

Diachronic Analysis of German Parliamentary Proceedings: Ideological Shifts through the Lens of Political Biases

Tobias Walter,¹ Celina Kirschner,¹ Steffen Eger,² Goran Glavaš,¹ Anne Lauscher,¹ Simone Paolo Ponzetto¹

¹*Data and Web Science Group, University of Mannheim, Mannheim, Germany*

²*Natural Language Learning Group, Technische Universität Darmstadt, Darmstadt, Germany*

celina.kirschner@hotmail.de, tobias.walter-ul@web.de,

{goran, anne, simone}@informatik.uni-mannheim.de, eger@aiphes.tu-darmstadt.de

Abstract—We analyze bias in historical corpora as encoded in diachronic distributional semantic models by focusing on two specific forms of bias, namely a political (i.e., anti-communism) and racist (i.e., antisemitism) one. For this, we use a new corpus of German parliamentary proceedings, DEUPARL, spanning the period 1867–2020. We complement this analysis of historical biases in diachronic word embeddings with a novel measure of bias on the basis of term co-occurrences and graph-based label propagation. The results of our bias measurements align with commonly perceived historical trends of antisemitic and anti-communist biases in German politics in different time periods, thus indicating the viability of analyzing historical bias trends using semantic spaces induced from historical corpora.

I. INTRODUCTION

Recent years have seen much work on the topic of bias in data-driven Artificial Intelligence [28] and Machine Learning [24]. In the case of Natural Language Processing, researchers have investigated the bias encoded within semantic spaces induced by both non-contextualized [5], [13] and contextualized embeddings [20], and a variety of methods has been developed to debias such embedding spaces in such a way as to make them fairer [4], [15], [22, *inter alia*].

However, despite much interest on studying and mitigating bias in semantic spaces, few works – [13], [33] being notable exceptions – have taken a historical perspective and looked at ways to study bias in diachronic corpora through the lens of distributional embedding models. In this paper, we create a new historical corpus of parliamentary proceedings spanning three different centuries and set to quantify different kinds of bias in the underlying historical periods. Our DEUPARL corpus covers the protocols of both the German Reichstag and the Bundestag from 1867 until 2020, thus spanning a large timeline that covers many crucial modern and contemporary events (e.g., two world wars): it consists of the digitized scans of the older German Reichstag and the digitally-born newer German Bundestag parliamentary proceedings. In this corpus, we focus on political and racist forms of bias, i.e., anti-communism and antisemitism, respectively. We achieve this by building upon previous work on quantifying the biases in textual corpora, that is: (1) we devise a number of term-based bias specifications that provide us with an operational

definition of antisemitic and anti-communist bias; (2) we quantify the degree of such biases using measures applied to word embeddings induced from different time slices of our corpus; (3) we complement this study of bias in dense semantic vectors with an analysis of bias in more ‘classic’ (i.e., sparse) distributional word representations using a label propagation algorithm applied to word co-occurrence graphs. As a result, we are able to provide a diachronic analysis of the bias trends in our historical corpus for different kinds of biases and word representations on the basis of our specifications.

The contributions of this work are the following ones:

- We present a new historical corpus of parliamentary proceedings built from the protocols of both German Reichstag and Bundestag between 1867 and 2020.
- We induce a variety of distributional semantic spaces, both dense and sparse, from different slices of our historical corpus and observe bias trends across time for different kinds of bias (i.e., antisemitism and anti-communism).
- We introduce a new measure of bias on the basis of term co-occurrences and graph-based label propagation, which makes it possible to quantify bias in sparse distributional spaces and thus in an arguably more interpretable way.

Our results are consistent with common historical interpretation of development of bias in German history, i.e., we identify an increase in antisemitism towards the NS period, with a peak in the Weimar Republic. This indicates the viability of studying the evolution of bias from a historical perspective on the basis of word embeddings, on associated tests, from historical corpora. We make all code and data (including DEUPARL) available at https://github.com/umanlp/crosstemporal_bias.

II. DEUPARL: GERMAN PARLIAMENTARY CORPUS

Our new historical corpus consisting of German parliamentary proceedings from 1867 until 2020, which we refer to as DEUPARL,¹ is built from two main sources.

¹We make DEUPARL available at: <https://tudatalib.ulb.tu-darmstadt.de/handle/tudatalib/2889>

TABLE I: Excerpt of a speech in a parliamentary session, from 1914. Left: OCR-scanned document. Right: Corrected version. Words containing OCR errors are underlined.

Gröber, Abgeordneter: Meine Herren, der Gedanke des Herrn Abgeordneten <u>Gothein</u> , der dem Beschluß des hohen Hauses in der letzten Sitzung zu Grunde lag, ist meines Erachtens zweifellos richtig. Es ist nicht er-	Gröber, Abgeordneter: Meine Herren, der Gedanke des Herrn Abgeordneten Gothein, der dem Beschluß des hohen Hauses in der letzten Sitzung zu Grunde lag, ist meines Erachtens zweifellos richtig. Es ist nicht er-
---	---

Reichstagsprotokolle. For parliamentary speeches before 1945, we use the German *Reichstagsprotokolle* which are accessible at <https://www.reichstagsprotokolle.de/>. The *Bay-erische Staatsbibliothek* distributes the digitized data upon request, arranged in three splits covering the years 1867–1895 (Norddeutscher Bund / Zollparlamente), 1895–1918 (Kaiserreich), and 1918–1942 (Weimarer Republik / Nationalsozialismus), respectively. These were OCR-scanned using the Abbyy FineReader software.² To check data quality, we extracted a sample of several thousand lines which we corrected by hand. In this sample, about 23% of all lines contained one or more OCR errors. A sample of the data is shown in Table I: we found the data to be of sufficiently high quality for our purposes, not urgently necessitating further OCR post-correction steps [11]. Overall, we extracted 5,446 individual protocols in the Reichstagsprotokolle, i.e., stenographic documents corresponding to the parliamentary sessions.

Bundestagsprotokolle. For parliamentary speeches after 1945, we use the *Bundestagsprotokolle* from <https://www.bundestag.de/protokolle/>.³ We extracted 4,260 protocols in 19 legislative periods from 1949 to March 2020.

Preprocessing. We conducted preprocessing steps for DEU-PARL (mostly affecting the Reichstagsprotokolle) including:⁴

- We fixed word segmentations at line breaks (e.g., “Poli- tik” → “Politik”);
- We lowercased and lemmatized all words (the latter using GermaLemma++ [29]) to increase statistical support for all subsequent models built on top of our data;
- We resolved cases such as “Luft- und Raumfahrt” into “Luftfahrt und Raumfahrt” and fixed erroneous word segmentations (e.g., where numbers and words were merged)
- We took care to keep only text from the individual parliamentary sessions, not considering various attachments to the protocols such as lists of names. To this end, we used regular expressions to extract the text within boundaries such as “Die Sitzung ist eröffnet” (“*the session is opened*”) and “Schluss der Sitzung hh Uhr mm Minute(n)” (“*end of session hh:mm*”)

Temporal splits. We defined our own temporal splits for a historical analysis of bias, taking historically accepted time

²<https://pdf.abbyy.com/>

³Note that, for the time period of the Division of Germany between 1949 and 1990, the Bundestagsprotokolle only cover the protocols of the Federal Republic of Germany (also known as “West Germany”).

⁴We also aimed at historical spelling normalization, but found it to lead to substantial amount of undesired text modifications.

TABLE II: Labels for time periods, range of time periods, and sizes of the corresponding data slices.

Period	Years	Tokens
KR1	1867-1890	40,585,912
KR2	1890-1918	77,175,976
WR	1918-1933	35,838,922
NS	1933-1942	230,018
CDU1	1949-1969	43,337,027
SPD1	1969-1982	35,208,879
CDU2	1982-1998	55,451,433
SPD2	1998-2005	28,614,189
CDU3	2005-2020	66,192,033

periods in modern German history as reference. These are: the Kaiserreich I (KR1, 1867–1890), Kaiserreich II (KR2, 1890–1918), Weimarer Republik (WR, 1918–1933), Nationalsozialismus (NS, 1933–1945). In the case of the Bundesrepublik (1949–2020), we divided the data into slices of contiguous time periods in which one of the two major German parties (Volksparteien, i.e., liberal-conservative CDU or social-democratic SPD) was in charge. Table II shows statistics of our data. For the Bundestagsprotokolle, slices range from between ~30 million to ~66 million tokens. For the Reichstagsprotokolle, the NS era is severely restricted in size, congruent with the observation that the parliamentary processes during the NS time were largely abolished.⁵

Party distribution. The number of parties in our historical time periods with active speakers ranges from 2 (during the NS period, including SPD and NSDAP) to 10 (during KR2).⁶ Of course, the party distribution may decisively affect our results, i.e., distribution of biases found. In this context, it is worth noting that certain parties were banned during the time span of our data. For example, the right-wing party NSDAP was forbidden after 1945. The communist party KPD was forbidden in 1956. Our data thus contains communist parties only during the time period of WR and CDU1.

Corpus analysis. In Figure 1, we show so-called ‘word clouds’ for each slice in our timeline (1867–2020).⁷ The word clouds graphically illustrate the frequency of words in text data, assigning more frequent words larger font sizes (we removed stop words and only consider words with a minimum frequency of 50 occurrences). As expected, vocabulary relevant for parliamentary debates is prominent in our textual material, e.g., ‘Abgeordneter’ (*representative*), ‘Herr’ (*mr.*),

⁵https://www.bundestag.de/parlament/geschichte/parlamentarismus/drittes_reich

⁶In the original data, party affiliation is attached to the speaker names.

⁷We use https://github.com/amueller/word_cloud for drawing.



Fig. 1: Word clouds showing relative frequency of words in text data for our nine periods: KR1, KR2, WR, NS (top, from left to right) and CDU1, SPD1, CDU2, SPD2, CDU3 (bottom).

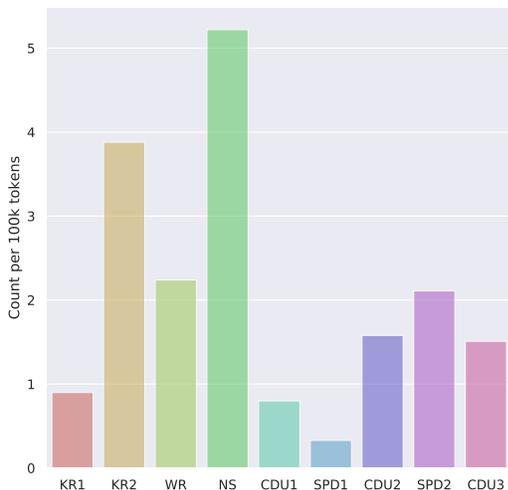


Fig. 2: Histogram summarizing the number of occurrence of the German word ‘Jude’ (*jew*) per 100k tokens across different time slices of our corpus.

‘Antrag’ (*petition*). Names for political parties (‘SPD’, ‘FDP’, ‘CDU/CSU’) become more frequent in more recent time periods. Interestingly, the word ‘deutsch’ (*German*) becomes more prominent from KR1 to KR2 to WR and then peaks in the NS period, where the word ‘Volk’ (ideologically-loaded *people* as nation) along with ‘Führer’ (*leader*, used as title for Nazi dictator Adolf Hitler) and ‘Berlin’ also become very prominent. The trend for ‘deutsch’ then reverses with the onset of the Bundesrepublik, i.e., after 1949, thus intuitively aligning with general understatement of national and cultural identity in Germany after World War II.

A limitation of word clouds is that they can highlight only a few frequent words from the corpus. To also analyze trends for lower frequency and highly relevant words, e.g., ‘Jude’ (*jew*) – due to its central relevance for one of the two kinds of bias that we focus on (i.e., antisemitism) – we show a temporal curve in Figure 2. Like ‘deutsch’, ‘Jude’ peaks in the NS period, with a frequency of almost 6 per 100k tokens, then sharply declines in the post war period to regain prominence again after 1982, though not on the same level as before 1949 (disregarding KR1). Finally, Figure 3 shows semantic shift [10] of the word ‘Judentum’ (*Judaism*) in our data. The plot

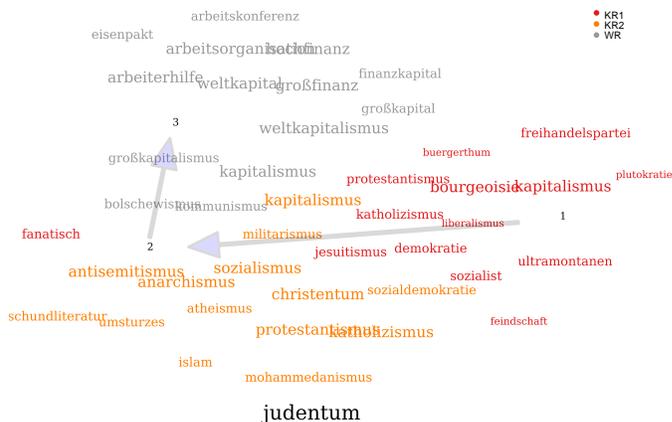


Fig. 3: The semantic representations of the word ‘Judentum’ (denoted by the numbers in chronological order, i.e., 1:KR1, 2:KR2, 3: WR) and their 15 nearest neighbors in the semantic space of each Reichstag period (red: KR1, orange:KR2, gray:WR).

is obtained by semantically embedding words from different time periods in a shared vector space using the model of [6] and then tracking in the semantic space how the word changes its fine-granular meaning over time periods using t-SNE [34] to visualize the embeddings (cf. also HistWords [17]). In this case, for instance, we see how ‘Judentum’ seemingly takes on (putatively negative) associations of ‘Bolschewismus’ (*Bolshevism*) and ‘Weltkapitalismus’ (*world capitalism*) in the time of the Weimar Republic.

III. CAPTURING POLITICAL BIASES: METHODOLOGY

Our bias evaluation pipeline is illustrated in Figure 4. The methodology consists of three main steps. We first slice DEU-PARL according to the dimension of interest: in this paper, our focus is on politically well-defined historical periods. We then induce computational text representations, namely (1) dense distributional word vector spaces (i.e., word embeddings) and (2) word co-occurrence graphs, for each of the corpus slices. Finally, given human-designed bias specifications for *antisemitism* and *anti-communism*, we compute bias scores with several measures on the basis of either word embedding spaces or word cocurrence graphs induced in the previous

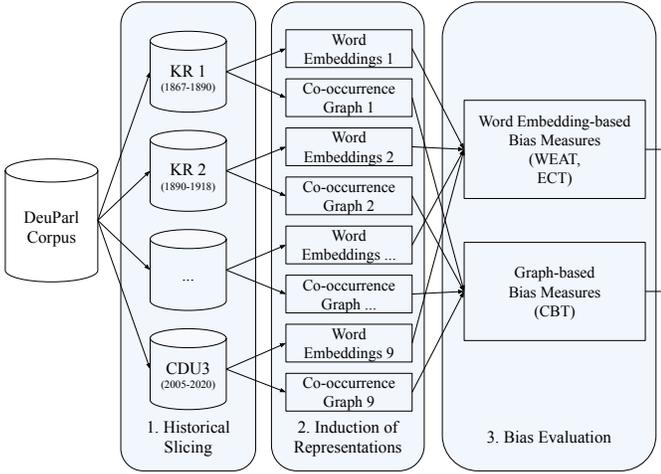


Fig. 4: Overview of our bias evaluation pipeline.

step. We next describe our bias specifications, followed by the detailed description of our word embedding-based and co-occurrence graph-based approaches for quantifying the bias effects across different time periods.

A. Bias Specifications

In order to quantify the degree of presence of some societal bias (e.g., racial, religious, political, or gender bias), one must first specify (i.e., formalize) that bias. In this work, we follow the established body of work on measuring biases in word embedding spaces [5], [9], [22] and adopt the so-called *explicit bias specification* $B_E = (T_1, T_2, A_1, A_2)$, consisting of two target term sets (T_1 and T_2), between which a significantly different association with respect to two sets of attribute terms (A_1 and A_2) is expected to exist. In this work, we investigate two dimensions of political biases, which arguably played the most important role in (most recent) German history: (1) antisemitic bias, i.e., the bias between the dominant Christian group (T_1) and the targeted Jewish group (T_2),⁸ and (2) anti-communist bias, as the bias between the mainstream political conservatism (T_1) and the targeted communism (T_2). We expect the diachronic analysis of the degree of presence (or lack thereof) of these two bias types to point to ideological shifts in German political history. In what follows, we describe concrete bias specifications for each of the two dimensions of analysis. Table III provides an overview on all bias specifications employed in our study.

Antisemitism. We quantify the presence of antisemitism as the difference in association between ‘Christian’ terms (T_1) and ‘Jewish’ terms (T_2) with respect to some attribute term sets A_1 (e.g., terms of positive sentiment like *love*) and A_2 (e.g., negative sentiment terms like *hate*). We obtain the terms for T_1 (i.e., Christian terms) by translating the Christian bias

⁸We decided to define Judaism along a religious dimension, rather than along a racial dimension, as often perceived e.g. during the NS times. Thus, we contrasted Jews to Christians rather than Germans or Aryans. We speculate that this could lead to an underestimation of antisemitic bias.

terms introduced by [13] from English to German (e.g., *church* → *Kirche*). We directly translated each term from the original English set and then judged the suitability of the translation for our use-case of diachronic bias analyses, resulting in minor modifications with respect to the original set. For instance, we discard the term ‘Kreuz’ (*cross*) due to its participation in the collocation ‘Eisernes Kreuz’, one of the highest Prussian military accolades. Besides, we added to T_1 the following terms: ‘Pfarrer’ (*pastor*), ‘Ostern’ (*Easter*) and ‘Bibel’ (*Bible*). We devised the target term set for Judaism from scratch, starting from central terms that clearly denote this religion, e.g., ‘jüdisch’ (*Jewish*) and ‘Jude’ (*Jew*).

We observe the differences between Christian and Jewish target sets through seven different *views*, each of which is defined with one pair of attribute sets (A_1, A_2). Views b–g) correspond to five bias lines identified by [33]:

- Sentiment.** We rely on the widely-used attribute set of pleasant (A_1) and unpleasant (A_2) terms introduced by [5] and translated to German by [21];
- Religious.** Terms capturing a religious view of what is “good” and “bad”, e.g., ‘Gläubige’ (*believers*, A_1);
- Economic.** Terms capturing the “jew as greedy” stereotype;
- Patriotic.** Nationally-loaded terms, e.g., ‘patriotisch’ (*patriotic*, A_1), ‘fremd’ (*foreign*, A_2);
- Racial.** Racist-related terminology capturing the idea of superior and inferior races;
- Conspiratorial.** Terms related to “jew world conspiracy” theories, e.g., ‘loyal’ (A_1), ‘betrügerisch’ (*deceitful*, A_2);
- Ethic.** Includes ethically-loaded terms describing what is “good” and “bad”, e.g., ‘tugendhaft’ (*virtuous*, A_1).

For each of these views, terms in A_2 have been recognized by [33] to be commonly used in constructions expressing negative stereotypes about Jews: for example, the negatively biased view of greedy and economically influential Jews symbolized by the Rothschilds (economic); the Jews involved in political machinations and secret plots symbolized by George Soros (conspiratorial); or the Nazi myth of Jews as distinct and inferior race (racial).

Anti-communism. In order to capture anti-communist bias, we devise two target term sets representing mainstream political conservatism (T_1) and the targeted communism (T_2). To this end, we resort to political and historical literature [32]. For both bias specifications, we account for linguistic changes across the analyzed range of 153 years by adapting the bias specifications after 1949 (starting with the German Bundestag). Here, we focus on three views:

- Sentiment.** We employ the same set of terms capturing general sentiment as above;
- Politics.** Terms relating to a political view of “good” and “bad”, e.g., ‘wirksam’ (*effective*, A_1);
- Propaganda.** Terms capturing anti-communist propaganda.

TABLE III: Bias specifications for representing antisemitism and anti-communism employed in our study. Superscript ⁻ indicates terms which we remove and superscript ⁺ indicates terms which we add for tests on historical slices starting from 1949 (German Bundestag).

Bias Type/View	Specification terms
Antisemitism	<i>T₁ (Christentum):</i> Taufe, Katholizismus, Christentum, evangelisch, Evangelium, Jesus, Christ, christlich, katholisch, Kirche, Pfarrer, Ostern, Bibel <i>T₂ (Judentum):</i> Rabbiner, Synagoge, kosher, Talmud, orthodox, Judentum, Jude, jüdisch, Mose, mosaisch, Israel, Abraham, zionistisch, israelitisch, Israelis
View: Sentiment	<i>A₁ (+):</i> streicheln, Freiheit, Gesundheit, Liebe, Frieden, Freude, Freund, Himmel, loyal, Vergnügen, Diamant, sanft, ehrlich, glücklich, Regenbogen, Geschenk, Ehre, Wunder, Sonnenaufgang, Familie, Lachen, Paradies, Ferien <i>A₂ (-):</i> Missbrauch, Absturz, Schmutz, Mord, Krankheit, Tod, Trauer, vergiften, stinken, Angriff, Katastrophe, Haß, verschmutzen, Tragödie, Scheidung, Gefängnis, Armut, häßlich, Krebs, töten, faul, erbrechen, Qual
View: Patriotism	<i>A₁ (+):</i> patriotisch, vaterländisch, reichstreu, nationalgesinnt, nationalstolz, königstreu, staatstreu, nationalgefühl, Volkstum ⁺ , Patriotismus ⁺ , Patriot ⁺ <i>A₂ (-):</i> nichtdeutsch, fremdländisch, fremd, undeutsch, vaterlandslos, reichsfeind, landesverräter reichsfeindlich, unpatriotisch, antideutsch, deutschfeindlich, unstürzler, staatsfeindlich, dissident, heimatlos, separatistisch ⁺ , staatenlos ⁺
View: Economy	<i>A₁ (+):</i> geben, Großzügigkeit, großzügig, selbstlos, genügsam, Großmut, uneigennützig, sparsam, Proletariat, Armut, Industriearbeiter <i>A₂ (-):</i> nehmen, Gier, gierig, egoistisch, habgierig, Habsucht, eigennützig, verschwenderisch, Bourgeoisie, Wohlstand, Bankier, Wucher ⁺
View: Conspiracy	<i>A₁ (+):</i> loyal, Kamerad, Ehrlichkeit, ersichtlich, aufrichtig, vertrauenswürdig, wahr, ehrlich, unschuldig, freundschaftlich, hell, zugänglich, machtlos, ohnmacht, untertan <i>A₂ (-):</i> illoyal, Spitzel, Verrat, geheim, hinterlistig, betrügerisch, unwahr, zweifelhaft, Verbrecher, bedrohlich, dunkel, geheimnis, einflussreich, weltmacht, herrschaft, verschwörung
View: Religion	<i>A₁ (+):</i> Gläubige, geistlich, Engel, heilig, fromm, geheiligt, göttlich, ehrwürdig, treu, gläubig, religiös <i>A₂ (-):</i> Atheist, weltlich, Teufel, irdisch, atheistisch, heidnisch, gottlos, verflucht, untreu, ungläubig, irreligiös, Gotteslästerung
View: Race	<i>A₁ (+):</i> normal, überlegenheit, gleichheit, angenehm, freundlich, ehrenwert, sympathie, akzeptiert, besser, national, rein, überlegen, sauber, ehrenhaft <i>A₂ (-):</i> seltsam, unterlegenheit, ungleichheit, unangenehm, boshaft, schändlich, hass, abgelehnt, schlechter, fremdländisch, unrein, unterlegen, schmutzig, verseucht, schädlich, niederträchtig
View: Ethics	<i>A₁ (+):</i> bescheiden, sittlich, anständig, tugendhaft, charakterfest, würdig, treu, moralisch, ehrlich, gesittet, gewissenhaft, vorbildlich <i>A₂ (-):</i> unbescheiden, unsittlich, unanständig, lüstern, korrupt, unwürdig, untreu, unmoralisch, unehrlich, verdorben, gewissenlos, barbarisch
Anti-communism	<i>T₁ (Conservatism):</i> Konservatismus, Tradition, Geschichte, Christentum ⁻ , Adel ⁻ , Monarchie ⁻ , Mittelalter ⁻ , Stände ⁻ , Werte, Moral, König ⁻ , Kaiser ⁻ , Hierarchie, Identität, Kontinuität, Sicherheit, Grundbesitz ⁻ , Autorität, Legitimität, Ordnung, Religion ⁻ , Kirche ⁻ , Erhaltung, Treue ⁻ , Tugend ⁻ , Bräuche, Sitten, Bewahrung, Gottesgnadentum ⁻ , Ständeordnung ⁻ , Restauration ⁻ , Bürger ⁺ , Bürgertum ⁺ , Regierung ⁺ , Wertordnung ⁺ , Bürgerlichkeit ⁺ , Stabilität ⁺ , Wohlstand ⁺ <i>T₂ (Communism):</i> Sozialismus ⁻ , Kommunismus, Proletariat, Arbeiter ⁻ , Klassengesellschaft, Klasse ⁻ , Revolution, Aufklärung ⁻ , Gemeinschaft ⁻ , Gerechtigkeit ⁻ , Armut ⁻ , Kapital, Gleichheit ⁻ , Chancen ⁻ , Freiheit ⁻ , Arbeiterklasse ⁻ , Solidarität ⁻ , Partei ⁻ , Verstaatlichung, Gewerkschaft ⁻ , Marx, Engels, Vergesellschaftung, Gemeineigentum, Widerstand, Kollektivierung, Arbeiterbewegung ⁻ , Aufstand ⁻ , Lenin ⁺ , Planwirtschaft ⁺ , Klassenkampf ⁺ , Proletariat ⁺ , Revolution ⁺ , Produktionsmittel ⁺ , Diktatur ⁺ , Bolschewiki ⁺ , Oktoberrevolution ⁺ , Räte ⁺ , Sowjetunion ⁺
View: Sentiment	<i>A₁ (+):</i> streicheln, Freiheit, Gesundheit, Liebe, Frieden, Freude, Freund, Himmel, loyal, Vergnügen, Diamant, sanft, ehrlich, glücklich, Regenbogen, Geschenk, Ehre, Wunder, Sonnenaufgang, Familie, Lachen, Paradies, Ferien <i>A₂ (-):</i> Missbrauch, Absturz, Schmutz, Mord, Krankheit, Tod, Trauer, vergiften, stinken, Angriff, Katastrophe, Haß, verschmutzen, Tragödie, Scheidung, Gefängnis, Armut, häßlich, Krebs, töten, faul, erbrechen, Qual
View: Politics	<i>A₁ (+):</i> sozial, progressiv, gemeinschaftlich, gemeinsam, zivilisiert, bewährt, wirksam, etabliert, demokratisch, hoch, möglich, fortschrittlich, gemäßigt, machbar, realistisch, früh, kontinuierlich, legitim, verlässlich, aufrichtig, intellektuell, sicher, Sicherheit, Fortschritt, pragmatisch, Vertrauen, Wandel ⁺ , sachlich ⁺ , Gewinn ⁺ , fähig ⁺ <i>A₂ (-):</i> unsozial, radikal, extrem, gefährlich, gefährdend, niedrig, nieder, unmöglich, undemokratisch, unrealistisch, spät, unlegitim, Gefahr, unehrlich, unaufrichtig, unintellektuell, unsicher, schwer, schwierig, Misstrauen, Stillstand ⁺ , Skandal ⁺ , skandalös ⁺ , Zukunft ⁺ , unsachlich ⁺ , Verlust ⁺ , unfähig ⁺
View: Propaganda	<i>A₁ (+):</i> Kamerad ⁻ , Kameraden ⁻ , Kameradschaft ⁻ , kameradschaftlich ⁻ , Vaterland ⁻ , Patriot ⁻ , Ehre, ehrlich, Einsatz, Untertan ⁻ , rein ⁻ , wir, Heimat ⁻ , deutsch, Deutschland ⁻ , Truppe ⁻ , Nationalstolz ⁻ , patriotisch, Volk, Befreiung ⁻ , Front ⁻ , Wahrheit, wahr, aufrichtig ⁺ , gemeinschaftlich ⁺ , Wertegemeinschaft ⁺ , Mitte ⁺ , Frieden ⁺ , Partnerschaft ⁺ , Integration ⁺ , Wandel ⁺ <i>A₂ (-):</i> Sabotage ⁻ , Saboteur ⁻ , Betrüger, Betrug, Gauner ⁻ , Schwindel ⁻ , Schwindler ⁻ , Parasit ⁻ , Volksfeind ⁻ , Reichsfeind ⁻ , undeutsch, unpatriotisch, reichsfeindlich, Volksverräter ⁻ , Spion ⁻ , Bolschewist ⁻ , fremd, unrein ⁻ , Kommunist, Spitzel ⁻ , anders, Lüge, Lügner, Dissident, Feind, Diktatur, Verschwörung ⁻ , verschwörerisch ⁻ , unehrlich ⁻ , feindlich ⁻ , Schmarotzer ⁺ , Elite ⁺ , Kriminelle ⁺ , kriminell ⁺

B. Measuring Bias in Semantic Spaces.

Most existing measures for capturing bias in text corpora (e.g., [4], [5], [9], [13], [15], [21], [22], *inter alia*), given an explicit bias specification, operate on word embeddings, namely dense vector representations of meaning in context [3], [25], [30]. In this work, we adopt two of the arguably most established bias measures based on word embeddings: Word Embedding Association Test (WEAT) [5] and Embedding Coherence Test (ECT) [9]. Training reliable word embeddings, however, requires substantial amounts of text. Given that some of the temporal slices in our DEUPARL corpus are of fairly limited size (see Table II), questions remain on whether the word embeddings induced from those slices can be sufficiently reliable. Because of this, we couple the embedding-based bias measures with a novel bias measure based on term co-occurrences and graph-based label propagation, which we dub Co-occurrence Bias Test (CBT).

Quantifying Bias in Word Embeddings. Word embeddings are dense numeric vector representations of words that aim to capture word meaning insofar that words with similar meaning get assigned vectors that are close to each other in the vector space. While there exist several different well-known algorithms for inducing word embeddings from a given corpus (see, e.g., [3], [25], [30]), they all rely on a distributional hypothesis [18] and, one way or the other, exploit the information about local word co-occurrences to derive dense semantic vectors of words, i.e., word embeddings.

Word Embedding Association Test (WEAT). Introduced by [5], the WEAT test is an adaptation of the Implicit Association Test (IAT) [26], [27]. Whereas IAT measures biases based on response times of human subjects to provided stimuli, WEAT quantifies the biases using semantic similarities between word embeddings of the same stimuli (i.e., terms in our bias specifications, see Table III). Given some word embedding space and an explicit bias specification $B_E=(T_1, T_2, A_1, A_2)$, the WEAT test statistic is computed as the differential of the associativity between T_1 and T_2 w.r.t. A_1 and A_2 , stemming from the mean similarity of their terms with terms in the attribute sets A_1 and A_2 , respectively:

$$s(B_E) = \sum_{t_1 \in T_1} s(t_1, A_1, A_2) - \sum_{t_2 \in T_2} s(t_2, A_1, A_2)$$

The association score s for an individual target term $t \in T_i$ is obtained as:

$$s(t, A_1, A_2) = \frac{1}{|A_1|} \sum_{a_1 \in A_1} \cos(\mathbf{t}, \mathbf{a}_1) - \frac{1}{|A_2|} \sum_{a_2 \in A_2} \cos(\mathbf{t}, \mathbf{a}_2)$$

where \mathbf{t} , \mathbf{a}_1 , \mathbf{a}_2 are word embeddings of terms t , a_1 , and a_2 , respectively and $\cos(\mathbf{x}, \mathbf{y}) \in [-1, 1]$ is the cosine of the angle enclosed by the vectors \mathbf{x} and \mathbf{y} . We note that $s(B_E)$ is large, for example, when terms in T_1 are positively associated with attributes in A_1 and negatively with attributes in A_2 and vice

versa for T_2 . In this case, there is a bias for T_1 to be associated with A_1 and T_2 to be associated with A_2 .

The significance of the above statistic is computed by comparing $s(B_E)$ with the scores $s(B_E^*)$ obtained for permutations of B_E , $B_E^* = (T_1^*, T_2^*, A_1, A_2)$, where T_1^* and T_2^* are equally sized partitions of $T_1 \cup T_2$. The p -value of the test is the probability of $s(B_E^*) > s(B_E)$. Finally, the *bias effect size*, i.e., the ‘‘amount’’ of bias, is computed as the normalized measure of separation between association distributions:

$$\frac{\mu(\{s(t_1, A_1, A_2)\}_{t_1 \in T_1}) - \mu(\{s(t_2, A_1, A_2)\}_{t_2 \in T_2})}{\sigma(\{s(t, A_1, A_2)\}_{t \in T_1 \cup T_2})}$$

with μ as the mean value and σ as the standard deviation.

Embedding Coherence Test (ECT). The ECT test, originally proposed by [9] and later adjusted by [22] quantifies the amount of explicit bias by comparing vectors of target sets T_1 and T_2 , obtained by averaging the embeddings of their constituent terms, with the vectors from a single attribute set $A = A_1 \cup A_2$. We first compute the embeddings for target sets T_1 and T_2 by averaging the vectors of their terms:

$$\mu_1 = \frac{1}{|T_1|} \sum_{t_1 \in T_1} \mathbf{t}_1; \quad \mu_2 = \frac{1}{|T_2|} \sum_{t_2 \in T_2} \mathbf{t}_2.$$

Next, for both μ_1 and μ_2 it computes the (cosine) similarities with vectors of all attribute terms $a \in A$:

$$\mathbf{s}_1 = [\cos(\mu_1, \mathbf{a}_i)]_{i=1}^{|A|}; \quad \mathbf{s}_2 = [\cos(\mu_2, \mathbf{a}_i)]_{i=1}^{|A|}.$$

The two resulting vectors of similarity scores, \mathbf{s}_1 (for T_1) and \mathbf{s}_2 (for T_2) are used to obtain the final ECT score: the Spearman’s correlation between the rank orders of \mathbf{s}_1 and \mathbf{s}_2 . Unlike for WEAT where a larger bias effect implies larger bias, the higher the ECT correlation, the lower the bias.

Measuring Bias with Co-occurrence Graphs. Both ECT and WEAT crucially depend on the quality of the semantic information encoded in the underlying word embedding space, which typically require large text corpora. We now propose a novel bias measure based on word co-occurrence graphs and graph-based label propagation, specifically designed to be applicable to small corpora as well.

We first compute scores between terms in the corpus indicating the level of their lexical association, based on their co-occurrence frequency. Concretely, we adopt Positive Pointwise Mutual Information (PPMI) as the measure of lexical association between terms. For a pair of words (w_1, w_2) , pointwise mutual information (PMI) is computed as the probability of joint occurrence $P(w_1, w_2)$ of w_1 and w_2 (in a co-occurrence window of fixed size), normalized with the product of the probabilities of individual terms, $P(w_1)$ and $P(w_2)$:

$$\text{PMI} = \log \frac{P(w_1, w_2)}{P(w_1) \cdot P(w_2)}$$

A PMI score of 0 (i.e., probability ratio of 1) means two words appear together exactly as frequently as one would

expect from their individual frequencies. Scores lower than zero imply a reliable lack of association between terms. This is why PPMI replaces negative PMI scores with 0, i.e., $\text{PPMI} = \max\{0, \text{PMI}\}$.

Co-occurrence Bias Test (CBT). The PPMI scores computed between all pairs of terms effectively define an undirected weighted graph which can then be used for semi-supervised label propagation of scores – from a small subset of nodes for which such scores are known to all other (unlabeled) nodes in the graph. We assign the labels to the nodes corresponding to the terms of the two attribute lists: the terms from the positive attribute list A_1 are assigned the score 1, whereas the terms from the negative attribute lists are labeled with 0. We next employ the semi-supervised graph-based label propagation algorithm named Harmonic Function Label Propagation (HFLP) [14], [37], [38] to induce the scores for all other nodes (i.e., terms) in the PPMI graph. Let \mathbf{W} be the weighted adjacency matrix of the PPMI graph with w_{ij} as the PPMI score between the i -th and j -th vocabulary term, and let \mathbf{D} be the corresponding diagonal degree matrix with entries $\mathbf{D}_{ii} = \sum_j w_{ij}$ and $\mathbf{D}_{ij} = 0$ for $i \neq j$. Assuming we index all labeled nodes (i.e., nodes corresponding to terms from A_1 and A_2) in W before all unlabeled nodes, the Laplacian of the graph, i.e., $\Delta = \mathbf{W} - \mathbf{D}$, can then be partitioned as:

$$\Delta = \begin{pmatrix} \Delta_{ll} & \Delta_{lu} \\ \Delta_{ul} & \Delta_{uu} \end{pmatrix}$$

Let \mathbf{f}_l be the binary vector of labels of labeled nodes (i.e., terms from A_1 and A_2). HFLP then offers a closed form solution for the scores of unlabeled nodes:

$$\mathbf{f}_u = -\Delta_{uu}^{-1} \Delta_{ul} \mathbf{f}_l.$$

The scores $f \in \mathbf{f}_u$ assigned to unlabeled nodes will then be in between 0 and 1. We are now interested in the scores assigned to the terms from the two target lists T_1 and T_2 , respectively. Concretely, we are interested in whether the mean of the scores for terms in T_1 statistically significantly differs from the mean of the scores assigned to terms in T_2 : if so, then we are observing a significant bias effect with respect to the attribute sets A_1 and A_2 . To this end, we apply the Student’s (unpaired, two-tailed) t -test and interpret the value of the t -statistic directly as the size bias effect:

$$t(T_1, T_2) = \frac{\mu_{T_1} - \mu_{T_2}}{s \cdot \sqrt{\frac{1}{|T_1|} + \frac{1}{|T_2|}}}$$

with $\mu(T)$ as the mean score that HFLP assigned to the terms in T and s as the estimator of the pooled standard deviation of the two sets of target scores:

$$s = \sqrt{\frac{(|T_1| - 1) \cdot s_{T_1}^2 + (|T_2| - 1) \cdot s_{T_2}^2}{|T_1| + |T_2| - 2}}$$

where $s_{T_1}^2$ and $s_{T_2}^2$ are the variances of the scores induced with HFLP for terms in T_1 and T_2 , respectively.

IV. EXPERIMENTS

We present our experiments on measuring bias in our DEUPARL corpus of German parliamentary proceedings.

A. Experimental Setup

For the embedding-based approach, we train Word2Vec CBOW models [25] on each historical slice independently. However, here, we exclude the NS slice due to its small size. We employ the *gensim* framework for inducing the embeddings with an embedding size of 200 and the window size as well as the minimum count for unigrams set to 5. For computing the PPMI matrices, we set the window size to 5 and the minimum term count to 10 with the exception of the NS slice for which we set the threshold to 2.

B. Results

The results for our embedding-based and the graph-based bias measurement approaches are depicted in Figures 5 and Figures 6, respectively. In the following, we first focus on antisemitic bias and then discuss the results for the anti-communist bias.

Antisemitic Bias. Figure 5a shows the mean and 95% confidence interval aggregated across the WEAT effect sizes for all bias views. By our specification, T_1 refers to the Christian group, T_2 to the Jewish group, A_1 to positively connotated words and A_2 to negatively connotated words. Thus, a large value of WEAT can be interpreted as antisemitic bias. As it can be seen, antisemitism measured by WEAT continuously rises starting from KR1 (Kaiserreich 1, 1867 – 1890) over KR2 (Kaiserreich 2, 1890 – 1918) until it reaches its peak in WR (1918 – 1942). Afterwards, except for a slight increase in CDU2 (1982 – 1998), the measured bias continuously drops and falls under the original level from KR1. Although there are slight differences, overall trends measured by ECT are similar (lower correlation indicates higher bias): the global minimum of the ECT mean can be found in the historical slice from the Weimar Republik (1918 – 1942), and the estimated bias from CDU2 (1982 – 1998) is smaller than the estimates obtained between KR1 and WR. The different bias views draw an interesting picture: for instance, while patriotic bias seems to drop, the bias measured towards attributes reflecting sentiment appears to stay consistent, even in the period of the Bundesrepublik. A reason for this is the fact that the discourse in the German parliament shifted towards discussing Germany’s responsibility for the Holocaust and the need for remembrance, which is *per se* not antisemitically biased, but carries a negative sentiment. Generally, the results obtained using the graph-based approach (CBT), shown in Figure 6b, are less significant but confirm the observations obtained using the embedding-based approach.

Anti-communist Bias. The mean of the aggregated WEAT scores for anti-communist bias obtained with different views is starting from a local optimum in KR1 (1867 – 1890), and then decreasing until they peak in CDU1 (1949 – 1969, see

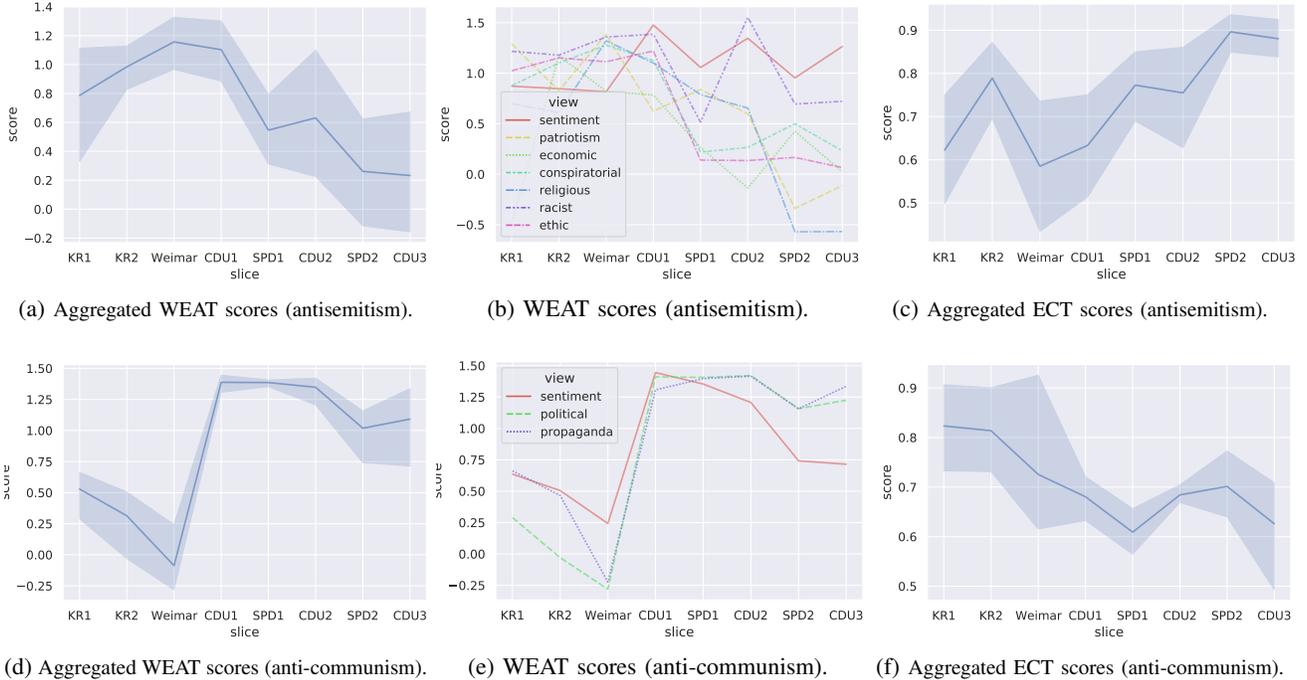


Fig. 5: Results of the embedding-based approach. WEAT: higher scores means lower bias; ECT: higher scores, lower bias.

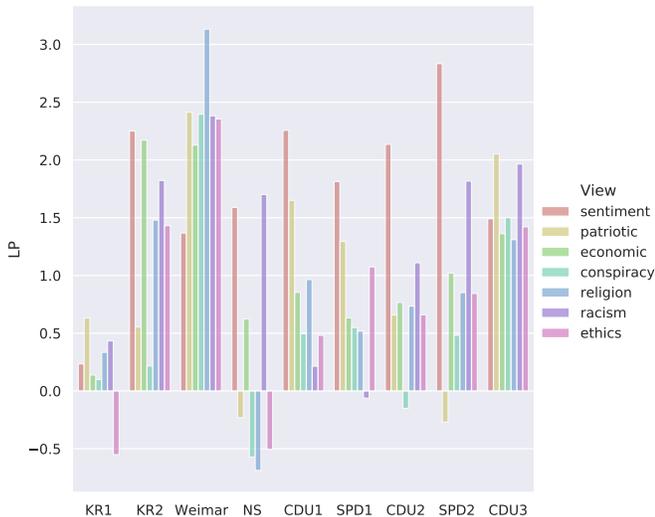
Figure 5d). Interestingly, compared to this finding, the ECT scores depicted in Figure 5f) show a different picture, while the overall trend remains similar: generally, the scores seem to vary more heavily and, measured with ECT (all scores are significant at $\alpha < 0.05$), the bias seems to steadily rise until SPD1, decrease a bit and then increase again to CDU3 (2005 – 2020). Intuitively, though the results are not entirely conclusive, all of these peaks as well as the shared trend of the increased bias over the years make sense: (i) We speculate that the peak in KR1 reflects the prominence of anti-socialist laws from that period; (ii) After the end of the monarchy in Germany, the communists formally founded the communist party *Kommunistische Partei Deutschlands* (KPD) and the SPD continued to establish itself as a moderate-leftist party, entering into the Weimarer Coalition. Until the 1932 election (6th legislative period in WR), the SPD remained the strongest political force in the Reichstag. Generally, the SPD gained a reputation for pursuing compromises between the different ideological factions while the KPD strongly advocated for abolishing the parliamentary system and even constituted an anti-parliamentary alliance with the NSDAP in the late WR years [16]. In a qualitative analysis, we found evidence indicating the moderate left’s desire to distance themselves from extremist communism. For instance, they denounce the bad reflection of the communist tactics on honest desire to improve conditions for the working class; (ii) after World War II (CDU1), the establishment of the Iron Curtain led to a clash of Western and Eastern ideologies [8]. In our study, we only include protocols from the Bundesrepublik, which reflect just the Western perspective. (iii) The reason for the peak in CDU3

may be attributed to the major parties’ general stance against populist and anti-democratic forces, which have been gaining momentum in recent years [19]. The different views indicate that, initially (KR1), the bias with respect to the sentiment attribute sets is higher, while later on, and especially in the last periods, it gets surpassed by the bias with respect to politics and propaganda terms.

Consistency across Measures. As already noticed by Lauscher et al. [22], different bias measurement approaches capture different aspects of bias in the distributional space and therefore sometimes deviate in terms of the amount of bias they reveal. In line with their findings, we note some inconsistencies: for instance, the amount of bias measured with CBT is generally lower than the one estimated with WEAT. This is not surprising, since WEAT has been shown to overestimate bias [12]. We therefore acknowledge the importance of running several measures in parallel in order to capture the most pervasive trends.

V. RELATED WORK

Measuring Bias in Semantic Spaces. Bolukbasi et al. [4] were the first to demonstrate that word embeddings induced from human-created corpora encode stereotypical human biases. More specifically, they showed that many distributional vector spaces allow for algebraically building biased analogies, such as the famous example “*man is to computer programmer as woman is to homemaker*”, an example of a sexist analogy. Following up on this, Caliskan et al. [5] presented



(a) CBT scores (antisemitism).



(b) CBT scores (anti-communism).

Fig. 6: Results for the graph-based approach.

WEAT, inspired from the Implicit Association Test [26] from psychology, which measures biases in terms of associative difference in similarity between term sets. The WEAT bias specifications were later translated to several languages, e.g., by Chaloner and Maldonado [7] and Lauscher et al. [21], [23]. We use the German sentiment attribute sets from the latter. A similar group of authors assembled a larger framework of bias evaluation measures [22] and proposed the notion of implicit and explicit bias specifications. This framework includes, for instance, the Embedding Coherence Test [9] and the Implicit Bias Test [15]. We adopt their notion of bias. Follow up research also proposed evaluation methods and resources for contextualized embedding methods [1], [20], [36, *inter alia*].

Natural Language Processing for the Analysis of Historical Corpora. The studies closest to ours were presented by Garg et al. [13] and Tripodi et al. [33]. Similar to us, Garg et al. quantify bias in historical corpora employing word vector spaces as proxies. However, they analyze gender and ethnic bias in U.S.-centric corpora, while we specifically focus on antisemitic and anti-communist bias in German parliamentary data, which allows us to capture ideological trends. Tripodi et al. focus, similar to us, on antisemitic bias, but they analyze French books and periodicals issues. We additionally propose a graph-based bias evaluation, which allows us to deal with smaller historical slices. Using word embeddings for studying change in historical corpora was introduced by Hamilton et al. [17] and Eger and Mehler [10]: however, while their contribution is on determining statistical laws of semantic change in diachronic word embeddings, we focus here instead on whether changes in bias measurements from dense and sparse embeddings align against known historical changes. Despite their enormous popularity, there exist alternatives to word embeddings as methodology to analyze historical corpora,

including topic modeling [2], [35] and event extraction [31].

VI. CONCLUSION

In this work, we have investigated how to employ the notion of bias in distributional word vector spaces for understanding ideological shifts in a large historical corpus of German parliamentary debates. The results obtained with the various available measures sometimes exhibit different amounts of bias. However, many shared trends can be found, which can, in turn, be attributed to events related to ideological shifts within the German history, e.g., the establishment of the iron curtain or peaks in antisemitism leading up to the NS time. Along with our study, we release DEUPARL hoping to fuel further research on computational understanding of German parliamentary proceedings.

Using the notion of bias, we are able to capture commonly accepted historical trends. While the embedding-based method has the advantage that a large range of bias evaluation measures for estimating the extend of stereotyping present in text representation models is available, it has the drawback in that the sizes of the slices need to be sufficiently large to induce reliable embeddings. In these cases, we proposed to ‘fall-back’ to more traditional meaning representations, i.e., co-occurrence graphs, and leverage label propagation algorithms to infer bias from sparse co-occurrence counts.

Our analysis shows that main ideological shifts can be read from historical corpora through the quantification of bias of semantic spaces induced from them. Our method has the potential to enable a bias-centric exploration of textual collections from experts and thus calls for the future integration of our quantitative analysis with qualitative ones in the manner of *distant reading* from historians. We also plan to tackle other languages than German, for a contrastive evaluation of multilingual cross-temporal bias, in future work.

ACKNOWLEDGMENTS

We thank the reviewers for insightful comments. We thank Niklas Friedrich for the independent reimplementations of the approaches and reproduction of the empirical results.

REFERENCES

- [1] C. Basta, M. R. Costa-jussà, and N. Casas, “Evaluating the underlying gender bias in contextualized word embeddings,” in *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, 2019, pp. 33–39.
- [2] D. M. Blei, “Topic modeling and digital humanities,” *Journal of Digital Humanities*, vol. 2, no. 1, 2012.
- [3] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching word vectors with subword information,” *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 135–146, 2017.
- [4] T. Bolukbasi, K. Chang, J. Y. Zou, V. Saligrama, and A. T. Kalai, “Man is to computer programmer as woman is to homemaker? debiasing word embeddings,” in *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, Eds., 2016, pp. 4349–4357.
- [5] A. Caliskan, J. J. Bryson, and A. Narayanan, “Semantics derived automatically from language corpora contain human-like biases,” *Science*, vol. 356, no. 6334, pp. 183–186, 2017.
- [6] V. D. Carlo, F. Bianchi, and M. Palmonari, “Training temporal word embeddings with a compass,” in *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019*, 2019, pp. 6326–6334.
- [7] K. Chaloner and A. Maldonado, “Measuring gender bias in word embeddings across domains and discovering new gender bias word categories,” in *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, 2019, pp. 25–32.
- [8] S. Kreuzberger and D. Hoffmann, “Antikommunismus und politische kultur in der bundesrepublik deutschland,” *Geistige Gefahr und Immunisierung der Gesellschaft, München*, 2014.
- [9] S. Dev and J. Phillips, “Attenuating bias in word vectors,” in *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, 2019, pp. 879–887.
- [10] S. Eger and A. Mehler, “On the linearity of semantic change: Investigating meaning variation via dynamic graph models,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2016, pp. 52–58.
- [11] S. Eger, T. von der Brück, and A. Mehler, “A comparison of four character-level string-to-string translation models for (OCR) spelling error correction,” *The Prague bulletin of mathematical linguistics*, vol. 105, no. 1, pp. 77 – 99, 2017.
- [12] K. Ethayarajh, D. Duvenaud, and G. Hirst, “Understanding undesirable word embedding associations,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 1696–1705.
- [13] N. Garg, L. Schiebinger, D. Jurafsky, and J. Zou, “Word embeddings quantify 100 years of gender and ethnic stereotypes,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 115, no. 16, pp. E3635–E3644, 2018.
- [14] G. Glavaš, F. Nanni, and S. P. Ponzetto, “Unsupervised cross-lingual scaling of political texts,” in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2017, pp. 688–693.
- [15] H. Gonen and Y. Goldberg, “Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds., 2019, pp. 609–614.
- [16] W. Halder, *Innenpolitik im Kaiserreich: 1871-1914*. Wissenschaft Buchgesellschaft, 2003.
- [17] W. L. Hamilton, J. Leskovec, and D. Jurafsky, “Diachronic word embeddings reveal statistical laws of semantic change,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016, pp. 1489–1501.
- [18] Z. S. Harris, “Distributional structure,” *Word*, vol. 10, no. 2-3, pp. 146–162, 1954.
- [19] U. Jun, *Die Parteien nach der Bundestagswahl 2017: Aktuelle Entwicklungen des Parteienwettbewerbs in Deutschland*. Springer-Verlag, 2020.
- [20] K. Kurita, N. Vyas, A. Pareek, A. W. Black, and Y. Tsvetkov, “Measuring bias in contextualized word representations,” in *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, 2019, pp. 166–172.
- [21] A. Lauscher and G. Glavaš, “Are we consistently biased? multidimensional analysis of biases in distributional word vectors,” in *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, 2019, pp. 85–91.
- [22] A. Lauscher, G. Glavaš, S. P. Ponzetto, and I. Vulić, “A general framework for implicit and explicit debiasing of distributional word vector spaces,” in *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020*, 2020, pp. 8131–8138.
- [23] A. Lauscher, R. Takieddin, S. P. Ponzetto, and G. Glavaš, “AraWEAT: Multidimensional analysis of biases in Arabic word embeddings,” in *Proceedings of the Fifth Arabic Natural Language Processing Workshop*, 2020, pp. 192–199.
- [24] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, “A survey on bias and fairness in machine learning,” 2019. [Online]. Available: <http://arxiv.org/abs/1908.09635>
- [25] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” in *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 2*, 2013, pp. 3111–3119.
- [26] B. A. Nosek, M. R. Banaji, and A. G. Greenwald, “Harvesting implicit group attitudes and beliefs from a demonstration web site,” *Group Dynamics: Theory, Research, and Practice*, vol. 6, no. 1, p. 101, 2002.
- [27] B. A. Nosek, A. G. Greenwald, and M. R. Banaji, “Understanding and using the implicit association test: Ii. method variables and construct validity,” *Personality and Social Psychology Bulletin*, vol. 31, no. 2, pp. 166–180, 2005.
- [28] E. Ntoutsi, P. Fafalios, U. Gadiraju, V. Iosifidis, W. Nejdl, M. Vidal, S. Ruggieri, F. Turini, S. Papadopoulos, E. Krasanakis, I. Kompatsiaris, K. Kinder-Kurlanda, C. Wagner, F. Karimi, M. Fernández, H. Alani, B. Berendt, T. Kruegel, C. Heinze, K. Broelemann, G. Kasneci, T. Tiropanis, and S. Staab, “Bias in data-driven AI systems - an introductory survey,” *CoRR*, vol. abs/2001.09762, 2020. [Online]. Available: <https://arxiv.org/abs/2001.09762>
- [29] K. Ortman, A. Roussel, and S. Dipper, “Evaluating off-the-shelf nlp tools for german,” in *KONVENS*, 2019.
- [30] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [31] M. Rovera, F. Nanni, and S. P. Ponzetto, “Event-based access to historical italian war memoirs,” *Journal on Computing and Cultural Heritage*, vol. 14, no. 1, 2021.
- [32] A. Schildt, *Konservatismus in Deutschland: von den Anfängen im 18. Jahrhundert bis zur Gegenwart*. Beck, 1998.
- [33] R. Tripodi, M. Warglien, S. Levis Sullam, and D. Paci, “Tracing antisemitic language through diachronic embedding projections: France 1789-1914,” in *Proceedings of the 1st International Workshop on Computational Approaches to Historical Language Change*, 2019, pp. 115–125.
- [34] L. van der Maaten and G. Hinton, “Visualizing data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.
- [35] T.-I. Yang, A. Torget, and R. Mihalcea, “Topic modeling on historical newspapers,” in *Proceedings of the 5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, 2011, pp. 96–104.
- [36] J. Zhao, T. Wang, M. Yatskar, R. Cotterell, V. Ordonez, and K.-W. Chang, “Gender bias in contextualized word embeddings,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 629–634.
- [37] X. Zhu, Z. Ghahramani, and J. D. Lafferty, “Semi-supervised learning using gaussian fields and harmonic functions,” in *Proceedings of the 20th International conference on Machine learning (ICML-03)*, 2003, pp. 912–919.
- [38] X. Zhu and A. B. Goldberg, “Introduction to semi-supervised learning,” *Synthesis lectures on artificial intelligence and machine learning*, vol. 3, no. 1, pp. 1–130, 2009.