
MigrationsKB: A KNOWLEDGE BASE OF PUBLIC ATTITUDES TOWARDS MIGRATIONS AND THEIR DRIVING FACTORS

A PREPRINT

Yiyi Chen, Harald Sack, Mehwish Alam

firstname.lastname@fiz-karlsruhe.de

FIZ Karlsruhe – Leibniz Institute for Information Infrastructure, Germany

Karlsruhe Institute of Technology, Institute AIFB, Germany

August 18, 2021

ABSTRACT

With the increasing trend in the topic of migration in Europe, the public is now more engaged in expressing their opinions through various platforms such as Twitter. Understanding the online discourses is therefore essential to capture the public opinion. The goal of this study is the analysis of social media platform to quantify public attitudes towards migrations and the identification of different factors causing these attitudes. The tweets spanning from 2013 to Jul-2021 in the European countries which are hosts to immigrants are collected, pre-processed, and filtered using advanced topic modeling technique. BERT-based entity linking and sentiment analysis, and attention-based hate speech detection are performed to annotate the curated tweets. Moreover, the external databases are used to identify the potential social and economic factors causing negative attitudes of the people about migration. To further promote research in the interdisciplinary fields of social science and computer science, the outcomes are integrated into a Knowledge Base (KB), i.e., MigrationsKB which significantly extends the existing models to take into account the public attitudes towards migrations and the economic indicators. This KB is made public using FAIR principles, which can be queried through SPARQL endpoint¹. Data dumps are made available on Zenodo².

Keywords Sentiment Analysis · Hate Speech Detection · Public Attitude · Knowledge Base · Social Media Analysis · Migration

1 Introduction

Social media has become one of the most widely used and important channel for people to express their opinions about the events happening around the globe. It is one of the most useful sources for measuring the attitudes of the public on important topics such as migration, climate change, green deal, etc. More specifically, migration has become one of the mainstream topics in Europe. Many European Projects have recently been focusing on the similar topic from different perspectives such as PERCEPTIONS³, METICOS⁴, etc. Many efforts have been put into studying the attitudes of public towards migrations from various perspectives (Hainmueller and Hopkins, 2014; Dennison and Dražanova, 2018; Helen Dempster and Hargrave, 2020). This study in particular focuses on analyzing social media platform in order to quantify and study public attitudes towards migrations and identify different factors which could be probable causes of these attitudes. This study utilizes advanced Artificial Intelligence (AI) methods based on knowledge graphs and neural networks such as contextual language models for analyzing these public attitudes on social media platform such as Twitter. This study aims at: (a) providing better understanding of public attitudes towards migrations, (b) explain possible reasons why these attitudes towards migrations are what they are, (c) define a Knowledge Base (KB) called as

¹<https://mgkb.fiz-karlsruhe.de/sparql/>

²<https://bit.ly/2W01Doc>

³<https://project.perceptions.eu/>

⁴<https://meticos-project.eu/>

MigrationsKB built by taking into account the semantics underlying this field of study, (d) publish this resource using FAIR principles (Wilkinson, 2016), i.e., make the resource Findable, Accessible, Interoperable, and Reusable.

Since the study mainly focuses on the analysis of social media platform such as Twitter, many kinds of challenges arise, i.e., millions of tweets in noisy natural language are being posted around the globe about a particular topic each day, which makes it impossible for the humans to process this information, leading to the necessity of automated processing. Many efforts have been conducted for creating KB for integrating twitter data, such as, TweetsKB (Fafalios et al., 2018), which is a KB of Twitter data in general, TweetsCOVID19 (Dimitrov et al., 2020), which uses COVID-19 related tweets from TweetsKB. One of the other analytical tool which is related to the topic of migration is MigrAnalytics (Alam et al., 2020), which analyzes tweets about migrations from TweetsKB.

Instead of using a subset of TweetsKB, MigrationsKB extends it by focusing on the specific aspects regarding the topic of migration and providing advanced deep learning based semantic annotations. It will help in facilitating further research in the field of social sciences, and to provide a viable corpus for further analysis in the domain of migration.

In order to analyze the public attitudes towards migrations in the destination countries in Europe, the geotagged tweets are extracted using migration related keywords. The irrelevant tweets are then filtered by using state-of-the-art neural network based topic modeling. It further utilizes contextualized word embeddings (Liu et al., 2020) and transfer learning for sentiment analysis and hate speech detection. Temporal and geographical dimensions are then explored for measuring the public attitudes towards migrations in a certain period of time in a certain region. Entity linking is applied to identify the entity mentions linked to Wikipedia for enabling easy search over the tweets related to a particular topic. In order to identify the potential social and economic factors driving the migration flows, external databases such as Eurostat⁵ and Statista⁶, are used to analyze the correlation between the public attitudes and the established economic indicators in a specific region in a certain time-period.

In order to enable the reusability of the results of this analysis, the outcome is then integrated into MigrationsKB which is an extension of the RDFS⁷ model as originally defined in TweetsKB. It is extended by defining new classes and entities to cover the Geo information of the tweets, the results of hate speech detection as well as to integrate the information about the Economic Indicators which could be the potential cause of negativity or hatred towards migrations. Finally, the competency questions are defined, the answers to which can be retrieved with the help of SPARQL⁸ queries. The source code has been made public for reproducibility reasons and is available through a GitHub repository⁹. Information related to MigrationsKB is available through the web page¹⁰. The SPARQL endpoint of MigrationsKB is also made publicly available to enable querying. The dump of annotated data is available through Zenodo.

This paper is structured as follows: In Section 2 discusses the related work. In Section 3, the details of the resource are presented. Section 4 provides details of MigrationsKB, the relevant competency questions are discussed. Finally, section 5 concludes the paper and discusses the future work.

2 Related Work

This section discusses studies which combine KBs and Twitter information from various domains in the first part. The second part of this section discusses the studies conducted for assessing the public attitudes towards migrations. The third part discusses the European projects involved in the topic of migration.

2.1 Knowledge Bases based on Twitter Data

Several studies have been conducted which provide a KB containing Tweets from a particular time span for making it more usable by researchers. TweetsKB (Fafalios et al., 2018) is one such KB which contains more than 1.5 billion tweets spanning almost 5 years, including entity and sentiment annotations, and provides a publicly available RDF dataset using established vocabularies to further facilitate different data exploration scenarios, such as entity-centric sentiment analysis and temporal entity analysis, etc. In the event of COVID-19 pandemic, TweetsCOVID19 (Dimitrov et al., 2020), deploying the RDF schema of TweetsKB, provides a knowledge base of COVID-19-related tweets, building

⁵<https://ec.europa.eu/eurostat>

⁶<https://www.statista.com/>

⁷<https://www.w3.org/TR/rdf-schema/>

⁸<https://www.w3.org/TR/rdf-sparql-query/>

⁹<https://github.com/MigrationsKB/MGKB>

¹⁰<https://MigrationsKB.github.io/MGKB/>

on a TweetsKB subset spanning from October 2019 to April 2020. The study applies the same feature extraction and data publishing methods as TweetsKB.

As a step forward in combining KB and Twitter information to the field of analyzing migration related data, MigrAnalytics (Alam et al., 2020) is introduced. It uses TweetsKB as a starting point to select data during the peak migration period from 2016 to 2017. MigrAnalytics analyzes tweets about migrations from TweetsKB including the hashtags and entities from the single seed word "Refugee" and then further combines European migration statistics to correlate with the selected tweets. MigrAnalytics enriches the keywords using WordNet, Wikipedia, and Word Embeddings. However, it uses a very naive algorithm for performing sentiment analysis. Moreover, it does not introduce any sophisticated way to remove the irrelevant tweets. In contrast, the methods used for generating MigrationsKB follow more advanced methods based on neural networks for sentiment analysis as well as hate speech detection and extends the RDFS model with relevant information. Also, due to the low Recall (i.e., 39%) of the entity linking (Dimitrov et al., 2020), the entities extracted in TweetsKB can not guarantee a comprehensive curation of migration related tweets. For the current study, three rounds of crawling are conducted (cf. Section 3.1).

2.2 Public Attitudes Towards Migrations

While the prominence of the topic of migration has risen sharply over the last decade in Europe, many efforts have been invested into analyzing the public attitudes towards migrations from various aspects. For instance, (Hainmueller and Hopkins, 2014) is based on the studies conducted during the last 2 decades explaining public attitudes on immigration policy in North America and Western Europe. The authors investigate the natives' attitudes towards immigration from perspectives of political economy and political psychology.

It is found that, attitudes towards immigration are shaped by sociotropic concerns about its cultural impacts - and to a lesser extent its economic impacts - on the nation as a whole. While in (Dennison and Dražanova, 2018), the authors explore the academic literature and the most up-to-date data across 17 countries on both sides of the Mediterranean. The study summarizes theoretical explanations for attitudes towards immigration including media effects, economic competition, contact and group threat theories, early life socialization effects, and psychological effects. It also concludes that in Europe, attitudes towards immigration are notably stable, rather than becoming more negative. More recently, (Helen Dempster and Hargrave, 2020) emphasizes on the factors of individuals' values and worldview. It states that individual factors (i.e., personality, early life norm acquisition, tertiary education, familial lifestyle and personal worldview) have more stable and strong impact on the person's attitudes towards immigration rather than the influence from politicians and media.

In (Dennison and Dražanova, 2018) and (Helen Dempster and Hargrave, 2020), the survey data is used exclusively, while for (Hainmueller and Hopkins, 2014) a comprehensive assessment of approximately 100 studies, including both survey and field experiment data, is conducted. However, for the current work, the analysis is performed on data from social media with automated approaches.

2.3 Projects on Migrations in Europe

Since migration has become one of the most popular and controversial topics in Europe, many projects have been established to gain perspectives and aid policy decisions regarding migrations. The PERCEPTIONS project aims to identify narratives, images, and perceptions of Europe abroad and to investigate how the discrepancies in different narratives lead to unrealistic expectations, problems, and security threats for both host countries and migrants. Eventually, it provides toolkits using all the above-mentioned information and measures to counteract the social issues. Meanwhile, METICOS project is mainly focused on creating a holistic solution to solve challenges for border management in the European Union.

3 Knowledge Base Construction of Migration Related Tweets

MigrationsKB is an extension over TweetsKB with the specific focus on the topic of migration (as the name depicts). The goal of this KB is two-fold: (i) to provide a semantically annotated, query-able resource about public attitudes on social media towards migrations, (ii) to provide an insight into which factors in terms of economic indicators are the cause of that attitude. In order to achieve these goals, an overall framework for constructing is shown in Figure 1. The first step is to define migration related keywords and perform keyword based extraction of geotagged tweets. The meta-data of the tweets is then extracted. Furthermore, topic modelling is performed for refining the tweets in case if irrelevant tweets are crawled in the tweet extraction phase. Contextual Embeddings are then used for performing sentiment analysis (i.e., tweets are classified into positive, negative and neutral) on the relevant set of tweets obtained after topic modeling. In order to further analyze the negative sentiments in terms of hate speech

Figure 1: Overall Framework for Constructing MigrationsKB.

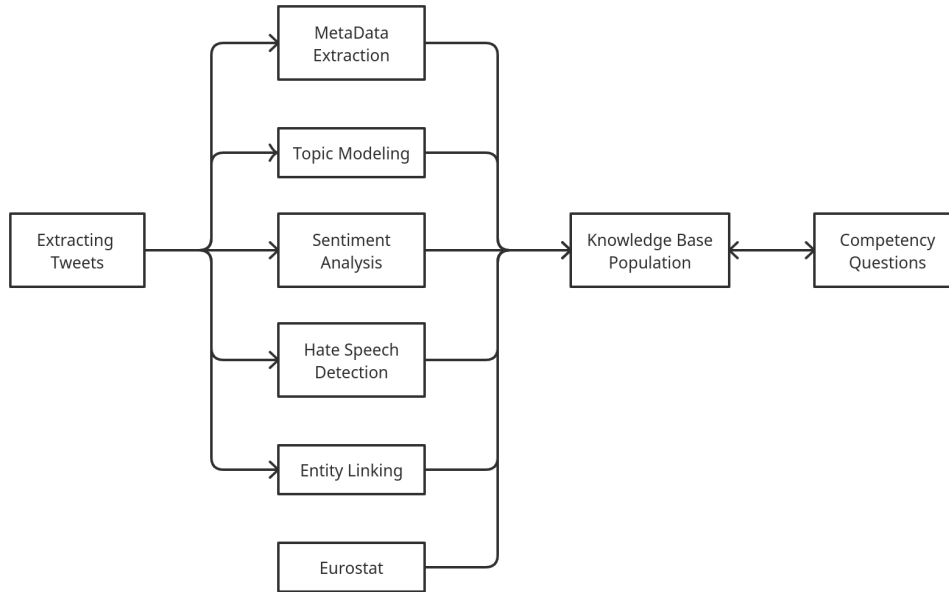


Table 1: Statistics of the EU countries with the most first time asylum applicants.

	2013	2014	2015	2016	2017	2018	2019	2020	SUM
Germany	126705	202645	476510	745160	222565	184180	165615	121955	2245335
Spain	4485	5615	14780	15755	36610	54050	117800	88530	337625
Poland	15240	8020	12190	12305	5045	4110	4070	2785	63765
France	66265	64310	76165	84270	99330	137665	151070	93470	772545
Sweden	54270	81185	162450	28795	26330	21560	26255	16225	417070
United Kingdom	30585	32785	40160	39735	34780	38840	46055	36041	298981
Austria	17500	28035	88160	42255	24715	13710	12860	14180	241415
Hungary	18895	42775	177135	29430	3390	670	500	115	272910
Switzerland	21305	23560	39445	27140	18015	15160	14195	10990	169810
Netherlands	13065	24495	44970	20945	18210	24025	25200	15255	186165
Italy	26620	64625	83540	122960	128850	59950	43770	26535	556850

against the immigrants/refugees, tweets are further classified into three classes, i.e., hate, offensive, and normal. In order to achieve the second goal of this study, an analysis of factors causing the negative sentiment or the hatred against immigrants/refugees is performed with the help of plots.

In order to make this information query-able with the help of SPARQL queries, MigrationsKB, a KB containing public attitudes towards migrations along with the factors driving these attitudes, is constructed. MigrationsKB is then populated with information extracted using the previously described framework. The statistics about these relevant factors such as unemployment rate, the gross domestic product growth rate (GDP), etc. are extracted from Eurostat, Statista, UK Parliament¹¹, and Office for National Statistics¹².

3.1 Collecting Migration Related Tweets

In order to identify the public attitudes towards migrations in the EU countries, the first step is to select a list of destination countries, i.e., the countries hosting the immigrants/refugees. In order to do so, the statistics about asylum applications (annual aggregated) present on Eurostat¹³ are used to obtain the countries with higher frequency of asylum applications during the period from 2013 to 2020, as shown in Table 1. The list of countries includes: Germany, Spain, Poland, France, Sweden, United Kingdom, Austria, Hungary, Switzerland, Netherlands, and Italy.

¹¹<https://www.parliament.uk/>

¹²<https://www.ons.gov.uk/>

¹³<https://ec.europa.eu/eurostat/databrowser/view/tps00191/default/table?lang=en>

Table 2: Statistics of Tweets.

	Germany	Spain	Poland	France	Sweden	UK	Austria	Hungary	Switzerland	Netherlands	Italy	SUM
1st round crawling	8,209	5,902	3,069	7,847	2,790	167,240	1,055	623	4,272	3,587	5,402	209,996
2nd round crawling	17,031	14,981	2,986	20,202	4,116	78,074	5,063	2,768	8,169	11,943	23,718	189,051
3rd round crawling	2,551	980	219	1,742	951	29,065	479	81	663	1,079	1,752	39,562
All (Unique)	26,892	21,392	6,187	29,049	7,556	265,448	6,394	3,355	12,062	16,095	30,023	424,453
Preprocessed (Unique)	25,498	20,240	5,764	26,514	7,263	248,580	6,027	3,226	11,658	15,346	27,223	397,423

In the second step, relevant tweets are extracted using keywords related to the topic of immigration and refugees using word embeddings. The words “immigration” and “refugee” are used as the seed words based on which top-50 most similar words are extracted using pretrained Word2Vec model on Google News¹⁴ as well as fastText embeddings¹⁵. These keywords are then manually filtered for relevance. Based on these keywords, a first round of crawling the tweets is performed. Then for the second round, most popular hashtags, i.e., the hashtags occurring in more than 100 crawled tweets are selected. These hashtags are also verified manually for relevance and then are used for crawling tweets in the second round. For both rounds, the crawled tweets span from 2013 to 2020. The third round of crawling is conducted, using both keywords and hashtags, to extract the tweets spanning from January to the end of July in 2021. The set of keywords, selected popular hashtags for filtering tweets are available on the GitHub repository¹⁶. The 10 most frequently occurring hashtags containing “refugee” and “immigrant” are shown in Table 8.

The extracted tweets are further filtered using their geographical information, i.e., only those tweets are selected which are geotagged with previously identified destination countries. About 66% of the crawled tweets have exact coordinates in their Geo metadata, the rest contain place names, such as “Budapest, Hungary”.

The tweets are then pre-processed by expanding contractions, removing the user mentions, reserved words (i.e., RT), emojis, smileys, numeric tokens, URLs, HTML tags. The punctuation marks are also removed. Moreover, the tokens except hashtags are lemmatized. Eventually, the “#” in hashtags are removed, while the tokens in hashtags are reserved. Finally, the stop-words are removed and only the tweets of sentence length greater than 1 are retained. More specifically for topic modeling (cf. Section 3.2), words with document frequency above 70% are removed. Table 2 shows the statistics of the extracted and preprocessed tweets.

3.2 Topic Modeling

In the previous step, the tweets are collected based on keywords, so it is inevitable that a lot of tweets are actually irrelevant to the topic of migration. For example, many tweets including the keyword “migrant” are about the topic “migrant birds”, and the tweets containing the keyword “exile” are about the “Japanese Band Exile”. Due to the large number of tweets collected (i.e., 397,423 preprocessed tweets), it is hard to manually filter out irrelevant tweets. In order to automate this process, topic modeling is performed.

Topic Modeling is used for extracting hidden semantic structures in the textual documents. One of the classical algorithms for topic modeling is Latent Dirichlet Allocation (LDA) (Blei et al., 2003) which represents each topic as the distribution over terms and each document as a mixture of topics. It is a very powerful algorithm, but it fails in case of huge vocabulary. Therefore, for the current study, the most recent topic modeling algorithm, Embedded Topic Model (ETM) (Dieng et al., 2020), is chosen. Similar to LDA, ETM models each document as a mixture of topics and the words are generated such that they belong to the topics (ranked according to their probability). It also relies on topic proportion and topic assignment. Topic proportion is the proportion of words in a document that belong to a topic, which are the main topics in the document. Topic assignment refers to important words in a given topic. In addition to that, ETM makes use of the embedding of each term and represents each topic in that embedding space. In word embeddings, the context of the word is determined by its surrounding words in a vector space but in case of ETM the context is defined based on the topic. The topic’s distribution over terms is proportional to the inner product of the topic’s embedding and each term’s embedding.

In the setting of the current study, the word and the topic embeddings are trained on tweets. First, the word embeddings are generated by training a Word2Vec skip-gram model on all the preprocessed tweets for 20 epochs, with minimal word frequency 2, dimension 300, negative samples 10 and window size 4. For obtaining the optimal training parameters for ETM, its performance is computed on a document completion task (Rosen-Zvi et al., 2012; Wallach et al., 2009). The parameters for which the highest performance is achieved, are selected and consequently ETM model is utilized. In order to obtain optimal parameters, the dataset is split into 85%, 10%, 5% for train, test, and validation set respectively.

¹⁴<https://github.com/mmihaltz/word2vec-GoogleNews-vectors>

¹⁵<https://fasttext.cc/docs/en/crawl-vectors.html>

¹⁶Keywords:<https://bit.ly/3APiKIw>, Hashtags: <https://bit.ly/2W1TMXk>

The size of the vocabulary of the dataset is 22850. To explore the optimal number of topics, the ETM is experimented with 25, 50, 75, and 100 topics. Initialized with the pretrained word embeddings, the ETM is trained on training data, with batch size 1000, Adam optimizer, and ReLU activation function. In order to select the best number of epochs for training ETM, the model is trained repeatedly by selecting 1-200 epochs and evaluated on the task of document completion (as described previously). The model performs the best on 172 number of epochs with 50 topics.

Table 3: The Results of ETM with Different Number of Topics. Bold values show the best results.

Nr. Topics	25	50	75	100
Val PPL	3329	3015	2920	2870
Best Epoch	185	172	176	178
Topic Coherence	0.0744	0.0777	0.0506	0.02
Topic Diversity	0.9696	0.9288	0.9056	0.7832
Topic Quality	0.0721	0.0721	0.0460	0.0157
Classified Nr. Of Topics	25	50	75	87

The metrics topic coherence and topic diversity are used for evaluation (Dieng et al., 2020). Topic coherence provides a quantitative measure of the interpretability of a topic (Mimno et al., 2011), which is the average point-wise mutual information of two words drawn randomly from the same tweet. A coherent topic would display words that are more likely to occur in the same tweet. In turn, the most likely words in a coherent topic should have high mutual information. In contrary, the topic diversity is defined as the percentage of unique words in the top 25 words of all topics. If there are topics that contain high percentage of words that overlap with the words in another topic, i.e., the diversity would be low, then the topics are redundant. If the diversity is close to 1, the topics are diverse. The results for models with different number of topics are shown in Table 3. The model with 100 topics has the lowest topic diversity and topic coherence, and only 87 topics are assigned to the tweets, which indicates redundancy of the topics. Comparing the models with 25, 50 and 75 topics, the model with 50 topics has a comparably the best topic quality and provides a wider range of topics. Therefore, the trained ETM model with 50 topics is used for classifying the topics for the tweets.

Table 4: Example of topic, terms belonging to the topics, an example Tweet, and its maximal probability score regarding the migration-related topics. The topics with * are chosen as migration-related.

Topic	Top Words	Preprocessed Tweet	Max. Probability Score
1*	refugee, seeker, kill, alien, enter	treatment refugee violate human rights dehumanize refugee endanger european value security argue group psychologist open letter	0.7195
2*	great, call, immigration, question, town	peddle lie interwoven thread brexit regional leave voter low exposure immigration easy scare foreigner queue town come assimilate quickly	1.1062
3*	work, refugee, covid, border, woman	yeah let corrupt nhs education system fine cause deport load hard work immigrant	0.8585
4	people, take, uk, health, hope	illegal immigrant get day uk free home cash health education maternity british national take fool katiehopkins	0.9598
5	stop, find, austria, future, country	proven liar self promote cheat allow uncounted unchecked immigration country cause current crisis	0.4782

The tweets are then refined based on topic embeddings. For each topic, the top 20 words (ranked by their probability) are selected as a representation of the topic¹⁷. These words are then manually verified based on their relevance to the topic of migration. The migration-related topics are chosen with the help of the probabilities associated to all the topics. Regarding the chosen topics, the maximal probability score for each tweet is extracted. Figure 2 shows the distribution of the maximal probability scores of the tweets regarding the chosen topics. By manual evaluation, the threshold for reserving the tweets by the probability score is set to 0.45. Figure 3 shows the distribution of the probability scores of tweets after filtering over the threshold. Eventually, out of the original 397423 tweets, 201555 are reserved for further analysis. Moreover, the topic of each tweet is defined by the maximal probability from all the topics. For example, in Table 4, the topics 1,2, and 3 belong to the chosen migration-related topics, while 4 and 5 are not. More specifically, for the tweet “illegal immigrant get day uk free home cash health education maternity british national take fool katiehopkins”, it is classified as topic 4, and has the maximal probability score 0.9598, which is over the threshold 0.45 and is reserved for the MigrationsKB. As shown in Figure 4 and Figure 5, the tweets are more evenly distributed before filtering, while there are more tweets with migration-related topics proportionally after filtering.

Figure 6 (left) shows the distribution of all the crawled tweets from destination countries from 2013 to Jul-2021. Most of the tweets are from the year 2019 geotagged with Germany, France, Netherlands, Italy, and Spain. Figure 6 (right)

¹⁷<https://bit.ly/2SZgpKb>

Figure 2: The distribution of maximal probability scores of tweets regarding migration-related topics before filtering.

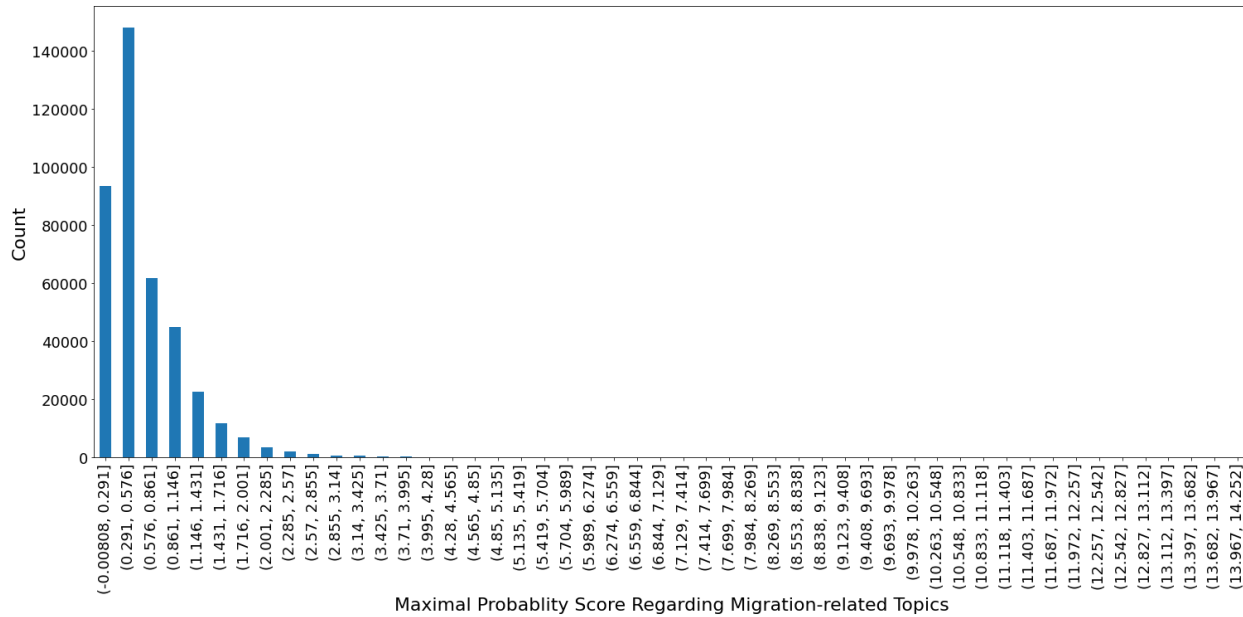
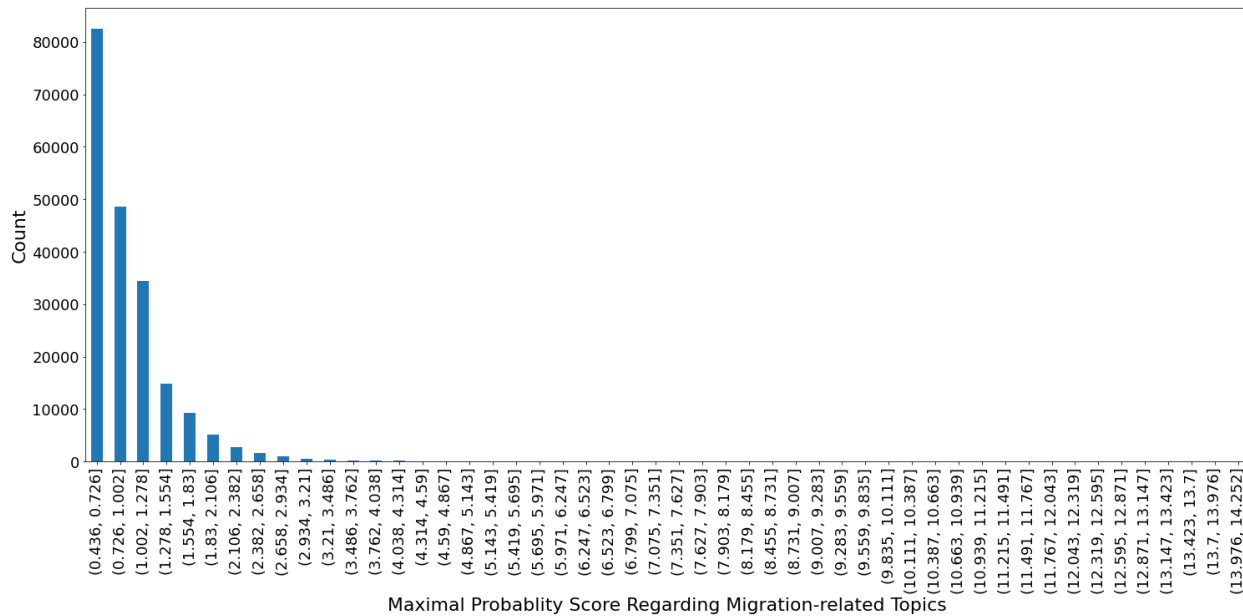


Figure 3: The distribution of maximal probability scores of tweets regarding migration-related topics after filtering.



shows the distribution of all the crawled tweets with the geotag UK within the time frame from 2013 to Jul-2021. Most of the tweets are from the year 2019 and 2020. UK is chosen because currently the focus of this study is English language and the majority of tweets are from there. Figure 7 shows the distribution of the filtered tweets after using ETM. For all the countries, the graph shows similar proportions/trends as Figure 6, but the number of tweets is obviously lower.

Figure 4: The distribution of topics before filtering.

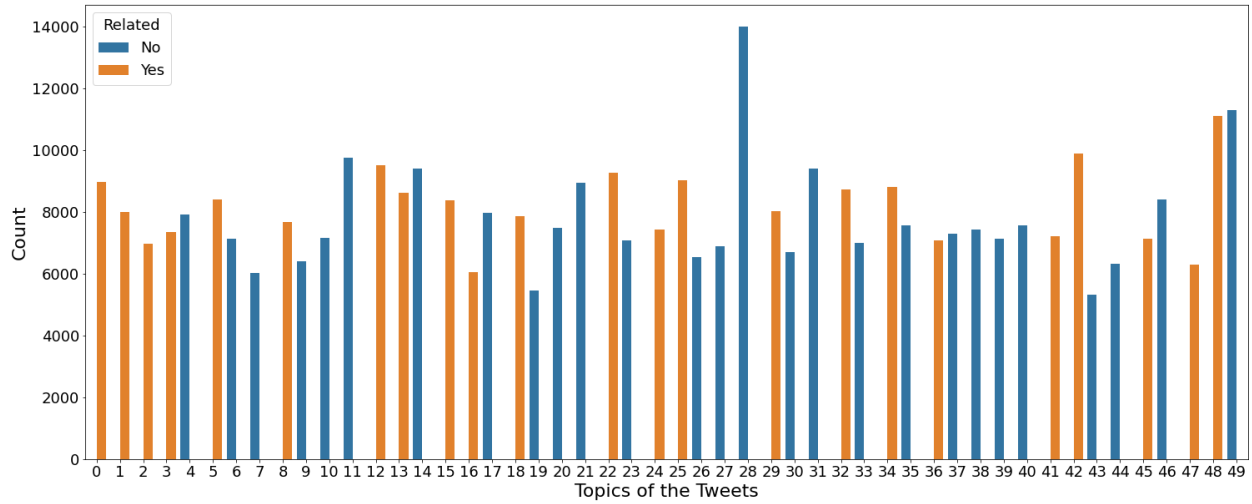
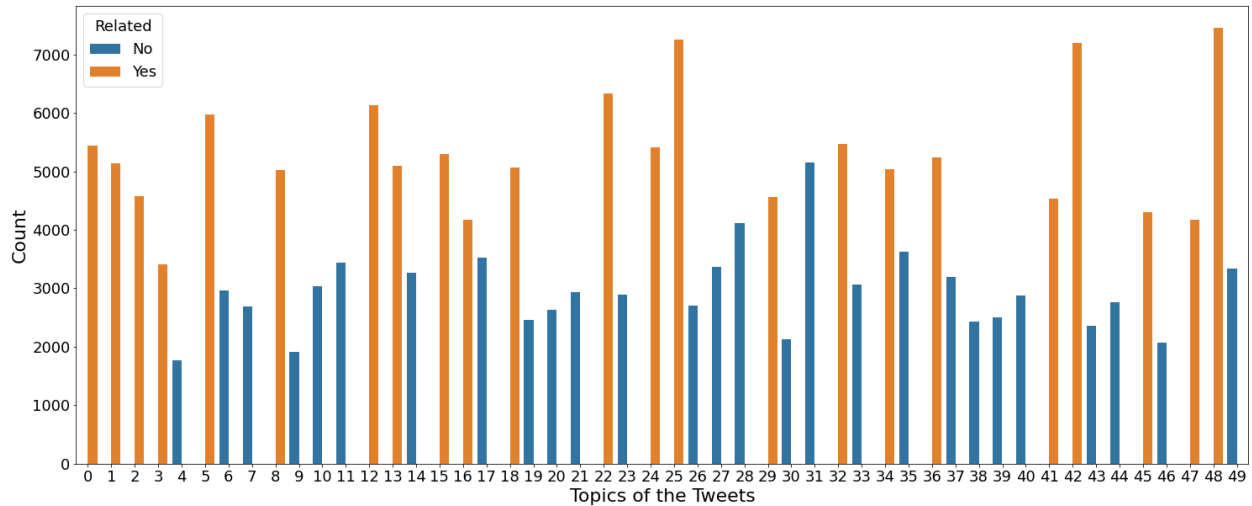


Figure 5: The distribution of topics after filtering.



3.3 Sentiment Analysis

In order to measure the public attitudes towards migrations, sentiment analysis is performed by classifying the tweets into positive, negative, and neutral sentiments. These public sentiments in the destination countries are then analyzed based on their geographic location and temporal information.

Training Dataset Selection. Since, there is a lack of datasets available for sentiment analysis particular to the domain of migration, the existing twitter datasets for sentiment analysis are used for fine-tuning language models for transfer learning on the collected data. Two twitter datasets for sentiment analysis are used most frequently, i.e., the Airline dataset¹⁸ and the SemEval2017 dataset¹⁹(Rosenthal et al., 2017). The Airline dataset focuses on travelers’ opinions on Twitter, which is domain specific. In comparison, the SemEval2017 dataset consists of broader topics of tweets including a range of named entities (e.g., Donald Trump, iPhone), geopolitical entities (e.g., Aleppo, Palestine), and other entities (e.g., Syrian refugees, gun control, and vegetarianism). The dataset is manually annotated using CrowdFlower. The language models are fine-tuned on both the datasets separately and on the combination. Table 5 shows the statistics of the datasets.

¹⁸<https://bit.ly/3u49hZT>

¹⁹<https://bit.ly/3v2jU0d>

Figure 6: Distribution of all the Crawled Tweets based on geographic location.

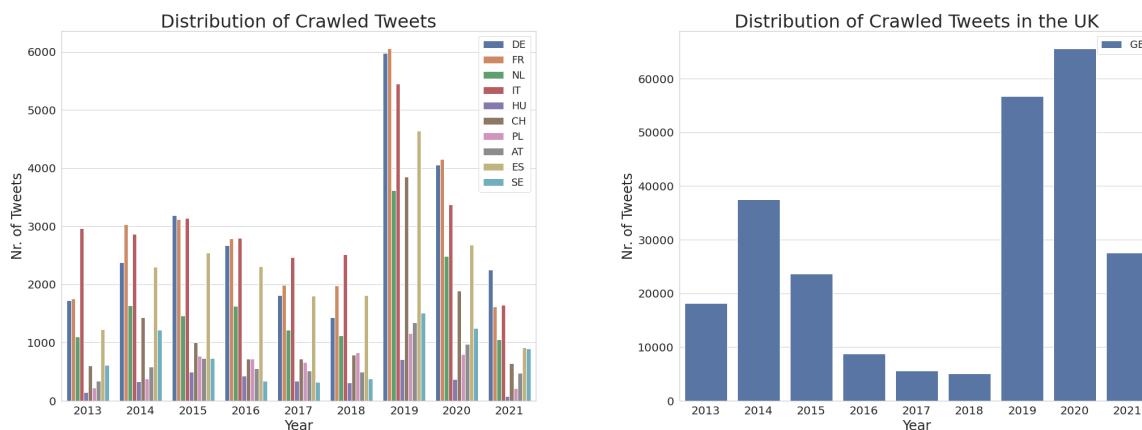
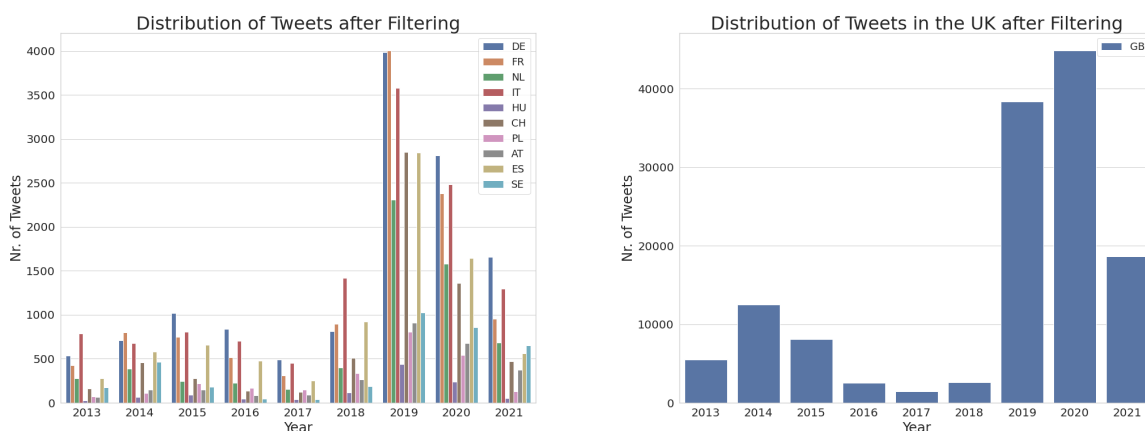


Figure 7: Distribution of Tweets based on geographic location after filtering using ETM.



Model Selection for Sentiment Analysis. For transfer learning, three contextual embedding models are chosen, i.e., BERT (Devlin et al., 2019), XLNet (Yang et al., 2019), and ULMFit (Howard and Ruder, 2018). These three models are fine-tuned as mentioned earlier. The fine-tuned models are then tested on the test set of their corresponding datasets, as well as on the test set of the other datasets. The performance of each of these models is measured. Because the curated tweets obtained from previous steps are not domain specific, the fine-tuned language model is required to perform well on a non-domain specific dataset. Therefore, all the models are evaluated on the test set of SemEval2017 dataset. The results of all the models are shown in Table 6. As shown in the Table 5, there are more neutral and positive tweets than negative ones in SemEval2017 which leads to class imbalance. The macro metrics are more robust to class imbalance and reflect the real performance classifying the minority classes compared to micro metrics, hence the macro F1 score, macro precision, macro recall and also standard accuracy are reported. Since, BERT fine-tuned on SemEval2017 training dataset renders the best results, which is also the state-of-the-art on this dataset (Rosenthal et al., 2017), it is chosen for transfer learning on the collected data for sentiment analysis.

Analyzing Public Sentiments Towards Migrations. In order to identify the public attitudes towards migrations, the sentiments of the tweets for each country in each year are aggregated. In Figure 8 (left), the public sentiments in the UK from 2013 to Jul-2021 are shown. It can be observed that, the total number of tweets about migration from 2013 to 2014, and the tweets with both positive and negative sentiments are increasing in the similar proportion. From 2016 to 2018, the topic of migration is less popular, while there is again a sharp increase in the tweets about migration from 2018 to 2020. Overall, the negative sentiment is much more significant towards migrations compared to the positive sentiment.

Table 5: Statistics of the Existing Datasets for Sentiment Analysis.

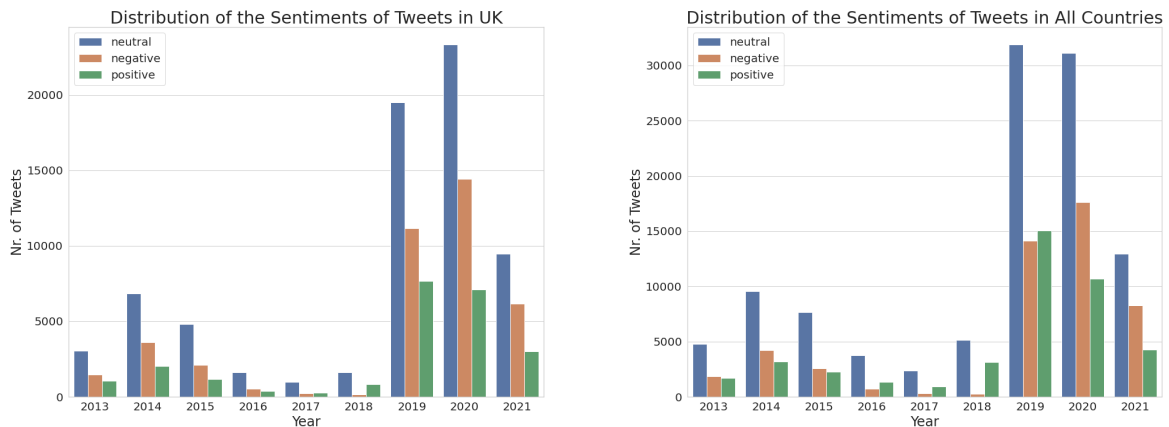
	Train	Validation	Test	All
SemEval2017				
negative	6291	752	766	7809
neutral	17981	2256	2287	22524
positive	15833	2006	1960	19799
Airline				
negative	7316	923	939	9178
neutral	2475	316	308	3099
positive	1921	225	217	2363
Combined				
negative	13608	1736	1643	16987
neutral	20516	2535	2572	25623
positive	17693	2207	2262	22162

Table 6: The Results of Contextual Embedding Models on SemEval2017 test dataset. Bold values show the best results.

Model	Fine-tuned	Acc	F1	Prec	Rec
XLNet	SemEval2017	0.7066	0.6851	0.6988	0.6719
XLNet	Airline	0.5565	0.5987	0.5965	0.601
XLNet	-	0.3718	0.3482	0.3348	0.3627
BERT	SemEval2017	0.7068	0.6949	0.7007	0.6892
ULMFiT	SemEval2017	0.6624	0.6365	0.6342	0.6388
ULMFiT	Airline	0.4709	0.5215	0.52	0.5231
BERT	Airline	0.5117	0.5831	0.5736	0.5929
BERT	-	0.5417	0.5722	0.5753	0.5692
BERT	Combined	0.6691	0.6627	0.6484	0.6776

As shown in Figure 8 (right), the public sentiment towards migrations in all 11 destination European countries follow similar trends as in the UK from 2013 to Jul-2021.

Figure 8: Temporal Distribution of the Sentiments of the Public towards Migrations. The left image shows the sentiments of the people towards migrations in the United Kingdom, and the right image shows the sentiments for all 11 destination countries in Europe.



3.4 Hate Speech Detection

To measure the negative attitude of the public towards migrations, hate speech detection is performed. The tweets are classified into one of the three classes hate, offensive, and normal. In order to perform transfer learning in this scenario,

all the hate speech detection models are trained on recently published manually annotated data about Hate Speech Detection, called as HateXplain (Mathew et al., 2021). Similar to previous studies on Hate Speech Detection, the sources of the dataset are Twitter (Waseem and Hovy, 2016; Davidson et al., 2017; Founta et al., 2018) and Gab (Mathew et al., 2019). For HateXplain Twitter dataset, the tweets are filtered from the 1% randomly collected tweets from Jan-2019 to Jun-2020 using the lexicons combined from (Davidson et al., 2017), (Ousidhoum et al., 2019), and (Mathew et al., 2018). The Gab dataset is originally introduced in (Mathew et al., 2018). All the data is annotated using Amazon Mechanical Turk (MTurk) where each text is annotated based on: (1) whether it is hate speech, offensive speech or normal; (2) the target communities in the text, including target groups such as Race, Religion, Gender, Sexual Orientation, and Miscellaneous; (3) if the text is considered as hate speech or offensive speech by majority, the annotators further annotate which parts of the text provide rationale for the given annotation (this ensures the explainability of manual annotation by the annotators).

HateXplain is split into train, validation, and test dataset by 80%, 10%, and 10%, for which the stratified split is performed to maintain class balance. BiRNN (Schuster and Paliwal, 1997) and BiRNN-Attention (Liu and Lane, 2016) are widely used for text classification task, and CNN-GRU (Zhang et al., 2018) is used for hate speech detection. In the current study, the experimentation is conducted using a combination of various models from CNN, BiLSTM, GRU, BiGRU, and an attention layer for selecting the best model for hate speech detection. For all the models, pre-trained on GloVe embeddings (Pennington et al., 2014) are used as reported in (Mathew et al., 2021). A dropout layer with dropout rate 0.3 is applied after the word embedding layer. For CNN models the convolutional layer has a filter size 100, and window sizes 2, 3, 4. The RNN models use hidden size 100. Finally, the softmax function is used for classifying the texts. The models with the highest accuracy on validation dataset (after training on the training dataset) are chosen for test on the test dataset, whose results are reported in Table 7. Eventually, the best performing pretrained model CNN+BiLSTM+Attention, which is comparable to the results of the best performing model from (Mathew et al., 2021), is used for transfer learning on the collected tweets.

Table 7: The statistics of the HateXplain Dataset and the results of different Hate Speech Detection models.

(a) The Statistics of HateXplain Dataset.

Dataset	Normal	Offensive	Hateful
Train	6251	4384	4748
Validation	781	548	593
Test	782	548	594

(b) The Results of Hate Speech Models on HateXplain. Bold values show the best results.

Model	Acc	F1	Prec	Rec
BiGRU	0.6533	0.6353	0.6343	0.6364
BiGRU+Attn	0.6445	0.6344	0.6297	0.6392
BiLSTM	0.6284	0.6211	0.6169	0.6253
BiLSTM+Attn	0.6512	0.6421	0.6386	0.6457
CNN+GRU	0.6544	0.6545	0.6541	0.6549
CNN+GRU+Attn	0.6450	0.6330	0.6372	0.6436
CNN+BiGRU	0.6575	0.6489	0.6461	0.6517
CNN+BiGRU+Attn	0.6606	0.6472	0.6444	0.6501
CNN+BiLSTM	0.6372	0.6496	0.6523	0.647
CNN+BiLSTM+Attn	0.6863	0.6751	0.6782	0.672

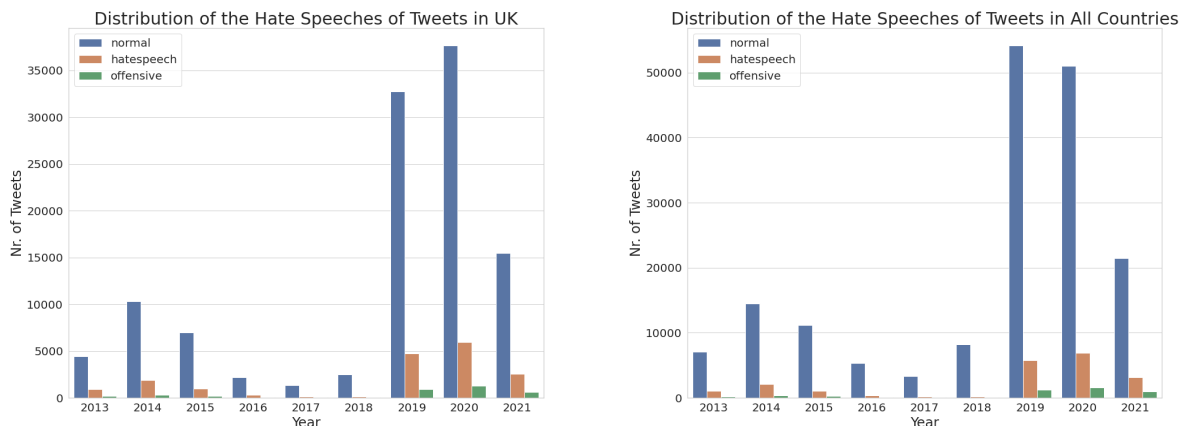
Analyzing Hate Speech in Public Opinions about Migration. The results for hate speech detection are aggregated temporally and geographically to identify the negative attitude of the public towards migrations. As shown in Figure 9, the number of hateful tweets is increasing from 2013 to 2014, decreasing from 2016 to 2018. It is then increasing again from 2019 to 2020 both in the UK and overall in 11 destination countries. In general, the proportion of offensive tweets and hateful tweets are always less than the tweets belonging to the normal class. In summary, the percentage of tweets classified as hate speech in the UK over 2013 to 2020 amounts to 12.98%, while in 11 destination countries it is about 9.36%.

3.5 Entity Linking

For entity linking BLINK (Wu et al., 2020) is used, which utilizes Wikipedia²⁰ as the target KB. Based on fine-tuned BERT, BLINK uses a two-stage approach. In the first stage, BLINK retrieves the candidates in a dense space defined by a bi-encoder that independently embeds the context of entity mention and the entity descriptions. Then in the second stage, each candidate is examined with a cross-encoder, that concatenates the entity mention and entity text.

²⁰The 2019/08/01 Wikipedia dump, which is downloadable in its raw format from <http://dl.fbaipublicfiles.com/BLINK/enwiki-pages-articles.xml.bz2>

Figure 9: Temporal distribution of tweets after hate speech detection. The left image shows the distribution of tweets from UK while the right image is for all the 11 EU countries.



BLINK outperforms state-of-the-art methods on several zero-shot benchmarks and also on established non-zero-shot evaluations such as TACKBP-2010²¹. Out of 201555 tweets in MigrationsKB, for 145747 tweets there is at least one entity mention detected using BLINK. For one tweet, the maximum number of detected entity mentions is 30. In total 89076 unique entities are detected. The detected entities are available online²². The 10 most frequently detected entities containing “refugee” is shown in Table 9.

3.6 Factors Effecting the Public Attitudes Towards Migrations

In order to learn the potential cause of the negative public attitudes towards migrations, the factors such as unemployment rate including youth unemployment rate and total unemployment rate, and GDP are studied. These factors are identified by the experts (Dennison and Drazanova, 2018) as the potential cause of negative attitude towards migrations. This data is collected from Eurostat, Statista, UK Parliament, and Office for National Statistics. Figure 10 shows the comparison between the factors (such as youth employment rate, total employment rate, and real GDP) and the negative attitudes (i.e., negative sentiment and hate speeches) in all the extracted tweets. On average in all 11 destination countries (see Figure 10) and individually in the UK (see Figure 11), the percentages of hate speech and negative sentiment of the public attitudes towards immigration are negatively correlated with the real GDP and positively correlated with total/youth unemployment rate, from 2013 to 2018 and from 2019 to 2020. In 2019, the percentages of hateful tweets and negative tweets are rapidly increased by more than 2% and 1% respectively as compared to 2018. The analysis for each of 11 countries are posted online²³.

4 MigrationsKB

This section discusses the extensions in the RDF/S model of TweetsKB for incorporating public attitudes towards migrations, as well as the economic indicators which drive these attitudes. Moreover, it also discusses some scenarios and competency questions (also translated to SPARQL) that could be asked by the experts to the MigrationsKB.

4.1 RDF/S Model of MigrationsKB

Figure 12 shows the RDF/S model of MigrationsKB. In order to fulfil the purpose of this study, several classes from already existing ontologies are re-used. A complete documentation on of this RDF/S data model is available online²⁴.

For incorporating the meta-data about the geographical location, following information is modified in the TweetsKB. The class `schema:Place` represents Geo information of a tweet. `schema:location` is used for associating a tweet

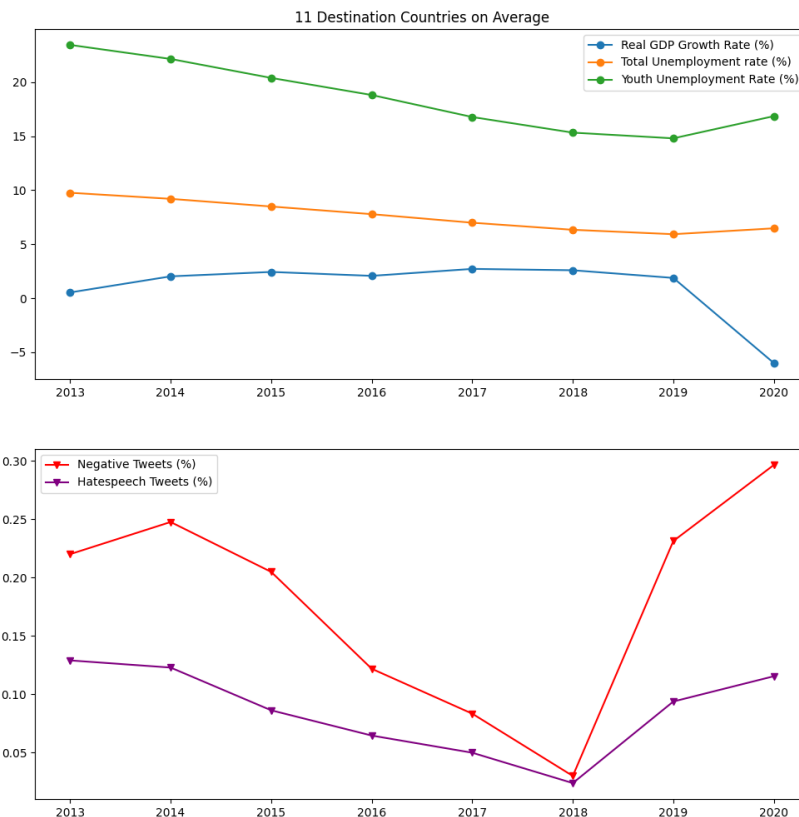
²¹<https://catalog.ldc.upenn.edu/LDC2018T16>

²²<https://bit.ly/3yPpa9D>

²³<https://migrationskb.github.io/MGKB/stats>

²⁴<https://bit.ly/3oERRSL>

Figure 10: The trend of Hate Speech against immigrants/refugees in all the identified destination countries from 2013 to Jul-2021.



(represented as `sioc:Post`) with schema:Place, i.e., its Geo information. `sioc:name` from SIOC²⁵ associates a place with its name represented as a text literal. `schema:addressCountry` specifies the country code of the geographic location of the tweet. `schema:latitude` and `schema:longitude` specify the coordinates of the Geo information.

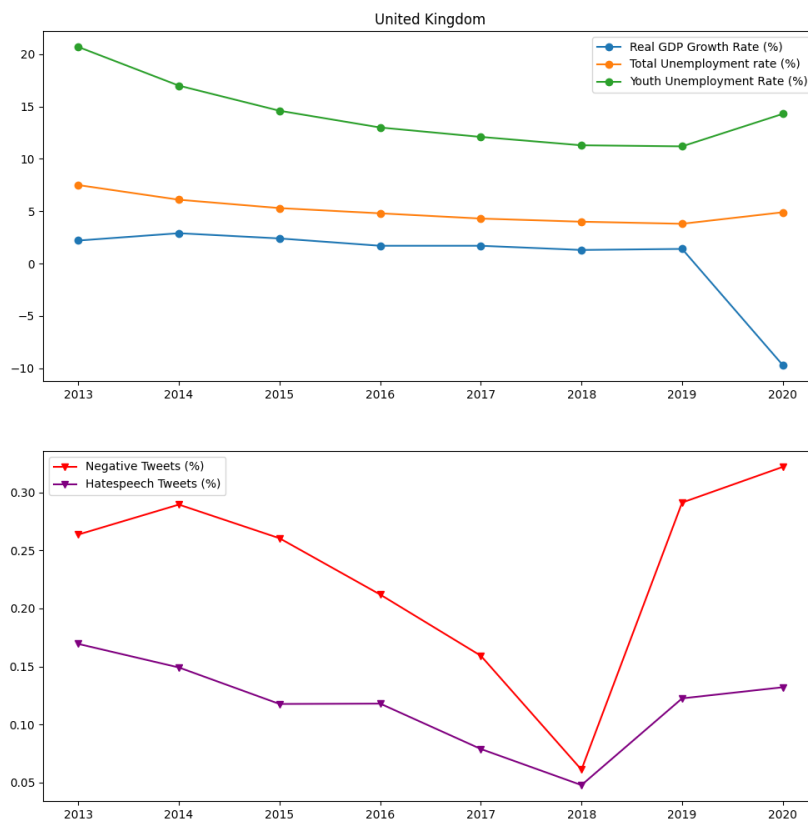
Representing Economic Indicators of European Union. To represent the economic indicators of the destination countries obtained from Eurostat, Financial Industry Business Ontology (FIBO)²⁶ is used.

- The class `fibo-ind-ei-ei:GrossDomesticProduct` represents the GDP of the country of the tweet in a certain year, which are specified by the properties `schema:addressCountry` and `dc:date` from DCMI, and the value of this indicator is represented by `fibo-ind-ei-ei:hasIndicatorValue`.
- The class `fibo-ind-ei-ei:UnemploymentRate` represents the unemployment rate in the country of the tweet in a certain year, represented with the help of the same properties, i.e., `schema:addressCountry`, `dc:date`, and `fibo-ind-ei-ei:hasIndicatorValue`.
- The class `fibo-ind-ei-ei:UnemployedPopulation` is used to specify the population of the unemployment rate.
- The class `fibo-fnd-dt-fd:ExplicitDate` represents the date when the statistics are last updated as a literal with the help of the property `fibo-fnd-dt-fd:hasDateValue`.

²⁵<http://rdfs.org/sioc/spec/>

²⁶<https://spec.edmcouncil.org/fibo/ontology>

Figure 11: The trend of hate speech against immigrants/refugees in the United Kingdom from 2013 to Jul-2021.



- The property `fibonacci-fnd-rel-rel:isCharacterizedBy` is used to associate a tweet with the economic indicators.

Representing Provenance Information. To represent the provenance information about the economic indicators, i.e., Eurostat, Statista, UK parliament, and Office of National Statistics, `PROV-O`²⁷ is used. The class `prov:Activity` defines an activity that occurs over a period of time and acts upon entities, which are defined by the class `prov:Entity`. The class `fibonacci-fnd-arr-asmt:AssessmentActivity` represents an assessment activity involving the evaluation and estimation of the economic indicators, which is a subclass of the class `prov:Activity`. The class `prov:Organization` represents a governmental organization or a company that is associated with the assessment activity, which is a subclass of the class `prov:Agent`. Further extensions are as follows:

- `dc:subject` represents a topic of a tweet resulting from topic modeling (see Section 3.2).
- `wana:neutral-emotion` represents the neutral sentiment of the tweet by applying sentiment analysis.
- `wana:hate`, `mgkb:offensive` and `mgkb:normal` represent the hate speeches, offensive speeches and normal speeches from hate speech detection of the tweets.
- `schema:ReplyAction` represents the action of reply regarding a tweet.
- `mgkb:EconomicIndicators` represents the economic indicators, which has the subclasses `fibonacci-ind-ei-ei:GrossDomesticProduct` and `fibonacci-ind-ei-ei:UnemploymentRate`.

²⁷<https://www.w3.org/TR/prov-o/>

- `mgkb:YouthUnemploymentRate` and `mgkb:TotalUnemploymentRate` represent the unemployment rates with respect to the population, i.e., the youth unemployment population and the total unemployed population.

4.2 Competency Questions.

MigrationsKB can further be used for answering the competency questions with the help of the SPARQL Queries, which can be used for querying information from MigrationsKB. The following query retrieves the top 10 hashtags which contain “refugee” and “immigrant”. The query result is shown in Table 8.

```
SELECT ?hashtagLabel (count(distinct ?tweet) as ?num) WHERE {
  ?tweet schema:mentions ?hashtag.
  ?hashtag a sioc_t:Tag ; rdfs:label ?hashtagLabel.
  FILTER( regex(?hashtagLabel, "refugee", "i") || lcase(str(?hashtagLabel))="refugee"
  || regex(?hashtagLabel, "immigrant", "i") || lcase(str(?hashtagLabel))="immigrant").
} GROUP BY ?hashtagLabel ORDER BY DESC(?num) LIMIT 10
```

Table 8: The Query Result of retrieving the top 10 hashtags which contain “refugee” and “immigrant”.

Hashtag	Nr. of Tweets
refugees	1354
RefugeesWelcome	1066
Refugees	638
refugee	493
RefugeeForum	372
refugeeswelcome	356
WorldRefugeeDay	333
immigrants	210
RefugeeWeek	183
WithRefugees	148

The following query retrieves the top 10 detected entities containing “refugee”. The query result is shown in Table 9.

```
SELECT ?entityLabel (count(?entityLabel) as ?numOfEntityMentions) where{
  ?tweet schema:mentions ?entity.
  ?entity a nee:Entity; nee:hasMatchedURI ?uri.
  ?uri a rdfs:Resource; rdfs:label ?entityLabel.
  FILTER( regex(?entityLabel, "refugee", "i") || lcase(str(?entityLabel))="refugee").
}GROUP BY ?entityLabel ORDER BY DESC(?numOfEntityMentions) LIMIT 10
```

Table 9: The Query Result of retrieving the top 10 entities which contain “refugee” and “immigrant”.

Entity	Nr. of Tweets
United Nations High Commissioner for Refugees	791
Refugee	183
Refugee Nation	149
Refugee Week	120
Convention Relating to the Status of Refugees	115
World Refugee Day	105
Refugee camp	46
European Refugee Fund	35
Refugee Council	34
Refugee Studies Centre	33

The following query identifies the sentiments and hatred of the people concerning refugees, i.e. search entities defining “refugee”. The query result is shown in Table 10.

```
SELECT ?EmotionCategory (count(distinct ?tweet) as ?numOfTweets) WHERE{
```

```

?tweet schema:mentions ?entity.
?entity a nee:Entity; nee:hasMatchedURI ?uri.
?uri a rdfs:Resource; rdfs:label ?x.
    FILTER( regex(?x, "refugee", "i") || lcase(str(?x))="refugee").
?tweet onyx:hasEmotionSet ?y.
?y a onyx:EmotionSet; onyx:hasEmotion ?z.
?z a onyx:Emotion; onyx:hasEmotionCategory ?EmotionCategory.
} GROUP BY ?EmotionCategory

```

Table 10: The Query Result of identifying the Sentiments and Hate Speech Emotions of the Public by searching entities defining “Refugees”.

Emotion Category	Nr. of Tweets
wna:neutral-emotion	1062
wna:posiive-emotion	714
wna:negative-emotion	253
mgkb:normal	1984
mgkb:offensive	8
wna:hate	37

The following query retrieves the GDPR indicator values and the number of tweets identified as hate speeches in the United Kingdom by year. The least amount of tweets occurs in 2017 and 2018 (shown in Figure 7), there is the least amount of the tweets classified in the “hate” class during hate speech detection. In the years 2019 and 2020, with the sharp decrease in the GDPR, there are many more hateful tweets than the previous years on yearly basis. The query result is shown in Table 11.

```

SELECT ?year ?IndValue (count(?tweet) as ?numOfTweets) where {
    ?tweet fibo_fnd_rel_rel:isCharacterizedBy ?gdpr.

    ?gdpr a fibo_ind_ei_ei:GrossDomesticProduct.
    ?gdpr schema:addressCountry "GB".
    ?gdpr dc:date ?year.
    ?gdpr fibo_ind_ei_ei:hasIndicatorValue ?IndValue.
    ?tweet onyx:hasEmotionSet ?y.
    ?y a onyx:EmotionSet; onyx:hasEmotion ?z.
    ?z a onyx:Emotion; onyx:hasEmotionCategory wna:hate.
}GROUP BY ?year ?IndValue ORDER BY DESC(?year)

```

Table 11: The GDPR and the Number of Hate Speeches in the United Kingdom.

Year	GDPR	Nr. of Tweet
2020	-9.7	5928
2019	1.4	4698
2018	1.3	126
2017	1.7	115
2016	1.7	299
2015	2.4	951
2014	2.9	1865
2013	2.2	934

5 Discussion and Future Work

In the current study, a Knowledge Base of migration related tweets is represented. The tweets are filtered and annotated using BERT-based sentiment analysis and attention-based hate speech detection. MigrationsKB extends the RDF/S model defined by TweetsKB by adding the Geo information of tweets, the statistics of economic indicators of European Union, and the results from hate speech detection algorithm. The corpus would assist further research in various fields such as social science by providing readily available information. With the integrated features, the relations between the public attitudes towards migrations and the economic factors have been made query-able.

As for now, the focus is solely on the English tweets, the distribution of the corpus is therefore highly skewed by the tweets from the United Kingdom. While focusing on the destination countries in Europe, there is already a wide variety of languages that need attention. Secondly, visualization tools and interfaces for querying will be created to help the experts in other fields such as social scientists to effectively interact with MigrationsKB. Finally, the MigrationsKB will be under continuous development by updating with newer relevant tweets. The more advanced approaches for topic modeling will be experimented, to analyze the change of topics in tweets in the temporal dimension.

Acknowledgement

This work is a part of ITFlows project²⁸. This project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 882986.

References

- Jens Hainmueller and Daniel J. Hopkins. Public attitudes toward immigration. *Annual Review of Political Science*, 17(1):225–249, 2014. doi:10.1146/annurev-polisci-102512-194818. URL <https://doi.org/10.1146/annurev-polisci-102512-194818>.
- James Dennison and Lenka Drazanova. Public attitudes on migration : rethinking how people perceive migration : an analysis of existing opinion polls in the euro-mediterranean region, 2018.
- Amy Leach Helen Dempster and Karen Hargrave. Public attitudes towards immigration and immigrants: What people think, why and how to influence them. 2020. URL <https://bit.ly/3xYcn3r>.
- Dumontier M. Aalbersberg I. et al. Wilkinson, M. The fair guiding principles for scientific data management and stewardship. *Scientific Data*, 2016. ISSN 2052-4463. URL <https://doi.org/10.1038/sdata.2016.18>.
- Pavlos Fafalios, Vasileios Iosifidis, Eirini Ntoutsis, and Stefan Dietze. Tweetskb: A public and large-scale RDF corpus of annotated tweets. *CoRR*, abs/1810.10308, 2018. URL <http://arxiv.org/abs/1810.10308>.
- Dimitar Dimitrov, Erdal Baran, Pavlos Fafalios, Ran Yu, Xiaofei Zhu, Matthäus Zloch, and Stefan Dietze. Tweetscov19 - a knowledge base of semantically annotated tweets about the covid-19 pandemic. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management, CIKM '20*, page 2991–2998, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450368599. doi:10.1145/3340531.3412765. URL <https://doi.org/10.1145/3340531.3412765>.
- M. Alam, G. A. Gesese, Z. Rezaie, and H. Sack. Migranalytics: Entity-based analytics of migration tweets. *CEUR workshop proceedings*, 2721:74–78, 2020. ISSN 1613-0073.
- Qi Liu, Matt J. Kusner, and Phil Blunsom. A survey on contextual embeddings. *CoRR*, abs/2003.07278, 2020. URL <https://arxiv.org/abs/2003.07278>.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, 2003. URL <http://jmlr.org/papers/v3/blei03a.html>.
- Adji Bousso Dieng, Francisco J. R. Ruiz, and David M. Blei. Topic modeling in embedding spaces. *Trans. Assoc. Comput. Linguistics*, 8:439–453, 2020. URL <https://transacl.org/ojs/index.php/tacl/article/view/2093>.
- Michal Rosen-Zvi, Thomas L. Griffiths, Mark Steyvers, and Padhraic Smyth. The author-topic model for authors and documents. *CoRR*, abs/1207.4169, 2012. URL <http://arxiv.org/abs/1207.4169>.
- Hanna M. Wallach, Iain Murray, Ruslan Salakhutdinov, and David Mimno. Evaluation methods for topic models. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, page 1105–1112, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605585161. doi:10.1145/1553374.1553515. URL <https://doi.org/10.1145/1553374.1553515>.
- David Mimno, Hanna M. Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. Optimizing semantic coherence in topic models. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, page 262–272, USA, 2011. Association for Computational Linguistics. ISBN 9781937284114.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. SemEval-2017 task 4: Sentiment analysis in Twitter. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 502–518, Vancouver, Canada, August 2017. Association for Computational Linguistics. doi:10.18653/v1/S17-2088. URL <https://www.aclweb.org/anthology/S17-2088>.

²⁸<https://www.itflows.eu/>

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. doi:10.18653/v1/n19-1423. URL <https://doi.org/10.18653/v1/n19-1423>.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 5754–5764, 2019.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 328–339. Association for Computational Linguistics, 2018.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. Hatexplain: A benchmark dataset for explainable hate speech detection. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence*. AAAI Press, 2021.
- Zeerak Waseem and Dirk Hovy. Hateful symbols or hateful people? predictive features for hate speech detection on Twitter. In *Proceedings of the NAACL Student Research Workshop*, pages 88–93, San Diego, California, June 2016. Association for Computational Linguistics. doi:10.18653/v1/N16-2013. URL <https://www.aclweb.org/anthology/N16-2013>.
- Thomas Davidson, Dana Warmsley, Michael W. Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language. *CoRR*, abs/1703.04009, 2017. URL <http://arxiv.org/abs/1703.04009>.
- Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. Large scale crowdsourcing and characterization of twitter abusive behavior. *CoRR*, abs/1802.00393, 2018. URL <http://arxiv.org/abs/1802.00393>.
- Binny Mathew, Anurag Illendula, Punyajoy Saha, Soumya Sarkar, Pawan Goyal, and Animesh Mukherjee. Temporal effects of unmoderated hate speech in gab. *CoRR*, abs/1909.10966, 2019. URL <http://arxiv.org/abs/1909.10966>.
- Nedjma Ousidhoum, Zizheng Lin, Hongming Zhang, Yangqiu Song, and Dit-Yan Yeung. Multilingual and multi-aspect hate speech analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4675–4684, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:10.18653/v1/D19-1474. URL <https://www.aclweb.org/anthology/D19-1474>.
- Binny Mathew, Ritam Dutt, Pawan Goyal, and Animesh Mukherjee. Spread of hate speech in online social media. *CoRR*, abs/1812.01693, 2018. URL <http://arxiv.org/abs/1812.01693>.
- M. Schuster and K.K. Paliwal. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681, 1997. doi:10.1109/78.650093.
- Bing Liu and Ian Lane. Attention-based recurrent neural network models for joint intent detection and slot filling, 2016.
- Ziqi Zhang, David Robinson, and Jonathan Tepper. Detecting hate speech on twitter using a convolution-gru based deep neural network. In Aldo Gangemi, Roberto Navigli, Maria-Esther Vidal, Pascal Hitzler, Raphaël Troncy, Laura Hollink, Anna Tordai, and Mehwish Alam, editors, *The Semantic Web*, pages 745–760, Cham, 2018. Springer International Publishing. ISBN 978-3-319-93417-4.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi:10.3115/v1/D14-1162. URL <https://www.aclweb.org/anthology/D14-1162>.
- Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. Scalable zero-shot entity linking with dense entity retrieval. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 6397–6407. Association for Computational Linguistics, 2020.

