WhatTheWikiFact: Fact-Checking Claims Against Wikipedia

Anton Chernyavskiy
HSE University
Moscow, Russia
aschernyavskiy_1@edu.hse.ru

Dmitry Ilvovsky HSE University Moscow, Russia dilvovsky@hse.ru Preslav Nakov Qatar Computing Research Institute, HBKU Doha, Qatar pnakov@hbku.edu.qa

ABSTRACT

The rise of Internet has made it a major source of information. Unfortunately, not all information online is true, and thus a number of fact-checking initiatives have been launched, both manual and automatic, to deal with the problem. Here, we present our contribution in this regard: WhatTheWikiFact, a system for automatic claim verification using Wikipedia. The system can predict the veracity of an input claim, and it further shows the evidence it has retrieved as part of the verification process. It shows confidence scores and a list of relevant Wikipedia articles, together with detailed information about each article, including the phrase used to retrieve it, the most relevant sentences extracted from it and their stance with respect to the input claim, as well as the associated probabilities. The system supports several languages: Bulgarian, English, and Russian.

CCS CONCEPTS

• Information systems \rightarrow Information retrieval; Information systems applications; Web searching and information discovery; • Computer systems organization \rightarrow Real-time systems; • Applied computing \rightarrow Document management and text processing; • Computing methodologies \rightarrow Natural language processing; Information extraction.

KEYWORDS

Fact-checking, Factuality, Veracity, Fake News, Disinformation, Misinformation, Document Retrieval, Sentence Retrieval, Stance Detection, Natural Language Inference, Wikipedia.

ACM Reference Format:

Anton Chernyavskiy, Dmitry Ilvovsky, and Preslav Nakov. 2021. What The WikiFact: Fact-Checking Claims Against Wikipedia. In Proceedings of . ACM, New York, NY, USA, 6 pages.

1 INTRODUCTION

Internet is abundant in platforms that allow users to share information online such as social networks, blogs, and forums. Unfortunately, not all information online is true, and there is a need to verify questionable claims that users encounter online.

As manual fact-checking is a complex and time-consuming task, it is important to develop tools that can help automate the process. Various task formulations have been proposed to automate fact-checking, e.g., SemEval RumourEval tasks [13, 20], the SemEval task on Fact Checking in Community Question Answering Forums [34], the CLEF CheckThat! lab [6, 16, 38, 40], the Fake News Challenge [21], and the FEVER task on Fact Extraction and VERification [3,

63, 64]. Here, we focus on the FEVER task, as it offers a large-scale training dataset, and enables explainable systems.

Interestingly, despite the popularity of research using the FEVER task formulation and its dataset, to the best of our knowledge, there are no publicly available running systems based on it. Here, we aim to bridge this gap with our WhatTheWikiFact system, which allows users to check claims against Wikipedia. It supports several languages (Bulgarian, English, and Russian) and uses a local Wikipedia for each language. The system first identifies relevant Wikipedia pages, and then finds relevant sentences within them. Then, it analyzes each claim submitted as an input against all extracted text fragments and makes a final verdict on the veracity of the claim: Truth, Lying, or Not Enough Info. What The Wiki Fact further displays all the information it used to make its decision, together with intermediate results for each verified claim, which includes (i) a list of the titles of the most relevant documents with links to the corresponding Wikipedia pages, and (ii) detailed information about each document as a table of retrieved text fragments (position in the document | text fragment), and a bar chart showing the confidence of the classifier in each label (Supports, Refutes or Not Enough Info) for each text fragment. This allows the user to quickly analyze the result.

The remainder of this paper is organized as follows: Section 2 discusses related work. Section 3 describes the dataset we used for training. Section 4 offers an overview of the system and its components. Section 4 discusses the core implementation details. Section 5 presents some evaluation results. Section 6 describes the system interface and its functionality, with some examples. Finally, Section 7 concludes and discusses future work.

2 RELATED WORK

Many task formulations have been proposed to address the spread of misinformation and disinformation online, and for each formulation, a number of approaches have been tried. Some good readings on the topic include surveys such as that by Shu et al. [59], who adopted a data mining perspective on "fake news" and focused on social media. Another survey [70] studied rumor detection in social media. The survey by Thorne and Vlachos [61] took a fact-checking perspective on "fake news" and related problems.

Li et al. [32] covered truth discovery in general. Lazer et al. [29] offered a general overview and discussion on the science of "fake news", while Vosoughi et al. [65] focused on the process of proliferation of true and false news online. Other recent surveys focused on stance detection [28], propaganda [12], social bots [17], false information [68] and bias on the Web [4]. Some very recent surveys focused on stance for misinformation and disinformation

¹The *WhatTheWikiFact* system is running online at https://www.tanbih.org/whatthewikifact

detection [23], on automatic fact-checking to assist human fact-checkers [39], on predicting the factuality and the bias of entire news outlets [42], and on multimodal disinformation detection [1].

Non-explainable fact-checking. The primary focus of our system and of this paper is fact-checking of claims. Relevant research includes work on credibility assessment in Twitter, which has been addressed using user-based, message-based, topic-based, and propagation-based features [8]. Rashkin et al. [51] analyzed the linguistic features used in the claims. Wang [66] presented the LIAR dataset, which focuses on fact-checking using only the input claim (its text and metadata). Lee et al. [30] found that misinformation can be discovered using perplexity analysis of the input claim, as perplexity is higher for false claims. A number of studies were also conducted on the feasibility of using language models for opendomain question answering and further as fact-checkers [31, 45, 54]. All this work was non-explainable.

Explainable fact-checking. This is a more relevant direction. Shaar et al. [55] developed two datasets for detecting previously fact-checked claims, which were extended and used for shared tasks as part of the CLEF CheckThat! lab in 2020 and 2021 [5, 6, 24, 40, 41, 56, 57]. Chen et al. [10] proposed table-based fact verification, and Gad-Elrab et al. [18] used knowledge graphs.

Stance detection as an element of fact-checking. This was the objective of the Fake News Challenge task [21, 53], as well as of the RumourEval tasks at SemEval in 2017 and 2019 [13, 20].

Fact-Checking Using Wikipedia. In our system, we use the FEVER dataset and task formulation, which enables Wikipedia-based explainable fact-checking [62]. The dataset was used in the FEVER shared tasks [63, 64], where most systems had the following components: (i) document retrieval, (ii) sentence retrieval, and (iii) natural language inference (NLI). Alonso-Reina et al. [2], Chakrabarty et al. [9], Hanselowski et al. [22], Otto [44] used a search API to retrieve relevant documents, while Yoneda et al. [67] used logistic regression. Word Mover's Distance [9], TF.IDF [33], ESIM [22], logistic regression [67], BERT [60] were used for sentence retrieval; and DAM [44], ESIM [25], Random Forest [52], LSTM [43] and BERT [60] were used for NLI. Here, we adopt a similar overall architecture.

System Demonstrations. Relevant demos include Hoaxy [58], for tracking misinformation in social networks and news sites, CredEye, [47] for credibility assessment, Tracy [18] for fact-checking using rules and knowledge graphs, Scrutinizer [27] for fact-checking statistical claims, STANCY [48] for stance detection using BERT and consistency constraints, Tanbih [69], which analyzes news articles and media outlets and predicts factuality of reporting, degree of propaganda, hyper-partisanship, political bias, framing, and stance with respect to various claims and topics, and FAKTA [35] for stance and evidence extraction from the Web.

Unlike these systems, we perform textual evidence-based fact-checking using Wikipedia, following the FEVER formulation of the task. The most related demo, FAKTA, is also trained on data from FEVER, but it focuses on the stance of a document, e.g., retrieved from the Web, with respect to the input claim, while our system

makes a prediction about a claim's factuality and gives sentencelevel evidence from Wikipedia to explain its decision. Moreover, our system supports several languages besides English.

3 DATA

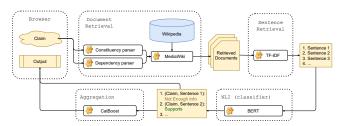
We train our system on the FEVER dataset [62], which includes a dump of 5.4M Wikipedia pages, and 220K claims labeled with one of the following three classes: Supports, Refutes or Not Enough Info. An example is shown in Figure 1. For the former two labels, there is also evidence provided, i.e., a sentence or a set of sentences that would allow one to make a verdict about the veracity of the input claim, while the Not Enough Info label indicates that there is no enough evidence in Wikipedia to prove or to refute the claim. There can be different sets of evidence for a given claim, but all of them would support the same label.



Figure 1: Manually labeled claim from FEVER.

4 SYSTEM OVERVIEW

Figure 2 shows the general architecture of our system, which is similar to the one described in [11]. First, the Document Retrieval (DR) module finds potentially relevant documents from Wikipedia. Then, the Sentence Retrieval (SR) module extracts the top-20 most relevant sentences from these documents. Afterwards, the Natural Language Inference (NLI) module classifies each claim–sentence pair as support/refute/NEI. Finally, the aggregation module makes a final prediction. We describe these steps below in more detail.



 $Figure\ 2: The\ architecture\ of\ our\ \textit{WhatTheWikiFact}\ system.$

Document Retrieval (DR). We use the Python MediaWiki API² to retrieve relevant documents from Wikipedia. We call the API with the results of our query generation module, which uses a

 $^{^2} http://wikipedia.readthedocs.io/en/latest/\\$

constituency parser [26] to extract noun phrases, which we use as separate queries. We further generate a query from the part of the claim up to the head word, which we extract using a dependency parser [15]. For each query, we retain the top-3 returned documents. As there can be many queries for an input claim, we further filter the results to improve the inference speed. Thus, we process the input claim and all retrieved titles shortened to the first bracket symbol using the Porter Stemmer [49], and we select them for the final set only if they are fully contained in the input.

Sentence Retrieval (SR). This component selects the top-20 most relevant sentences from the documents retrieved by DR component. The relevance is estimated as the cosine between the TF.IDF representations of the claim and of the document titles. We apply a variant of TF.IDF that uses binarization of the non-zero term counts, and considers words with a document fraction higher than 0.85 as stopwords. Note that there are no thresholds for the relevance score, and thus all retrieved sentences can potentially have a score of zero, e.g., if all top-20 documents happen to be irrelevant.

Natural Language Inference (NLI). This component takes as input the source claim, a separator, a sentence retrieved by the SR component, another separator, and the title of the document. Here, we use all retrieved sentences independently. For this component, we use BERT [14], which we fine-tune using a balanced training dataset. In particular, we use instances of the classes Supports and Refutes from the training set, and also the top-3 retrieved sentences for each claim, which we label as Not Enough Info.

Aggregation. We use CatBoost gradient boosting [50] to aggregate the sentence-level predictions of the NLI component. We train the model on part of the validation set using the stacked label probabilities predicted by the NLI model. Thus, we use a 60-dimensional (3 scores for 20 sentences) feature vector, potentially padded by zeros in case of lack of sentences in the retrieval phase.

We further report the maximal probability of the Supports (Refutes) class among all sentences scaled by the percentage of such sentences among the possibly relevant ones, i.e., excluding Not Enough Info, as the confidence score for Truth (Lying) predictions and the minimal probability of the Not Enough Info, otherwise.

Overall Architecture. WhatTheWikiFact has a server and a client parts, which are connected via a REST API. The client part is implemented using the Streamlit framework³—it is the GUI part, that is, the interface in the browser. The server part uses the Fast API⁴. We use Allen NLP library [19] for the constituency and the dependency parsers for texts in English, and Natasha⁵ and Stanza⁶ libraries for texts in Russian and Bulgarian, respectively. We use the official repository for BERT.⁷ We preload these models on the server, and we serve them using POST requests made by the client.

The parsers receive a piece of text as input, which is then tokenized and stemmed using the NLTK library [7] for English and Russian, and using BulStem [36, 37] for Bulgarian. BERT receives a list of sentence pairs, preprocessed using the BERT tokenizer.

Our NLI component is English-only, and thus we use the Google Translator API⁸ to translate the input into English first.

5 EVALUATION

The machine learning model we use in our *WhatTheWikiFact* system achieves an accuracy of 73.22 and a FEVER score of 67.44 on the test part of the FEVER dataset. Note that these scores are better than those for the best system at the FEVER shared task [63], which had an accuracy of 68.21, and a FEVER score of 64.21. They are also on par with the best system from the builder phase of the FEVER2.0 shared task [64], which had a FEVER score of 68.46. We should note that there have been some better results reported in the literature since then. In fact, we also had stronger results in our offline experiments. However, we use this particular model, as we need real-time execution, which requires certain compromise in terms of accuracy for the sake of improved speed of execution, which is essential for a real-time system like ours.

Our analysis shows that the biggest fraction of the system classification errors are in distinguishing Not Enough Info sentences from the rest. At the same time, the accuracy of the retrieval component is almost 91%, which means that in 91% of the cases, we do retrieve correct evidence for the final set of potentially relevant sentences. Therefore, as the system classification errors are not very frequent, such cases can be easily analyzed by the user in our system's output.

6 USER INTERFACE

Figure 3 shows a snapshot of *WhatTheWikiFact*'s output for an input claim. We can see that the system offers an overview of the verification results, which includes a verdict, the system's confidence in that verdict, and a list of possibly relevant documents: title and a link to the corresponding Wikipedia page.

Figure 3 shows that the system has retrieved three articles for the input claim "Napoleon Bonaparte declared Joan of Arc a national symbol." However, this does not mean that the system considers all of these articles as relevant. It just means that they contain some the top-20 most relevant sentences.

By clicking on *Show info*, the user can expand the panel with details about each document. As a result, the following information will be shown (illustrated on Figure 4):

- The part of the input claim used to retrieve the document.
- A bar chart of the predicted stance labels for the input claim with respect to each retrieved sentence. The stance is expressed as one of the classes Supports (SUP), Refutes (REF), or Not Enough Info (NEI). The chart further shows the class probability, which is also represented as the bar height, sentence number, and label, which is also indicated with the corresponding color. Note that there are three bars for each sentence, i.e., one for each label. Moreover, the bars are ordered (grouped) by labels or optionally by sentences (it is specified in the *Output options* section).
- A table of the most relevant sentences. For each sentence in that table, we show its position in the document as well as the document length, thus reflecting also the relative position and allowing for easy matching with the bar chart.

³http://streamlit.io/

⁴http://fastapi.tiangolo.com/

⁵http://github.com/natasha/natasha

⁶http://stanfordnlp.github.io/stanza/

⁷http://github.com/google-research/bert

 $^{^8} https://github.com/nidhaloff/deep-translator\\$

WhatTheWikiFact: Fact-Checking Claims Against Wikipedia This tool performs fact checking by comparing information with Wikipedia data. Select the language English Enter claims here (one per row, max 5) Napoleon Bonaparte declared Joan of Arc a national symbol. Output options + Verify Claim Verdict Confidence Retrieved articles: Joan of Arc Show info

Figure 3: Screenshot of what our WhatTheWikiFact system outputs when verifying the claim: "Napoleon Bonaparte declared Joan of Arc a national symbol."

For example, Figure 4 shows that the article "Joan of Arc" is relevant for fact-checking the input claim, as it includes a sentence that supports it. Other retrieved sentences in this document have almost 100% probability of a Not Enough Info label, and are thus irrelevant. Manual analysis by the user for the remaining two documents—"Napoleon", which was retrieved by the phrase "Napoleon Bonaparte", and "National symbol", which was retrieved by the phrase "a national symbol"—could confirm that they are indeed irrelevant. With these retrieval results, the user can manually inspect the boundary classification cases, especially for veracity prediction with low confidence or in case of Not Enough Info.

Show info

7 CONCLUSION AND FUTURE WORK

Napoleon

We have presented WhatTheWikiFact, a system for automatic claim verification using Wikipedia. The system reports the veracity for each input claim supplemented with evidence retrieved during the verification process, thus offering explainability. It also shows confidence scores and a set of relevant Wikipedia articles. Moreover, it allows the user to obtain detailed information about each article, including the exact phrase that was used to retrieve it, a list of the most relevant sentences with to the input claim that it contains, and their stance probabilities regarding the input claim. The system supports several languages: Bulgarian, English, and Russian.

In future work, we plan to implement a more accurate model by distilling knowledge from a larger model. Another direction we want to explore is to add additional languages for verification using local Wikipedias, which can be implemented without additional training, e.g., using multilingual BERT or language adapters [46].

ACKNOWLEDGMENTS

This research was done within the framework of the HSE University Basic Research Program.



Figure 4: Screenshot showing details about "Joan of Arc", a Wikipedia article retrieved when verifying the claim "Napoleon Bonaparte declared Joan of Arc a national symbol."

It is also part of the Tanbih mega-project (http://tanbih.qcri.org/), which is developed at the Qatar Computing Research Institute, HBKU, and aims to limit the impact of "fake news," propaganda, and media bias by making users aware of what they are reading.

REFERENCES

- Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fabrizio Silvestri, Dimiter Dimitrov, Giovanni Da San Martino, Shaden Shaar, Hamed Firooz, and Preslav Nakov. 2021. A Survey on Multimodal Disinformation Detection. arXiv/2103.12541 (2021).
- [2] Aimée Alonso-Reina, Robiert Sepúlveda-Torres, Estela Saquete, and Manuel Palomar. 2019. Team GPLSI. Approach for automated fact checking. In Proceedings of the Second Workshop on Fact Extraction and VERification (FEVER '19). Association for Computational Linguistics, Hong Kong, China, 110–114.
- [3] Rami AÎy, Zhijiang Guo, Michael Šejr Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. FEVEROUS: Fact Extraction and VERification Over Unstructured and Structured information. In Proceedings of the 35th Conference on Neural Information Processing Systems Datasets and Benchmarks Track (NeurIPS '21).
- [4] Ricardo Baeza-Yates. 2018. Bias on the Web. Commun. ACM 61, 6 (2018), 54-61.
- [5] Alberto Barrón-Cedeño, Tamer Elsayed, Preslav Nakov, Giovanni Da San Martino, Maram Hasanain, Reem Suwaileh, Fatima Haouari, Nikolay Babulkov, Bayan Hamdan, Alex Nikolov, Shaden Shaar, and Zien Sheikh Ali. 2020. Overview of CheckThat! 2020 Automatic Identification and Verification of Claims in Social Media. In Proceedings of the 11th International Conference of the CLEF Association: Experimental IR Meets Multilinguality, Multimodality, and Interaction (CLEF '20). Springer, Thessaloniki, Greece, 215–236.
- [6] Alberto Barrón-Cedeño, Tamer Elsayed, Preslav Nakov, Giovanni Da San Martino, Maram Hasanain, Reem Suwaileh, and Fatima Haouari. 2020. CheckThat! at CLEF 2020: Enabling the Automatic Identification and Verification of Claims in Social Media. In Proceedings of the 42nd European Conference on Information Retrieval (ECIR '20). Springer, Lisbon, Portugal, 499–507.
- [7] Steven Bird and Edward Loper. 2004. NLTK: The Natural Language Toolkit. In Proceedings of the ACL Interactive Poster and Demonstration Sessions (ACL '04). Association for Computational Linguistics, Barcelona, Spain, 214–217.
- [8] Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information Credibility on Twitter. In Proceedings of the 20th International Conference on World Wide Web (WWW '11). Association for Computing Machinery, Hyderabad, India, 675–684.
- [9] Tuhin Chakrabarty, Tariq Alhindi, and Smaranda Muresan. 2018. Robust Document Retrieval and Individual Evidence Modeling for Fact Extraction and Verification. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium, 127–131.
- [10] Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. 2020. TabFact: A Large-scale Dataset for Table-based Fact Verification. In Proceedings of the 8th International Conference on Learning Representations (ICLR '20). OpenReview.net, Addis Ababa, Ethiopia.
- [11] Anton Chernyavskiy and Dmitry Ilvovsky. 2019. Extract and Aggregate: A Novel Domain-Independent Approach to Factual Data Verification. In Proceedings of the Second Workshop on Fact Extraction and VERification (FEVER). Association for Computational Linguistics, Hong Kong, China, 69–78.
- [12] Giovanni Da San Martino, Stefano Cresci, Alberto Barrón-Cedeño, Seunghak Yu, Roberto Di Pietro, and Preslav Nakov. 2020. A Survey on Computational Propaganda Detection. In Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI-PRICAI '20). International Joint Conferences on Artificial Intelligence Organization, Yokohama, Japan (online), 4826-4832.
- [13] Leon Derczynski, Kalina Bontcheva, Maria Liakata, Rob Procter, Geraldine Wong Sak Hoi, and Arkaitz Zubiaga. 2017. SemEval-2017 Task 8: RumourEval: Determining rumour veracity and support for rumours. In Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval '17). Association for Computational Linguistics, Vancouver, Canada, 69-76.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT '19). Association for Computational Linguistics, Minneapolis, Minnesota, USA, 4171– 4186.
- [15] Timothy Dozat and Christopher D. Manning. 2017. Deep Biaffine Attention for Neural Dependency Parsing. In Proceedings of the 5th International Conference on Learning Representations (ICLR '17). OpenReview.net, Toulon, France.
- [16] Tamer Elsayed, Preslav Nakov, Alberto Barrón-Cedeño, Maram Hasanain, Reem Suwaileh, Pepa Atanasova, and Giovanni Da San Martino. 2019. CheckThat! at CLEF 2019: Automatic Identification and Verification of Claims. In Proceedings of the 41st European Conference on Information Retrieval (ECIR '19). CEUR-WS.org, Cologne, Germany, 309–315.
- [17] Emilio Ferrara, Onur Varol, Clayton Davis, Filippo Menczer, and Alessandro Flammini. 2016. The Rise of Social Bots. Commun. ACM 59, 7 (2016), 96–104.
- [18] Mohamed H. Gad-Elrab, Daria Stepanova, Jacopo Urbani, and Gerhard Weikum. 2019. Tracy: Tracing Facts over Knowledge Graphs and Text. In *Proceedings of the World Wide Web Conference (WWW '19)*. Association for Computing Machinery, San Francisco, CA, USA, 3516–3520.

- [19] Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. AllenNLP: A Deep Semantic Natural Language Processing Platform. In Proceedings of Workshop for NLP Open Source Software (NLP-OSS '18). Association for Computational Linguistics, Melbourne, Australia, 1–6.
- [20] Genevieve Gorrell, Elena Kochkina, Maria Liakata, Ahmet Aker, Arkaitz Zubiaga, Kalina Bontcheva, and Leon Derczynski. 2019. SemEval-2019 Task 7: RumourEval, Determining Rumour Veracity and Support for Rumours. In Proceedings of the 13th International Workshop on Semantic Evaluation (SemEval '19). Association for Computational Linguistics, Minnesola, USA, 845–854.
- [21] Andreas Hanselowski, Avinesh PVS, Benjamin Schiller, Felix Caspelherr, Debanjan Chaudhuri, Christian M. Meyer, and Iryna Gurevych. 2018. A Retrospective Analysis of the Fake News Challenge Stance-Detection Task. In Proceedings of the 27th International Conference on Computational Linguistics (COLING '18). Association for Computational Linguistics, Santa Fe, New Mexico, USA, 1859–1874.
- [22] Andreas Hanselowski, Hao Zhang, Zile Li, Daniil Sorokin, Benjamin Schiller, Claudia Schulz, and Iryna Gurevych. 2018. UKP-Athene: Multi-Sentence Textual Entailment for Claim Verification. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium, 103–108.
- [23] Momchil Hardalov, Arnav Arora, Preslav Nakov, and Isabelle Augenstein. 2021. A Survey on Stance Detection for Mis- and Disinformation Identification. arXiv/2103.00242 (2021).
- [24] Maram Hasanain, Fatima Haouari, Reem Suwaileh, Zien Sheikh Ali, Bayan Hamdan, Tamer Elsayed, Alberto Barrón-Cedeño, Giovanni Da San Martino, and Preslav Nakov. 2020. Overview of CheckThat! 2020 Arabic: Automatic Identification and Verification of Claims in Social Media. In Working Notes of CLEF 2020—Conference and Labs of the Evaluation Forum (CLEF '2020). CEUR-WS.org, Thessaloniki, Greece.
- [25] Christopher Hidey and Mona Diab. 2018. Team SWEEPer: Joint Sentence Extraction and Fact Checking with Pointer Networks. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics. Brussels, Belgium, 150–155.
- [26] Vidur Joshi, Matthew Peters, and Mark Hopkins. 2018. Extending a Parser to Distant Domains Using a Few Dozen Partially Annotated Examples. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL '18). Association for Computational Linguistics, Melbourne, Australia, 1190–1199.
- [27] Georgios Karagiannis, Mohammed Saeed, Paolo Papotti, and Immanuel Trummer. 2020. Scrutinizer: A Mixed-Initiative Approach to Large-Scale, Data-Driven Claim Verification. Proc. VLDB Endow. 13, 12 (2020), 2508–2521.
- [28] Dilek Küçük and Fazli Can. 2020. Stance Detection: A Survey. ACM Comput. Surv. 53, 1, Article 12 (2020), 37 pages.
- [29] David M.J. Lazer, Matthew A. Baum, Yochai Benkler, Adam J. Berinsky, Kelly M. Greenhill, Filippo Menczer, Miriam J. Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, Michael Schudson, Steven A. Sloman, Cass R. Sunstein, Emily A. Thorson, Duncan J. Watts, and Jonathan L. Zittrain. 2018. The science of fake news. Science 359, 6380 (2018), 1094–1096.
- [30] Nayeon Lee, Yejin Bang, Andrea Madotto, and Pascale Fung. 2021. Towards Few-shot Fact-Checking via Perplexity. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT '21). Association for Computational Linguistics, Online, 1971–1981.
- [31] Nayeon Lee, Belinda Li, Sinong Wang, Wen-tau Yih, Hao Ma, and Madian Khabsa. 2020. Language Models as Fact Checkers?. In Proceedings of the Third Workshop on Fact Extraction and VERification (FEVER '20). Association for Computational Linguistics, Online, 36–41.
- [32] Yaliang Li, Jing Gao, Chuishi Meng, Qi Li, Lu Su, Bo Zhao, Wei Fan, and Jiawei Han. 2016. A Survey on Truth Discovery. SIGKDD Explor. Newsl. 17, 2 (2016), 1–16.
- [33] Christopher Malon. 2018. Team Papelo: Transformer Networks at FEVER. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium, 109–113.
- [34] Tsvetomila Mihaylova, Georgi Karadzhov, Pepa Atanasova, Ramy Baly, Mitra Mohtarami, and Preslav Nakov. 2019. SemEval-2019 Task 8: Fact Checking in Community Question Answering Forums. In Proceedings of the 13th International Workshop on Semantic Evaluation (SemEval '19). Association for Computational Linguistics, Minneapolis, Minnesota, USA, 860–869.
- [35] Moin Nadeem, Wei Fang, Brian Xu, Mitra Mohtarami, and James Glass. 2019. FAKTA: An Automatic End-to-End Fact Checking System. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations) (NAACL-HLT '19). Association for Computational Linguistics, Minneapolis, Minnesota, 78–83.
- [36] Preslav Nakov. 2003. Building an Inflectional Stemmer for Bulgarian. In Proceedings of the 4th International Conference on Computer Systems and Technologies (CompSysTech '03). Sofia, Bulgaria, 419–424.
- [37] Preslav Nakov. 2003. BulStem: Design and Evaluation of Inflectional Stemmer for Bulgarian. In Proceedings of the Workshop on Balkan Language Resources and Tools. Thessaloniki, Greece.

- [38] Preslav Nakov, Alberto Barrón-Cedeño, Tamer Elsayed, Reem Suwaileh, Lluís Màrquez, Wajdi Zaghouani, Pepa Atanasova, Spas Kyuchukov, and Giovanni Da San Martino. 2018. Overview of the CLEF-2018 CheckThat! Lab on Automatic Identification and Verification of Political Claims. In Proceedings of the Ninth International Conference of the CLEF Association: Experimental IR Meets Multilinguality, Multimodality, and Interaction (CLEF '18) (Lecture Notes in Computer Science). Springer, Avignon, France, 372–387.
- [39] Preslav Nakov, David Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021. Automated Fact-Checking for Assisting Human Fact-Checkers. In Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI '21). International Joint Conferences on Artificial Intelligence Organization, Online, 4551–4558.
- [40] Preslav Nakov, Giovanni Da San Martino, Tamer Elsayed, Alberto Barrón-Cedeño, Rubén Míguez, Shaden Shaar, Firoj Alam, Fatima Haouari, Maram Hasanain, Nikolay Babulkov, Alex Nikolov, Gautam Kishore Shahi, Julia Maria Struß, and Thomas Mandl. 2021. The CLEF-2021 CheckThat! Lab on Detecting Check-Worthy Claims, Previously Fact-Checked Claims, and Fake News. In Proceedings of the 43rd European Conference on Information Retrieval (ECIR '21). Springer, Lucca, Italy, 639-649.
- [41] Preslav Nakov, Giovanni Da San Martino, Tamer Elsayed, Alberto Barrón-Cedeño, Rubén Miguez, Shaden Shaar, Firoj Alam, Fatima Haouari, Maram Hasanain, Watheq Mansour, Bayan Hamdan, Zien Sheikh Ali, Nikolay Babulkov, Alex Nikolov, Gautam Kishore Shahi, Julia Maria Struß, Thomas Mandl, Mucahid Kutlu, and Yavuz Selim Kartal. 2021. Overview of the CLEF-2021 CheckThat! Lab on Detecting Check-Worthy Claims, Previously Fact-Checked Claims, and Fake News. In Proceedings of the 12th International Conference of the CLEF Association: Information Access Evaluation Meets Multiliguality, Multimodality, and Visualization (CLEF '21). Bucharest, Romania (online).
- [42] Preslav Nakov, Husrev Taha Sencar, Jisun An, and Haewoon Kwak. 2021. A Survey on Predicting the Factuality and the Bias of News Media. arXiv/2103.12506 (2021).
- [43] Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining Fact Extraction and Verification with Neural Semantic Matching Networks. In Proceedings of the AAAI Conference on Artificial Intelligence (AAAI '19). 6859–6866.
- [44] Wolfgang Otto. 2018. Team GESIS Cologne: An all in all sentence-based approach for FEVER. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium, 145– 149.
- [45] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language Models as Knowledge Bases?. 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 '19). Association for Computational Linguistics, Hong Kong, China, 2463–2473.
- [46] Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. 2020. AdapterHub: A Framework for Adapting Transformers. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP '20). Association for Computational Linguistics, Online, 46–54.
- [47] Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2018. CredEye: A Credibility Lens for Analyzing and Explaining Misinformation. In Companion Proceedings of the The Web Conference 2018 (WWW '18). International World Wide Web Conferences Steering Committee, Lyon, France, 155–158.
- [48] Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2019. STANCY: Stance Classification Based on Consistency Cues. 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 '19). Association for Computational Linguistics, Hong Kong, China, 6413– 6418.
- [49] M. Porter. 1980. An algorithm for suffix stripping. Program 40 (1980), 211–218.
- [50] Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. 2018. CatBoost: Unbiased Boosting with Categorical Features. In Proceedings of the 32nd International Conference on Neural Information Processing Systems (NIPS'18). Curran Associates Inc., Montréal, Canada, 6639-6649.
- [51] Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. 2017. Truth of Varying Shades: Analyzing Language in Fake News and Political Fact-Checking. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP '17). Association for Computational Linguistics, Copenhagen, Denmark, 2931–2937.
- [52] Aniketh Janardhan Reddy, Gil Rocha, and Diego Esteves. 2018. DeFactoNLP: Fact Verification using Entity Recognition, TFIDF Vector Comparison and Decomposable Attention. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium. 132–137.
- [53] B. Riedel, Isabelle Augenstein, Georgios P. Spithourakis, and S. Riedel. 2017. A simple but tough-to-beat baseline for the Fake News Challenge stance detection

- task. ArXiv abs/1707.03264 (2017).
- [54] Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How Much Knowledge Can You Pack Into the Parameters of a Language Model?. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP '20). Association for Computational Linguistics, Online, 5418–5426.
- [55] Shaden Shaar, Nikolay Babulkov, Giovanni Da San Martino, and Preslav Nakov. 2020. That is a Known Lie: Detecting Previously Fact-Checked Claims. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL '20). Association for Computational Linguistics, Online, 3607–3618.
- [56] Shaden Shaar, Fatima Haouari, Watheq Mansour, Maram Hasanain, Nikolay Babulkov, Firoj Alam, Giovanni Da San Martino, Tamer Elsayed, and Preslav Nakov. 2021. Overview of the CLEF-2021 CheckThat! Lab Task 2 on Detecting Previously Fact-Checked Claims in Tweets and Political Debates. In Working Notes of CLEF 2021—Conference and Labs of the Evaluation Forum (CLEF '21). CEUR-WS, Bucharest, Romania (online).
- [57] Shaden Shaar, Alex Nikolov, Nikolay Babulkov, Firoj Alam, Alberto Barrón-Cedeño, Tamer Elsayed, Maram Hasanain, Reem Suwaileh, Fatima Haouari, Giovanni Da San Martino, and Preslav Nakov. 2020. Overview of CheckThat! 2020 English: Automatic Identification and Verification of Claims in Social Media. In Working Notes of CLEF 2020—Conference and Labs of the Evaluation Forum (CLEF '20). CEUR-WS, Thessaloniki, Greece.
- [58] Chengcheng Shao, Giovanni Luca Ciampaglia, Alessandro Flammini, and Filippo Menczer. 2016. Hoaxy: A Platform for Tracking Online Misinformation. In Proceedings of the 25th International Conference Companion on World Wide Web (WWW '16 Companion). International World Wide Web Conferences Steering Committee, Montréal, Québec, Canada, 745–750.
- [59] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake News Detection on Social Media: A Data Mining Perspective. SIGKDD Explor. Newsl. 19, 1 (2017), 22–36.
- [60] Dominik Stammbach and Guenter Neumann. 2019. Team DOMLIN: Exploiting Evidence Enhancement for the FEVER Shared Task. In Proceedings of the Second Workshop on Fact Extraction and VERification (FEVER '19). Association for Computational Linguistics, Hong Kong, China, 105–109.
- [61] James Thorne and Andreas Vlachos. 2018. Automated Fact Checking: Task Formulations, Methods and Future Directions. In Proceedings of the 27th International Conference on Computational Linguistics (COLING '18). Association for Computational Linguistics, Santa Fe, New Mexico, USA, 3346–3359.
- [62] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a Large-scale Dataset for Fact Extraction and VERification. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT '18). Association for Computational Linguistics, New Orleans, Louisiana, USA, 809– 819.
- [63] James Thorne, Andreas Vlachos, Oana Cocarascu, Christos Christodoulopoulos, and Arpit Mittal. 2018. The Fact Extraction and VERification (FEVER) Shared Task. In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium, 1-9.
- [64] James Thorne, Andreas Vlachos, Oana Cocarascu, Christos Christodoulopoulos, and Arpit Mittal. 2019. The FEVER2.0 Shared Task. In Proceedings of the Second Workshop on Fact Extraction and VERification (FEVER '19). Association for Computational Linguistics, Hong Kong, China, 1–6.
- [65] Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. Science 359, 6380 (2018), 1146–1151.
- [66] William Yang Wang. 2017. "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL '17). Association for Computational Linguistics, Vancouver, Canada, 422–426.
- [67] Takuma Yoneda, Jeff Mitchell, Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. UCL Machine Reading Group: Four Factor Framework For Fact Finding (HexaF). In Proceedings of the First Workshop on Fact Extraction and VERification (FEVER '18). Association for Computational Linguistics, Brussels, Belgium, 97–102.
- [68] Savvas Zannettou, Michael Sirivianos, Jeremy Blackburn, and Nicolas Kourtellis. 2019. The Web of False Information: Rumors, Fake News, Hoaxes, Clickbait, and Various Other Shenanigans. J. Data and Information Quality 11, 3 (2019), 10:1–10:37.
- [69] Yifan Zhang, Giovanni Da San Martino, Alberto Barrón-Cedeño, Salvatore Romeo, Jisun An, Haewoon Kwak, Todor Staykovski, Israa Jaradat, Georgi Karadzhov, Ramy Baly, Kareem Darwish, James Glass, and Preslav Nakov. 2019. Tanbih: Get To Know What You Are Reading. 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 '19): System Demonstrations. Association for Computational Linguistics, Hong Kong, China, 223–228.
- [70] Arkaitz Zubiaga, Ahmet Aker, Kalina Bontcheva, Maria Liakata, and Rob Procter. 2018. Detection and Resolution of Rumours in Social Media: A Survey. ACM Comput. Surv. 51, 2, Article 32 (2018), 36 pages.