# Exploring Approaches for Measuring Risk in the News

Emanuele **Di Buccio**[1,2,3,*], Federico **Neresini**[2]

[1]*Department of Information Engineering, University of Padova, Italy*
[2]*Department of Philosophy, Sociology, Education and Applied Psychology, University of Padova, Italy*
[3]*Department of Statistical Sciences, University of Padova, Italy*

#### Abstract

This paper presents a preliminary investigation into automatic approaches that rely solely on content-based features to compute indicators such as the "risk indicator", which aims to provide a measure of the extent to which risk is present/evoked in a set of informative resources. We built a dataset comprising English newspaper articles, labeled by ten experts in Social Sciences and Humanities using a three-level scale to denote the 'degree' of risk suggested by each article. The study reports on experimental results obtained by considering different instantiations of the indicator, specifically exploiting normalized term frequencies and the Concept Mover Distance.

#### Keywords

Information retrieval and access, Computational Methods for Social Sciences, Risk

## 1. Introduction

The availability of a vast amount of digitized content has provided novel opportunities for investigating research questions in Humanities and Social Sciences, for instance, for carrying out investigations on heterogeneous and longitudinal sets of informative resources. One such research question is how some public issues are discussed in different informative resources, e.g., different media streams. Let us consider, for instance, issues such as "nuclear power" or "artificial intelligence". How are those issues discussed in the newspapers, social media, or other media streams? Can we measure if the narrative is close to concepts such as risk and if that changes through time? Or if the "amount" of risk is larger when discussing issues in future scenarios? The same questions are also relevant to other *indicators*, e.g., for measuring the extent to which the discourse on a public issue is or becomes conflictual over time. If several indicators are available, we can use them as features to identify, for instance, the extent to which a topic covered a a set of documents is or is becoming controversial.

Because of the large amount of data potentially involved in this kind of investigations, we should devise completely automatic methods for computing such indicators or, at least, methods to support expert users in their investigations and reduce their cognitive effort by providing a subset of large corpora where risk is "more present."

This paper reports a preliminary experimental evaluation of methods based on content-based features to compute indicators such as "risk". We will rely on two existing approaches. The first approach exploits the frequency of terms in a controlled vocabulary to compute the indicator on a per document basis; the approach, proposed in [1] and inspired by the idea proposed in [2], can include information other than term frequency, like the document length and the average document length in the document corpus. The second approach is based on the Concept Mover's Distance, which was proposed [3] and exploits distributed representations of words to measure the extent to which a focal concept, in our case "risk", is present in texts.

The experimental evaluation has been carried out on a labeled dataset of news articles gathered from English online newspapers. Those articles were manually labeled by ten assessors, specifically

researchers in Humanities and Social Sciences, using a 3-level scale. The presence of risk in the articles measured by the two considered approaches is compared with the labels assigned by the users.

The analysis of the obtained results shows a modest but positive relationship between the measures and the manually assigned labels.

## 2. Related Works

Works on measuring the presence of risk(s) or the extent to which documents evoke the concept of risk span a variety of research fields and tasks.

The *eRisk* series of Workshops,[1] which started in 2017, aims at the early detection of risk on the Internet. Detection techniques investigated in eRisk aim at the identification of different situations related to health and safety, which include depression, harassment, signs of anorexia, self-harm, and pathological gambling. The dataset adopted in the evaluation initiative relies on a text collection gathered from Reddit [4]. Another work that focuses on the detection of risks, but in a different setting and with a different goal, is [5]. In that work Latent Dirichlet Allocation (LDA) [6] was adopted to identify and then monitor risks due to beeswax adulteration; the adopted corpus consists in news articles collected through the Medical Information System from the Europe Media Monitor (EMM/MEDISYS) [7, 8]. However, the work reported in this paper considers a different task: our objective is to measure the degree to which a document, particularly a newspaper article, evokes the concept of risk, thus having a "measure" of the extent to which a newspaper reader perceives risk when reading an article.

The work reported in [9] investigates dense representations of words obtained using Word2Vec, to analyze diverse variables that might affect perception of risk. Tasks investigated include predicting risk perception, which is framed as a regression task, investigating word associations with risk sources, and predicting ratings of risk dimensions. Risk sources might include technologies, activities and occupations, and geopolitical forces. One objective of [9] was to gain insights into the opportunities provided by dense vector representations, for instance, on their capability to provide the same results that can be obtained through instruments such as surveys. Indeed, as discussed in that paper, these instruments have some limitations. The results obtained in a survey generally refer to a specific risk source or a set of risk sources; however, the obtained insights do not necessarily generalize to other (novel) risk sources. Moreover, it is not possible to carry out a "retrospective" survey, which might be beneficial to gain insights into the risk perception of a particular source in the past. In [9], the risk sources were presented to a set of participants that needed to rate them in terms of riskiness, using a scale from -100 to 100, and evaluate them on nine dimensions of risk perception; another task was, given a risk dimension, to list the first three words that came into their mind. The obtained results show the effectiveness of Word2Vec-based representation. The work reported in this paper shares with [9] the adoption of dense vector representations but has a different focus: not risk sources but documents. In [9] the idea of relying on newspapers to gain a retrospective perception of risk is mentioned as a possible research direction, but it is not explicitly explored.

Research works on automatic detection and analysis of Media Frames might also be helpful for the task addressed in this paper. A possible definition of framing is that proposed in [10]: "To frame is to select some aspects of a perceived reality and make them more salient in a communicating text, in such a way as to promote a particular problem definition, causal interpretation, moral evaluation, and/or treatment recommendation for the item described". Even if this is not a research direction pursued in this paper, the relationship between framing theory and the task of "measuring" risk and the usefulness of past contributions on automatic techniques to detect and analyze frames [11, 12, 13, 14] is worth investigating.

---

[1]https://erisk.irlab.org/

## 3. Risk Indicator Instantiations

As mentioned in Section 1, we will consider two different instantiations of the risk indicator.

The first instantiation shares the same intuition underlying the simple solution discussed in [3], where the occurrences of the focal concept, in our case "risk", are counted and normalized over the number of words in the document. As stressed in that paper, one of the limitations of this approach is that a positive value for the indicator instantiation is obtained only if the word "risk" occurs in the document. The indicator instantiation proposed in [1] mitigates this issue not using a single keyword but a set of keywords to describe the focal concept. For instance, in the case of risk, the set of keywords can be: "risk", "danger", "hazard", "damage", and "harm". These keywords can be identified with the help of experts, through automatic techniques, or using a hybrid strategy where the expert is supposed to use automatic extraction techniques to identify a set of candidates. Besides this limitation, this instantiation does not consider other information, such as the length of the document. A possible solution to take into consideration this additional document feature is to rely on the pivoted document length normalization [15]; moreover, to mitigate the effect of the number of occurrences on the indicator value, a possible option is to use the "saturated" term frequency as in BM25. Therefore, given a set of keywords $\mathcal{K}$, the value of the indicator for the document $d$ can be computed as:

$$\mathcal{I}_{\mathcal{K}}(d) \quad = \quad \frac{1}{|\mathcal{K}|} \sum_{w \in \mathcal{K}} \frac{n_L(w,d)/B}{n_L(w,d)/B + k_1} \tag{1}$$

where $n_L(w,d)$ is the frequency of the term $w$ in the document $d$; $n_L(w,d)$ is normalized by $B$: $(1-b) + b\frac{\mathrm{dl}(d)}{\mathrm{avgdl}(\mathcal{C})}$ where $\mathrm{dl}(d)$ is the length of the document $d$ and $\mathrm{avgdl}(\mathcal{C})$ is the average document length in the corpus $\mathcal{C}$; $b \in [0,1]$ is a parameter that controls the weight assigned to the document length normalization. The $k_1$'s values control the term frequency's effect on a document's indicator value. Differently from the original proposal [1], we did not count only occurrences for the exact keywords in $\mathcal{K}$, but we looked for occurrences of words containing the keyword as prefix — all the words starting with "risk", which include "riskiness" and "risky". This instantiation of the risk indicator and other variants are currently implemented and integrated in the TIPS Platform [16].

The other instantiation of the risk indicator relies on the Concept Mover's Distance (CMD) [3]. CMD is a way to measure the extent to which a concept, in our case "risk", is present in a document. CMD builds on Word Mover's Distance (WMD), which measures the dissimilarity between two documents as "the minimum distance that the embedded words of one document need to "travel" to reach the embedded words of another document." [17]. In [3], a pseudo-document composed of the term or the terms denoting a specified concept is built, and then the WMDs are computed on the corpus including the pseudo-document; the CMD is inverted in order to get a measure of "closeness". A detailed description of WMD and CMD can be found respectively in [17] and [3]. One of the advantages of the second indicator instantiation is that WMD, and consequently CMD, can rely on different methods to compute word embeddings, thus providing flexibility in terms of word representation.

## 4. Evaluation

We performed an experimental evaluation to gain some preliminary insights into the capability of the two indicator instantiations in measuring the engagement of the concept of risk in a document. Section 4.1 will reports details on how the two approaches were used, e.g., specific values of the parameters and of the vocabularies. Section 4.2 will describe the dataset that we built to carry out the evaluation. Section 4.3 will report and discuss the adopted evaluation measures and the obtained results.

### 4.1. Experimental Settings

As for the first approach, we considered two possible instantiations of the risk indicator using Equation 1. In the first, the indicator assign to each word in $\mathcal{K}$ the term frequency normalized over the document

length, in line with the simple approach discussed in [3]. In the second instantiation, we relied on Eq.1 and we set $b = 0.75$ and $k_1 = 1.2$, which are the default values adopted in BM25. We implement the first approach in Java using the functionalities of Apache Lucene for processing texts, e.g., for tokenization. The average document length was computed over all the articles and was approximated to the next integer, namely 420. The vocabulary used to compute the indicators are those reported above: "risk*", "danger*", "hazard*", "damag*" and "harm*", where the asterisk notations is used to specify that the match in the document is a prefix match — all the words beginning with that word are considered a match and are counted to compute the statistics needed by the indicator instantiation. Hereafter, we will refer to these two instantiations using the label kwnf, since the basic idea is to rely on the KeyWord(s) Normalized Frequency. $\mathtt{kwnf}_1$ is the "simple" instantiation, $\mathtt{kwnf}_2$ refers to the instantiations of Equation 1 with with $b = 0.75$ and $k_1 = 1.2$. In TIPS the computation of the indicator is performed at indexing time and therefore it is stored along with other document metadata, accessible at query time. We did not rely on a stemmer to the represent document to capture the (prefix) terms in $\mathcal{K}$ in order to have a complete control on the adopted vocabulary.

As for the CMD, we relied on the implementation available in the text2map library.[2] Since CMD requires distributional representation, we explored two possible options: Glove trained on the dataset and fastText English Word Vectors trained on the Common Crawl. As for Glove, we used 50 as word vector dimension and 100 iterations. For each distributional representation, we considered two possible instantiations. The first instantiation, hereafter denoted as $\mathrm{CMD}_{risk}$, uses on the work *risk* to describe the focal concept. The second instantiation, hereafter denoted by $\mathrm{CMD}_{riskset}$, "defines" risk using compound concept.[3] The compound concept was specified using the same words adopted for the first approach: "risk danger hazard damage harm".

### 4.2. Dataset

As for the datasets, we relied on the corpora indexed in the TIPS Platform [16]. TIPS is a research platform designed and developed to assist researchers in the fields of Social Sciences and Humanities to carry out investigations of the discourse of Science and Technology in the Media. It provides search and analytics, e.g., the monitor the presence of technoscientific documents in the newspapers over time or in different section categories, to compute indicator values and extract topics and follow their trend/presence over time [16].

Through the TIPS platform, we sampled 1800 articles published from 2015 to 2019 in the online versions of five English Newspapers. Then we removed duplicates and near-duplicates using the approaches described in [16] which rely respectively on MD5 hashing of content and metadata and on Locality Sensitive Hashing for Minhash Signatures [18]. After duplicates and near-duplicates removal, we obtained 1736 articles. Those articles were distributed to a set of 10 assessors, which consisted of researchers in sociology, humanities, and linguistics; each of them was assigned 200 articles and had to specify if in the article

- no risk was suggested
- risk was suggested only marginally
- high risk suggested

Some articles were assessed in terms of risk by two researchers; if the assessment was not consistent among the two, we did not considered the document for the evaluation. Other articles were removed because the content no more accessible via the stored URL. The statistics on the resulting dataset are reported in Table 1.

The content of the articles was extracted using some regular expressions manually devised in the course of the TIPS Project and the newspaper3k library.[4]

---

[2]https://culturalcartography.gitlab.io/text2map/
[3]https://culturalcartography.gitlab.io/text2map/articles/CMDist-concept-movers-distance.html
[4]https://github.com/codelucas/newspaper

**Table 1**

Dataset used for the experimental evaluation.

| Class | Number of documents |
|---|---|
| No risk suggested | 1036 |
| Risk suggested only marginally | 185 |
| Risk highly suggested | 184 |

**Table 2**

Experimental results in terms of NDCG.

| Method | NDCG@10 | NDCG@20 |
|---|---|---|
| $\mathrm{kwnf}_1$ | 0.6255 | 0.4718 |
| $\mathrm{kwnf}_2$ | 0.4804 | 0.4860 |
| $\mathrm{CMD}_{\mathrm{glove+risk}}$ | 0.4137 | 0.3647 |
| $\mathrm{CMD}_{\mathrm{fasttext+risk}}$ | 0.7217 | 0.5918 |
| $\mathrm{CMD}_{\mathrm{glove+riskset}}$ | 0.4515 | 0.4214 |
| $\mathrm{CMD}_{\mathrm{fasttext+riskset}}$ | 0.7652 | 0.6402 |

We opted for considering articles gathered from the TIPS platform because the assessment of the level of risk was carried out in conjunction with another task: assessing relevance of the articles with respect to science and technology. Indeed, TIPS is equipped with text classifiers that rely on Supervised Machine Learning techniques; therefore, labeled data is needed both to train the classifiers and to evaluate their effectiveness. The sample was extracted to check the effectiveness of the classifiers, originally trained and evaluated on a different labeled set of documents. When labeling this additional sample, we ask the assessors to assess also the "level" of risk.

### 4.3. Results

In order to evaluate the effectiveness of the approaches described in Section 3, we exploited two measures. One is the Spearman Rank Correlation between the "true" labels – the one obtained as described in Section 4.2 – and the scores assigned by the different considered approaches. The other measure adopted is the Normalized Discounted Cumulative Gain (NDCG) [19], specifically the variant adopted in [20]. This measure provides an indication of the approach capability to rank highly relevant (in our case "risky") documents at high rank positions. The gain adopted to compute the NDCG were:

- 0: no risk suggested
- 1: risk suggested only marginally
- 2: risk highly suggested

We computed both NDCG@k with two different cutoff: $k = 10$ and $k = 20$.

Table 2 reports the results in terms of NDCG. The most promising approaches, both for NDCG@10 and NDCG@20, are the CMD instantiations that rely on fasttext. Among these, the approach using the controlled vocabulary to represent the "compound" concept proves to be the most effective. More simple approaches such as $\mathrm{kwnf}_1$, which only relies on normalized term frequency, is one of the most effective in terms of NDCG@10.

Table 3 reports the results obtained for the two instantiations of Equation 1 and those based on CMD and fasttext, since they were the most effective CMD instantiations. Results are reported in terms of Spearman Rank Correlation. The results show a modest but positive correlation between the true values and the computed scores. Also for this measure, the results suggest that the $\mathrm{CMD}_{\mathrm{riskset}}$ variant, which exploits the compound concept to describe risk, is the most promising.

**Table 3**
Experimental results in terms of Spearman Rank Correlation.

| Method | Spearman Rank Correlation |
| --- | --- |
| $\text{kwnf}_1$ | 0.2269 |
| $\text{kwnf}_2$ | 0.2284 |
| $\text{CMD}_{\text{fasttext+risk}}$ | 0.3018 |
| $\text{CMD}_{\text{fasttext+riskset}}$ | 0.3469 |

## 5. Final Remarks

In this paper, we reported on an ongoing investigation we are carrying out on possible instantiations of the *risk indicator*. We built a dataset with the collaboration of expert users in Social Sciences and Humanities. We studied the effectiveness of two previously proposed approaches and the correlation of the predicted "risk indicator" value with the gathered labels. The correlation is positive but weak, suggesting further work is needed.

Indeed, the current work should be considered as a preliminary investigation since there are a number of limitations that must be addressed.

One limitation is the modest size of the dataset, which includes only 1405 documents. As part of a teaching activity carried out in a Master Degree Course at the University of Padova, we have asked students in the course to label additional documents using the considered three-level risk scale. The result of this activity will be adopted to increase the size of the dataset. We are planning to redistribute some of the articles in the labeled dataset in order to analyze the consistency of the assessment among the experts.

Another limitation is the type of considered approaches. The term frequency-based normalization was adopted to investigate the effectiveness of a simple approach. CMD was considered in order to benefit from distributional representation. However, other approaches besides Glove and Fasttext might be adopted, for example, representations rooted in transformer-based architectures such as BERT and its subsequent extensions/variants, that were shown to be particularly promising — see for instance [21]. These representations can be adopted in conjunction with the CMD or with novel strategies, e.g., based on subspace-based representations.

The approach considered in this work relies solely on the content of the documents. We are pursuing this research direction due to the nature of the informative resources under consideration—namely, newspapers—where reader interactions, such as comments, are not necessarily available. These types of approaches can also be applied to other kinds of documents, such as parliamentary data [22].

However, content-based approaches can be complemented with other sources of evidence, such as reader interactions or information gathered from other media streams. User interaction, in particular, has shown to be promising in tasks like the analysis and prediction of controversies in social media streams, such as Reddit [23].

As mentioned in Section 2, another research direction that we intend to investigate is the relationship with Media Frames and the effectiveness of automatic approaches for frame detection and analysis [11, 24, 13] for the task considered in this paper. The focus will be on approaches such as that proposed in FrameFinder [14] for frame label prediction, which models the problem as a zero-shot prediction task.

# References

[1] E. Di Buccio, A. Lorenzet, M. Melucci, F. Neresini, Unveiling latent states behind social indicators, in: R. Gavaldà, I. Zliobaite, J. Gama (Eds.), Proceedings of the First Workshop on Data Science for Social Good, SoGood@ECML-PKDD 2016, Riva del Garda, Italy, September 19, 2016, volume 1831 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2016.

[2] F. Neresini, A. Lorenzet, Can media monitoring be a proxy for public opinion about technoscientific controversies? The case of the Italian public debate on nuclear power., Public understanding of science (Bristol, England) (2014).

[3] D. S. Stoltz, M. A. Taylor, Concept mover's distance: measuring concept engagement via word embeddings in texts, Journal of Computational Social Science 2 (2019) 293–313. doi:`10.1007/s42001-019-00048-6`.

[4] D. E. Losada, F. Crestani, A Test Collection for Research on Depression and Language Use, 2016, pp. 28–39. doi:`10.1007/978-3-319-44564-9_3`.

[5] A. Rortais, F. Barrucci, V. Ercolano, J. Linge, A. Christodoulidou, J.-P. Cravedi, R. Garcia-Matas, C. Saegerman, L. Svečnjak, A topic model approach to identify and track emerging risks from beeswax adulteration in the media, Food Control 119 (2021) 107435. doi:`10.1016/j.foodcont.2020.107435`.

[6] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent dirichlet allocation, Journal of Machine Learning Research 3 (2003). doi:`10.5555/944919.944937`.

[7] J. P. Linge, R. Steinberger, T. Weber, R. Yangarber, E. van der Goot, D. Al Khudhairy, N. Stilianakis, Internet surveillance systems for early alerting of health threats, Eurosurveillance 14 (2009) 19162.

[8] R. Steinberger, B. Pouliquen, E. van der Goot, An introduction to the Europe Media Monitor family of applications, in: Proceedings of the SIGIR 2009 Workshop on Information Access in a Multilingual World, volume 43, 2009. `arXiv:1309.5290`.

[9] S. Bhatia, Predicting risk perception: New insights from data science, Management Science 65 (2019) 3800–3823. doi:`10.1287/mnsc.2018.3121`.

[10] R. M. Entman, Framing: Toward clarification of a fractured paradigm, Journal of Communication 43 (1993) 51–58. doi:`10.1111/j.1460-2466.1993.tb01304.x`.

[11] A. E. Boydstun, D. Card, J. Gross, P. Resnick, N. A. Smith, Tracking the Development of Media Frames within and across Policy Issues (2018). doi:`10.1184/R1/6473780.v1`.

[12] M. Ali, N. Hassan, A survey of computational framing analysis approaches, in: Y. Goldberg, Z. Kozareva, Y. Zhang (Eds.), Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 2022, pp. 9335–9348. doi:`10.18653/v1/2022.emnlp-main.633`.

[13] J. Piskorski, N. Stefanovitch, G. Da San Martino, P. Nakov, SemEval-2023 task 3: Detecting the category, the framing, and the persuasion techniques in online news in a multi-lingual setup, in: A. K. Ojha, A. S. Doğruöz, G. Da San Martino, H. Tayyar Madabushi, R. Kumar, E. Sartori (Eds.), Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023), Association for Computational Linguistics, Toronto, Canada, 2023, pp. 2343–2361. doi:`10.18653/v1/2023.semeval-1.317`.

[14] M. Reiter-Haas, B. Klösch, M. Hadler, E. Lex, Framefinder: Explorative multi-perspective framing extraction from news headlines, in: CHIIR 2024 - Proceedings of the 2024 Conference on Human Information Interaction and Retrieval, Association for Computing Machinery, Inc, 2024, pp. 381–385. doi:`10.1145/3627508.3638308`.

[15] A. Singhal, C. Buckley, M. Mitra, Pivoted document length normalization, in: Proceedings of the 19th annual international ACM SIGIR conference on Research and development in information retrieval - SIGIR '96, ACM Press, 1996, pp. 21–29. doi:`10.1145/243199.243206`.

[16] E. Di Buccio, A. Cammozzo, F. Neresini, A. Zanatta, TIPS: search and analytics for social science research, in: L. Tamine, E. Amigó, J. Mothe (Eds.), Proceedings of the 2nd Joint Conference of the Information Retrieval Communities in Europe (CIRCLE 2022), Samatan, Gers, France, July 4-7, 2022, volume 3178 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2022.

[17] M. J. Kusner, Y. Sun, N. I. Kolkin, K. Q. Weinberger, From word embeddings to document distances, in: Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15, JMLR.org, 2015, p. 957–966.

[18] A. Z. Broder, Identifying and filtering near-duplicate documents, in: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), volume 1848, 2000, pp. 1–10. doi:10.1007/3-540-45123-4_1.

[19] K. Järvelin, J. Kekäläinen, Cumulated gain-based evaluation of IR techniques, ACM TOIS 20 (2002) 422–446.

[20] C. L. A. Clarke, N. Craswell, I. Soboroff, G. V. Cormack, Overview of the TREC 2010 Web Track, in: Proceedings of The Nineteenth Text REtrieval Conference, TREC 2010, Gaithersburg, Maryland, USA, November 16-19, 2010, volume Special Publication 500-294, NIST, 2010.

[21] O. Khattab, M. Zaharia, Colbert: Efficient and effective passage search via contextualized late interaction over bert, in: Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20, Association for Computing Machinery, New York, NY, USA, 2020, p. 39–48. doi:10.1145/3397271.3401075.

[22] T. Erjavec, M. Ogrodniczuk, P. Osenova, N. Ljubešić, K. Simov, A. Pančur, M. Rudolf, M. Kopp, S. Barkarson, S. Steingrímsson, c. Çöltekin, J. de Does, K. Depuydt, T. Agnoloni, G. Venturi, M. C. Pérez, L. D. de Macedo, C. Navarretta, G. Luxardo, M. Coole, P. Rayson, V. Morkevičius, T. Krilavičius, R. Darundefinedis, O. Ring, R. van Heusden, M. Marx, D. Fišer, The parlamint corpora of parliamentary proceedings, Language Resources and Evaluation 57 (2022) 415–448. doi:10.1007/s10579-021-09574-0.

[23] P. Koncar, S. Walk, D. Helic, Analysis and prediction of multilingual controversy on reddit, in: 13th ACM Web Science Conference 2021, ACM, 2021, pp. 215–224. doi:10.1145/3447535.3462481.

[24] D. Card, A. E. Boydstun, J. H. Gross, P. Resnik, N. A. Smith, The media frames corpus: Annotations of frames across issues, in: C. Zong, M. Strube (Eds.), Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 438–444. doi:10.3115/v1/P15-2072.