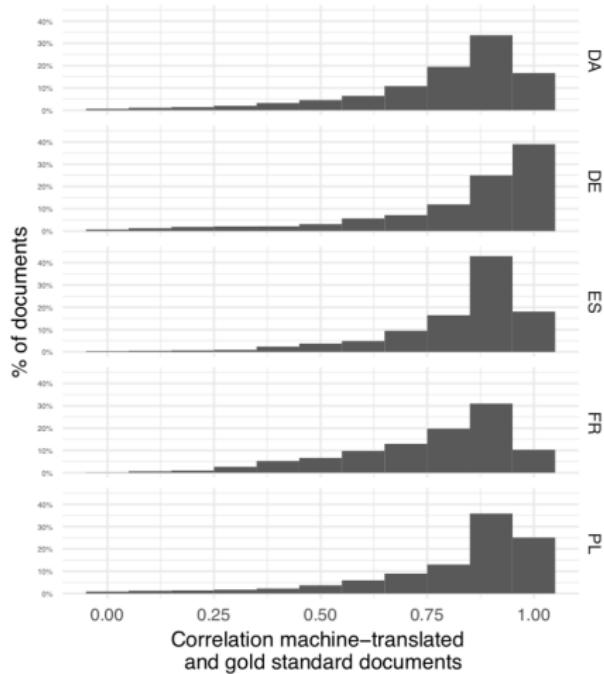
# Results at document level

Figure 4: Unique TDM features for gold standard and machine-translated corpora



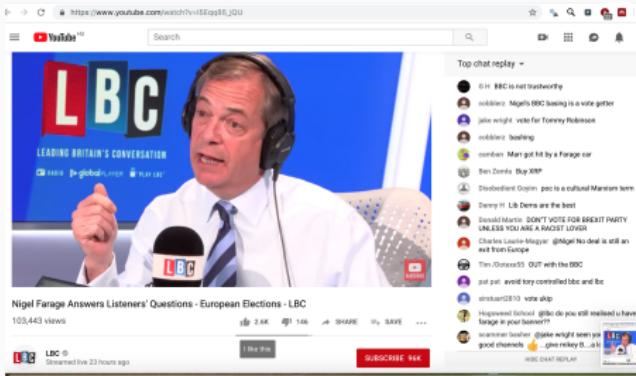
# Results at document level

Figure 5: Similarity of document-level topical prevalence with equal number of topics



# Automated speech recognition

- Lots of political speech occurs in talk shows, on Youtube channels, etc.
- But much of it not transcribed
- Proksch et al. (2019) test the quality of Automatic Speech Recognition (ASR) systems for bag of words models





## Testing the Validity of Automatic Speech Recognition for Political Text Analysis

Sven-Oliver Proksch<sup>①</sup>, Christopher Wratil<sup>②,1</sup>  
and Jens Wäckerle<sup>①</sup>

<sup>1</sup> Cologne Center for Comparative Politics, University of Cologne, Germany. Email: [so.proksch@uni-koeln.de](mailto:so.proksch@uni-koeln.de)

<sup>2</sup> Minda de Gunzburg Center for European Studies, Harvard University, Cambridge, MA 02138, USA

### Abstract

The analysis of political texts from parliamentary speeches, party manifestos, social media, or press releases forms the basis of major and growing fields in political science, not least since advances in “text-as-data” methods have rendered the analysis of large text corpora straightforward. However, a lot of sources of political speech are not regularly transcribed, and their on-demand transcription by humans is prohibitively expensive for research purposes. This class includes political speech in certain legislatures, during political party conferences as well as television interviews and talk shows. We showcase how scholars can use automatic speech recognition systems to analyze such speech with quantitative text analysis models of the “bag-of-words” variety. To probe results for robustness to transcription error, we present an original “word error rate simulation” (WERSIM) procedure implemented in *R*. We demonstrate the potential of automatic speech recognition to address open questions in political science with two substantive applications and discuss its limitations and practical challenges.

**Keywords:** Google, YouTube, text analysis, transcriptions, automatic speech recognition, campaigns

# Automated speech recognition

**Corpus:** 57 speeches (French, German, English) from the EP's State of the Union plenary (SOTEU)

---

	Human transcriptions	Google Speech API	Youtube Video
Speech 1			
Speech 2			
Speech 3			

---

# Automated speech recognition

Measure: word error rate (WER):

$$WER = \frac{S + D + I}{N}$$

$N$  = number of words in the gold standard transcriptions

$S$  = number words with inaccurate ASR transcription  
("substitutions")

$D$  = number of words that are missing in the ASR transcription  
("deletions")

$I$  = number of words in the ASR transcription that are not in the  
gold standard text ("insertions")

# WER

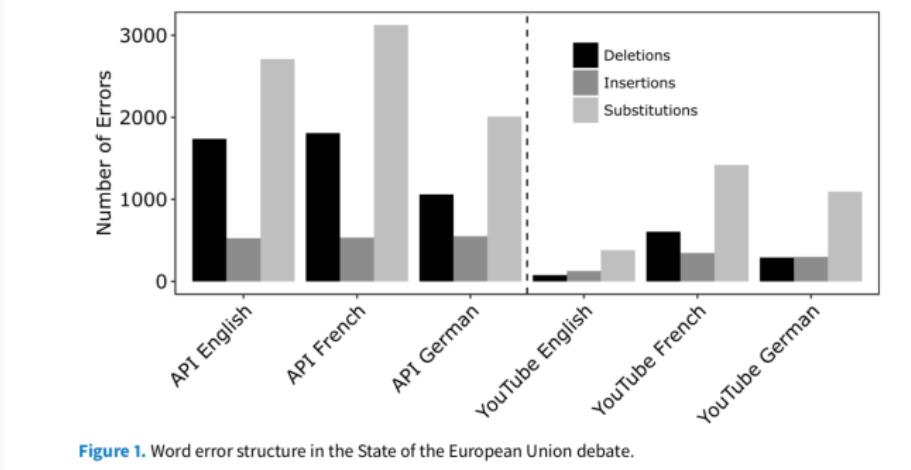


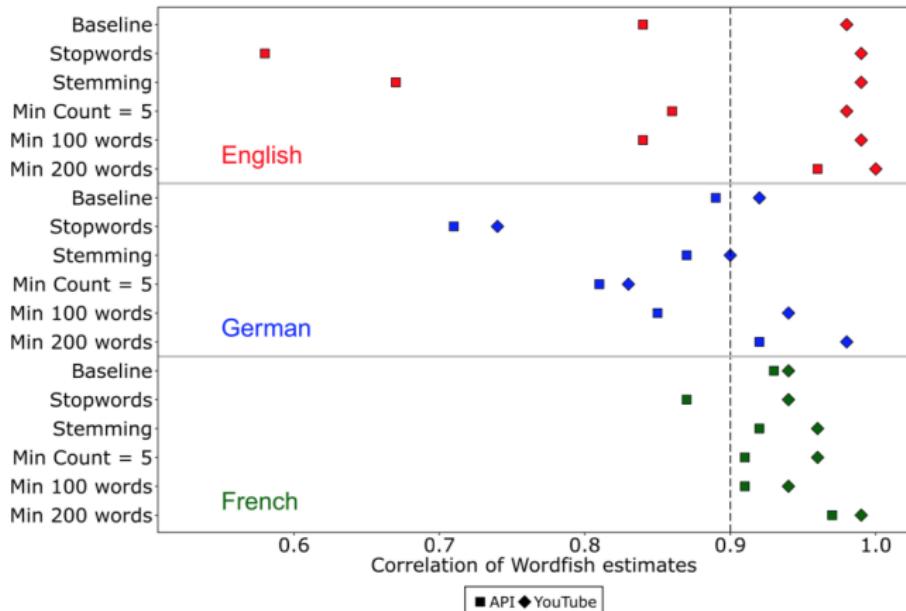
Figure 1. Word error structure in the State of the European Union debate.

Average WER **Youtube**: 0.03 (English), 0.12 (French), 0.10 (German)

Average WER **Google API**: 0.21 (English), 0.26 (French), 0.21 (German)

# Applications: scaling

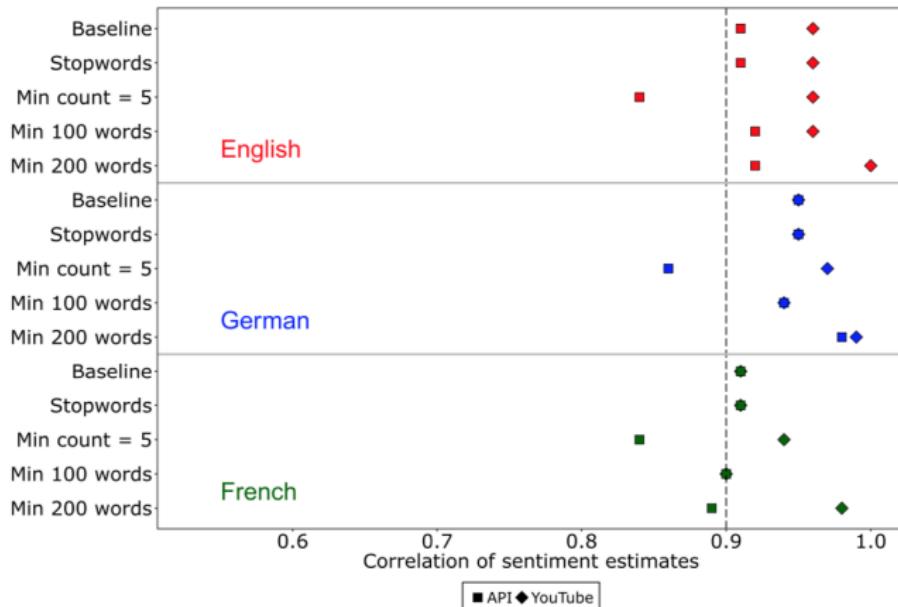
(a) Correlation of Wordfish estimates from human and ASR transcriptions



Wordfish scaling of speeches on a single dimension

# Applications: sentiment analysis

(b) Correlation of sentiment estimates from human and ASR transcriptions



## Take away

---

Automatic Speech Recognition generates meaningful data for bag-of-words text models

This will make transcription radically cheaper

Authors introduce a procedure (WERSIM) to simulate impact of increases in WER on quantities of interest

## Take away

---



## Challenges

---

# Challenges

---

1. Sampling
2. Preprocessing
3. Contextual meaning
4. Measurement
5. Causality

# Sampling

---

# Sampling

---

Speeches are different from Tweets are different from Facebook posts are different Op-Eds are different from party manifestos

# Sampling

---

Speeches are different from Tweets are different from Facebook posts are different Op-Eds are different from party manifestos

Politicians are strategic

## Sampling

---

Speeches are different from Tweets are different from Facebook posts are different Op-Eds are different from party manifestos

Politicians are strategic

Case in point: national politicians and the European Union (Rauh et al. 2019)

# Sampling

## Winning words: Individual differences in linguistic style among U.S. presidential and vice presidential candidates ★

Richard B. Slatcher <sup>a</sup>, , Cindy K. Chung <sup>a</sup>, James W. Pennebaker <sup>a</sup>, Lori D. Stone <sup>b</sup>

 Show more

<https://doi.org/10.1016/j.jrp.2006.01.006>

[Get rights and content](#)

### Abstract

The present study examines the personalities and psychological states of the 2004 candidates for U.S. president and vice president through their use of words. The transcripts of 271 televised interviews, press conferences, and campaign debates of John Kerry, John Edwards, George W. Bush, and Dick Cheney between January 4 and November 3, 2004 were analyzed using a computerized text analysis program. Distinct **linguistic** styles were found among these four political candidates, as well as differences between political parties and candidate types. Drawing on previous research linking word use and personality characteristics, the results suggest that the candidates had unique linguistic styles variously associated with cognitive complexity, **femininity**, depression, aging, presidentiality, and honesty.



Use LIWC to measure presidential candidates' **linguistic styles**

# Sampling

Table 2. Mean standardized scores for **linguistic** measures by speaker

Linguistic measure	Bush	Cheney	Kerry	Edwards	F	$\eta^2_p$
Cognitive Complexity	-.19 <sub>a,b</sub>	1.16 <sub>b</sub>	.05 <sub>a,b</sub>	-.32 <sub>a</sub>	2.60*	.03
Femininity	.45 <sub>b</sub>	-1.54 <sub>a</sub>	-.70 <sub>a</sub>	.78 <sub>b</sub>	17.93*	.10
Depression	-.85 <sub>a</sub>	-.98 <sub>a</sub>	1.19 <sub>b</sub>	.58 <sub>b</sub>	7.73*	.08
Age	1.34 <sub>b</sub>	-.59 <sub>a</sub>	-1.02 <sub>a</sub>	-.73 <sub>a</sub>	3.15*	.04
Presidentiality	.66 <sub>c</sub>	.87 <sub>c</sub>	.01 <sub>b</sub>	-1.47 <sub>a</sub>	6.79*	.07
Honesty	-.23 <sub>b</sub>	.90 <sub>c</sub>	-.80 <sub>a</sub>	.72 <sub>c</sub>	8.88*	.12

Note. Means for each measure are means of scores that have been standardized across the entire sample ( $N = 271$ ). Because of unequal sample sizes of speech samples for each of the four speakers, the average of the means for each respective measure does not equal zero. Means with different subscript levels are significantly different from one another at  $p < .05$  using Bonferroni post hoc comparison tests.  $\eta^2_p$ 's are conservative estimates of effect sizes for the overall differences among speakers for each linguistic measure, controlling for the effects of speech source type.

\*

$p < .05$ .

# Sampling

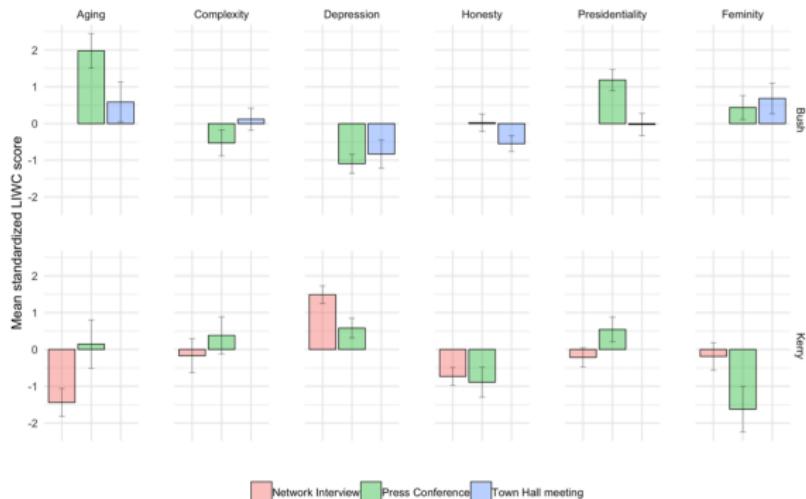


Figure 1. Linguistic style of George W. Bush and John Kerry on six linguistic style dimensions.

Note. This figure displays standardized LIWC scores for George W. Bush and John Kerry on six linguistic style dimensions (aging, complexity, depression, honesty, presidentiality, femininity) for various text sources for which we have at least 10 observations: network interviews (Kerry:  $n = 44$ ), press conferences (Bush:  $n = 57$ ; Kerry:  $n = 21$ ), and town hall meetings (Bush:  $n = 38$ ) (for more information, see Slatcher et al., 2007).

Linguistic style depends on corpus under study (Schoonvelde et al., 2019)

# Preprocessing

---

# Preprocessing

---

- Recent paper by Matthew Denny and Arthur Spirling (2018)
- “For just seven possible (binary) preprocessing steps, there would be  $2^7 = 128$  possible models to run and analyze”
- Possibility of ‘heading down “forking paths of inference”’ (Gelman and Loken 2014)
  - End result may crucially depend on arbitrary steps

# Preprocessing steps

---

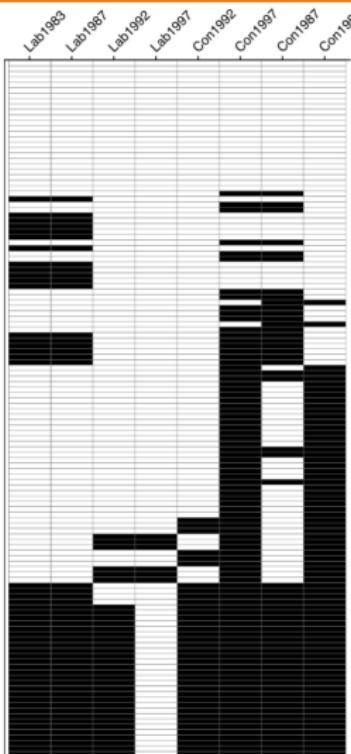
1. Punctuation
2. Numbers
3. Lowercasing
4. Stemming
5. W Stopword removal
6. 3 n-gram inclusion
7. Infrequently Used Terms

## Application: Wordfish on UK manifestos

- Wordfish model (Slapin and Proksch 2008) Labour and Conservative manifestos over 4 elections (1983 - 1997)
  - Estimates latent ideological position for each text
- Prior belief of order of manifestos:

$\text{Lab}_{1983} < \text{Lab}_{1987} < \text{Lab}_{1992} < \text{Lab}_{1997} < \text{Con}_{1992} < \text{Con}_{1997} < \text{Con}_{1987} < \text{Con}_{1983}$ .

# Wordfish positions of UK manifestos



**Figure 1.** Wordfish results for the 128 different preprocessing possibilities. Each row of the plot represents a different specification. A white bar implies that the manifesto for that year is in the correct place as regards our priors. A black bar implies it was misplaced.

# Preprocessing matters

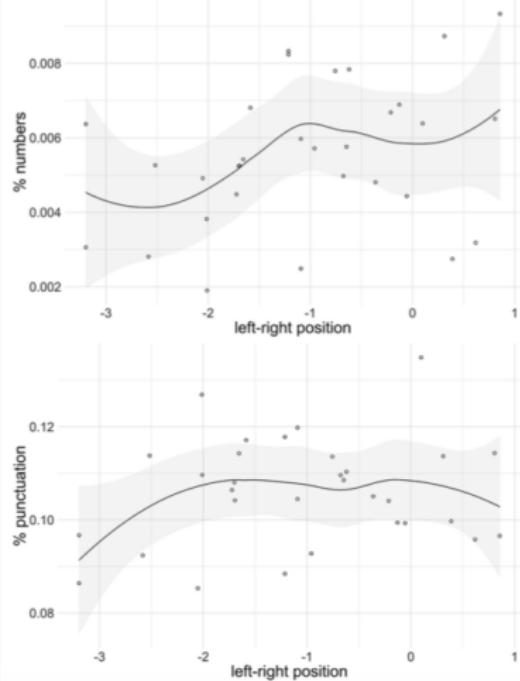
---

Preprocessing steps matter?

But not clear why they matter . . .

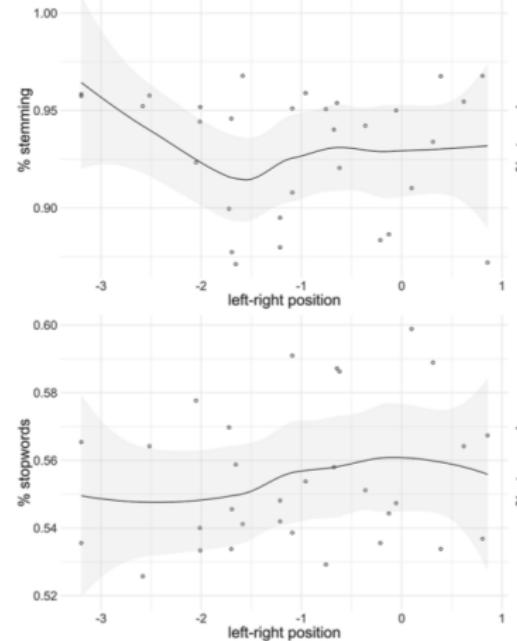
Maybe preprocessing is **not ideologicallyl neutral?**

# Preprocessing matters



Based on 2100 speeches by 31 MEPs (Schoonvelde, Schumacher & Bakker, 2019)

# Preprocessing matters



Based on 2100 speeches by 31 MEPs (Schoonvelde, Schumacher & Bakker, 2019)

# Preprocessing matters

*Political Psychology*



*Political Psychology*, Vol. xx, No. xx, 2016  
doi: 10.1111/pops.12327

## On the Grammar of Politics—or Why Conservatives Prefer Nouns

Aleksandra Cichocka  
*University of Kent*

Michał Bilewicz  
*University of Warsaw*

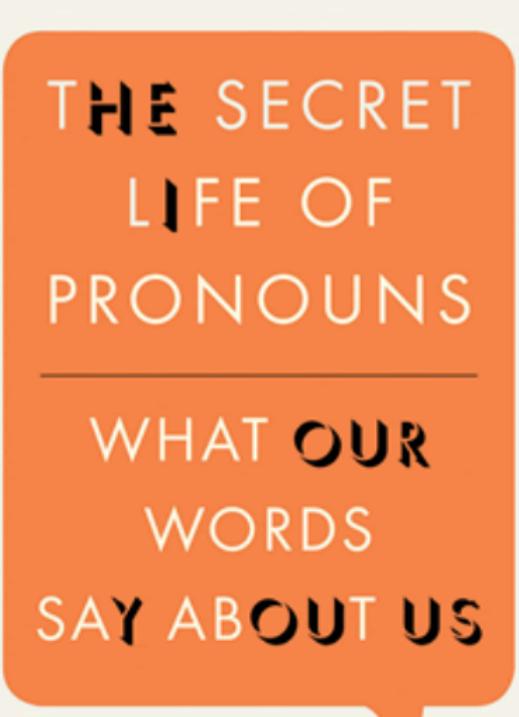
John T. Jost  
*New York University*

Natasza Marrouch  
*University of Connecticut*

Marta Witkowska  
*University of Warsaw*

“Nouns convey greater permanence, stability of subjects and objects, as well as categorical perceptions of social actors and the world at large. As such, they are likely to address conservatives? greater needs for order, certainty, and predictability.”

## Preprocessing matters



**THE SECRET  
LIFE OF  
PRONOUNS**

---

**WHAT OUR  
WORDS  
SAY ABOUT US**

JAMES W. PENNEBAKER

Finds associations between **personality type and linguistic patterns**: Introverted speakers prefer rich vocabulary (Oberlander & Gill, 2006), use more negations (e.g., Pennebaker & King, 1999). Neurotic speakers use more first-person singular and negative emotion words (Pennebaker & King, 1999). People high on openness to experience use more tentative words, such as 'maybe' or 'perhaps', and longer words (Pennebaker & King, 1999).

## **Contextual meaning**

---

# Contextual meaning

---

Bags of words ignore context

Polysemy poses a problem – “Low unemployment” versus “low economic growth”

Solutions:

- multigrams
- word embeddings

## Word embeddings

---



## Contextual meaning

---

But just like bag of words word embeddings require choices  
(Rodriguez & Spirling, 2019)

- Window size
- Size of the embedding space
- Local or pre-trained models

Impact of these choices not always clear

Same goes for pinpointing change in meaning using word embeddings models (Rodman, *Political Analysis, Forthcoming*)

# **Measurement**

---

# Measurement

Construct bycatch – when we mine text for one construct (latent position), we usually catch multiple (e.g., position and sentiment)

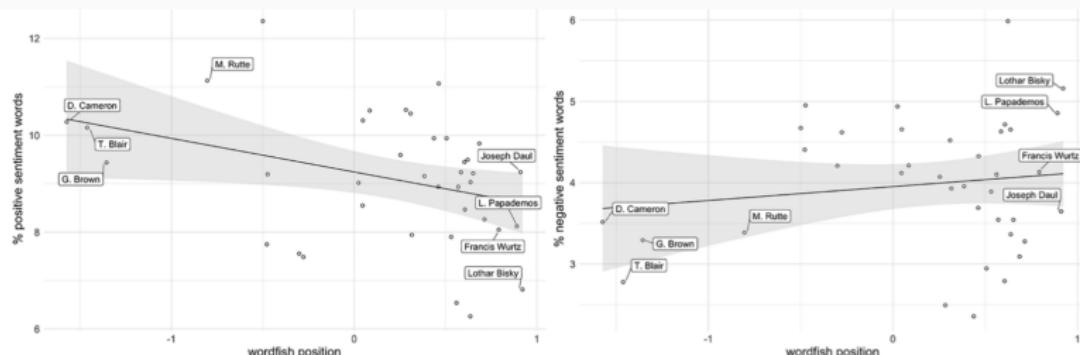


Figure 4. Use of sentiment and estimated *Wordfish* positions of speeches by politicians in the European Union.

Note. These scatterplots denote the average use of negative sentiment (left) and average use of positive sentiment (right) over the range of the estimated average *Wordfish* position of heads of government and MEP group leaders. It shows that *Wordfish* scores and the use of positive sentiment are correlated with each other.

Schoonvelde, Schumacher & Bakker (2019)

# Measurement

**Trait versus state:** is text a manifestation of some latent characteristic of the speaker or driven by external circumstances?

Disciplinary differences in, for example, psychology and political science



Journal Information  
Journal TOC

[Search APA PsycNET](#)

## PsyARTICLES: Journal Article

[The language of well-being: Tracking fluctuations in emotion experience through everyday speech.](#)

[© Request Permissions](#)

**Sun, Jessie, Schwartz, H. Andrew, Son, Youngseo, Kern, Margaret L., Vazire, Simine**  
*Journal of Personality and Social Psychology*, Apr 04 , 2019, No Pagination Specified

The words that people use have been found to reflect stable psychological traits, but less is known about the extent to which everyday fluctuations in spoken language reflect transient psychological states. We explored within-person associations between spoken words and self-reported state emotion among 185 participants who wore the Electronically Activated Recorder (EAR; an unobtrusive audio recording device) and completed experience sampling reports of their positive and negative emotions 4 times per day for 7 days (1,579 observations). We examined language using the Linguistic Inquiry and Word Count program (LIWC; theoretically created dictionaries) and open-vocabulary themes (clusters of data-driven semantically-related words). Although some studies give the impression that LIWC's positive and negative emotion dictionaries can be used as indicators of emotion experience, we found that when computed on spoken language, LIWC emotion scores were not significantly associated with self-reports of state emotion experience. Exploration of other categories of language variables suggests a number of hypotheses about substantive everyday correlates of momentary positive and negative emotion that can be tested in future studies. These findings (a) suggest that LIWC positive and negative emotion dictionaries may not capture self-reported subjective emotion experience when applied to everyday speech, (b) emphasize the importance of establishing the validity of language-based measures within one's target domain, (c) demonstrate the potential for developing new hypotheses about personality processes from the open-ended words that are used in everyday speech, and (d) extend perspectives on intraindividual variability to the domain of spoken language. (PsycINFO Database Record (c) 2019 APA, all rights reserved)

# Causality

---

# Causality

---

Only very few text papers take causal inference seriously

Little overlap between the text as data and the causal inference crowd

# Causality

“In 2015, a legal ruling forced the German government to declassify all its public opinion research. Our causal identification strategy exploits the demonstrably exogenous timing of the reports’ dissemination to cabinet members within a window of a few days. We find that exposure to the public opinion reports leads elites to change their rhetoric markedly.”

Hager & Hilbig (Forthcoming,  
*American Journal of Political  
Science*)



# Causality

---

Hager & Hilbig also pre-registered their hypotheses

Down the line this will become a necessity when doing applied text work

## Some best practices when analyzing political speech

---

### 1. Use similar text sources to the extent possible

- When multiple text sources are the only option, account for this in the analysis
- For example, by adding meta data to your model (e.g., structural topic models)

## Some best practices when analyzing political speech

1. Use similar text sources to the extent possible
  - When multiple text sources are the only option, account for this in the analysis
  - For example, by adding meta data to your model (e.g., structural topic models)
2. Get to **know your data**. How did a speech come about? Who was involved? Incorporate this information in the analysis

## Some best practices when analyzing political speech

1. Use similar text sources to the extent possible
  - When multiple text sources are the only option, account for this in the analysis
  - For example, by adding meta data to your model (e.g., structural topic models)
2. Get to **know your data**. How did a speech come about? Who was involved? Incorporate this information in the analysis
3. Consider in what ways preprocessing steps can correlate with stable speaker characteristics.
  - Average results across preprocessing steps, particularly when working with a small corpus

# Some best practices when analyzing political speech

1. Use similar text sources to the extent possible
  - When multiple text sources are the only option, account for this in the analysis
  - For example, by adding meta data to your model (e.g., structural topic models)
2. Get to **know your data**. How did a speech come about? Who was involved? Incorporate this information in the analysis
3. Consider in what ways preprocessing steps can correlate with stable speaker characteristics.
  - Average results across preprocessing steps, particularly when working with a small corpus
4. Consider **construct by-catch**

# Some best practices when analyzing political speech

1. Use similar text sources to the extent possible
  - When multiple text sources are the only option, account for this in the analysis
  - For example, by adding meta data to your model (e.g., structural topic models)
2. Get to **know your data**. How did a speech come about? Who was involved? Incorporate this information in the analysis
3. Consider in what ways preprocessing steps can correlate with stable speaker characteristics.
  - Average results across preprocessing steps, particularly when working with a small corpus
4. Consider **construct by-catch**
5. Pre-register your hypotheses

## Moving forward

---

# Speech versus communication

---

We often analyze text, but what we are interested in is  
**communication**

Communication is both verbal and non-verbal

It is not just about **what** politicians say but **how** they say it

Promising new directions in analysis of speech combine text with  
other data sources / other forms of communication

# Speech and voice pitch

## Audio-as-data

- In political science: Work by Bryce Dietrich and colleagues

Automated text analysis to detect topics and **voice pitch** to detect emotional intensity

- Understands voice pitch as sincere and text as strategic

Women more than men increase their voice pitch when talking about women's issues

- Correlates with relevant behavioral patterns

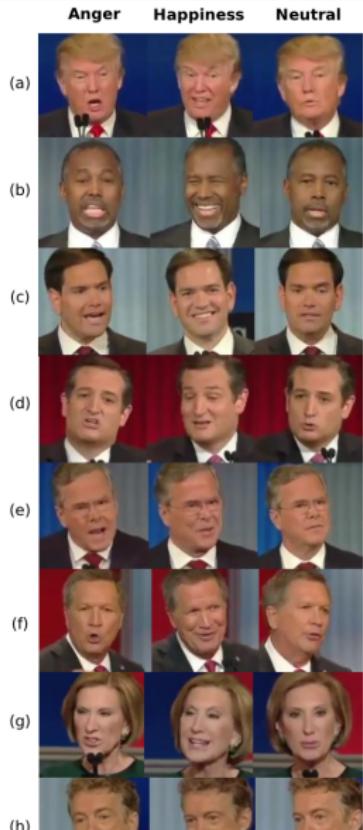
## Pitch Perfect: Vocal Pitch and the Emotional Intensity of Congressional Speech\*

Bryce J. Dietrich<sup>†</sup>      Matthew Hayes<sup>‡</sup>      Diana Z. O'Brien<sup>§</sup>

### Abstract

Though audio archives are available for a number of political institutions, the data they provide receive scant attention from researchers. Yet, audio data offer important insights, including information about speakers' emotional states. Using one of the largest collections of natural audio ever compiled—74,158 Congressional floor speeches—we introduce a novel measure of legislators' emotional intensity: small changes in vocal pitch that are difficult for speakers to control. Applying our measure to MCs' floor speeches about women, we show that female MCs speak with greater emotional intensity when talking about women as compared to both their male colleagues and their speech on other topics. Our two supplementary analyses suggest that increased vocal pitch is consistent with legislators' broader issue commitments, and that emotionally intense speech may affect other lawmakers' behavior. More generally, by demonstrating the utility of audio-as-data approaches, our work highlights a new way of studying political speech.

# Speech and images



Boussalis and Coan (2019) use Microsoft Face API to score stills from Republican contenders on basic emotions

Match these real time with approval ratings from a focus group

Find that displays of anger **increase approval ratings** among the public

# Speech and images

---

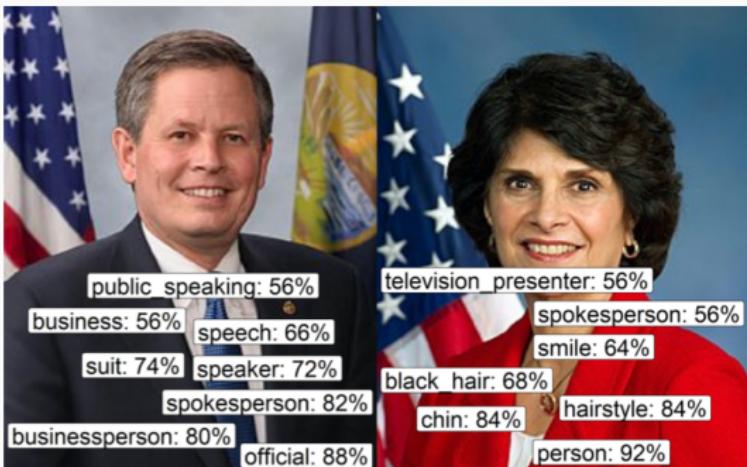
But how **biased** is  
off-the-shelf image  
recognition  
software?



Wikipedia image annotations for Steve Daines and Lucille Roybal-Allard

# Speech and images

But how **biased** is  
off-the-shelf image  
recognition  
software?



Wikipedia image annotations for Steve Daines and Lucille Roybal-Allard

## Speech and video

---



← → C https://www.youtube.com/channel/UC-IHJZR3Gqxm24\_Vd\_AJ5Yw

YouTube  M

Library History Watch later Liked videos

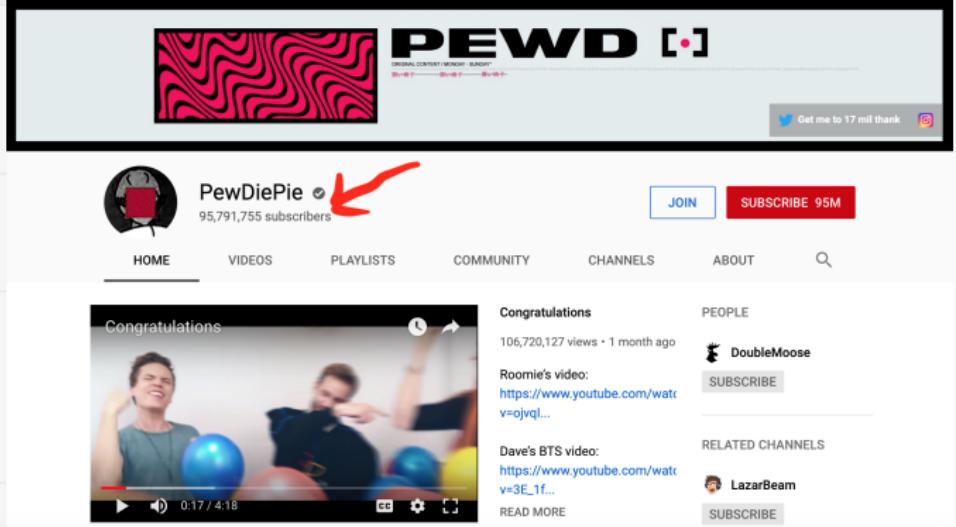
SUBSCRIPTIONS

- Justin Esarey
- Browse channels

MORE FROM YOUTUBE

- YouTube Premium
- YouTube Movies
- Gaming
- Live

Settings



The image shows the YouTube channel page for PewDiePie. At the top, there's a banner with a pink and black wavy pattern and the text "PEWD [.]". Below the banner, the channel name "PewDiePie" is displayed with a verified badge (a blue circle with a white checkmark). A red arrow points to this badge. The subscriber count is listed as "95,791,755 subscribers". To the right of the channel name are "JOIN" and "SUBSCRIBE 95M" buttons. Below the channel name, there are tabs for "HOME", "VIDEOS", "PLAYLISTS", "COMMUNITY", "CHANNELS", and "ABOUT".  
  
The main content area features a video thumbnail for a video titled "Congratulations" showing two men laughing. Below the video are controls for play, volume, and settings, along with the duration "0:17 / 4:18".  
  
To the right of the video, there's a "Congratulations" section with a video thumbnail of two men, a "PEOPLE" section with a user named "DoubleMoose", and a "RELATED CHANNELS" section with a user named "LazarBeam".  
  
At the bottom of the page, there's a "READ MORE" link.

## LTTA: Linguistic Temporal Trajectory Analysis

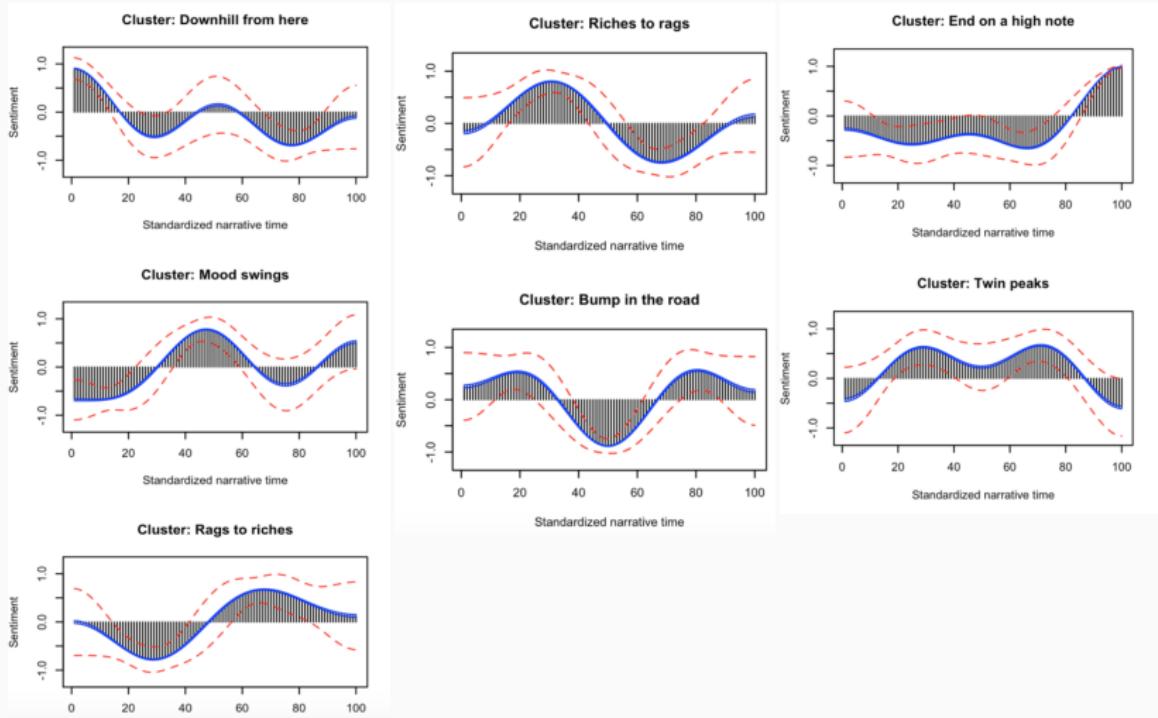
### Identifying the sentiment styles of YouTube's vloggers

**Bennett Kleinberg**  
Department of Psychology  
University of Amsterdam  
  
Department of Security  
and Crime Science  
University College London  
b.a.r.kleinberg@uva.nl

**Maximilian Mozes**  
Department of  
Informatics  
Technical University  
of Munich  
mozes@cs.tum.edu

**Isabelle van der Vegt**  
Department of Security and  
Crime Science  
University College London  
isabelle.vegt.17@ucl.ac.uk

- Method to inspect how linguistic markers like sentiment shift over time in a text
  - sentiment detection in these using a **sliding window** – standardize within fixed time periods
- Corpus: selection of vlogs from the most popular vloggers on Youtube
- Obtain transcripts from all vlogs produced



Cluster	Family	Female	Male
Downhill from here	2.23	1.26	-2.88*
Mood swings	-2.31	1.96	1.25
Rags to riches	2.13	-1.95	-1.08
Riches to rags	-2.05	4.88*	-0.56
Bump in the road	1.69	-1.12	-1.08
End on a high note	-5.16*	-6.03*	8.32*
Twin peaks	3.83*	2.25	-4.99*

Table 3. Standardized residuals for the cluster-by-gender association.

Lots of cool political applications possible!

# Conclusion

---

Speech is a great source of political data

But we have to be careful when analysing it

Keep in mind work in other disciplines

**Thank you!**

**Twitter: @hjms**