

# **Beyond Word Clouds**

## **Combining Entities and Topics for Fine-Grained Analyses of Historical (and Political) Texts**



**Federico Nanni**, Hiram Kümper and Simone Paolo Ponzetto  
University of Mannheim

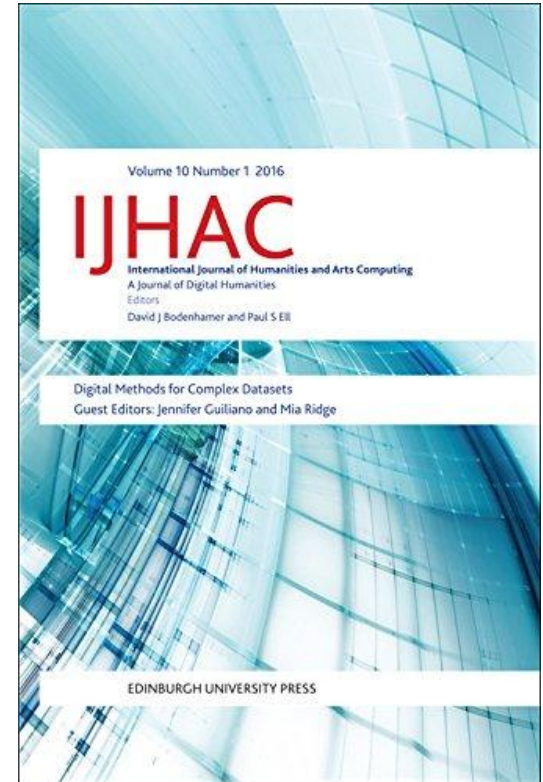
# Background

- Bachelor and Master in **Contemporary History**
- Final year PhD student in **Digital Humanities** at UniBo
  - Thesis: “The Web as a Historical Corpus”
- Researcher at **Data and Web Science Group**, UniMannheim

# Cfp: special issue of IJHAC

The **future of digital methods** for complex datasets.

Guest editors: Jennifer Giuliano and Mia Ridge



# **The future of digital history**

Historians are dealing with large collections:

# The future of digital history

Historians are dealing with large collections:



# The future of digital history

Historians are dealing with large collections:



- Text exploration
- Information Retrieval
- Quantification

# NLP and digital history

Voyant Tools: Reveal Your Texts

http://voyanttools.org/corpus=1321892653071.7045

Highrise MELCamp Stephens QT and HTML SYSJU DRMC Garmin Louisa Doane Doane Reference Pages HLU-Outlook Web Saved Tabs Read Later Instapaper Yojimbo

Voyant Tools: Reveal Your Texts

**Corpus**

Summary

- There are 127 documents in this corpus with a total of **224,647 words** and **14,640 unique words**.
- Longest documents** (by words): *Black Columbus Transcripts/Everything You Learn Comes Together; Jackson, Reggie.txt* (11,011). Shortest documents: *Overcoming Struggle...* (101), *JON LAMPLEY NARRATIVE...* (108). All...
- Highest vocabulary density** (by words per token): *Overcoming Struggle...* (873.3), *JON LAMPLEY NARRATIVE...* (857.4). Lowest density: *Steppen' Out on Faith...* (125.2), *Finding Myself - Singleton...* (142.0). All...
- Most frequent words** in the corpus: *like* (1,589), *just* (1,263), *read* (1,214), *know* (1,194), *school* (870). More...
- Words with **notable peaks in frequency** across the corpus: *really* (774), *like* (1,589), *just* (1,263), *read* (1,214), *know* (1,194), *school* (870). More...
- Distinctive words** (compared to the rest of the corpus):
  - Route 66 Averil, Julia...: *like* (1), *just* (3), *read* (2), *know* (1)

**Words in the Entire Corpus**

**Corpus Reader**

REGGIE: Teaches over there from time to time. So yeah probably back then, when she found that out about my reading level, I'm sure she pushed me pretty hard at home in addition to whatever I was doing in school at the time.

[Media skips forward]

REGGIE: Not really, basically in fifth grade, no, no, no, I think it was second, no that's not right, third or fourth grade I think it was, they had a general music class where like all the kids play recorders.

SPEAKER: Oh yeah.

[laughing]

REGGIE: And then from there, I think it was fifth grade, I played trumpet for a year and then in middle school and on I went and actually started playing drums, although I was actually playing the drums way before that. I think I started when I was three.

SPEAKER: Where, how?

REGGIE: At home, my mom bought me a little drum set and I would literally play, as she tells me-

[laughing]

REGGIE: ...play along with the radio and records, you know, and then I started playing in **church** when I was like six or seven. I've been playing ever since. As far as informative study

**Word Trends**

Relative Frequencies

church

Relative

Segments Search

**Keywords in Context**

Left	Keyword	Right
then I started playing in	church	when I was like six
four, when we were at	church	I would make my way
when I wasn't playing in	church	what I would always do
grew up just playing in	church	and I kind of played

Page 1 of 1 Context Preview

**Words in Documents**

Type	Count	Relative	Trend
church	5	11.08	

**Corpus**

127 documents with 224,647 tokens and 14,640 types

Document Label	Tokens	Types	Densh
1) Black Columbus Transcripts/Route 66 Averil, Julia (Transcript.txt)	211	125	592.4
2) Black Columbus Transcripts/Everything You Learn Comes Together; Jacks...	4,514	816	180.8
3) Black Columbus Transcripts/Determination: St. Cyr, Frantz (Transcript.txt)	307	147	478.8
4) Black Columbus Transcripts/Mar'Jesna's Narrative Transcript.txt	995	316	317.6
5) Black Columbus Transcripts/The Power of Words Dunn, Victoria (2010-03-1...	2,210	672	304.1
6) Black Columbus Transcripts/Lewis Black's Influence on My Writing-D. Fabio...	578	212	366.8

Page 1 of 1 Reset Search 1 of 1

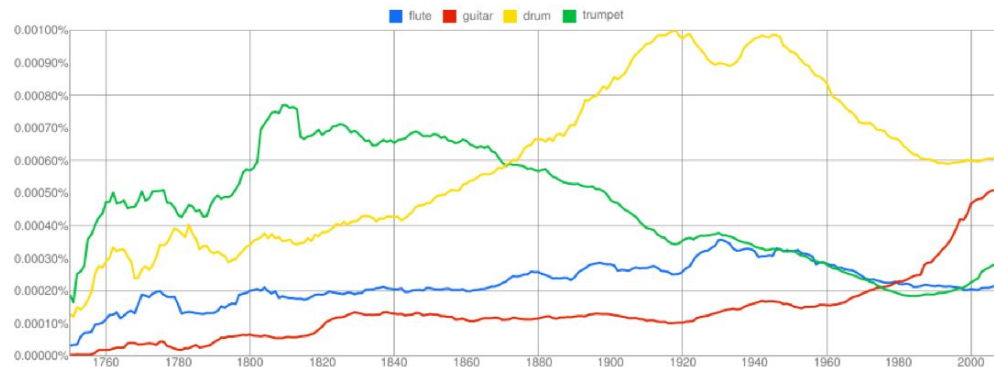
Voyant Tools, Stefan Sinclair & Geoffrey Rockwell (©2011) v. 1.0 beta (4302)

# NLP and digital history

Google labs Books Ngram Viewer

Graph these case-sensitive comma-separated phrases:  between  and  from the corpus  with smoothing of .

[Search lots of books](#)



Search in Google Books:

<a href="#">1750 - 1793</a>	<a href="#">1794 - 1923</a>	<a href="#">1924 - 1940</a>	<a href="#">1941 - 1981</a>	<a href="#">1982 - 2008</a>	<a href="#">flute</a>
<a href="#">1750 - 1825</a>	<a href="#">1826 - 1979</a>	<a href="#">1980 - 1992</a>	<a href="#">1993 - 2001</a>	<a href="#">2002 - 2008</a>	<a href="#">guitar</a>
<a href="#">1750 - 1777</a>	<a href="#">1778 - 1801</a>	<a href="#">1802 - 1817</a>	<a href="#">1818 - 1952</a>	<a href="#">1953 - 2008</a>	<a href="#">trumpet</a>
<a href="#">1750 - 1801</a>	<a href="#">1802 - 1911</a>	<a href="#">1912 - 1927</a>	<a href="#">1928 - 1983</a>	<a href="#">1984 - 2008</a>	<a href="#">drum</a>

Run your own experiment! Raw data is available for download [here](#).



# NLP and digital history

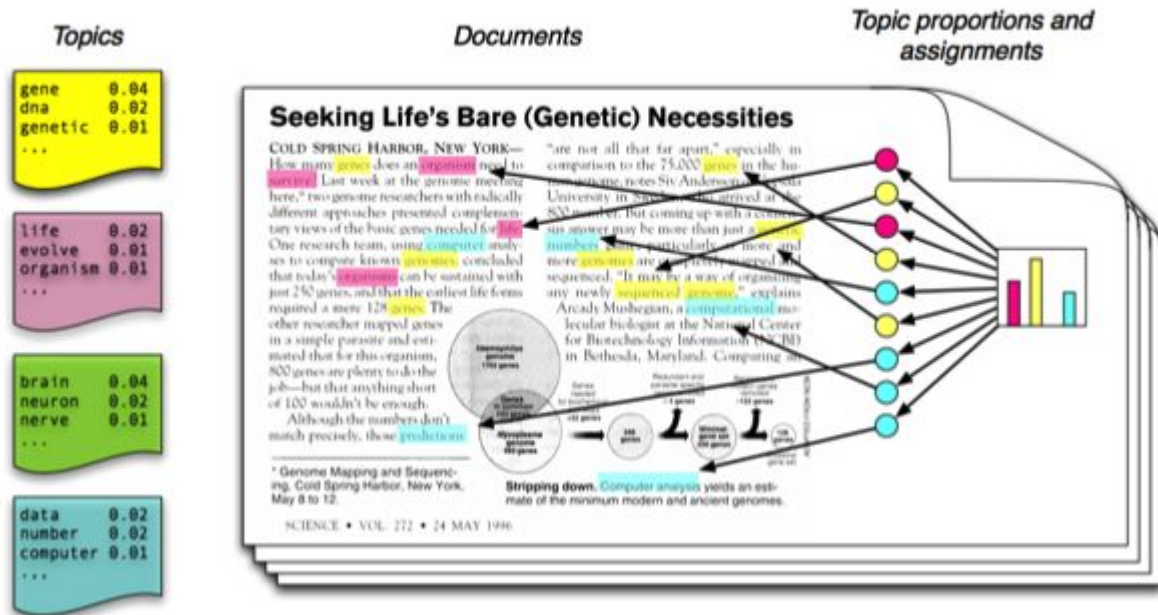
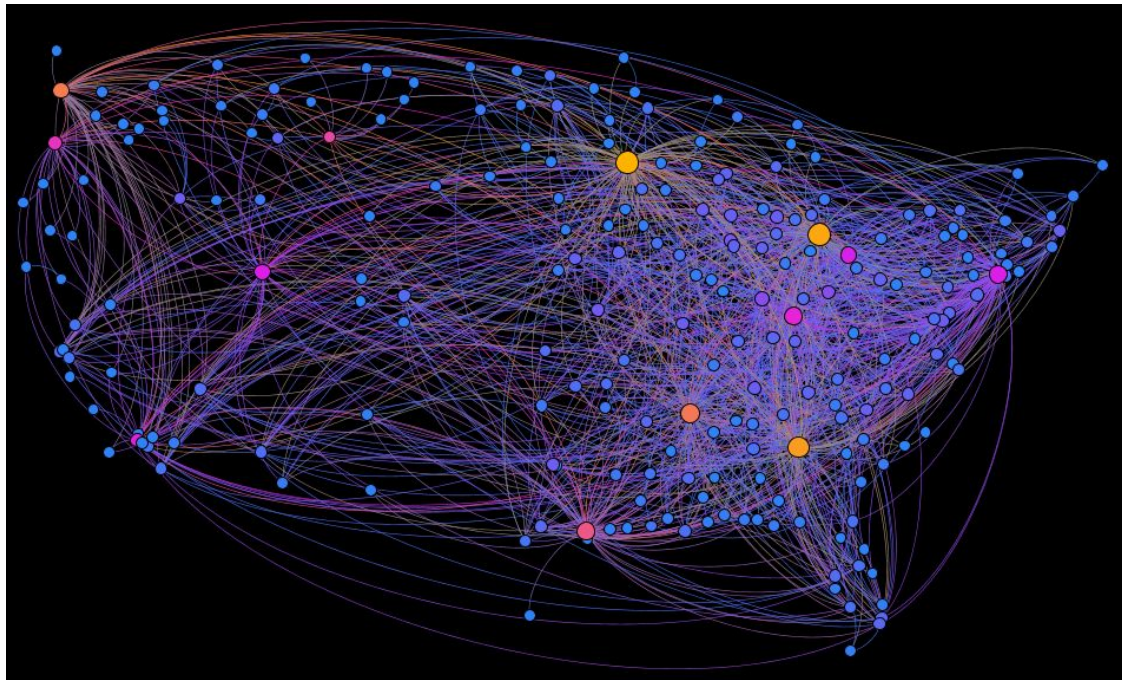
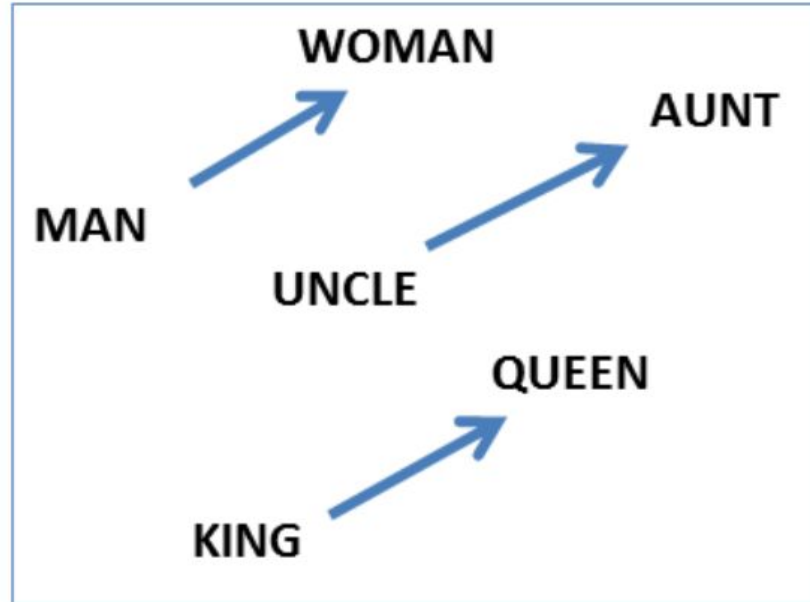


Figure source: Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4), 77-84.

# NLP and digital history



# NLP and digital history



# Current problems

It's very difficult to:

- Find the **right approach** for a specific question
- Establish its **reliability**
- Move **beyond** text exploration
- **Answer** a humanities research question

# Our goal at Data and Web Science Group

Sustain **hypothesis-testing** analyses

Offer both tool implementations and **evaluation platforms**

Train students with experience both in **data science and humanities / social sciences**

Nanni, Kümper and Ponzetto, “Semi-supervised Textual Analysis and Historical Research Helping Each Other”, IJHAC, 2016.

# Today's talk

Overview of three researches we recently conducted at DWS:

- 1) Entities as Topic Labels
- 2) Building Entity-based Collections of Global Events
- 3) Topic-Based Analysis of Political Positions

# Entities as Topic Labels

---

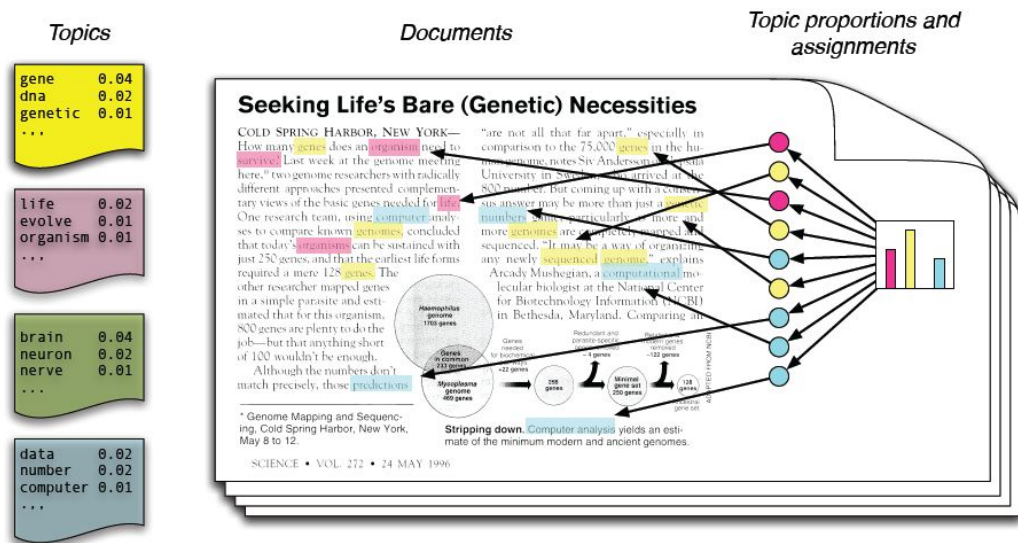
Anne Lauscher<sup>1</sup>, **Federico Nanni**<sup>1</sup>, Pablo Ruiz Fabo<sup>2</sup> and Simone Paolo Ponzetto<sup>1</sup>

<sup>1</sup>Data and Web Science Group, University of Mannheim

<sup>2</sup>LATTICE Lab, École Normale Supérieure

# Why topic models are awesome

They are able to identify the most important topics in a collection of documents.





# Why topic models are **NOT** awesome

The topics obtained are difficult to interpret and evaluate.

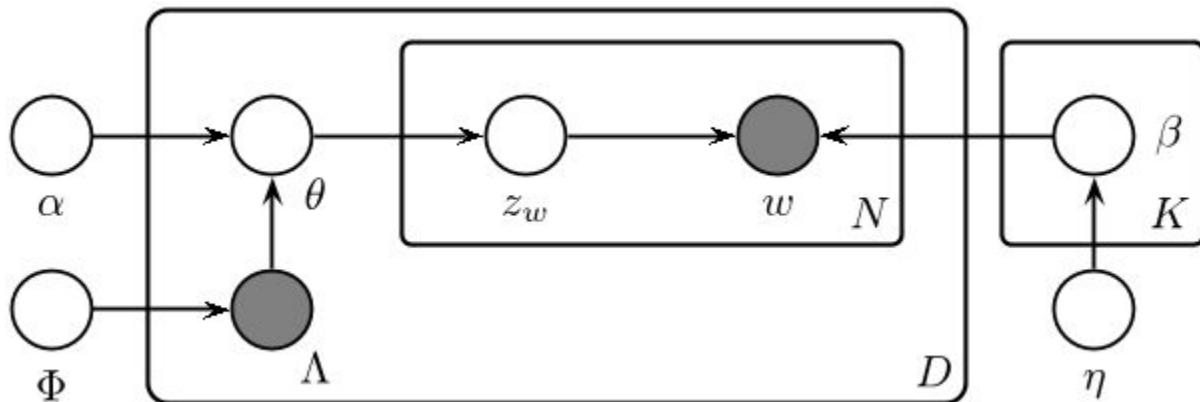
# Why topic models are **NOT** awesome

The topics obtained are difficult to interpret and evaluate.



# Labeled LDA

Each **document** is described with one or more **labels**, each label is associated with a specific **topic** (Ramage et al., 2009).

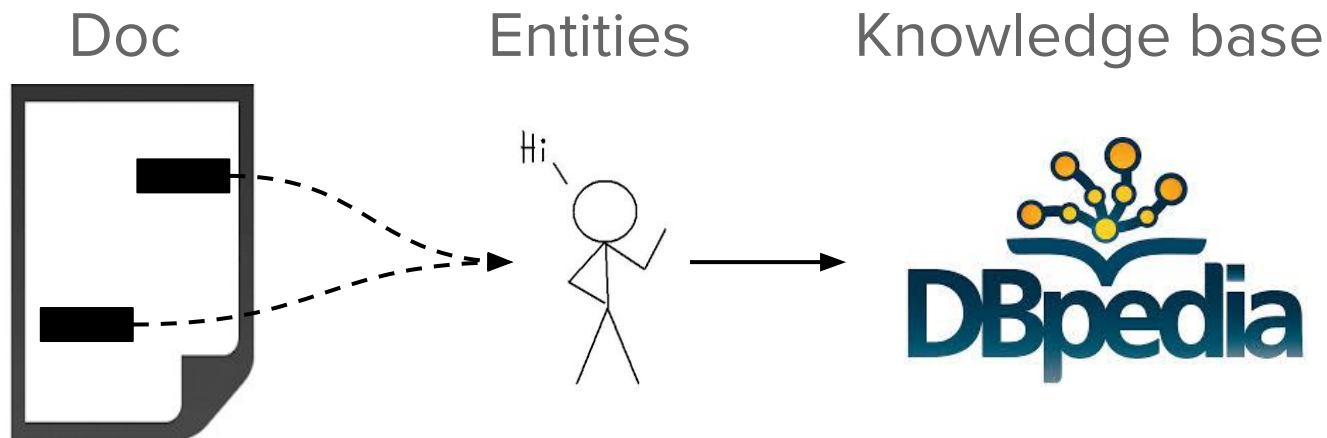


# How to automatically obtain labels?

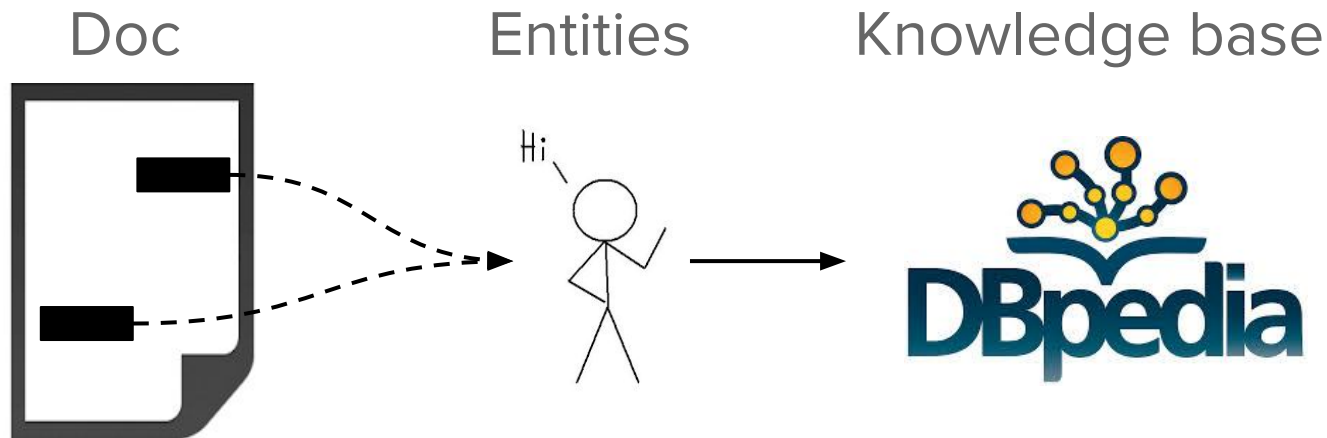
Different approaches:

- Keyphrase digger (FBK - Trento)
- Labeling the obtained topics (Hulpus et al., 2014)

# Our approach: entities

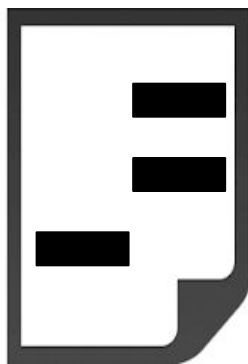


# Our approach: entities



# Our approach: entity ranking

Doc1



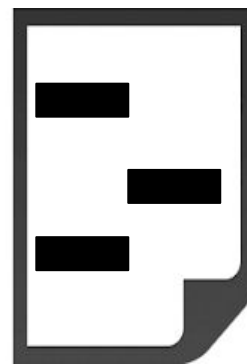
Entity1  
Entity2  
Entity3

Doc2



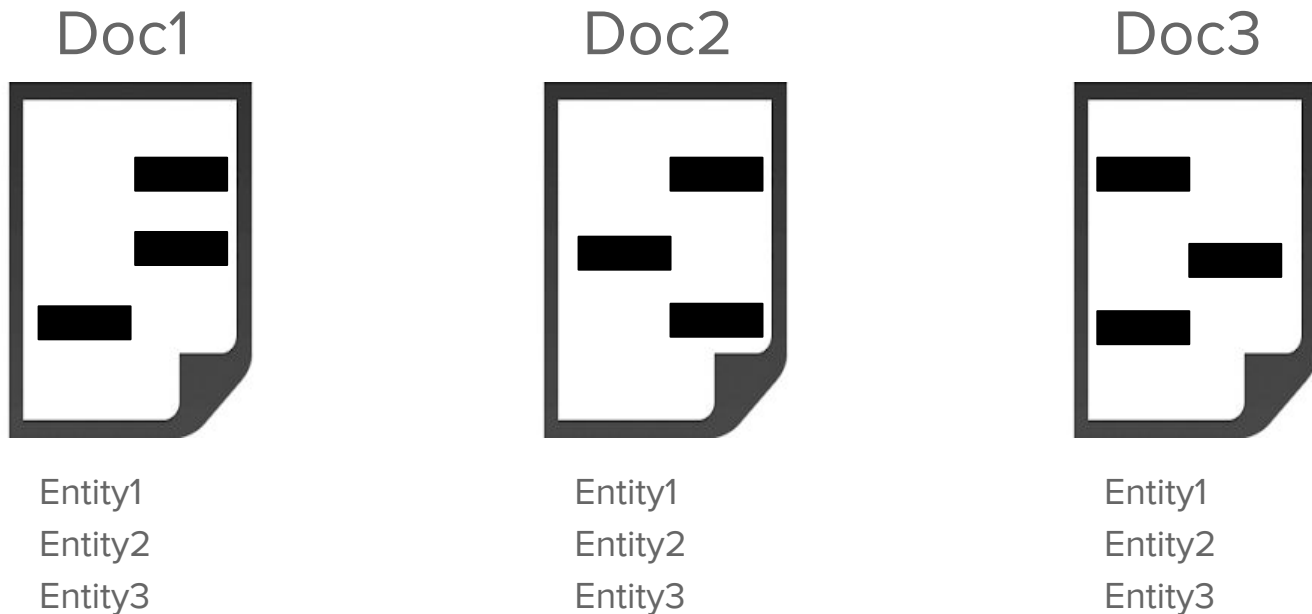
Entity1  
Entity2  
Entity3

Doc3



Entity1  
Entity2  
Entity3

# Our approach: entities as topic labels



**Labeled LDA!**



# Different case-study

Transcripts from **European Parliament's** fifth term (1999-2004).



Threads in the **Enron Corpus** (600.000 emails, 158 employees).



Discussions in the **Hillary Clinton Email Dataset**, a collection of redacted versions of emails (available on Kaggle).



# Europarl corpus

Examined most relevant topics addressed by each party in the European Parliament's fifth term (1999-2004).



## Les Verts (France)

	Label: Consumer (47%)	Label: GMO (34%)
Topic words	product directive consumer safety law market	human health agreement food measure sustainable

# Europarl corpus

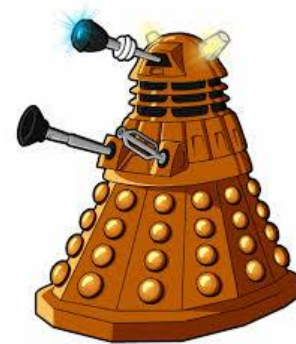
Examined most relevant topics addressed by each party in the European Parliament's fifth term (1999-2004).



## Les Verts (France)

	Label: Consumer (47%)	Label: GMO (34%)
Topic words	product directive consumer safety law market	human health agreement food measure sustainable

EVALUATE!  
EVALUATE!  
EVALUATE!



# How to evaluate it?

1. Label selection

**Les Verts** (France)

**Consumer**

**GMO**

# How to evaluate it?

1. Label selection
2. Label ranking

## Les Verts (France)

Consumer (**47%**)

GMO (**34%**)

# How to evaluate it?

1. Label selection
2. Label ranking
3. Label-topic relation

## Les Verts (France)

Consumer (47%)



product  
directive  
consumer  
safety  
law  
market

GMO (34%)



human  
health  
agreement  
food  
measure  
sustainable

# How to evaluate it?

1. Label selection
2. Label ranking
3. Label-topic relation

Following slides introduce  
the ongoing master thesis work  
of **Anne Lauscher** (Univ. Mannheim)

## Les Verts (France)

Consumer (47%)



product  
directive  
consumer  
safety  
law  
market

GMO (34%)



human  
health  
agreement  
food  
measure  
sustainable

# Label selection and ranking

Country Weather  
Limiting magnitude  
Marginal cost  
Economic growth  
Brand management  
Barriers to entry  
Fax  
Three Mile Island accident  
George W. Bush  
Cambridge Energy Research Associates  
Macroeconomics  
Customer relationship management  
Loving Every Minute (album)  
Wave function  
Drought

Is the country facing sustained tight markets?:

In general, capacity margins have been falling over the last few years as utilities have refrained from building baseload capacity, and others have focused on developing primarily peaking capacity. The last large building boom of coal plants ended in early 1970's and nuclear boom dropped sharply after Three Mile Island incident in 1979. During 1990's, projected capacity margins have fallen from the 15 - 20 % toward 10% and below. Outside of the California situation, NYC also poses a potential risk for this summer.

If general, however, by 2002 - 2003, the amount of capacity proposed in each region more than covers normal load growth for meeting peak hour demand. Remaining question relates to performance of existing coal and nuclear stacks, also what happens during periods of persistent drought.

Will volatility and prices remain high or lessen?:

To the extent that more capacity becomes merchant oriented focusing in marginal cost economics, and transmission congestion persists, increased volatility will continue, especially in ISO/pool type environments. Prices will reflect primary fuel dynamics, especially the interplay between gas and oil. The case for lower prices will reflect an overbuild scenario beyond this year, coupled with slowing economic activity and low incidence of extreme weather events.

How fast can generation be added and what returns should be expected?

The variability in development of greenfield capacity depends on time to permit at state or local levels. Construction time is fairly constant. In general 18 months is reasonable time frame from concept to first fire. To the extent that power plants are project financed with minimum 30 equity, returns should be consistent with other comparable project financed opportunities available to fund managers. We do not expect any more fully debt financed facilities in the near term.

Can companies make major profits owning generation long term?



# Label selection and ranking

Country Weather  
Limiting magnitude  
Marginal cost  
Economic growth  
Brand management  
Barriers to entry  
Fax  
Three Mile Island accident  
George W. Bush  
Cambridge Energy Research Associates  
Macroeconomics  
Customer relationship management  
Loving Every Minute (album)  
Wave function  
Drought



Mean Average Precision

Is the country facing sustained tight markets?:

In general, capacity margins have been falling over the last few years as utilities have refrained from building baseload capacity, and others have focused on developing primarily peaking capacity. The last large building boom of coal plants ended in early 1970's and nuclear boom dropped sharply after Three Mile Island incident in 1979. During 1990's, projected capacity margins have fallen from the 15 - 20% toward 10% and below. Outside of the California situation, NYC also poses a potential risk for this summer.

If general, however, by 2002 - 2003, the amount of capacity proposed in each region more than covers normal load growth for meeting peak hour demand. Remaining question relates to performance of existing coal and nuclear stacks, also what happens during periods of persistent drought.

Will volatility and prices remain high or lessen?:

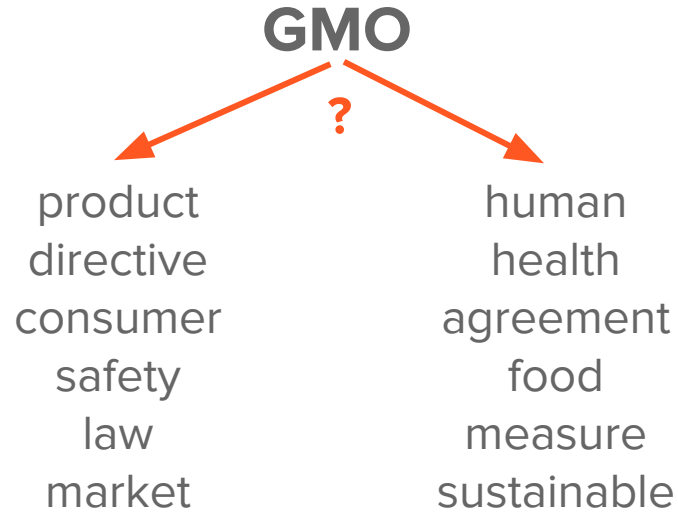
To the extent that more capacity becomes merchant oriented focusing in marginal cost economics, and transmission congestion persists, increased volatility will continue, especially in ISO/pool type environments. Prices will reflect primary fuel dynamics, especially the interplay between gas and oil. The case for lower prices will reflect an overbuild scenario beyond this year, coupled with slowing economic activity and low incidence of extreme weather events.

How fast can generation be added and what returns should be expected?

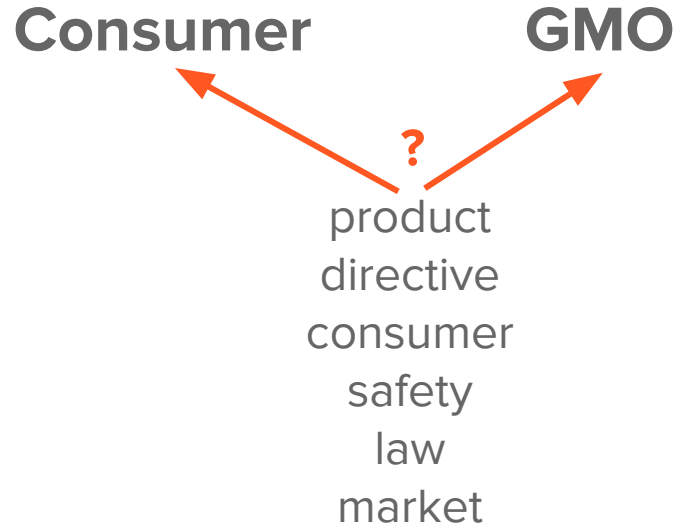
The variability in development of greenfield capacity depends on time to permit at state or local levels. Construction time is fairly constant. In general 18 months is reasonable time frame from concept to first fire. To the extent that power plants are project financed with minimum 30 equity, returns should be consistent with other comparable project financed opportunities available to fund managers. We do not expect any more fully debt financed facilities in the near term.

Can companies make major profits owning generation long term?

# Topic-label relation



# Topic-label relation



# First results - 50 docs labeled

	<b>Avg number of label selected</b>	<b>Recall on user selection</b>	<b>Docs with 0 annotations</b>
EuroParl	4	0.88	2
EnronCorpus	4	0.91	7
ClintonCorpus	4	0.95	1

# First results - ranking

	MAP TF-IDF Ranking	MAP LLDA Ranking
RandomBaseline	0.30	0.30
EuroParl	0.51	<b>0.54</b>
EnronCorpus	0.40	<b>0.41</b>
ClintonCorpus	0.48	<b>0.52</b>

# Conclusion

Entity-labels could **drastically improve** topic interpretability.

However it is necessary to **always evaluate** them.

We will release:

- The pipeline for labeling topics with entities
- A tool for evaluating each step of the process

# Building Entity-based Collections of Global Events



**Federico Nanni**, Simone Paolo Ponzetto and Laura Dietz

Data and Web Science Group

University of Mannheim

{federico,simone,dietz}@informatik.uni-mannheim.de

# Global events





# Global events



Identified by a common name: the **Wall Street Crash of 1929**.

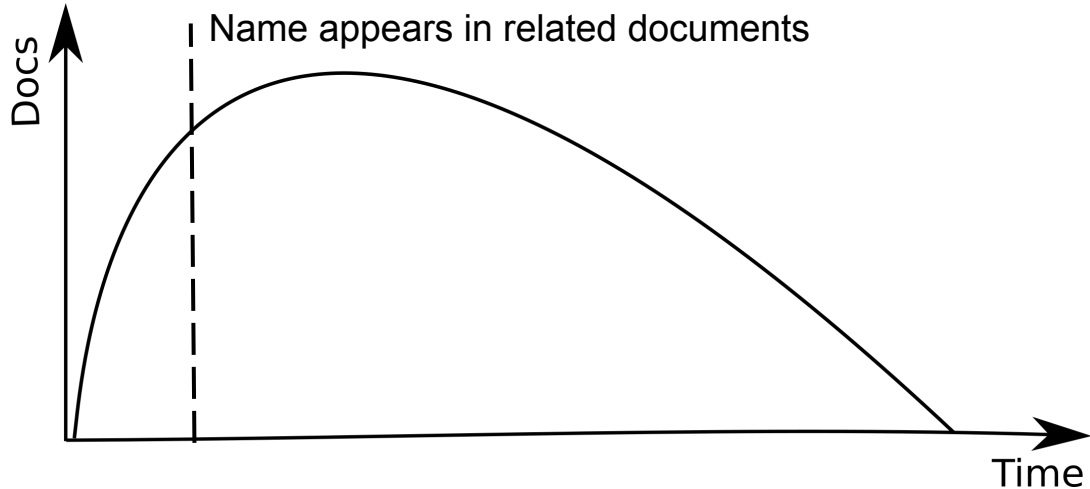
Also known as: the **Black Tuesday**, the **Great Crash**, or the **Stock Market Crash**.

# Global events



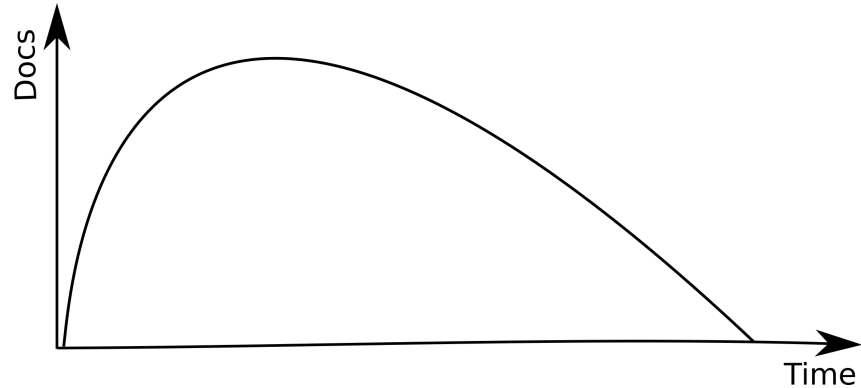
Identified by a common name: the **Wall Street Crash of 1929**.

Also known as: the **Black Tuesday**, the **Great Crash**, or the **Stock Market Crash**.



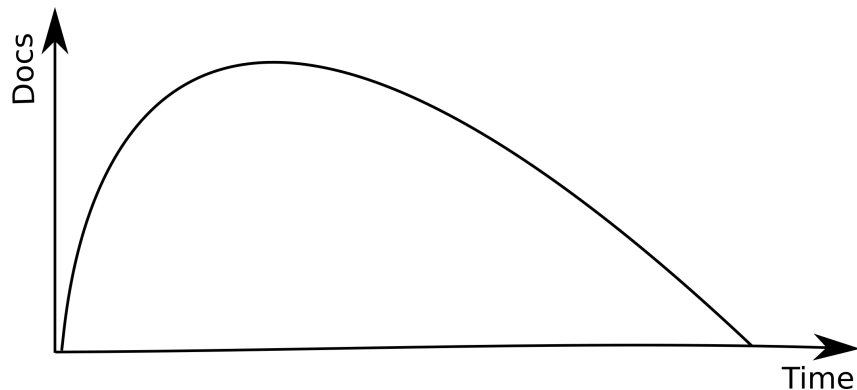
# Retrospective analyses

Studying how **events are perceived by society** is a fundamental research task for humanists and social scientists.



# Retrospective analyses

Studying how **events are perceived by society** is a fundamental research task for humanists and social scientists.



# Solutions: doc filtering, focused crawling

Collect the documents that contain all query words (e.g. Kedzie et al., 2014).

E.g. query: "wall street crash tuesday stock market 1929"

The FIRST with the LATEST Full wired Press leased wire

People's Paper Santa Ana Daily Register Orange County

VOL. XXIV. NO. 285. SANTA ANA, CALIFORNIA, MONDAY, OCTOBER 28, 1929. 20 PAGES 3c Per Copy. 65c Per Month

## BILLIONS LOST AS STOCKS CRASH

### Pantages In Jail Awaits Sentence To State Prison

#### GUILTY SAYS JURY AFTER MANY HOURS

Theater Man Declares He Got Raw Deal and Did Not Ever Have Chance

MISS PRINGLE GLAD In Statement Says Not Sure Whether She Will Continue Her Stage Career

OF ANGELES, Oct. 28.—(AP)—Dressed in the blue uniform of the Los Angeles county jail, Pantages was "banged" today and sentenced to state prison for about his conviction that night of an assault upon Katherine Pringle.

"I believe I got a raw deal," Pantages declared as he walked away with the jury from the court house in a van.

"I don't know what I did to her," he said, "but I know she was a fine girl."

The indictment returned in the case against Miss Pringle was for the assault upon her in the lobby of the Grand Hotel.

The indictment returned in the case against Miss Pringle was for the assault upon her in the lobby of the Grand Hotel.

The indictment returned in the case against Miss Pringle was for the assault upon her in the lobby of the Grand Hotel.

#### SENATE "FARM BLOC" STARTS BATTLE FOR HIGHER TARIFFS

Friend Brings Man To Jail On Rum Charge

#### BINGHAM TAKES EXCEPTION TO SENATE LOBBY

Charges Committee With Unfair Political Tactics During Probe

#### WESTERN REPUBLICANS GO TO WASHINGTON TO INCREASE RATES

Washington, Oct. 28.—(AP)—Changing the Senate tariff investigation committee will make possible higher rates on many goods, including iron and steel, according to a statement issued today by the committee.

The committee, headed by Senator Bingham, Republican, Oklahoma, will make a study of the tariff on iron and steel, and will recommend a 10 per cent increase on the iron and steel tariff.

The committee also will make a study of the tariff on copper, and will recommend a 10 per cent increase on the copper tariff.

The committee also will make a study of the tariff on zinc, and will recommend a 10 per cent increase on the zinc tariff.

#### THREE U. S. COMMITMENTS ARE CRASH VICTIMS

Auto Planes Into Engine As Routers Return from Big Football Game

Los Angeles, Calif., Oct. 28.—(AP)—Two deaths in a Tulsa bus

#### JURY BELIEVES HER STORY

Banking Support Unable to Prevent Break Accompanied by Wild Trading

#### TICKERS FAR BEHIND

General Electric, U. S. Steel, Other Leading Issues Caught In Crash

#### MAINSTAYS OF MARKET IN PLUNGE

Banking Support Unable to Prevent Break Accompanied by Wild Trading

#### TICKERS FAR BEHIND

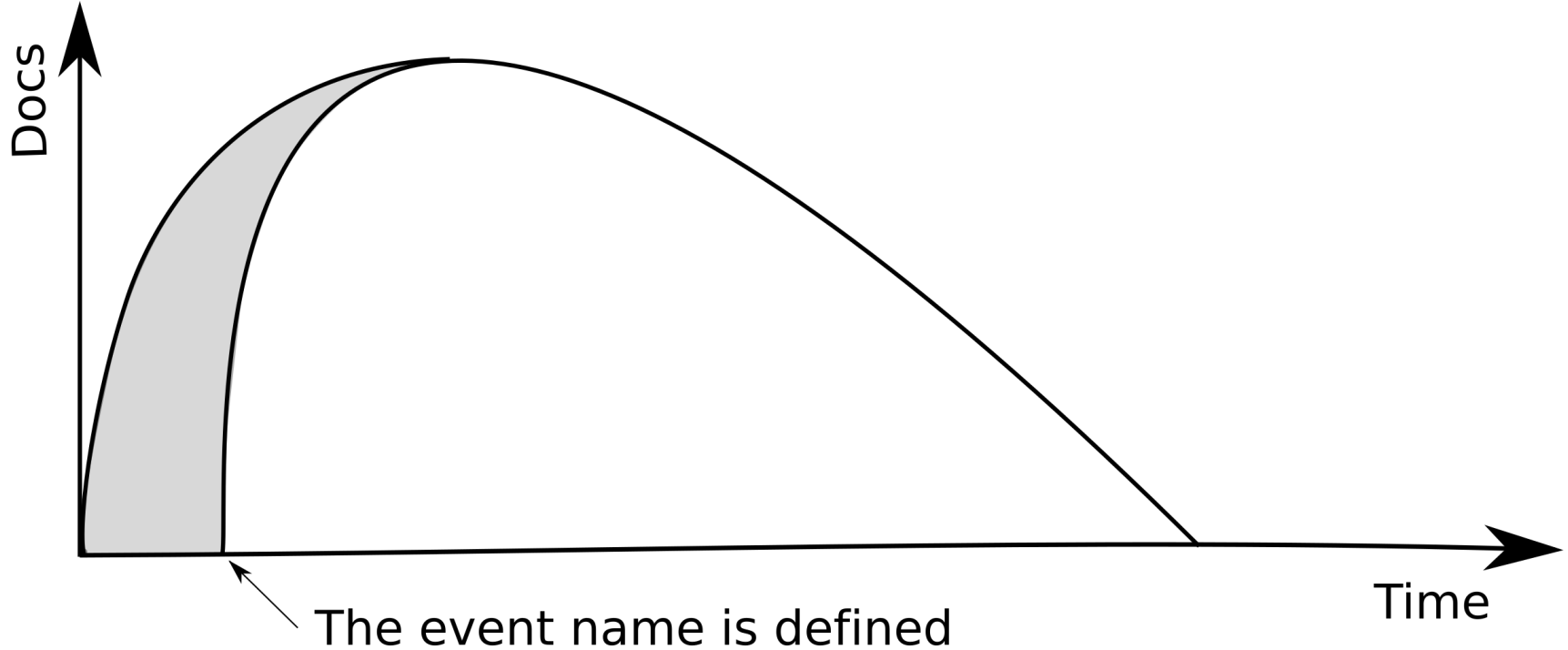
General Electric, U. S. Steel, Other Leading Issues Caught In Crash



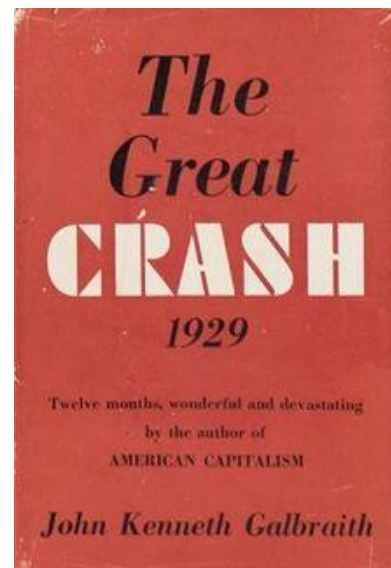
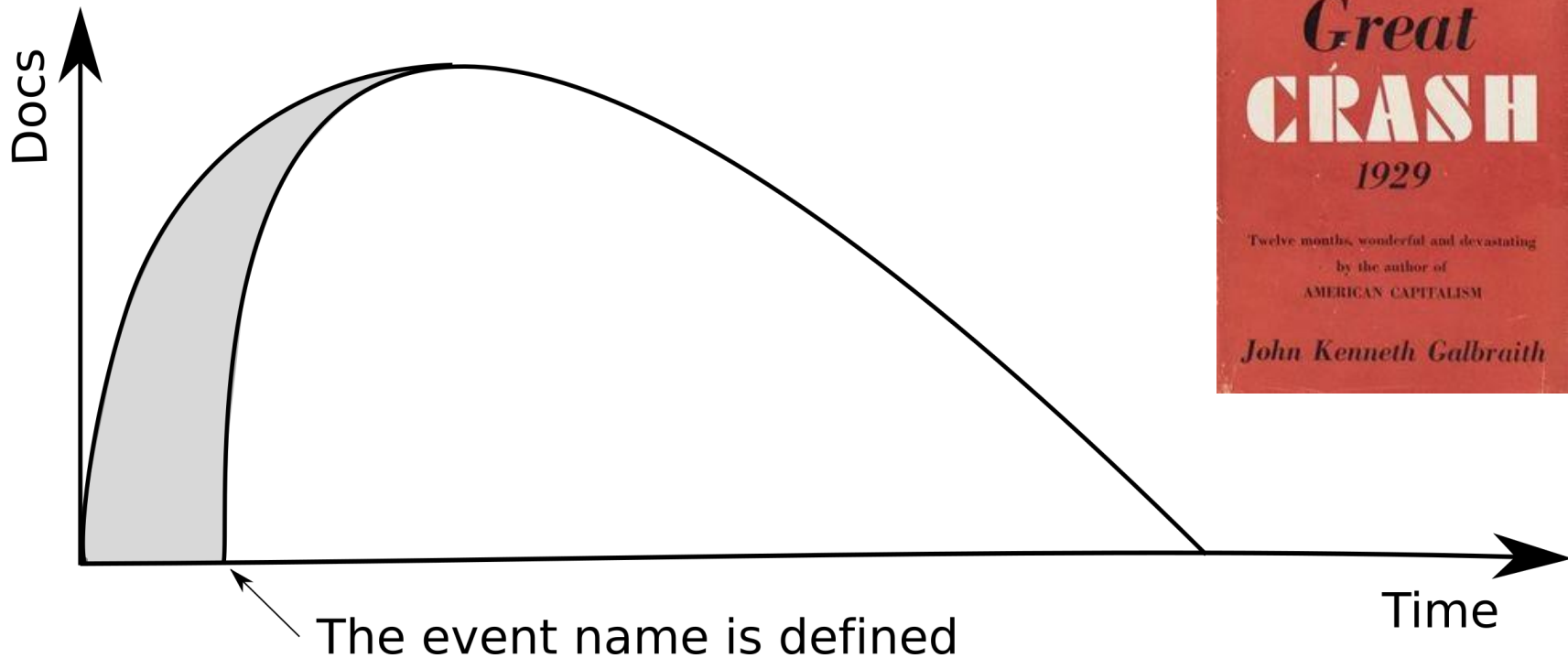
#### CHRISTIANITY MUST DESTROY WAR OR WAR WILL OVERTHROW CHURCH, SAYS PASTOR SCHROCK

"WAR is a problem in arithmetic," said the Rev. Perry F. Schrock, in his sermon at the First Christian church last night. He selected verse out of a

# Missing the early stages



# Missing the early stages



# A contemporary example



**@ReallyVirtual**

Sohaib Athar

Helicopter hovering above Abbottabad at 1AM (is a rare event).

18 hours ago via TweetDeck



# A contemporary example

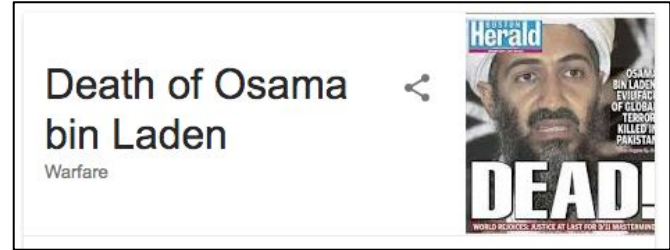


**@ReallyVirtual**

Sohaib Athar

Helicopter hovering above Abbottabad at 1AM (is a rare event).

18 hours ago via TweetDeck



# Brief idea for finding early stories

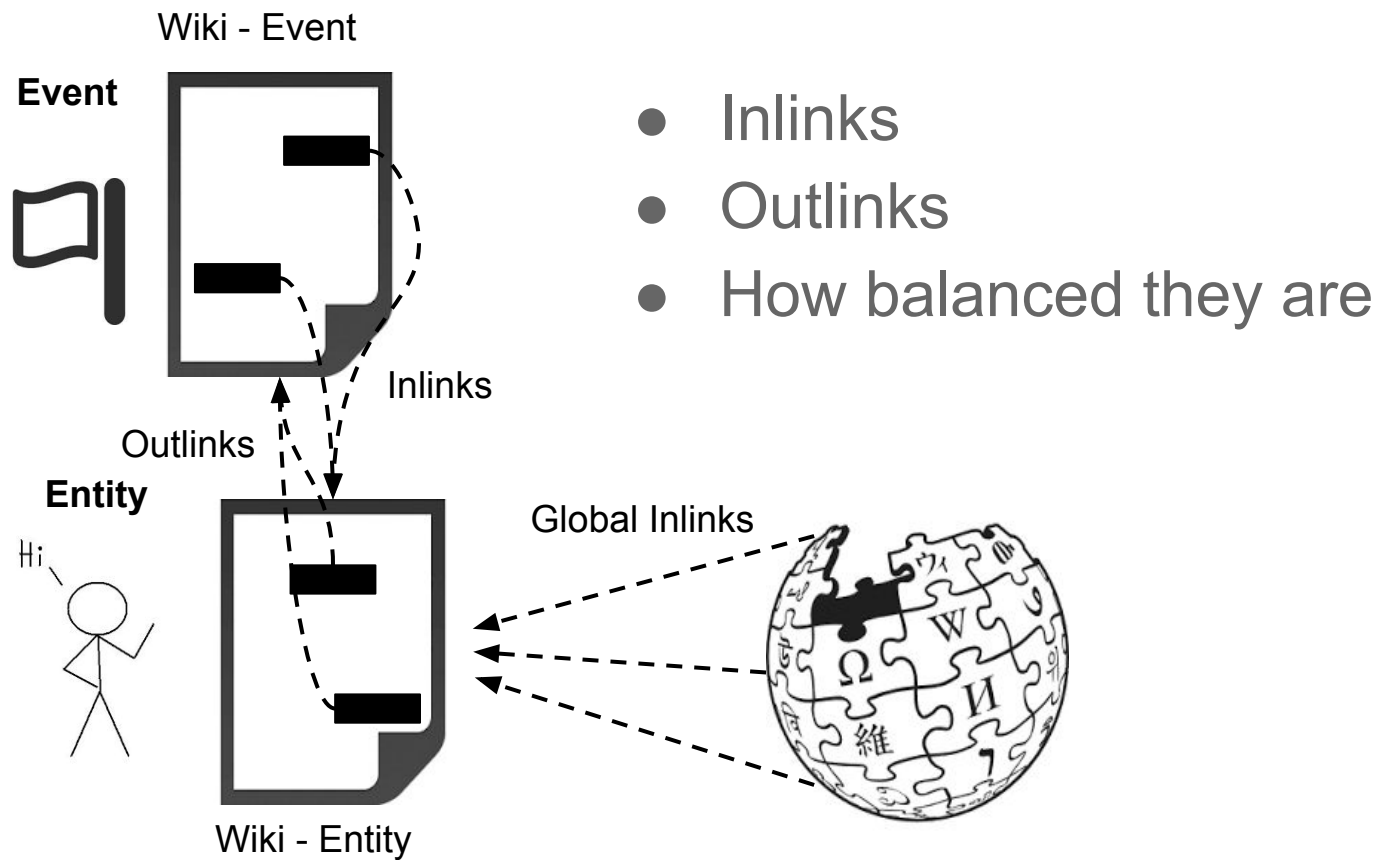
Given a named event:

1. Identify **related entities**
2. Retrieving **text passages** with entity in context
3. Building a **language model** for each entity

=> Entity event-query expansion!

**How do we identify related entities?**

# Simple graph-based entity relatedness



# Identifying related entities

---

**System**

---

Stics

---

Wiki2Vec

---

WikipediaRanking

---

Eventipedia (our)

---

**Gold standard:** 10 global events (between 2012 and 2014).

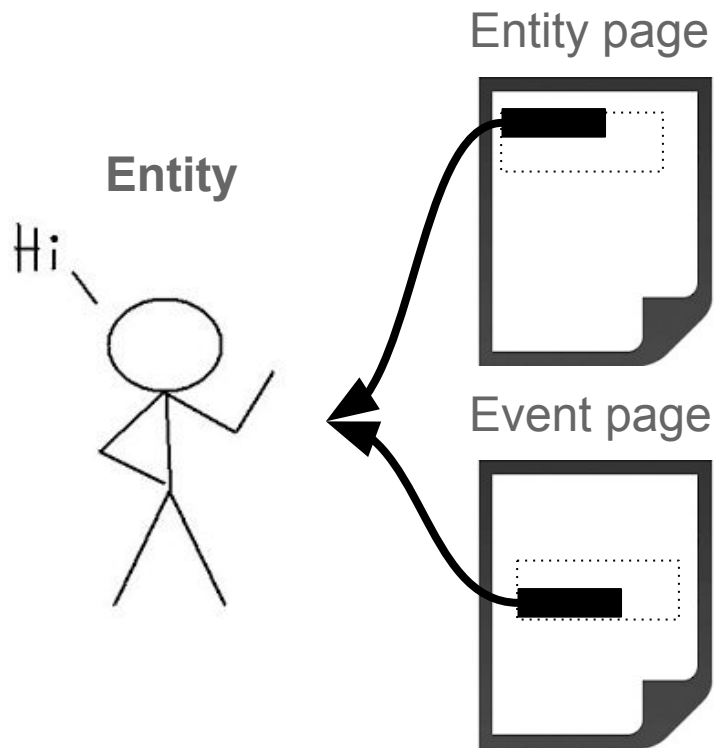
Human annotators assess the relevance of retrieved entities on a binary scale.

# Identifying related entities

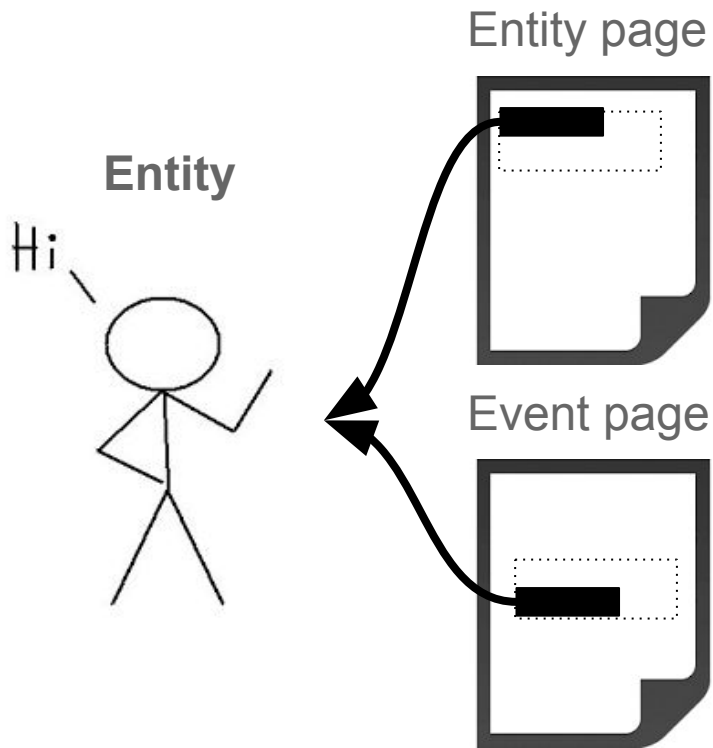
System	MAP@10	Micro-Prec@10
Stics	0.54 ± 0.07	0.59 ± 0.05
Wiki2Vec	0.59 ± 0.11	0.64 ± 0.04
WikipediaRanking	0.66 ± 0.09	0.71 ± 0.05
Eventipedia (our)	<b>0.74 ± 0.05</b>	<b>0.81 ± 0.04</b>

**Gold standard:** 10 global events (between 2012 and 2014).  
Human annotators assess the relevance of retrieved entities on a binary scale.

# Retrieving explanatory passages



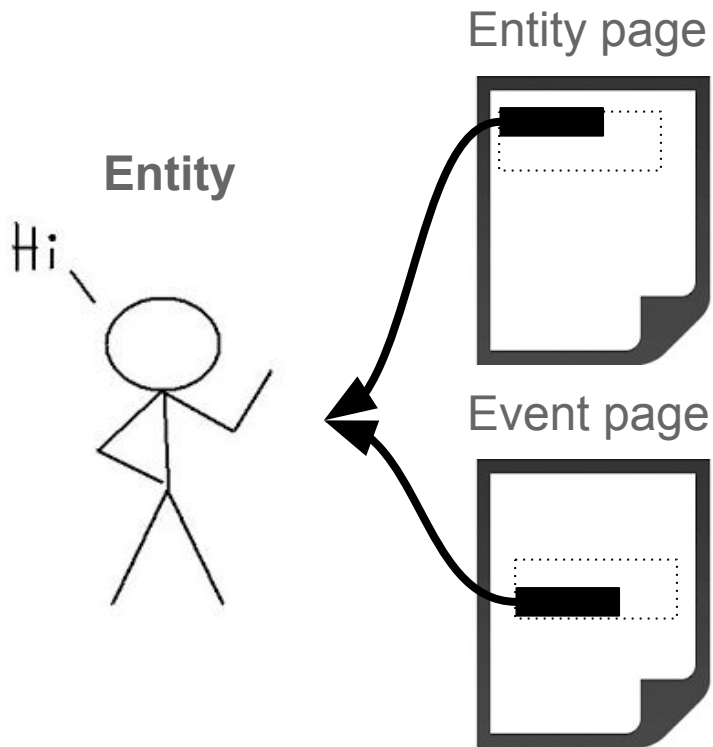
# Retrieving explanatory passages



“**Russia**, also officially known as the Russian Federation, is a sovereign state in northern Eurasia. It is a federal semi-presidential republic.”



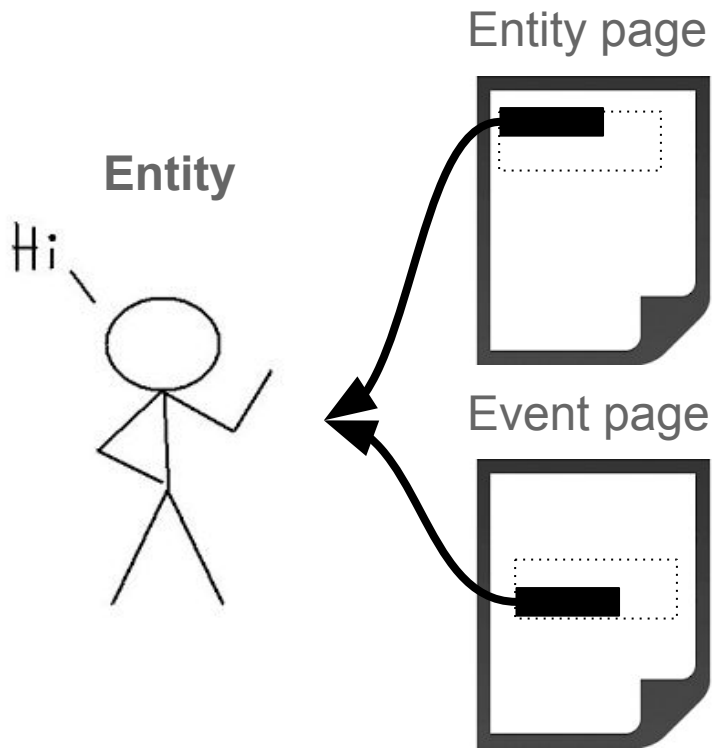
# Retrieving explanatory passages



“**Russia**, also officially known as the Russian Federation, is a sovereign state in northern Eurasia. It is a federal semi-presidential republic.”

“From the early stages, the Syrian government received technical, financial, military and political support from **Russia**, Iran and Iraq.”

# Retrieving explanatory passages



	<b>% Good</b>
Entity page snippet	45%
Event page snippet	<b>68%</b>

# Conclusions

A simple **entity relatedness approach** is useful for identifying entities related to an event.

**Entity** needs to be considered **in the context of the event** for building the language model.

Next step: evaluate “early stories detection”.

# Topic-Based Analysis of Political Positions in US Electoral Campaigns

---

**Federico Nanni**, Caecilia Zirn, Goran Glavas  
Jason Eichorst and Simone Paolo Ponzetto

Data and Web Science Group  
Collaborative Research Center SFB 884  
University of Mannheim

{federico,caecilia,goran,simone}@informatik.uni-mannheim.de

# Introduction

- Campaigns designed to **inform voters** on candidates ideas

# Introduction

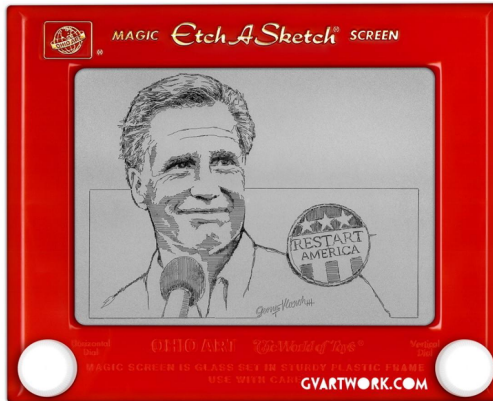
- Campaigns designed to **inform voters** on candidates ideas
- **Converging on a position** is an interesting process because candidates must satisfy
  - Interests of party members and groups during primaries
  - Interests of voters in general elections

# Introduction

- Campaigns designed to **inform voters** on candidates ideas
- **Converging on a position** is an interesting process because candidates must satisfy
  - Interests of party members and groups during primaries
  - Interests of voters in general elections
- Often, there is a **notable shift** in candidate positions between primaries and general elections

# Mitt Romney's Etch-a-Sketch

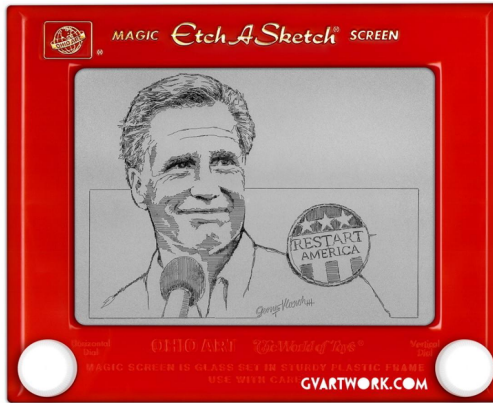
“Everything changes. It’s **almost like an Etch-a-Sketch**. You can kind of shake it up and restart all over again.”





# Mitt Romney's Etch-a-Sketch

“Everything changes. It’s **almost like an Etch-a-Sketch**. You can kind of shake it up and restart all over again.”



Starting Hypothesis: For certain **topics** these changes might be **more prominent** than for others

# Topics

Coarse-grained topics from **Comparative Manifesto Project**:

1. External relations
2. Freedom and Democracy
3. Political System
4. Economy
5. Welfare and Quality of Life
6. Fabric of Society
7. Social Groups

# Speech dataset

Speeches **manually labeled** with topics at paragraph level.

Two annotators, moderate IAA of 0.55 (Cohen's kappa).

Contains altogether **9 speeches** from 2008, 2012, and 2016 elections (around 1k paragraphs).

# Topical classification

Topic classifier: SVM with **lexical** and **semantic** features.

Two experimental settings:

1. **Domain transfer setting**: training the model on manifestos, testing on speeches
2. **Pure speeches setting** – Folded cross-validation on speeches

Baseline model: SVM with bag-of-words features.

# Topical classification

Classification results (in terms of F1 score):

- Domain transfer setting (training on manifestos): 36.2%
- Baseline model (BoW SVM, CV on speeches): 71.2%
- Pure speech setting (CV on speeches): **78.6%**

Conclusion: **Transfer learning doesn't work** between different domains of political text.

# Political scaling

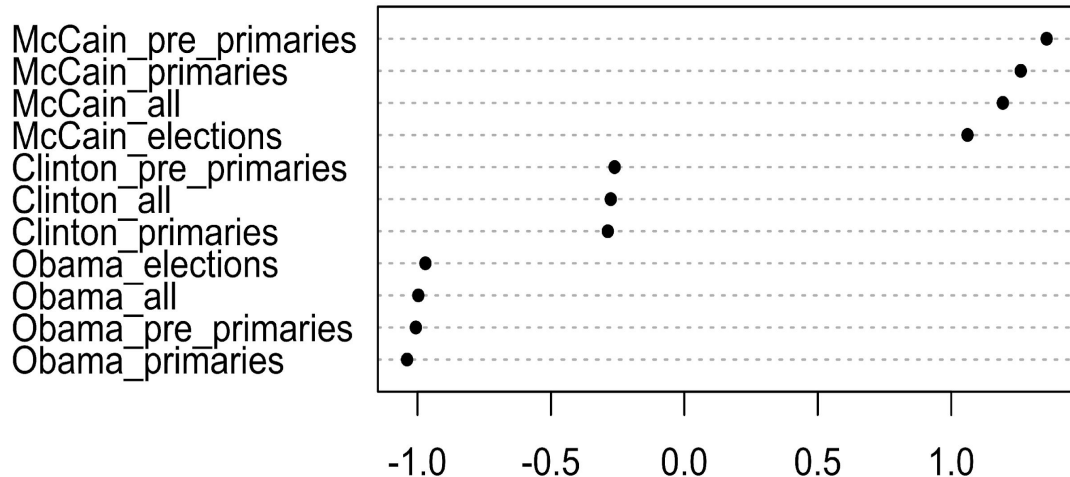
Best performing model made **topic predictions** on the speeches.

For each candidate, we concatenate all paragraphs with **same topic and same phase** (pre-primaries, primaries, elections).

Finally, we feed each phase-topic slice to the **Wordfish** tool.

# Coarse-grained analysis (2008)

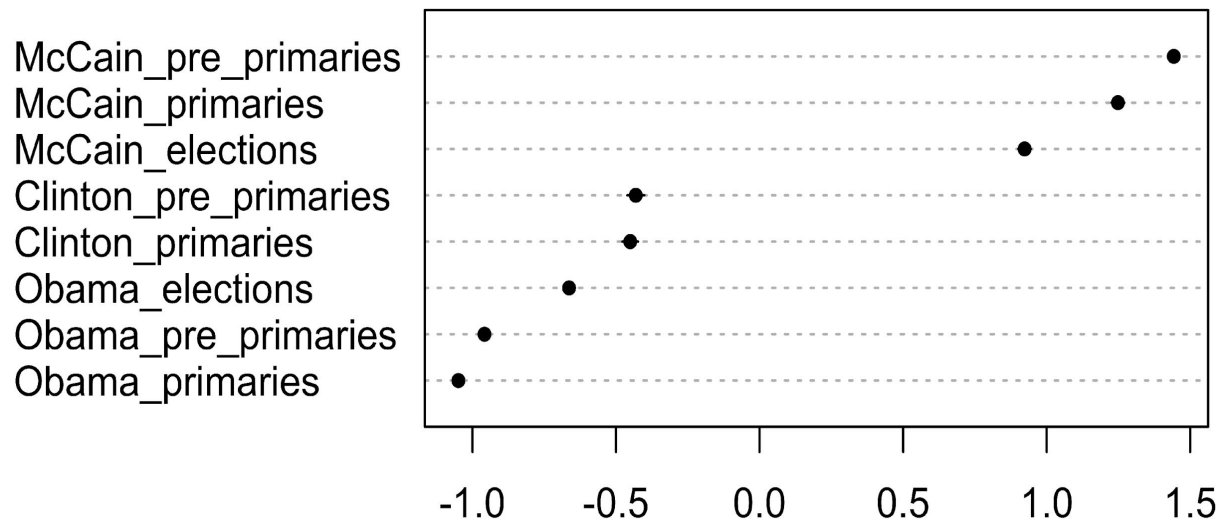
First we analyzed the **general positions** (Wordfish on entire speeches), as a baseline for analysis.



# Fine-grained analysis (2008)

Next, we analyzed topic-specific positions:

## External Relations

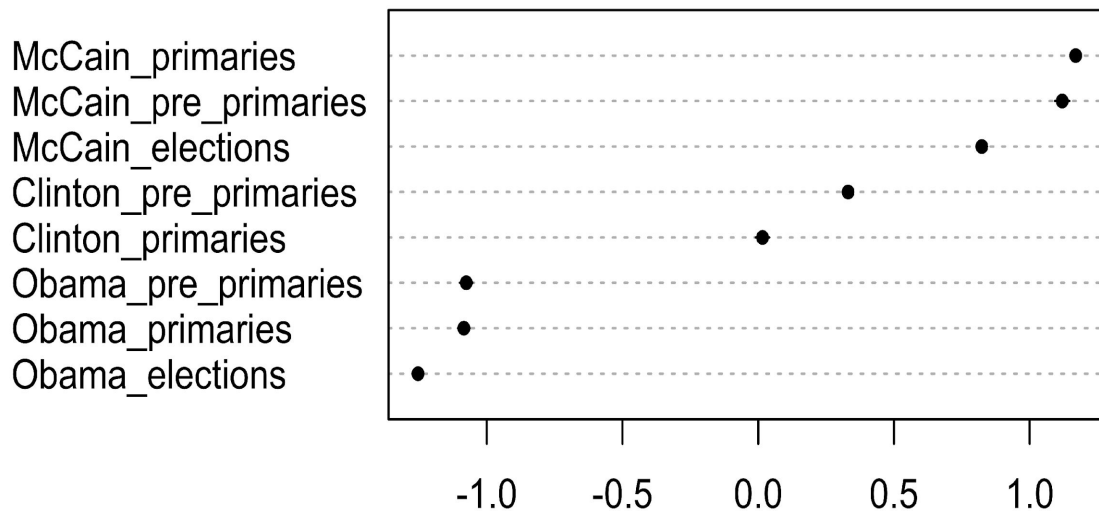




# Fine-grained analysis (2008)

Next, we analyzed topic-specific positions:

## Welfare and Quality of Life



# Conclusion

Topic-based position analysis could offer **fine-grained perspectives** on political campaigns.

Important **evaluating** approaches for the task.

We are working on a **Python** implementation of the pipeline!

# Final wrap-up

Our approach to DH has a **strong NLP approach**.

Focus on:

1. Hypothesis-testing analyses
2. Tool evaluation
3. Train researchers with highly interdisciplinary profiles

**Thanks!**

## **Data and Web Science Group**

<http://dws.informatik.uni-mannheim.de>

{federico,simone}@informatik.uni-mannheim.de

## **Historisches Institut**

<http://www.geschichte.uni-mannheim.de>

hiram.kuemper@uni-mannheim.de