# Evaluation of Fake News Detection with Knowledge-Enhanced Language Models

**Chenxi Whitehouse, Tillman Weyde, Pranava Madhyastha, Nikos Komninos**

City, University of London {chenxi.whitehouse, t.e.weyde, pranava.madhyastha, nikos.komninos.1}@city.ac.uk

## Abstract

Recent advances in fake news detection have exploited the success of large-scale pre-trained language models (PLMs). The predominant state-of-the-art approaches are based on fine-tuning PLMs on labelled fake news datasets. However, large-scale PLMs are generally not trained on structured factual data and hence may not possess priors that are grounded in factually accurate knowledge. The use of existing knowledge bases (KBs) with rich human-curated factual information has thus the potential to make fake news detection more effective and robust. In this paper, we investigate the impact of knowledge integration into PLMs for fake news detection. We study several state-of-the-art approaches for knowledge integration, mostly using Wikidata as KB, on two popular fake news datasets - `LIAR`, a politics-based dataset, and `COVID-19`, a dataset of messages posted on social media relating to the COVID-19 pandemic. Our experiments show that knowledge-enhanced models can significantly improve fake news detection on `LIAR` where the KB is relevant and up-to-date. The mixed results on `COVID-19` highlight the reliance on stylistic features and the importance of domain specific and current KBs. The code is available at https://github.com/chenxwh/fake-news-detection.

## Introduction

The world is witnessing a growing epidemic of fake news, which includes misinformation, disinformation, rumours, hoaxes, and other forms of rapidly spread and factually inaccurate information (Sharma et al. 2019). Fake news has been observed to severely impact political processes because of the wide reach of social media (Allcott and Gentzkow 2017). Misinformation related to medical issues, such as the COVID-19 pandemic, can cost lives (O'Connor and Murphy 2020). Automated methods for fake news detection and mitigation are a critical yet technically challenging problem (Thorne and Vlachos 2018).

In this paper, we focus on content-based fake news detection: methods that assess the truthfulness of news items based only on the text without using metadata. State-of-the-art models for this task are driven by advances in large-scale pre-trained language models (PLMs) (e.g. Liu et al. 2019a; Kaliyar, Goswami, and Narang 2021), which are trained on

vast amounts of raw web-based text using self-supervised methods (Rogers, Kovaleva, and Rumshisky 2020). A major limitation of these models is the lack of explicit grounding to real world entities and relations, which makes it difficult to recover factual knowledge (Bender et al. 2021). On the other hand, knowledge bases (KBs) provide a rich source of structured and human-curated factual knowledge, often complementary to what is found in raw text. This has recently led to the development of KB-augmented language models. Fake news detection can particularly benefit from the integration of KBs, making such models less dependent and reliant on surface level linguistic features.

In this study, we empirically analyse the impact of recent state-of-the-art knowledge integration methods, which enhance PLMs with KBs, for content-based fake news detection tasks. We evaluate ERNIE (Zhang et al. 2019), Know-Bert (Peters et al. 2019), KEPLER (Wang et al. 2021b) and K-ADAPTER (Wang et al. 2021a) on two distinct publicly available datasets: `LIAR` (Wang 2017), a politically oriented dataset, and `COVID-19` (Patwa et al. 2021), a dataset related to the recent pandemic. We find that integrating knowledge can improve fake news detection accuracy, given that the knowledge bases are relevant and up-to-date. Our experiments are not designed to find new state-of-the-art models for these datasets, but to investigate the effect of knowledge base integration into PLMs.

Our contributions are as follows: we evaluate multiple KB integration methods for fake news detection, we investigate model and data aspects that can prevent KB integration from being effective or from being effectively measured, and we discuss the potential for real-world applications.

In the following sections, we present a brief overview of four state-of-the-art methods that integrate KBs with PLMs studied in this paper. We then introduce and compare the datasets, the experiments with different knowledge-enhanced models, and the effectiveness of entity linking. We discuss our findings with respect to the necessary conditions for KB integration to be effective and how to assess its effect in application scenarios. Finally, we discuss the challenges in fake news detection and promising future directions.

## Method

In this section, we introduce the models with KB integration, and describe the datasets and our experimental setup.

## Knowledge Integration for PLM

Standard deep learning models obtain information from predicting and classifying text as they are trained, but have no prior knowledge of, or interaction with, world knowledge. Although PLMs can effectively characterise linguistic patterns from text to generate high-quality context-aware representations, they are limited in their grasp of knowledge, concepts, and relations, which are essential for some Natural Language Processing (NLP) tasks, including assessing the truthfulness of a news item.

On the other hand, KBs like Wikidata (https://www.wikidata.org) and WordNet (Miller 1995) contain rich curated information about the world. Thus, they could greatly complement PLMs if effective integration methods were available. Several efforts have been made to integrate KBs into PLMs. In this paper, we study the following models:

**ERNIE** injects knowledge into BERT (Devlin et al. 2019) by pre-training a language model on both large corpora and KBs. It uses TAGME (Ferragina and Scaiella 2010) to link entities to Wikidata. TAMGE identifies entity mentions in input text and links them to associated entity embeddings, which are then fused into the corresponding positions of the text. The knowledge-based learning objective is to predict the correct token-entity alignment. ERNIE has enhanced performance over BERT in entity typing and relation classification (Zhang et al. 2019).

**KnowBert** incorporates KBs into BERT using a knowledge attention and re-contextualisation mechanism. It identifies entity spans in input text and incorporates an integrated entity linker in the model to retrieve entity embeddings from a KB. The entity linker is responsible for entity disambiguation, which considers 30 entity candidates and uses their weighted average embedding. Knowledge-enhanced entity-span representations are then re-contextualised with a word-to-entity attention technique. KnowBert has shown improvement over BERT in relationship extraction, entity typing and word sense disambiguation (Peters et al. 2019).

**KEPLER** integrates factual knowledge into PLMs by adding a knowledge embedding objective with the supervision from a KB, and optimising it jointly with language modelling objectives. KEPLER is trained to encode the entities from their contextual descriptions, which enhances the ability of PLMs to extract knowledge from text. By keeping the original structures of PLMs, KEPLER can be used in general downstream NLP tasks without additional inference overhead. It is shown that KEPLER improves performance over RoBERTa (Liu et al. 2019b) in relationship extraction, entity typing and link prediction (Wang et al. 2021b).

**K-ADAPTER** retains the PLMs unchanged, but adds learnable adapter features that are trained in a multi-task setting on relation prediction and dependency-tree prediction. Two kinds of knowledge adapters have been developed by Wang et al. (2021a): factual knowledge obtained from automatically aligned text triples on Wikipedia and Wikidata, and linguistic knowledge obtained via dependency parsing. Both have been found to improve relation classification, entity typing and question answering (Wang et al. 2021a).

## Datasets

In our experiments, we use `LIAR` and `COVID-19` to study fake news detection. They both consist of short statements, but with different content, time of collection, linguistic and stylistic features.

**LIAR** was collected in 2017 from Politifact (https://www.politifact.com). It includes 12.8k human-labelled short statements about US politics from various contexts, i.e. news releases, TV interviews, campaign speeches, etc. Each statement has been rated for truthfulness by a Politifact editor using a six grade scale: "pants-fire", "false", "barely-true", "half-true", "mostly true", and "true". `LIAR` also provides metadata (e.g. speaker, context), which we do not use in our experiments. While Wang (2017) has been widely cited, we only found three other results for our specific task (no metadata, six classes) (Alhindi, Petridis, and Muresan 2018; Liu et al. 2019a; Chernyavskiy and Ilovsky 2020), the latter has the current best accuracy of 34.5%.

**COVID-19** was collected in 2020 after the COVID-19 outbreak. It consists of 10.5k posts related to the pandemic which are obtained from different social-media sites including Twitter, Facebook, and Instagram. The fake posts were collected from various fact-checking websites, i.e. Politifact and NewsChecker (https://newschecker.in), and the real posts were from Twitter using verified Twitter handles. Each post has a label, "real" or "fake". It was used as a shared task in the CONSTRAINT 2021 workshop (Chakraborty 2021) with the best reported accuracy of 98.69%.

## Experimental Setup

We use an empirical approach to study the effect of knowledge integration on fake news detection, to understand how knowledge is used by the model, and to evaluate the quality of the entity linker to the KB.

ERNIE and KnowBert are built on BERT-base, whereas KEPLER and K-ADAPTER are enhanced from RoBERTa-base and RoBERTa-large, respectively. We follow the concept of an ablation study to investigate the influence of the external knowledge by comparing the performance of each knowledge-enhanced PLM with the corresponding baseline model. We note that ERNIE and KnowBert incorporate entity embeddings that are linked to the input. Therefore we visualise the entities linked that contribute to the fake news detection task in ERNIE, and design experiments to investigate the impact of entity disambiguation of KnowBert.

We evaluate the performance of the models on fake news detection by fine-tuning the knowledge-enhanced PLMs on the training set with the same hyper-parameter settings. The input text is fed first to the PLM, and followed by a dropout ($p = 0.1$) and a linear layer. The output is then passed to a softmax layer for classification. We use AdamW optimiser (Loshchilov and Hutter 2019) (learning rate of $5 \times 10^{-6}$) and cross entropy as the loss function. Maximum input length is set to 128, and the batch size is 4. We train for 10 epochs and usually observe convergence after five. We perform five runs for each experiment and report the average accuracy with the standard deviation. Both `LIAR` and `COVID-19` are already

1426

(a) Word count per statement     (b) POS, punctuation, numbers in `LIAR`     (c) POS, punctuation, numbers, https in `COVID-19`
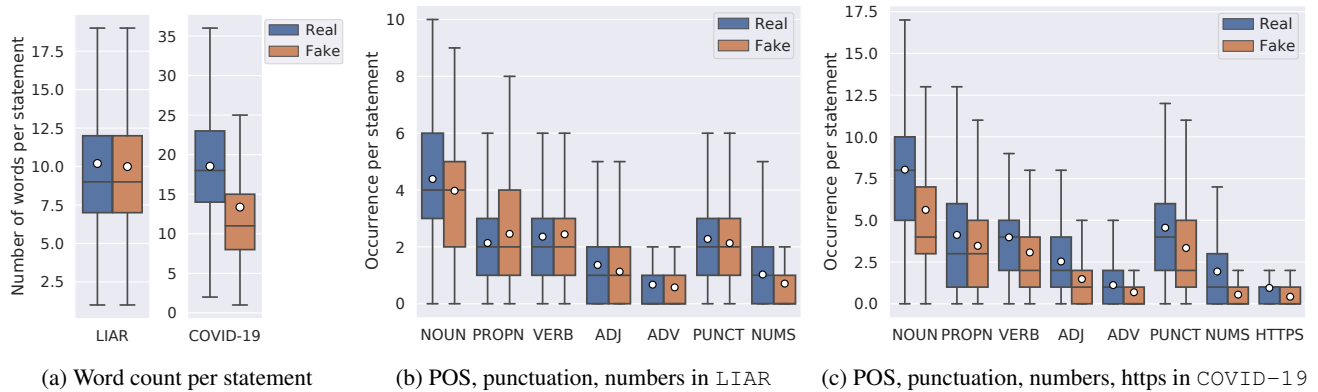
Figure 1: Number of words, POS tags, punctuation and numbers per statement in real and fake news in `LIAR` and `COVID-19`, and number of https-links per statement in `COVID-19`. The mean values are shown as white filled circles in the plot.

divided into train, validation, and test splits, which we use in our experiments as provided.

**Linguistic Feature Analysis** We also perform linguistic feature analysis following the work in Horne and Adali (2017) to investigate stylistic difference between real and fake news in the datasets. We use spaCy (https://spacy.io) to parse the statements and get the Part-of-Speech (POS) tags.

For `LIAR`, we group "pants-fire", "false", and "barely-true" as fake and "half-true", "mostly true", and "true" as real. We compare the distribution of different words, POS tags (NOUN, PROPN, VERB, ADJ, ADV), punctuation, and number-like words in each statement in Figure 1.

The length of posts is quite different between the two classes in `COVID-19`, with average of 32 and 22 for real and fake statements, respectively, as shown in Figure 1a. `LIAR`, on the other hand, has similar statement length, with 18 words per statement for real and 17 for fake.

In general, `COVID-19` has distinct linguistic features between classes whereas `LIAR` shows more similar features. In particular, `COVID-19` contains links, mostly https links, which are listed as a separate category in Figure 1c, showing a very skewed distribution.

## Experiments and Results

For our experiments, we use ERNIE, three pre-trained KnowBert models with different KBs (Wiki, WordNet, W+W), KEPLER, and K-ADAPTER with three adapters (F, L, F-L) in the published implementation, fine-tune the models to our task, and compare the result with the baseline models - BERT-base, RoBERTa-base, and RoBERTa-large.

**Detection Accuracy** The detection accuracy of the knowledge-enhanced PLMs and the corresponding baselines is shown in Table 1. On `LIAR`, all knowledge-enhanced methods improve over the baseline with KnowBert-W+W reaches the best overall result (improvement of $+2.59$ over BERT-base), whereas on `COVID-19`, only three of eight models show improvement, and only by a small margin.

The computational cost varies per approach. KEPLER retains the baseline PLM architecture, thus there is no

| MODEL | BASE | LIAR | COVID-19 |
|---|---|---|---|
| **B**ERT-**B**ase (BB) | - | $26.36_{\pm0.58}$ | $97.51_{\pm0.19}$ |
| **R**oBERTa-**B**ase (RB) | - | $26.71_{\pm0.93}$ | $97.61_{\pm0.26}$ |
| **R**oBERTa-**L**arge (RL) | - | $\mathbf{27.36}_{\pm0.79}$ | $\mathbf{97.92}_{\pm0.17}$ |
| ERNIE | BB | $27.53_{\pm0.13}$ | $97.30_{\pm0.18}$ |
| KnowBert-Wiki | BB | $27.64_{\pm0.09}$ | $97.37_{\pm0.09}$ |
| KEPLER | RB | $26.77_{\pm1.15}$ | $97.58_{\pm0.15}$ |
| K-ADAPTER-F | RL | $\mathbf{28.63}_{\pm0.90}{}^{*}$ | $\mathbf{97.92}_{\pm0.10}$ |
| KnowBert-WordNet | BB | $26.95_{\pm0.45}$ | $97.00_{\pm0.06}$ |
| KnowBert-W+W | BB | $\mathbf{28.95}_{\pm0.64}{}^{*}$ | $97.56_{\pm0.15}$ |
| K-ADAPTER-L | RL | $28.46_{\pm0.87}{}^{*}$ | $98.07_{\pm0.09}$ |
| K-ADAPTER-F-L | RL | $27.45_{\pm0.78}$ | $\mathbf{98.11}_{\pm0.14}$ |

Table 1: Detection accuracy results (average of five runs). The first section corresponds to the baseline models. Models in the second section use Wikidata KB. The third section shows models using other KBs and features. The best values within each section per dataset are marked in bold. The subscript numbers with $\pm$ show the standard deviation. Results with $*$ indicate statistically significant improvements over the baseline, both for mean (t-test, one-sided, $p < .05$) and median (Wilcoxon signed rank test, one-sided, $p < .05$).

overhead compared to RoBERTa-base. K-ADAPTER also freezes the RoBERTa-large layers, but there is an overhead of 9-23% from the adapters, while the overhead for Know-Bert is 40-87% and 111-131% for ERNIE.

**KB Linking** ERNIE and KnowBert create links between the text and KB entities at runtime and the quality of this linking influences the output. ERNIE uses TAGME and selects only one entity candidate per text span. In Figure 2 we show the 50 most frequently selected KB entities for each dataset. We can see that in `COVID-19`, the most frequent entities are not content-related ("https", "twitter") while "COVID-19", the most frequent relevant term in the dataset, is missing in the linked entities. For `LIAR`, on the other hand, the linked entities seem relevant. Since `LIAR` was collected three years earlier, it is apparently a better match for
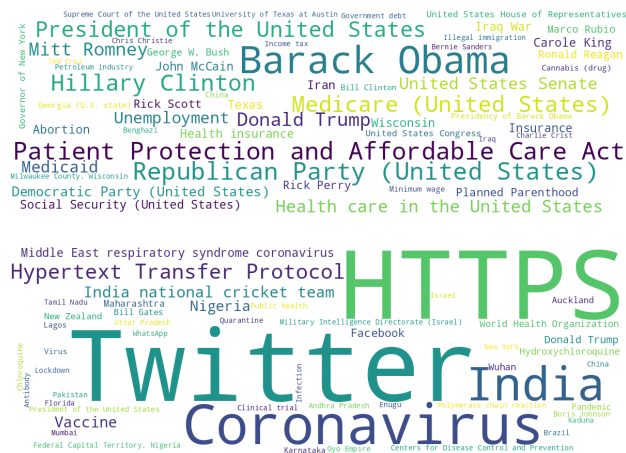
Figure 2: Word clouds for the 50 most frequent entities linked by ERNIE in `LIAR` (top) and `COVID-19` (bottom).

the entity linker and the KB used. Another potential influence on the effectiveness of KB integration is the number of linked entities. In contrast to ERNIE, KnowBert selects 30 most probable entities per text span. In a sensitivity study, we restrict KnowBert-W+W to only one entity, which reduces the accuracy on `LIAR` from 28.95% to 27.31%, below the accuracy of ERNIE (27.53%).

## Discussion

The reliable improvement of detection accuracy on `LIAR` by integrating PLMs with Wikidata shows the potential of knowledge integration exceeding the results obtained by integrating multiple types of metadata by Wang (2017). On the other hand, the improvements are good but not dramatic for `LIAR` and not consistent for `COVID-19`. We can identify two aspects contributing to the result which are relevant to the effective use of knowledge-enhanced models:

1) Currentness and relevance of the KB: as `COVID-19` was collected after most of the PLMs were trained, some terms such as "COVID-19" are not in the KB;

2) Quality of the dataset: the `COVID-19` dataset contains confounders that provide strong cues, overshadowing the impact of the knowledge base. The most important one is the occurrence of https links, which appear in 95.3% of the real posts but only 42.3% of the fake posts.

There is also potential to achieve more explainability and interpretability with direct KB integration at runtime. Take this statement from `COVID-19`: *"DNA Vaccine: injecting genetic material into the host so that host cells create proteins that are similar to those in the virus against which the host then creates antibodies"* as an example, KnowBert-W+W correctly classifies it as "real", whereas BERT-base fails. We observe most mention spans in the statement, i.e. *"DNA"*, *"injecting"*, *"genetic"*, *"genetic material"*, *"host"*, *"cells"*, etc. are correctly linked to entities *"DNA"*, *"Injection_(medicine)"*, *"Genetics"*, *"Genome"*, *"Host_(biology)"*, *"Cell_(biology)"*, respectively, therefore it seems that the entity links may have contributed to

KnowBert-W-W for this classification. However, the level of explainability is still limited.

**Application Aspects** Automatic fake news detection in practice adds two dynamic application aspects, which are difficult to test with static datasets as our experiment on `COVID-19` has shown:

(1) Dynamic adaptation: it is necessary to update the system to the changing characteristics of real and fake news (Silva and Almeida 2021). Knowledge-enhanced models that use KBs at runtime offer an opportunity to update the KB independent of the model. This has the advantage that fake news can be recognised as contradicting the KB before there are any fake news examples.

(2) Adversarial robustness: fake news authors are very likely to take evasive action. Adapting the text style is relatively easy and could be automated, which makes the detection with stylistic features difficult (see Zellers et al. 2019; Schuster et al. 2020).

Deployment of fake news detection in social media will also need human verification, e.g. when a user challenges actions taken against them. Here, KB integration can offer the advantage of insight into knowledge that has been used in the detection for better explainability.

## Related Work

In recent years large-scale PLMs i.e. BERT and RoBERTa have dominated NLP tasks, including some content-based fake news detection (Kaliyar, Goswami, and Narang 2021). Most fake news detection approaches either combine text with metadata (e.g. Ding, Hu, and Chang 2020), or focus only on the source of the text (e.g. Nørregaard, Horne, and Adali 2019). For `LIAR`, Alhindi, Petridis, and Muresan (2018) extend the data with evidence sentences in a new dataset called `LIAR-PLUS` to improve detection. Chernyavskiy and Ilvovsky (2020) introduce a Deep Averaging Network to model the discursive structure of the text and use Siamese models on the extended text data. Liu et al. (2019a) predict labels at two levels of granularity. For `COVID-19`, there are a number of results from the CONSTRAINTS 2021 workshop (Chakraborty 2021) which use a wide variety of traditional and neural NLP models. None of these approaches uses external knowledge, so they could all benefit from KB integration.

## Conclusion and Future Work

In this paper, we study the effectiveness of enhancing PLMs with knowledge bases for fake news detection. We find that integrating knowledge with PLMs can be beneficial on a static dataset but it depends on suitable KBs and the quality of the data. On the modelling level there are many routes for improvement. For practical application, more insight on what knowledge is used would be useful as well as dynamic adaptation of the models and the KBs. Integrating KBs with PLMs offers potentially more robust and timely fake news detection. However, a new evaluation approach, i.e. a testing scenario that models dynamic knowledge as well as adversarial and automatic fake news generators, is needed to assess the true potential of knowledge integration.

# References

Alhindi, T.; Petridis, S.; and Muresan, S. 2018. Where is Your Evidence: Improving Fact-checking by Justification Modeling. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, 85–90. Brussels, Belgium: ACL.

Allcott, H.; and Gentzkow, M. 2017. Social Media and Fake News in the 2016 Election. *Journal of Economic Perspectives*, 31(2): 211–36.

Bender, E. M.; Gebru, T.; McMillan-Major, A.; and Shmitchell, S. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623.

Chakraborty, T. 2021. *Combating Online Hostile Posts in Regional Languages during Emergency Situation: First International Workshop, CONSTRAINT 2021, Collocated with AAAI 2021, Virtual Event, February 8, 2021, Revised Selected Papers*. Springer Nature.

Chernyavskiy, A.; and Ilvovsky, D. 2020. Recursive Neural Text Classification Using Discourse Tree Structure for Argumentation Mining and Sentiment Analysis Tasks. In *ISMIS*, 90–101. Springer.

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*, 4171–4186. Minneapolis, Minnesota: ACL.

Ding, J.; Hu, Y.; and Chang, H. 2020. BERT-based Mental Model, a Better Fake News Detector. In *Proceedings of the 2020 6th international conference on computing and artificial intelligence*, 396–400.

Ferragina, P.; and Scaiella, U. 2010. TAGME: On-the-fly Annotation of Short Text Fragments (by Wikipedia Entities). *Proceedings of the 19th ACM international conference on Information and knowledge management*.

Horne, B.; and Adali, S. 2017. This Just In: Fake News Packs A Lot in Title, Uses Simpler, Repetitive Content in Text Body, More Similar to Satire Than Real News. *Proceedings of the International AAAI Conference on Web and Social Media*, 11(1): 759–766.

Kaliyar, R. K.; Goswami, A.; and Narang, P. 2021. Fake-BERT: Fake News Detection in Social Media with a BERT-based Deep Learning Approach. *Multimedia tools and applications*, 80(8): 11765–11788.

Liu, C.; Wu, X.; Yu, M.; Li, G.; Jiang, J.; qing Huang, W.; and Lu, X. 2019a. A Two-Stage Model Based on BERT for Short Fake News Detection. In *KSEM*.

Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019b. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR*, abs/1907.11692.

Loshchilov, I.; and Hutter, F. 2019. Decoupled Weight Decay Regularization. In *7th International Conference on Learning Representations, ICLR 2019*.

Miller, G. 1995. WordNet: A Lexical Database for English. *Commun. ACM*, 38: 39–41.

Nørregaard, J.; Horne, B. D.; and Adali, S. 2019. NELA-GT-2018: A Large Multi-Labelled News Dataset for the Study of Misinformation in News Articles. *Proceedings of the International AAAI Conference on Web and Social Media*, 13(01): 630–638.

O'Connor, C.; and Murphy, M. 2020. Going Viral: Doctors Must Tackle Fake News in the Covid-19 Pandemic. *Bmj*, 369(10.1136).

Patwa, P.; Sharma, S.; Pykl, S.; Guptha, V.; Kumari, G.; Akhtar, M. S.; Ekbal, A.; Das, A.; and Chakraborty, T. 2021. Fighting an Infodemic: COVID-19 Fake News Dataset. In *CONSTRAINT@AAAI*.

Peters, M. E.; Neumann, M.; Logan, R.; Schwartz, R.; Joshi, V.; Singh, S.; and Smith, N. A. 2019. Knowledge Enhanced Contextual Word Representations. In *EMNLP/IJCNLP*, 6086–6093. Hong Kong, China: ACL.

Rogers, A.; Kovaleva, O.; and Rumshisky, A. 2020. A Primer in BERTology: What We Know about How BERT Works. *Transactions of the ACL*, 8: 842–866.

Schuster, T.; Schuster, R.; Shah, D. J.; and Barzilay, R. 2020. The Limitations of Stylometry for Detecting Machine-Generated Fake News. *Computational Linguistics*, 46(2): 499–510.

Sharma, K.; Qian, F.; Jiang, H.; Ruchansky, N.; Zhang, M.; and Liu, Y. 2019. Combating Fake News. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10: 1 – 42.

Silva, R. M.; and Almeida, T. A. 2021. How Concept Drift can Impair the Classification of Fake News. In *Anais do IX Symposium on Knowledge Discovery, Mining and Learning*, 121–128. SBC.

Thorne, J.; and Vlachos, A. 2018. Automated Fact Checking: Task Formulations, Methods and Future Directions. In *Proceedings of the 27th International Conference on Computational Linguistics*, 3346–3359. Santa Fe, New Mexico, USA: ACL.

Wang, R.; Tang, D.; Duan, N.; Wei, Z.; Huang, X.; Ji, J.; Cao, G.; Jiang, D.; and Zhou, M. 2021a. K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 1405–1418. Online: ACL.

Wang, W. Y. 2017. "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection. In *Proceedings of the 55th Annual Meeting of the ACL (Volume 2: Short Papers)*, 422–426. Vancouver, Canada: ACL.

Wang, X.; Gao, T.; Zhu, Z.; Zhang, Z.; Liu, Z.; Li, J.; and Tang, J. 2021b. KEPLER: A Unified Model for Knowledge Embedding and Pre-trained Language Representation. *Trans. Assoc. Comput. Linguistics*, 9: 176–194.

Zellers, R.; Holtzman, A.; Rashkin, H.; Bisk, Y.; Farhadi, A.; Roesner, F.; and Choi, Y. 2019. Defending Against Neural Fake News. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Zhang, Z.; Han, X.; Liu, Z.; Jiang, X.; Sun, M.; and Liu, Q. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1441–1451. Florence, Italy: ACL.