

Multimodal Fusion with BERT and Attention Mechanism for Fake News Detection

Nguyen Manh Duc Tuan*
Toyo University
Tokyo, Japan
ductuan024@gmail.com

Pham Quang Nhat Minh
Aimesoft JSC
Hanoi, Vietnam
minhpham@aimesoft.com

Abstract—Fake news detection is an important task for increasing the reliability of the information on the internet since fake news is spreading fast on social media and has a negative effect on our society. In this paper, we present a novel method for detecting fake news by fusing multi-modal features derived from textual and visual data. Specifically, we proposed a scaled dot-product attention mechanism to capture the relationship between text features extracted by a pre-trained BERT model and visual features extracted by a pre-trained VGG-19 model. Experimental results showed that our method improved against the current state-of-the-art method on a public Twitter dataset by 3.1% accuracy.

Index Terms—Fake News Detection, Multimodal, BERT, Attention Mechanism

I. INTRODUCTION

Recently, fake news detection has received much attention in both NLP and the data mining research community. Fake news can negatively affect society, and spreads like a real virus, especially via social media. An example of how fake news can affect us is the 2016 and 2020 U.S. presidential elections where tons of fake news were spread on Facebook and Twitter. That makes fake news detection a crucial task.

Our work is inspired by the idea that different modalities can show different aspects of news and they complement each other in checking the reliability of information [20]. Fake news is usually created by both images and texts. In fake news, images may not be related to the content of the post [3]. Figure 1 is an image of two kids and they may be siblings, but the image is neither about *Vietnamese people* nor the *earthquake in Nepal in 2015*.

In this paper, we present a novel multimodal model for identifying fake news. We use neural networks to obtain feature representations from different modalities. Multimodal features are fused using the attention mechanism and put into a sigmoid layer for classification. Specifically, we use the BERTweet model [17] to obtain feature representations from texts and a pre-trained VGG-19 model to obtain feature representations from images. We propose a scaled dot-product



Fig. 1: An example of an image used in a fake news about two Vietnamese siblings at the 2015 Nepal earthquake.

attention mechanism on both texts and images, and also a self-attention mechanism on images because we see that in non-photoshopped images, all parts of images are related. Textual representations and visual representations along with three attention outputs are combined to improve the accuracy of fake news detection. Our proposed model obtained 80.8% of accuracy and 80% of F1-score.

The remainder of the paper is structured as follows. In Section II, we discuss the existing work on fake news detection, focusing on multimodal approaches. In Section III, we present our proposed method for fake news detection and in Section IV, we describe the dataset that we used in experiments. Section V gives our experimental results and results analysis. Finally, in Section VI, we conclude the paper and give some remarks.

II. RELATED WORK

The majority of previous work on fake news detection used text and user metadata features. Specifically, textual features can be obtained by applying convolutional neural networks (CNN) [12], [22]. Conneau et al. [6] has shown that a deep stack of local operations can help a model to learn the high-level hierarchical representation of a sentence and that increasing the depth leads to the performance improvement. Furthermore, deeper CNNs with residual connections can help

*This work was done when the author was an internship student at Aimesoft JSC

to avoid over-fitting and solve the vanishing gradient problem [12]. Textual features can also be manually designed from word clues, patterns, or other linguistic features of texts such as their writing styles [10], [24], [26]. We can also analyze unreliable news based on the sentiment analysis [23].

Fake news can be detected by analyzing social network information including user-based features and network-based features. User-based features were extracted from user profiles in [9], [15], [18]. For example, the number of followers, the number of friends, and registration ages are useful features to determine the credibility of a user post [4]. Network-based features can be extracted from the propagation of posts or tweets on graphs [16], [28].

Recent studies have shown that visual features and correlations between modalities are useful factors in detecting fake news [11], [13], [20], [25]–[27]. Many of them have used word embeddings, which is low-level embeddings for getting textual representations. Different from those, we applied the BERTweet model for extracting sentence embeddings.

Jin et al. [11] created an end-to-end network called att-RNN. It used attention mechanisms to combine text, image, and social-related features. In that network, texts and social contexts were concatenated together and passed into an LSTM network. Image features were extracted from a pre-trained VGG-19 model and outputs of the LSTM are used with an attention mechanism for fusing the visual features. Finally, the average of outputs of the LSTM was joint with visual features for classification.

Wang et al. [24] introduced a fake news detection model that makes use of multimodal features in adversarial neural networks. That model contains two major parts: event discriminator and fake news detector. For an input text, they used word embedding vectors as input and used CNN layers on it to generate the textual representation. The representation for an image was extracted using a pre-trained VGG-19 model. Finally, both textual and visual representations were concatenated and passed into two fully connected layers on top of the above two major parts.

Khattar et al. [13] built a similar architecture called Multimodal Variational Autoencoder for Fake News Detection (MVAE). Similar to EANN, they used word-embedding vectors, but instead of using CNN, bi-directional LSTMs were used to extract textual features. Visual features were extracted by VGG-19. These two vectors were concatenated to create the latent vectors which were passed into a decoder for reconstructing original samples. The latent vectors are then passed into two fully connected layers for fake news detection.

Singhal et al. [20] created a survey to show the essence of texts and images in fake news detection and they built an architecture called Spotfake. It uses BERT to obtain text feature representations and a pre-trained VGG-19 network to obtain feature representations from images. Then two single models were combined to make the prediction. While that method obtained good results, it did not make use of the correlation between modalities, which is useful for fake news detection. In the survey, Singhal et al. [20] showed that 81.4%

of the people, who took the survey, could distinguish fake and real news when both images and texts are given, while the number is 38.4% if only texts are given and 32.6% if only images are given. That indicated that multiple modalities provide more information and useful for fake news detection.

III. METHODOLOGY

In this section, we describe our approach to fake news detection. As shown in Figure 2 our proposed model contains four parts. The first part is the textual feature extractor in which we use BERT [8] and CNN layers to extract the contextual text features. The second part of the model is the visual feature extractor that extracts the visual features from a post. The third part is the common feature extractor in which we proposed an attention mechanism to extract the features from both texts and images. Finally, the last part is a multiple feature combination component that merges the representations derived from different components to obtain the feature representation for the entire post.

A. Textual Feature Extractor

Before extracting textual features, we performed the following pre-processing steps on texts.

- We converted words and tokens that have been lengthened into short forms. For example, “Cooooool” into “Cool”.
- There are some emojis written in text format, such as “:)”, “:(”, etc. We changed those emojis into sentiment words “happy” or “sad”.
- We did tokenization, word normalization, word segmentation with ekphrasis [2], a text analysis tool for social media.

We use BERT, the state-of-the-art model in many NLP tasks, to extract the feature representation of a tweet. BERT model has been shown to be effective in many NLP tasks, including text classification. In this task, We use BERTweet [17], a BERT model pre-trained on Tweet data. Cheema et al. [5] and Devlin et al. [8] have suggested that different hidden layers of BERT can capture different kinds of semantic information of the text, and the last four hidden layers of BERT are good for extracting information in a feature-based approach. Thus, we concatenate the last 4 hidden layers as contextual embeddings of tokens. After that, we use 1D-CNN layers [14] with filter sizes 2, 3, 4, and 5 to extract more information from different sets of word vectors for prediction. After passing the embedding vectors through 1D-CNN layers, we stack those output vectors vertically and pass them into two additional 1D-CNN layers with residual connections and obtain the output T_m in which the number of filters in each CNN layer is $d_T = 768$. Finally, we flatten T_m , pass it through two fully connected layers and obtain $d_T = 32$ as the final vector size of the textual representation T_f .

B. Visual Feature Extractor

We use VGG-19 model pre-trained on ImageNet dataset [19] for visual feature extraction. We take the output of the pre-trained VGG-19’s second last layer and

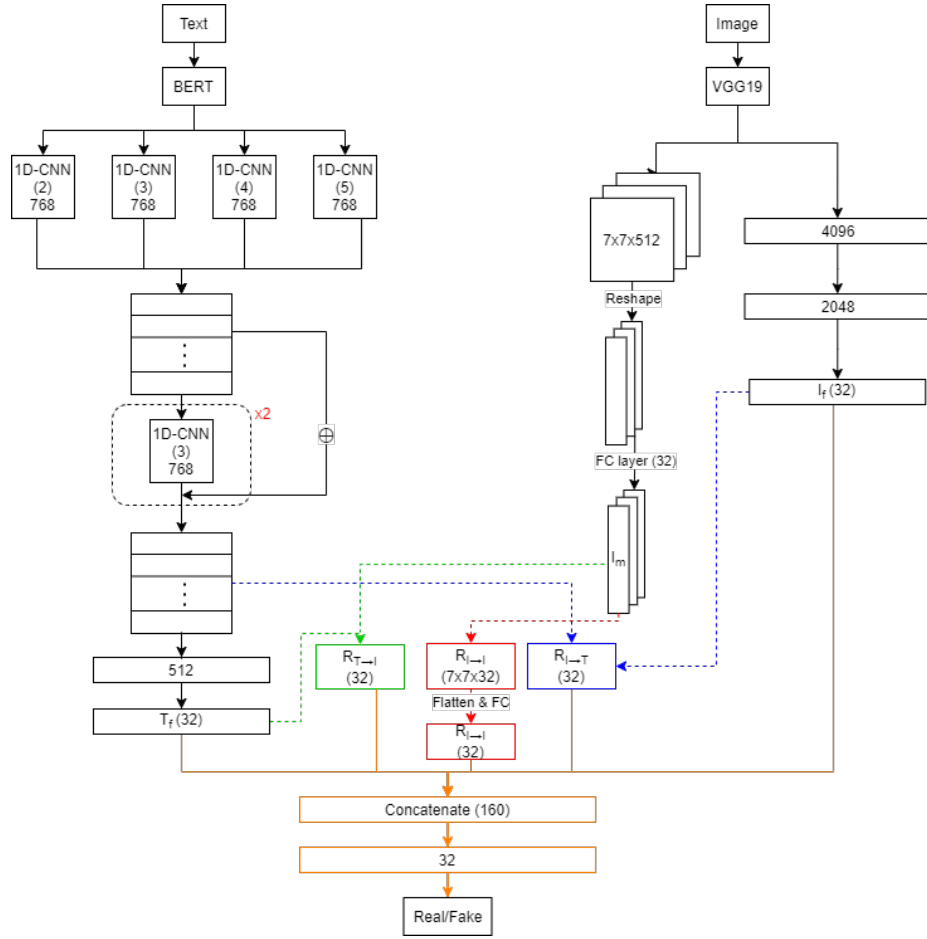


Fig. 2: Model Architecture

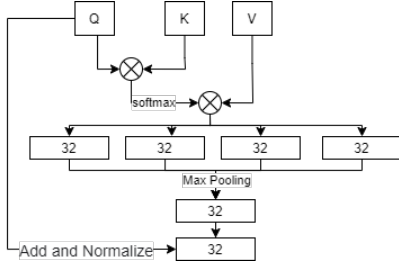


Fig. 3: Scaled dot-product attention mechanism

pass it through two feed-forward layers to decrease the dimension to a vector size to $d_I = 32$ as the final visual representation I_f . We also extract the third last layer output I_m that we will use later in the common feature extracting part. I_m is reshaped from 3D tensor to 2D tensor of shape $regions \times d_{I_m}$ where $regions$ is the number of pixels in the output of the pre-trained VGG-19's second final layer. There are $7 \times 7 = 49$ regions for an image of size 224×224 . Finally, we put it into a fully connected layer and obtain the visual feature representation of size $d_{I_m} = 32$.

C. Common Feature Extractor

In this part, we proposed to apply the scaled dot-product attention mechanism on visual features and textual features (I_f , I_m , T_f , T_m) to capture how well text and images are related to each other in a post. We applied the attention mechanism for text and visual features in two directions and we also used the self-attention mechanism on the images since we believe that the contents in an unmanipulated image should relate to each other. Figure 3 shows the general design of our proposed attention mechanism.

When we use information from a text to make the comparison to an image, or we use the text vector representation T_f as Query and image regions' features I_m as Key and Value, the three terms Query, Key and Value are defined as follows.

$$Q = T_f \times W_Q, K = I_m \times W_K, V = I_m \times W_V$$

$$W_Q \in R^{d_T \times d}, W_K \in R^{d_{I_m} \times d}, W_V \in R^{d_{I_m} \times d}$$

where $d = 32$, W_Q, W_K, W_V are weight matrices and \times is matrix-multiplication operation.

The output matrix of the scaled dot-product attention applied on Q, K, and V is calculated as follows.

$$Att_{T \rightarrow I} = softmax(\frac{Q \times K^T}{\sqrt{d}}) \times V$$

where $Att_{T \rightarrow I}$ is the attention's output matrix when we use text vector representation T_f as Query and image regions' features I_m as Key and Value.

Similarly, when we use information from the image to make the comparison to the text, we obtain $Att_{I \rightarrow T}$, the attention's output matrix in which the visual vector representation I_f is used as Query and T_m is used as Key and Value.

The output matrix $Att_{I \rightarrow I}$ of the self-attention mechanism applied on an image is calculated as follows.

$$Q = I_m \times W_Q, K = I_m \times W_K, V = I_m \times W_V \\ W_Q \in R^{d_{I_m} \times d}, W_K \in R^{d_{I_m} \times d}, W_V \in R^{d_{I_m} \times d}$$

After obtaining three attention's output matrices $Att_{T \rightarrow I}$, $Att_{I \rightarrow T}$ and $Att_{I \rightarrow I}$, we pass each of them into four different fully connected layers with the size of 32, which is the same as d_T , d_I and d_{I_m} . Then we take a maximum of 4 vectors, pass it into another fully connected layer and add a residual connection into it. Finally, we obtain three vectors $R_{T \rightarrow I}$, $R_{I \rightarrow T}$ and $R_{I \rightarrow I}$. We also use layer normalization [1] to the output of each attention block.

D. Multiple Feature Combination

In this step, $R_{I \rightarrow I}$ is flattened and pass into a fully connected layer of size 32 and we obtain $R'_{I \rightarrow I}$. Finally, we concatenate 5 outputs: T_f , I_f , $R_{T \rightarrow I}$, $R_{I \rightarrow T}$ and $R'_{I \rightarrow I}$ and pass the concatenated tensor through a fully connected layer with 32 neural units. In summary, we have 2 vector representations for textual features, which are T_f , $R_{I \rightarrow T}$ and 3 vector representations for visual features, which are I_f , $R_{T \rightarrow I}$, $R'_{I \rightarrow I}$.

IV. DATASET

We evaluated our proposed model on MediaEval 2016 data. The dataset was released for the Verifying Multimedia Use challenge at MediaEval 2016 [3], [7]. The data contains tweets and images associated with tweets. There are 17,000 unique tweets about various events. We used the same train/test split as provided in MediaEval 2016 data in which 15,000 news tweets (9,000 fake-news and 6,000 real-news) were used for training and 2,000 news tweets were used for evaluation. Some tweets have attached videos, but we only kept the data with texts and attached images and removed samples with attached videos.

V. EXPERIMENTS AND RESULTS

A. Experimental setup

In experiments, we used one image in a post as the input for visual feature extraction. We randomly chose one image from the post if the post contains multiple images.

We used BERTweet model [17] for text feature extraction. We chose the maximum sequence length of 32 for the BERT

model and used padding for texts whose lengths are shorter or longer than 32. In our proposed model, we kept weights of pre-trained BERT and VGG-19 fixed and used them as feature extractors because in preliminary experiments, we found that fine-tuning BERT and VGG-19 did not improve the performance of our model.

Hyper-parameters used in experiments are as follows. The number of filters of each 1D-CNN layer is 768, and we used the pooling size of 3 in max-pooling layers. The hidden size of the fully connected layer before T_f is 768.

We resized all the images to 224x224x3. The second last output of VGG-19 has the size of 4096 and the third last output has the shape of 7x7x512. The hidden size of the fully connected layer before I_f is 2048 and the hidden size of the layer before I_m is 32.

After each fully connected layer in our proposed model, we applied the Dropout technique [21] and set dropout rate to 0.3. We trained the model with 10 epochs, the batch size of 256, and Adam optimizer with the learning rate of 1e-4. Also, after each feed-forward layer and CNN layer, we used a batch normalization layer.

B. Results

In experiments, we compared our proposed model with some baseline uni-modal models and multimodal models as follows. Baseline multimodal models have been described in Section II.

- **TextLSTM** is an LSTM network that includes a bidirectional LSTM layer, a feed-forward layer, and a softmax layer. We used Google pre-trained word embeddings with the dimension of 32. The model used only textual features.
- **Textual (BERTweet)** used only textual features derived from the pre-trained BERTweet model. The model corresponds to the left component in Figure 2.
- **Visual** used only visual features derived by VGG-19. The model corresponds to the right component in Figure 2.
- **att-RNN** is the multimodal model with attention mechanism presented in [11]. We did not use social features in **att-RNN** for a fair comparison.
- **EANN** is the Event Adversarial Neural Networks presented in [24].
- **MVAE** [13] is a multimodal model based on Variational Autoencoder.
- **Spotfake** [20] combined multimodal features extracted by BERT and VGG-19.

Table I shows the results of the baselines and our proposed method on the MediaEval 2016 dataset. Our proposed model outperformed baseline uni-modal models and multi-modal models. In comparison with Spotfake - a model which combines multimodal features derived from a pre-trained BERT model and a pre-trained VGG-19 model, our model improved 3.1% in accuracy and 4.4% in Macro-F1.

We also compared two versions of BERT in the proposed model: BERTweet and BERT base in Table II. The results indicated that on MediaEval 2016 data, BERTweet outperforms

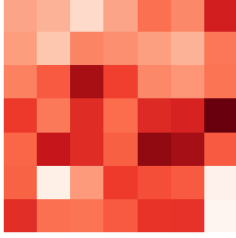


(a) Real News



(b) Fake News

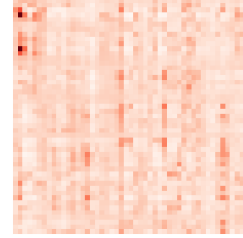
Fig. 4: The image in the news contents.



(a) Text as Q and Image as K

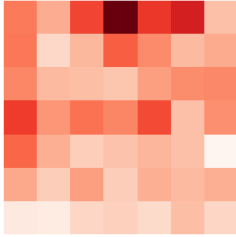


(b) Image as Q and Text as K



(c) Self-attention on Image

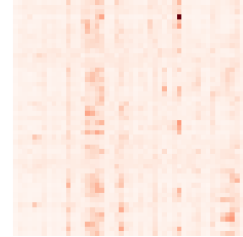
Fig. 5: Attention outputs on real news.



(a) Text as Q and Image as K



(b) Image as Q and Text as K



(c) Self-attention on Image

Fig. 6: Attention outputs on fake news.

Model	Accuracy	Fake News			Real News		
		Precision	Recall	F1-score	Precision	Recall	F1-score
TextLSTM	0.526	0.586	0.553	0.569	0.469	0.526	0.496
Textual (BERTweet)	0.666	0.667	0.840	0.743	0.664	0.430	0.522
Visual	0.596	0.695	0.518	0.593	0.524	0.7	0.599
att-RNN	0.664	0.749	0.615	0.676	0.589	0.728	0.651
EANN	0.648	0.810	0.498	0.617	0.584	0.759	0.660
MVAE	0.745	0.801	0.719	0.758	0.689	0.777	0.730
Spotfake	0.777	0.751	0.900	0.82	0.832	0.606	0.701
Proposed model	0.812	0.813	0.874	0.843	0.810	0.728	0.767

TABLE I: Performance comparison between models

Model	Accuracy	Fake News			Real News		
		Precision	Recall	F1-score	Precision	Recall	F1-score
BERT base	0.669	0.769	0.607	0.678	0.585	0.753	0.659
BERT large	0.788	0.806	0.832	0.819	0.762	0.728	0.744
BERTweet	0.812	0.813	0.874	0.843	0.810	0.728	0.767

TABLE II: Performance in different versions of BERT

BERT models trained on general domain corpus. A plausible explanation is that BERTweet was trained on 850 million

English Tweets which have the same domain as MediaEval 2016 data.

In order to show the effectiveness of our proposed attention mechanism, we plot the attention weight heat maps for a piece of real news and a piece of fake news in Figure 4a and Figure 4b, respectively. The real news (Figure 4a) comes with text content: *RT @saserief: “@Conflicts: French military deployed on the streets of #Paris large scale terror attack taking place - @lepoint https://t... and the fake news (Figure 4b) has content: *RT @londonorganiser: Ever get the feeling you’re being had Europe? http://t.co/shIVo0S7RZ*. There is clear discrimination between the two images through the heat maps. In the fake news’ image, regions are less related to each other since it is a manipulated image from 4 different images with added words. Figure 5b shows that the picture can clearly express the content of the tweet (and vice versa), which is about *military deployed on the streets*, meanwhile, there is no connection between the text content and images in*

the fake news, so it is very ambiguous when using the picture to express the content of the tweet.

VI. CONCLUSION

We have presented a multimodal approach for fake news identification on social media. We combined textual features derived from a pre-trained BERT model and visual features derived from VGG-19 pre-trained on ImageNet data. Especially, we proposed a novel attention mechanism to learn the correlation between modalities. Experimental results confirmed the effectiveness of our method in the fake news detection task. In future work, we plan to use more than one image for identifying fake news. We also intend to use user-related features such as the number of friends, the number of followers, and other metadata features from the tweets such as the number of likes and comments.

REFERENCES

- [1] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," 2016.
- [2] C. Baziotis, N. Pelekis, and C. Doulkeridis, "Datastories at semeval-2017 task 4: Deep lstm with attention for message-level and topic-based sentiment analysis," in *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*. Association for Computational Linguistics, August 2017, pp. 747–754.
- [3] C. Boididou, K. Andreadou, S. Papadopoulos, D.-T. Dang-Nguyen, G. Boato, M. Riegler, Y. Kompatsiaris *et al.*, "Verifying multimedia use at mediaeval 2015," *MediaEval*, vol. 3, no. 3, p. 7, 2015.
- [4] C. Castillo, M. Mendoza, and B. Poblete, "Information credibility on twitter," in *Proceedings of the 20th International Conference on World Wide Web*. Association for Computing Machinery, 2011, p. 675–684.
- [5] G. S. Cheema, S. Hakimov, and R. Ewerth, "Tib's visual analytics group at mediaeval '20: Detecting fake news on corona virus and 5g conspiracy," 2021.
- [6] A. Conneau, H. Schwenk, L. Barrault, and Y. Lecun, "Very deep convolutional networks for text classification," 2017.
- [7] Detection and visualization of misleading content on Twitter, "Boididou, christina and papadopoulos, symeon and zampoglou, markos and apostolidis, lazaros and papadopoulou, olga and kompatsiaris, yiannis," *International Journal of Multimedia Information Retrieval*, vol. 7, no. 1, pp. 71–86, 2018.
- [8] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "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, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186.
- [9] X. Duan, E. Naghizade, D. Spina, and X. Zhang, "RMIT at PAN-CLEF 2020: Profiling Fake News Spreaders on Twitter," in *CLEF 2020 Labs and Workshops, Notebook Papers*, L. Cappellato, C. Eickhoff, N. Ferro, and A. Névéol, Eds. CEUR Workshop Proceedings, Sep. 2020.
- [10] S. Ghosh and C. Shah, "Towards automatic fake news classification," *Proceedings of the Association for Information Science and Technology*, vol. 55, no. 1, pp. 805–807, 2018.
- [11] Z. Jin, J. Cao, H. Guo, Y. Zhang, and J. Luo, "Multimodal fusion with recurrent neural networks for rumor detection on microblogs," in *Proceedings of the 25th ACM International Conference on Multimedia*. Association for Computing Machinery, 2017, p. 795–816.
- [12] R. K. Kaliyar, A. Goswami, P. Narang, and S. Sinha, "Fndnet – a deep convolutional neural network for fake news detection," *Cogn. Syst. Res.*, vol. 61, no. C, p. 32–44, Jun. 2020.
- [13] D. Khattar, J. S. Goud, M. Gupta, and V. Varma, "Mvae: Multimodal variational autoencoder for fake news detection," in *The World Wide Web Conference*. Association for Computing Machinery, 2019, p. 2915–2921.
- [14] Y. Kim, "Convolutional neural networks for sentence classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 1746–1751.
- [15] S. Krishnan and M. Chen, "Identifying tweets with fake news," in *2018 IEEE International Conference on Information Reuse and Integration (IRI)*, 2018, pp. 460–464.
- [16] J. Ma, W. Gao, and K.-F. Wong, "Rumor detection on twitter with tree-structured recursive neural networks," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, Jul. 2018, pp. 1980–1989.
- [17] D. Q. Nguyen, T. Vu, and A. T. Nguyen, "BERTweet: A pre-trained language model for English Tweets," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2020, pp. 9–14.
- [18] K. Shu, X. Zhou, S. Wang, R. Zafarani, and H. Liu, "The role of user profile for fake news detection," 2019.
- [19] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015.
- [20] S. Singhal, R. R. Shah, T. Chakraborty, P. Kumaraguru, and S. Satoh, "Spotfake: A multi-modal framework for fake news detection," in *2019 IEEE Fifth International Conference on Multimedia Big Data (BigMM)*, 2019, pp. 39–47.
- [21] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 56, pp. 1929–1958, 2014.
- [22] L. Tian, X. Zhang, Y. Wang, and H. Liu, "Early detection of rumours on twitter via stance transfer learning," in *Advances in Information Retrieval*, J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, and F. Martins, Eds. Cham: Springer International Publishing, 2020, pp. 575–588.
- [23] L. Wang, Y. Wang, G. de Melo, and G. Weikum, "Five shades of untruth: Finer-grained classification of fake news," in *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 2018, pp. 593–594.
- [24] Y. Wang, F. Ma, Z. Jin, Y. Yuan, G. Xun, K. Jha, L. Su, and J. Gao, "Eann: Event adversarial neural networks for multi-modal fake news detection," *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018.
- [25] K. Wu, S. Yang, and K. Q. Zhu, "False rumors detection on sina weibo by propagation structures," *2015 IEEE 31st International Conference on Data Engineering*, pp. 651–662, 2015.
- [26] Y. Yang, L. Zheng, J. Zhang, Q. Cui, Z. Li, and P. S. Yu, "Ti-cnn: Convolutional neural networks for fake news detection," 2018.
- [27] X. Zhou, J. Wu, and R. Zafarani, "Safe: Similarity-aware multi-modal fake news detection," 2020.
- [28] X. Zhou and R. Zafarani, "Network-based fake news detection: A pattern-driven approach," 2019.