



Fake news outbreak 2021: Can we stop the viral spread?

Tanveer Khan^{a,*}, Antonis Michalas^{a,b}, Adnan Akhunzada^c

^a Tampere University, Finland

^b RISE AB

^c Technical University of Denmark, Denmark

ARTICLE INFO

Keywords:

Fake news
Fact checking
Machine learning
Tools
Datasets

ABSTRACT

Social Networks' omnipresence and ease of use has revolutionized the generation and distribution of information in today's world. However, easy access to information does not equal an increased level of public knowledge. Unlike traditional media channels, social networks also facilitate faster and wider spread of disinformation and misinformation. Viral spread of false information has serious implications on the behaviours, attitudes and beliefs of the public, and ultimately can seriously endanger the democratic processes. Limiting false information's negative impact through early detection and control of extensive spread presents the main challenge facing researchers today. In this survey paper, we extensively analyze a wide range of different solutions for the early detection of fake news in the existing literature. More precisely, we examine Machine Learning (ML) models for the identification and classification of fake news, online fake news detection competitions, statistical outputs as well as the advantages and disadvantages of some of the available data sets. Finally, we evaluate the online web browsing tools available for detecting and mitigating fake news and present some open research challenges.

1. Introduction

The popularity of Online Social Networks (OSNs) has rapidly increased in recent years. Social media has shaped the digital world to an extent it is now an indispensable part of life for most of us (Gazi and Çetin, 2017). Rapid and extensive adoption of online services is influencing and changing how we access information, how we organize to demand political change and how we find partners. One of the main advantages and attractions of social media is the fact that it is fast and free. This technology has dramatically reshaped the news and media industries since becoming a dominant and growing source of news and information for hundreds of millions of people (Kaplan, 2015). In the United States today more people are using social media as a news source than ever before (Suciu). Social media has progressively changed the way we consume and create news. The ease of producing and distributing news through OSNs has also simultaneously sharply increased the spread of fake news.

Fake news is not a new phenomenon; it existed long before the arrival of social media. However, following the 2016 US presidential election it has become a buzzword (Albright, 2016). There are numerous examples of fake news through history. A notable one from the antiquity is the Mark Anthony smear campaign circa 44 BC (Posetti and Matthews,

2018). In more recent times, examples include the anti-German campaign, German corpse factory in 1917 (Neander and Marlin, 2010) and the Reich Ministry of Public Enlightenment and Propaganda established in 1933 by the Nazis to spread Nazi ideology and incite violence against Jews (Bytwerk, 2010).

Although propaganda campaigns and spread of fabricated news may have been around for centuries, their fast and effective dissemination only became possible by means of a modern technology such as the internet. The internet revolutionized fake news, regardless of how the misinformation is manifested: whether we are talking about a rumor, disinformation, or biased, sloppy, erroneous reporting. In a recent study (Wong, 2016), it was found that almost 50 percent of traffic taken from Facebook is fake and hyperpartisan, while at the same time, news publishers relied on Facebook for 20 percent of their traffic. In another study, it was found that 8 percent of 25 million Universal Resource Locator (URLs) posted on social media were indicative of malware, phishing and scams (Thomas, 2013).

Researchers in Germany conducted a study regarding fake news distribution in the country and people's attitudes and reactions towards it (Christian Reuter et al., 2019). Based on the published results, 59 percent of participants stated that they had encountered fake news; in some regions, this number increased to almost 80 percent (G. L. R. D.,

* Corresponding author.

E-mail address: tanveer.khan@tuni.fi (T. Khan).

2019). Furthermore, more than 80 percent of participants agreed fake news poses a threat and 78 percent strongly believed it directly harms democracy. Government institutions and powerful individuals use it as a weapon against their opponents (Carson, 2019). In the 2016 US presidential election, a significant shift in how social media was used to reinforce and popularize narrow opinions was observed. In November of the same year, 159 million visits to fake news websites were recorded (Allcott and Gentzkow, 2017), while the most widely shared stories were considered to be fake (Silverman, 2016). Similarly, it is believed that the distribution of fake news influenced the UK European Union membership referendum (Grice, 2017).

However, fake news is not only about politics. During the recent fires in Australia, several maps and pictures of Australia's unprecedented bushfires spread widely on social media. While users posted them to raise awareness, the result was exactly the opposite since some of the viral maps were misleading, spreading disinformation that could even cost human lives (Rannard, 2020). The recent COVID-19 pandemic accelerated the rise of conspiracy theories in social media. Some were alleging that the novel coronavirus is a bio-weapon funded by Bill Gates to increase the selling of vaccines (M. for minds). Undoubtedly fake news threaten multiple spheres of life and can bring devastation not only to economic and political aspects but peoples' wellbeing and lives.

1.1. An overview of this survey

The main motivation behind our study was to provide a comprehensive overview of the methods already used in fake news detection as well as bridge the knowledge gap in the field, thereby helping boost interdisciplinary research collaboration. This work's main aim is to provide a general introduction to the current state of research in the field.

We performed an extensive search of a wide range of existing solutions designed primarily to detect fake news. The studies used deal with identification of fake news based on ML models, network propagation models, fact-checking methods etc. More precisely, we start by examining how researchers formulate ML models for the identification and classification of fake news, which tools are used for detecting fake news and conclude by identifying open research challenges in this domain.

Comparison to Related Surveys. In a related work by Vitaly Klyuev (2018), an overview of the different semantic methods by concentrating on Natural Language Processing (NLP) and text mining techniques was provided. In addition, the author also discussed automatic fact-checking as well as the detection of social bots. In another study, Oshikawa et al. (2018) focused on the automatic detection of fake news using NLP techniques only. Two studies can be singled out as being the closest to our work. First, Study by Collins et al. (2020) which examined fake news detection models by studying the various variants of fake news and provided a review of recent trends in combating malicious contents on social media. Second, a study by Shu et al. (2020a) which mostly focused on various forms of disinformation, factors influencing it and mitigating approaches.

Although some similarities are inevitable, our work varies from the aforementioned ones. We provide a more detailed description of some of the approaches used and highlight the advantages and limitations of some of the methods. Additionally, our work is not limited to NLP techniques, but also examines types of detection models available, such as, knowledge-based approaches, fact-checking (manual and automatic) and hybrid approaches. Furthermore, our approach considers how the NLP techniques are used for the detection of other variants of fake news such as rumors, clickbaits, misinformation and disinformation. Finally, it also examines the governmental approaches taken to combat fake news and its variants.

1.2. Organization

The rest of this paper is organised as follows: Section 2 discusses the

most important methods for detecting fake news, in Section 3, we detailed both the automatic and manual assessment of news and analysed different ways of measuring the relevance, credibility and quality of sources. To automate the process of fake news detection, the analysis of comprehensive data sets is of paramount importance. To this end, in Section 4, we first discuss the characteristics of online tools used for identifying fake news and then compare and discuss different data sets used to train ML algorithms to effectively identify fake news. The classification of existing literature, identified challenges, future directions and existing governmental strategies to tackle the problem of fake news detection are discussed in Section 5. Finally, the concluding remarks are given in Section 6.

2. Fake news analysis

People are heavily dependent on social media for getting information and spend a substantial amount of time interacting on it. In 2018, the Pew Research Center revealed that 68 percent of Americans (Matsa and Shearer, 2018) used social media to obtain information. On average, 45 percent of the world's population spend 2 h and 23 min per day on social media and this figure is constantly increasing (Asano, 2017). The biggest problem with information available on social media is its low quality. Unlike the traditional media, at the moment, there is no regulatory authority checking the quality of information shared on social media. The negative potential of such unchecked information became evident during the 2016 US presidential election.¹ In short, it is of paramount importance to start considering fake news as a critical issue that needs to be solved.

In spite of the overwhelming evidence supporting the need to detect fake news, there is, as yet, no universally accepted definition of fake news. According to (Lazer et al., 2018), "fake news is fabricated information that mimics news media content in form but not in organizational process or intent". In a similar way, fake news is defined as "a news article that is intentionally and verifiable false" (Shu et al., 2017). Some articles also associate fake news with terms such as deceptive news (Allcott and Gentzkow, 2017), satire news (Rubin et al., 2015), clickbait (Chen et al., 2015), rumors (Zubiaga et al., 2018), misinformation (Kucharski, 2016), and disinformation (Kshetri and Voas, 2017). Hence, these terms are used interchangeably in this survey.

The following forms of misuse of information have been considered as variants of fake news in the existing literature (Tandoc et al., 2018; Rubin et al., 2015):

- **Clickbait:** Snappy headlines that easily capture user attention without fulfilling user expectations since they are often tenuously related to the actual story. Their main aim is to increase revenue by increasing the number of visitors to a website.
- **Propaganda:** Deliberately biased information designed to mislead the audience. Recently, an increased interest has been observed in propaganda due to its relevance to the political events (Rubin et al., 2015).
- **Satire or Parody:** Fake information published by several websites for the entertainment of users such as "The Daily Mash" website. This type of fake news typically use exaggeration or humor to present audiences with news updates.
- **Sloppy Journalism:** Unreliable and unverified information shared by journalists that can mislead readers.
- **Misleading Headings:** Stories that are not completely false, but feature sensationalist or misleading headlines.

¹ <https://www.independent.co.uk/life-style/gadgets-and-tech/news/tumblr-russian-hacking-us-presidential-election-fake-news-internet-research-agency-propaganda-bots-a8274321.html>.

- **Slanted or Biased News:** Information that describes one side of a story by suppressing evidence that supports the other side or argument.

For years, researchers have been working to develop algorithms to analyze the content and evaluate the context of information published by users. Our review of the existing literature is organised in the following way: [subsection 2.1](#), examines approaches to identifying different types of user accounts such as bots, spammers and cyborgs. It is followed by [subsection 2.2](#), where different methods used for identifying rumors and clickbaits are discussed. In [subsection 2.3](#), the users' content and context features are considered while in [subsection 2.4](#), different approaches for the early detection of fake news by considering its propagation are discussed.

2.1. User account analysis

According to a report published in 2021 Twitter alone has 340 million users, 11.7 million registered apps, delivers 500 million tweets a day and 200 billion tweets a year ([Aslam, 2021](#)). It's popularity has made it an ideal target for bots, or automated programs ([Kaur et al., 2018](#)). Recently, it was reported that around 5–10 percent of Twitter accounts are bots and responsible for the generation of 20–25 percent of all tweets ([U. of Eastern Finland, 2019](#)). Some of the bots are legitimate, comply with Twitter objectives, and can generate a substantial volume of benign tweets like blogs and news updates. Other bots, however, can be used for malicious purposes such as a malware that gathers passwords or a spam that adds random users as friends and expects to be followed back ([IONOS, 2018](#)). Such bots have a more detrimental effect particularly when spreading fake news. The significance of differentiating the legitimate bots from the malicious ones emerged from the fact that malicious bots can also be used to mimic human behaviour in a negative way.

Researchers examined bots, in a number of existing publications ([Cresci et al., 2017](#); [Gibert et al., 2020](#); [Edwards et al., 2014](#); [Lee et al., 2011](#); [Wu et al., 2013](#); [Stone-Gross et al., 2009](#)). [Gilani et al. \(2017\)](#) focused on classifying Twitter accounts into "human" and "bots" and analyzing the impact each has on Twitter. The proposed technique was based on previous work by "Stweeler" ([gilani, 2018](#)) for the collection, processing, analysis, and annotation of data. For the identification of bots, human annotation was used, where participants differentiated bots from humans and generated a reliable data set for classification. The process provided an in-depth characterization of bots and humans by observing differences and similarities. The finding stated that the bots' removal from Twitter causes serious repercussions for content production and information dissemination and also indicated that bots count on re-tweeting, redirecting users via URLs, and uploading media. However, the imprecision in the existing algorithm revealed by the authors and the manual collection of data limited the ability to analyze accounts.

Similarly, [Giachanou et al. \(Giachanou and Ghanem, 2019\)](#) investigated whether the Twitter account author is human or a bot and further determined the gender of a human account. For this purpose, a linear Support Vector Machines (SVM) classifier was trained to analyze words, character grams, and stylistic features. For the identification of human gender, a stochastic gradient descent classifier was used to assess the sentiment of tweets, words, and character grams and point wise mutual information features – the importance of terms per gender. The data set used consisted of tweets in English and Spanish. The experiments illustrated the accuracy of bot detection, i.e. 0.906 for bots in English and 0.856 for Spanish. Similarly, for the identification of gender, the accuracy for English tweets amounted to 0.773 and 0.648 for Spanish tweets. In the long run, the bot detection model outperformed the gender detection model.

Another account type that can be generated on Twitter is a Cyborg. Cyborg refers to a human-assisted bot or bot-assisted human ([Chu et al., 2010](#)). Cyborgs have characteristics of both human-generated and

bot-generated accounts and as such require a level of human engagement. These accounts facilitate posting various information more frequently, rapidly and long-term ([DFRLab, 2016](#)). Differentiating a cyborg from a human can be a challenging task. The automated turing test ([Von Ahn et al., 2004](#)) used to detect undesirable or bot programs is not capable of differentiating cyborgs from humans. However, [Jeff Yan \(2006\)](#) proposed that a cyborg might be differentiated by comparing the characteristics of a machine and human elements of a cyborg. Similarly, [Chu et al. \(2012\)](#) differentiate between bot, cyborg and human accounts by taking into account tweet content, tweeting behaviour and features of the account.

OSNs also serve as platforms for the rapid creation and spread of spam. Spammers act similarly to bots and are responsible for posting malicious links, prohibited content and phishing sites ([Perez, 2018](#); [Michalakis and Murray, 2017](#)). Traditional methods of detecting spammers that utilize network structure are classified into three categories:

- Link-based, where the number of links is used as a measure of trust. These links are considered to be built by legitimate users ([Lee et al., 2011](#)).
- Neighbor-based, which treats links as a measure of homophily, the tendency for linked users to share similar beliefs and values ([Hu et al., 2013](#); [Rayana and Akoglu, 2015](#); [Li et al., 2016](#)).
- Group-based, which recognizes that spammers often work in groups to coordinate attacks ([Jindal and Liu, 2007](#)). Group-based methods detect spammers by taking advantage of the group structure hidden in the social network. Additionally, spammers behave differently from legitimate users so they can be treated as outliers ([Gao et al., 2010](#); [Akoglu et al., 2015](#)).

Current efforts for detection of social spammers utilize the structure and behavioural patterns of social spammers in an attempt to discover how their behaviour can be differentiated from legitimate users ([Lim et al., 2010](#); [Chu et al., 2010](#); [Ye et al., 2016](#); [Li et al., 2015](#); [Xue et al., 2013](#); [Yang et al., 2014](#)). However, spammers often find ways to create a link with legitimate users, making it more difficult to detect specific spamming patterns. [Wu et al. \(2017a\)](#) tackled this problem by taking into account both content and network structure. They proposed "Sparse Group Modeling for Adaptive Spammer Detection (SGASD)" that can detect both types of spammers – those within a group and individuals.

Another challenging task is detection of camouflaged content polluters on OSNs. Content polluters – spammers, scanners and fraudsters – first establish links with a legitimate user and then merge the malicious with real content. Due to insufficient label information available for camouflaged posts in online media, the use of these manipulated links and contents as camouflage makes detecting polluters very difficult. In order to tackle this challenge, [Wu et al. \(2017b\)](#) studied how camouflaged content polluters can be detected and proposed a method called "Camouflaged Content Polluters using Discriminate Analysis (CCPDA)" which can detect content polluters using the patterns of camouflaged pollution.

[Chris et al. \(Grier et al., 2010\)](#) spam detection analysis juxtaposed two different types of Twitter accounts – a "professional spamming account" whose sole purpose is to distribute spam, versus "accounts compromised by spammers". The authors found that accounts currently sending spam had been compromised by spammers; once legitimate, they became controlled by spammers. Furthermore, to detect spam activity on Twitter, a directed social graph model ([Wang, 2010](#)) based on friend and follower relationships was proposed. Different classifier techniques were used to distinguish between the spammer and normal behaviour and determined that the Naive Bayes classifier performs better with respect to F-measure.

Huge momentum has been observed where user-generated content is exploited in micro-blogs for predicting real-world phenomena such as prices and traded stock volume on financial markets ([Cresci et al., 2016](#)). Research efforts in this domain targeted sentiment metrics as a predictor

for stock prices (Bollen et al., 2011; Chen et al., 2014; Gabrovšek et al., 2017), company tweets and the topology of the stock network (Mao et al., 2012; Ruiz et al., 2012) and used weblogs pointing to the relationship between companies (Kharratzadeh and Coates, 2012). Cresci et al. (2019) demonstrated the use of twitter stock micro-blogs as a platform for bots and spammers to practice cash-tag piggybacking – an activity for promoting low-value stocks by exploiting the popularity of real high-value stocks. They employed a spambot detection algorithm to detect accounts that issue suspicious financial tweets. Nine million tweets from five main US financial markets, which presented stocks with unusually high social significance compared to their low financial relevance, were investigated with respect to their social and financial significance. These tweets were compared with financial data from Google finance. The results indicated that 71 percent of users were classified as bots and that high discussion of low-value financial stocks was due to a massive number of synchronized tweets.

Twitter currently has no defined policy for addressing automated malicious programs operating on its platform. However, it is expected that these malicious accounts will be deleted in the near future (Schwartz, 2018). A survey of the literature has identified numerous studies (Ergahin et al., 2017; Lee et al., 2010; Gilani et al., 2017; Davis et al., 2016; Antoniadis et al., 2015; Weimer et al., 2007; Adewole et al., 2017; Wanas et al., 2008; Weerkamp and De Rijke, 2008; Morris et al., 2012) that describe the important characteristics which can be used for the identification of bots on Twitter. Despite these attempts, limitations still exist in employing these characteristics for detecting fake news, especially, early detection of fake news during its propagation. Other methods, such as network propagation, have to be utilized for this purpose.

2.2. Identifying rumors and clickbaits

Social media is like a blank sheet of paper on which anything can be written (Yaraghi, 2019), and people easily become dependent on it as a channel for sharing information. This exactly is the reason why social media platforms (e.g. Twitter and Facebook) are highly scrutinized for the information shared on them (Haralabopoulos et al., 2015). These platforms have undertaken some efforts to combat the spread of fake news but have largely failed to minimize its effect. For instance, in the United States, 60 percent of adults who depend on social media for news consumption are sharing false information (Wong, 2019). In April 2013, two explosions during the Boston Marathon gained tremendous notoriety in the news and on social media, and the tragedy was commented on in millions of tweets. However, many of those tweets were rumors (controversial factual claims) and contained fake information, including conspiracy theories. Similarly, a survey published by Kroll – a business intelligence and investigating firm – states that 84 percent of companies feel threatened by the rise of rumors and fake news fuelled by social media (Binham, 2019). On Weibo, rumors were detected in more than one-third of trending topics (Zhao et al., 2015). The spread of rumors on social media has also become an important issue for companies worldwide. Still, there is no clear policy defined by social media administrators to verify shared content. Below, we discuss different techniques that have been proposed by researchers to address this problem.

Traditionally, human observers have been used to identify trending rumors. Currently, research is focused on building an automated rumor identification tool. For this purpose, a rumor detection technique (Zhao et al., 2015) was designed. In this technique, two types of clusters were generated: posts containing words of inquiry such as “Really”, “What”, “Is it true?” were grouped into one cluster. These inquiries were then used to detect rumor clusters. Similarly, posts without words of inquiry were grouped into another cluster. Similar statements were extracted from both clusters. The clusters were then ranked, based on their likelihood of containing these words. Later, the entire cluster was scanned for disputed claims. These experiments, performed with Twitter data, resulted in earlier and effective detection of rumors (almost 50 rumor

clusters were identified). However, there is still considerable space to improve these results (Zhao et al., 2015). For instance, the manual collection of inquiry words could be improved by training a classifier and the process of ranking could be improved by exploring more features for the rumor cluster algorithm.

People can share fake information on social media for various reasons. One of those is to increase readership, which is easily achievable by using clickbait. Clickbait is a false advertisement with an attached hyperlink. It is specifically designed to get users to view and read the contents inside the link (Anna Escher, 2016). These advertisements attract users with catchy headlines but contain little in the way of meaningful content. A large number of users are lured by clickbait. Monther et al. (Aldwairi and Alwahedi, 2018) provided a solution to protect users from clickbait in the form of a tool that filters and detects sites containing fake news. In categorizing a web page as a source of fake news, they considered several factors. The tool navigates the content of a web page, analyzes the syntactical structure of the links and searches for words that might have a misleading effect. The user is then notified before accessing the web page. In addition, the tool searches for the words associated with the title in the links and compares it with a certain threshold. It also monitors punctuation marks such as question and exclamation marks used on the web page, marking it as a potential clickbait. Furthermore, they examined the bounce rate factor–percentage of visitors who leave a website, associated with the web page. Where the bounce rate factor was high, the content was marked as a potential source of misleading information.

A competition was organised with the aim of building a classifier rating the extent to which a media post can be described as clickbait. In the clickbait competition, the data set was generated from Twitter and consisted of 38,517 Twitter posts from 27 US news publishers (Potthast et al., 2018). Out of 38,517 tweets, 19,538 were available in the training set and 18,979 were available for testing. For each tweet, a clickbait score was assigned by five annotators from Amazon Mechanical Turk. The clickbait scores assigned by human evaluators were: 1.0 heavily clickbaity, 0.66 considerably clickbaity, 0.33 slightly clickbaity and 0.0 not clickbaity. The goal was to propose a regression model that could determine the probable clickbaitiness of a post. The evaluation metric used for the competition was Mean Squared Error (MSE). In this competition, Omidvar et al. (2018) proposed a model using the deep learning method and won the challenge. They achieved the lowest MSE for clickbait detection by using a bi-directional Gated Recurrent Unit (biGRU). Instead of solving the clickbait challenge using a regression model, Yiewi Zhou (Zhou, 2017) reformulated the problem as a multi-classification. On the hidden state of biGRU, a token level self-attentive mechanism was applied to perform multi-classification. This self attentive Neural Network (NN) was trained without performing any manual feature engineering. They used 5 self-attentive NNs with a 80-20 percent split and obtained the second lowest MSE value. Similarly, Alexey Grigorev (2017) proposed an ensemble of Linear SVM models to solve the clickbait problem and achieved the third lowest MSE value. In addition to the given data set, they gathered more data from multiple Facebook groups that mostly contained clickbait posts by using the approach described in “identified clickbaits using ML” (Thakur, 2016).

2.3. Content and context analysis for fake news

The rapid dissemination of fake news is so pernicious that researchers resolved towards trying to automate the process by using ML techniques such as Deep Neural Networks (DNN). However, the Black box problem – a lack of transparency in decision-making in the NN – obscures reliability. Nicole et al. (O’Brien et al., 2018) addressed the deep learning “black-box problem” for fake news detection. A data set composed of 24,000 articles was created, consisting of 12,000 fake and 12,000 genuine articles. The fake articles were collected from Kaggle while genuine ones were sourced from The Guardian and The New York

Table 1
Detailed Summary of the Studies used in Bots, Clickbaits and Rumors.

Approach	Data set and Features	Evaluation Metrics	Finding or Outcomes	Weaknesses	Platform
URL blacklisting for Spam detection (Grier et al., 2010)	400 million tweets, 25 million URLs	Click through, measuring delay	8% of 25 million URLs indicative of phishing, scams and malware	Inaccurate when used for web services (Thomas et al., 2011)	Twitter
Bayesian classifier for Spam detection (Wang, 2010)	25 K Users, 500 K tweets, 49 million follower/friend	Precision	Web crawler, directed social graph model, 89% precision in spam detection	Manual analysis of collected data (Lee and Kim, 2013)	Twitter
SVM for Spam detection (Benevenuto et al., 2010)	54 million users, 1.9 billion links, 1.8 billion tweets	Precision, recall, Micro F-1, Macro F-1	Correctly classified: 70% spammers, 96% non-spammers	Manually labelled data set (Ghosh et al., 2012)	Twitter
Naive bayes for account classification (Ersahin et al., 2017)	501 fake accounts, 499 real accounts, profile and tweets	ROC curve, F1 score, confusion matrix	Accuracy 90.9%	Manually labelled data set (Alothali et al., 2018)	Twitter
Ranking model for rumor detection (Zhao et al., 2015)	10,240,066 tweets, keyword search	Precision, detection time	Clustering and classification performs effectively, earlier detection of rumors	More features can be explored for rumor clustering	Twitter
SVM-rank for account classification (Gilani et al., 2017)	60 million tweets, tweets frequency	Classification accuracy, precision, recall and F1 measure	Develop and evaluate a mechanism to classify automated agents and human users	Cannot check relative frequency of any particular URL (Nasim et al., 2018)	Twitter
SVM and Stochastic Gradient Descent for bots and gender profiling (Giachanou and Ghanem, 2019)	English and Spanish Tweets, textual, stylistic	Accuracy	Words and char grams are important feature for gender and bot detection	–	Twitter
SGASD for spam detection (Wu et al., 2017a)	TwitterS, TwitterH, Network, content	Precision, recall, F1	Present SGASD framework for spammer detection	Network information focuses on user instead of information	Twitter
Logistic Regression for stance detection (Ferreira and Vlachos, 2016)	300 rumors, 2595 news articles, headlines	Accuracy, precision, recall	Emergent Dataset used for a variety of NLP tasks	Data set (cannot learn all nuances of tasks)	Emergent project

Times. The study concluded that DNNs can be used to detect the language patterns in fabricated news. Additionally, the algorithm can also be used for detecting fake news in novel topics.

Another technique to tackle the deep learning “black-box problem” in fake news detection is CSI (capture, score and integrate) – a three-step system which incorporates the three basic characteristics of fabricated news (Ruchansky et al., 2017). These characteristics include text, source, and the response provided by users to articulate missing information. In the first step, a Recurrent Neural Network (RNN) is used to capture the momentary pattern of user activity. The second step estimates the source of suspicions related to user behaviour. The third, hybrid step involves integration of steps one and two and is used to predict fake articles. The experiments were performed on real-world data sets and demonstrated a high level of accuracy in predicting fake articles. Still, a bottleneck in using a computationally intensive model is posed by the lack of a manually labelled fake news data set. William Yang Wang (2017) addressed the limited availability of labelled data sets for combating fake news using statistical approaches and chose a contemporary publicly available data set called LIAR. This data set was utilized to investigate fabricated news using linguistic patterns. The results were based on an evaluation of several approaches, including Logistic Regression (LR), the Convolution Neural Network (CNN) model, Long Short-Term Memory (LSTM) networks and SVM. They concluded that combination of meta-data with text significantly improves the detection of fake news. According to the authors, this body of information can also be used to detect rumors, classify stance and carry out topic modeling, argument mining, and political Natural Language Processing (NLP) research. Table 1 present a summary of the different approaches proposed for both the account as well as content and context analysis of fake news.

In 2017 a competition named Fake News Challenge (FNC) was held with the aim to use Artificial Intelligence (AI) techniques to combat the problem of fake news. During the initial phase, stance detection was used. It refers to the relative stance to any issue, claim or topic made by two pieces of relevant text (what other organizations say about the topic). A two-level scoring system was applied – 25 percent weight was assigned if the text was deemed to be related or unrelated to its headline and 75 percent weight was assigned on the basis of labelling the related

pairs as agrees, disagrees, discusses or unrelated. In this competition, the top team submitted an ensemble model for a Deep Convolution Neural Network (DCNN) and Gradient-Boosted Decision Tree (GBDT) with a weighted average of 50/50 (Sean Baird, 2017). The DCNN and GBDT separately did not achieve perfect accuracy. However, the combination of both approaches correctly detected the stance of each headline with a score of 82.01. Similarly, approach proposed by team Athene (Hanselowski, 2017) achieved a score of 81.97 and won second place in the competition. They used an ensemble approach involving multi-layer perception and applied MLP and Bag-of-Word (BoW) features to the challenge. The team in third place, Riedel et al. (2017), proposed a stance detection system for FNC Stage 1. For the input text, they used two BoW representations. A MLP classifier was used with one hidden layer having 100 units. For the hidden layer, a Rectified Linear Unit (RELU) activation function was used while the final linear layer utilized a soft-max. They achieved an accuracy of 81.72.

At a different competition named Web Search and Data Mining (WSDN) 2019, fake news was detected by classifying the titles of articles. Using a given title for any fake news article ‘A’ and a title for another incoming news article ‘B’, people were asked to classify the incoming article into one of three categories: agrees, disagrees and unrelated (Risdal, 2017). The winner of this competition Lam Pham (2019), who achieved 88.298 percent weighted accuracy on the private leader boards and 88.098 percent weighted accuracy on the public leader boards. This ensemble approach incorporated NNs and gradient boosting trees. In addition, Bidirectional Encoder Representation from Transformer (BERT) was used for encoding news title pairs, transforming and incorporating them into a new representational space. The approach by Liu et al. won a second place (Liu et al., 2019) by proposing a novel ensemble framework based on the Natural Language Interference (NLI) task. Their proposed framework for fake news classification consisted of three-level architecture with a 25 BERT model along with a blending ensemble strategy in the first level followed by 6 ML models and finally a single LR for the final classification. Yang et al. (2019) also considered this problem as a NLI task and considered both the NLI model as well as the BERT. They trained the strongest NLI models, Dense RNN, Dense CNN, ESIM, Gate CNN (Dauphin et al., 2017) and decomposable attention, and achieved an accuracy of 88.063 percent.

Table 2
Detailed Summary of the Studies used in Network as well as Content and Context Analysis.

Approach	Data set and Features	Evaluation Metrics	Finding or Outcomes	Weaknesses	Platform
RNN and CNN for fake news detection (Liu and Wu, 2018)	Weibo (Ma et al., 2016), Twitter 15 and Twitter 16 (Ma et al., 2017), user profiles	Effectiveness, efficiency	Outperforms a state-of-the-art model in terms of both effectiveness and efficiency	Problem with computational efficiency and interpretability (Zhou and Zafarani, 2019)	Twitter, Weibo
Bayesian inference for fake news detection (Tschitschholek et al., 2018)	4039 users, 88,234 edges, users and spammers	Utility, engagement, robustness	Outperforms NO-LEARN and RANDOM algorithms	Trustworthiness of news sources is ambiguous	Facebook
Diffusion of network information for classification (Vosoughi et al., 2018)	126,000 stories, 3 million users, 4.5 million tweets, retweets, users	Diffusion dynamics	Fake news spreads rapidly and deeply and is more novel than true news	Information cascades (Constantinides et al., 2018)	Twitter
LSTM-RNN for fake news detection (Wu and Liu, 2018)	3600 fake news, 68,892 real news, network information	Micro-F1, Macro-F1	Trace miner: classifying social media messages	Only considers network information	Twitter
Network flow model for fact-checking (Shiralkar et al., 2017)	Synthetic corpora, real-world data set, Edge capacity	AUROC	Network flow techniques are promising for fact-checking	Limitation of content-based approach (Pan et al., 2018)	WSDM-Cup 2017 (Hannah Bastl, 2017)
Bow-Tie and D-core decomposition for user analysis (Garcia et al., 2017)	40 million users, 1.47 billion follower links	Reputation, social influence, popularity	Global metrics are more predictive than local	Theory-driven approach (Hasani-Mavriqi et al., 2018)	Twitter
Hybrid CNN for fake news detection (Wang, 2017)	LIAR 12,836 short statements, Metadata and text features	Accuracy	LIAR data set, Integrate text and metadata	Justification and evidence are ignored in experiments	Politifact
RNN and user behaviour for fake news detection (Ruchansky et al., 2017)	Two real-world data sets (Twitter and Weibo) (Ma et al., 2016), text	Classification accuracy	More accurate in fake news classification	No assumptions about user distribution behaviour	Twitter, Weibo
DNN for fake news detection (O'Brien et al., 2018)	12,000 fake and 12,000 real news, language patterns	Accuracy	Observes subtle differences in language patterns of real and fake news	Only predicts truthfulness of claims	Different websites

2.4. Network propagation and detection of fake news

One of the OSNs main strong points is facilitating the propagation of information between users. The information of interest to users is further shared with relatives, friends, etc (Kong et al., 2019). In order to detect the propagation of fake news at its early stage (Liu and Wu, 2020), it is crucial to be able to understand and measure the information propagation process. The influence of propagation on OSNs and their impact on network structure was studied in (Saxena et al., 2015; Hong et al., 2011). Ye et al. (Ye and Wu, 2010) study revealed that more than 45.1 percent of information shared by a user on social media is further propagated by his/her followers. Furthermore, approximately 37.1 percent of the information shared is propagated up to 4 hops from the original publisher.

Liu and Wu (2018) used the data network features and introduced a popular network model for the early detection of fake news. They addressed the limitation of low accuracy of early fake news detection by classifying news propagation paths as a multivariate time series. Characteristics of each user involved in spreading news were represented by a numerical vector. Then a time series classifier was built by combining CNN and RNN. This classifier was used for fake news detection by capturing the local and global variations of observed characteristics along the propagation path. This model is considered as more robust, as it relies on common user characteristics which are more reliable and accessible in the early stage of news propagation. The experiments were performed on two real-world data sets based on Weibo (Ma et al., 2016) and Twitter (Ma et al., 2017). The proposed model detected fake news within 5 min of its spread with 92 percent accuracy for Weibo and 85 percent accuracy for Twitter data sets.

Sebastian et al. (Tschitschholek et al., 2018) examined the ways to minimize the spread of fake news at an early stage by stopping its propagation in the network. They aggregated user flagging, a feature introduced by Facebook that allows users to flag fake news. In order to utilize this feature efficiently, the authors developed a technique called 'DETECTIVE' which uses Bayesian Inference to learn flagging accuracy. Extensive experiments were performed by using a publicly available data set (Leskovec and Mcauley, 2012) from Facebook. The results

indicated that even with minimum user engagement DETECTIVE can leverage crowd signals to detect fake news. It delivered better results in comparison to existing algorithms, i.e. NO-Learn and RANDOM.

The dissemination of misinformation on OSNs has a particularly undesirable effect when it comes to public emergencies. Dynamic Linear Threshold (DLT) model (Litou et al., 2016) was developed to attempt and limit this type of information. It analyzes the user's probability, based on an analysis of competing beliefs, of propagating either credible or non-credible news. Moreover, an optimization problem based on DLT was formulated to identify a certain set of users that could be responsible for limiting the spread of misinformation by initiating the propagation of credible information.

A study by Garcia et al. (2017) focused on examining reputation (Dimitriou and Michalas, 2012, 2014; Michalas and Komninos, 2014), popularity and social influence on Twitter using digital traces from 2009 to 2016. They evaluated the global features and specific parameters that make users more popular, keep them more active and determine their social influence. Global measures of reputation were calculated by taking into account the network information for more than 40 million users. These new global features of reputation are based on the D-core decomposition method (Giatsidis et al., 2013) and The Twitter's bow-tie structure (Broder et al., 2000) in 2009. The results indicated that social influence is more related to popularity than reputation, and global network metrics such as social reputation are more accurate predictors for social influence than local metrics such as followers, etc.

Soroush et al. (Vosoughi et al., 2018) collected and studied twitter data from 2006 to 2007 in order to classify it as true or false news. News is classified as true or false based on information collected from six independent fact-checking organizations. They generated a data set that consisted of approximately 126,000 tweets, tweeted by 3 million twitter users approximately 4.5 million times. They found that fake news was more novel and inspired surprise, fear, and disgust in replies, while true news inspired trust, sadness, anticipation, and joy. As people prefer to share novel information, false news spreads more rapidly, deeply and broadly than true news. According to Panos et al. (Constantinides et al., 2018), rapid dissemination of information on social media is due to information cascade. Liang Wu and Huan Liu (Wu and Liu, 2018) also

classified twitter messages using diffusion network information. Instead of using content features, they focused on the propagation of Twitter messages. They proposed trace miner, a novel approach that uses diffusion network information to classify social media messages. Trace miner accepts message traces as inputs and outputs its category. Table 2 presents a detailed summary of the studies used in network as well as content and context analysis.

After reviewing the studies discussed above, it became evident there is no 'one size fits all' when it comes to fake news detection. Extensive research is still required to fully understand the dynamic nature of this problem.

3. Fact checking

The rapid spread of fraudulent information is a big problem for readers who fail to determine whether a piece of information is real or fake. Since fake news is a big threat to society and responsible for spreading confusion, it is necessary to have an efficient and accurate solution to verify information in order to secure the global content platform. To address the problem of fake news, the American media education agency Poynter established the International Fact-Checking Network (IFCN) in 2015, which is responsible for observing trends in fact-checking as well as providing training to fact-checkers. A great deal of effort has already been devoted to providing a platform where fact-checking organizations around the world can use a uniform code of principles to prevent the spread of fake news. Two fact-checking organizations, Snopes and Politifact, developed a fake news detection tool useful in classifying fake news levels in stages. However, this tool requires a lot of manual work. There is a profound need for a model that can automatically detect fake news.

Giovanni et al. reduced the complex manual fact-checking task to a simple network analysis problem (Ciampaglia et al., 2015), as such problems are easy to solve computationally. The proposed approach was evaluated by analyzing tens of thousands of statements related to culture, history, biographical and geographical information using a public knowledge graph extracted from Wikipedia. They found that true statements consistently receive higher support in comparison to false ones and concluded that applying network analytics to large-scale knowledge repositories provides new strategies for automatic fact-checking. Below, we examine two facets of fact checking problem. In 3.1 we look into computational approaches to automatic fact checking, whereas in 3.2, we concentrate on the issue of trust and credibility of the information and the source providing it.

3.1. Towards automatic fact checking

Computational approaches to fact-checking are considered key to tackling the massive spread of misinformation. These approaches are scalable and effective in evaluating the accuracy of dubious claims. In addition, they improve the productivity of human fact-checkers.

One of the proposed approaches is an unsupervised network flow-based approach (Shiralkar et al., 2017), which helps to ascertain the credibility of a statement of fact. The statement of fact is available as a set of three elements that consist of the subject entity, the object entity, and the relation between them. First, the background information of any real-world entity is viewed on a knowledge graph as a flow of the network. Then, a knowledge stream is built by computational fact-checking which shows the connection between the subject and object of a set. The authors evaluated network flow model on actual and customized fact data sets and found it to be quite effective in separating true and false statements.

A study by Baly et al. (2018) examined on the factuality and bias of claims across various news media. They collected features from articles of the target news websites, their URL structures, the web traffic they attract, their twitter accounts (where applicable) as well as Wikipedia pages. These features were then used to train the SVM classifier for bias

and factuality separately. The evaluation, showed that the articles' features achieved the best performance on factuality and bias, Wikipedia features were somewhat useful for bias but not for factuality, and Twitter and URL features faired better in factuality than bias.

A different approach for an automatic fake news detection (Pérez-Rosas et al., 2017) was based on several exploratory analyses to identify the linguistic differences between legitimate and fake news. It involved the introduction of two novel data sets, the first collected using both manual and crowdsourcing annotation, and the second generated directly from the web. Based on these, first several exploratory analyses were performed to identify the linguistic properties most common for fake news. Secondly, a fake news detector model based on these extracted linguistic features was built. They concluded that the proposed system performed better than humans in certain scenarios with respect to more serious and diverse news sources. However, human beings outperformed the proposed model in the celebrity domain.

OSNs are also used as a vector for the diffusion of hoaxes. Hoaxes spread uncontrollably as propagation of such news depends on very active users. At the same time, news organizations devote a great deal of time and effort to high-quality fact-checking of information online. Eugenio et al. (Tacchini et al., 2017) used two classification algorithms: LR and Boolean Crowd Sourcing (BCS) for classifying Facebook posts as hoaxes or non-hoaxes based on users who "liked" it. On a data set of 15, 500 posts and 909,236 users, they obtained a classification accuracy of more than 99 percent. The proposed technique even worked for users who "liked" both hoax and non-hoax posts. Similarly, Kumar et al. (2016) studied the presence of hoaxes in Wikipedia articles based on a data set consisting of 20K hoaxes explicitly and manually labelled by Wikipedia editors. According to their findings, hoaxes have very little impact and can be easily detected. A multi-modal hoax detection system that merges diverse modalities – the source, text, and the image of a tweet was proposed by Maigrot et al. (2016). Their findings suggested that using only source or text modality ensures high performance in comparison to using all the modalities. Marcella et al. (Tambuscio et al., 2015) focused on the diffusion of hoaxes on OSNs by considering hoaxes as viruses in which a normal user, once infected, behaves as a hoax-spreader. The proposed stochastic epidemic model can be interpreted as a Susceptible-Infected-Susceptible (SIS) or Susceptible-Infected-Recovered (SIR) model – the infected user can either be a believer (someone who believes the fake news) or a fact-checker (checking the news before believing it). The model was implemented and tested on homogeneous, heterogeneous and real networks. Based on a wide range of values and topologies, the fact-checking activity was analysed and then a threshold was defined for fact-checking probability (verifying probability). This threshold was used to achieve the complete removal of fake news based on the number of fact-checkers considering the news as fake or real. A study by Shao et al. focused on the temporal relation between the spread of misinformation and fact-checking, and the different ways in which both are shared by users. They proposed Hoaxy (Shao et al., 2016) – a model useful in the collection, detection, and analysis of this type of misinformation. They generated a data set by collecting data from both fake news (71 sites, 1, 287,768 tweets, 171,035 users and 96,400 URLs) and fact-checking (6 sites, 154,526 tweets, 78,624 users and 11,183 URLs) sources. According to their results, fact-checking data sharing lags behind misinformation by 10–20 h. They suggested that social news observatories could play an important role by providing the dynamics of real and fake news distribution and the associated risks.

3.2. Trust and credibility

The ease of sharing and discovering information on social media results in a huge amount of content published for target audiences. Both participants (those who share and consume) must check the credibility of shared content. Social media also enables its users to act simultaneously as content producers and consumers. The content consumer has

Table 3
Detailed summary of the studies used in Fact-checking, Trust and Credibility.

Approach	Data set and Features	Evaluation Metrics	Finding or Outcomes	Weaknesses	Platform
Structural modeling (Torres et al., 2018)	541 users, Age, gender, network size	Reliability, validity	Development of a research model based on perceptions	Ignores news connection with recipient	Social networking sites
Measuring user trust for fake news detection (Shu et al., 2018)	Two data sets, Explicit and implicit profile features	Follower to following counts ratio	Expert and naive user features differ	Does not consider the bias and credibility of users	Twitter
SVM-rank for credibility assessment (Gupta et al., 2014)	10,074,150 tweets, 4,996,448 users, features obtained from high impact crisis events	Response time, usability, effectiveness,	TweetCred browser extension	Results influenced by context of tweets and personalization	Twitter
Automated learning and Standard collaborative filtering (Ghavipour and Meybodi, 2018a)	Advogato, Observer, Apprentice, Journeyer, Master	Coverage, prediction accuracy	Efficient and accurate trust path discovery	Does not consider the dynamic nature of trust (Ghavipour and Meybodi, 2018b)	Advogato
Spam and bot detection (Cresci et al., 2019)	9 million tweets, 30,032 companies, Market capitalization, industrial classification	Cashtag	Uncovering malicious practices–cashtag piggybacking	–	Twitter
Supervised learning for misinformation detection (Antoniadis et al., 2015)	80,294 tweets, 59,660 users	Accuracy, precision, recall, F-Measure	Accuracy of timely identification of misinformation at 70%	Undefined intentions (Balestrucci et al., 2019)	Twitter
Stochastic epidemic model for fact-checking (Tambuscio et al., 2015)	Network with 1000 nodes, Spreading rate, forgetting probability	Probability	Define a fact-checking probability for hoaxes	Does not consider the heterogeneity of agents	Facebook
Hoaxy for fact-checking (Shao et al., 2016)	Fake news and fact-checking sources, data volume, time series	Keyword correlation	Propagation of fake news is dominated by active users	Fake news makes more of a contribution to data set generation	Twitter
Random forest classifier for fake news detection (Potthast et al., 2017)	1627 articles, Writing style	Accuracy, precision, recall, F1	Distinguished hyperpartisan and mainstream	Not applicable for fake news detection	Facebook
Linear SVM classifier for fake news detection (Pérez-Rosas et al., 2017)	100 fake and 100 legitimate articles	Accuracy, precision, recall, F1 measures	Two data sets, accuracy comparable to humans in detecting fake news	Humans perform better in celebrity domain	Web
LR, BCS algorithm for classification (Tacchini et al., 2017)	15,500 posts, 909,236 users, likes	Accuracy	Classification accuracy 99% for hoaxes and non-hoaxes	Limited conspiracy theories in data set (Shu et al., 2020b)	Facebook
SVM classifier for predicting factuality (Baly et al., 2018)	1066 news websites, URL, article, account	Accuracy, F_1 score, MAE and its variant	Predicting the factuality of reports and bias of news media	Limiting sharing of false content is challenging (Paschen, 2019)	Entire news medium

more flexibility in what content to follow. For the content producer, it is necessary to check and evaluate the source of information. If a user is interested in receiving information regarding a particular topic of interest from a specific source, his primary task is to check the credibility, relevance, and quality of that source. Different ways of checking credibility include:

- Looking for other users who have subscribed to such information (Canini et al., 2011).
- Assessing both the expertise (support and recommendations from other professionals) (Ericsson et al., 2018; Wang et al., 2015) and user credibility.
- Assessing the credibility of the sources (examining the content and peer support) (Rieh and Danielson, 2007).

Researchers have proposed different techniques for identifying credible and reputable sources of information. Canini et al. (2011) proposed an algorithm based on both the content and social status of the user. Weng et al. merged the web page ranking technique and topic modeling to compute the rank of a Twitter user (Weng et al., 2010). Cha et al. (2010) studied the factors that specify user influence. A random walk approach (Perozzi et al., 2014) was proposed for separating credible sources from malicious ones by performing network feature analysis.

TweetCred, a real-time web-based system, was developed to evaluate the credibility of tweets (Gupta et al., 2014). It assigns a credibility score to each tweet on a user time line rating from 1 (low credibility) to 7 (high credibility). The credibility score is then computed using a

semi-supervised ranking algorithm trained on a data set consisting of an extensive set of features collected from previous work (Pérez-Rosas et al., 2017) and manually labelled by humans. The TweetCred evaluation was performed based on its usability, effectiveness, and response time. An 80 percent credibility score was calculated and displayed within 6 s. Additionally, 63 percent of users either agreed or disagreed with the generated score by 1–2 points. Irrespective of its effectiveness, the results were still influenced by user personalization and the context of tweets which did not involve factual information.

A different research model was developed – based on perceptions related to news authors, news sharers, and users – to test verification behaviours of users. (Torres et al., 2018). The aim was to study the validation of content published by users on Social Networking Sites (SNSs). The results were assessed using a three-step analysis to evaluate the measurement model, structural model, and common method bias. It focused on the epistemology of declarations of interpersonal trust to examine factors that influence user trust in disseminated news on SNSs. To test the full model, the researchers used SmartPLS 2.0. The evaluation showed that the variety in social ties on SNSs increases trust among network participants and trust in the network reduces news verification behaviours. However, the evaluation disregards the importance of the nature of news connected with the recipient.

Trust is an important factor to be considered when engaging in social interaction on social media. When measuring trust between two unknown users, the challenging task is the discovery of a reliable trust path. In (Ghavipour and Meybodi, 2018a), Ghavipour et al. addressed the problem of reliable trust paths by utilizing a heuristic algorithm built on learning automata called DLATrust. They proposed a new approach

Table 4
Detailed summary of the online web browsing tools.

Tools	Availability	Proposed	Technique	Input	Output	Source
SurfSafe	Browser extension	Robhat labs	Comparison and textual analysis	Images and text	Safe, warning, unsafe	100 fact-checking, trusted organizations
Trusted News	Browser extension	Trusted News	–	Website content	Trustworthy, biased, satire	MetaCert protocol
Fake News Detector	Browser extension	Robhino	Crowd sourcing, ML	News content	Fake news, clickbait, extremely biased	Feedback by other tools
Fake News Guard	Browser extension	Fake News Guard	AI, network analysis, fact-checking	Webpages, links	Fake or not	Fact checkers
Decodex	Browser extension	Laurent’s team	–	Pieces of information	Satire, info, no information	600 websites
BS Detector	Browser extension	Daniel Sieradski,	Comparison model	URLs	Fake news, conspiracy theory, clickbait, extremely biased etc	Data set of unreliable domains
TrustyTweet	Browser extension	Katrin Hartwig and Christian Reuter (Hartwig and Reuter, 2019)	Media literacy	Tweets	Politically neutral, transparent and intuitive warnings	Potential indicators from previous studies
TweetCred	Browser extension	–	Semi-supervised ranking model	Tweets	Credibility score	Twitter data
FiB	Browser extension	DEVPOST	Text analysis, image analysis, web scraping	Facebook posts	Trust score	Verification using keyword extraction, image recognition, source verification
BotOrNot	Website, REST API	Clayton et al.	Classification algorithm	Twitter screen name	Bot likelihood score	Accounts for recent history including tweet mentions
LiT.RL News Verification	Web browser	Rubin et al.	NLP, support vector machine	Language used	Satirical news, clickbait, falsified news	Lexico-syntactic features in text

Table 5
Detailed summary of the available data sets in the existing literature.

Dataset	Statistics	Observations	Goal	Approach	Sources	Limitations
CRED-BANK	60 million tweets, 1049 real events	Manually annotated	Credibility assessment	Media events are linked to a human credibility judgement	Twitter	Collected tweets are not related to fake news articles (Shu et al., 2020b)
LIAR	12,836 statements	Manually annotated	Fact-checking	Assessment of truthfulness of claim	TruthO-Meter, fact-checking	Instead of entire article based on short statements (Shu et al., 2020b)
FAKE NEWS Net	Social context, content, spatiotemporal information	Fake News Tracker	Analyzing and visualizing fake news	Fake news diffusion, user engagement	PolitiFact, Twitter	Social engagement of the articles
Memetracker 9	90 million documents	22 million district phrases are extracted	Temporal patterns	Tracking ideas, new topic memes	1.65 million sites	–
BuzzFeed	left-wing and right-wing articles	Rates post as “true”, “mixture of true and false”, “false”	Fact-checked	Facebook engagement number	Facebook	Based on headlines and text only (Shu et al., 2020b)
BS Detector	–	BS Detector assigned labels	News veracity	Manually compiled list of domains	Web pages	Instead of a human expert, a tool is used for news veracity
BuzzFace	2263 news articles, 1.6 million comments	“true”, “mixture of true and false”, “false” and “no factual comments”	Veracity assessment	Extension of BuzzFeed including comments	Facebook	Context and content information but no temporal information (Shu et al., 2020b)
Facebook HOAX	15,500 posts, 32 pages, 2,300,000 likes	Scientific pages are non-hoaxes, conspiracy pages are hoaxes	Post classification into hoaxes and non-hoaxes	Number of likes per post and per user, relation between pages	Facebook	Few instances of news and conspiracy theories (Shu et al., 2020b)
Higgs-Twitter	527,496 users, 985,590 tweets	632,207 geo-located tweets	User behaviour accuracy	Analysis of spatial and temporal user activity	Twitter	No labelling of fake news
Trust and Believe	50,000 users	Manually annotated, Active learning	Influence score	Active learning approach	Twitter	Small dataset

for aggregating the trust values from multiple paths based on a standard collaborative filtering mechanism. The experiments performed on Advogato – a well-known network data set for trust, showed the efficiency and high accuracy in predicting the trust of reliable paths between two indirectly connected users.

Liu and Wu et al. (Shu et al., 2018) studied the correlation between user profiles and the fake news shared on social media. A real-world data set comprising social context and news content was built for categorizing users based on measuring their trust in fake news. Representative groups of both experienced users (able to differentiate between real and fake news) and naive users (unable to differentiate between real and

fake news) were selected. They proposed that the features relevant to these users could be useful in identifying fake news. The results for identified user groups showed that the distribution satisfies power-law distribution (Clauset et al., 2009) with high R^2 scores. This result indicated a significant difference between features of experienced and naive users. However, the paper left unexplored the credibility and political bias of experienced users before characterizing them for fake news detection.

The timely detection of misinformation and sharing of credible information during emergency situations are of utmost importance. The challenge of distinguishing useful information from misinformation

Table 6
Classification of the Studies Surveyed based on the Platform Used – Facebook and Twitter.

Platform	Research Papers
Twitter	Grier et al. (Grier et al., 2010), Ye et al. (Ye and Wu, 2010), Chengcheng et al. (Shao et al., 2016), Soroush et al. (Vosoughi et al., 2018), Liang Wu and Huan liu (Wu and Liu, 2018), Hartwig et al. (Hartwig and Reuter, 2019), Davis et al. (Davis et al., 2016), Cha et al. (Cha et al., 2010), Weng et al. (Weng et al., 2010), Canini et al. (Canini et al., 2011), Thomas Kurt (Thomas, 2013), Holton et al. (Holton and Lewis, 2011), Antoniadis et al. (Antoniadis et al., 2015), Alessandro et al. (Balestrucci et al., 2019), Zhao et al. (Zhao et al., 2015), Gupta et al. (Gupta et al., 2014), Khan and Michalas (Khan and Michalas, 2020)
Facebook	Tambuscio et al. (Tambuscio et al., 2015), Joon Ian Wong (Wong, 2016), Monther et al. (Aldwairi and Alwahedi, 2018), Potthast et al. (Potthast et al., 2017), Alexey Grigorev (Grigorev, 2017), Fake News Guard ^a , BuzzFace (Santia et al., 2018), FacebookHoax (Tacchini et al., 2017; Shu et al., 2020b), Sebastian et al. (Tschitschholek et al., 2018), Detective (Leskovec and Mcauley, 2012)

^a <https://www.eu-startups.com/directory/fake-news-guard/>.

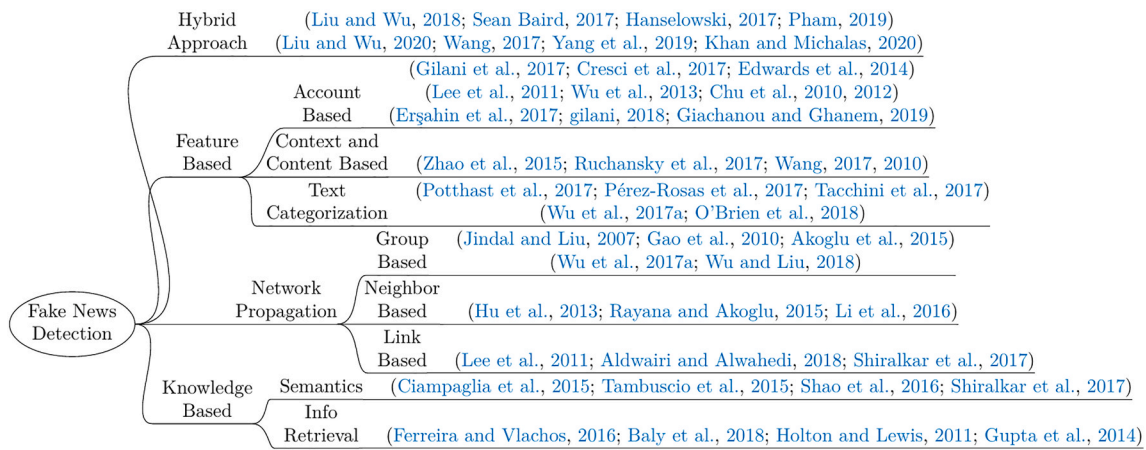


Fig. 1. Classification of the existing literature based on four paradigms – hybrid approach, feature based, network propagation, and knowledge based.

during these events is however still significant. Moreover, a lack of know-how about social networks makes it even more challenging to discern the credibility of shared information (Webwise, 2019). Antoniadis et al. (2015) developed a detection model to identify misinformation and suspicious behavioural patterns during emergency events on the Twitter platform. The model was based on a supervised learning technique using the user’s profile and tweets. The experiments were performed on a data set consisting of 59,660 users and 80,294 tweets. The authors filtered 81 percent of the tweets and claimed that more than 23 percent were misrepresentations. Although the proposed technique makes no distinction between intentional and unintentional information (Balestrucci et al., 2019), it successfully achieved timely detection.

In Table 3, we analyze trust and reputation models in terms of the mechanism used, data set as well as the outcomes and weaknesses of each model. In the existing literature, insufficient importance is given to the sources responsible for spreading the fake news. Evaluating the source is not straightforward process, as there are multiple variables to be considered in source verification, such as affiliation and reputation of the source, expertise in the domain, agreement or disapproval of other sources etc. Moreover, the absence of a source makes information unreliable, regardless of whether it is generated by an authentic source or not. Hence, fake news evaluation requires a model capable of performing source tracking, verification and validation.

4. Tools and data resources

Social media popularity, the availability of the internet, the extreme growth of user-generated website content, the lack of quality control and poor governance all provide fertile ground for sharing and spreading false and unverified information. This has led to continuous deterioration of information veracity. As the significance of the fake news problem is growing, the research community is proposing increasingly robust and accurate solutions. Some of the proposed solutions are discussed below and their characteristics are provided in Table 4.

- **BS-Detector²**: Available as a browser extension for both Mozilla and Chrome. It searches for all the links available on a webpage that are linked to unreliable sources and checks these links against a manually compiled list of domains. It can classify the domains as fake news, conspiracy theory, clickbait, extremely biased, satire, proceed with caution, etc. The BS detector has been downloaded and installed around about 25,000 times (Griswold, 2016).

- **FiB³**: The distribution of content is as important as its creation, FiB takes both post creation as well as distribution into account. It verifies the authenticity of a post in real time using AI. The AI uses keyword extraction, image recognition and source verification to check the authenticity of posts and provide a trust score. In addition, FiB tries to provide true information for posts that are deemed false (Figueira and Oliveira, 2017).
- **Trusted News add-on⁴**: Built in conjunction with MetaCertProtocol powered by the Metacert organization to help users spot suspicious or fake news. It is used to measure the credibility of website content and flags content as good, questionable or harmful. It gives a wider set of outputs, including marking website contents as malicious, satirical, trustworthy, untrustworthy, biased, clickbait and unknown (Walsh, 2019).
- **SurfSafe⁵**: There are different ways to analyze fake news such as textual analysis, image analysis, etc. Ash Bhat and Rohan Phadte focused on the analysis of fake news using images and generated a data set which consists of images collected from 100 fact-checking and trusted new sites. They developed a plug-in that checks the

² <https://gitlab.com/bs-detector/bs-detector>.

³ <https://devpost.com/software/fib>.

⁴ <https://trusted-news.com/>.

⁵ <https://chrome.google.com/webstore/detail/surfsafe-join-the-fight-a/hbpa-gabeiphkfhbboacggckhkipgdmh?hl=en>.

Table 7
Approaches taken by governments to tackle the problem of fake news.

Country	Focus	Approach/Action
Argentina	Fact-checking resources for public	– Commission created to verify fake news during national election campaign; – Imposing sanctions for spreading fake news.
Sweden	Foreign disinformation campaign	– Media broadcasts and publications are governed by law; – Educating citizens
Canada	Foreign disinformation campaign	– No specific law developed to prohibit the spread of fake news. Laws related to the criminal code or broadcasting distribution regulation may be relevant to spreading fake news.
China	Election misinformation	– Spreading fake news is a crime under China's criminal law; – Imposition of a fine and imprisonment; – Reliable information is published to systematically rebut fake news
Egypt	Media regulation	– Three domestic laws have been passed to regulate information distribution and its accuracy; – Imposing sanctions for spreading fake news
France	Election misinformation	– No specific law but there is general legislation against fake news; – Imposing sanctions for spreading fake news
Germany	Hate speech	– A number of civil and criminal laws exist for fake news; – Network enforcement act specific for fighting fake news
Israel	Foreign disinformation campaign	– High-level committee appointed by the president to examine the current law for threats and find ways to address them; – Imposing sanctions for spreading fake news
Japan	Media regulation	– A law exists to counter fake news; – Ministry of Communication and Internal Affairs work jointly to counter fake news
Kenya	Election misinformation	– Computer misuse and cyber-crime act has been passed, not yet in force; – Educating citizens
Malaysia	Election misinformation	– Malaysian anti-fake News Act 2018; – A fact-checking portal is operated by government agencies; – Imposing sanctions for spreading fake news
Nicaragua	Media regulation	– No specific law available, however some provisions can be found within the penal code and election law
Russia	Election misinformation	– Passed legislation that addresses the spread of fake news; – Imposing sanctions for spreading fake news
Brazil	Election misinformation	– No law but the topic is under discussion in congress; – Fines and imprisonment
United Kingdom	Foreign disinformation campaign	– No legislation to scrutinize or validate news on social media; – Reliable information is published to systematically rebut fake news
United Arab Emirates	Election misinformation	– Sharing misinformation is a crime by law; – Imposition of a fine
United States	disinformation, misinformation	– Proposed a federal law; – State media literacy initiatives

images against a generated data set. The main idea is to check each new image against the generated image data set. If the image is used in a fake context or modified, the information as a whole is considered fake (Clark).

- **BotOrNot**: A publicly available service used to assign a classification score to a Twitter account. This score is assigned to an account on the basis of the similarity it exhibits to the known characteristics of social bots. This classification system leverages more than 1000 features extracted from contents, interaction patterns and available metadata (Davis et al., 2016). These features are further grouped into six sub-classes:
 - Network features - built by extracting the statistical features for mentions, retweets, and hashtag co-occurrence.
 - User features - based on twitter metadata such as creation time of account, languages, locations.
 - Friend features - dependent on the statistics of social contacts such as number of followers, posts and so on.
 - Temporal features - recording the timing pattern for content generation and distribution.
 - Content features - based on part-of-speech tagging.
 - Sentiment features - built by using a sentiment analysis algorithm that takes into account happiness, emotion scores, etc.
- **Decodex**⁶: An online fake news detection tool that alerts the user to the potential of fake news by labeling the information as 'satire', 'info' and 'no information' (Gielczyk et al., 2019).
- **TrustyTweet**⁷: TrustyTweet is a browser plug-in, proposed for twitter users to assess and increase media literacy. It shifts the focus from fake news detection by labelling to supporting users to make their own assessment by providing transparent, neutral and intuitive hints when dealing with the fake news. TrustyTweet is based on gathering the potential indicators for fake news, already identified and proven to be promising in previous studies (Hartwig and Reuter, 2019).
- **Fake News Detector**⁸: The Fake News Detector is an open source project used for flagging news. A user can flag news as either fake news, extremely biased or clickbait. The user flagging activity is visible to other fake news detector users who may flag it again. Once the news is flagged, it is saved in the repository and accessible to Robhino – an ML robot trained on the inputs provided by humans that flags news automatically as clickbait, fake news or extremely biased news.
- **Fake News Guard**⁹: Available as a browser extension, it can verify the links displayed on Facebook or any page visited by the user. There is insufficient information about the way this tool works, however the key idea is that the "Fake news guard uses the AI technique along with network analysis and fact-checking".
- **TweetCred**¹⁰: A web browser tool used for assessing the credibility of tweets by using a supervised ranking algorithm trained on more than 45 features. TweetCred assigns a credibility score for each tweet on the user time line. Over the course of three months, TweetCred was installed 1127 times and computed the credibility score for 5.4 million tweets (Gupta et al., 2014).
- **LiT.RL News Verification**¹¹: A research tool that analyses the language used on web pages. The core functionality of the News Verification browser is textual data analysis using NLP and automatic classification using a SVM. It automatically detects and highlights website news as clickbait, satirical fake news and fabricated news (Rubin et al., 2019).

⁶ <https://chrome.google.com/webstore/detail/decodex/kbpkclapffgmndlaifaaalgkaagkfdod?hl=fr>.

⁷ <https://peasec.de/2019/trustytweet/>.

⁸ <https://github.com/fake-news-detector/fake-news-detector/tree/master/robhino>.

⁹ <http://www.eu-startups.com/directory/fake-news-guard/>.

¹⁰ <http://twitdigest.iiitd.edu.in/TweetCred/>.

¹¹ <https://victoriarubin.fims.uwo.ca/2018/12/19/release-for-the-lit-rl-news-verification-browser-detecting-clickbait-satire-and-falsified-news/>.

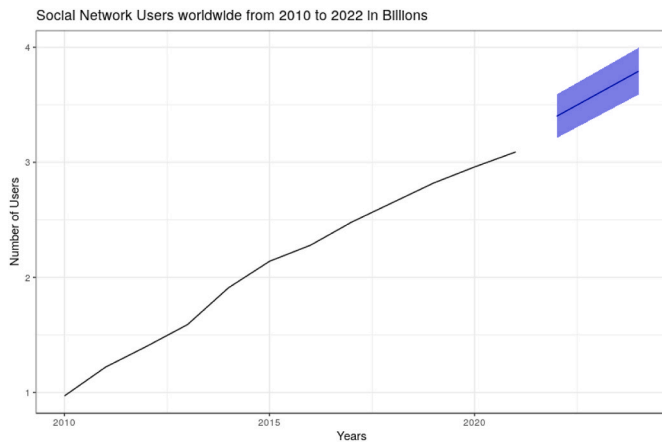


Fig. 2. Number of social users.

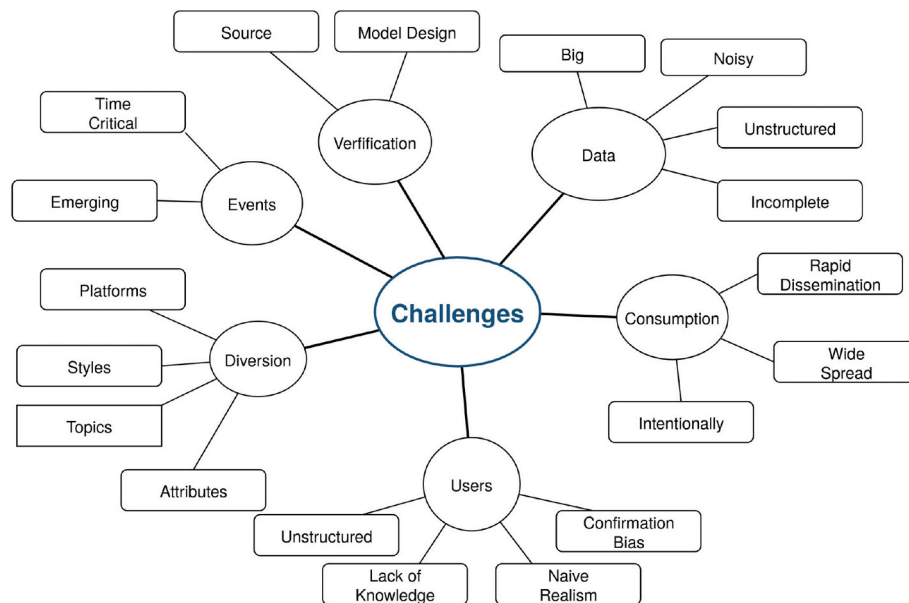


Fig. 3. Future challenges.

Detecting fake news on social media poses many challenges as most fake news is intentionally written. Researchers are considering different information, such as user behaviour, the engagement content of news, etc. to tackle the problem. However, there is no data set available that could provide the information on how fake news propagates, how different users interact with fake news, how to extract temporal features which could help to detect it and what the impact of fake news truly is.

In the previous section, we discussed the automatic detection of fake news using ML models. ML models require high quality data set to be efficient. This continues to be a major challenge when it comes to social media data due to its unstructured nature, high dimensionality, etc. In order to facilitate research in this field, a comprehensive guide to existing data sets is required. Below we present the details for some of the more widely used ones:

- **CredBank**: Collected by tracking more than 1 billion tweets between October 2014 and February 2015 (Mitra, 2016). It consists of tweets, events, topics and an associated credibility judgment assigned by humans. The data set comprises 60 million tweets which are further categorized into 1049 real-world events. Further, the data is spread into a streaming tweet file, topic file, credibility annotation file and searched tweet file (Mitra and Gilbert, 2015).

- **LIAR**: A publicly available fake news detection data set (Fernandes, 2019) that can be used for fact-checking. It consists of 12,836 short statements labelled manually by humans. In order to verify their truthfulness, each statement is evaluated by the editor of **POLITIFACT.COM**. Each statement is labelled in any of the following six categories: true, mostly-true, half-true, barely-true, false, pants on fire (Wang, 2017).
- **Memetracker9**: This data set (Leskovec et al., 2009a) recorded social media activity and online mainstream content over a three-month period. They used a Spinn3rAPI and collected 90 million documents from 165 million different websites (Leskovec et al., 2009b). The data set they generated is 350 GB in size. First, they extracted 112 million quotes, which were further refined and from which 22 million distinct phrases were collected.
- **FakeNewsNet¹²**: A multi-dimensional data repository consisting of social context, content and spatiotemporal information (Shu et al., 2020b). The data set was constructed using FakeNewsTracker, a tool used for collecting, analyzing as well as visualizing fake news. In the

given data set, the content consists of news, articles and images while context consists of information related to the user, post, response and network. The spatiotemporal information consists of spatial (user profile with location, tweets with location) and temporal information (timestamp for news and responses).

- **BuzzFeedNews¹³**: This data set, recorded all the news published by 9 news agencies on Facebook regarding the US election. The articles and news were fact-checked by journalists from BuzzFeed. It contains 1627 articles and 826 streams from hyperpartisan Facebook pages which publish misleading and false information at an alarming rate (News, 2016).
- **BS Detector (Vieira, 2017)**: This data set was collected by using BS Detector, a web browser extension for both Chrome and Mozilla. It is used to search all the links linked to unreliable sources on a given web page. These links are checked across a manually compiled list of domains.

¹² <https://github.com/KaiDMML/FakeNewsNet>.

¹³ <https://www.buzzfeednews.com/article/craigsilverman/partisan-fb-pages-analysis>.

- **BuzzFace:** This data set (Santia et al., 2018) consists of 2263 news articles and 1.6 million comments. Buzzface is based on extending the BuzzFeed data set by adding comments related to Facebook news articles. The news articles were categorized as “mostly true”, “mixture of true and false”, “mostly false” and “no factual comments”.
- **FacebookHoax:** In this data set (Tacchini et al., 2017), facebook graph API is used for the collection of data. It consists of 15,500 posts, of which 8923 are hoaxes and the rest are non-hoaxes. These posts were collected from 32 pages: 14 conspiracy and 18 scientific pages. In addition, the data set also includes the number of likes, which exceeds 2.3 millions.
- **Higgs-Twitter:** This data set (De Domenico et al., 2013) consists of Twitter posts related to the discovery of the new Higgs boson particle. The tweets were collected using the Twitter API. It consists of all the tweets that contain one of the following hashtags or keywords: cern, higgs, lhc, boson. The data set consists of 527,496 users and 985,590 analysed tweets of which 632,207 were geo-located tweets.
- **Trust and Believe:** This data set consists of 50000 Twitter users, all of whom were politicians (Khan, 2021). For each user, a unique profile is created containing 19 features. A total of 1000 user was manually annotated, with the rest being classified using an active learning approach.

Table 5 presents a detailed summary of the available data sets used for fake news detection in existing literature. Most are either small in size or contain mainly uni-modal data. The existing multi-modal data sets, unfortunately, still can't be used as a benchmark for training and testing models for fake news detection (Jindal et al., 2019). The next step is to generate large and comprehensive data sets that would include resources from which all relevant information could be extracted.

5. Discussion and challenges

Solving the problem of fake news detection and minimizing their impact on society is one of the important issues being considered in the research community. In this review, we analysed different studies using varying methods for detecting fake news. With the aim of aiding future research we provide their classification based on the social media platform used in Table 6.

Similarly, a study of the current literature on false news identification can be divided into four paradigms: hybrid approach, feature-based, network propagation and knowledge-based. The hybrid approaches employ both human and ML approaches for the detection of fake news. In the feature-based method, multiple features associated with a specific social media account are used to detect fake news. This paradigm can further be divided into three sub-categories – account-based, context and content-based and Text categorization. These methods are explicitly discussed in section 2. The third paradigm, network propagation, describes the potential methods for discovering, flagging and stopping the propagation of fake news in its infancy. The final paradigm entails supplementing AI models with human expert knowledge for decision-making (see section 2). An overview of these paradigms is given in Fig. 1.

Identifying and mitigating the spread of fake news and its variants presents a set of unique challenges. Fake news dissemination is a part of coordinated campaigns targeting a specific audience with the aim of generating a plausible impact on either local or global level. Many companies as well as entire countries were faced with the need to start building mechanisms to protect citizens from fake news. In September 2019, Facebook announced it was contributing \$10 million to a fund to improve deepfake detection technologies while several governments have taken different initiatives to defeat this problem (Funke Daniel, 2019; Rusu and Herman, 2019; G. L. R. D., 2019). Educational institutions and non-profit organizations have also tried to mitigate the problem through advocacy and literacy campaigns. Specifically, these

institutions in collaboration with technology companies have designed various techniques for detecting, flagging, and reporting fake news (Carlson, 2017; Northman, 2019; Sardarizadeh, 2019; Read, 2019).

Table 7 summarizes the actions that have been taken by governments around the world in order to battle the spread of fake news.

The greatest obstacle in fake news detection is that the information spreads through social media platforms like forest fire (especially if it's polarizing) which when not addressed, becomes viral in a matter of milliseconds (Stahl, 2018). The implications of this instantaneous consumption of information, on the other hand, are long-lasting. As a result, fake news becomes indistinguishable from real information, and the ongoing trends are difficult to recognize. We believe that fake news propagation can only be successfully controlled through early detection (see section 2.4). Another significant problem is that the rise in the influence of social media is closely connected to the increase in the number of users. According to Fig. 2, there are currently more than 3 billion users and by 2024 this number is expected to exceed 4 billion, a development that will eventually lead to an exponential rise in data (Tankovska, 2021). This data is most likely to be potentially uncertain due to inconsistencies, incompleteness, noise and unstructured nature. This complexity increases the velocity, variety, and amount of data and will most probably jeopardize the legitimacy of the results of any standard analytic processes and decisions that would be based on them. Analysis of such data requires tailor-made advanced analytical mechanisms. Designing techniques that could efficiently predict or evaluate future courses of action with high precision thus remains very challenging.

To summarize, humans are susceptible to becoming victims of false information due to their intrinsic way of processing and interpreting information being influenced by cognitive biases – namely, by the Truth Bias, Naive Realism and Confirmation Bias (Stahl, 2018). Consequently, all fake information floating around can lead to false information which is capable of ruining the “balance of news ecosystem”. The main challenge is that most users do not pay more attention to the manipulated information, while those who are manipulating it are systematically trying to create more confusion. The outcome of this process is that the people's ability to decipher real from false information is further impeded (Shu et al., 2017; Rubin, 2017).

Can we stop the viral spread?, the answer obviously is *Not yet* and it is because of the critical challenges surrounding the detection of fake news (see Fig. 3). Several efforts, however, have been put in place to help limit it such as media literacy. Media literacy comprised of practices that enable people to access and critically evaluate content across different media seems like the only valid solution. Although this is, and always was a challenging task, a coherent understanding, proper education, training, awareness and responsible media engagement could change this (Bulger and Davison, 2018). In the mean time, resisting disinformation and “fake news” culture should be promoted and encouraged. In addition, cross-disciplinary collaboration (i.e., social psychology, political science, sociology, communication studies etc.) can help and streamline findings across diverse disciplines to devise a holistic approach for understanding the media environment structure and how it operates.

6. Conclusion

Today, OSNs can be seen as platforms where people from all over the world can instantly communicate with strangers and even influence people's actions. Social media has shaped the digital world to an extent that they now seem like an indispensable part of our daily lives. However, social networks' ease of use has also revolutionized the generation and distribution of fake news. This prevailing trend has had a significant impact on our societies.

In this survey paper, we studied the problem of fake news detection from two different perspectives. Firstly, to assist users in identifying who they are interacting with, we looked at different approaches in existing

literature used for the identification and classification of user accounts. To this end, we analysed in depth both the users' context (anyone) and content (anything). For the early identification and mitigation of fake news, we studied different approaches that focus on data network features. Recently proposed approaches for measuring the relevance, credibility, and quality of sources were analysed in detail.

Secondly, we approached the problem of automating fake news detection by elaborating on the top three approaches used during fake news detection competitions and looked at the characteristics of more robust and accurate web-browsing tools. We also examined the statistical outputs, advantages, and disadvantages of some of the publicly available data sets. As the detection and prevention of fake news presents specific challenges, our conclusion identified potential challenges and promising research directions.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgement

This research has received funding from the EU research projects ASCLEPIOS (No. 826093) and CYBELE (No 825355).

References

- Adewole, K.S., Anuar, N.B., Kamsin, A., Varathan, K.D., Razak, S.A., 2017. Malicious accounts: dark of the social networks. *J. Netw. Comput. Appl.* 79, 41–67.
- Akoglu, L., Tong, H., Koutra, D., 2015. Graph based anomaly detection and description: a survey. *Data Min. Knowl. Discov.* 29 (3), 626–688.
- Albright, J., 2016. The #election2016 Micro-propaganda Machine. <https://medium.com/@d1gi/the-election2016-micro-propaganda-machine-383449cc1fba>.
- Aldwairi, M., Alwahedi, A., 2018. Detecting fake news in social media networks. *Procedia Computer Science* 141, 215–222.
- Allcott, H., Gentzkow, M., 2017. Social media and fake news in the 2016 election. *J. Econ. Perspect.* 31 (2), 211–236.
- Alothali, E., Zaki, N., Mohamed, E.A., Alashwal, H., 2018. Detecting social bots on twitter: a literature review. In: 2018 International Conference on Innovations in Information Technology (IIT). IEEE, pp. 175–180.
- Anna Escher, A.H., 2016. Wtf Is Clickbait? <https://techcrunch.com/2016/09/25/wtf-is-clickbait/>.
- Antoniadis, S., Litou, I., Kalogeraki, V., 2015. A model for identifying misinformation in online social networks. In: OTM Confederated International Conferences" on the Move to Meaningful Internet Systems". Springer, pp. 473–482.
- Asano, E., 2017. How Much Time Do People Spend on Social Media? *Social Media Today*, pp. 290–306.
- Aslam, S., 2021. Twitter by the Numbers: Stats, Demographics & Fun Facts. <https://www.omniaagency.com/twitter-statistics/>.
- Balestrucci, A., De Nicola, R., Inverso, O., Trubiani, C., 2019. Identification of credulous users on twitter. In: Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing, pp. 2096–2103.
- Baly, R., Karadzov, G., Alexandrov, D., Glass, J., Nakov, P., 2018. Predicting Factuality of Reporting and Bias of News Media Sources arXiv preprint arXiv:1810.01765.
- Benevenuto, F., Magno, G., Rodrigues, T., Almeida, V., 2010. Detecting spammers on twitter. In: Collaboration, Electronic Messaging, Anti-abuse and Spam Conference (CEAS), vol. 6, p. 12.
- Binham, C., 2019. Companies Fear Rise of Fake News and Social Media Rumours. <http://www.ft.com/content/4241a2f6-e080-11e9-9743-db5a370481bc>.
- Bollen, J., Mao, H., Pepe, A., 2011. Modeling public mood and emotion: twitter sentiment and socio-economic phenomena. In: Fifth International AAAI Conference on Weblogs and Social Media.
- Broder, A., Kumar, R., Maghoul, F., Raghavan, P., Rajagopalan, S., Stata, R., Tomkins, A., Wiener, J., 2000. Graph structure in the web. *Comput. Network.* 33 (1–6), 309–320.
- Bulger, M., Davison, P., 2018. The Promises, Challenges, and Futures of Media Literacy. Bytwerk, R.L., 2010. Grassroots Propaganda in the Third Reich: the Reich Ring for National Socialist Propaganda and Public Enlightenment. *German Studies Review*, pp. 93–118.
- Canini, K.R., Suh, B., Pirolli, P.L., 2011. Finding credible information sources in social networks based on content and social structure. In: 2011 IEEE Third International Conference on Privacy, Security, Risk and Trust and 2011 IEEE Third International Conference on Social Computing. IEEE, pp. 1–8.
- Carlson, E., 2017. Flagging Fake News. <https://niemanreports.org/articles/flagging-fake-news/>.
- Carson, J., 2019. Fake News: what Exactly Is it and How Can You Spot it? <https://www.telegraph.co.uk/technology/0/fake-news-exactly-has-really-had-influence/>.
- Cha, M., Haddadi, H., Benevenuto, F., Gummadi, K.P., 2010. Measuring user influence in twitter: the million follower fallacy. In: Fourth International AAAI Conference on Weblogs and Social Media.
- Chen, H., De, P., Hu, Y.J., Hwang, B.-H., 2014. Wisdom of crowds: the value of stock opinions transmitted through social media. *Rev. Financ. Stud.* 27 (5), 1367–1403.
- Chen, Y., Conroy, N.J., Rubin, V.L., 2015. Misleading online content: recognizing clickbait as "false news". In: Proceedings of the 2015 ACM on Workshop on Multimodal Deception Detection, pp. 15–19.
- Christian Reuter, J.K., Hartwig, Katrin, Schlegel, N., 2019. Fake news perception in Germany: a representative study of people's attitudes and approaches to counteract disinformation. In: 14th International Conference on Wirtschaftsinformatik.
- Chu, Z., Gianvecchio, S., Wang, H., Jajodia, S., 2010. Who is tweeting on twitter: human, bot, or cyborg?. In: Proceedings of the 26th Annual Computer Security Applications Conference. ACM, pp. 21–30.
- Chu, Z., Gianvecchio, S., Wang, H., Jajodia, S., 2012. Detecting automation of twitter accounts: are you a human, bot, or cyborg? *IEEE Trans. Dependable Secure Comput.* 9 (6), 811–824.
- Ciampaglia, G.L., Shiralkar, P., Rocha, L.M., Bollen, J., Menczer, F., Flammini, A., 2015. Computational fact checking from knowledge networks. *PLoS One* 10 (6), e0128193.
- Clark, B., SurfSafe Offers a Browser-Based Solution to Fake News (201b). <https://thextweb.com/insider/2018/08/21/surfsafe-offers-a-browser-based-solution-to-fake-news/>.
- Clauset, A., Shalizi, C.R., Newman, M.E., 2009. Power-law distributions in empirical data. *SIAM Rev.* 51 (4), 661–703.
- Collins, B., Hoang, D.T., Nguyen, N.T., Hwang, D., 2020. Trends in combating fake news on social media—a survey. *Journal of Information and Telecommunication* 1–20.
- Constantinides, P., Henfridsson, O., Parker, G.G., 2018. Introduction-platforms and Infrastructures in the Digital Age.
- Cresci, S., Di Pietro, R., Petrocchi, M., Spognardi, A., Tesconi, M., 2016. Dna-inspired online behavioral modeling and its application to spambot detection. *IEEE Intell. Syst.* 31 (5), 58–64.
- Cresci, S., Di Pietro, R., Petrocchi, M., Spognardi, A., Tesconi, M., 2017. The paradigm-shift of social spambots: evidence, theories, and tools for the arms race. In: Proceedings of the 26th International Conference on World Wide Web Companion, pp. 963–972.
- Cresci, S., Lillo, F., Regoli, D., Tardelli, S., Tesconi, M., 2019. Cashtag piggybacking: uncovering spam and bot activity in stock microblogs on twitter. *ACM Trans. Web* 13 (2), 11.
- Dauphin, Y.N., Fan, A., Auli, M., Grangier, D., 2017. Language modeling with gated convolutional networks. In: Proceedings of the 34th International Conference on Machine Learning, vol. 70. JMLR.org, pp. 933–941.
- Davis, C.A., Varol, O., Ferrara, E., Flammini, A., Menczer, F., 2016. Botnot: a system to evaluate social bots. In: Proceedings of the 25th International Conference Companion on World Wide Web. International World Wide Web Conferences Steering Committee, pp. 273–274.
- De Domenico, M., Lima, A., Mougel, P., Musolesi, M., 2013. The anatomy of a scientific rumor. *Sci. Rep.* 3, 2980.
- DFRLab, 2016. Human, Bot or Cyborg? <https://medium.com/@DFRLab/human-bot-or-cyborg-41273cdbl1e17>.
- Dimitriou, T., Michalas, A., 2012. Multi-party trust computation in decentralized environments. In: 2012 5th International Conference on New Technologies, Mobility and Security. NTMS, pp. 1–5. <https://doi.org/10.1109/NTMS.2012.6208686>.
- Dimitriou, T., Michalas, A., 2014. Multi-party trust computation in decentralized environments in the presence of malicious adversaries. *Ad Hoc Netw.* 15, 53–66. <https://doi.org/10.1016/j.adhoc.2013.04.013>.
- Edwards, C., Edwards, A., Spence, P.R., Shelton, A.K., 2014. Is that a bot running the social media feed? testing the differences in perceptions of communication quality for a human agent and a bot agent on twitter. *Comput. Hum. Behav.* 33, 372–376.
- Ericsson, K.A., Hoffman, R.R., Kozbelt, A., 2018. *The Cambridge Handbook of Expertise and Expert Performance*. Cambridge University Press.
- Erşahin, B., Aktaş, Ö., Kılınc, D., Akyol, C., 2017. Twitter fake account detection. In: 2017 International Conference on Computer Science and Engineering (UBMK). IEEE, pp. 388–392.
- Fernandes, T., 2019. Liardataset. https://github.com/thiagorainmaker77/liar_dataset.
- Ferreira, W., Vlachos, A., 2016. Emergent: a novel data-set for stance classification. In: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1163–1168.
- Figueira, Á., Oliveira, L., 2017. The current state of fake news: challenges and opportunities. *Procedia Computer Science* 121, 817–825.
- M. for minds. Spread of coronavirus fake news causes hundreds of deaths. <https://www.dw.com/en/coronavirus-misinformation/a-54529310>.
- Funke Daniel, F.D., 2019. A Guide to Anti-misinformation Actions Around the World. <https://www.poynter.org/ifcn/anti-misinformation-actions/>.
- Gabrovšek, P., Aleksovski, D., Mozetič, I., Grčar, M., 2017. Twitter sentiment around the earnings announcement events. *PLoS One* 12 (2).
- Gao, J., Liang, F., Fan, W., Wang, C., Sun, Y., Han, J., 2010. On community outliers and their efficient detection in information networks. In: Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, pp. 813–822.
- Garcia, D., Mavrodiev, P., Casati, D., Schweitzer, F., 2017. Understanding popularity, reputation, and social influence in the twitter society. *Pol. Internet* 9 (3), 343–364.
- Gazi, M.A., Çetin, M., 2017. The research of the level of social media addiction of university students. *International Journal of Social Sciences and Education Research* 3 (2), 549–559.

- Ghaviipour, M., Meybodi, M.R., 2018a. Trust propagation algorithm based on learning automata for inferring local trust in online social networks. *Knowl. Base Syst.* 143, 307–316.
- Ghaviipour, M., Meybodi, M.R., 2018b. A dynamic algorithm for stochastic trust propagation in online social networks: learning automata approach. *Comput. Commun.* 123, 11–23.
- Ghosh, S., Viswanath, B., Kooti, F., Sharma, N.K., Korlam, G., Benevenuto, F., Ganguly, N., Gummadi, K.P., 2012. Understanding and combating link farming in the twitter social network. In: *Proceedings of the 21st International Conference on World Wide Web*, pp. 61–70.
- Giachanou, A., Ghanem, B., 2019. Bot and gender detection using textual and stylistic information. *Pan* 16, 5.
- Giatsidis, C., Thilikos, D.M., Vazirgiannis, M., 2013. D-cores: measuring collaboration of directed graphs based on degeneracy. *Knowl. Inf. Syst.* 35 (2), 311–343.
- Gibert, D., Mateu, C., Planes, J., 2020. The rise of machine learning for detection and classification of malware: research developments, trends and challenges. *J. Netw. Comput. Appl.* 153, 102526.
- Gielczyk, A., Wawrzyniak, R., Choraś, M., 2019. Evaluation of the existing tools for fake news detection. In: *IFIP International Conference on Computer Information Systems and Industrial Management*. Springer, pp. 144–151.
- gilani, Z., 2018. STCS - Streaming Twitter Computation System. <https://github.com/zafargilani/stcs>.
- Gilani, Z., Kochmar, E., Crowcroft, J., 2017. Classification of twitter accounts into automated agents and human users. In: *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*. ACM, pp. 489–496.
- Grice, A., 2017. Fake News Handed Brexiters the Referendum and Now They Have No Idea what They're Doing. <https://www.independent.co.uk/voices/michael-gove-boris-johnson-brexit-euro-sceptic-press-theresa-may-a7533806.html>.
- Grier, C., Thomas, K., Paxson, V., Zhang, M., 2010. @ spam: the underground on 140 characters or less. In: *Proceedings of the 17th ACM Conference on Computer and Communications Security*. ACM, pp. 27–37.
- Grigorev, A., 2017. Identifying Clickbait Posts on Social Media with an Ensemble of Linear Models *arXiv preprint arXiv:1710.00399*.
- Griswold, A., 2016. Facebook Warned People that a Popular Fake News Detector Might Be "Unsafe". <https://qz.com/851894/facebook-said-bs-detector-a-plugin-to-detect-fake-news-might-be-unsafe/>.
- Gupta, A., Kumaraguru, P., Castillo, C., Meier, P., 2014. Tweetcred: real-time credibility assessment of content on twitter. In: *International Conference on Social Informatics*. Springer, pp. 228–243.
- Hannah Bastl, E.H., 2017. Bjorn Buchhold, Triple Scoring. <https://www.wsdm-cup-2017.org/triple-scoring.html>.
- Hanselowski, A., 2017. Team Athene on the Fake News Challenge. <https://medium.com/@andre134679/team-athene-on-the-fake-news-challenge-28a5cf5e017b>.
- Haralabopoulos, G., Anagnostopoulos, L., Zeadally, S., 2015. Lifespan and propagation of information in on-line social networks: a case study based on reddit. *J. Netw. Comput. Appl.* 56, 88–100.
- Hartwig, K., Reuter, C., 2019. Trustytweet: an indicator-based browser-plugin to assist users in dealing with fake news on twitter. In: *Proceedings of the International Conference on Wirtschaftsinformatik (WI)*.
- Hasani-Mavriqi, I., Kowald, D., Helic, D., Lex, E., 2018. Consensus dynamics in online collaboration systems. *Computational social networks* 5 (1), 2.
- Holton, A.E., Lewis, S.C., 2011. Journalists, social media, and the use of humor on twitter. *Electron. J. Commun.* 21 (1/2), 1–22.
- Hong, L., Dan, O., Davison, B.D., 2011. Predicting popular messages in twitter. In: *Proceedings of the 20th International Conference Companion on World Wide Web*, pp. 57–58.
- Hu, X., Tang, J., Zhang, Y., Liu, H., 2013. Social spammer detection in microblogging. In: *Twenty-Third International Joint Conference on Artificial Intelligence*.
- IONOS, D.G., 2018. Social Bots – the Technology behind Fake News. <https://www.ionos.com/digitalguide/online-marketing/social-media/social-bots/>.
- Jindal, N., Liu, B., 2007. Review spam detection. In: *Proceedings of the 16th International Conference on World Wide Web*. ACM, pp. 1189–1190.
- Jindal, S., Sood, R., Singh, R., Vatsa, M., Chakraborty, T., 2019. Newsbag: A Multimodal Benchmark Dataset for Fake News Detection.
- Kaplan, A.M., 2015. Social media, the digital revolution, and the business of media. *Int. J. Media Manag.* 17 (4), 197–199.
- Kaur, R., Singh, S., Kumar, H., 2018. Rise of spam and compromised accounts in online social networks: a state-of-the-art review of different combating approaches. *J. Netw. Comput. Appl.* 112, 53–88.
- Khan, T., Jan, 2021. Trust and Believe – Should We? Evaluating the Trustworthiness of Twitter Users. <https://doi.org/10.5281/zenodo.4428240>.
- Khan, T., Michalas, A., 2020. Trust and believe - should we? evaluating the trustworthiness of twitter users. In: *2020 IEEE 19th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*. IEEE Computer Society, Los Alamitos, CA, USA, pp. 1791–1800. <https://doi.org/10.1109/TrustCom50675.2020.00246>.
- Kharratzadeh, M., Coates, M., 2012. Weblog analysis for predicting correlations in stock price evolutions. In: *Sixth International AAAI Conference on Weblogs and Social Media*.
- Klyuev, V., 2018. Fake news filtering: semantic approaches. In: *2018 7th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions)(ICRITO)*. IEEE, pp. 9–15.
- Kong, X., Shi, Y., Yu, S., Liu, J., Xia, F., 2019. Academic social networks: modeling, analysis, mining and applications. *J. Netw. Comput. Appl.* 132, 86–103.
- Kshetri, N., Voas, J., 2017. The economics of "fake news". *IT Professional* 19 (6), 8–12.
- Kucharski, A., 2016. Study epidemiology of fake news. *Nature* 540 (7634), 525, 525.
- Kumar, S., West, R., Leskovec, J., 2016. Disinformation on the web: impact, characteristics, and detection of wikipedia hoaxes. In: *Proceedings of the 25th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, pp. 591–602.
- Lazer, D.M., Baum, M.A., Benkler, Y., Berinsky, A.J., Greenhill, K.M., Menczer, F., Metzger, M.J., Nyhan, B., Pennycook, G., Rothschild, D., et al., 2018. The science of fake news. *Science* 359 (6380), 1094–1096.
- Lee, S., Kim, J., 2013. Warningbird: a near real-time detection system for suspicious urls in twitter stream. *IEEE Trans. Dependable Secure Comput.* 10 (3), 183–195.
- Lee, K., Caverlee, J., Webb, S., 2010. Uncovering social spammers: social honeypots+ machine learning. In: *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, pp. 435–442.
- Lee, K., Eoff, B.D., Caverlee, J., 2011. Seven months with the devils: a long-term study of content polluters on twitter. In: *Fifth International AAAI Conference on Weblogs and Social Media*.
- Leskovec, J., McAuley, J.J., 2012. Learning to discover social circles in ego networks. In: *Advances in Neural Information Processing Systems*, pp. 539–547.
- Leskovec, J., Backstrom, L., Kleinberg, J., 2009a. Meme-tracking and the dynamics of the news cycle. In: *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, pp. 497–506.
- Leskovec, J., Backstrom, L., Kleinberg, J., 2009b. Meme-tracking and the dynamics of the news cycle. In: *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 497–506.
- Li, H., Chen, Z., Mukherjee, A., Liu, B., Shao, J., 2015. Analyzing and detecting opinion spam on a large-scale dataset via temporal and spatial patterns. In: *Ninth International AAAI Conference on Web and Social Media*.
- Li, J., Hu, X., Wu, L., Liu, H., 2016. Robust unsupervised feature selection on networked data. In: *Proceedings of the 2016 SIAM International Conference on Data Mining*. SIAM, pp. 387–395.
- Lim, E.-P., Nguyen, V.-A., Jindal, N., Liu, B., Lauw, H.W., 2010. Detecting product review spammers using rating behaviors. In: *Proceedings of the 19th ACM International Conference on Information and Knowledge Management*. ACM, pp. 939–948.
- Litou, I., Kalogeraki, V., Katakis, I., Gunopulos, D., 2016. Real-time and cost-effective limitation of misinformation propagation. In: *2016 17th IEEE International Conference on Mobile Data Management (MDM)*, vol. 1. IEEE, pp. 158–163.
- Liu, Y., Wu, Y.-F.B., 2018. Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. In: *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Liu, Y., Wu, Y.-F.B., 2020. Fned: a deep network for fake news early detection on social media. *ACM Trans. Inf. Syst.* 38 (3), 1–33.
- Liu, S., Liu, S., Ren, L., 2019. Trust or suspect? an empirical ensemble framework for fake news classification. In: *Proceedings of the 12th ACM International Conference on Web Search and Data Mining*, Melbourne, Australia, pp. 11–15.
- Ma, J., Gao, W., Mitra, P., Kwon, S., Jansen, B.J., Wong, K.-F., Cha, M., 2016. Detecting Rumors from Microblogs with Recurrent Neural Networks.
- Ma, J., Gao, W., Wong, K.-F., 2017. Detect Rumors in Microblog Posts Using Propagation Structure via Kernel Learning. *Association for Computational Linguistics*.
- Maigrot, C., Claveau, V., Kijak, E., Sicre, R., 2016. Mediaeval 2016: A Multimodal System for the Verifying Multimedia Use Task.
- Mao, Y., Wei, W., Wang, B., Liu, B., 2012. Correlating s&p 500 stocks with twitter data. In: *Proceedings of the First ACM International Workshop on Hot Topics on Interdisciplinary Social Networks Research*, pp. 69–72.
- Matsa, K.E., Shearer, E., 2018. News Use across Social Media Platforms 2018. *Pew Research Center*, 10.
- Michalas, A., Komninos, N., 2014. The lord of the sense: a privacy preserving reputation system for participatory sensing applications. In: *Computers and Communication (ISCC), 2014 IEEE Symposium*. IEEE, pp. 1–6.
- Michalas, A., Murray, R., 2017. Keep pies away from kids: a raspberry pi attacking tool. In: *Proceedings of the 2017 Workshop on Internet of Things Security and Privacy, IoTS&P '17*. ACM, New York, NY, USA, pp. 61–62. <https://doi.org/10.1145/3139937.3139953>.
- Mitra, E.G., 2016. Tanushree, CREDBANK-Data. <https://github.com/compsocial/CREDBANK-data>.
- Mitra, T., Gilbert, E., 2015. Credbank: a large-scale social media corpus with associated credibility annotations. In: *Ninth International AAAI Conference on Web and Social Media*.
- Morris, M.R., Counts, S., Roseway, A., Hoff, A., Schwarz, J., 2012. Tweeting is believing?: understanding microblog credibility perceptions. In: *Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work*. ACM, pp. 441–450.
- Nasim, M., Nguyen, A., Lothian, N., Cope, R., Mitchell, L., 2018. Real-time detection of content polluters in partially observable twitter networks. In: *Companion Proceedings of the the Web Conference 2018*, pp. 1331–1339.
- Neander, J., Marlin, R., 2010. Media and propaganda: the northcliffe press and the corpse factory story of world war i. *Global Media J.: Canadian Edition* 3 (2).
- News, B., 2016. Fact-Checking Facebook Politics Pages. <https://github.com/BuzzFeedNews/2016-10-facebook-fact-check>.
- Northman, T., 2019. Instagram Is Removing "Fake News" from the Platform. <https://hypebae.com/2019/8/instagram-fake-news-removing-tool-flagging-misinformation>.
- Omidvar, A., Jiang, H., An, A., 2018. Using neural network for identifying clickbaits in online news media. In: *Annual International Symposium on Information Management and Big Data*. Springer, pp. 220–232.
- Oshikawa, R., Qian, J., Wang, W.Y., 2018. A Survey on Natural Language Processing for Fake News Detection *arXiv preprint arXiv:1811.00770*.

- O'Brien, N., Latessa, S., Evangelopoulos, G., Boix, X., 2018. The Language of Fake News: Opening the Black-Box of Deep Learning Based Detectors.
- Pan, J.Z., Pavlova, S., Li, C., Li, N., Li, Y., Liu, J., 2018. Content based fake news detection using knowledge graphs. In: International Semantic Web Conference. Springer, pp. 669–683.
- Paschen, J., 2019. Investigating the emotional appeal of fake news using artificial intelligence and human contributions. *J. Prod. Brand Manag.*
- Perez, S., 2018. Twitter's Spam Reporting Tool Now Lets You Specify Type, Including if It's a Fake Account. <https://techcrunch.com/2018/10/31/twitters-spam-reporting-tool-now-lets-you-specify-type-including-if-its-a-fake-account/>.
- Pérez-Rosas, V., Kleinberg, B., Lefevre, A., Mihalcea, R., 2017. Automatic Detection of Fake News arXiv preprint arXiv:1708.07104.
- Perozzi, B., Al-Rfou, R., Skiena, S., 2014. Deepwalk: online learning of social representations. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 701–710.
- Pham, L., 2019. Transferring, Transforming, Ensembling: the Novel Formula of Identifying Fake News.
- Posetti, J., Matthews, A., 2018. A short guide to the history of 'fake news' and disinformation. *International Center for Journalists* 7, 2018–07.
- Potthast, M., Kiesel, J., Reinartz, K., Bevendorf, J., Stein, B., 2017. A Stylometric Inquiry into Hyperpartisan and Fake News arXiv preprint arXiv:1702.05638.
- Potthast, M., Gollub, T., Wiegmann, M., Stein, B., Hagen, M., Komlossy, K., Schuster, S., Fernandez, E.P.G., Jun. 2018. Webis Clickbait Corpus 2017 (Webis-clickbait-17). <https://doi.org/10.5281/zenodo.3346491>. URL <https://doi.org/10.5281/zenodo.3346491>.
- Rannard, G., 2020. Australia Fires: Misleading Maps and Pictures Go Viral.
- Rayana, S., Akoglu, L., 2015. Collective opinion spam detection: bridging review networks and metadata. In: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, pp. 985–994.
- Read, D., 2019. Social Media News: Fake News Flagging Tool, Clear Facebook History and More. <https://skedsocial.com/blog/fake-news-flagging-tool/>.
- Riedel, B., Augenstein, I., Spithourakis, G., Riedel, S., 2017. A Simple but Tough-To-Beat Baseline for the Fake News Challenge Stance Detection Task corr abs/1707.03264.
- Rieh, S.Y., Danielson, D.R., 2007. Credibility: a multidisciplinary framework. *Annu. Rev. Inf. Sci. Technol.* 41 (1), 307–364.
- Risdal, M., 2017. Getting Real about Fake News. <https://www.kaggle.com/mrisdal/fake-news>.
- Rubin, V.L., 2017. Deception detection and rumor debunking for social media. In: The SAGE Handbook of Social Media Research Methods. Sage, p. 342.
- Rubin, V.L., Chen, Y., Conroy, N.J., 2015. Deception detection for news: three types of fakes. In: Proceedings of the 78th ASIS&T Annual Meeting: Information Science with Impact: Research in and for the Community. American Society for Information Science, p. 83.
- Rubin, V., Brogly, C., Conroy, N., Chen, Y., Cornwell, S.E., Asubiaro, T.V., 2019. A news verification browser for the detection of clickbait, satire, and falsified news. *The Journal of Open Source Software* 4 (35), 1.
- Ruchansky, N., Seo, S., Liu, Y., 2017. Csi: a hybrid deep model for fake news detection. In: Proceedings of the 2017 ACM on Conference on Information and Knowledge Management. ACM, pp. 797–806.
- Ruiz, E.J., Hristidis, V., Castillo, C., Gionis, A., Jaimes, A., 2012. Correlating financial time series with micro-blogging activity. In: Proceedings of the Fifth ACM International Conference on Web Search and Data Mining, pp. 513–522.
- Rusu, M.-L., Herman, R.-E., 2019. Legislative measures adopted at the international level against fake news. In: International Conference KNOWLEDGE-BASED ORGANIZATION, vol. 25. Sciencdo, pp. 324–330.
- Santia, G.C., Williams, J.R., Buzzface, 2018. A news veracity dataset with facebook user commentary and egos. In: Twelfth International AAAI Conference on Web and Social Media.
- Sardarizadeh, S., 2019. Instagram Fact-Check: Can a New Flagging Tool Stop Fake News? <https://www.bbc.com/news/blogs-trending-49449005>.
- Saxena, A., Iyengar, S., Gupta, Y., 2015. Understanding spreading patterns on social networks based on network topology. In: Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015, pp. 1616–1617.
- Schwartz, O., 2018. Your favorite Twitter bots are about die, thanks to upcoming rule changes. <https://qz.com/1422765/your-favorite-twitter-bots-are-about-die-thanks-to-upcoming-rule-changes/>.
- Sean Baird, Y.P., 2017. Doug Sibley, Talos Targets Disinformation with Fake News Challenge Victory. <https://blog.talosintelligence.com/2017/06/talos-fake-news-challenge.html>.
- Shao, C., Ciampaglia, G.L., Flammini, A., Menczer, F., 2016. Hoaxy: a platform for tracking online misinformation. In: Proceedings of the 25th International Conference Companion on World Wide Web. International World Wide Web Conferences Steering Committee, pp. 745–750.
- Shiralkar, P., Flammini, A., Menczer, F., Ciampaglia, G.L., 2017. Finding streams in knowledge graphs to support fact checking. In: 2017 IEEE International Conference on Data Mining (ICDM). IEEE, pp. 859–864.
- Shu, K., Sliva, A., Wang, S., Tang, J., Liu, H., 2017. Fake news detection on social media: a data mining perspective. *ACM SIGKDD Explorations Newsletter* 19 (1), 22–36.
- Shu, K., Wang, S., Liu, H., 2018. Understanding user profiles on social media for fake news detection. In: 2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR). IEEE, pp. 430–435.
- Shu, K., Bhattacharjee, A., Alatawi, F., Nazer, T.H., Ding, K., Karami, M., Liu, H., 2020a. Combating disinformation in a social media age. *Wiley Interdisciplinary Reviews: Data Min. Knowl. Discov.* 10 (6), e1385.
- Shu, K., Mahudeswaran, D., Wang, S., Lee, D., Liu, H., 2020b. 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.
- Silverman, C., 2016. This Analysis Shows How Viral Fake Election News Stories Outperformed Real News on Facebook <https://www.buzzfeednews.com/article/craigsilverman/viral-fake-election-news-outperformed-real-news-on-facebook\#.hVVK8Br7G>.
- Stahl, K., 2018. Fake News Detection in Social Media, vol. 6. California State University Stanislaus.
- Stone-Gross, B., Cova, M., Cavallaro, L., Gilbert, B., Szydowski, M., Kemmerer, R., Kruegel, C., Vigna, G., 2009. Your botnet is my botnet: analysis of a botnet takeover. In: Proceedings of the 16th ACM Conference on Computer and Communications Security. ACM, pp. 635–647.
- P. Suci, More Americans Are Getting Their News From Social Media.
- Tacchini, E., Ballarin, G., Della Vedova, M.L., Moret, S., de Alfaro, L., 2017. Some like it Hoax: Automated Fake News Detection in Social Networks arXiv preprint arXiv:1704.07506.
- Tambuscio, M., Ruffo, G., Flammini, A., Menczer, F., 2015. Fact-checking effect on viral hoaxes: a model of misinformation spread in social networks. In: Proceedings of the 24th International Conference on World Wide Web. ACM, pp. 977–982.
- Tandoc Jr., E.C., Lim, Z.W., Ling, R., 2018. Defining “fake news” a typology of scholarly definitions. *Digital Journalism* 6 (2), 137–153.
- Tankovska, H., 2021. Number of Social Network Users Worldwide from 2017 to 2025. <https://www.statista.com/statistics/278414/number-of-worldwide-social-network-users/>.
- Thakur, A., 2016. Identifying Clickbaits Using Machine Learning. <https://www.linkedin.com/pulse/identifying-clickbaits-using-machine-learning-abhishek-thakur/>.
- G. L. R. D., 2019. The Law Library of Congress, 53K Rumors Spread in Egypt in Only 60 Days, Study Reveals. <https://www.loc.gov/law/help/fake-news/counter-fake-news.pdf>.
- Thomas, K., 2013. The Role of the Underground Economy in Social Network Spam and Abuse. Ph.D. thesis. UC Berkeley.
- Thomas, K., Grier, C., Ma, J., Paxson, V., Song, D., 2011. Design and evaluation of a real-time url spam filtering service. In: 2011 IEEE Symposium on Security and Privacy. IEEE, pp. 447–462.
- Torres, R., Gerhart, N., Negahban, A., 2018. Combating fake news: an investigation of information verification behaviors on social networking sites. In: Proceedings of the 51st Hawaii International Conference on System Sciences.
- Tschiatschhokle, S., Singla, A., Gomez Rodriguez, M., Merchant, A., Krause, A., 2018. Fake news detection in social networks via crowd signals. In: Companion Proceedings of the the Web Conference 2018. International World Wide Web Conferences Steering Committee, pp. 517–524.
- U. of Eastern Finland, 2019. New Application Can Detect Twitter Bots in Any Language. <https://phys.org/news/2019-06-application-twitter-bots-language.html>.
- Vieira, T., 2017. Bs-detector-dataset. <https://github.com/thiagovas/bs-detector-dataset>.
- Von Ahn, L., Blum, M., Langford, J., 2004. Telling humans and computers apart automatically. *Commun. ACM* 47 (2), 56–60.
- Vosoughi, S., Roy, D., Aral, S., 2018. The spread of true and false news online. *Science* 359 (6380), 1146–1151.
- Walsh, P., 2019. Factmata Trusted News Chrome Add-On Has Been Turned off until Further Notice. https://medium.com/@Paul_Walsh/factmata-trusted-news-chrome-add-on-has-been-turned-off-until-further-notice-7566f7312f86.
- Wanas, N., El-Saban, M., Ashour, H., Ammar, W., 2008. Automatic scoring of online discussion posts. In: Proceedings of the 2nd ACM Workshop on Information Credibility on the Web. ACM, pp. 19–26.
- Wang, A.H., 2010. Don't follow me: spam detection in twitter. In: 2010 International Conference on Security and Cryptography (SECRYPT). IEEE, pp. 1–10.
- Wang, W.Y., 2017. Liar, Liar Pants on Fire”: A New Benchmark Dataset for Fake News Detection arXiv preprint arXiv:1705.00648.
- Wang, Y., Yin, G., Cai, Z., Dong, Y., Dong, H., 2015. A trust-based probabilistic recommendation model for social networks. *J. Netw. Comput. Appl.* 55, 59–67.
- Webwise, 2019. Explained: what Is Fake News? <https://www.webwise.ie/teachers/wh-at-is-fake-news/>.
- Weerkamp, W., De Rijke, M., 2008. Credibility improves topical blog post retrieval. In: Proceedings of ACL-08. HLT, pp. 923–931.
- Weimer, M., Gurevych, I., Mühlhäuser, M., 2007. Automatically assessing the post quality in online discussions on software. In: Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions. Association for Computational Linguistics, pp. 125–128.
- Weng, J., Lim, E.-P., Jiang, J., He, Q., 2010. Twitterrank: finding topic-sensitive influential tweeters. In: Proceedings of the Third ACM International Conference on Web Search and Data Mining. ACM, pp. 261–270.
- Wong, J.I., 2016. Almost All the Traffic to Fake News Sites Is from Facebook. new data show. <https://qz.com/848917/facebook-fb-fake-news-data-from-jumpshot-its-the-biggest-traffic-referrer-to-fake-and-hyperpartisan-news-sites/>. (Accessed 1 January 2020).
- Wong, Q., 2019. Fake News Is Thriving Thanks to Social Media Users, Study Finds. <https://www.cnet.com/news/fake-news-more-likely-to-spread-on-social-media-study-finds/>.
- Wu, L., Liu, H., 2018. Tracing fake-news footprints: characterizing social media messages by how they propagate. In: Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining. ACM, pp. 637–645.
- Wu, X., Feng, Z., Fan, W., Gao, J., Yu, Y., 2013. Detecting marionette microblog users for improved information credibility. In: Joint European Conference on Machine Learning and Knowledge Discovery in Databases. Springer, pp. 483–498.

- Wu, L., Hu, X., Morstatter, F., Liu, H., 2017a. Adaptive spammer detection with sparse group modeling. In: Eleventh International AAAI Conference on Web and Social Media.
- Wu, L., Hu, X., Morstatter, F., Liu, H., 2017b. Detecting camouflaged content polluters. In: Eleventh International AAAI Conference on Web and Social Media.
- Xue, J., Yang, Z., Yang, X., Wang, X., Chen, L., Dai, Y., 2013. Votetrust: leveraging friend invitation graph to defend against social network sybils. In: 2013 Proceedings IEEE INFOCOM. IEEE, pp. 2400–2408.
- Yan, J., 2006. Bot, cyborg and automated turing test. In: International Workshop on Security Protocols. Springer, pp. 190–197.
- Yang, Z., Wilson, C., Wang, X., Gao, T., Zhao, B.Y., Dai, Y., 2014. Uncovering social network sybils in the wild. *ACM Trans. Knowl. Discov. Data* 8 (1), 2.
- Yang, K.-C., Niven, T., Kao, H.-Y., 2019. Fake News Detection as Natural Language Inference.
- Yaraghi, N., 2019. How Should Social Media Platforms Combat Misinformation and Hate Speech? <https://www.brookings.edu/blog/techtank/2019/04/09/how-should-social-media-platforms-combat-misinformation-and-hate-speech/#cancel>.
- Ye, S., Wu, S.F., 2010. Measuring message propagation and social influence on twitter. com. In: International Conference on Social Informatics. Springer, pp. 216–231.
- Ye, J., Kumar, S., Akoglu, L., 2016. Temporal opinion spam detection by multivariate indicative signals. In: Tenth International AAAI Conference on Web and Social Media.
- Zhao, Z., Resnick, P., Mei, Q., 2015. Enquiring minds: early detection of rumors in social media from enquiry posts. In: Proceedings of the 24th International Conference on World Wide Web. International World Wide Web Conferences Steering Committee, pp. 1395–1405.
- Zhou, Y., 2017. Clickbait Detection in Tweets Using Self-Attentive Network arXiv preprint arXiv:1710.05364.
- Zhou, X., Zafarani, R., 2019. Network-based fake news detection: a pattern-driven approach. *ACM SIGKDD Explorations Newsletter* 21 (2), 48–60.
- Zubiaga, A., Aker, A., Bontcheva, K., Liakata, M., Procter, R., 2018. Detection and resolution of rumours in social media: a survey. *ACM Comput. Surv.* 51 (2), 32.



Tanveer Khan received the Master Degree in Information Security from COMSATS University Islamabad Pakistan. After his Master's, he worked as a Data Analyst on the project CYBER Threat Intelligence Platform at COMSATS University, Islamabad, Pakistan. He also worked as a Junior analyst at Trillium Infosec, Pakistan. Currently, he is working as a Ph.D., Researcher at the Department Computing Sciences, at Tampere University, Finland. He is also a member of Network and Information Security Group (NISEC) at Tampere University, Finland. His interest is in privacy-preserving machine learning, fake news detection in social networks, cyber security, digital forensics and malware analysis.



Prof. Antonis Michalakis received his PhD in Network Security from Aalborg University, Denmark and he is currently working as an Assistant Professor at the Department Computing Sciences, at Tampere University, Finland where he also coleads the Network and Information Security Group (NISEC). The group comprises Ph.D., students, professors and researchers. Group members conduct research in areas spanning from the theoretical foundations of cryptography to the design and implementation of leading edge efficient and secure communication protocols. Apart from his research work at NISEC, as an assistant professor he is actively involved in the teaching activities of the University. Finally, his role expands to student supervision and research projects coordination. Furthermore, Antonis has published a significant number of papers in field related journals and conferences and has participated as a speaker in various conferences and workshops. His research interests include private and secure e-voting systems, reputation systems, privacy in decentralized environments, cloud computing, trusted computing and privacy preserving protocols in eHealth and participatory sensing applications.



Adnan Akhuzada is a Senior Member, IEEE with extensive 12 years of R&D experience both in ICT industry and academia having proven track record of high impact published research (i.e., US Patents, Journals, Transactions, Books, Reputable Magazines, Chapters, Conference and Conference Proceedings) and research funding. His experience as an educator & researcher is diverse that includes work as a Lecturer, a Senior Lecturer, a Year Tutor, as an Assistant Professor at COMSATS University, as a Senior Researcher at RISE SICs Vasteras AB, Sweden, and currently as a research fellow & WP lead at Technical University of Denmark (DTU) having mentor-ship of graduate students, and supervision of academic and R&D projects both at UG and PG level. He has been rigorously involved in international accreditation such as Accreditation Board for Engineering and Technology (ABET), and curriculum development according to the latest guidelines of ACM/IEEE. He also served as an In-charge BS Software Engineering (SE) and Telecommunication Networks (TN) Programme at Comsats University. He is currently involved in various EU and Swedish funded projects of cyber security. He is a member of technical programme committee of varied reputable conferences and editorial boards. He is presently serving as an associate editor of IEEE Access.