

How Does the Language of ‘Threat’ Vary Across News Domains?

A Semi-Supervised Pipeline for Understanding Narrative Components in News Contexts

Igor Ryazanov¹ and Johanna Björklund¹

Abstract—By identifying and characterising the narratives told in news media we can better understand political and societal processes. The problem is challenging from the perspective of natural language processing because it requires a combination of quantitative and qualitative methods. This paper reports on work in progress, which aims to build a human-in-the-loop pipeline for analysing how the variation of narrative themes across different domains, based on topic modelling and word embeddings. As an illustration, we study the language associated with the threat narrative in British news media.

I. INTRODUCTION

Due to the drastic changes in news distribution over the past decades, considerable attention has been given to how ongoing events are framed in news reporting. In the realm of digital news media, concerns include increasing polarisation and a decrease in the relative share of political reporting [1]. News reporting has a direct effect on the political landscape because alternative news framing translates into competing public discourses and, by extension, electoral results [2]. Studying the framings and narratives in the media is, therefore, vital for understanding political processes. Extensive qualitative analysis of large-scale news corpora is, however, expensive and can hardly be feasible. This provides the motivation to apply natural language processing to both facilitate qualitative research at scale and enable quantitative approaches to narrative understanding. In this paper, we propose a pipeline for descriptive multi-domain analysis of narrative subcomponents in the news media.

Advancements in natural language processing (NLP) methods provide a variety of tools for political communication analysis. Some of these are found in applications within or adjacent to the area of narrative understanding, such as stance detection [3] and sentiment analysis [4]. Algorithms based on neural networks have been shown to discern the difference in published texts by different partisan actors [5] and predict the ideological alignment of social media posters [6]. The majority of these methods can be described as classification algorithms, relying on supervised learning on either narrow domain-specific annotated datasets or fine-tuning large language models (LLMs) which have been trained on huge general-domain corpora. They are

primarily used for quantitative studies and practical applications where the target phenomena are well-defined and the domain shift is limited. The methods, however, are not always suitable to assist qualitative and descriptive research. The low level of interpretability of the machine learning methods in general, and of the LLMs in particular, is also a factor.

We are studying the usefulness of NLP for fine-grained narrative structure analysis, e.g. extracting narrative substructures and revealing context-specific language. The proposed pipeline, which is still a work in progress, serves to describe the contextual use of narrative themes within different overarching topics. As an example, we study the language used by news publishers to express the notion of threat and risk. We choose this type of semantic relation because its presence in a news article almost guarantees a degree of partiality which affects the readers’ perception of the issue.

The pipeline consists of the following steps:

- Applying semi-supervised or unsupervised topic modelling to find latent topics in a text corpus
- Training contextual embeddings for each discovered topic
- Computing the closest terms in the embedding space to describe the notion of threat in each topic

In the system presented here, the embeddings are produced by Word2Vec, and topics are derived through Correlation Explanation (CorEx) where clusters are shaped around user-provided anchor words [7]. Depending on the available domain knowledge and discoveries from unsupervised clustering, the anchor words can guide the model to find crisper topics in a semi-supervised fashion. The output of the pipeline is a collection of descriptions of a selected concept (in our case, threat) for each of the generated topics.

II. BACKGROUND

In the context of NLP applications, interpretations and definitions of narrative structures and elements can vary greatly, depending on specific tasks and domains. Here we relate our problem to several of these approaches. While our goal to investigate the language abstract notions in different contexts (here exemplified by threat) does not match them exactly, it shares many similarities with e.g. stance detection and narrative discovery.

¹ I. Ryazanov and J. Björklund are with the Department of Computing Science, Faculty of Science and Technology, Umeå University, Umeå, Sweden igorr@cs.umu.se; johanna@cs.umu.se. This work is supported by Marianne and Marcus Wallenberg’s Foundation, the Swedish Research Council, and WASP-HS.

A. Opinion mining and stance detection

As demonstrated in a number of studies, certain narrative-like notions can be captured by large language models trained on the document level. The documents or sentences are labelled by human annotators as containing such notions, and the definition of the notion is left to expert judgement. The assumption is that the representation of the document is rich enough that it captures narrative elements regardless of form. This is very prominent in, e.g. hate speech detection, where the key challenge comes from the fact that the hateful intent can take misleading forms and does not rely on any specific device to be conveyed [8], [9]. Similarly, it can be the case for the stance detection task, where the stance towards an issue cannot always be determined by positive or negative vocabulary and sentiment analysis is not enough to draw meaningful conclusions [3].

The downside is that the model trained to classify entire documents would be able to identify, e.g. a stance towards a specific political issue, but not necessarily what constitutes the narrative within the text. For example, an article can be shown to include the notion of ‘threat’, but explaining what constitutes ‘threat’ beyond the label becomes problematic. Since the algorithmic decision applies to an entire document, for narrative detection purposes, these models are more applicable to shorter messages and higher-level narratives (‘pro-abortion’ as opposed to ‘threat’ or ‘success’).

B. Narrative extraction

In the field of computational narratology, a common approach to narratives involves determining key entities and their relations. The inspiration for such methods comes from the structural interpretations of stories in formalist folklore studies [10]. Character or role detection can take various forms but often includes assigning a fixed set of archetypal roles or narrative frames (‘villain’, ‘protagonist’, ‘victim’, etc.) to specific entities in the story. In the news article domain, abstract role detection has been realised, for example, through the combination of entity extraction and sentiment analysis [11]. There has been, for example, some success in applying these methods in computational studies of conspiracy theories classifying entities within the publication as ‘insiders’ or ‘outsiders’ [12].

A more complex approach involves constructing the relations of the extracted entities, where the resulting narrative representation usually takes the form of a graph. This method has been applied, for instance, to conspiracy theory discovery. Relating this directly to our research, the authors used the presence of the notion of ‘threat’ encoded as subject-verb-copula triplets as the main criterion to detect specific theory elements [13]. We, however, are interested in how the narrative component of ‘threat’ is different in different news contexts and not in which contexts it defines.

C. Perspective extraction

Recently, Minnema et al. [14] put forth a framework based on Frame Semantics. They apply a FrameNet [15] parser LOME [16] to analyse perspectives in news media event

description. Instead of building a graph, the focus is on analysing linguistic frames invoked by specific texts. While the purpose of the model is similar to our task, it is focused on the analysis of the specific events or topics (e.g. femicide reporting in Italy [17]) rather than comparing the contexts.

III. PRE-STUDY

A. Semi-supervised topic modelling for contextual ‘threat’ understanding

In our goal to keep the pipeline as robust and explainable as possible, we investigated the possibility of extracting contextual descriptions of ‘threat’ purely by applying semi-supervised topic modelling, which has the advantage of being interpretable in terms of probabilities, unlike Word2Vec word embeddings. Semi-supervised topic modelling has been used, for example, to investigate the presence of gendered latent topics in different contexts [18]. We explored the option of taking a similar approach, presupposing that there exists a specific cluster of news articles or their fragments centred around the target concept of threat. All of the pre-study has been performed on the same dataset as the rest of the paper and its thematical subsets: sports and politics. If the subsets would each contain a threat-related cluster of articles, we would have been able to compare their content and, therefore, the definitions of threat in these contexts.

B. Experiments with pSSLDA and CorEx

In the first series of experiments, we applied Latent Dirichlet Allocation (LDA) with z -labels (pSSLDA) [19]. At the initialisation step, it assigns additional weight to predefined seed words for specific topics. After initialisation, the algorithm proceeds in an unsupervised fashion. Through our experiments, we initialised clusters with various threat-related word combinations, as well as tested other similar notions, such as ‘success’. The resulting topical distribution remained near-identical to the output of the unsupervised LDA model. Moreover, the topic order remained unstable even with the seeding, somewhat counterintuitively: one could have expected the military conflict-related news topic to be consistently initiated by seed words, such as ‘danger’.

The second series of preliminary experiments using the more restrictive CorEx also displayed negative results: the topics initialised with the threat-related anchor (seed) words did not seem to be immediately humanly interpretable and had significant overlaps with other clusters. Our interpretation is that in a news dataset, the event-specific language dominates all other vocabulary particularities, making event-based topics very easily separable. So even if clusters of text corresponding to the notions, such as ‘threat’ exist, they remain statistically insignificant in comparison. While this may not be the case for other abstract topics, it is reasonable to expect a co-occurrence-based method to find clusters based on topic-specific terms rather than the presence of a higher-level semantic construct. Thus, we rejected using purely semi-supervised topic modelling for our task.

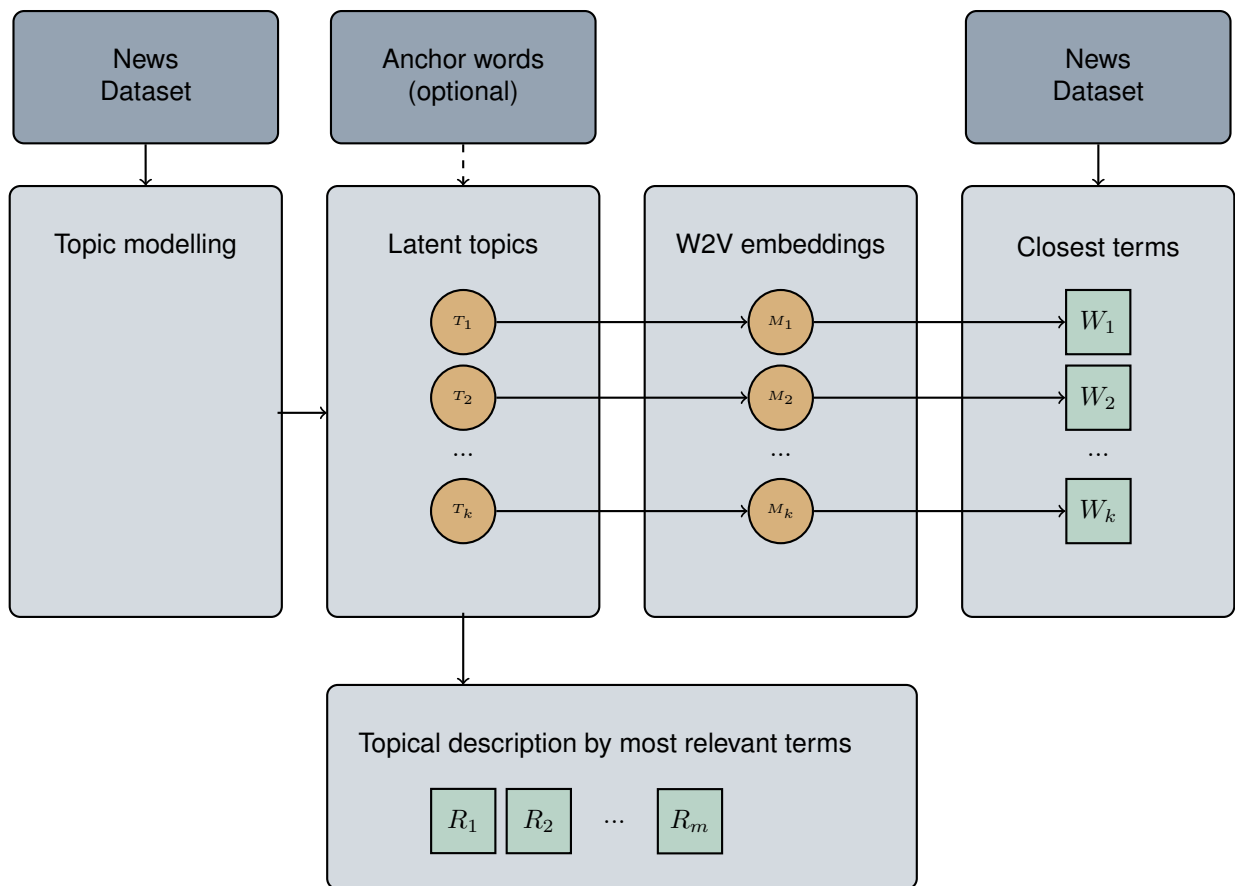


Fig. 1. In the proposed pipeline, the user provides a corpus for topic modelling, a parameter k declaring the expected number of topics, and, if desired, a number of anchor words to guide the topic formation. The result is k topics, represented by topical descriptions R_i consisting of the m terms most relevant to each. If the user is not satisfied, they may update the anchor words and repeat the clustering until the topics are as desired. In the next step, semantic embeddings M_i are computed for each topic T_i , and based on these, the closest set of closest terms W_i to the target concept are extracted.

IV. METHODS

A. Topic Modelling

To demonstrate the chosen approach, we apply it to study the language of threat in news media. The initial topic modelling is done with the help of CorEx, and in our analysis, we equate these topics with ‘contexts’. Since the purpose of the pipeline is to perform exploratory analysis and assist in qualitative studies, it is vital to have some control over topic distribution based on domain knowledge. CorEx is based on mutual information between words and topics and offers more restrictive and flexible semi-supervised functionality compared to, e.g. Latent Dirichlet Allocation with z -labels [19]. It is also valuable for corpus exploration because it does not require transfer learning: Topic modelling based on neural networks can perform better on specific tasks [20], but also introduces additional bias from out-of-context corpora. Due to the nature of the models, this bias cannot be easily separated from the properties of the target datasets, which is a disadvantage in exploratory analysis.

B. Word Embeddings

After the topics have been established, we train Word2Vec embeddings on texts in each individual topic. Word em-

beddings preserve some degree of semantic relations from natural language [21] and have been used for investigating the definitions of concepts, as well as the narratives surrounding them. For example, Papisavva et al. [22] apply Word2Vec to find words associated with QAnon. It has also been used to compare semantic contexts: e.g. the language use of parliamentary motions of the opposing Swedish political parties [23]. While our goal is somewhat different, the approach is similar. We use a set of keywords as a representation for the target concept: in this case, threat.

Since our principal interest is investigating concrete media contexts, we avoid language models that require pre-training. Such models would inject implicit out-of-context bias, making comparisons harder. Additionally, the document cluster for each topic is relatively small, and it has been shown that Word2Vec can outperform transformer-based language models on smaller datasets when trained from scratch [24].

V. EXPERIMENTAL SETUP

A. Dataset

To evaluate the pipeline (outlined in Figure 1), we experiment with a collection of news articles from mainstream free-access British news media. The dataset was collected

between May and early August 2022 and contains 57,996 unique articles out of 100,000 in total. Top-5 most frequent news sources are *The Sun*, *The Independent*, *The Daily Express*, *The Daily Mirror* and *The Daily Star* together constituting approximately 23% of the articles in the dataset.

For each article, the following information was collected: the title (headline), the preamble, the body, the URL to the article, and the publication date and time. For the purposes of the experiments, the headline, the preamble and the body are concatenated into single text entries. The articles are tokenised and processed into the matrix of token counts with the Scikit-learn library.

B. Threat Definition

In this experiment, our goal is to define the notion of threat through the common language terms that encapsulate the relation of one entity presenting a threat to another. In that, our definition is similar to the definition of the frames in FrameNet: a frame is “a script-like conceptual structure that describes a particular type of situation, object, or event along with its participants and props” [15]. The semantically closest frame of FrameNet is ‘Risky_situation’ – “A particular Situation is likely (or unlikely) to result in a harmful event befalling an Asset.”, so we choose to use the nouns in the FrameNet ontology that invoke the ‘Risky_situation’ (*threat*, *danger*, *risk*) as a triple of keywords.

VI. EXPERIMENTS

A. Unsupervised Topic Modelling

We assume minimal domain knowledge and task CorEx to identify news topics without anchor words to detect latent topics in the dataset. The number of topics is initially chosen to be below 10 to limit the scope of further analysis and after the preliminary experimentation set to 5. We evaluated the top-20 most relevant terms for each cluster, and in four cases out of five, the topics seem to have a clear focus (top-3 terms listed in parenthesis):

- T_0 : (*league*, *season*, *premier*)
- T_1 : (*government*, *cost*, *crisis*)
- T_2 : (*police*, *court*, *officers*)
- T_3 : (*love*, *instagram*, *star*)
- T_4 : (*like*, *just*, *think*)

The final topic T_4 has the lowest topic correlation value. It is based on non-topic-specific words and seems to include articles not matching with the rest of the clusters.

B. Semi-supervised Topic Modelling

Based on the initial result above, we can expect that the four topics (loosely defined as ‘Sports’, ‘Costs crisis’ ‘Police’, and ‘TV and celebrities’) are present in the data, but it would be unreasonable to assume that all (or even most) articles would belong to one of them. Gallagher et al. [7], when investigating the performance of CorEx, set the number of clusters as high as 50 for a news dataset while initialising some of them with anchors to create crisp topics. The four topics described above are already shown to be present in the data. Even with minimal domain knowledge, we can also

expect significant coverage of the Russia-Ukraine war in the summer of 2022, which is a potentially valuable context for threat interpretation. To isolate the five chosen topics, we restrict them with three anchor words each and set the total number of topics to 20.

The anchor words and the resulting topics are shown in Table I. Each of them is described by the top 10 most relevant terms. The list includes anchor words used for initialisation (in bold). Corex is a discriminative model and allows articles to belong to several topics.

C. Word Embeddings

For the articles in each topic, we train a Word2Vec model and extract the terms in the embedding space that are the closest to the three keywords: *threat*, *danger* and *risk*. We use the Gensim implementation of Word2Vec with the following parameters for all individual models: ignoring unique words (frequency 2 or more), word window size – 10, vector size – 100. The results are presented in Table II. For each context–keyword pair, we show seven words with the closest vector representation measured by cosine distance. Words unique to the contexts are highlighted in bold. As we can see, these neighbourhoods vary greatly between the contexts.

D. Analysis

At this step in the pipeline, the automatic analysis could be complemented with expert knowledge to add a qualitative element to the study. However, even by analysing the output of the models superficially and without specialised knowledge, we notice certain peculiarities. One such thing is that the language of the ‘cost crisis’ topic is as strong if not stronger than the language of the ‘war’ topic. Another observation is that the word ‘fetus’ is likely present together with ‘court’ because the Roe v. Wade ruling was overturned by the US Supreme Court within the time frame. With greater domain knowledge, one could choose better anchors and obtain crisper clusters. It is, for example, likely the term ‘TV’ caused the last topic to skew towards war instead of the cultural sphere as was intended.

VII. DISCUSSION

A. Limitations

One immediate limitation of the pipeline is the need for human guidance in the clustering process. While annotated data is not necessary, clustering does benefit greatly from domain knowledge, as the unsupervised version is unlikely to produce meaningful results. Another related disadvantage is the lack of ‘ground truth’ knowledge to test the results with the problem framed as it is. We, however, work on the assumption that a media researcher would not look for purely quantitative output in their application, using this framework as a facilitator for qualitative research instead.

Another potential weakness is the need to define any potential target concept with the keywords. While it is not in itself problematic with a domain expert’s input, we so far lack the evidence to judge what kind of concept definition is preferable in a general case.

TABLE I
FIVE TOPICS INITIATED WITH ANCHOR WORDS FROM THE ORIGINAL DATASET.

Id	Size	Top-10 terms
T_0	4018	war, ukraine, russia , russian, invasion, military, putin, ukrainian, forces, vladimir
T_1	13,594	league, player, sport , players, football, squad, goals, champions, clubs, winger
T_2	8249	cost, crisis, economy , living, struggling, cuts, poverty, spiralling, poorest, unemployment
T_3	10,134	police, court, officers , arrested, incident, investigation, crime, victim, judge, guilty
T_4	8070	tv, celebrity, singer, series , ekin, davide, su, island, luca, sanclimenti

TABLE II
THREAT, RISK, AND DANGER IN FIVE CONTEXTS (IDENTIFIED LATENT TOPICS), REPRESENTED BY THE CLOSEST TERMS. TERMS UNIQUE TO EACH TOPIC IN THIS SELECTION ARE HIGHLIGHTED IN BOLD.

Context	Key word	Top-7 closest terms
'War'	threat	geopolitical, escalation, strategic, warfare, threats, counter, risks
	danger	context, grave, disaster, dangers, height, catastrophe, breadth
	risk	risks, situation, safety, disaster, consequences, escalation, threat
'Sports'	threat	creativity, backline, pressing, presence, attack, defence, possession
	danger	trouble, threat, territory, control, possession, fall, lines
	risk	risks, cause, reduce, potentially, size, financial, prevent
'Cost crisis'	threat	escalation, conflict, threats, grave, geopolitical, danger, nuclear
	danger	threat, midst, fear, prospect, consequences, serious, concern
	risk	risks, damage, serious, consequences, expense, causes, cause
'Police'	threat	aggression, conflict, wider, issue, escalation, terrorist, outrage
	danger	risk, unnecessary, fetus, fear, cause, distress, potentially, harm
	risk	risks, potentially, danger, effective, level, cause, levels, nature
'TV & celebrities'	threat	nuclear, opposition, economic, missiles, conflict, kremlin, russia
	danger	circumstances, saturated, elements, doom, arteries, cynicism, turmoil
	risk	levels, height, safety, consumers, increases, damage, assistance

B. Future work

In our ongoing project, we have set ourselves several goals that would expand on existing experiments:

- We plan to study how the frequency and distribution of keywords or their combinations affect the results. As we propose to use the pipeline to study relatively high-level and not explicitly domain-specific concepts, we would like to at least outline how unspecific the keywords need to be. Based on this, we hope to provide a high-level recommendation on how to define the concepts. Similarly, we plan to formulate a recommendation on how to guide topic modelling with anchor words.
- Then, we aim to repeat the experimental scenario for other abstract narrative concepts, such as 'success' or 'failure', and compare the model's performances.
- The next step is to extend the experiments to an analogous dataset of the Swedish media. While it is reasonable to expect topic modelling and Word2Vec to

work similarly well in another Germanic language, the news media culture is different, which is likely to cover the use of language.

Moving further, we can see this pipeline being used in comparative studies of news publishing language in different contexts, not only limiting them to event-based topics. The contexts can include, e.g. different types of publications (mainstream vs tabloids vs new media) or political alignments ('left wing' vs 'right wing'). Another potential use case is comparing the language of the same publication over a time period to investigate how language shifts within particular news contexts. Finally, an even more challenging task requiring more topical expertise would be drawing comparisons between the same concepts in the news media of different countries in their respective languages.

C. Conclusion

We have implemented a semi-supervised pipeline to analyse the expression of narrative themes in different media contexts. Previous studies have used word embeddings to describe terms and stances, and we extend this to a more abstract notion and produce a complete pipeline to perform comparative analysis. Our next steps include applying the pipeline to other languages and comparing the performance and results to English-language media. We also reach out to media researchers to identify other relevant applications. We believe that when there is sufficient domain knowledge to guide topic formation, this mixed-method approach can be an effective tool for narrative analysis.

ACKNOWLEDGEMENT

This work was supported by the Wallenberg AI, Autonomous Systems and Software Program – Humanities and Society (WASP-HS), as well as the Marianne and Marcus Wallenberg Foundation and the Swedish Research Council.

We are thankful to Hannah Devinney and Anton Eklund for their helpful comments on an early version of this manuscript that did much to improve the content. We are also grateful to Adlede AB for their support in the data collection.

REFERENCES

- [1] P. Van Aelst, J. Strömbäck, T. Aalberg, F. Esser, C. De Vreese, J. Matthes, D. Hopmann, S. Salgado, N. Hubé, A. Stepińska, et al., “Political communication in a high-choice media environment: a challenge for democracy?,” *Annals of the International Communication Association*, vol. 41, no. 1, pp. 3–27, 2017.
- [2] L. Alonso-Muñoz and A. Casero-Ripollés, “Populism against europe in social media: The eurosceptic discourse on twitter in spain, italy, france, and united kingdom during the campaign of the 2019 european parliament election,” *Frontiers in communication*, vol. 5, p. 54, 2020.
- [3] D. Küçük and F. Can, “Stance detection: A survey,” *ACM Computing Surveys (CSUR)*, vol. 53, no. 1, pp. 1–37, 2020.
- [4] M. Birjali, M. Kasri, and A. Beni-Hssane, “A comprehensive survey on sentiment analysis: Approaches, challenges and trends,” *Knowledge-Based Systems*, vol. 226, p. 107134, 2021.
- [5] T. Fagni and S. Cresci, “Fine-grained prediction of political leaning on social media with unsupervised deep learning,” *Journal of Artificial Intelligence Research*, vol. 73, pp. 633–672, 2022.
- [6] W. Chen, X. Zhang, T. Wang, B. Yang, and Y. Li, “Opinion-aware knowledge graph for political ideology detection,” in *IJCAI*, vol. 17, pp. 3647–3653, 2017.
- [7] R. J. Gallagher, K. Reing, D. Kale, and G. Ver Steeg, “Anchored correlation explanation: Topic modeling with minimal domain knowledge,” *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 529–542, 2017.
- [8] F. Poletto, V. Basile, M. Sanguinetti, C. Bosco, and V. Patti, “Resources and benchmark corpora for hate speech detection: a systematic review,” *Language Resources and Evaluation*, vol. 55, pp. 477–523, 2021.
- [9] A. Schmidt and M. Wiegand, “A survey on hate speech detection using natural language processing,” in *Proceedings of the fifth international workshop on natural language processing for social media*, pp. 1–10, 2017.
- [10] V. I. Propp, *Morphology of the Folktale*, vol. 9. University of Texas Press, 1968.
- [11] D. Gomez-Zara, M. Boon, and L. Birnbaum, “Who is the hero, the villain, and the victim? detection of roles in news articles using natural language techniques,” pp. 311–315, 2018.
- [12] P. Holur, T. Wang, S. Shahsavari, T. Tangherlini, and V. Roychowdhury, “Which side are you on? Insider-Outsider classification in conspiracy-theoretic social media,” pp. 4975–4987, 2022.
- [13] S. Shahsavari, P. Holur, T. Wang, T. R. Tangherlini, and V. Roychowdhury, “Conspiracy in the time of corona: Automatic detection of emerging covid-19 conspiracy theories in social media and the news,” *Journal of computational social science*, vol. 3, no. 2, pp. 279–317, 2020.
- [14] G. Minnema, S. Gemelli, C. Zanchi, T. Caselli, and M. Nissim, “SocioFillmore: A Tool for Discovering Perspectives,” pp. 240–250, 2022.
- [15] C. F. Baker, C. J. Fillmore, and J. B. Lowe, “The berkeley framenet project,” in *COLING 1998 Volume 1: The 17th International Conference on Computational Linguistics*, 1998.
- [16] P. Xia, G. Qin, S. Vashishtha, Y. Chen, T. Chen, C. May, C. Harman, K. Rawlins, A. S. White, and B. Van Durme, “LOME: Large Ontology Multilingual Extraction,” pp. 149–159, 2021.
- [17] G. Minnema, S. Gemelli, C. Zanchi, V. Patti, T. Caselli, and M. Nissim, “Frame semantics for social nlp in italian: Analyzing responsibility framing in femicide news reports,” in *Italian Conference on Computational Linguistics*, 2021.
- [18] H. Devinney, J. Björklund, and H. Björklund, “Semi-supervised topic modeling for gender bias discovery in english and swedish,” in *GeBNLP2020, COLING’2020–The 28th International Conference on Computational Linguistics, December 8-13, 2020, Online*, pp. 79–92, Association for Computational Linguistics, 2020.
- [19] D. Andrzejewski and X. Zhu, “Latent dirichlet allocation with topic-inset knowledge,” in *Proceedings of the NAACL HLT 2009 Workshop on Semi-Supervised Learning for Natural Language Processing*, pp. 43–48, 2009.
- [20] M. R. Grootendorst, “Bertopic: Neural topic modeling with a class-based tf-idf procedure,” *ArXiv*, vol. abs/2203.05794, 2022.
- [21] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” *Advances in neural information processing systems*, vol. 26, 2013.
- [22] A. Papasavva, J. Blackburn, G. Stringhini, S. Zannettou, and E. D. Cristofaro, ““is it a coincidence?”: An exploratory study of QAnon on Voat,” in *Proceedings of the Web Conference 2021*, pp. 460–471, 2021.
- [23] A. Fredén, M. Johansson, P. Kisić Merino, and D. Saynova, *A Comparison of Language Processing Models in Political Analysis: Evidence from Sweden*. Oct. 2021.
- [24] A. Edwards, J. Camacho-Collados, H. De Ribaupierre, and A. Preece, “Go simple and pre-train on domain-specific corpora: On the role of training data for text classification,” in *Proceedings of the 28th international conference on computational linguistics*, pp. 5522–5529, 2020.