
USER PERCEPTIONS OF AUTOMATIC FAKE NEWS DETECTION: CAN ALGORITHMS FIGHT ONLINE MISINFORMATION?

Bruno Tafur
University of Cambridge
Cambridge
bt403@cantab.ac.uk

Advait Sarkar
University of Cambridge
Cambridge
advait@microsoft.com

ABSTRACT

Fake news detection algorithms apply machine learning to various news attributes and their relationships. However, their success is usually evaluated based on how the algorithm performs on a static benchmark, independent of real users. On the other hand, studies of user trust in fake news has identified relevant factors such as the user’s previous beliefs, the article format, and the source’s reputation. We present a user study (n=40) evaluating how warnings issued by fake news detection algorithms affect the user’s ability to detect misinformation. We find that such warnings strongly influence users’ perception of the truth, that even a moderately accurate classifier can improve overall user accuracy, and that users tend to be biased towards agreeing with the algorithm, even when it is incorrect.

Keywords fake news, machine learning, misinformation, user trust

1 Introduction

False information can be rapidly published and spread online, reaching massive audiences [1, 2, 3]. Social networks have catalyzed the spread of misinformation, affecting people worldwide. In response, research has focused on fake news detection using machine learning [4, 5, 6, 7, 8]. These algorithms typically rely on analysis of content, information source, or propagation patterns [5]. High accuracy models use collections of features [5, 9, 10].

Studies have measured users’ ability to detect fake news [11], and users’ change in trust with respect to changes in user and expert reputation ratings or changes in content structure and format [12, 13, 14, 15, 16]. Researchers have proposed behavioural “nudges”, causing users to distrust or critically evaluate news before sharing [17, 18]. Such nudges could be generated by a fake news detection algorithm, or reputation ratings.

Some studies have evaluated the use of warnings in news articles, generally focusing on a subset of fake or disputed articles. [16, 19, 20]. The focus is mostly on analysing the effect of textual information and the design of the articles when showing the warnings.

However, the user perspective and interaction with fake news detection algorithms is under-explored. What might it be like for a user to browse the web, a mixture of true and false information, with accuracy ratings issued by an algorithm? Moreover, what if the algorithm, as is usually the case, makes mistakes?

In this paper, we explore user interaction with fake news detection algorithms, by measuring how frequently users agree with algorithm warnings in different scenarios. We evaluate if an imperfect algorithm, on the whole, can help users better identify misinformation. We make the following contributions:

- We present an empirical study with 40 participants, demonstrating that using fake news detection algorithms can improve user ability to discern fake news from real news.
- We demonstrate that user *accuracy* when detecting fake and real news can vary depending on the truthfulness of the article and the algorithm’s correctness.

- We present evidence that users' *agreement* with a fake news algorithm is affected by its correctness, and that because users tend to agree with system advice, incorrect warnings may be worse than no warnings.
- We confirm previous results that suggest that users' accuracy in identifying fake and real news is relatively low [16, 14, 21]. Specifically, in our case, participants had 61% accuracy without algorithm support.

2 Related Work

Discussions of fake news are made complex due to its various definitions [1]. Related topics include misinformation and disinformation. Misinformation refers to distributed information that is false or inaccurate, while disinformation is the distribution of misinformation on purpose with a deceptive intent [2]. Fake news can include concepts such as rumours, satire and disinformation [1]. In our work, we make the simplifying assumption that fake news is articles with false factual statements, where falsity is determined by a label in our dataset (which in turn derives from the perceptions of data annotators, who attempt to reflect scientific consensus), regardless of the intent of the publisher.

2.1 Fake news detection algorithms

Fake news detection algorithms apply supervised machine learning methods, commonly based on content analysis, source analysis or propagation patterns [5]. Propagation-based models often use graph neural networks [22, 23, 24], analysing the patterns of news spread. Other models focus on content, classifying based on the style of the text or comparison with a factual knowledge base [23, 25, 4, 26].

The most effective models have used mixed approaches [9, 5, 10]. For example, Antoun et al. [5] use features from the headline, the body and the relationship with the top five searches in Google. Their tests employ transformer architectures, such as RoBERTa and XLNet, and reached an F1-score of up to 98% when detecting fake news with data from the Qatar International Cybersecurity Contest. The model developed by Shu et al. [10] achieved an 88% F1-score in Politifact data and 87% in Buzzfeed data. These studies focus on the performance of algorithms on static benchmarks, and not on the user experience.

2.2 User studies

Research has measured users' detection accuracy of fake news, taking into consideration factors such as news format, topic, structure and reputation [12, 13, 14, 15, 16]. There is a diversity of terminology and measurements: studies refer variously to user trust, believability, credibility or user engagement [11].

Kim & Dennis [12] analysed how changing the presentation format of an article affects users' trust and further actions on social media, finding that showing the source of an article increased scepticism and having low source ratings affected the article's believability. Similarly, Kim et al. [13] analysed the impact of user, source and expert reputation ratings on article credibility. They found that all three kinds of ratings influenced trust. Both studies found confirmation bias: users tend to believe more in articles aligned with their previous beliefs.

A study by Spezzano et al. [15] examined how having an image, title, excerpt, and source affected user's trust. The study then compared human-level detection to the performance of an algorithm. Human accuracy results showed a 62% detection accuracy when title and image were included versus 53% when they were just exposed to an extract, but their detection algorithm achieved 83% accuracy.

Another study analysed how credibility was affected by the topic of the news and by the number of likes on social media [14]. Their results showed an average user detection accuracy of 51%, and users had a higher detection accuracy with fake news than real news. The authors suggested that this difference could be due to a deception-bias, i.e., a presumption of deception by the users. They also found a relationship between Facebook likes and the higher credibility.

Some studies have evaluated applying warnings to articles [16, 19, 20, 27, 21] or informational checklists [28]. Kirchner & Reuter [16] tested warnings on fake news, finding that users prefer to be warned about fake news and that the warnings affected the perceived accuracy of false headlines. Similarly, Pennycook et al. [19] showed that applying warnings to a subset of fake news impacts the perceived accuracy of fake news headlines. However, the articles that did not show a warning were negatively affected as users tend to assume that the absence of a warning made the news article more truthful. This is interpreted as an "implied truth" effect and suggests warnings should be balanced between fake and real news. Also, Clayton et al. [20] found that using "Rated False" and "Disputed" tags reduced the belief in fake news, with "Rated False" having a more significant effect. Therefore, there could also be an impact generated by the type of warning used.

Seo et al. [27] evaluated effects of 3 types of warnings on participants' recognition, detection and sharing of fake news. They found that giving explanations inside the warning had a positive effect on the user's decision. The study did not cover possible scenarios of algorithm incorrectness and its impact on users' perception of the truth, and focused on adding warnings to disputed articles but not to non-disputed ones.

A close precedent to our work is Snijders et al. [21], who explore effects of individual confidence on trust in the advice of an algorithmic fake news detector, finding that participants are less likely to accept algorithmic advice for news stories about which they are themselves confident. Crucially, the algorithm used in their study was not fully accurate. Another study by Lu et al. [29] also analysed nudges by an AI algorithm and analysed the effect related to news spread. However, in both studies, their report does not differentially analyse user agreement in the cases where the algorithmic advice was correct versus erroneous.

In summary, previous research has developed accurate, yet imperfect algorithms for fake news detection. Studies of users have shown that cues such as format, topic, prior beliefs, source, author, images, etc. all contribute to the perceived credibility of a news article. Furthermore, studies have shown that warnings can help users evaluate misinformation. However, no previous user study has contrasted the full set of scenarios that may arise with an imperfect fake news detector in practice, including true information that has been incorrectly flagged as false, or false information that has been incorrectly flagged as true. Our study fills this gap.

3 Method

We designed a study to investigate the impact of an imperfect fake news detection algorithm on users' ability to detect fake and real news.

Participants. We recruited 40 English-speaking adults via the Amazon Mechanical Turk platform. We selected participants from English-speaking countries, who had carried out more than 1,000 surveys on the platform, and with an approval rate of more than 99%.

Dataset. We manually curated a dataset 40 news articles.¹ Twenty articles were compiled manually from recent publications in Snopes and Politifact, fact-checking websites that have been used in previous research to generate similar datasets [30, 31, 32]. A further 10 articles were picked from the MisInfoText dataset [31, 30] and a final 10 articles were selected from the Fake and Real dataset [33, 34]. Each article consisted of a headline and a single paragraph of content. We chose this short length to reduce reading time, so that participants could be exposed to a large number of articles during the study. We selected articles based on the following criteria:

- Articles needed to be easily established as true or false based on publicly available information from authoritative sources. We removed articles where we could not determine truth or falsity with high certainty. We acknowledge that this is a subjective step.
- Articles needed to be current at the time of the study.
- The dataset needed to be balanced, having the same number of true and false articles, and having a diverse range of topics that were roughly equally represented. We chose articles representing the following broad topics: food & health, politics, climate change, world news, and COVID-19.

Protocol. We adopted a counterbalanced, within-subjects design. All participants were shown 40 articles, and asked to rate them as fake or real. Each user rated half of the articles with algorithm support (with warnings) and half of the articles without warnings. The articles with algorithm support displayed a warning based on the classification output of a (hypothetical) fake news detection algorithm with imperfect accuracy. Examples of the articles with and without warnings, and an example of a full question shown in the survey, can be seen in Figure 1. Sample headlines for each topic can be seen in Table 2.

To counterbalance the stimuli used in the with and without-warnings conditions, the 40 articles were divided into two batches of twenty articles, containing 10 real and 10 fake articles each. The 40 participants were randomly divided into two equal groups. Group A of the participants was shown Batch 1 without warnings and Batch 2 with the algorithm detection warnings. Group B was exposed to Batch 1 with warnings and Batch 2 without warnings. This is illustrated in Table 1. For each participant, the items within Batch 1 and Batch 2 were randomly shuffled, to mitigate order effects.

In practice, fake news detection algorithms report accuracies between 60% and 90% [10, 23]. We therefore chose to make our simulated algorithm 70% accurate, aligning with previous work [21]. Thus, in 30% of the cases, the algorithm

¹Archived on the Open Science Framework at <https://osf.io/bjn2q/>

Table 1: Participant groups and condition assignments

Participant group A	Participant group B
Batch 1, without warnings	Batch 1, with warnings
Batch 2, with warnings	Batch 2, without warnings



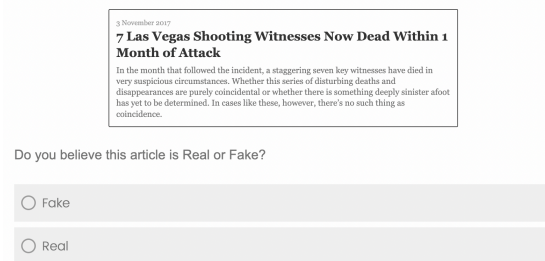
(a) Without warning



(b) With fake article warning



(c) With real article warning



(d) Sample question shown in the survey (no warning shown)

Figure 1: Sample news articles shown in the survey

either incorrectly flagged a true article as false, or a false article as true. In practice we randomly selected 3 (out of 10) real and 3 fake stories per batch to show incorrect warnings for; each participant thus encountered 14 correct and 6 incorrect warnings in the with-warnings condition. This falls in the range of possible accuracies of a real algorithm, and allows sufficient instances of misclassification to observe interesting data about user responses in that range.

The study did not collect any personal or identifiable information, and was approved by our institutional ethics committee.

4 Results

We define user *accuracy* as the percentage of articles correctly classified by the user. We define user *agreement* with the algorithm as the percentage of articles where the user rating (real or fake) matches the algorithm rating.

4.1 User accuracy

4.1.1 Effects of algorithm warnings

The mean user accuracy with and without warnings is shown in Table 3. Accuracy increased from 60.5% without warnings to 68% with warnings. This difference is a medium effect size (Cohen’s $D = 0.65$) and is statistically significant (t-test, $p = 7.78 \cdot 10^{-12}$). The distributions can be seen in Figure 2.

4.1.2 Effects of truth, falsity, and algorithm correctness

Table 3 breaks down user accuracy by fake and real articles. The average user accuracy for detecting fake news is 69%, while for real news it is 59%. This difference is statistically significant (Wilcoxon signed-rank test, $z = 2.75, p = 0.006$. A non-parametric test is used here as the data are not normally distributed).

Table 2: Sample article headlines per topic.

Category	Fake	True
Food & Health	FBI Issues Horrifying Warning to Frequent Grocery Shoppers	“Rattlesnake selfie” results in a \$153K medical bill
Politics	Trump Just Got Banned From The Place He Proposed To Melania	Trump tweets cartoon of train hitting CNN reporter
Climate Change	Facebook Spamming Climate Posts with “Climate Science Center” Propaganda	Amid higher global temperatures, sea ice at record lows at poles
World	Billion Dollar Company Tells Employees How They’re Allowed To React On Social Media To 46% Pay Cuts	Nigeria says U.S. agrees delayed \$593 million fighter plane sale
COVID	1000 Peer Reviewed Studies Questioning Covid-19 Vaccine Safety	Trump booed after revealing he got a Covid booster shot

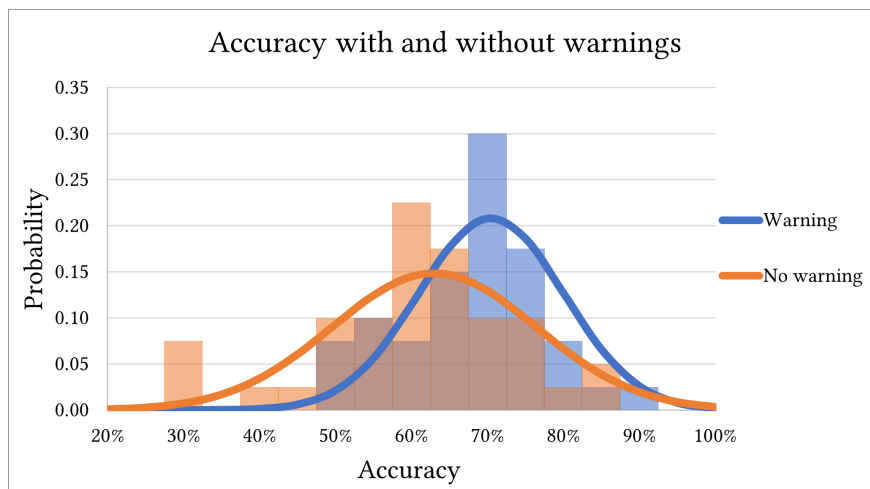


Figure 2: Normalised distribution of user accuracy with and without algorithm warnings. Accuracy improves with warnings.

Warnings in Table 3 are subdivided into correct and incorrect. User accuracy with correct warnings was considerably higher than with incorrect warnings, going from 42.5% to 78.9%. This difference is statistically significant (Wilcoxon signed-rank test, $z = 4.14, p = 4 \cdot 10^{-5}$).

4.1.3 Effect of news topics

Figure 3 (left) shows user accuracy by topic, restricted just trials where there was no warning, or a correct warning (our sample was not large enough to investigate incorrect warnings by topics, which we leave for future work). The difference between topics is statistically significant (Friedman test, $p = 0.007$), but this topic needs further investigation as the number of news articles in each topic is relatively small. We include this preliminary analysis here as a starting point for such investigation.

4.2 User agreement with the algorithm

We examined the frequency with which users agreed with the algorithm ratings, whether fake or real. On average, the users agreed with the algorithm decision 72.5% of the time.

Table 4 summarizes user agreement with the warnings, when the underlying article is fake or real, and the warning is correct or wrong. Users agreed with the algorithm 74.8% of the time for fake articles and 70.3% of the time for real articles, but this difference is not statistically significant. Nonetheless, this can be loosely interpreted as the users showing slightly more trust in the algorithm when the warnings were related to fake news.

Table 3: User accuracy for fake and real news, with and without algorithm warnings

	Acc. Fake (σ)	Acc. Real (σ)	Mean Acc. (σ)
With warning	75.3% ($\pm 13.5\%$)	60.8% ($\pm 15.2\%$)	68.0% ($\pm 9.5\%$)
Correct warning	85.7%	72.1%	78.9%
Incorrect warning	50.8%	34.2%	42.5%
Without warning	63.0% ($\pm 20.5\%$)	58.0% ($\pm 17.2\%$)	60.5% ($\pm 13.2\%$)
Overall mean	69.1% ($\pm 15.1\%$)	59.4% ($\pm 12.5\%$)	

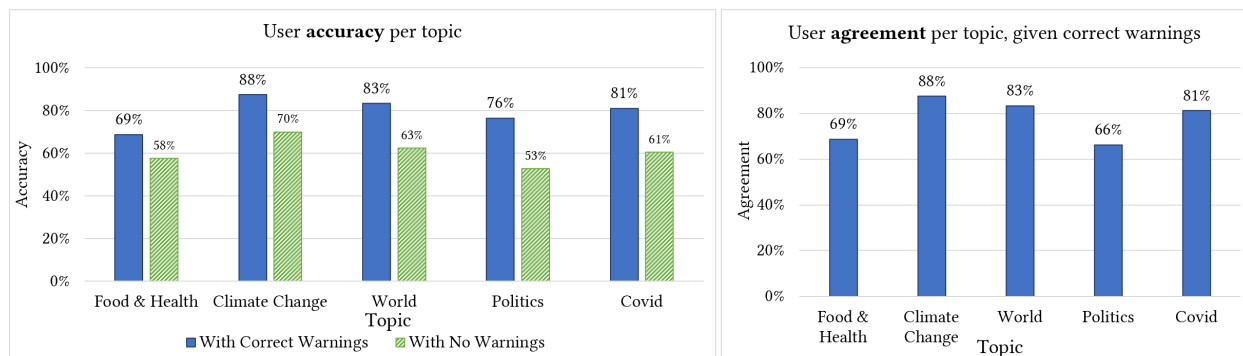


Figure 3: Left: Comparing user accuracy per topic with correct warnings, versus without warnings. Right: user agreement with algorithm per topic, with correct warnings.

The users agreed with the algorithm recommendations 78.9% of the time when the algorithm was correct, and the users agreed with the algorithm 57.8% of the time when it was incorrect. This difference is statistically significant (Wilcoxon signed-rank test, $z = 4.17$, $p = 10 \cdot 10^{-5}$). Thus, users are considerably more likely to agree with the algorithm when it is giving the correct advice. Nonetheless, even when the algorithm is wrong, users agree with the recommendations more than half the time.

4.2.1 Effect of news topics

Finally, the percentage agreement in specific news topics can be seen in Figure 3 (right). The difference between topics was found to be significant (Friedman test, $p < 0.01$). Therefore, it is possible that users are more likely to follow more the algorithm in certain topics than in others. However, similar to our analysis of the topics in user accuracy, these results may require further studies with a greater variety of news articles per topic.

5 Limitations

Participants recruited through Amazon Mechanical Turk have known limitations, such as perverse incentives to complete the task as quickly as possible without regard to correctness. We attempted to mitigate these by setting a high threshold for ratings when selecting participants, and instituting time checks and analysis of response distributions to ensure that participants were not providing random responses.

Table 4: Fake vs. real news user agreement with and without algorithm warnings

	Agreement Fake (σ)	Agreement Real (σ)	Mean Agreement (σ)
Correct warnings	85.7%	72.1%	78.9% ($\pm 15.7\%$)
Incorrect warnings	49.2%	65.8%	57.8% ($\pm 31.3\%$)
Mean Agreement	74.8% ($\pm 18.0\%$)	70.3% ($\pm 23.1\%$)	72.5% ($\pm 18.1\%$)

Our analysis of topics was post-hoc; we did not intend to compare differences between topics from the outset, thus our dataset was too small to include many instances of each topic. Future work could continue along the lines suggested in our preliminary analysis by expanding the dataset greatly, to create representative sets for different topics.

This study did not cover getting qualitative insights from the users. This feedback could be important to understand users' opinions of the system warnings and the nature of their mental processes when critically evaluating the truthfulness of an article in the context of a system warning (which may be wrong). Future studies could gather more qualitative data and contrast them with quantitative findings.

Some other variables could be modified as well to analyse their impact. We fixed the algorithm accuracy of 70%; future studies could test the impact of different levels of accuracy. We did not vary the formatting of the articles or warnings, but future studies might explore showing warnings with different text messages. Possible changes to the text messages include stating different verbs for the warning "Our system believes...", for example, "classifies" or "thinks", or changing the wording overall. Article attributes could also be changed to give the users more information, such as showing the news source in the extract. Future studies may explore the interaction of these presentational factors with algorithm warnings.

6 Discussion

We find that a fake news detection algorithm, even one that is only 70% accurate, can help users distinguish fake and real news. Overall, showing warnings helped improve user accuracy by 7.5 percentage points, and this improvement was statistically significant.

User accuracy in detecting fake news was greater than detecting real news. This may be deception bias, which has also been found in other studies [14], where users tend to be more skeptical when their attention is drawn to the possibility of encountering false articles. Fake news article headlines tend to be less subtle, which might make them easier to identify.

In contrast to Snijders et al., [21], our study measured a larger improvement in accuracy when using algorithmic advice. This is likely due to multiple methodological differences. Our design has greater ecological validity: our participants did not choose when to receive advice, they were not shown the ground truth in order to help calibrate trust in the model, and they did not see the same article repeatedly and therefore could not retroactively update their judgement after having seen the ground truth.

At the time of this study, we inhabit an information culture where most web users are not consciously critical of the truth of articles in their day-to-day reading. To increase the ecological validity of our findings, a future study could observe people's behaviour in a more naturalistic, day-to-day setting. For example, this might take the form of a longitudinal diary study, complemented with logs of their interactions in a social media platform or news website.

Algorithm correctness was also a significant factor. Users agreed with the algorithm's advice 72.5% of the time. Even when the algorithm gave wrong advice, users agreed with their recommendations more than half the time. When the system was correct, users had higher accuracy than when the system was wrong; possibly indicating an unwillingness to go against the algorithm's decision, which can be viewed as "overtrust" or "inappropriate trust" [35]. For example, four users followed the algorithm's decision 100% of the time, placing total trust in the algorithm. These results suggest that it would be better not to give a warning, than to give an incorrect warning. At a high level of trust, increasing the algorithm's accuracy is also likely to increase users' accuracy.

7 Conclusion

Even though the performance of fake news detection algorithms has improved in recent research, and they are essential tools in the fight against misinformation, they are still not (and may never be) fully accurate. It is therefore necessary to evaluate the effects of using an imperfect algorithm on users' trust, to test whether such systems are worth using despite their possibility for error.

In our study, the user accuracy in classifying news as real or fake increased by 7.5 percentage points when assisted by an algorithm that itself was only 70% accurate. This accuracy can be affected by other factors such as the nature of the article, the algorithm correctness and the news topics.

We further found that user agreement with the algorithm was 78.9% when the algorithm was correct, and 57.8% when it was wrong. A high level of user trust in the algorithm, even when incorrect, indicates that no warnings may be better than incorrect warnings, and that any improvement in the accuracy of the system should translate directly into improvement in the user's ability to detect fake news.

References

- [1] Xinyi Zhou and Reza Zafarani. A Survey of Fake News: Fundamental Theories, Detection Methods, and Opportunities. *ACM Computing Surveys*, 53(5):109:1–109:40, September 2020.
- [2] James Thorne and Andreas Vlachos. Automated Fact Checking: Task Formulations, Methods and Future Directions. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3346–3359, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics.
- [3] Jana Lasser, Segun Taofeek Aroyehun, Almog Simchon, Fabio Carrella, David Garcia, and Stephan Lewandowsky. Social media sharing of low quality news sources by political elites. *PNAS Nexus*, 2022.
- [4] William Yang Wang. "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection. *arXiv:1705.00648 [cs]*, May 2017. arXiv: 1705.00648 version: 1.
- [5] Wissam Antoun, Fady Baly, Rim Achour, Amir Hussein, and Hazem Hajj. State of the Art Models for Fake News Detection Tasks. In *2020 IEEE International Conference on Informatics, IoT, and Enabling Technologies (ICIoT)*, pages 519–524, February 2020.
- [6] Robiert Sepúlveda-Torres, Marta Vicente, Estela Saquete, Elena Lloret, and Manuel Palomar. Exploring Summarization to Enhance Headline Stance Detection. In Elisabeth Métais, Farid Meziane, Helmut Horacek, and Epaminondas Kapetanios, editors, *Natural Language Processing and Information Systems*, Lecture Notes in Computer Science, pages 243–254, Cham, 2021. Springer International Publishing.
- [7] Valeria Mazzeo, Andrea Rapisarda, and Giovanni Giuffrida. Detection of Fake News on COVID-19 on Web Search Engines. *Frontiers in Physics*, 9, 2021.
- [8] Monther Aldwairi and Ali Alwahedi. Detecting Fake News in Social Media Networks. *Procedia Computer Science*, 141:215–222, January 2018.
- [9] Jiawei Zhang, Bowen Dong, and Philip S. Yu. FakeDetector: Effective Fake News Detection with Deep Diffusive Neural Network. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, pages 1826–1829, April 2020. ISSN: 2375-026X.
- [10] Kai Shu, Suhang Wang, and Huan Liu. Beyond News Contents: The Role of Social Context for Fake News Detection. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM '19*, pages 312–320, New York, NY, USA, January 2019. Association for Computing Machinery.
- [11] Victoria Vziatyshva, Yadviga Sinyavskaya, Alexander Porshnev, Maxim Terpilovskii, Sergey Koltcov, and Kirill Bryanov. Testing Users' Ability to Recognize Fake News in Three Countries. An Experimental Perspective. In Gabriele Meiselwitz, editor, *Social Computing and Social Media: Experience Design and Social Network Analysis*, Lecture Notes in Computer Science, pages 370–390, Cham, 2021. Springer International Publishing.
- [12] Antino Kim and Alan R. Dennis. Says Who? The Effects of Presentation Format and Source Rating on Fake News in Social Media. SSRN Scholarly Paper ID 2987866, Social Science Research Network, Rochester, NY, August 2018.
- [13] Antino Kim, Patricia Moravec, and Alan R. Dennis. Combating Fake News on Social Media with Source Ratings: The Effects of User and Expert Reputation Ratings. SSRN Scholarly Paper ID 3090355, Social Science Research Network, Rochester, NY, February 2019.
- [14] Mufan Luo, Jeffrey T. Hancock, and David M. Markowitz. Credibility Perceptions and Detection Accuracy of Fake News Headlines on Social Media: Effects of Truth-Bias and Endorsement Cues. *Communication Research*, May 2020. Publisher: SAGE Publications Inc.
- [15] Francesca Spezzano, Anu Shrestha, Jerry Alan Fails, and Brian W. Stone. That's fake news! reliability of news when provided title, image, source bias & full article. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):109:1–109:19, April 2021.
- [16] Jan Kirchner and Christian Reuter. Countering Fake News: A Comparison of Possible Solutions Regarding User Acceptance and Effectiveness. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW2):1–27, October 2020.
- [17] Loukas Konstantinou, Ana Caraban, and Evangelos Karapanos. Combating Misinformation Through Nudging. In David Lamas, Fernando Loizides, Lennart Nacke, Helen Petrie, Marco Winckler, and Panayiotis Zaphiris, editors, *Human-Computer Interaction – INTERACT 2019*, volume 11749, pages 630–634. Springer International Publishing, Cham, 2019. Series Title: Lecture Notes in Computer Science.
- [18] Farnaz Jahanbakhsh, Amy X. Zhang, Adam J. Berinsky, Gordon Pennycook, David G. Rand, and David R. Karger. Exploring lightweight interventions at posting time to reduce the sharing of misinformation on social media. 5:1–42.

- [19] Gordon Pennycook, Adam Bear, Evan T. Collins, and David G. Rand. The Implied Truth Effect: Attaching Warnings to a Subset of Fake News Headlines Increases Perceived Accuracy of Headlines Without Warnings. *Management Science*, 66(11):4944–4957, November 2020. Publisher: INFORMS.
- [20] Katherine Clayton, Spencer Blair, Jonathan A. Busam, Samuel Forstner, John Glance, Guy Green, Anna Kawata, Akhila Kovvuri, Jonathan Martin, Evan Morgan, Morgan Sandhu, Rachel Sang, Rachel Scholz-Bright, Austin T. Welch, Andrew G. Wolff, Amanda Zhou, and Brendan Nyhan. Real Solutions for Fake News? Measuring the Effectiveness of General Warnings and Fact-Check Tags in Reducing Belief in False Stories on Social Media. *Political Behavior*, 42(4):1073–1095, December 2020.
- [21] Chris Snijders, Rianne Conijn, Evie de Fouw, and Kilian van Berlo. Humans and algorithms detecting fake news: Effects of individual and contextual confidence on trust in algorithmic advice. *International Journal of Human-Computer Interaction*, pages 1–12, 2022.
- [22] Tian Bian, Xi Xiao, Tingyang Xu, Peilin Zhao, Wenbing Huang, Yu Rong, and Junzhou Huang. Rumor Detection on Social Media with Bi-Directional Graph Convolutional Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01):549–556, April 2020. Number: 01.
- [23] Yi Han, Shanika Karunasekera, and Christopher Leckie. Graph Neural Networks with Continual Learning for Fake News Detection from Social Media. *arXiv:2007.03316 [cs]*, August 2020. arXiv: 2007.03316.
- [24] Federico Monti, Fabrizio Frasca, Davide Eynard, Damon Mannion, and Michael M. Bronstein. Fake News Detection on Social Media using Geometric Deep Learning. *arXiv:1902.06673 [cs, stat]*, February 2019. arXiv: 1902.06673.
- [25] Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra Lefevre, and Rada Mihalcea. Automatic Detection of Fake News. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3391–3401, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics.
- [26] Sadia Afroz, Michael Brennan, and Rachel Greenstadt. Detecting Hoaxes, Frauds, and Deception in Writing Style Online. In *2012 IEEE Symposium on Security and Privacy*, pages 461–475, May 2012. ISSN: 2375-1207.
- [27] Haeseung Seo, Aiping Xiong, and Dongwon Lee. Trust It or Not: Effects of Machine-Learning Warnings in Helping Individuals Mitigate Misinformation. In *Proceedings of the 10th ACM Conference on Web Science, WebSci '19*, pages 265–274, New York, NY, USA, June 2019. Association for Computing Machinery.
- [28] Hendrik Heuer and Elena Leah Glassman. A comparative evaluation of interventions against misinformation: Augmenting the WHO checklist. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22, pages 1–21. Association for Computing Machinery.
- [29] Zhuoran Lu, Patrick Li, Weilong Wang, and Ming Yin. The effects of AI-based credibility indicators on the detection and spread of misinformation under social influence. 6:461:1–461:27.
- [30] Fatemeh Torabi Asr and Maite Taboada. Big Data and quality data for fake news and misinformation detection. *Big Data & Society*, 6(1):2053951719843310, January 2019. Publisher: SAGE Publications Ltd.
- [31] Fatemeh Torabi Asr and Maite Taboada. The Data Challenge in Misinformation Detection: Source Reputation vs. Content Veracity. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 10–15, Brussels, Belgium, November 2018. Association for Computational Linguistics.
- [32] Kai Shu, Deepak Mahudeswaran, Dongwon Lee, and Huan Liu. FakeNewsNet: A Data Repository with News Content, Social Context, and Spatiotemporal Information for Studying Fake News on Social Media. *Big Data*, 8(3):171–188, June 2020. Publisher: Mary Ann Liebert, Inc., publishers.
- [33] Hadeer Ahmed, Issa Traore, and Sherif Saad. Detecting opinion spams and fake news using text classification. *Security and Privacy*, page 15, 2018.
- [34] Hadeer Ahmed, Issa Traore, and Sherif Saad. Detection of online fake news using n-gram analysis and machine learning techniques. In *International conference on intelligent, secure, and dependable systems in distributed and cloud environments*, pages 127–138. Springer, 2017.
- [35] Fumeng Yang, Zhuanyi Huang, Jean Scholtz, and Dustin L Arendt. How do visual explanations foster end users' appropriate trust in machine learning? In *Proceedings of the 25th International Conference on Intelligent User Interfaces*, pages 189–201, 2020.