
BERT Model Analysis of Dutch Perceptions on Countries during Political Upheaval

Annet Prudon (*5910447*)

Master's Thesis
MSc Applied Data Science
Faculty of Science
Utrecht University

Supervised by Dr. E. J. Meijers, and Dr. M. J. Smit

28 June 2024



**Utrecht
University**

Abstract

This thesis explores how perceptions of countries are shaped and specifically influenced by political upheaval. Perceptions of places are often not created by direct experiences but through media, art, news, and social discourse. These perceptions are dynamic, continuously evolving due to the narratives and discourses generated by a wide range of actors. To investigate the perception of regions and the impact of political upheaval, this study utilizes BERT models (BERTopic and robBERT v2) for topic modelling and sentiment analysis on parliamentary debate transcripts from the Dutch House of Representatives. The findings indicate that rapid changes, such as the Russian invasion of Ukraine, are more easily detectable in terms of topics and sentiment compared to gradual changes like Brexit. However, the results are preliminary, suggesting a need for further research with refined methodologies. Future work should include applying the analysis to the entire dataset and comparing it with Ukraine and UK filtered datasets. Enhancing data preprocessing steps is also essential to improve topic interpretability. Additionally, sentiment analysis should be conducted using a robBERT v2 model trained specifically on parliamentary debates. These improvements aim to strengthen the conclusions of this thesis and facilitate broader research on the impact of political upheaval in other countries.

TABLE OF CONTENTS

1	Introduction	1
2	Data and Methods	3
2.1	Data	3
2.2	Preprocessing	4
2.3	Methods	5
2.3.1	Topic modelling	5
2.3.2	Sentiment Analysis	6
3	Results	6
3.1	Topic modelling	6
3.1.1	Ukraine	6
3.1.2	United Kingdom	8
3.2	Sentiment Analysis	10
3.2.1	Ukraine	10
3.2.2	United Kingdom	11
4	Discussion	13
5	Conclusion	14
	References	15
A	Appendix: Extra analysis output	17
B	Appendix: Programming code	23
B.1	XML to DataFrame code	23
B.2	Preprocessing code	25
B.3	BERTopic (topic modelling) implementation code	28
B.4	robbERT (sentiment analysis) implementation code	30

1 Introduction

Since World War II there have been multiple initiatives to create an unified Europe, to fulfill the need for a cohesive European identity. The identity of a place and its people is a well-researched subject. However, most research overlooks the perspective of one country from another, which is essential for understanding the perception of other nations and fostering a sense of unity (Peng et al., 2020). There is significant interest in this subject because the perception of one country by another greatly influences its development in various fields, such as economic growth, migration trends, and tourism dynamics. Additionally, it impacts the willingness to provide aid during difficult times (Agostini and van Zomeren, 2021, Chayinska et al., 2017). These perceptions, however, are not static; they are particularly susceptible to transformation as a result of political upheaval (Ivanov et al., 2017). This sense of unity is crucial for supporting each other during difficult times, such as in the event of a conflict.

Moreover, human perception is always determined by the mental shortcuts made to save time in comprehending our surroundings. As a result, we humans form stereotypes about every subject, which influence our thoughts and actions. With a similar mechanism we humans try to grasp the world and perceive its varieties. Our perception of countries is often sculpted not by firsthand experiences, but through the lens of media, art, news, and social conversations. These perceptions are in constant change, continually reshaped by the evolving narratives and discourses from an extensive variety of actors.

To gain a better understanding of how these perceptions of other countries are formed, Natural Language Processing (NLP) can be of service. Over the past decade, NLP has taken a lift. This expansion is driven by the availability of extensive text-based resources and the development of various NLP models. As a result, a substantial number of research papers have been published, focusing on methods to extract and interpret textual information using NLP models. These studies explore different data sources and model architectures. The former varies from news articles to legal documents, and Twitter discussions to political speeches. The latter varies from logistic regression to decision trees, and neural networks to transformers. This research specifically uses BERT type models to analyse transcripts of political debates.

Regarding data sources, most research on measuring people's sentiment or opinion on specific subjects relies on Twitter data (Ruz et al., 2020, Chaudhry et al., 2021, Ilyas et al., 2020). However, since the take-over and reorganisation of X, formerly Twitter, the data is no longer accessible and individual expressions have become more extreme. Therefore, the source of a new, reliable, and well-documented text-based source is needed. In addition to social media data, political speeches and legislative documents are increasingly being analysed using NLP models (Abercrombie and Batista-Navarro, 2020, Müller-Hansen et al., 2021, Greene and Cross, 2017, Gillioz et al., 2020, Silva et al., 2021). These types of research mainly seek information about the political agenda and focus of the relevant government. This thesis focuses on the parliamentary debates of the House of Representatives in the Netherlands. Despite the availability of transcribed Dutch debates, which are publicly accessible through an API, there is limited existing research on this topic.

Regarding model architectures, there are numerous new techniques that use (slightly) different assumptions and methods to extract information from textual data. One of the most well-known and widely used models is BERT, a deep learning model introduced by Google in 2018. Consequently, it has undergone extensive testing and various model variations have been developed. Specifically, a topic modelling technique BERTopic (Grootendorst, 2022) and a large pre-trained general Dutch language model robBERT (Delobelle et al., 2020). The former can implement Dynamic Topic Modelling (DTM), which is required for this research due to the dependency between topics over time. The latter is trained on Dutch text and outperforms other sentiment analysis (BERT) models.

Therefore, this research aims to use the described BERT models to extract information from transcripts of parliamentary debates in the House of Representatives in the Netherlands. The goal is to gain insight into the evolution of public perception. Specifically, to gain a deeper understanding of how politicians discuss a particular country and how this discourse changes after political upheavals. This will be a first step towards obtaining knowledge about the imaging of a region using BERT models.

Research question: *What is the effect of political upheaval on the way of speaking about a country by analysing changes in discussed topics and sentiments through political debates?*

To be specific, this research focuses on the defining moments in Ukraine and United Kingdom (UK), respectively the Russian invasion of Ukraine and the Brexit. The former involves rapid change with a military operation where life and death are at stake. The latter is a slower process with a greater focus on financial and trade matters. Despite their inherent differences, both topics are highly relevant in the Netherlands and are therefore frequently discussed in the House of Representatives.

This research contributes to the knowledge gap in multiple ways. First, parliamentary debate transcripts are rarely used in scientific research, despite being easily accessible in the Netherlands. It is indeed a very different data source compared to Twitter data. In political debates, representatives speak with careful consideration and substantiated opinions, unlike the impulsive and extreme rhetoric often found on Twitter. Additionally, the House of Representatives should provide a good representation of the general public opinion since its members are democratically elected. Second, testing BERT models on Dutch textual data is beneficial because most NLP research primarily focuses on English texts. Especially since the two recently developed models used in this research (BERTopic and robBERT) are easily applicable and fitting for the Dutch language (more on this in the Methods section). Third, these BERT models are not often tested on parliamentary debate data. Fourth, this research helps to take the first steps in understanding the perception of regions and the impact of political upheaval on these perceptions.

This thesis is further organized as follows. In the next section, Data, the data and its source are discussed. Thereafter, in Preprocessing the preprocessing of the data is described. In Methods the implementation of the BERT models, and evaluation measures are discussed. Then, the results of the models are presented in Results, after which the implications and limitations of this research are discussed in Discussion. The

overall conclusion is given in the Conclusion part of the thesis. In the Appendix, additional results (A) and the programming code (B) can be found.

2 Data and Methods

2.1 Data

The data is obtained from the Open Data Portal from the House of Representatives, where all parliamentary debates between 25 June 2013 until 7 March 2024 are transcribed (Tweede Kamer der Staten-Generaal, n.d.-b). By using the OData API (Tweede Kamer der Staten-Generaal, n.d.-a), all reports were downloaded as XML file and transformed to a Python Pandas DataFrame with the following attributes (see Table 1). All rows with no text or title are deleted.

Table 1

Attributes extracted from the reports.

Attribute	Content
<i>Text</i>	Text said by speaker in one turn
<i>Speaker</i>	Speaker
<i>Date</i>	Date of debate
<i>Activity</i>	Debate part
<i>Title</i>	Title of the debate part (subject of debate)
<i>Speaking turn</i>	Index of speaking turn for the speaker
<i>Party</i>	Party of the speaker
<i>Cabinet</i>	Boolean to mark when speaker is (prime) minister or Secretary of State

Several key moments occurred throughout the data collection period in both Ukraine and the United Kingdom. The analysis results will be compared between these moments for both subjects (see Table 2).

Table 2

Defining moments happening in the subjects Ukraine and UK.

Ukraine events	Date	UK events	Date
Annexation Krim	<i>18 March 2014</i>	Referendum	<i>23 June 2016</i>
MH17	<i>17 July 2014</i>	Brexit	<i>31 January 2020</i>
Referendum	<i>6 April 2016</i>		
Invasion by Russia	<i>24 February 2022</i>		

2.2 Preprocessing

For data preprocessing, the following steps are taken. To optimize runtime efficiency within the limited timeframe of this research, an initial filtration is performed to reduce the dataset size. The Ukraine dataset retains the part (based on *Title*) of the debate (based on *Date*) where the text contains the string "Oekra" (not "Oekraïne" due to the diaeresis). Similarly, the United Kingdom (UK) dataset retains the debate parts containing the strings "Brexit" or "Verenigd Koninkrijk" (both to ensure that all discussions about Brexit that do not use the word *Brexit* are kept). Additionally, only debates where *Oekraïne* or *Brexit/Verenigd Koninkrijk* appear more than three times are retained. This is done to eliminate the possibility of the term being used randomly rather than being the actual subject of discussion. Next, to clean the data, all irrelevant political technicalities without substantive discussion are filtered out from the Ukraine and UK datasets: "Opening", "Mededelingen", "Stemmingen", "Afscheid", "Hamerstukken", "Regeling van werkzaamheden". At last, all three datasets (all reports, Ukraine reports, and UK reports) are saved as an csv file. The Ukraine dataset consists of 251 debate parts, and the UK dataset of 151 debate parts (refer to Table 3 for distribution across the years).

Table 3

Number of debate parts in both Ukraine and UK dataset.

Year	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
Ukraine	4	17	20	21	12	11	14	4	10	92	38	13
UK	1	11	11	16	17	21	21	16	20	12	3	2

To prepare the data for topic modelling and sentiment analysis, the following steps have been taken. First, all text is tokenized. Specifically for the UK dataset, the strings *Verenigd Koninkrijk* are replaced with *VERENIGDKONINKRIJK* to prevent them for vanishing during tokenization. Second, all punctuation is deleted, except for hyphens, *ā*, and *i* due to its occurrence in relevant words (for example, Ukraine is written as *OekraĀine*). Third, all stop words are deleted, and to complete the stop words list, an additional 191 words have been incorporated (Countwordsfree, n.d.). The following words are also included to the stop words list due to their frequent use in political politeness and terminology: *heer*, *mevrouw*, *staatssecretaris*, *minister*, *voorzitter*, *kamer*, *kabinet*, *motie*, *vraag*. Fourth, with the use of part of speech tagging from *spacy* (specifically the Dutch trained pipeline *nl_core_news_sm*), all auxiliaries are removed due to their frequent use but meaninglessness. Fifth, all empty documents are deleted. Sixth, all reports are grouped by *Title* and *Date* such that each row is one debate part (i.e. document).

2.3 Methods

2.3.1 Topic modelling

For topic modelling in this type of research, dynamic topic modelling (DTM) is recommended (Blei and Lafferty, 2006, Müller-Hansen et al., 2021, Silva et al., 2021). Unlike Latent Dirichlet Allocation (LDA), which assumes that topics are independent (i.e. uncorrelated), DTM is better suited for comparing topics over time. This DTM is a dynamic extension to LDA (Blei and Lafferty, 2006) and addresses the issue of topic evolution by allowing topics to change and evolve over time. This research uses BERTopic model for the DTM, which has proven to create coherent topics (Hadiat, 2022). Specifically, BERTopic uses pre-trained transformer-based language models to create document embeddings, cluster them, and build topic representations using the class-based TF-IDF technique (Grootendorst, 2022).

To implement BERTopic, the model `bertopic` is imported and the language is set on Dutch. The parameter fine-tuning and specifications used are as follows (Grootendorst, 2024):

- `verbose` is set on `True` to be able to track the stages of the model.
- `global_tuning` is turned off, because this will create a global representation of the topics which is not needed for this research. Only `evolutionary_tuning` is turned on since this research aims to gain insight in how the topics evolve over time.
- `Date` is reformed to a `Year` attribute such that no timestamp bins are needed. These bins will distribute all documents into equal-sized segments. However, because the data is not evenly balanced across time, using bins may lead to timestamps that are difficult to interpret. Therefore, `Year` is used instead of `Date`.
- `min_topic_size` is set on 5 because the input is a relatively small dataset.
- `nr_topics` is set to "auto" because the model already reduces the calculated topics using a Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) and determining the number of topics beforehand is tricky.
- To diversify the topic representations, `MaximalMarginalRelevance` (MMR) from `bertopic.representations` is used, which is useful for the interpretation of the topics. However, using MMR does slightly reduce the coherence score. Nevertheless, this difference is not significant, and it enhances the interpretability of topics, thereby facilitating analysis. Therefore, the diversity parameter is set to 0.1.

To be able to answer the research question, the topics will be plotted against the years along with their frequency across those years. To evaluate the performance of the model, the *U_{mass}* coherence score (i.e. how often two words occur together in the corpus) is calculated, which is best when closest to zero (Mimno et al., 2011).

2.3.2 Sentiment Analysis

For the sentiment analysis, the robBERT v2 model (Delobelle et al., 2020) is used. The robBERT v2 model is one of the best performing pre-trained language models in Dutch and outperforms other versions of BERT models, especially on sentiment analysis (Delobelle et al., 2020). The pre-trained version `DTAI-KULeuven/robbert-v2-dutch-sentiment` on Huggingface (DTAI-KULeuven, 2022) is trained on book reviews, to classify them positive or negative with a sentiment score. So, for this research, this model is utilized to classify the debate parts as positive or negative and compute the sentiment scores. To evaluate the performance of the analysis, a random debate segment is selected from each year, and three random strings within that segment are assigned sentiment labels and scores. This approach allows for checking the sentiment evaluation against intuitive feelings associated with each string.

3 Results

3.1 Topic modelling

3.1.1 Ukraine

The BERTopic model generated 25 topics for the Ukraine dataset, see Table 9 in Appendix to view all topics and its coherence scores. The overall coherence score ($U_{mass} = -0.791$) is quite good, since it is smaller than -1 and therefore close to zero.

To analyse the generated topics, all topics are grouped on similarity. For example, topics 17 and 18 are grouped under *Condemning Russian Actions* since the following words are used for these topics respectively: *Russia, Ukraine, sanctions, Netherlands, criminal*, and *sanctions, Netherlands, European*. Likewise, some topics are grouped together in *Irrelevant* because they do not seem to relate to any defining moments for Ukraine. For instance, topic 1 contains the following words: *farmers, nature, animals*, and *agriculture*. All topic groups can be seen in Table 4.

To lay these groups next to the timeline of defining moments in Ukraine, the words per year linked to the topic are checked. First, *Condemning Russian actions*, in which topic 17 is the most interesting because in 2015 and 2022 the words *gas* and *Ukraine* are linked. However, it does not have the best coherence score ($U_{mass} = -2.091$). Second, in *Safety* no linked words appear to have anything to do with the defining moments (*ATMs* and *virus*). Third, for *Refugees*, topic 9 ($U_{mass} = -0.209$) contains different countries from which refugees cause problems in the Netherlands before 2022 and in 2022 *Ukraine* appears, which could be due to the Russian invasion in 2022 (see Table 5). Fourth, in group *Energy crisis* in 2017 the words *sanction* and *list of sanctions* appear which could be related to the sanctions for Russia. Additionally,

Table 4

Topic groups for Ukraine dataset.

Group	Topics
<i>Condemning Russian actions</i>	17, 18
<i>Safety</i>	6, 20, 21, 22
<i>Refugees</i>	9, 11
<i>Energy crisis</i>	19
<i>Trade</i>	15, 23
<i>Military</i>	8, 13
<i>Referendum</i>	5
<i>Irrelevant</i>	1, 2, 3, 4, 7, 10, 12, 14, 16, 24, 25

there is a gap from 2018 to 2022 (see Figure 1) where the topic is not discussed whereas in 2022 and 2023 words related to the energy crisis caused by the lack of trading with Russia appear (*fuel, price gap*). Fifth, in *Trade* there are no words that can be linked to the defining moments, except for one report in 2015 mentioning *Ukraine, Russia, and Europe*. Sixth, in *Military* in 2022 and 2024 *Ukraine* is linked to the topics as expected. However, in 2023 only words related to Israel are linked to the topics which is highly unexpected. Seventh, for *Referendum* it is quite unexpected that there are respectively 13 and 9 debates about this topic (see Figure 1). Lastly, the words linked in the *Irrelevant* group are clearly irrelevant to the defining moments (*ISIS, media law, biodiversity*). However, these irrelevant topics do have relatively high coherence scores, ranging from -0.083 (Topic 3) to -1.164 (Topic 16).

Table 5

Words linked to topic 9 (Refugees) per year.

Year	Frequency	Words
2014	1	kinderpardon, kinderen, criterium, gemeentetoezicht, regering
2015	1	turkije, europese, vluchtelingen, nederland, libanon
2016	2	griekenland, vluchtelingen, nederland, maatregelen, relocatie
2017	1	asielzoekers, afgewezen, landen, marokko, illegaliteit
2018	2	azmani, politie, migratie, pact, nederland
2021	1	school, gemeentes, inspectie, ouders, nederland
2022	6	oekraïne, opvang, mensen, vluchtelingen, gemeenten
2023	4	zorg, dementie, nederland, aantal, zorgkosten

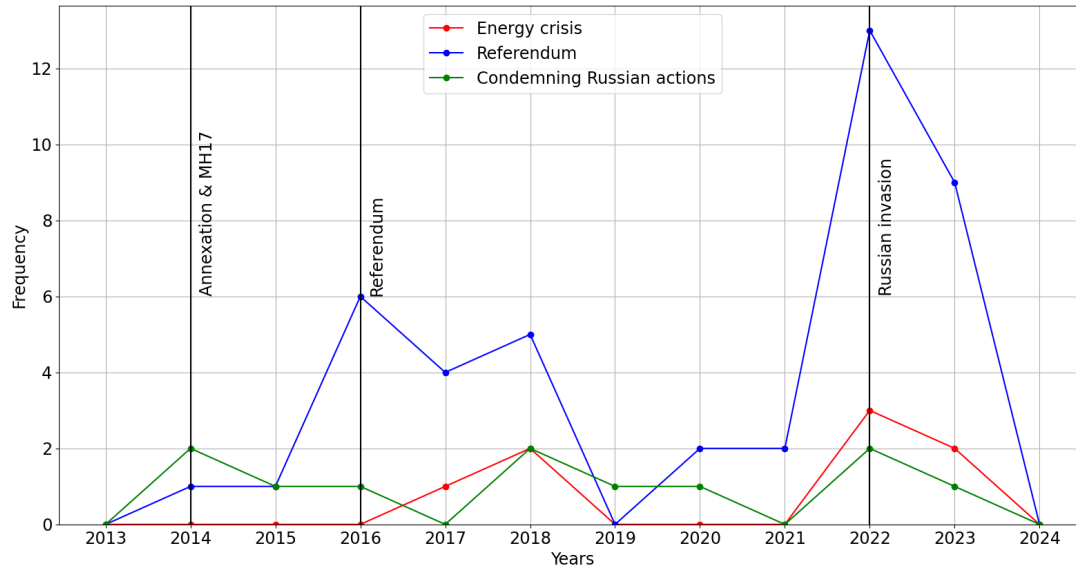


Figure 1. Occurrences of topics 19 (Energy crisis, $U_{mass} = -0.457$), 5 (Referendum, $U_{mass} = -0.280$), and 17 (Condemning Russian actions, $U_{mass} = -2.091$) per year

3.1.2 United Kingdom

The BERTopic model generated 23 topics for the UK dataset, see Table 10 in Appendix to view all topics and its coherence scores. The overall coherence score ($U_{mass} = -0.910$) is quite good, since it is smaller than -1 and therefore close to zero.

Like the Ukraine dataset, group are created on similarity. For instance, topics 2 and 17 are grouped in *International conflicts* due to the words being related to conflicts outside of the Netherlands: *Assad, Israel, Ukraine, and criminal*. On the other hand, topic 4 is grouped as irrelevant since it contains less meaningful words: *pact, naturally, think*. All topic groups can be seen in Table 6.

To analyse all topics, the linked words are viewed and determined whether they fit into the timeline of defining moments. First, *Health care and COVID*, topic 3 is linked to words such as *Ebola, children, and parents* before 2020, whereas in 2020 and 2021 it is as expected linked to *virus and testing* with way more reports (see Figure 2). After 2021 it is linked to *questions, postcovid, and report* (see Table 7). Topic 5 is focused on the organisation of health care by being linked to *medicine, health insurance, and antibiotics*. Second, topic 6 is mostly linked to words like *farmers, animals, and nature*, with a significantly higher number of reports in 2021 (see Figure 2). Additionally, topic 10 is focused on words related to sustainability: *plastic, biomass, and waste*. Third, topic 2 is focused on the turmoil in the Middle East (*Israel, Syria, Iraq*), and *China and Turkey*. Whereas topic 17 is focused on the turmoil in *Ukraine and Russia*, only appearing

Table 6

Topic groups for UK dataset.

Group	Topics
<i>Health care and COVID</i>	3,5
<i>Agriculture and climate</i>	6, 10
<i>International conflicts</i>	2, 17
<i>Financial</i>	11, 13, 15, 21, 22
<i>Infrastructure</i>	8, 12
<i>Media and Journalism</i>	19, 20
<i>Energy</i>	16
<i>Irrelevant</i>	1, 4, 7, 9, 14, 18, 23

between 2015 and 2019 (see Figure 2). Fourth, for *Financial* group, it is most interesting to see that there is a peak in 2017, linked to words as *companies*, *shareholders*, and *sector*. Fifth, for group *Infrastructure*, it is quite uninterpretable due to randomness of linked words to both topics: *spy*, *airport*, *school advise*, and *asphalt*, *nuclear power plant*, *Vietnamese*. Sixth, for group *Media and Journalism*, topic 19 is very random: *education*, *radio*, *sport*, whereas topic 20 focuses on *Turkey* and *freedom of speech* (see Figure 2). Seventh, topic 16 is also linked to quite random words, *mining*, *ProRail*, and *energy*. Lastly, for the irrelevant group, all topics are quite evenly distributed and do not change over time as one would expect along side the defining moments in the UK.

Table 7

Topic 3 over time.

Year	Frequency	Words
2014	1	fritsma, overleg, ebola, geboorteakte, kinderen
2015	2	advies, kinderen, inspectie, salarissen, banken
2019	1	vaccinatiegraad, ouders, gezondheidsraad, vaccineren, vaccinaties
2020	20	testen, virus, vragen, aantal, natuurlijk
2021	18	natuurlijk, zorg, vaccineren, testen, vragen
2022	2	dank, besmettingen, vragen, aantal, zorg
2023	3	postcovid, maatregelen, rapport, minister, vragen

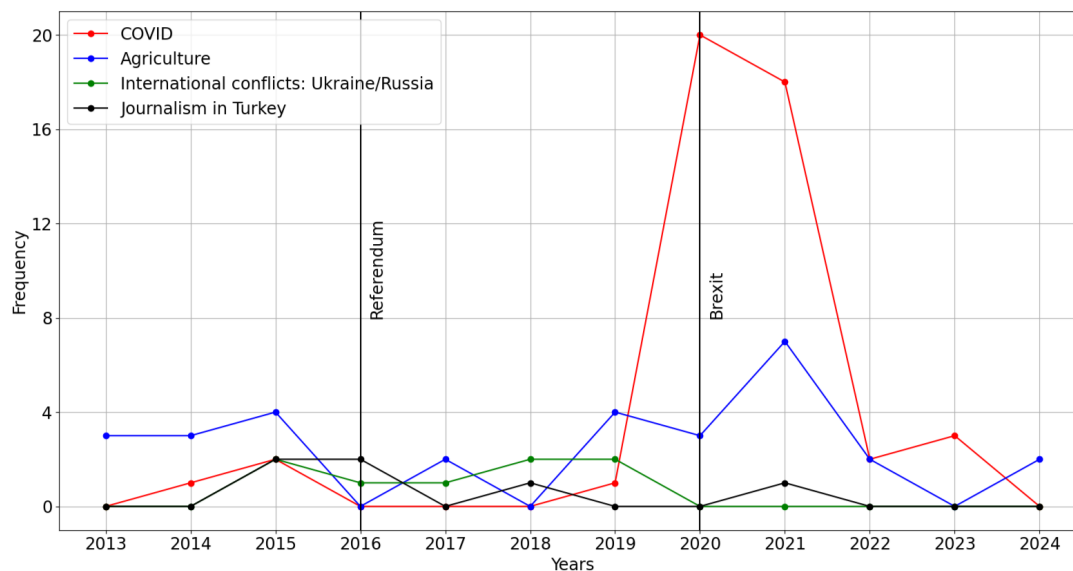


Figure 2. Occurrences of topics 3 (COVID, $U_{mass} = -0.630$), 6 (Agriculture, $U_{mass} = -0.586$), 17 (International Conflicts, $U_{mass} = -1.433$), and 20 (Journalism in Turkey, $U_{mass} = -0.319$) per year

3.2 Sentiment Analysis

3.2.1 Ukraine

The pre-trained version DTAI-KULeuven/robbert-v2-dutch-sentiment of the robbERT v2 model calculated the label (*Positive* or *Negative*) and its sentiment score for each debate part of the Ukraine dataset. The positive labelled sentiment score is clearly higher after 2020 with a similar or smaller standard deviation (see Table 8). The negative labelled sentiment score does not show a similar change. In Figure 3, there are two interesting things to see. First, only in 2016, the year of the referendum, there are more debate parts labelled as negative than positive. Second, there is a big growth of number of debate parts as well as number of positive labelled debate parts in 2022, the year of the Russian invasion. To evaluate the performance, the sentiment label and score is calculated for a random sample of words (see Table 11 in Appendix). Some words (*behaalde* and *geschript*) are evaluated as expected (respectively *Positive* and *Negative*), whereas some words (*belemmeren* and *verwoest*) are not intuitively evaluated (both *Positive*).

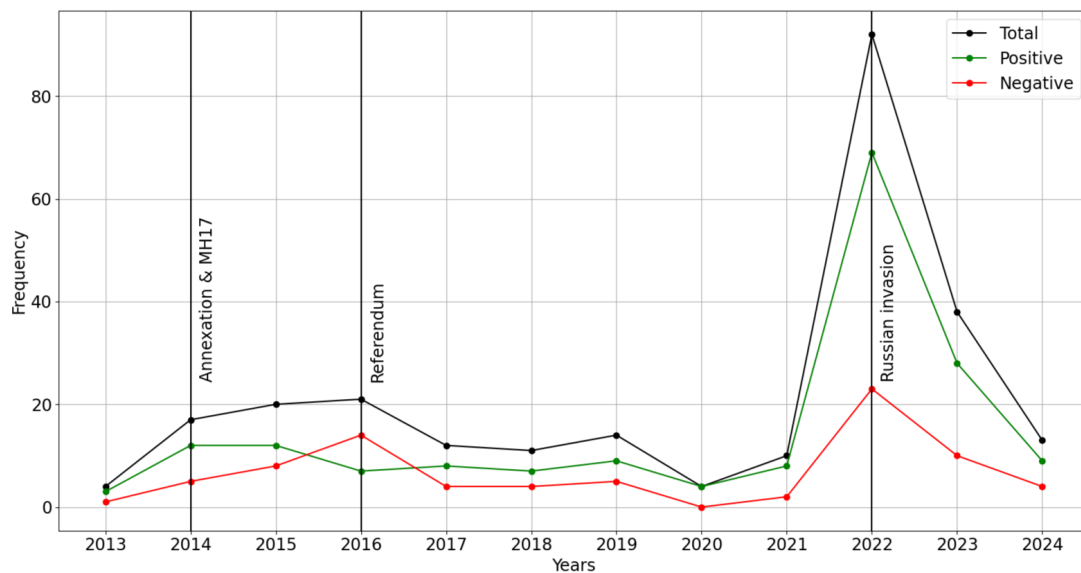


Figure 3. Frequency of sentiment labels per year for the Ukraine dataset.

3.2.2 United Kingdom

Correspondingly, the robbERT v2 model calculated the label and its sentiment score for each debate part of the UK dataset. The negative labelled sentiment score is the highest in 2015 ($\mu = 0.975$, $\sigma = 0.017$) and 2022 ($\mu = 0.995$, $\sigma = 0.002$) (see Table 8). In Figure 4, it is interesting to see that in 2016 the amount of positive labelled debate parts is the same as negative labelled debated parts, where everywhere else there are more debate parts labelled positive. Especially since the number of total debate parts has quite a lift. To evaluate the performance, the sentiment label and score is calculated for a random sample of words (see Table 12 in Appendix). Some words (*schrok* and *migrantencrisis*) are evaluated as expected (both *Negative*), whereas some words (*laagste* and *optreden*) are not intuitively evaluated (both *Positive*).

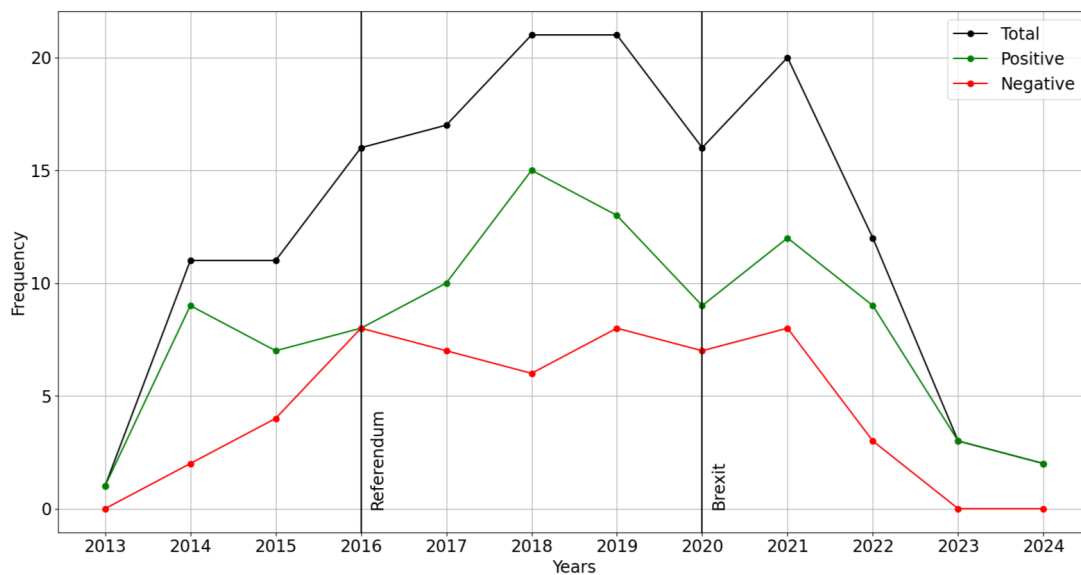


Figure 4. Frequency of sentiment labels per year for the UK dataset.

Table 8

Mean and standard deviation per year per sentiment label for both datasets.

Year	Ukraine				UK			
	Positive sentiment		Negative sentiment		Positive sentiment		Negative sentiment	
	Mean	Std Dev	Mean	Std Dev	Mean	Std Dev	Mean	Std Dev
2013	0.853	0.220	0.868	NA	0.996	NA	NA	NA
2014	0.905	0.149	0.836	0.210	0.892	0.164	0.787	0.278
2015	0.891	0.162	0.861	0.215	0.854	0.222	0.975	0.017
2016	0.872	0.160	0.803	0.210	0.931	0.071	0.896	0.156
2017	0.953	0.075	0.921	0.034	0.897	0.123	0.810	0.168
2018	0.888	0.140	0.982	0.022	0.882	0.166	0.852	0.125
2019	0.881	0.160	0.894	0.211	0.924	0.127	0.886	0.126
2020	0.883	0.184	NA	NA	0.912	0.085	0.855	0.204
2021	0.924	0.142	0.863	0.181	0.921	0.139	0.836	0.158
2022	0.944	0.100	0.888	0.137	0.961	0.058	0.995	0.002
2023	0.928	0.125	0.890	0.180	0.857	0.103	NA	NA
2024	0.968	0.055	0.896	0.130	0.916	0.062	NA	NA

4 Discussion

To recap the results, for the topic modelling most of the coherence scores are quite good (i.e. close to zero). However, the interpretation of the topics does not always make sense in the context of Ukraine or UK. For the sentiment analysis, the single sample words are partly not intuitively evaluated. However, the frequency of the classification in the context of Ukraine or UK is quite as expected.

Despite the good coherence scores, this U_{mass} score is a small representation of a topic model's performance, so to look at the topics and see if they make sense is still the best way to evaluate its performance. However, some topics seem quite random and not all words are coherent in meaning. Therefore, a better, more extensive data cleaning process is needed to be able to check the modelled topics with the timelines of Ukraine and UK. For example, more irrelevant words (*thanks, for example, as*) could be deleted as stop words. In this research however, there was too little time to extent this process.

Additionally, the use of evaluating a sample of words with a sentiment to check the performance of robBERT v2 is not the best way. The pre-trained version (DTAI-KULeuven, 2022) is trained to classify a complete text, not separate words, which gives the model a very small amount of information to classify its sentiment. A better way to test its performance is to train the version on similar labelled text, other transcripts of parliamentary debates for example, and to test this on a validation set (Delobelle et al., 2020). Unfortunately, in this research there was a lack of labelled transcripts and too little time to label.

Apart from the method, there are also a few limitations in the data. First, this study utilized a small dataset of 256 and 151 debate parts for respectively the Ukraine and UK dataset. These are relatively few documents for especially topic modelling. Second, the imbalance in the Ukraine dataset may have impacted the topic modelling results, suggesting that a more balanced dataset could yield more accurate insights. The UK dataset was more in balance in number of debate parts over time.

Lastly, BERT and similar NLP models operate as "black boxes," meaning their internal workings are not easily interpretable. This lack of transparency makes it difficult to understand how these models generate outputs from given inputs. Additionally, it poses challenges in assessing their performance and understanding their behaviour, which is particularly problematic when traditional metrics are applied without clear insight into the model's decision-making process (in this research the coherence score). Due to this black box characteristic, the performance in this research is difficult to indicate.

All these limitations seem to make this research not of essence. Nevertheless, applying these new models contributes to the evolution of using BERT models in a scientific analysis. Additionally, applying it to a new sort of textual data, parliamentary debates instead of Twitter or review based, creates a new testing field for these models and whether they work as well as other sources. Moreover, the impact on the theoretical base is also of interest since it gains insight in how the perceptions on a region are affected by political upheaval. It is clear to see a difference in the Ukraine analysis between the defining moment, especially the invasion of Russia. For the UK it is more gradual change and might need a more extensive method to be able to see differences between the defining moments.

To conclude, further research on this subject with a revised method is definitely necessary. To be specific, it would be interesting to apply the same analysis on the whole dataset and compare those to the Ukraine and UK filtered datasets. In this analysis, the pre-processing of the data should be more extensive, for example deleting more meaningless words. Also, the sentiment analysis should be done with a trained robBERT v2 on parliamentary debates. After this, the focus could be on different countries with political upheaval to test the generalizability. Lastly, the dataset could be way larger since the parliamentary debate transcripts go further back than 2013. Enlarging the data gives more insight, which was not possible in this research due to limited time.

5 Conclusion

To generally conclude, the modelled topics have a good coherence, however, the topics are not as interpretable as hoped. This interpretability is the best way to test the performance of the topic modelling, so it is key to extend the preprocessing steps in further research. Subsequently, the sentiment analysis evaluates a bit as predicted, but there is no valid way to check the performance of the model. Therefore, it is essential to train the model on similar data to be able to check the performance and to improve the analysis.

For the theoretical background, this research is a first step towards implementing a new method to analyse the perspective of an external party on the political upheaval. Specifically, the topic modelling and sentiment analysis on the Ukraine data is mostly as predicted, so it is clear that a big change in a country is reflected on the data analysis. To answer the research question: *What is the effect of political upheaval on the way of speaking about a country by analysing changes in discussed topics and sentiments through political debates?* There is further research needed with a revised method in which time is too limited to be able to fully answer the question. However, since there is more of an effect to see in the Ukraine dataset than the UK dataset, one could say that the more sudden the political upheaval, the clearer a change in the analysis. Additionally, in the Ukraine defining moments it is a matter of life or death whereas the UK defining moments are more financially focused. This might be of effect in the analysis. Therefore, widening the view and including more countries with political upheaval will be of essence to state these results with more certainty.

References

- Abercrombie, G., & Batista-Navarro, R. T. (2020). Parlvote: A corpus for sentiment analysis of political debates. *Proceedings of the Twelfth Language Resources and Evaluation Conference, LREC 2020*, 5073–5078.
- Agostini, M., & van Zomeren, M. (2021). Toward a comprehensive and potentially cross-cultural model of why people engage in collective action: A quantitative research synthesis of four motivations and structural constraints. *Psychological Bulletin*, 147(7), 667.
- Blei, D. M., & Lafferty, J. D. (2006). Dynamic topic models. *Proceedings of the 23rd International Conference on Machine Learning*, 113–120.
- Chaudhry, H. N., Javed, Y., Kulsoom, F., Mehmood, Z., Khan, Z. I., Shoaib, U., & Janjua, S. H. (2021). Sentiment analysis of before and after elections: Twitter data of us election 2020. *Electronics*, 10(17), 2082. <https://doi.org/https://doi.org/10.3390/electronics10172082>
- Chayinska, M., Minescu, A., & McGarty, C. (2017). Political solidarity through action (and inaction): How international relations changed intracultural perceptions in ukraine. *Group Processes & Intergroup Relations*, 20(3), 396–408.
- Countwordsfree. (n.d.). *The list of stop words*. <https://countwordsfree.com/stopwords/dutch> Accessed May 2024.
- Delobelle, P., Winters, T., & Berendt, B. (2020). RobBERT: A Dutch RoBERTa-based Language Model. *Findings of the Association for Computational Linguistics: EMNLP 2020*, 3255–3265. <https://doi.org/10.18653/v1/2020.findings-emnlp.292>
- DTAI-KULeuven. (2022). RobBERT v2: Dutch Sentiment Analysis Model. RobBERT finetuned for sentiment analysis on DBRD [Online; accessed June 12, 2024]. <https://huggingface.co/DTAI-KULeuven/robbert-v2-dutch-sentiment>
- Gillioz, A., Casas, J., Mugellini, E., & Abou Khaled, O. (2020). Overview of the transformer-based models for nlp tasks. *2020 15th Conference on Computer Science and Information Systems (FedCSIS)*, 179–183.
- Greene, D., & Cross, J. P. (2017). Exploring the political agenda of the european parliament using a dynamic topic modeling approach. *Political Analysis*, 25(1), 77–94.
- Grootendorst, M. (2022). Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint, arXiv:2203.05794*.
- Grootendorst, M. (2024). *BERTopic: Topic modeling with bert* [[Online; accessed June 12, 2024]]. <https://maartengr.github.io/BERTopic/index.html>
- Hadiat, A. R. (2022). Topic modeling evaluations: The relationship between coherency and accuracy. *Doctoral dissertation*.
- Ilyas, S., Soomro, Z., Anwar, A., Shahzad, H., & Yaqub, U. (2020). Analyzing brexit’s impact using sentiment analysis and topic modeling on twitter discussion. *The 21st Annual International Conference*

- on *Digital Government Research (dg.o '20)*, 6. <https://doi.org/https://doi.org/10.1145/3396956.3396973>
- Ivanov, S., Gavrilina, M., Webster, C., & Ralko, V. (2017). Impacts of political instability on the tourism industry in ukraine. *Journal of Policy Research in Tourism, Leisure and Events*, 9(1), 100–127.
- Mimno, D., Wallach, H., Talley, E., Leenders, M., & McCallum, A. (2011). Optimizing semantic coherence in topic models. *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, 262–272.
- Müller-Hansen, F., Callaghan, M. W., Lee, Y. T., Leipprand, A., Flachsland, C., & Minx, J. C. (2021). Who cares about coal? analyzing 70 years of german parliamentary debates on coal with dynamic topic modeling. *Energy Research & Social Science*, 72. <https://doi.org/https://doi.org/10.1016/j.erss.2020.101869>
- Peng, J., Strijker, D., & Wu, Q. (2020). Place identity: How far have we come in exploring its meanings? *Frontiers in Psychology*, 11. <https://doi.org/10.3389/fpsyg.2020.00294>
- Ruz, G. A., Henríquez, P. A., & Mascareño, A. (2020). Sentiment analysis of twitter data during critical events through bayesian networks classifiers. *Future Generation Computer Systems*, 106, 92–104. <https://doi.org/https://doi.org/10.1016/j.future.2020.01.005>
- Silva, N. F. D., Silva, M. C. R., Pereira, F. S., Tarrega, J. P. M., Beinotti, J. V. P., Fonseca, M., & de Carvalho, A. C. D. L. (2021). Evaluating topic models in portuguese political comments about bills from brazil's chamber of deputies. *Intelligent Systems: 10th Brazilian Conference, BRACIS 2021, Virtual Event, November 29–December 3, 2021, Proceedings, Part II*, 10, 104–120.
- Tweede Kamer der Staten-Generaal. (n.d.-a). *Open Data Portaal: OData API*. <https://opendata.tweedekamer.nl/documentatie/odata-api> Accessed April 2024.
- Tweede Kamer der Staten-Generaal. (n.d.-b). *Open Data Portaal: Verslag*. <https://opendata.tweedekamer.nl/documentatie/verslag> Accessed April 2024.

A Appendix: Extra analysis output

Table 9

All topics for Ukraine, with coherence score per topic.

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
-0.389	-0.688	-0.112	-0.083	-0.280
boeren	europese	nederland	nederland	referendum
natuur	nederland	jaar	komen	wetsvoorstel
jaar	president	natuurlijk	zeggen	grondwet
maken	minister	denk	dank	voorstel
komen	europa	landen	gezegd	democratie
dieren	natuurlijk	dank	jaar	zeggen
nederland	commissie	regering	regering	amendement
dank	landen	europese	debat	vragen
landbouw	dank	staat	partijen	nederland
mee	oekraïne	geld	natuurlijk	jaar
Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
-0.193	-0.123	-0.649	-0.209	-0.889
politie	nederland	israël	mensen	gemeenten
vragen	syrië	nederland	opvang	overheid
openbaar	vragen	palestijnse	oekraïne	koning
ministerie	zeggen	zeggen	nederland	minister
zorgen	nederlandse	staat	vluchtelingen	president
nederland	groepen	jaar	asielzoekers	zeggen
bericht	informatie	regering	landen	natuurlijk
veiligheid	gezegd	kernwapens	dank	echt
geval	zoals	nodig	zorgen	democratie
zoals	debat	veiligheid	bijvoorbeeld	koninklijk
Topic 11	Topic 12	Topic 13	Topic 14	Topic 15
-0.386	-0.172	-0.960	-0.230	-5.049
werken	zorg	defensie	dank	verdrag
werk	agema	krijgsmacht	constaterende	europese
arbeidsmarkt	vragen	militairen	vragen	handel
werkgevers	natuurlijk	nederland	verzoekt	nederland
denk	tekorten	samenwerking	economie	landen
groep	patiënten	nederlandse	overwegende	nederlandse
sociale	bijvoorbeeld	dijk	voorgesteld	canadese
zorgen	jaar	oekraïne	beraadslaging	boeren
werknemers	zorgen	boswijk	regering	europa
bijvoorbeeld	zoals	militaire	gesprek	handelsverdrag

Topic 16 <i>-1.164</i>	Topic 17 <i>-2.091</i>	Topic 18 <i>-0.584</i>	Topic 19 <i>-0.457</i>	Topic 20 <i>-0.251</i>
pensioen nieuwe pensioenen pensioenfondsen pensioenstelsel gepensioneerden pensioenfonds amendement pensioenakkoord fondsen	rusland russische oekraïne sancties nederland regering landen staatsaansprakelijkheid nederlandse strafrechtelijk	hof sancties nederlandse nederland verdrag rechten rechter mensenrechten europes europese	gas huishoudens natuurlijk nederland betalen energie prijzen maatregelen geld contracten	veiligheid nationale nederland onderzoek bijvoorbeeld verkeer zorgen bedrijven zaken geval
Topic 21 <i>-0.722</i>	Topic 22 <i>-0.925</i>	Topic 23 <i>-2.209</i>	Topic 24 <i>-0.493</i>	Topic 25 <i>-0.468</i>
zorg testen maatregelen natuurlijk virus echt vragen denk eigenlijk horeca	kinderen nederland veiligheid ouders politie onderwijs journaTopicen school justitie scholen	nederland landen handel ontwikkelingssamenwerking wereld nederlandse vrouwen europa buitenlandse organisaties	president minister vragen dividendbelasting nederland vragenuur debat natuurlijk premier zeggen	marokko vragen nederland politie marokkanen zorg huisartsen aantal zoals onderzoek

Table 10

All topics for the UK dataset, with coherence score per topic.

Topic 1 <i>-0.044</i>	Topic 2 <i>-0.112</i>	Topic 3 <i>-0.630</i>	Topic 4 <i>-0.023</i>	Topic 5 <i>-0.597</i>
nederland jaar europese zeggen natuurlijk dank komt zoals landen graag	syrië nederland nederlandse regering steun israël assad informatie turkije debat	maatregelen natuurlijk testen vragen aantal virus denk besmettingen week eigenlijk	pact nederland natuurlijk zeggen komen denk gezegd zegt zoals eigenlijk	zorg patiënten zorgverzekeraars zorgverzekeraar kwaliteit nodig nederland dank belangrijk verzekeraars
Topic 6 <i>-0.586</i>	Topic 7 <i>-0.613</i>	Topic 8 <i>-0.487</i>	Topic 9 <i>-1.531</i>	Topic 10 <i>-0.594</i>
boeren dieren natuur landbouw nederland maken jaar sector vissers commissie	politie kwetsbaarheden wetsvoorstel seksuele rechter staat seksueel nederland aantal zoals	airport nr luchtvaart luchthavens luchtruim verzoekt regering orde vluchten dag	studenten onderwijs scholen wetsvoorstel instellingen student amendement leraren docenten leerlingen	plastic circulaire stoffen water verzoekt regering orde voorgesteld drinkwater blijkt
Topic 11 <i>-1.531</i>	Topic 12 <i>-4.629</i>	Topic 13 <i>-0.523</i>	Topic 14 <i>-1.385</i>	Topic 15 <i>-1.125</i>
nederland belasting belastingontwijking country belastingdienst belastingverdragen ontwikkelingslanden fiscale nederlandse vliegbelasting	apk rijden verkeersveiligheid kentekenplicht km vvd voertuigen regering automobilist wegen	bedrijven bedenktijd nederland natuurlijk belang publieke zoals overheid aandeelhouders dank	wetsvoorstel beveiligers nederland verdrag landen mensensmokkel particuliere protocol europese rechten	banken financiën sector amendement wetsvoorstel bonussen beloningsbeleid bank toezichthouder bedrijven

Topic 16 <i>-1.625</i>	Topic 17 <i>-1.433</i>	Topic 18 <i>-0.154</i>	Topic 19 <i>-1.787</i>	Topic 20 <i>-0.319</i>
amendement	rusland	vluchtelingen	publieke	turkije
schade	oekraïne	begroting	regionale	turkse
mijnbouwwet	russische	nederland	radio	persvrijheid
gas	landen	turkije	concessiebeleidsplan	nederland
bewijslast	staatsaansprakelijkheid	geld	commissariaat	nederlandse
amendementen	nederland	griekenland	transparantie	journalisten
richtlijn	internationale	miljard	toezicht	regering
veiligheid	buitenlandse	europese	media	vrijheid
banken	strafrechtelijk	regering	omroepverenigingen	europese
gaswinning	nederlandse	financiã	wetsvoorstel	buitenlandse
Topic 21 <i>-0.990</i>	Topic 22 <i>-0.531</i>	Topic 23 <i>-0.206</i>		
pensioen	amendement	europese		
pensioenfondsen	belastingplan	commissie		
pensioenfonds	fiscale	verzoekt		
pensioenen	belasting	orde		
pensioenstelsel	belastingdienst	parlement		
gepensioneerden	miljard	albanie		
fondsen	betalen	regering		
risico	voorzitter	voorgesteld		
nabestaandenpensioen	maatregelen	beraadslaging		
werkgevers	miljoen	turkije		

Table 11

Random sample of words and its sentiment evaluation for the Ukraine dataset.

Year	Word	Sentiment score	Sentiment label
2013	incentivepotje	0.982	Positive
	behaalde	0.973	Positive
	uitgebreidere	0.808	Positive
2014	aanscherpen	0.977	Positive
	uitkomsten	0.960	Negative
	projectie	0.890	Negative
2015	naadloos	0.653	Negative
	thematiek	0.948	Negative
	geschrapt	0.999	Negative
2016	Investigation	0.863	Positive
	daarbij	0.882	Negative
	deel	0.943	Negative
2017	herindelingen	0.971	Positive
	representatieve	0.996	Negative
	werking	0.559	Negative
2018	drie	0.856	Positive
	belemmeren	0.998	Positive
	effecten	0.868	Negative
2019	informatieve	0.999	Positive
	gekkigheid	0.966	Negative
	Aanhangers	0.535	Positive
2020	orgaan	0.905	Positive
	vergaande	0.994	Negative
	opkomt	0.705	Positive
2021	verwoest	0.965	Positive
	opdringen	0.931	Positive
	indringend	0.989	Positive
2022	woorden	0.562	Positive
	gek	0.997	Negative
	komen	0.882	Positive
2023	gevallen	0.944	Negative
	groepsdenken	0.998	Positive
	achteren	0.992	Positive
2024	verwachting	0.527	Negative
	Valstar	0.998	Positive
	vaststelling	0.930	Positive

Table 12

Random sample of words and its sentiment evaluation for the UK dataset.

Year	Word	Sentiment score	Sentiment label
2013	doorrekeningen	0.972	Negative
	vormgeven	0.910	Positive
	onderkennen	0.999	Positive
2014	voorzorgsprincipe	0.832	Positive
	realiseer	0.994	Positive
	hernieuwbaar	0.998	Positive
2015	dienen	0.985	Negative
	verwachting	0.527	Negative
	migrantencrisis	0.587	Negative
2016	daarmee	0.865	Negative
	meenemen	0.964	Positive
	liggende	0.984	Positive
2017	gebaseerd	0.774	Positive
	schrok	0.913	Negative
	recht	0.997	Positive
2018	genueanceerd	0.993	Positive
	laagste	0.951	Positive
	2021-2027	0.962	Positive
2019	verengelt	0.576	Negative
	verstrekken	0.559	Negative
	onthouden	0.987	Negative
2020	dekken	0.966	Positive
	standaardkarikatuur	0.951	Negative
	praat	0.803	Positive
2021	afspraken	0.636	Positive
	schrift	0.881	Negative
	eenheid	0.913	Positive
2022	intens	0.998	Positive
	Zij	0.969	Positive
	landbouwgrond	0.613	Positive
2023	diensten	0.781	Positive
	democratie	0.938	Positive
	verzetten	0.952	Negative
2024	oordeelt	0.870	Negative
	deelname	0.931	Positive
	optreden	0.998	Positive

B Appendix: Programming code

B.1 XML to DataFrame code

```
1 import os
2 from bs4 import BeautifulSoup
3 import pandas as pd
4
5 # Initialize lists to store data
6 texts = []
7 speakers = []
8 dates = []
9 indices = []
10 activities = []
11 titles = []
12
13 file_path = "C:/Users/annet/Documents/Applied Data Science/Thesis/data/verslagen"
14 file_names = os.listdir(file_path)
15
16 for file_name in file_names:
17     with open(os.path.join(file_path, file_name), 'r', encoding='utf-8') as text:
18         data = text.read()
19
20     # Parse the XML
21     soup = BeautifulSoup(data, 'xml')
22     elements = soup.find_all('tekst')
23     date = file_name.split('_')[0]
24
25     for index, element in enumerate(elements):
26         alineaitems = element.find_all("alineaitem")
27
28         if alineaitems:
29             speaker = alineaitems[0].get_text().strip()
30             text = element.get_text()
31             text_clean = text.replace(speaker, "")
32
33             # Find subject within debate
34             parent_activiteit = None
35             for parent in element.parents:
36                 if parent.name == 'activiteit':
37                     parent_activiteit = parent
38                     break
39
40             if parent_activiteit:
41                 parent = parent_activiteit['soort']
42                 title = parent_activiteit.find('titel').text
```

```

43     else:
44         parent = None
45         title = None
46
47     # Append data to lists
48     texts.append(text_clean)
49     speakers.append(speaker)
50     dates.append(date)
51     indices.append(index)
52     activities.append(parent)
53     titles.append(title)
54     else:
55         print("No aline items found in file:", file_name)
56     print("File " + file_name + " done!")
57
58 # Create DataFrame from lists
59 df = pd.DataFrame({
60     'text': texts,
61     'speaker': speakers,
62     'date': dates,
63     'activity': activities,
64     'title': titles,
65     'speaking turn': indices,
66     'party': None,
67     'cabinet': None
68 })
69
70 # Cleaning
71 # Extract party of speaker
72 df["party"] = df["speaker"].str.extract(r'\(([^\)]+)\):')
73 # Extract whether the cabinet is speaking ((prime)minister/secretary of state)
74 df["cabinet"] = df["speaker"].str.match(r'^(Minister|Staatssecretaris) (\s+(\w+))*:$').
75     astype(bool)
76
77 # Delete all empty text and title rows
78 df.loc[:, ['text', 'title']] = df.loc[:, ['text', 'title']].fillna("")
79 df_filtered = df[df['text'].str.contains(r'[a-zA-Z0-9]') & df['title'].str.contains(r'[a-
80     zA-Z0-9]')]
81
82 # Save DataFrame to CSV
83 df_filtered.to_csv("C:/Users/annet/Documents/Applied Data Science/Thesis/data/verslagen.
84     csv", sep=';', encoding='utf-8', index=False)
85
86 # Oekraine filtering
87 # Filter rows where the text contains "Oekra"
88 df_oekra = df_filtered[df_filtered['text'].str.contains("Oekra")]
89 # Extract unique dates and titles from filtered rows

```

```

87 dates_with_oeakra = df_oeakra['date'].unique()
88 titles_with_oeakra = df_oeakra['title'].unique()
89 # Filter the original DataFrame based on both dates and titles
90 df_oeakra = df_filtered[(df_filtered['date'].isin(dates_with_oeakra)) & (df_filtered['title
    '].isin(titles_with_oeakra))]
91
92 # Irrelevant activities removed
93 activities_to_remove = ["Opening", "Mededelingen", "Stemmingen", "Afscheid", "
    Hamerstukken", "Regeling van werkzaamheden"]
94 df_oeakra = df_oeakra[~df_oeakra['activity'].isin(activities_to_remove)]
95
96 # Save DataFrame to CSV
97 df_oeakra.to_csv("C:/Users/annet/Documents/Applied Data Science/Thesis/data/
    verslagen_oeakra.csv", sep=';', encoding='utf-8', index=False)
98
99 # VK+Brexit filtering
100 df_vk = df_filtered[df_filtered['text'].str.contains(pat = '(Brexit|Verenigd\sKoninkrijk)
    ', regex = True)]
101 # Extract unique dates from filtered rows
102 dates_with_vkbrexit = df_vk['date'].unique()
103 titles_with_vkbrexit = df_vk['title'].unique()
104 # Filter the original DataFrame based on both dates and titles
105 df_vk = df_filtered[(df_filtered['date'].isin(dates_with_vkbrexit)) & (df_filtered['title
    '].isin(titles_with_vkbrexit))]
106
107 # Irrelevant activities removed
108 df_vk = df_vk[~df_vk['activity'].isin(activities_to_remove)]
109
110 # Save DataFrame to CSV
111 df_vk.to_csv("C:/Users/annet/Documents/Applied Data Science/Thesis/data/
    verslagen_vkbrexit.csv", sep=';', encoding='utf-8', index=False)

```

B.2 Preprocessing code

```

1 import os
2 import pandas as pd
3 from nltk.corpus import stopwords
4 from nltk.tokenize import word_tokenize
5 import nltk
6 nltk.download('stopwords')
7 import plotly.io as pio
8 pio.renderers.default='browser'
9 import spacy
10

```

```

11 path = "C:/Users/annet/Documents/Applied Data Science/Thesis/data/"
12 os.chdir(path)
13
14 df_oe = pd.read_csv('verslagen_vkbrexit.csv', sep = ";")
15 df_vk = pd.read_csv('verslagen_vkbrexit.csv', sep = ";")
16
17 # Create stop words list
18 stop_words = set(stopwords.words('dutch'))
19 unique_items = set()
20 with open('C:/Users/annet/Documents/Applied Data Science/Thesis/stop_words_dutch.txt', 'r
    ') as file:
21     for line in file:
22         unique_items.add(line.strip()) # Strip any whitespace characters and add the line
            to the set
23 extra_stop_words = unique_items.union(stop_words)
24 extra_stop_words.update(["heer", "mevrouw", "staatssecretaris", "minister", "voorzitter",
            "kamer", "kabinet", "motie", "vraag"])
25
26 # Preprocessing function
27 def preprocess(text):
28     def is_valid_word(word):
29         return any(char.isalpha() for char in word) or '\tilde{a}' in word or '\ddot{i}' in
            word or '-' in word
30
31     # Replace "Verenigd Koninkrijk" with a unique placeholder
32     placeholder = "VERENIGDKONINKRIJK"
33     text = text.replace("Verenigd Koninkrijk", placeholder)
34     # Split text into words while preserving special characters
35     tokens = word_tokenize(text)
36     # Filter out words that do not meet the criteria
37     tokens = [word for word in tokens if is_valid_word(word)]
38     # Remove stopwords in a case-insensitive manner
39     tokens = [word for word in tokens if word.lower() not in extra_stop_words]
40     return tokens
41
42 nlp = spacy.load("nl_core_news_sm")
43
44 # Function to process and tag text using SpaCy
45 def tag_text_filt(text_list):
46     tokens = []
47     for text in text_list:
48         doc = nlp(text)
49         tokens.extend([token.text for token in doc if token.pos_ != 'AUX'])
50     return tokens
51
52 # Function to create a string by joining elements of list
53 def to_tokenize(text):

```

```

54     return text.split()
55
56 def filter_and_group_documents(df, name):
57     # Apply preprocessing and tagging
58     df["pros_token"] = df["text"].apply(preprocess)
59     df["pros_token_filt"] = df["pros_token"].apply(tag_text_filt)
60
61     # Filter out empty documents
62     df = df[df['pros_token_filt'].map(len) > 0]
63
64     # Group the tokens per debate part (title - date combination)
65     df_grouped = df.groupby(['date', 'title']).agg(
66         {'pros_token_filt': lambda x: ' '.join([item for sublist in x for item in sublist])
67         }
68     ).reset_index()
69
70     # Make string/text column
71     df_grouped['pros_string_filt'] = df_grouped['pros_token_filt'].apply(to_tokenize)
72
73     if name == "UK":
74         # Count occurrences of "VERENIGDKONINKRIJK" in strings (> 3)
75         df_grouped['vk_count'] = df_grouped['pros_token_filt'].str.count("
76             VERENIGDKONINKRIJK")
77         df_filt = df_grouped[df_grouped['vk_count'] > 3]
78     else:
79         # Count occurrences of "Oekra" in strings (> 3)
80         df_grouped['oek_count'] = df_grouped['pros_token_filt'].str.count("Oekra")
81         df_filt = df_grouped[df_grouped['oek_count'] > 3]
82
83     # Convert 'date' column to datetime format
84     df_filt['date'] = pd.to_datetime(df_filt['date'], format='%Y%m%d')
85     # Extract year from 'date' column
86     df_filt['year'] = df_filt['date'].dt.year
87
88     return df_filt
89
90 df_oek_filt = df_oek.apply(filter_and_group_documents(df_oek, "Ukraine"))
91 df_oek_filt.to_csv("C:/Users/annet/Documents/Applied Data Science/Thesis/data/
92     processed_oek.csv", sep=';', encoding='utf-8', index=False)
93
94 df_vk_filt = df_oek.apply(filter_and_group_documents(df_oek, "UK"))
95 df_vk_filt.to_csv("C:/Users/annet/Documents/Applied Data Science/Thesis/data/processed_vk
96     .csv", sep=';', encoding='utf-8', index=False)

```

B.3 BERTopic (topic modelling) implementation code

```

1 import os
2 import pandas as pd
3 import bertopic
4 import plotly.io as pio
5 pio.renderers.default='browser'
6 import gensim.corpora as corpora
7 from gensim.models.coherencemodel import CoherenceModel
8 from bertopic.representation import MaximalMarginalRelevance
9
10 path = "C:/Users/annet/Documents/Applied Data Science/Thesis/data/"
11 os.chdir(path)
12
13 df_oeK = pd.read_csv('processed_oeK.csv', sep = ";")
14 df_vk = pd.read_csv('processed_vk.csv', sep = ";")
15
16 def topic_modelling(df):
17     representation_model = MaximalMarginalRelevance(diversity=0.1)
18
19     # BERTopic model
20     model = bertopic.BERTopic(language="Dutch", min_topic_size=5, verbose=True, nr_topics=
21         'auto', representation_model=representation_model)
22     topics, probs = model.fit_transform(df['pros_token_filt'])
23     topics_over_time = model.topics_over_time(df['pros_token_filt'], df['year'],
24         global_tuning=False)
25
26     model.visualize_topics_over_time(topics_over_time)
27
28     # Save topics_over_time plot and model
29     fig = model.visualize_topics_over_time(topics_over_time)
30     fig.write_html("topics_VK_plot.html")
31
32     model.save("BERTopic_vk", serialization="pickle")
33
34     # Coherence score
35     # Initialize model and docs
36     docs = df['pros_token_filt']
37
38     # Extract vectorizer and tokenizer from BERTopic
39     vectorizer = model.vectorizer_model
40     tokenizer = vectorizer.build_tokenizer()
41
42     # Extract features for Topic Coherence evaluation
43     words = vectorizer.get_feature_names_out()
44     tokens = [tokenizer(doc) for doc in docs]

```

```
43 dictionary = corpora.Dictionary(tokens)
44 corpus = [dictionary.doc2bow(token) for token in tokens]
45 topic_words = [[words for words, _ in model.get_topic(topic)]
46                 for topic in range(len(set(topics))-1)]
47
48 # Evaluate coherence scores
49 coherence_model = CoherenceModel(topics=topic_words,
50                                 texts=tokens,
51                                 corpus=corpus,
52                                 dictionary=dictionary,
53                                 coherence='u_mass')
54 coherence_score = coherence_model.get_coherence()
55
56 cs_per_topic = coherence_model.get_coherence_per_topic()
57
58 return model, topics_over_time, coherence_score, cs_per_topic
59
60 df_oe_k_tm = df_oe_k.apply(topic_modelling)
61 df_vk_tm = df_vk.apply(topic_modelling)
```


B.4 roBBERT (sentiment analysis) implementation code

```

1 import os
2 import pandas as pd
3 from transformers import AutoTokenizer, pipeline, AutoModelForSequenceClassification
4 import random
5 import ast
6
7 tokenizer = AutoTokenizer.from_pretrained("DTAI-KULeuven/robbert-v2-dutch-sentiment")
8 model = AutoModelForSequenceClassification.from_pretrained("DTAI-KULeuven/robbert-v2-
9     dutch-sentiment")
10 sentiment_pipeline = pipeline("sentiment-analysis", model=model, tokenizer=tokenizer,
11     truncation=True)
12
13 path = "C:/Users/annet/Documents/Applied Data Science/Thesis/data/"
14 os.chdir(path)
15
16 # Function to calculate sentiment per debate part
17 def sentiment_calculation(input_df):
18     senties = []
19     for index, row in input_df.iterrows():
20         sentiment = sentiment_pipeline(row['pros_token_filt'])
21         senties.append({
22             'title': row['title'],
23             'date': row['date'],
24             'year': row['year'],
25             'sentiment_score': sentiment[0].get('score'),
26             'sentiment_label': sentiment[0].get('label')
27         })
28     return pd.DataFrame(senties)
29
30 # Function to analyse sentiment per year
31 def analyse_sentiment(df):
32     result = df.groupby('year').agg(
33         avg_pos_senti_score=('sentiment_score', lambda x: x[df['sentiment_label'] == '
34             Positive'].mean()),
35         std_pos_senti_score=('sentiment_score', lambda x: x[df['sentiment_label'] == '
36             Positive'].std()),
37         avg_neg_senti_score=('sentiment_score', lambda x: x[df['sentiment_label'] == '
38             Negative'].mean()),
39         std_neg_senti_score=('sentiment_score', lambda x: x[df['sentiment_label'] == '
40             Negative'].std()),
41         total_count=('sentiment_score', 'count'),
42         count_neg=('sentiment_label', lambda x: (x == 'Negative').sum()),
43         count_pos=('sentiment_label', lambda x: (x == 'Positive').sum())
44     ).reset_index()

```

```

39     return result
40
41 df_oe_k = pd.read_csv('processed_oe_k.csv', sep = ";")
42 df_oe_k_filt = df_oe_k[df_oe_k['oe_kra_count'] > 3]
43 df_oe_k_result = sentiment_calculation(df_oe_k_filt)
44 analyse_oe_k_result = analyse_sentiment(df_oe_k_result)
45
46 df_vk = pd.read_csv('processed_vk.csv', sep = ";")
47 df_vk_filt = df_vk[df_vk['vk_count'] > 3]
48 df_vk_result = sentiment_calculation(df_vk_filt)
49 analyse_vk_result = analyse_sentiment(df_vk_result)
50
51 # Performance test for the sentiment analysis
52 def performance_test(df):
53     # Create a subset dataframe and extract the words
54     subset = df[['pros_string_filt', 'year']].rename(columns={"pros_string_filt": "string"
55     })
56     subset["words"] = subset["string"].apply(ast.literal_eval)
57
58     # Sample one row per year
59     years = range(2013, 2025)
60     sampled_subset = pd.concat([subset[subset['year'] == year].sample(n=1) for year in
61     years])
62
63     # Sample three unique words per row if possible
64     sampled_subset['sampled_words'] = sampled_subset['words'].apply(lambda words: random.
65     sample(list(set(words)), min(3, len(set(words)))))
66
67     # Collect sentiment analysis results
68     senties = [
69         {
70             'year': row['year'],
71             'word': word,
72             'sentiment_score': sentiment[0]['score'],
73             'sentiment_label': sentiment[0]['label']
74         }
75         for _, row in sampled_subset.iterrows()
76         for word in row['sampled_words']
77         for sentiment in [sentiment_pipeline(word)]
78     ]
79     return pd.DataFrame(senties)
80
81 perform_oe_k = performance_test(df_oe_k_filt)
82 perform_vk = performance_test(df_vk_filt)

```