

# Creating an Austrian Language Polarity Dictionary with the Crowd

Thomas E. Kolb, Katharina Sekanina, Bettina M. J. Kern,  
Julia Neidhardt, Andreas Baumann and Tanja Wissik

NLP Seminar @ TU Wien

Funded by:



**Stadt  
Wien**

Grant number:  
MA7-737909/19

**ÖAW**

ÖSTERREICHISCHE  
AKADEMIE DER  
WISSENSCHAFTEN

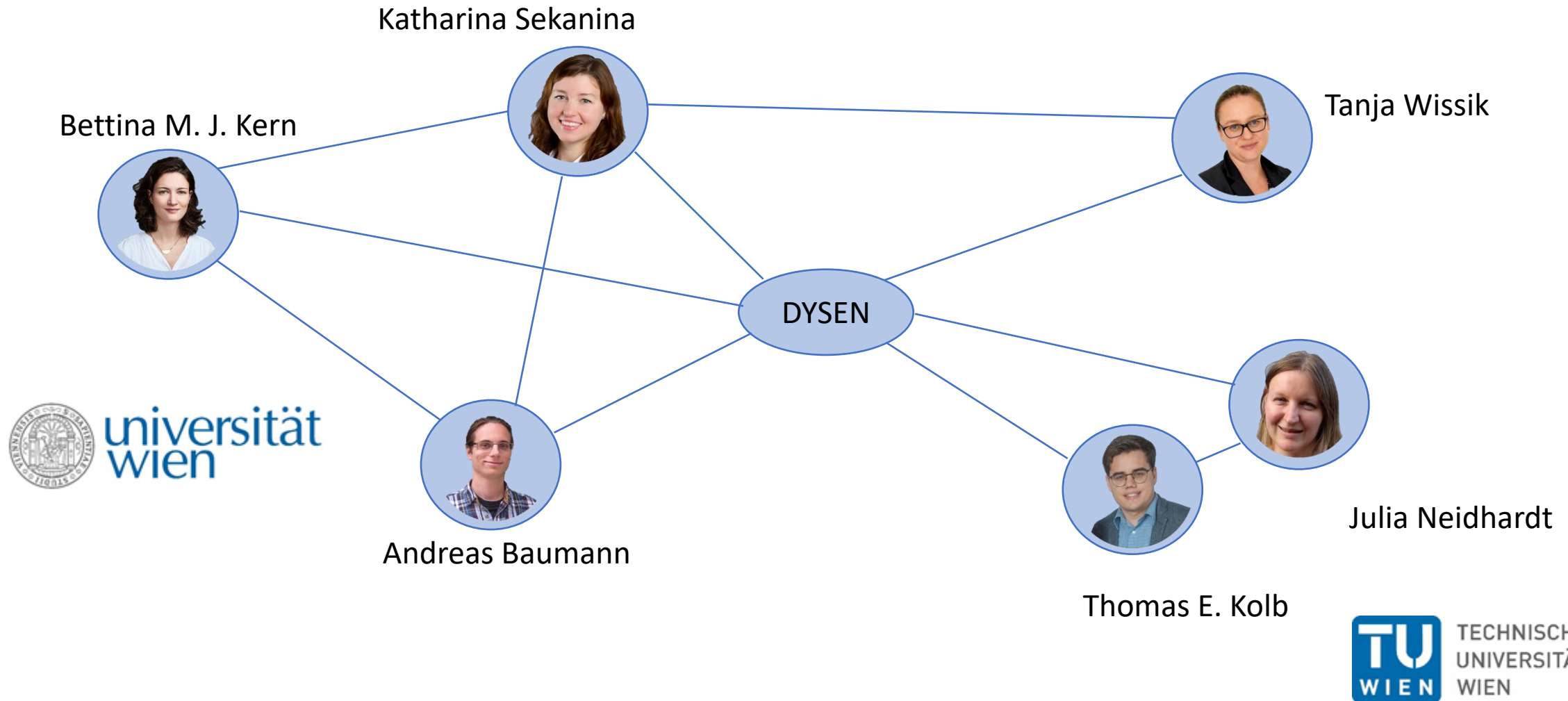


universität  
wien



TECHNISCHE  
UNIVERSITÄT  
WIEN

# Project Team



# DYSEN Project

## Dynamic Sentiment Analysis as Emotional Compass for the Digital Media Landscape

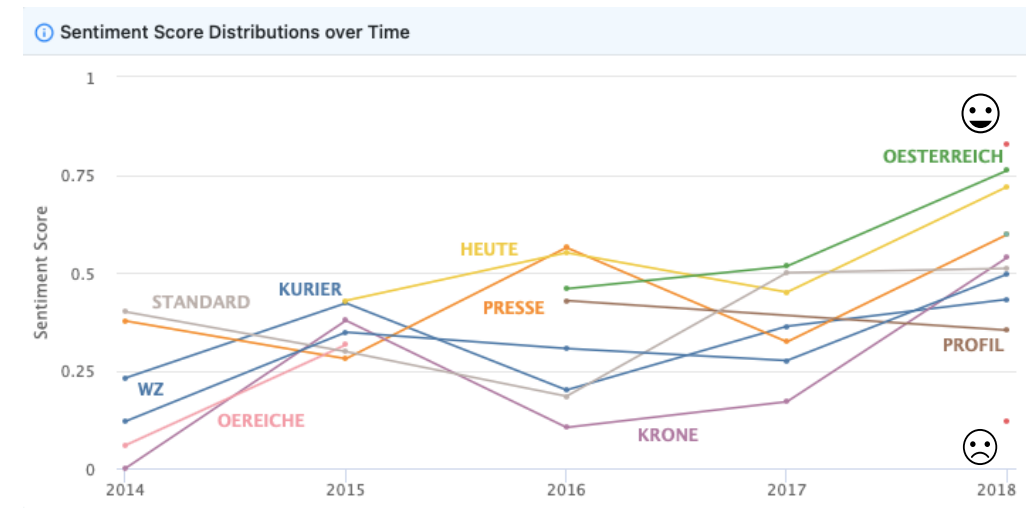
 **Research question:** How do print media report about the Viennese politicians?

 **Aim of the project:** Develop a tool that can detect change of emotional polarization of politicians in Austrian Newspapers

Funded by:



<https://dysen.acdh.oew.ac.at/dysen/>

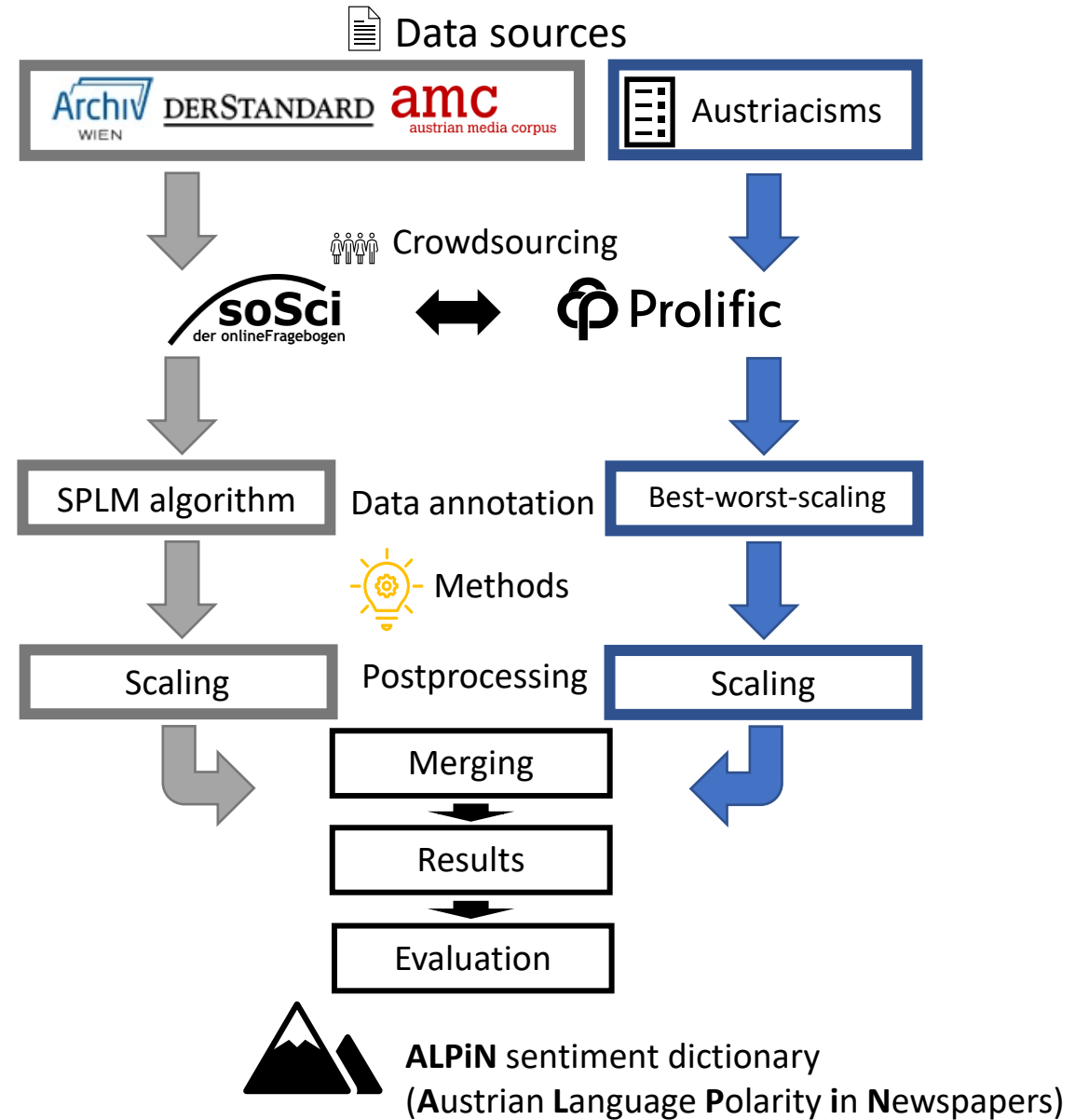


# Problem Statement

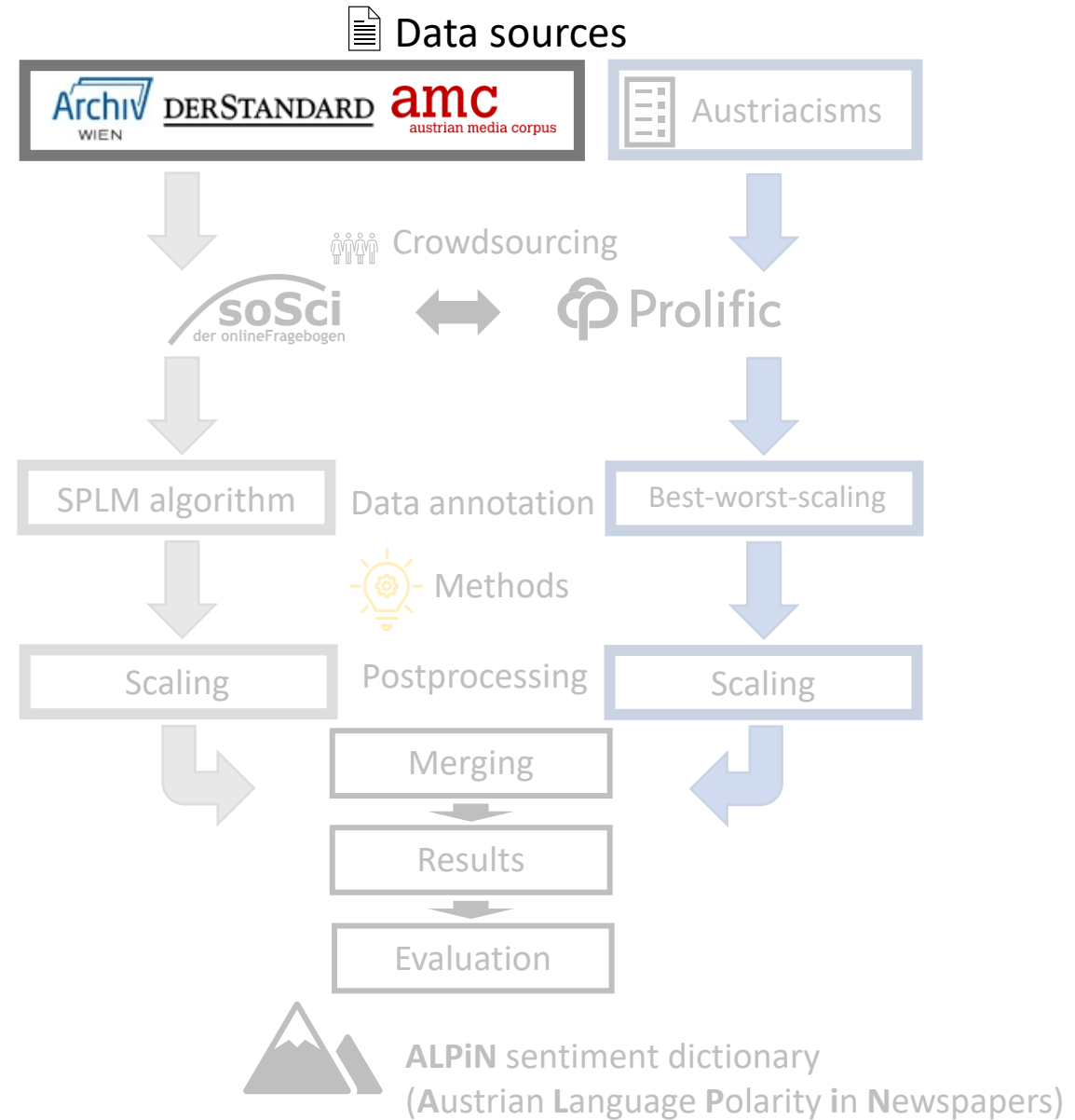
✘ There is no sentiment dictionary for Austrian German in this domain

🎯 **Goal:** Create Austrian German language resource in the domain of news media and politics

# Content



# Content



# Data sources: Viennese politicians

Retrieved from the politician archive of Vienna (POLAR) of the Vienna City and State Archives.

**Definition of “Viennese politician”:** All members of the

- Vienna City Council
- Vienna City Senate
- Vienna State Parliament
- Vienna State Government

who were active between the 13<sup>th</sup> and the 20<sup>th</sup> parliamentary term (1983 to 2020) = **487 politicians**

# Data sources: DERSTANDARD (1 Million Posts Corpus)

(Schabus et al., 2017)

- Forum posts of 12 months from 2015 to 2016
- 3599 posts labelled for sentiment by professional forum moderators

	ID_Post	Body	Category		
0	3326	Top qualifizierte Leute verdienen auch viel.	SentimentNeutral		
1	5321	Gott sei dank ist für sie eine Umfrage alles, ...	SentimentNegative	SentimentNeutral	1865
2	5590	Sorry, aber die FPÖ tut eigentlich gar nichts ...	SentimentNeutral	SentimentNegative	1691
3	6015	Weil es dein meisten Leuten verständlicherweis...	SentimentNegative	SentimentPositive	43
4	8213	Na wer weis was da vorgefallen ist...	SentimentNeutral		



# Data sources: **amc** Austrian Media Corpus

austrian media corpus (Ransmayr et al., 2017)

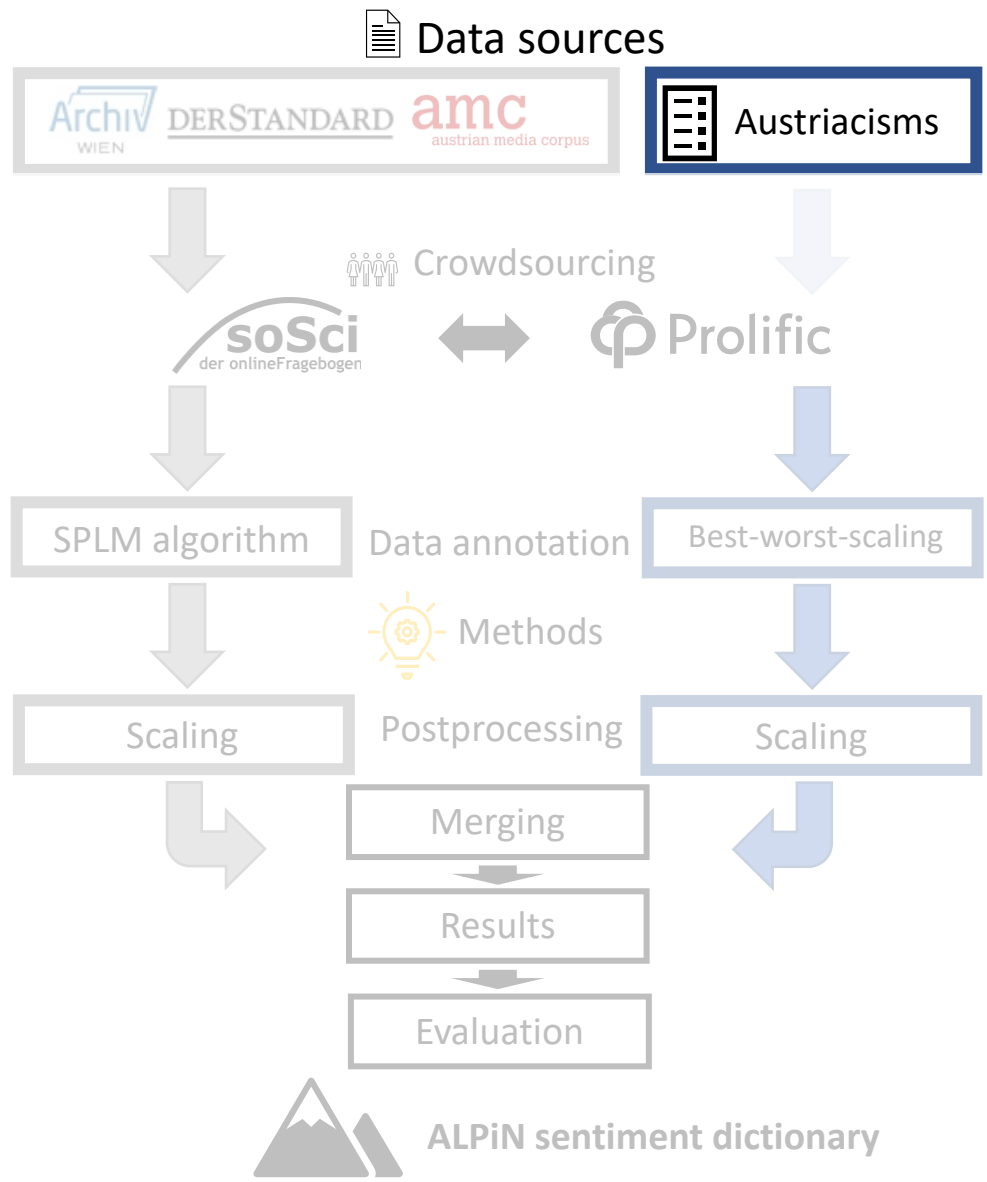
- Contains Austrian print media
- Preprocessed and linguistically annotated
- Yearly updates

## Our data:

- We use print media related to Vienna between 1996 and 2017
- Excluded APA<sup>1</sup> and OTS<sup>2</sup> articles ("*Presseaussendungen*")
- Text snippets of 60 tokens around the politicians' name were extracted

1. <https://apa.at/>

2. <https://www.ots.at/>



# Data sources: Austriacisms

## Based on:

- „Variantenwörterbuch des Deutschen“ (VWB; words specific to Austria) (Bickel et al.,2015)
- Austriacism list of Wikipedia<sup>1</sup>

## Restrictions:

The combined list is manually checked and cleaned up by linguist experts of our project team = **538 remaining words** (*pos tagged with: noun, adjective, verb*)

1. [https://de.wikipedia.org/wiki/Liste\\_von\\_austriacismen](https://de.wikipedia.org/wiki/Liste_von_austriacismen)

# Crowdsourcing

**Aim:** Attain sentiment annotations from the crowd for:

**amc**  
austrian media corpus

Text snippets from the amc data that mention Viennese politicians



Austriacism list

**By using:**



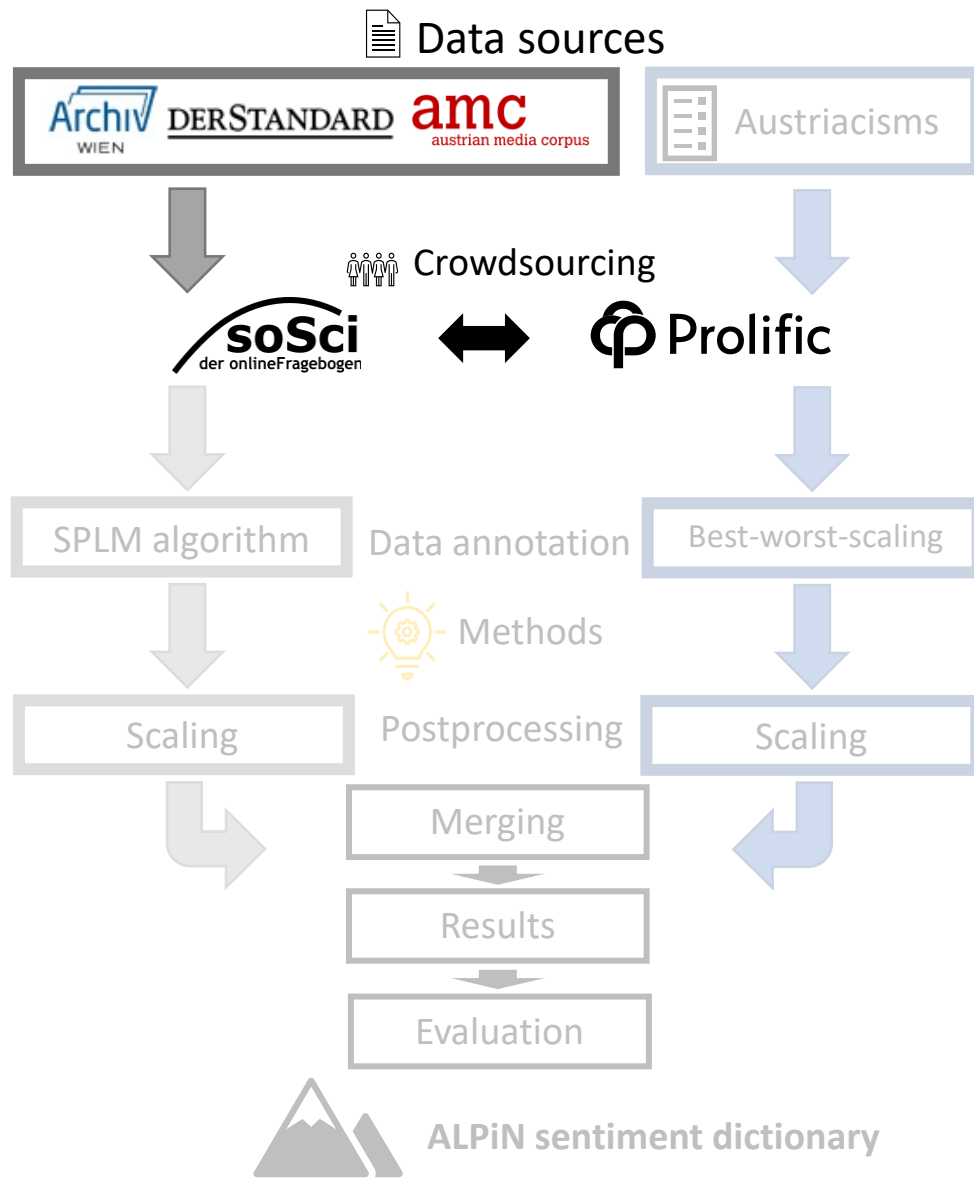
SoSci Survey<sup>1</sup>: platform for designing surveys



Prolific<sup>2</sup>: platform to find research participants who fill out the survey

1. <https://www.soscisurvey.de/>

2. <https://www.prolific.co/>



# Crowd sourcing: **amc** Austrian Media Corpus

austrian media corpus

- Each item labelled  $\geq 3$  times
- Majority vote (equal number per class = rated as neutral)
- Three classes: positive, neutral, negative
- Survey:
  - 100 randomly selected text snippets
  - +24 items for quality control ( $\geq 75\%$  correct)

## **Restricted annotators by:**

- Current Country of Residence (Germany, Austria, Switzerland)
- Nationality (Germany, Austria, Switzerland)
- First Language (German)

# Crowd sourcing: **amc** Austrian Media Corpus austrian media corpus

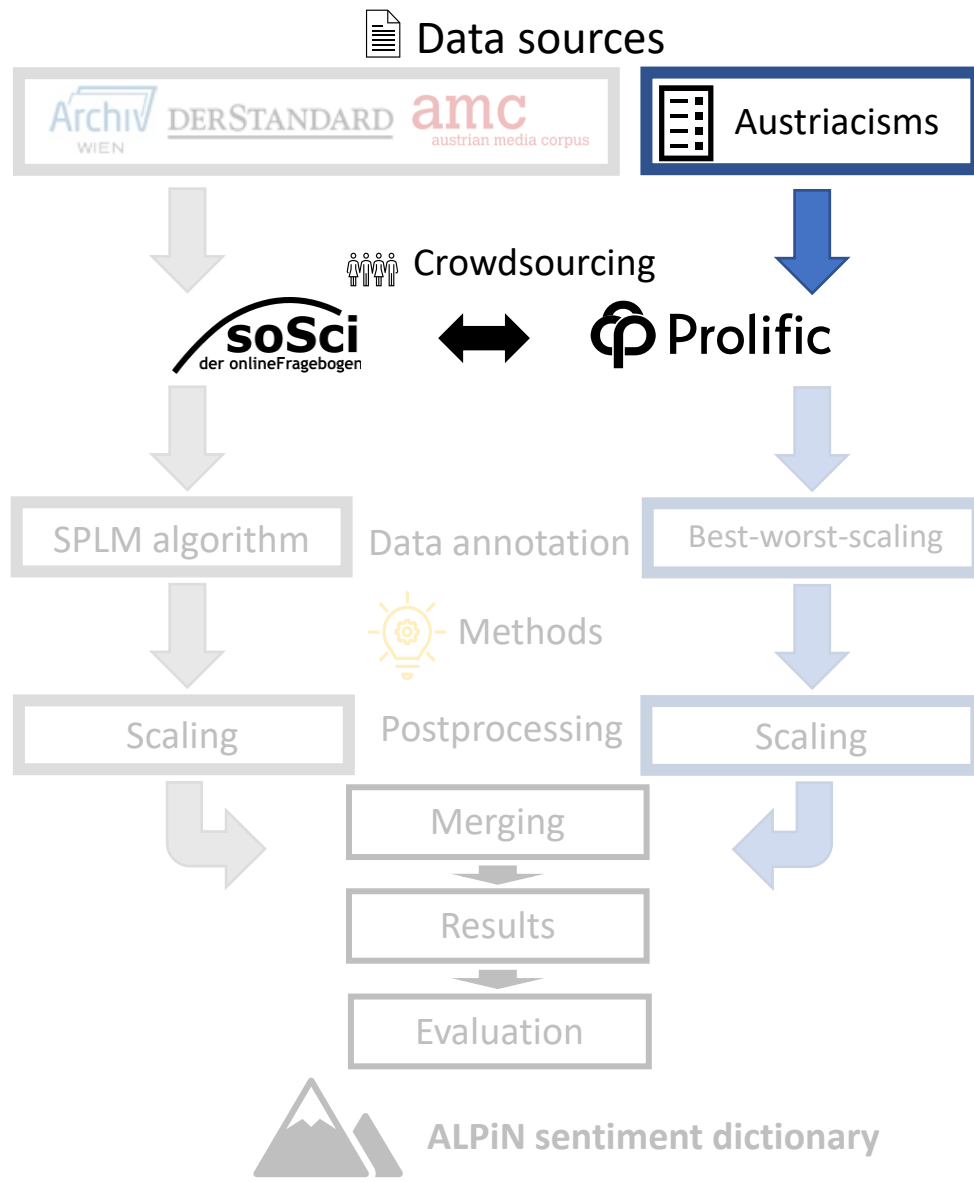
## 1st annotation run (70 annotators after excluding the 14 bad ones)

- 2376 items
  - Fleiss-Kappa: 0.295 (fair inter-annotator agreement)
- |          |      |
|----------|------|
| neutral  | 1202 |
| positive | 598  |
| negative | 576  |

## 2nd annotation run (88 annotators after excluding the 15 bad ones)

- 2970 items
  - Fleiss-Kappa: 0.283 (fair inter-annotator agreement)
- |          |      |
|----------|------|
| neutral  | 1492 |
| positive | 787  |
| negative | 691  |

**Output:** 5346 labelled text snippets including Viennese politicians





# Crowd sourcing: Austriacisms (Survey 1)

## Survey 1 (Preselection):

- Over 1 600 words in total
- 500 words per survey
- +25 words for quality control
- Four options (positive, neutral, negative, unknown)

## Restricted annotators by:

- Current Country of Residence (Austria)
- Nationality (Austria)
- First Language (German)

	negativ	neutral	positiv	unbekannt
lebensbejahend	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Seuche	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Vernaderer	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Gewand	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

# Crowd sourcing: Austriacisms (Survey 2)

## Survey 2:

- Best-worst-scaling (BWS) method<sup>1</sup> (Kiritchenko & Mohammad, 2017)
- 1074 tuples
- 130 tuples per survey
- +20 tuples for quality control ( $\geq 75\%$  correct)
- **Restricted annotators by:**
  - Current Country of Residence (Austria)
  - Nationality (Austria)
  - First Language (German)

5. Bitte wählen Sie das positivste und negativste Wort aus der Liste.

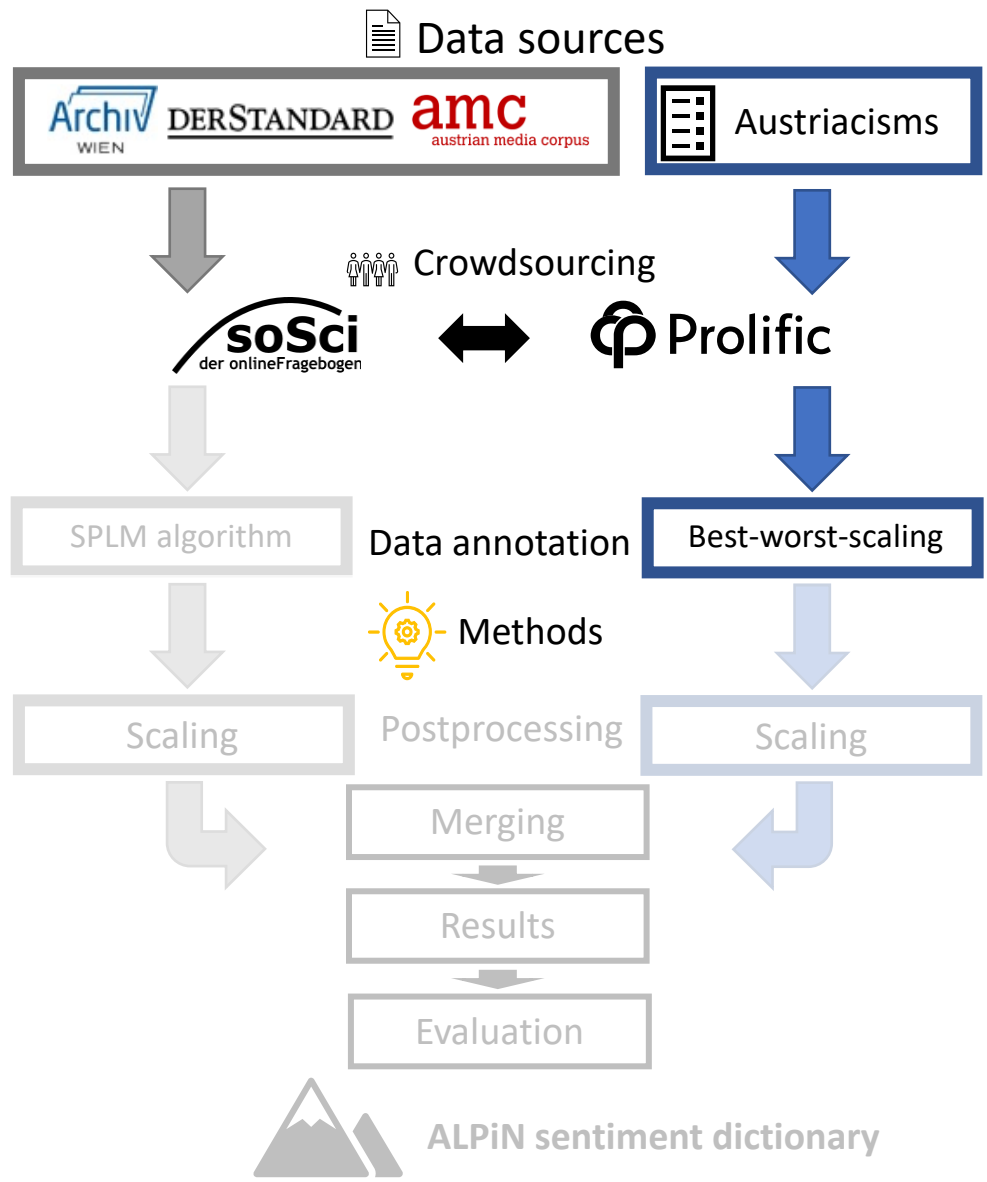
Ohrwaschel
waschelnass
großgoschert
Sanktus

am positivsten
am negativsten

# Crowd sourcing: Austriacisms

- 34 annotators after excluding the 6 bad ones
- **Output:** 4417 tuples (BestItem, WorstItem)

	Item1	Item2	Item3	Item4	BestItem	WorstItem
0	Rodel	Knödelakademie	Keiler	Gelenksbeschwerden	Rodel	Gelenksbeschwerden
1	brennheiß	Stornoversicherung	Scherz(e)l	sich ausgehen	sich ausgehen	brennheiß
2	Steireranzug	Causa	Pönale	Lokalausweis	Lokalausweis	Steireranzug
3	Alumnat	Beiwagerl	Servus	kiefeln	Servus	kiefeln
4	Patschenkino	Aufnahmestopp	Straßenerhalter	Marmeladinger	Straßenerhalter	Aufnahmestopp
...	...	...	...	...	...	...
4412	ferten	Ermäßigungsausweis	Halbpreisspass	versumpfern	Ermäßigungsausweis	versumpfern
4413	Zuhause	Bramburi	Mistbauer	Beiwagerl	Zuhause	Mistbauer
4414	Oja!	ludeln	Rettung	gar	Oja!	ludeln
4415	Stützlehrer	Mascherl	Einspänner	grauslich	Mascherl	grauslich
4416	Jausenbrot	enthaften	versperren	Schubhaft	Jausenbrot	Schubhaft



# Methods: Data Annotation (Atriacisms)

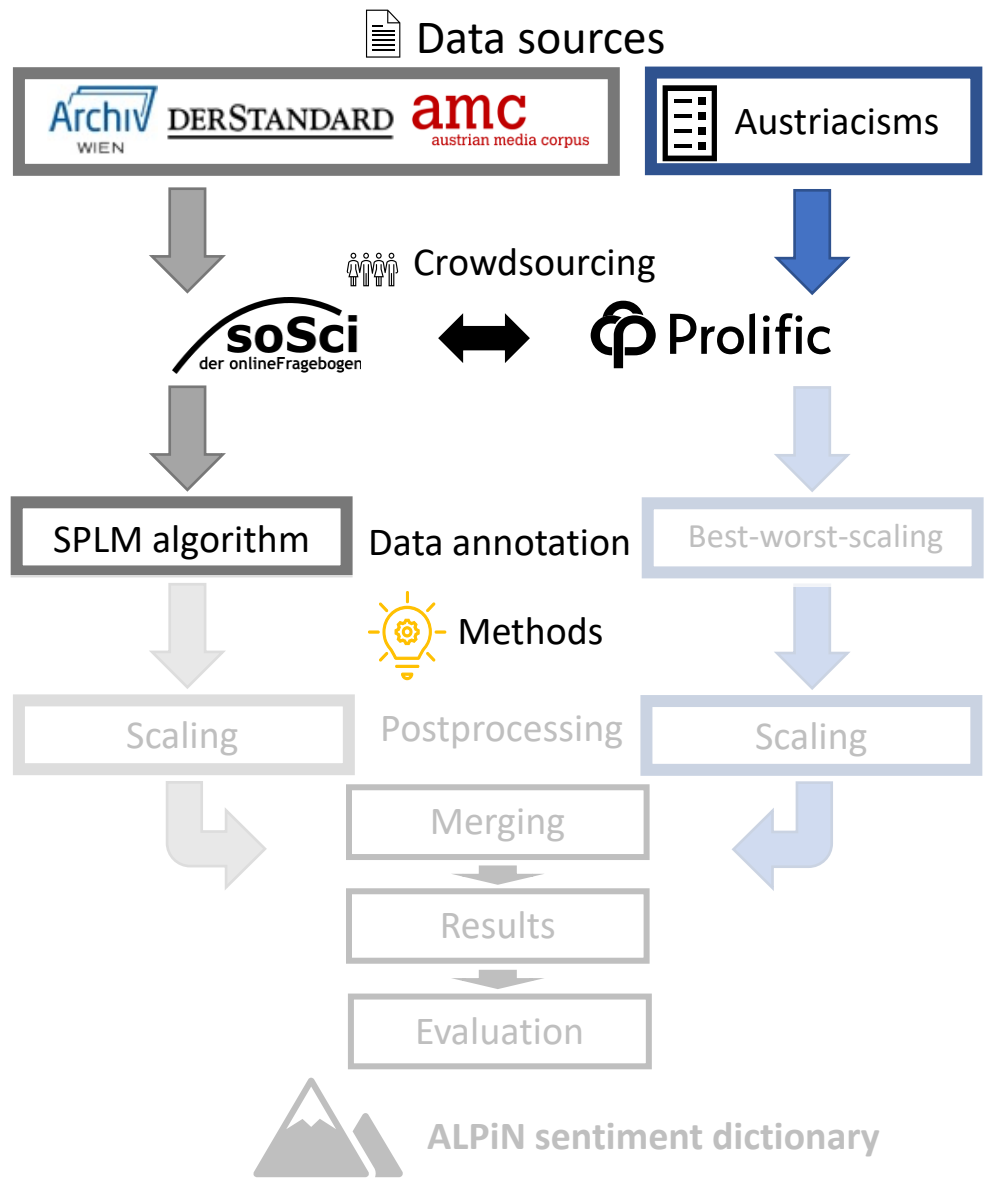
**Best-worst-scaling (BWS) method** (Kiritchenko & Mohammad, 2017)

## split-half reliability:

- Spearman correlation: 0.9159 +/- 0.0051
- Pearson correlation: 0.9164 +/- 0.0049

**Output: 537 words**

	word	tag	short-tag	score	scaled
0	fesch	ADJ	a	0.882	0.910217
1	Zuckerl	NOUN	n	0.879	0.907121
2	Topfenpalatschinke	NOUN	n	0.857	0.884417
3	leiwand	ADJ	a	0.853	0.880289
4	Ersparnis	NOUN	n	0.844	0.871001
...	...	...	...	...	...
533	Schussattentat	NOUN	n	-0.844	-0.871001
534	Exekution	NOUN	n	-0.848	-0.875129
535	speiben	VERB	v	-0.875	-0.902993
536	Brandleger	NOUN	n	-0.879	-0.907121
537	Fotze	NOUN	n	-0.969	-1.000000



# Methods: Data Annotation (amc, derStandard)

**SPLM method** (Almatarneh & Gamallo, 2018)

Algorithm to generate a sentiment score based on labelled text items.

**Remark:** “neutral” sentiment labels of the derStandard dataset were converted to “positive”. This was required to the high imbalance in the dataset.

SentimentNeutral	1865
SentimentNegative	1691
SentimentPositive	43

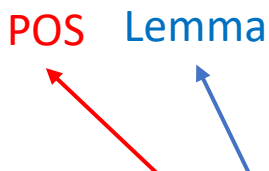
# Methods: Data Annotation (amc, derStandard) preprocessing (1)

## amc

„**Lemma**“ and „**POS**“ based on  
the amc corpora<sup>1</sup> (RFTagger/Tiger corpus )

## derStandard

„**Lemma**“ and „**POS**“ based on  
spacy („de\_core\_news\_sm“)



	text	polarity
0	[(X, Top), (ADJ, qualifizieren), (NOUN, Leute)...	neutral
1	[(NOUN, Gott), (AUX, sein), (ADP, danken), (AU...	negative
2	[(PUNCT, "), (DET, der), (PROPN, FPÖ), (AUX, w...	neutral
3	[(SCONJ, Weil), (PRON, ich), (DET, mein), (DET...	negative
4	[(INTJ, Na), (PRON, wer), (NOUN, weis), (PRON,...	neutral

example based on the derStandard dataset



# Methods: Data Annotation (amc, derStandard) preprocessing (2)

Map result of „POS“ tagging to „wordnet“ tags to reduce the number of tags:

```

amc
tag = {
  'ADJA':wn.ADJ, # attributive adjectives
  'ADJD':wn.ADV, # adjective with predicative or adverbial usage
  'ADV':wn.ADV, # adverbs
  'N':wn.NOUN, # Noun
  'VFIN':wn.VERB, # finite verb
  'VIMP':wn.VERB, # imperative verbs
  'VINF':wn.VERB, # infinitival verb
  'VPP':wn.VERB # participle verb
}

derStandard
# https://universaldependencies.org/u/pos/
tag = {
  'ADJ':wn.ADJ, # adjective
  'ADV':wn.ADV, # adverbs
  'NOUN':wn.NOUN, # noun
  'PRON':wn.NOUN, # proper noun
  'PROPN':wn.NOUN, # proper noun
  'VERB':wn.VERB # verb
}
    
```

	polarity	wordnet tag	text
0	neutral		[(so, r), (eine, ), (Ansinnen, n), (scheinen, ...
1	negative		[(Roland, n), (Sperk, n), (Vorsitzende, n), (d...
2	neutral		[(feiern, v), (werden, v), (in, ), (Szenelokal...
3	neutral		[(in, ), (die, ), (ÖVP, n), (machen, v), (sich...
4	neutral		[(bei, ), (eine, ), (Erfolg, n), (die, ), (Vol...
...	...		...
8909	positive		[(Russland, n), (ist, ), (in, ), (wk1, ), (vor...
8910	positive		[(Was, ), (tendenziell, r), (kein, ), (schlech...
8911	positive		[(Was, ), (Unsinn, n), (Der, ), (Linguistik, n...
8912	negative		[(wien, n), (verschreckt, v), (investoren, v),...
8913	negative		[(Früher, r), (haben, ), (sie, ), (ein, ), (vi...

8914 rows × 2 columns

# Methods: Data Annotation (amc, derStandard) result (1)

	word	Tag	D
8293	auch	r	1.807385e-03
9541	ich	n	1.655992e-03
3533	sehr	r	1.094594e-03
6266	geben	v	8.939163e-04
1139	Frau	n	8.637304e-04
...	...	...	...
3729	Quote	n	1.707558e-06
3758	notwendig	r	1.542211e-06
1252	klar	r	1.542211e-06
2394	überhaupt	r	1.211518e-06
7530	brauchen	v	8.808251e-07

4675 rows × 3 columns

positive words

	word	Tag	D
4949	haben	v	-2.215640e-03
2837	sein	v	-1.959662e-03
1552	hier	r	-9.608509e-04
5185	Flüchtling	n	-9.493369e-04
4123	nur	r	-8.721120e-04
...	...	...	...
7736	Anton	n	-1.653465e-07
10822	David	n	-1.653465e-07
1595	Spital	n	-1.653465e-07
3735	Einkommen	n	-1.653465e-07
9023	Typus	n	-1.653465e-07

4392 rows × 3 columns

negative words

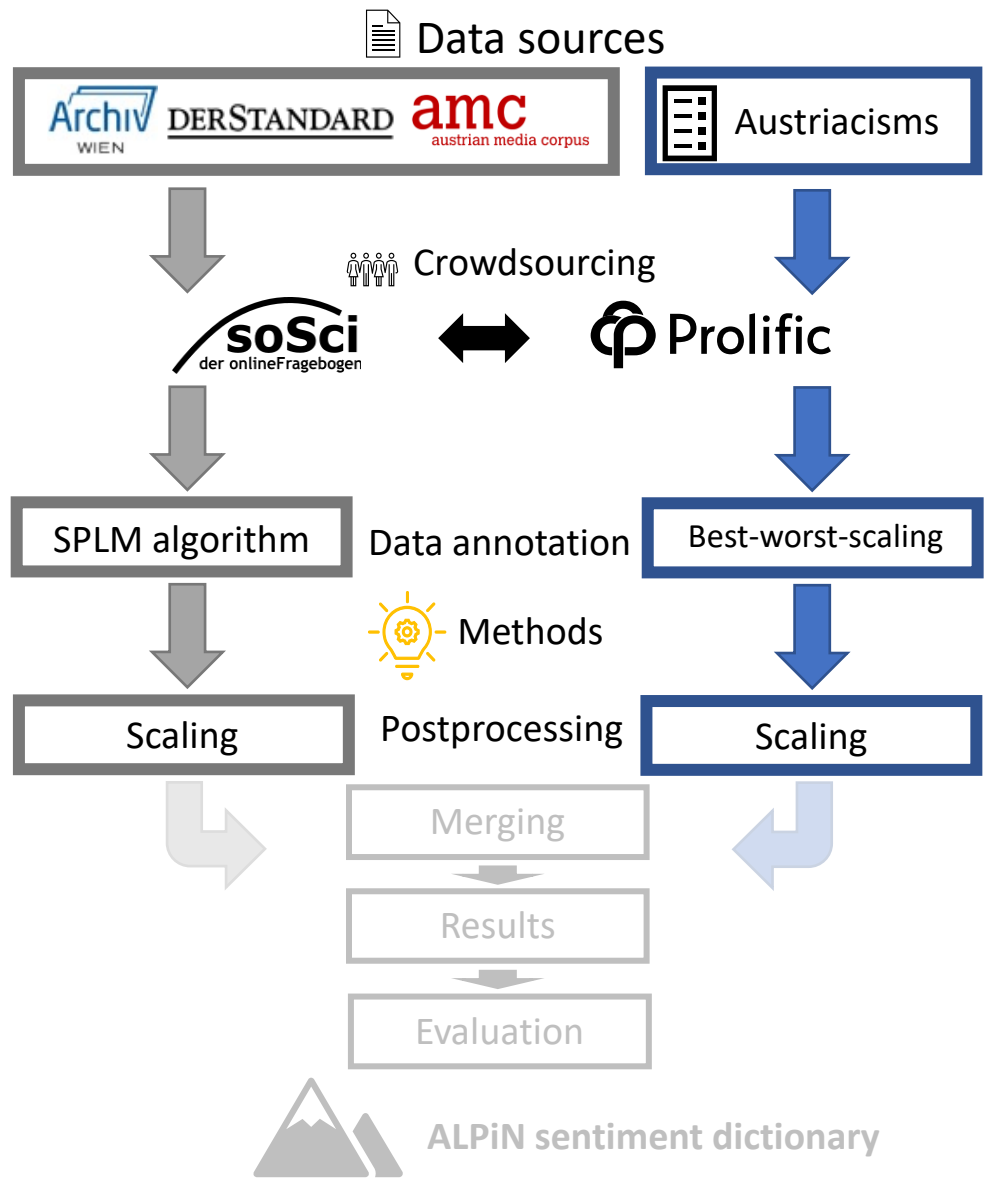
$D(w)$ : sentiment score  
 $D(w) [-1;+1]$

# Methods: Data Annotation (amc, derStandard) result (2)

	word	Tag	D
0	auch	r	0.001807
1	ich	n	0.001656
2	sehr	r	0.001095
3	geben	v	0.000894
4	Frau	n	0.000864
...	...	...	...
9062	nur	r	-0.000872
9063	Flüchtling	n	-0.000949
9064	hier	r	-0.000961
9065	sein	v	-0.001960
9066	haben	v	-0.002216

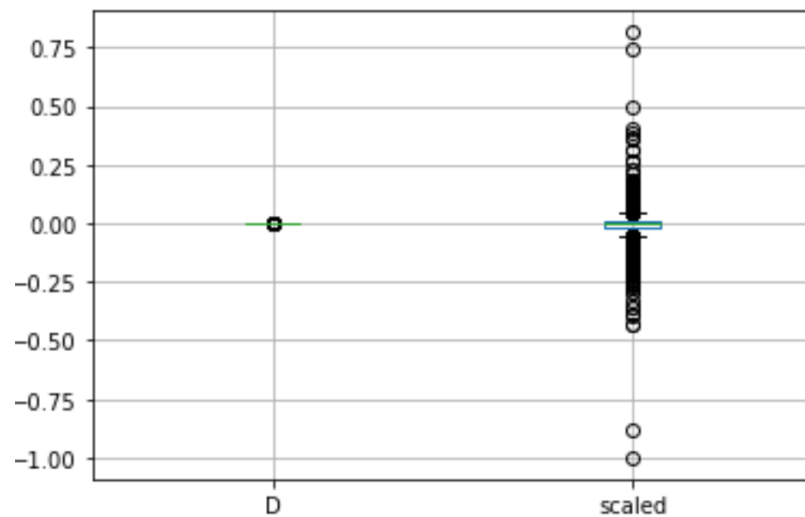
$D(w)$  = sentiment score

9067 rows × 3 columns

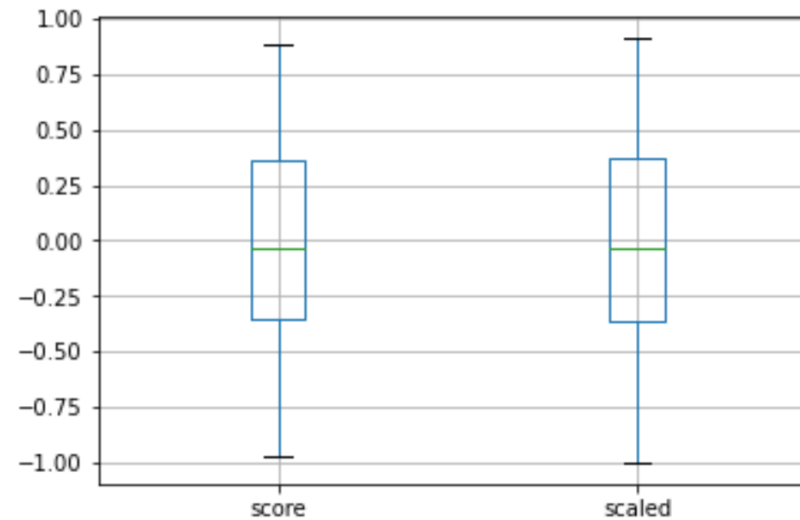


# Methods: Postprocessing (1)

Scaling to  $[-1,+1]$  with „max\_abs\_scaler of sklearn“<sup>1</sup> before merging the dictionaries

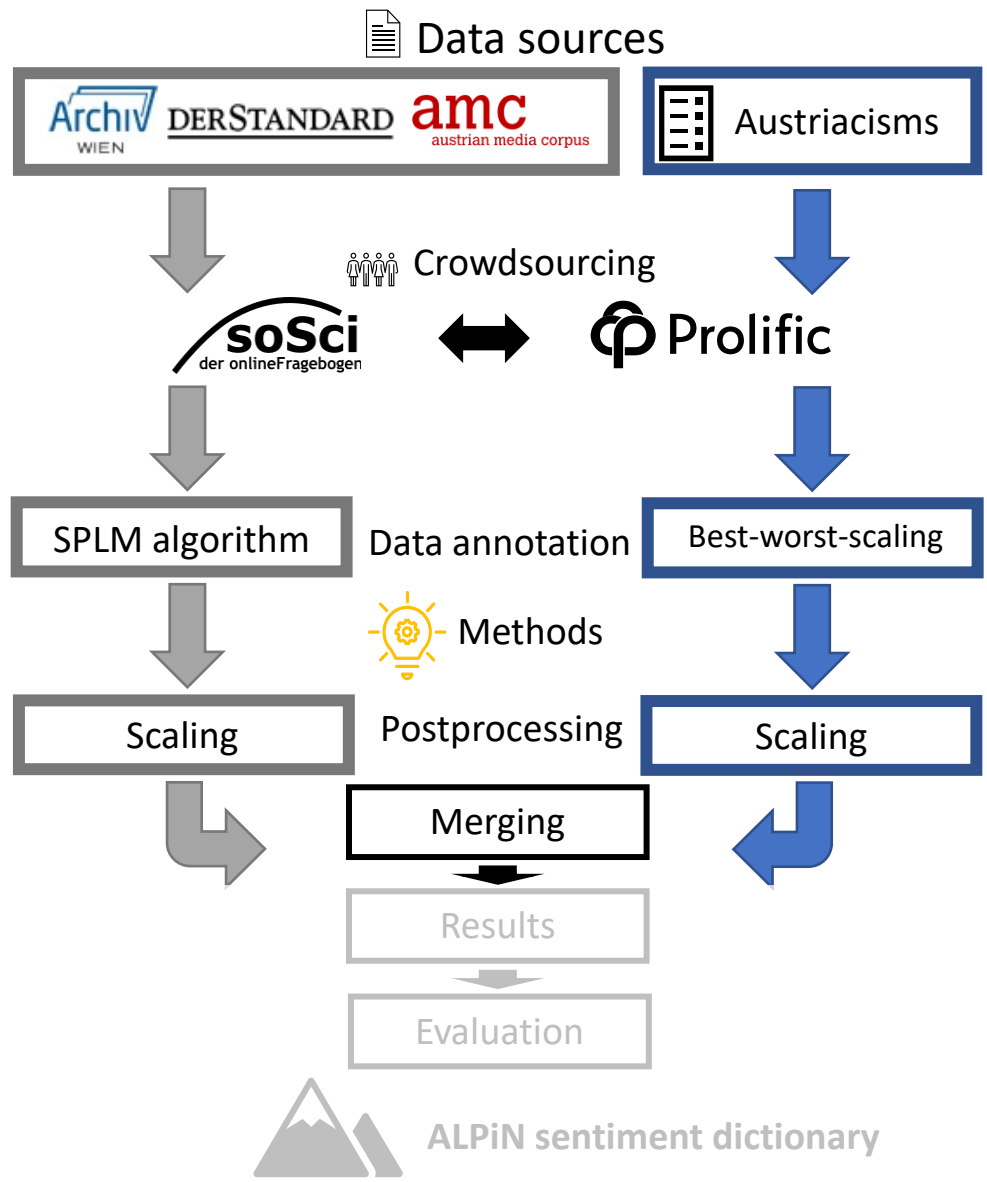


amc with derStandard after applying SPLM



Austriacisms after applying BWS

1. <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MaxAbsScaler.html>



# Methods: Postprocessing (2)

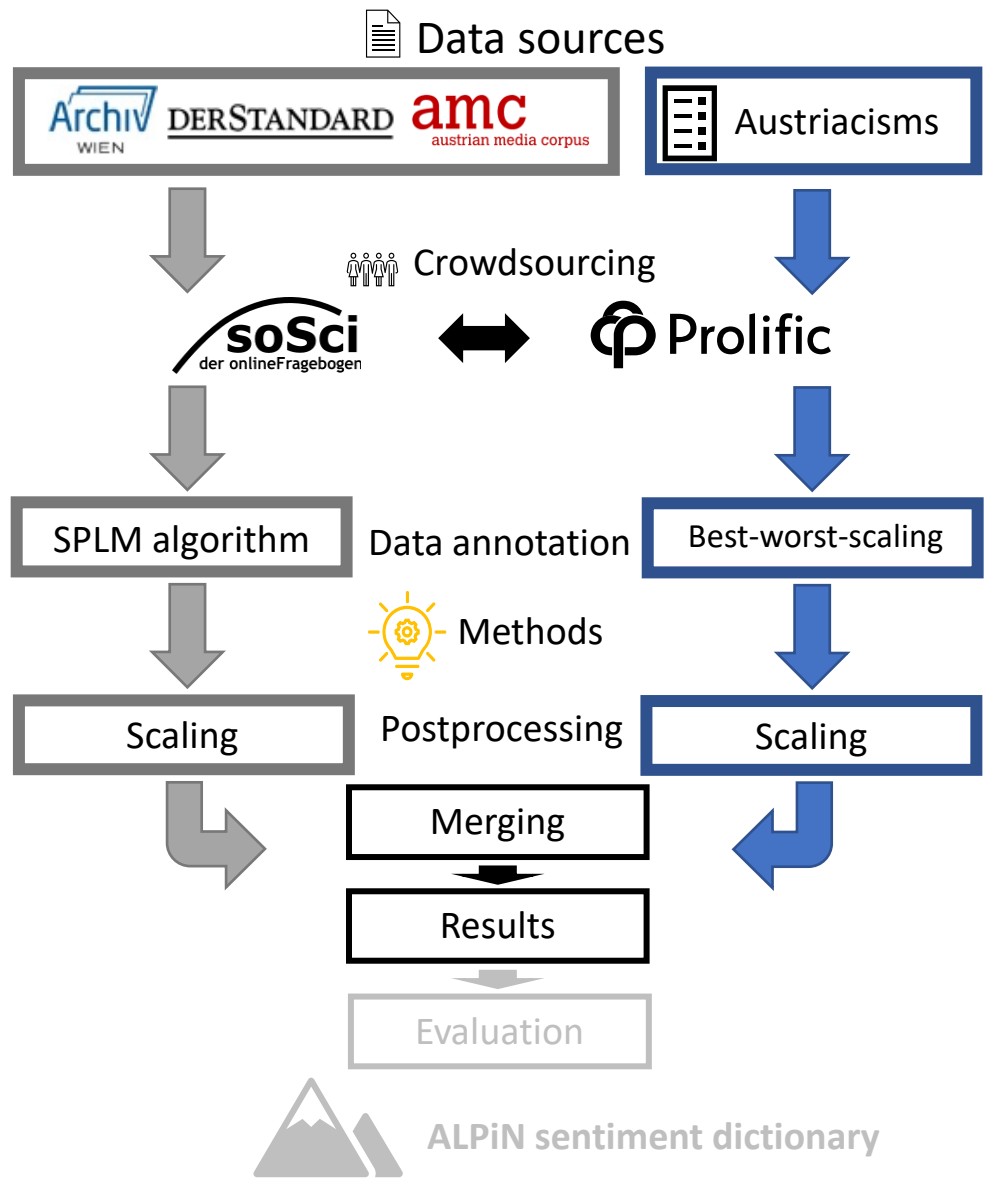
Comparison of words which occur in both dictionaries (amc+derStandard vs austriacisms):

	word	short-tag	sentiment_austriacism	sentiment_dysen_with_derstandard
0	Wiese	n	0.750258	0.011685
1	Karenz	n	0.742002	0.017528
2	Angelobung	n	0.728586	-0.021754
3	Ehrenzeichen	n	0.710010	0.029213
4	Gehalt	n	0.644995	0.091790
5	aufrecht	a	0.625387	-0.013376
6	maturieren	v	0.625387	0.011685
7	ÖAMTC	n	0.562436	-0.013376
8	einbringen	v	0.547988	-0.052732
9	Team	n	0.515996	0.072572

20	Abgang	n	0.000000	-0.033439
21	Klappe	n	-0.226006	-0.013376
22	klagen	v	-0.312693	-0.023445
23	angreifen	v	-0.343653	-0.001765
24	Fleck	n	-0.375645	-0.013376
25	Einvernahme	n	-0.386997	-0.013376
26	Freunderwirtschaft	n	-0.437564	-0.013376
27	versperren	v	-0.486068	-0.013376
28	Mist	n	-0.594427	-0.013376
29	sekkieren	v	-0.688338	-0.013376
30	exekutieren	v	-0.837977	-0.001691
31	Exekution	n	-0.875129	0.011685

## Restrictions:

During merging duplicates will be removed by using the Austriacism words prioritized.





# Results

## amc data only

	word	Tag	D
<b>0</b>	neu	a	0.002108
<b>1</b>	Wien	n	0.002040
<b>2</b>	Wiener	a	0.001465
<b>3</b>	Jahr	n	0.001432
<b>4</b>	Michael	n	0.001307
...	...	...	...
<b>4863</b>	Pilz	n	-0.001522
<b>4864</b>	Westenthaler	n	-0.001664
<b>4865</b>	Peter	n	-0.001664
<b>4866</b>	sein	v	-0.002586
<b>4867</b>	haben	v	-0.003756

4868 rows × 3 columns

## amc with derStandard

	word	Tag	D
<b>0</b>	auch	r	0.001807
<b>1</b>	ich	n	0.001656
<b>2</b>	sehr	r	0.001095
<b>3</b>	geben	v	0.000894
<b>4</b>	Frau	n	0.000864
...	...	...	...
<b>9062</b>	nur	r	-0.000872
<b>9063</b>	Flüchtling	n	-0.000949
<b>9064</b>	hier	r	-0.000961
<b>9065</b>	sein	v	-0.001960
<b>9066</b>	haben	v	-0.002216

9067 rows × 3 columns

# Results

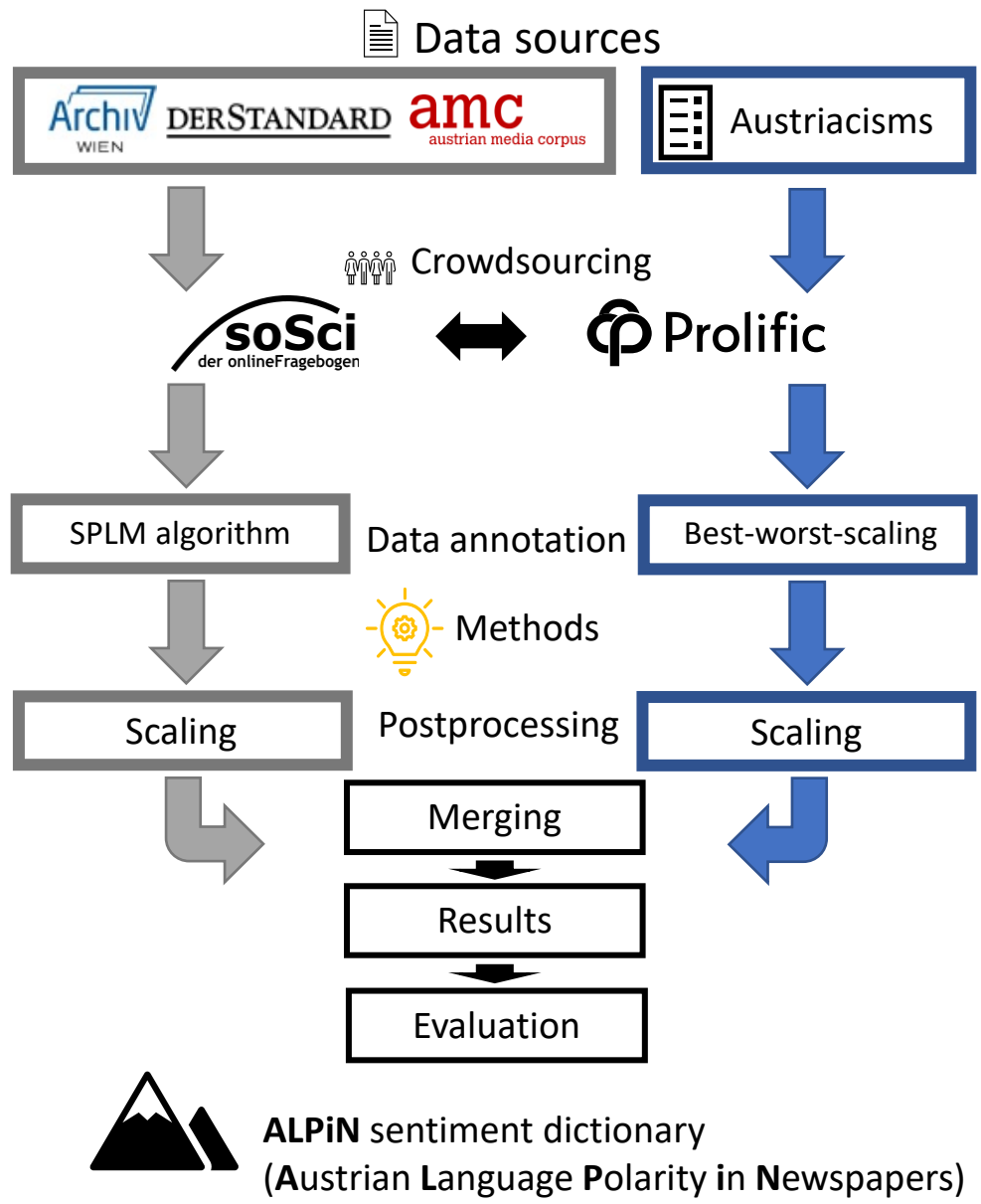
**amc + derStandard + austriacisms**

Scaled to [-1,+1] with  
„max\_abs\_scaler of sklearn“<sup>1</sup>

	word	short-tag	scaled
<b>0</b>	fesch	a	0.910217
<b>1</b>	Zuckerl	n	0.907121
<b>2</b>	Topfenpalatschinke	n	0.884417
<b>3</b>	leiwand	a	0.880289
<b>4</b>	Ersparnis	n	0.871001
...	...	...	...
<b>9568</b>	sein	v	-0.884468
<b>9569</b>	speiben	v	-0.902993
<b>9570</b>	Brandleger	n	-0.907121
<b>9571</b>	Fotze	n	-1.000000
<b>9572</b>	haben	v	-1.000000

9573 rows × 3 columns

1. <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MaxAbsScaler.html>



# Evaluation (1)

**Method:** Kfold (5 folds), cross\_validation, SVC(kernel='linear')

## Features:

- Count of positive words in text-item
- Count of negative words in text-item
- proportion

	text	polarity	count_pos	count_neg	proportion
0	[(qualifizieren, a), (Leute, n), (verdienen, v...	positive	5	0	5.000000
1	[(Gott, n), (ich, n), (Umfrage, n), (alle, n),...	negative	7	2	3.500000
2	[(FPÖ, n), (Rohr, n), (schießen, v), (Regierun...	positive	5	6	0.833333
3	[(ich, n), (Leute, n), (verständlicherweise, r...	negative	6	4	1.500000
4	[(wer, n), (weis, n), (was, n), (da, r), (vorf...	positive	3	2	1.500000

dataset after feature calculation

# Evaluation (2)

Evaluate the dictionary which is based on amc, derStandard and the austriacism list against "derStandard" and "DYSEN":

**1<sup>st</sup>** against derStandard only

```
fit_time 0.70582594871521
score_time 0.03183770179748535
test_accuracy 0.7532595425745635
test_precision 0.7655951442582838
test_recall 0.7740803341990627
test_f1_score 0.7688922021025179
```

**2<sup>nd</sup>** against amc only

```
fit_time 0.17100081443786622
score_time 0.01651768684387207
test_accuracy 0.8150271692254615
test_precision 0.8322847276249622
test_recall 0.8117444005270092
test_f1_score 0.8186862563698238
```

# Discussion (1)

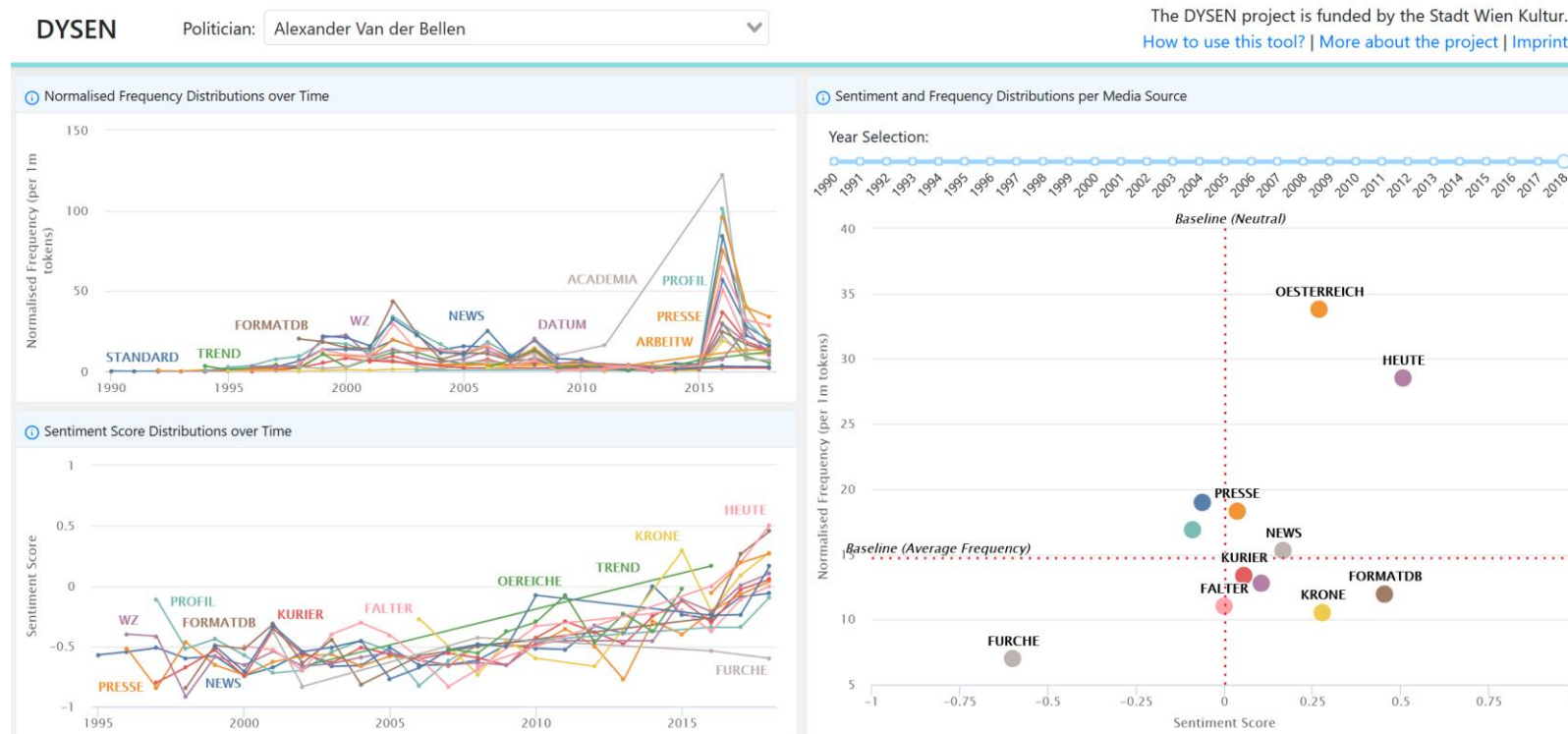
- Difficult to label news media (mainly “neutral” texts), as a result the inter-annotator agreement is not as high as in other domains by using similar methods
- Limited text length
- No external dataset for evaluation

## **Future work:**

- Improvement of the text extraction by using Aspect-based sentiment analysis
- Investing more money to label a bigger dataset
- Expanding the scope of the project to all politicians and media in Austria

# Discussion (2)

Tool created as part of the DYSEN project which uses the ALPiN dict.:



<https://dysen-tool.acdh-dev.oeaw.ac.at/> (work in progress)



# ALPiN Dictionary (1)

***„Austrian Language Polarity in Politics and Newspapers“***

 Current research topic in our DYSEN project

 Currently there is no dictionary based on Austrian-German in the domain of news media and politics

**amc**  
austrian media corpus

Based on the „Austrian Media Corpus“ phrases related to Viennese politicians of the last 20 years

 Labelled dataset created via crowd-sourcing (prolific) by Austrian German native speakers

 Dictionary generated by applying the SPLM (Almatarneh & Gamallo, 2018) algorithm





# ALPiN Dictionary (2)

*„Austrian Language Polarity in Politics and Newspapers“*



Incorporation of Austriacisms by using the best-worst-scaling (BWS) to improve the quality of the labels (Kiritchenko & Mohammad, 2017)

DERSTANDARD

Incorporation of derStandard (popular Austrian news media) forum posts



Diverse independent data-sources (amc, derStandard, austriacisms)



Resulting resource and paper will be submitted/published by the end of the year

# Thank you very much!

*thomas.kolb@tuwien.ac.at*

Funded by:



**Stadt  
Wien**

*Grant number:  
MA7-737909/19*

**ÖAW**

ÖSTERREICHISCHE  
AKADEMIE DER  
WISSENSCHAFTEN



**universität  
wien**



**TECHNISCHE  
UNIVERSITÄT  
WIEN**

# References

Dietmar Schabus, Marcin Skowron, and Martin Trapp. 2017. One Million Posts: A Data Set of German Online Discussions. In Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '17). Association for Computing Machinery, New York, NY, USA, 1241–1244. DOI:<https://doi.org/10.1145/3077136.3080711>

Ransmayr, Jutta, Karlheinz Mörth, und Matej Ďurčo (2017): AMC (Austrian Media Corpus) – Korpusbasierte Forschungen zum österreichischen Deutsch. In Digitale Methoden der Korpusforschung in Österreich (= Veröffentlichungen zur Linguistik und Kommunikationsforschung Nr. 30), Hrsg. C. Resch und W. U. Dressler, 27-38. Wien: Verlag der Österreichischen Akademie der Wissenschaften.

Bickel, H., Hofer, L., & Suter, S. (2015). 22. Variantenwörterbuch des Deutschen (VWB)–NEU. In Regionale Variation des Deutschen (pp. 541-562). De Gruyter.

Kiritchenko, S., & Mohammad, S. M. (2017). Capturing Reliable Fine-Grained Sentiment Associations by Crowdsourcing and Best-Worst Scaling.

Kiritchenko, S., & Mohammad, S. (2017). Best-Worst Scaling More Reliable than Rating Scales: A Case Study on Sentiment Intensity Annotation. Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), 465–470. <https://doi.org/10.18653/v1/P17-2074>

Almatarneh, S., & Gamallo, P. (2018). Automatic Construction of Domain-Specific Sentiment Lexicons for Polarity Classification. 175–182. [https://doi.org/10.1007/978-3-319-61578-3\\_17](https://doi.org/10.1007/978-3-319-61578-3_17)

Rouces, J., Tahmasebi, N., Borin, L., & Eide, S. R. (2018). Generating a Gold Standard for a Swedish Sentiment Lexicon. LREC.