

# Fake news: an algorithmic perspective on fact-checking

Martijn B.J. Schouten  
11295562

Bachelor thesis  
Credits: 12 EC

Bachelor's degree Information Science

University of Amsterdam  
Faculty of Science  
Science Park 904  
1098 XH Amsterdam

*Supervisor*  
Dr. M. J. Marx

ILPS, IvI  
Faculty of Science  
University of Amsterdam  
Science Park 904  
1098 XH Amsterdam

2019-06

Add a  
final date

20

## Abstract

21

Add an  
abstract

22	<b>Contents</b>	
23	<b>1 Introduction</b>	<b>4</b>
24	<b>2 Related Work</b>	<b>5</b>
25	2.1 Automatic fake news detection . . . . .	5
26	2.2 Pre-trained textual embeddings . . . . .	5
27	2.3 Pooling . . . . .	6
28	2.4 Padding . . . . .	6
29	2.5 Neural text classifiers . . . . .	7
30	<b>3 Methodology</b>	<b>7</b>
31	3.1 Description of the data . . . . .	7
32	3.2 Data preprocessing and cleaning . . . . .	8
33	3.2.1 Filtering statements . . . . .	8
34	3.2.2 Reducing labels . . . . .	9
35	3.3 Methods . . . . .	9
36	3.4 Applying embedding techniques . . . . .	9
37	3.4.1 RQ1 . . . . .	10
38	3.4.2 RQ2 . . . . .	10
39	3.4.3 RQ3 . . . . .	10
40	<b>4 Evaluation</b>	<b>10</b>
41	4.1 RQ1 . . . . .	10
42	4.2 RQ2 . . . . .	10
43	4.3 RQ3 . . . . .	10
44	<b>5 Conclusions</b>	<b>12</b>
45	5.1 Acknowledgements . . . . .	12

# 1 Introduction

The ability to broadcast information on a large scale has been in the hands of large publishing organizations in the pre-Internet era, but nowadays everyone can share news via social media [9]. This introduces risks on validity and authenticity of news, as social media and digital platforms can speed up the spread of falsehoods without much effort from the author [6].

As a matter of fact, 63% of adults in the United States prefer to read their news on the Internet. Young adults take the lead: 76% of adults between the ages 18 and 49 get their primary news consumption via the web, compared to just 43% for adults of 50 years and older [16]. As time passes by, social media is slowly becoming the primary source of news for more and more people.

The main danger of this development is that human perception is often skewed with regards to objectivity of facts. Naïve realism let consumers of news belief that their perception is right, while other's perceptions are uninformed. Furthermore, confirmation bias results in consumers preferring information that confirms beliefs they already have [25]. This makes consumers vulnerable for the spread of misinformation or fake news.

According to the European Commission, "*disinformation - or fake news - consists of verifiably false or misleading information that is created, presented and disseminated for economic gain or to intentionally deceive the public, and may cause public harm*" [6]. The answer to the problem of fake news as of recently has been to manually fact-check statements on validity, but, as Shu et al. underlines, one of the downsides to this approach is that fake news typically relates to newly emerging, time-critical events. This means the real news may not be fully verified by proper knowledge bases due to a lack of contradicting claims [25]. An automated approach would both help in solving the problem of human subjectivity and the speed at which false information is spread in the current news spreading landscape. Furthermore, such an approach can help human fact-checking by targetting statements that are most likely to be false.

Natural language processing has been in rapid development over the past years. With the releases of OpenAI's GPT-2 model in February of this year and Google's BERT in the autumn of 2018, state-of-the-art pre-trained textual embedding techniques have shown promising results on various classification tasks [20][8]. Although fake news classification has been attempted before [28][11], performance has been rather low. However, these new pre-trained textual embeddings have not yet been used in the fight against disinformation.

This thesis is focussed on the following research question: *what is the performance of combinations of pre-trained embedding techniques with machine learning algorithms when classifying fake news?* This main question will be answered through the results of the following subquestions:

**RQ1** Which way of pooling vectors to a fixed length works best for classifying fake news?

**RQ2** At what maximum sequence length do neural networks hold the highest accuracy when classifying fake news?

**RQ3** How well do neural network classification architectures classify fake news compared to non-neural classification algorithms?

## 93 2 Related Work

### 94 2.1 Automatic fake news detection

95 There have been several attempts in the past to create classifiers for automatic  
96 detection of lies and fake news. Wang used both linear and neural classifiers to  
97 classify statements from the Liar dataset into 6 possible gradations of truthfulness.  
98 Furthermore, he added speaker metadata to improve the result of his classifications.  
99 Both with and without introducing speaker metadata, the best performing  
100 architecture was found to be a convolutional neural network. With an accuracy  
101 of 27% without, and 27,4% with metadata on the test set, Wang was able to  
102 perform 6,2% and 6,6% better than the majority baseline of 20,8%[28].

103 From the same dataset, Khurana extracted linguistic features such as n-  
104 grams, sentiment, number of capital letters and POS tags to classify the data  
105 into 3 labels instead of the original 6 labels. For classification, she used a set of  
106 non-neural classifiers. Her best performing classifier, using gradient boosting,  
107 obtained an accuracy of 49,03%, which performed around 5% better than the  
108 majority baseline of 44,28% [11].

109 The British factchecking organization Full Fact has developed an architecture  
110 that is able to monitor and factcheck statements from the British Parliament  
111 and major media outlets in the United Kingdom. It can automatically factcheck  
112 the accuracy of statistical claims, for example [4]. For detecting factual claims  
113 from texts, the organization uses InferSent, which is a way of transfer learning  
114 that has been proved to perform well for the use case of Full Fact [13].

115 Various tools with regards to fake news detection are also available. Faker  
116 Fact is a tool which can classify texts into a set of categories ranging from satire  
117 to agenda-driven, the former identifying humorous intent, the latter identifying  
118 manipulation [1]. Hoaxy, on the other hand, allows for the visualization of  
119 unverified claims through Twitter networks [23].

### 120 2.2 Pre-trained textual embeddings

121 Traditionally, feature representation for text classification is often based on the  
122 bag-of-words model, containing linguistic patterns as features, such as unigrams,  
123 bigrams or n-grams. However, these approaches completely ignore contextual  
124 information or word order in texts, and are unable to capture semantics of words.  
125 As a result, classifiers may be unable to correctly identify patterns, affecting the  
126 classification accuracy [14].

127 As an answer to these problems, pre-trained text embeddings have been  
128 rising in popularity, both in use and in research. Before classification is possible,  
129 text data needs to be transformed into numbers to be able to be interpreted by  
130 classifier algorithms. Fundamentally, text embeddings are vector representations  
131 of linguistic structures, allowing for usage of text in classifiers. The process of  
132 turning words into these embeddings is typically powered by statistics gathered  
133 from large unlabeled corpora of text data [15].

134 In 2017, Vaswani et al. proposed a novel architecture for embedding text  
135 data called the Transformer. With the main aim originally being translation  
136 from one sentence in some language to another sentence in another language,  
137 Transformers are based on an encoder-decoder model. These models take the  
138 sequence of input words, convert it to an intermediate representation, after

139 which the decoder creates an output sequence.

140 The main strength of the Transformer architecture is its focus on attention  
141 to create the intermediate representation. The encoder receives a sequence of  
142 inputs, and reads it at once, as opposed to sequentially (either from left to right  
143 or from right to left, as humans do). This allows the encoder to learn the context  
144 of a word based on all of the surrounding text [27].

145 Inspired by the Transformer architecture, numerous new text embedding  
146 techniques have been developed which give the possibility to classify texts by  
147 making use of the intermediate representation of the Transformer model. As  
148 shown by the Bidirectional Encoder Representations from Transformer (BERT)  
149 model by Devlin et al., these techniques have beaten existing benchmarks in  
150 natural language processing, underlining the importance of context in text  
151 embeddings [8].

## 152 2.3 Pooling

153 Linear classifiers need data in a two-dimensional shape to be able to perform  
154 calculations. In the case of raw text data, sentences in the dataset have variable  
155 word lengths, resulting in a different vector length when turning the text into a  
156 vector representation. To turn the vector representations into a uniform length,  
157 we can either cut off the vectors at a fixed length (*padding*), or we can perform  
158 calculations to reduce the length of the vectors (*pooling*).

159 In computer vision, feature pooling is used to reduce noise in data. The  
160 goal of this step is to transform joint feature representations into a new, more  
161 usable one that preserves important information while discarding irrelevant  
162 details. Pooling techniques such as max pooling and average pooling perform  
163 mathematical operations to reduce several numbers into one [5]. In the case of  
164 transforming the shape of the data, we can reduce vectors to the smallest vector  
165 in the dataset to create a uniform shape.

166 Scherer et al. compared performance of two pooling operations on a convolutional  
167 neural network architecture. The first pooling method extracted maximums  
168 and the second one was primarily based on working with averages. They have  
169 shown that a max pooling operation is vastly superior for capturing invariances  
170 in image-like data [22].

171 Shen et al. noted that in text classification, only a small number of key  
172 words contribute to the final prediction. As a result, simple pooling operations  
173 are surprisingly effective for representing documents [24]. Lai et al., Hu et al.  
174 and Zhang et al. use a max pooling layer in a (recurrent) convolutional neural  
175 network for identifying key features in text classification [14][10][30]. In the case  
176 of text classification, max pooling strategies seem to be the most popular.

## 177 2.4 Padding

178 When padding a sequence, a list of sequences is transformed to a specific length.  
179 Sequences longer than the desired length will be truncated to fit the requirement,  
180 while sequences shorter than the desired length will be padded by a specified  
181 value [2]. To fill the sequences, a value of zero is often used. Hu et al. also use  
182 zero values for padding their sequences [10].

183 Apart from controlling the size of the feature dimension, padding has other  
184 uses as well. Simard et al. make use of sequence padding for convolutional

neural networks to center feature units, and concluded it did not impact the performance of the classifier significantly [26]. Wen et al. apply padding to convolutional network models to prevent dimension loss [29].

## 188 2.5 Neural text classifiers

Wang has shown that neural networks perform slightly better on classifying fake news than linear classifiers. In his research, he compared accuracies on support vector machines, logistic regressions, bidirectional LSTMs and convolutional neural networks with each other. With his 6 label classification, his support vector machine implementation was the best performing linear classifier, but the performance was slightly worse than the best performing neural network (25.8% for the former, and 26% for the latter) [28].

Wang used two neural network architectures both well known for their robustness and performance when it comes to text classification. The first model, the bidirectional Long Short Term Memory (LSTM) networks, are specifically tailored at keeping track of information for a long period of time. This makes those models able to keep track of the context in a more intelligent way when compared to a standard non-neural classification algorithm [17].

The second architecture, the convolutional neural network, apply a set of filters in its layers to local features. These models are shown to be effective in numerous natural language processing applications, such as semantic parsing, search query retrieval, sentence modeling and other traditional NLP tasks [12].

### 206 3 Methodology

### 207 3.1 Description of the data

For classifying fake news, Wang’s Liar dataset will be used [28]. The Liar dataset contains 12,791 short statements from Politifact.com, which are labeled manually by a Politifact.com editor on truthfulness. The statements are an average of 18 tokens long, and the topics vary from different political subjects, as can be seen in figure 1. Truthfulness is evaluated by assigning one of 6 labels, ranging from *pants-on-fire* to *true*. The distribution of statements across the original 6 labels can be seen in table 2.

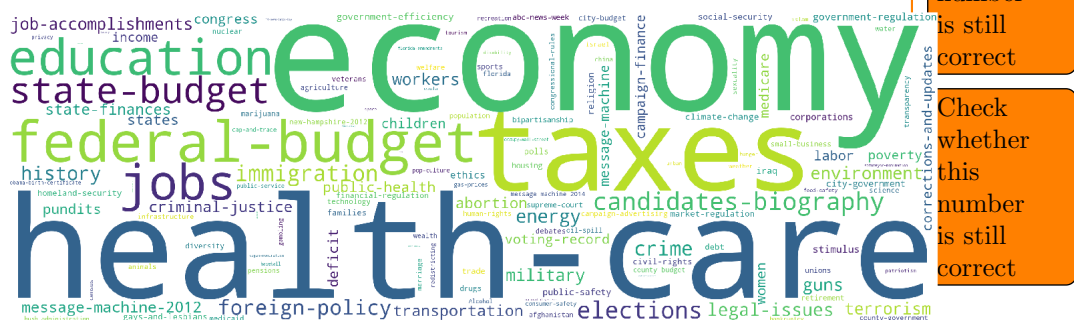


Figure 1: An overview of all statement topics in the Liar dataset.

id	11044.json
label	pants-fire
statement	The Mexican government forces many bad people into our country.
subjects	foreign-policy,immigration
speaker	donald-trump
speaker_job	President-Elect
state	New York
party	republican
context	an interview with NBC's Katy Tur
mostly_true_count	37
half_true_count	51
barely_true_count	63
false_count	114
pants_on_fire_count	61

Table 1: An example entry in the Liar dataset.

For each statement, the dataset contains an id, a label, a subject, a speaker, the function title of the speaker, the affiliated state and political affiliation, the context of the statement and a vector with a truthfulness history. An example of such a data entry can be seen in table 1. Wang introduced this truthfulness history to boost the prediction scores, as speakers with a track record of lying are expected to have a lower chance of speaking the truth when classifying new statements. However, for our application we are only interested in the statement itself and its corresponding label. Due to cheapness and spreadability, a large amount of fake news is spread over social media [25]. This means author information and metadata will not readily be available in real world circumstances.

The original dataset has been split beforehand into a test, train and validation set. The train set contains 80% of the total amount of statements, while the test and validation set both contain approximately 10% of the statements.

Check whether this number is still correct

## 3.2 Data preprocessing and cleaning

### 3.2.1 Filtering statements

The original dataset contained statements ranging from 1 sentence to 19. On closer inspection of statements with the high amounts of sentences, it was found that not all statements were processed from source files into dataframes correctly. As a result, records of some different statements were joined together, forming a single string. To combat this, the following regular expression was used to filter those statements out:

```
\\.json\\t(mostly-true|true|half-true|false|barely-true|pants-fire)\\
```

After applying this regular expression, the total amount of sentences in the statements were reduced from a maximum of 19 to a maximum of 11.



6 labels	3 labels	2 labels
true (16.1%)	true (35,3%)	true (55,8%)
mostly-true (19.2%)		
half-true (20.5%)	half-true (20.5%)	
barely-true (16.4%)	false (44,19%)	false (44,19%)
false (19.6%)		
pants-fire (8.19%)		

Table 2: Distribution of labels from the original label distribution when reducing the amount of labels.

### 3.2.2 Reducing labels

Wang’s main objective was to classify fake news into a fine-grained category of fakeness [28]. For our main research question, we aim to predict whether the statements are fake news or not. This means the statements do not necessarily need to be distinguished into these fine-grained categories. Because of this, the classifiers used to predict fake news in this research will be trained on the original 6 labels, Khurana’s division into three labels [11], and a binary classification. The division from the original 6 labels into the lesser amounts of labels can be seen in table 2. This way, we can better compare performance of pre-trained embeddings to existing research on this dataset.

## 3.3 Methods

### 3.4 Applying embedding techniques

As our main research question is focussed on pre-trained word embeddings, the first step in the classification process is to turn the statements of the Liar dataset into vectors. For this purpose, the Flair framework will be used. Flair contains interfaces for turning words into embeddings, built on the PyTorch platform [21][19]. Using Flair, we have access to the following 5 state-of-the-art pre-trained embedding techniques:

- ELMo (Embeddings from Language Models) [18];
- BERT [8];
- Generative Pre-Training (GPT) [?];
- Transformer-XL [7];
- Flair [3].

265 **3.4.1 RQ1**

266 **3.4.2 RQ2**

267 **3.4.3 RQ3**

## 268 **4 Evaluation**

### 269 **4.1 RQ1**

270 **Which way of pooling vectors to a fixed length works best for classifying**  
271 **fake news?**

272 The answer to this question will be given by comparing accuracy of a logistic  
273 regression with different pooling strategies (max, min, average).

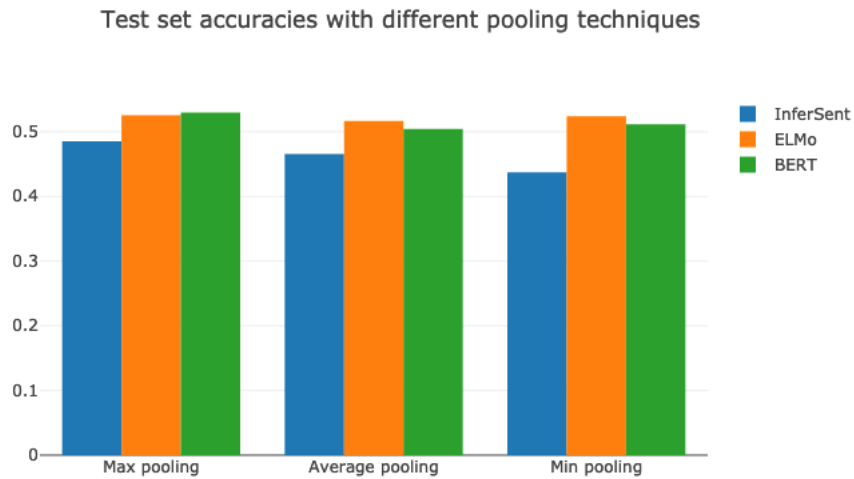


Figure 2: Comparing different pooling strategies.

### 274 **4.2 RQ2**

275 **At what padding sequence length do neural networks hold the highest**  
276 **accuracy when classifying fake news?**

277 The answer to this question will be given by comparing accuracy of a bidirectional  
278 LSTM architecture with variable padding lengths.

279 By comparing all embedding techniques, I plan on getting a peak padding  
280 length with the best overall performance shared by all embedding methods. At  
281 this moment, the optimal length averages at around a length of 22.

### 282 **4.3 RQ3**

283 **How well do neural network classification architectures classify fake**  
284 **news compared to non-neural classification algorithms?**

Test set accuracy of padded datasets with variable maximum lengths

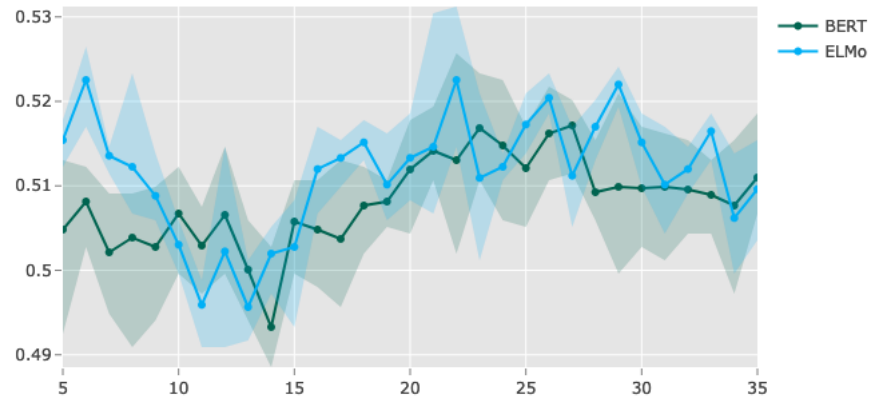


Figure 3: Comparing different maximum padding lengths.

285 The answer of this question will be given by comparing two linear classifiers  
286 with two neural classifiers.

Test set accuracy of machine learning models for each embedding technique

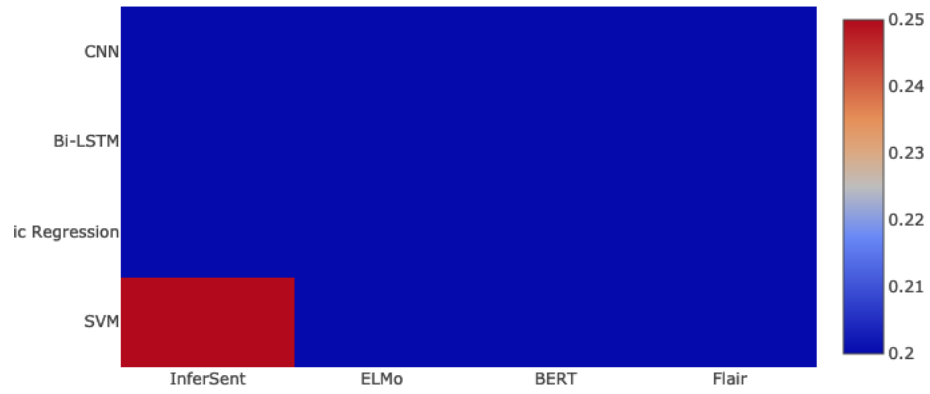


Figure 4: Comparing linear classifiers with neural classifiers.

## 287 5 Conclusions

### 288 5.1 Acknowledgements

## References

- [1] About fakerfact.
- [2] Sequence preprocessing.
- [3] Alan Akbik, Tanja Bergmann, and Roland Vollgraf. Pooled contextualized embeddings for named entity recognition. In *NAACL 2019, 2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, page to appear, 2019.
- [4] Mevan Babakar and Will Moy. The state of automated factchecking. Technical report, Full Fact, 2016.
- [5] Y-Lan Boureau, Jean Ponce, and Yann LeCun. A theoretical analysis of feature pooling in visual recognition. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 111–118, 2010.
- [6] European Commission. Fake news and online disinformation. <https://ec.europa.eu/digital-single-market/en/fake-news-disinformation>, 2018. Retrieved on 8th of April, 2019.
- [7] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *CoRR*, abs/1901.02860, 2019.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [9] Lee Howell et al. Digital wildfires in a hyperconnected world. *WEF Report*, 3:15–94, 2013.
- [10] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. Convolutional neural network architectures for matching natural language sentences. In *Advances in neural information processing systems*, pages 2042–2050, 2014.
- [11] Urja Khurana. The linguistic features of fake news headlines and statements, 2017.
- [12] Yoon Kim. Convolutional neural networks for sentence classification. *CoRR*, abs/1408.5882, 2014.
- [13] Lev Konstantinovskiy. Sentence embeddings for automated factchecking - lev konstantinovskiy.
- [14] Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. Recurrent convolutional neural networks for text classification. In *Twenty-ninth AAAI conference on artificial intelligence*, 2015.
- [15] Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhersch, and Armand Joulin. Advances in pre-training distributed word representations. *CoRR*, abs/1712.09405, 2017.

- [16] Amy Mitchell and Hannah Klein. Americans still prefer watching to reading the news – and mostly still through television. 2018.
- [17] Christopher Olah. Understanding lstm networks. <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>, August 2015.
- [18] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *CoRR*, abs/1802.05365, 2018.
- [19] PyTorch. From research to production. <https://pytorch.org/>. Retrieved on 10th of June 2019.
- [20] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [21] Zalando Research. Flair. <https://github.com/zalando-research/flair>, 2019.
- [22] Dominik Scherer, Andreas Müller, and Sven Behnke. Evaluation of pooling operations in convolutional architectures for object recognition. In *International conference on artificial neural networks*, pages 92–101. Springer, 2010.
- [23] Chengcheng Shao, Giovanni Luca Ciampaglia, Alessandro Flammini, and Filippo Menczer. Hoaxy: A platform for tracking online misinformation. In *Proceedings of the 25th International Conference Companion on World Wide Web, WWW ’16 Companion*, pages 745–750, Republic and Canton of Geneva, Switzerland, 2016. International World Wide Web Conferences Steering Committee.
- [24] Dinghan Shen, Guoyin Wang, Wenlin Wang, Martin Renqiang Min, Qinliang Su, Yizhe Zhang, Chunyuan Li, Ricardo Henao, and Lawrence Carin. Baseline needs more love: On simple word-embedding-based models and associated pooling mechanisms. *CoRR*, abs/1805.09843, 2018.
- [25] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. Fake news detection on social media: A data mining perspective. *CoRR*, abs/1708.01967, 2017.
- [26] Patrice Y Simard, David Steinkraus, John C Platt, et al. Best practices for convolutional neural networks applied to visual document analysis. In *Icdar*, volume 3, 2003.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- [28] William Yang Wang. ”liar, liar pants on fire”: A new benchmark dataset for fake news detection. *CoRR*, abs/1705.00648, 2017.
- [29] L. Wen, X. Li, L. Gao, and Y. Zhang. A new convolutional neural network-based data-driven fault diagnosis method. *IEEE Transactions on Industrial Electronics*, 65(7):5990–5998, July 2018.

- 369 [30] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional  
370 networks for text classification. In *Advances in neural information*  
371 *processing systems*, pages 649–657, 2015.