

A Survey on Predicting the Factuality and the Bias of News Media

Anonymous ACL submission

Abstract

The present level of proliferation of fake, biased, and propagandistic content online has made it impossible to fact-check every single suspicious claim or article, either manually or automatically. Thus, many researchers are shifting their attention to higher granularity, aiming to profile entire news outlets, which makes it possible to detect likely “fake news” the moment it is published, by simply checking the reliability of its source. Source factuality is also an important element of systems for automatic fact-checking and “fake news” detection, as they need to assess the reliability of the evidence they retrieve online. Political bias detection, which in the Western political landscape is about predicting left-center-right bias, is an equally important topic, which has experienced a similar shift towards profiling entire news outlets. Moreover, there is a clear connection between the two, as highly biased media are less likely to be factual; yet, the two problems have been addressed separately.

In this survey, we review the state of the art on media profiling for factuality and bias, arguing for the need to model them jointly. We further discuss interesting recent advances in using different information sources and modalities, which go beyond the text of the articles the target news outlet has published. Finally, we discuss current challenges and outline future research directions.

1 Introduction

The rise of the Web has made it possible for anybody to create a website and to become a *news medium*. This was a hugely positive development as it elevated freedom of expression to a whole new level, allowing anybody to have their voice heard. With the subsequent rise of social media, anybody could potentially reach out to a vast audience, something that until recently was only possible for major news outlets. One of the consequences was a *trust crisis*: with traditional news media stripped off

their gate-keeping role, the society was left unprotected against potential manipulation.

The issue became a general concern in 2016, a year marked by micro-targeted online disinformation at an unprecedented scale in connection to Brexit and the US Presidential election. These developments gave rise to the term “fake news.”

In an attempt to solve the trust problem, several initiatives, such as PolitiFact, Snopes, FactCheck, and Full Fact, have been launched to fact-check suspicious claims manually. However, given the scale of the proliferation of false information online, it was unfeasible to fact-check every single suspicious claim, even when this was done automatically, not only for computational reasons but also due to timing. In order to fact-check a claim manually or automatically, it is required to verify the stance of mainstream media concerning that claim and/or the reaction of users on social media. Accumulating this evidence takes time, and delay means more potential sharing of the malicious content. A study has shown that, for some very viral claims, more than 50% of the sharing happens within the first ten minutes after posting the micro-post on social media (Zaman et al., 2014), and thus timing is of utmost importance. Moreover, an extensive recent study has found that “fake news” spreads six times faster and reaches much farther than real news (Vosoughi et al., 2018).

A much more promising alternative is to profile the medium that initially published the news article with a suspicious claim. Since media that have published fake or biased content in the past are more likely to do so in the future, profiling media in advance makes it possible to detect likely “fake news” the moment it is published by simply checking the reliability of its source.

Estimating the reliability of a news source is important for claim fact-checking (Nguyen et al., 2018), and it also gives an important prior when solving article-level tasks such as “fake news”

084 and click-bait detection (Hardalov et al., 2016; 134
085 Karadzhov et al., 2017a; De Sarkar et al., 2018; 135
086 Pérez-Rosas et al., 2018; Brill, 2001; Finberg et al., 136
087 2002; Pan et al., 2018; Nguyen et al., 2022). 137

088 There have been several surveys on fake 138
089 news (Shu et al., 2017; da Silva et al., 2019; Zhou 139
090 and Zafarani, 2020), mis/dis-information (Islam 140
091 et al., 2020; Alam et al., 2022; Hardalov et al., 141
092 2022a), fact-checking (Thorne and Vlachos, 2018a; 142
093 Kotonya and Toni, 2020; Nakov et al., 2021; Guo 143
094 et al., 2022a), truth discovery (Li et al., 2016), and 144
095 propaganda detection (Martino et al., 2020). How- 145
096 ever, they have focused on claims or articles, while 146
097 here we survey research on profiling entire news 147
098 outlets for factuality and bias. 148

099 2 Factuality 150

100 Veracity of information has been studied at dif- 151
101 ferent levels: (i) claim-level (e.g., *fact-checking*), 152
102 (ii) article-level (e.g., *“fake news” detection*), 153
103 (iii) user-level (e.g., *hunting for trolls*), and 154
104 (iv) medium-level (e.g., *source reliability estima-* 155
105 *tion*). Our primary interest here is in the latter. 156

106 At the claim-level, significant effort has been 157
107 paid to fact-checking and rumor detection using 158
108 information from social media, i.e., how users re- 159
109 ply to the claim (Canini et al., 2011; Castillo et al., 160
110 2011; Ma et al., 2015, 2016; Zubiaga et al., 2016; 161
111 Ma et al., 2017; Dungs et al., 2018; Kochkina et al., 162
112 2018; Hardalov et al., 2022b; Nguyen et al., 2022), 163
113 but there is a need for more comprehensive ap- 164
114 proaches (Thorne and Vlachos, 2018b; Guo et al., 165
115 2022b). A set of web pages and snippets from 166
116 search engines have also been used as a source of 167
117 information (Mukherjee and Weikum, 2015; Popat 168
118 et al., 2016, 2017; Karadzhov et al., 2017b; Mi- 169
119 haylova et al., 2018; Baly et al., 2018b). In either 170
120 case, the most important information for the claim- 171
121 level tasks are *stance* (does a tweet or a news article 172
122 agree or disagree with the claim?) and *source reli-* 173
123 *ability* (do we trust the user who posted the tweet 174
124 or the medium that published the news article?). 175

125 The problem of source reliability remains largely 176
126 under-explored. In the case of social media and 177
127 community fora, it concerns modeling the user, 178
128 e.g., there has been research on finding opinion ma- 179
129 nipulation *trolls* (Mihaylov and Nakov, 2016), *sock-* 180
130 *puppets* (Maity et al., 2017), *Internet water army* 181
131 (Chen et al., 2013), and *seminar users* (Darwish 182
132 et al.). In the case of the Web, it is about source 183
133 trustworthiness (the URL domain, the medium).

In early work, the source reliability of news me-
dia has often been estimated automatically based
on the general stance of the target medium with
respect to known true/false claims, without access
to gold labels about the overall medium-level fac-
tuality of reporting (Dong et al., 2015; Mukherjee
and Weikum, 2015; Popat et al., 2016, 2017, 2018).

More recent work has addressed the task as one
on its own right. Baly et al. (2018a) used gold
labels from Media Bias/Fact Check, and rich infor-
mation sources: articles published by the medium,
what is said about it on Wikipedia, metadata from
its Twitter profile, URL structure, and traffic infor-
mation. In follow-up work, Baly et al. (2019)
used the same representation to jointly predict me-
dia factuality and bias on an ordinal scale, using
a multi-task ordinal regression setup. Then, Baly
et al. (2020b) extended the information sources
to include Facebook followers and speech signals
from the news medium’s channel on YouTube (if
any). Hounsel et al. (2020) proposed to use domain,
certificate, and hosting information of the website
infrastructure. Finally, Panayotov et al. (2022) used
audience overlap and graph neural networks.

3 Bias 158

159 Compared to factuality, which can be objectively 160
161 determined by whether a piece of information is 161
162 true or not, media bias has more complex dimen- 162
163 sions. For the last few decades, many scholars have 163
164 conceptualized media bias in different ways. For 164
165 instance, a bias can be defined as “imbalance or 165
166 inequality of coverage rather than as a departure 166
167 from truth” (Stevenson et al., 1973). A departure 167
168 from truth, however, can be measured only when 168
169 the accurate record of the event is available (e.g., 169
170 trial transcript and reporting).

A different definition, namely “any systematic
slant favoring one candidate or ideology over an-
other” (Waldman and Devitt, 1998), is proposed
to capture various dimensions rather than coverage
imbalance, such as favorability conveyed in visual
representations (i.e., news photos). E.g., smiling,
speaking at the podium, cheering crowd, and eye-
level shots are preferred over frowning, sitting, be-
ing alone, and shots from above, respectively.

D’Alessio and Allen (2000) reviewed 59 studies
about partisan media bias in presidential elections.
They proposed to categorize media bias into the
following three types: (i) *gatekeeping bias*, where
editors and journalists ‘select’ the stories to report,

(ii) *coverage bias*, where the amount of news coverage (e.g., the length of newspapers articles, or the time given on television) each party receives is systematically biased to one party at the expense of the other one, and (iii) *statement bias*, where news media interject their attitudes or opinions in the news reporting. Groeling (2013) proposed a more relaxed concept of media bias, which is “a portrayal of reality that is significantly and systematically (not randomly) distorted,” to take a variety of media bias dimensions into account. In particular, he focused on two main forms of media bias—*selection bias* (i.e., what to cover) and *presentation bias* (i.e., how to cover it)—driven by the choices of newsmakers.

Selection bias or *gatekeeping* bias has been studied in various ways, including qualitative interviews or surveys of journalists and editors about the decision-making process they use to select the stories in their newsroom (Tandoc Jr, 2014). Here, news selection is not necessarily confined to political context. News reporting about any news items can be considered as the unit of analysis.

Data-driven research on selection bias commonly follows three steps: (i) collect news articles (for newspapers or online news) or transcripts (for TV news) for a target period, (ii) conduct content analysis to find the news coverage of politicians, parties, or events. Optionally, study the tone of the news articles (e.g., negative news are more frequently reported) (Soroka, 2012), and (iii) identify systematic biases by comparing news coverage. An exhaustive database of news stories is thus essential for selection bias research. While commercial databases, such as Lexis Nexis, have been widely used (Soroka, 2012; Padgett et al., 2019; Gilens and Hertzman, 2000; Boykoff and Boykoff, 2004), publicly available datasets, such as GDELT, start to get attention (Boudemagh and Moise, 2017; Kwak and An, 2014; Boudemagh and Moise, 2017) and are getting validated by comparing multiple sources (Kwak and An, 2016; Weaver and Bimber, 2008; Kwak and An, 2016).

Presentation bias has been characterized from diverse perspectives, including framing (Entman, 2007), visuals (Barrett and Barrington, 2005), sources (Baum and Groeling, 2008), tone (Soroka, 2012), and more. Particularly, framing bias has been actively studied in many disciplines.

Framing Bias refers to a bias that highlights a certain aspect of an event or an issue more than the others (Entman, 1993). Emphasizing a particular aspect can deliver a distorted view toward the issue even without the use of biased expressions.

Framing biases have been typically studied at issue level (Kim and Johnson, 2022). Researchers collect news articles about an issue or an event, conduct manual content analysis, and build a frame detection model (Baumer et al., 2015). Open-source tools to help the analysis have been proposed (Bhattachia et al., 2021; Morstatter et al., 2018). While this approach can characterize diverse frames, it is not trivial to compare framing across issues.

The Media Frames Corpus (MFC) was proposed to address this limitation (Card et al., 2015). It contains articles annotated with 15 generic frames (including *others*) across three policy issues. Several studies have demonstrated reasonable prediction performance of the general media frames with different datasets (Field et al., 2018; Kwak et al., 2020). These 15 general frames were also used for analyzing political discourse on social media (Johnson et al., 2017). These frames are often customized to a specific issue by adding issue-specific frames (Liu et al., 2019), even though doing so somewhat contradicts the original motivation of general media frames, namely to be able to compare frames across various issues.

News slant was proposed to characterize how framing in news reports favors one side over the other (Entman, 2007). The media-level slant thus could differ across issues (Ganguly et al., 2020).

A variety of methods have been proposed to quantify the extent of news slant in traditional news media by (i) linking media outlets to politicians with known political positions, (ii) directly analyzing news content, and (iii) using shared audience among media outlets. Groseclose and Milyo (2005) assigned an ADA (Americans for Democratic Action) score for each media outlet by investigating co-citations of think-tanks by members of Congress and media outlets. Gentzkow and Shapiro (2010) proposed an ideological slant index of news media in a seminal study. The news slant is measured by the extent of phrases in news coverage that are more frequently used by one political party (i.e., Democratic or Republican) congress members than by another one in the 2005 Congress Record. Their frequency-based approach successfully finds politically charged phrases such as *death tax* or *war on*

284 *terror* by Republicans and associated media and
285 *estate tax* or *war in Iraq* by Democrats and associ-
286 ated media, and they further computed media slant
287 index for 433 newspapers. The choice of words
288 by political party members and news media is con-
289 sidered framing because they purposely highlight
290 some aspect of the issue over other ones.

291 An et al. (2012) proposed a method to compute
292 media slant scores by measuring distances between
293 media sources by their mutual followers on Twitter
294 (An et al., 2011, 2012). Stefanov et al. (2020)
295 identified the political leanings of media outlets and
296 influential people on Twitter based on their stance
297 on controversial topics. They built clusters of users
298 around core vocal ones based on their behaviour
299 on Twitter such as retweeting, using a procedure
300 proposed in (Darwish et al., 2020).

301 **Left-center-right bias (or left-right bias)** was
302 studied based on media-level annotation from spe-
303 cialized online platforms, such as News Guard,
304 AllSides, and Media Bias/Fact Check, where jour-
305 nalists use carefully designed guidelines to make
306 the judgments. Researchers have then trained sys-
307 tems to predict this bias using a variety of informa-
308 tion sources such as analyzing the corresponding
309 YouTube channels (Dinkov et al., 2019), and using
310 information from the articles the target news outlet
311 has published, what is there about them in social
312 media and in Wikipedia (Baly et al., 2020b).

313 There has also been work on predicting the left-
314 center-right bias of articles, which is somewhat
315 relevant here as it can be an element of media-level
316 analysis. Such systems are typically trained us-
317 ing distant supervision, projecting the label from a
318 medium to each article from that medium, which is
319 an easy way to obtain large datasets, needed to train
320 contemporary deep learning models. For example,
321 Kulkarni et al. (2018) used site-level annotations
322 from the AllSides website for political bias detec-
323 tion. The same approach was used to study hy-
324 perpartisanship, i.e., extremely one-sided reporting
325 (Potthast et al., 2018), as part SemEval-2019 task 4
326 on Hyper-partisan News Detection (Kiesel et al.,
327 2019). More recent work has demonstrated the
328 dangers of distant supervision and has introduced a
329 dataset for left-center-right bias with proper manual
330 article-level annotations (Baly et al., 2020a).

331 4 Joint Modeling

332 There is a well-known connection between factual-
333 ity and bias. For example, hyper-partisanship (high

334 bias) is often linked to low trustworthiness (Pot-
335 thast et al., 2018), e.g., appealing to emotions rather
336 than sticking to the facts, while center media tend
337 to be generally more impartial and also more trust-
338 worthy. Moreover, some of the datasets used for
339 the two tasks have media-level annotations for both
340 factuality and bias. Thus, it makes sense to model
341 factuality and bias jointly.

342 Yet, joint modeling of the two tasks remains
343 severely underexplored. In fact, there has been
344 a single attempt at doing so to date: Baly et al.
345 (2019) proposed a multi-task learning formulation.
346 They further took into account the ordinal nature
347 of the labels for both tasks, noting that classifying
348 an *extreme right* medium as *extreme-left* is a huge
349 error, while classifying it as a *center* is a smaller
350 one, and predicting *right* is an even smaller error.
351 Similarly, predicting a *high-factuality* label for a
352 *low-factuality* medium is a bigger mistake than pre-
353 dicting *mixed factuality*. Thus, they proposed a
354 multi-task ordinal regression model, copula ordi-
355 nal regression (Walecki et al., 2016), which jointly
356 predicts factuality and bias on ordinal scales. They
357 further used several auxiliary tasks, modeling cen-
358 trality, hyper-partisanship, as well as left-vs.-right
359 bias on a coarse-grained scale.

360 This is challenging as it requires understanding
361 the interactions between the two dimensions. Al-
362 though the relationship between extreme bias and
363 low factuality follows intuition, uncovering the con-
364 nection between being factual but biased or non-
365 factual but unbiased requires more detailed insights.
366 For news media that exhibit a mixed behavior in
367 both aspects, this poses an even greater difficulty.

368 5 Basis of Prediction

369 5.1 Textual Content

370 5.1.1 Representation

371 The most natural representation for a source is as
372 a sample of articles it has published, which in turn
373 can be represented using linguistic features or as
374 continuous representations.

375 *Linguistic Features* focus on language use, and
376 they have been shown to be useful for detecting
377 fake articles, as well as for predicting the political
378 bias and the factuality of reporting of news media
379 (Horne et al., 2018; Baly et al., 2018a). For ex-
380 ample, Horne and Adali (2017) showed that “fake
381 news” pack a lot of information in the title (as many
382 people do not read beyond the title, e.g., in social
383 media), and use shorter, simpler, and repetitive con-

384 tent in the body (as writing fake information takes a
385 lot of effort). Such features can be calculated based
386 on the Linguistic Inquiry and Word Count (LIWC)
387 lexicon and used to distinguish articles from trusted
388 sources vs. hoaxes vs. satire vs. propaganda
389 (pen). They can be also modeled using linguistic
390 markers (Mihaylova et al., 2018) such as *factives*
391 from (Hooper, 1975), *assertives* from (Hooper,
392 1975), *implicatives* from (Karttunen, 1971), *hedges*
393 from (Hyland, 2005), *Wiki-bias* terms from (Re-
394 casens et al., 2013), *subjectivity* cues from (Riloff
395 and Wiebe, 2003), and *sentiment* cues from (Liu
396 et al., 2005). There are 141 such features in the
397 NELA toolkit (Horne et al., 2018):

- **Style:** part-of-speech tags, use of specific words (function words, pronouns, etc.), and features for clickbait title classification;
- **Complexity:** type-token ratio, readability, number of cognitive process words (identifying discrepancy, insight, certainty, etc.);
- **Bias:** features modeling bias using lexicons (Recasens et al., 2013; Mukherjee and Weikum, 2015) and subjectivity, calculated by pre-trained classifiers (Horne et al., 2017);
- **Affect:** sentiment scores from lexicons (Recasens et al., 2013; Mitchell et al., 2013) and full systems (Hutto and Gilbert, 2014);
- **Morality:** features based on the Moral Foundation Theory (Graham et al., 2009) and lexicons (Lin et al., 2018);
- **Event:** features modeling time and location.

415 *Embedding representations:* An alternative way to
416 represent an article is to use embedding representa-
417 tions, typically based on large pre-trained language
418 models, such as BERT (Devlin et al., 2019). This
419 can be done without fine-tuning, e.g., by encod-
420 ing an article (possibly truncated, e.g., BERT can
421 take up to 512 tokens as an input) and then av-
422 eraging the word representations extracted from
423 the second-to-last layer. Alternatively, one can
424 use pre-trained sentence encoders such as Sentence
425 BERT (Reimers and Gurevych, 2019). Finally, one
426 can obtain representations that are relevant to the
427 target task, e.g., by fine-tuning BERT to predict
428 the label (bias or factuality) of the medium that an
429 article comes from, in the form of distant supervi-
430 sion (Baly et al., 2020b). One issue with distant
431 supervision is that the model can end up learning
432 to detect the source of the target news article in-
433 stead of predicting its factuality and bias, which
434 can be fixed using adversarial media adaptation and

a specially adapted triplet loss (Baly et al., 2020a).

5.1.2 Aggregation

In order to obtain a representation/prediction for an entire medium, there is a need to aggregate the representations/predictions for its articles.

Averaging article-level representations: One could average the representations for all articles to obtain a representation for a medium, which can then be used in a medium-level classifier. Using arithmetic averaging is a good idea as it captures the general trend of articles in a medium, while limiting the impact of outliers. For instance, if a medium is known to align with left-wing ideology, this should not change if it published a few articles that align with right-wing ideology.

Aggregating posterior probabilities: Alternatively, each article can be represented by a \mathcal{C} -dimensional vector that corresponds to its posterior probabilities of belonging to each class c_i , $i \in \{1, \dots, \mathcal{C}\}$ of the given task, whether it is predicting the political bias or the factuality of the target news medium. Finally, these article-level posterior probabilities are averaged in order to aggregate them at the medium level.

5.2 Multimedia Content

Nowadays, almost all news websites heavily rely on multimedia content. This dependence, however, also makes multimedia a very effective means for dispensing an intended, and even manipulated, messages. The increasing availability of automated and AI-powered multimedia editing and synthesis tools, combined with massive computational power, makes such capabilities accessible to everyone.

Given that the multimedia editors of a news site typically follow a defined workflow when creating, acquiring, editing, and curating content for their pages, this pattern adds a crucial dimension to profiling the factuality and the bias of a news source. In fact, questions around the origin and the veracity of photographic images and videos have long been the subject of multimedia forensics research (Sencar and Memon, 2013; Sencar et al., 2022). There has been research on verifying metadata integrity (Yang et al., 2020; Kee et al.; Iuliani et al., 2018; Yang et al., 2020), digital integrity (Cozzolino and Verdoliva, 2018; Korus, 2017; Cozzolino and Verdoliva, 2018), physical integrity (Matern et al., 2020; O'Brien et al., 2012; Iuliani et al., 2017; Matern et al., 2020; Riess et al., 2017; Peng et al., 2017) identification of processing

traces (Hadwiger et al., 2019), and discrimination of synthesized (i.e., GAN generated) media (Agarwal et al., 2020; Li et al., 2018; Agarwal et al., 2020; Verdoliva, 2020). However, these capabilities have only been sparsely explored in the context of predicting factuality and bias.

Existing work mainly considered characteristics of images appearing in trustworthy vs. unreliable sources. It was proposed to use visual characteristics (Jin et al., 2016), deep-learning representations (Qi et al., 2019; Khattar et al., 2019; Qi et al., 2019; Singhal et al., 2019), image provenance information from reverse image search (Zlatkova et al., 2019), and self-consistency with respect to metadata (Huh et al., 2018). Overall, multimedia characteristics have a strong potential that is yet to be fully used for news media profiling.

5.3 Audience Homophily

The well-known homophily principle, “birds of a feather flock together,” crucially asserts that similar individuals interact with each other at a higher rate. Therefore, audience representation could be another approach to describe a news media outlet whereby an overall, descriptive characteristic of followers of the outlet is obtained. Then, by evaluating the similarity of audience-centric representations with previously categorized news media, its factuality and bias can be inferred.

Ribeiro et al. (2018) used Facebook’s targeted advertising tool to infer the ideological leaning of online media based on the political leaning of the users who interacted with these media. An et al. (2012) relied on follow relationships on Twitter to ascertain the ideological leaning of news media and users. Wong et al. (2013) studied retweet behavior to infer the ideological leanings of online media sources and of popular Twitter accounts. Barberá (2015) proposed a model based on the follower relationships to media sources and Twitter personalities to estimate their ideological leaning.

Stefanov et al. (2020) predicted the political leaning of media with respect to a topic by observing the users of which side of the debate on a polarizing topic were sharing content from which media in support of their position. They constructed a user-media graph and then used label propagation and graph neural networks to derive representations for media, which they used for classification. They further aggregated the leanings across several polarizing topics to come up with a left-center-right

polarization prediction.

Following a similar approach, (Baly et al., 2020b) considered three social media platforms for audience characterization. On Twitter, they proposed to use self-descriptions in publicly accessible profiles of users following the account of a medium. For each medium, a representation is obtained by encoding the biographic descriptions of Twitter followers and averaging the resulting textual representations. The second characterization involves how the audience of the medium’s YouTube channel responds to each video in terms of number of comments, views, likes and dislikes. By averaging these statistics over all videos, a medium-level representation is obtained. The last audience representation is obtained using Facebook’s advertising platform, which is used to obtain demographic information for the audience interested in each medium. This data is used to obtain the audience distribution over the political spectrum. The distribution is then divided into five categories to label each medium accordingly: very conservative, conservative, moderate, liberal, and very liberal.

5.4 Infrastructure Characteristics

Beyond textual, visual, and audience features, news sites also exhibit distinct characteristics that relate to the underlying infrastructure and technological components deployed to serve their content online. In this regard, the prediction problem is analogous to a well-studied one in the cybersecurity domain where the goal has been to identify infrastructure characteristics of malicious domains (Anderson et al., 2007; Invernizzi et al., 2014) that are used for malware distribution (Wang et al., 2013; Invernizzi et al., 2014), phishing (Purwanto et al., 2020; James et al., 2013; Mohammad et al., 2012, 2014; Purwanto et al., 2020), online scams (Alrwais et al., 2017; Konte et al., 2009; Hao et al., 2016), and spamming (Anderson et al., 2007; Hao et al., 2009). Since establishing the infrastructure of a news medium involves several decisions with respect to technological aspects, it is plausible to expect that news media with varying IT practices and different levels of access to IT resources will differ in their characteristics.

There has been very little work on network, web design, and data elements of a news website to characterize new sites for factuality and bias. At the network level, (Hounsel et al., 2020) aimed to distinguish disinformation websites vs. authentic web-

sites vs. sites not related to news or politics, and found that features related to a website’s domain name, registration, and DNS configuration work best. Concerning the web design aspect, [Castelo et al. \(2019\)](#) introduced a web page classifier based on several features that govern the structure and the style of a page in addition to three categories of linguistic features. Their binary classification results (real vs. fake news) on several datasets showed that the web-markup features consistently perform well and are complementary to linguistic ones.

Finally, at the data level [Fairbanks et al. \(2018\)](#) examined the source of web pages to identify shared data objects, such as mutually linked sites, scripts, and images, across web sites. This information is then used to create a shared data object graph. By comparing the content-level features to the structural properties of the graph, they found that the use of mutually shared objects yields better performance in predicting both factuality and bias for a site, especially for factuality. Overall, a major advantage of using infrastructure features is their content- and audience-agnostic nature. This allows making reliable predictions when only limited textual and visual content is available and without an established audience interest in a news medium.

6 Lessons Learned

Factuality and bias have some commonalities as they exert negative influences on the public by delivering information that is deviated from the truth. Not surprisingly, some news media purposely take a biased position in the political landscape and appeal to partisan audiences. This trend becomes apparent in recent years mainly because the news industry becomes more and more competitive. Many journalists and editors, however, have concerns about their biases in news selection and reporting and try to be neutral or at least report diverse perspectives of an issue.

As the bias can be conveyed by different means—text, photos, and videos—, media bias can get subtle in many dimensions. Among them, ideological bias is an important conceptualization due to the importance of media bias in the political context. In the US context, the ideological bias could be broadly defined as conservative, center, and liberal. Then, the (ideological) bias prediction task is formulated as predicting whether a given news story, including both text and visual elements, favors one party over the other. Reported results so far show

that accurate prediction of this ideological bias of a news medium is a far more easier task than assessing factuality. This is, in fact, not surprising as evaluation of the factuality ultimately depends on the authenticity and the objectivity of the particular claims stated in a news story, essentially requiring verification from other sources and observations.

Although more sophisticated analysis of the text style and multimedia characteristics may be expected to improve the achievable accuracy, it is evident that there is a big need to complement the textual and visual elements of a news medium with others. In this regard, recent studies have demonstrated the potential of audience homophily and the medium’s infrastructure characteristics in bridging the existing performance gap. The content-agnostic nature of these characteristics make them useful in the early discovery and categorization of news media even in the absence of sufficient content.

7 Major Challenges

Ordinal scales: While the ideological bias (news slant) is typically modeled as left-center-right, there exists a spectrum within each bias based on bias intensity. A hyperpartisan bias prediction task has been tested to differentiate far-right from right and far-left from left, but it does not model the political bias using an ordinal scale. Difficulties in labeling the bias (i.e., creating ground-truth datasets) by experts or crowdsourcing is a major hurdle for modeling ideological bias as an ordinal variable.

Multimodality: In news reporting, a photo typically gets high attention, and readers can sometimes understand a news story from news photos only, even without reading the text. Indeed, news text and photos are strongly coupled together and deliver relevant information about news stories to readers. Thus, there should be a benefit from modeling news text and photos together to understand their bias and factuality ([Alam et al., 2022](#)), and potential harmfulness ([Sharma et al., 2022](#)).

Evaluation granularity: The label of a news medium is inferred from a sample of observations. This can introduce a measurement bias when a news medium does not exhibit the same reporting behavior with all news items it publishes. This is especially the case for media that have a particular stance in only certain issues ([Ganguly et al., 2020](#)). Thus, reliable estimation of factuality and bias labels require analyzing a relatively large amount of

684	content covering a range of issues.	
685	Variability in factuality & bias ratings: These	
686	ratings are inherently not static and may change	
687	over time when a news medium takes corrective ac-	
688	tion to address issues raised by fact-checkers. Thus,	
689	the ground truth needed for building a learning ap-	
690	proach varies, triggering the need for re-evaluating	
691	the performance of the proposed approaches. Thus,	
692	there is a need to take into account the sensitivity of	
693	a learning approach to such small but nevertheless	
694	inevitable variations.	
695	Dataset size: The datasets for media-level fac-	
696	tuality and bias are relatively small, typically of	
697	a few hundred examples. They are derived from	
698	sites, such as Media Bias/Fact Check and AllSides,	
699	where domain experts perform manual analysis.	
700	Annotation vs. modeling: One problem is that	
701	human annotators judge the factuality of reporting	
702	and the bias of media based on criteria that are	
703	not easy to automate or based on information that	
704	may not be accessible to automatic systems. For	
705	example, if a news outlet is judged to be of mixed	
706	factuality based on it having failed just 2-3 fact-	
707	checks, for an automatic system to arrive at the	
708	same conclusion using the same idea, it would have	
709	to select for analysis the exact same articles where	
710	the false claims were made.	
711	Data availability: Primarily due to copyright	
712	issues, there are only a few publicly available	
713	datasets of the full text of news for research pur-	
714	poses. Instead, indexed data (e.g., GDELT dataset ¹)	
715	by mentioned actors, events, locations, sources, or	
716	tones are available and have been analyzed in many	
717	studies. A set of news headlines collected from	
718	news websites or aggregated websites (e.g., All-	
719	Sides) are also shared more actively for research	
720	purposes. Considering the importance of social	
721	media channels in news dissemination, researchers	
722	collect and analyze social media posts of official	
723	accounts of news media. As social media posts are	
724	relatively more informal than news articles to fit	
725	for social media audience (Park et al., 2021), more	
726	studies are required for understanding their biases	
727	and factuality correctly.	
728	8 Future Forecasting	
729	Support for non-English corpora and different	
730	political systems: Most of the studies we review	
	are conducted for English. More research on bias	731
	and factuality for other languages thus is expected.	732
	Recently, various approaches are proposed to accel-	733
	erate NLP research for resource-scarce languages,	734
	such as multilingual word embeddings. We believe	735
	that those efforts help conduct bias and factuality re-	736
	search for non-English corpora. One non-technical	737
	issue here is that not all the countries have US-like	738
	left-center-right political biases. For example, there	739
	might exist a multiparty system in some countries.	740
	In that case, understanding relevant political biases	741
	should be the first step in media bias research.	742
	Incorporation of video content: TV news ac-	743
	counts for significant portions of the news industry.	744
	Also, the presence of news media becomes strong	745
	in video-driven social media platforms over time.	746
	To get high user engagements, news media outlets	747
	upload short video clips curated for social media	748
	use, particularly on existing social media. Most pre-	749
	vious studies on bias in video news have analyzed	750
	their transcripts instead of analyzing video directly.	751
	Commercial databases, such as Lexis Nexis, or	752
	open-source libraries to create subtitles are used	753
	to analyze news transcripts. We expect that more	754
	studies on analyzing video contents in an end-to-	755
	end manner will be presented to fully understand	756
	the bias and factuality of video news.	757
	Bringing practical implications: Since the fac-	758
	tuality and the bias of news media largely influence	759
	the public, it is crucial to implement working sys-	760
	tems, so that readers can benefit from a rich stream	761
	of research. Several stand-alone websites, such	762
	as Media Bias/Fact Check, AllSides, and Tanbih	763
	(Zhang et al., 2019), aim to make media bias and	764
	factuality transparent to end-users, thus promoting	765
	media literacy. We expect new tools and services	766
	to support more media and languages.	767
	9 Conclusion	768
	We reviewed the state of the art on media profil-	769
	ing for factuality and bias, arguing for the need to	770
	model them jointly. We further discussed interest-	771
	ing recent advances in exploiting different infor-	772
	mation sources and different modalities, which go	773
	beyond the text of the articles the target news outlet	774
	has published. Finally, we discussed current chal-	775
	lenges and outlined promising research directions.	776

¹<https://www.gdeltproject.org/>

777
778
779
780
781
782

783
784
785
786
787
788
789
790

791
792
793
794
795
796

797
798
799
800

801
802
803
804

805
806
807
808

809
810
811
812
813
814
815

816
817
818
819
820
821

822
823
824
825
826
827
828
829

830
831

References

Shruti Agarwal, Hany Farid, Ohad Fried, and Maneesh Agrawala. 2020. Detecting deep-fake videos from phoneme-viseme mismatches. In *CVPR*, pages 2814–2822. IEEE.

Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fabrizio Silvestri, Dimiter Dimitrov, Giovanni Da San Martino, Shaden Shaar, Hamed Firooz, and Preslav Nakov. 2022. A survey on multimodal disinformation detection. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6625–6643, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Sumayah Alrwais, Xiaojing Liao, Xianghang Mi, Peng Wang, Xiaofeng Wang, Feng Qian, Raheem Beyah, and Damon McCoy. 2017. Under the shadow of sunshine: Understanding and detecting bulletproof hosting on legitimate service provider networks. In *IEEE SP*, pages 805–823.

Jisun An, Meeyoung Cha, Krishna Gummadi, Jon Crowcroft, and Daniele Quercia. 2012. Visualizing media bias through twitter. In *AAAI ICWSM*, volume 6.

Jisun An, Meeyoung Cha, P. Krishna Gummadi, and Jon Crowcroft. 2011. Media landscape in twitter: A world of new conventions and political diversity. In *AAAI ICWSM*.

David S. Anderson, Chris Fleizach, Stefan Savage, and Geoffrey M. Voelker. 2007. *Spamscatter: Characterizing internet scam hosting infrastructure*. Ph.D. thesis.

Ramy Baly, Giovanni Da San Martino, James Glass, and Preslav Nakov. 2020a. We can detect your bias: Predicting the political ideology of news articles. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4982–4991, Online. Association for Computational Linguistics.

Ramy Baly, Georgi Karadzhov, Dimitar Alexandrov, James Glass, and Preslav Nakov. 2018a. Predicting factuality of reporting and bias of news media sources. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3528–3539, Brussels, Belgium.

Ramy Baly, Georgi Karadzhov, Jisun An, Haewoon Kwak, Yoan Dinkov, Ahmed Ali, James Glass, and Preslav Nakov. 2020b. What was written vs. who read it: News media profiling using text analysis and social media context. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3364–3374, Online. Association for Computational Linguistics.

Ramy Baly, Georgi Karadzhov, Abdelrhman Saleh, James Glass, and Preslav Nakov. 2019. Multi-task

ordinal regression for jointly predicting the trustworthiness and the leading political ideology of news media. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2109–2116, Minneapolis, Minnesota. 832
833
834
835
836
837
838

Ramy Baly, Mitra Mohtarami, James Glass, Lluís Màrquez, Alessandro Moschitti, and Preslav Nakov. 2018b. Integrating stance detection and fact checking in a unified corpus. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 21–27, New Orleans, Louisiana. 839
840
841
842
843
844
845
846

Pablo Barberá. 2015. Birds of the same feather tweet together: Bayesian ideal point estimation using Twitter data. *Political Analysis*, 23(1):76–91. 847
848
849

Andrew W Barrett and Lowell W Barrington. 2005. Bias in newspaper photograph selection. *Political Research Quarterly*, 58(4):609–618. 850
851
852

Matthew A Baum and Tim Groeling. 2008. New media and the polarization of american political discourse. *Political Communication*, 25(4):345–365. 853
854
855

Eric Baumer, Elisha Elovic, Ying Qin, Francesca Polletta, and Geri Gay. 2015. Testing and comparing computational approaches for identifying the language of framing in political news. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1472–1482, Denver, Colorado. Association for Computational Linguistics. 856
857
858
859
860
861
862
863
864

Vibhu Bhatia, Vidya Prasad Akavoor, Sejin Paik, Lei Guo, Mona Jalal, Alyssa Smith, David Assefa Tofu, Edward Edberg Halim, Yimeng Sun, Margrit Betke, Prakash Ishwar, and Derry Tanti Wijaya. 2021. OpenFraming: Open-sourced tool for computational framing analysis of multilingual data. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 242–250, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 865
866
867
868
869
870
871
872
873
874

Emina Boudemagh and Izabela Moise. 2017. News media coverage of refugees in 2016: a GDELT case study. In *AAAI ICWSM*. 875
876
877

Maxwell T Boykoff and Jules M Boykoff. 2004. Balance as bias: Global warming and the us prestige press. *Global environmental change*, 14(2):125–136. 878
879
880

Ann M Brill. 2001. Online journalists embrace new marketing function. *Newspaper Research Journal*, 22(2):28–40. 881
882
883

Kevin R. Canini, Bongwon Suh, and Peter L. Pirolli. 2011. Finding credible information sources in social networks based on content and social structure. In *Proceedings of the IEEE International Conference* 884
885
886
887

888	<i>on Privacy, Security, Risk, and Trust, and the IEEE International Conference on Social Computing, SocialCom/PASSAT '11</i> , pages 1–8.	<i>Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota.	943 944 945
891	Dallas Card, Amber E. Boydston, Justin H. Gross, Philip Resnik, and Noah A. Smith. 2015. The media frames corpus: Annotations of frames across issues . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 438–444, Beijing, China. Association for Computational Linguistics.	Yoan Dinkov, Ahmed Ali, Ivan Koychev, and Preslav Nakov. 2019. Predicting the leading political ideology of Youtube channels using acoustic, textual and metadata information. In <i>Proceedings of the 20th Annual Conference of the International Speech Communication Association, INTERSPEECH '19</i> .	946 947 948 949 950 951
892			
893			
894			
895			
896			
897			
898			
899			
900	Sonia Castelo, Thais Almeida, Anas Elghafari, Aécio Santos, Kien Pham, Eduardo Nakamura, and Juliana Freire. 2019. A topic-agnostic approach for identifying fake news pages. In <i>WWW Companion</i> , pages 975–980.	Xin Luna Dong, Evgeniy Gabrilovich, Kevin Murphy, Van Dang, Wilko Horn, Camillo Lugaresi, Shaohua Sun, and Wei Zhang. 2015. Knowledge-based trust: Estimating the trustworthiness of web sources . <i>Proc. VLDB Endow.</i> , 8(9):938–949.	952 953 954 955 956
901			
902			
903			
904			
905	Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information credibility on twitter . In <i>Proceedings of the 20th International Conference on World Wide Web, WWW 2011, Hyderabad, India, March 28 - April 1, 2011</i> , pages 675–684.	Sebastian Dungs, Ahmet Aker, Norbert Fuhr, and Kalina Bontcheva. 2018. Can rumour stance alone predict veracity? In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 3360–3370, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	957 958 959 960 961 962
906			
907			
908			
909			
910	Cheng Chen, Kui Wu, Srinivasan Venkatesh, and Xudong Zhang. 2013. Battling the internet water army: detection of hidden paid posters . In <i>Advances in Social Networks Analysis and Mining 2013, ASONAM '13, Niagara, ON, Canada - August 25 - 29, 2013</i> , pages 116–120.	Robert M. Entman. 1993. Framing: Toward clarification of a fractured paradigm . <i>Journal of Communication</i> , 43(4):51–58.	963 964 965
911			
912			
913			
914			
915			
916	Davide Cozzolino and Luisa Verdoliva. 2018. Camera-based image forgery localization using convolutional neural networks. In <i>EUSIPCO</i> , pages 1372–1376.	Robert M Entman. 2007. Framing bias: Media in the distribution of power. <i>Journal of communication</i> , 57(1):163–173.	966 967 968
917			
918			
919	Fernando Cardoso Durier da Silva, Rafael Vieira, and Ana Cristina Bicharra Garcia. 2019. Can machines learn to detect fake news? a survey focused on social media. In <i>HICSS</i> .	James Fairbanks, Natalie Fitch, Nathan Knauf, and Erica Briscoe. 2018. Credibility assessment in the news: do we need to read. In <i>MIS2 Workshop</i> , pages 799–800.	969 970 971 972
920			
921			
922			
923	Dave D’Alessio and Mike Allen. 2000. Media bias in presidential elections: A meta-analysis. <i>Journal of communication</i> , 50(4):133–156.	Anjalie Field, Doron Kliger, Shuly Wintner, Jennifer Pan, Dan Jurafsky, and Yulia Tsvetkov. 2018. Framing and agenda-setting in Russian news: a computational analysis of intricate political strategies . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3570–3580, Brussels, Belgium. Association for Computational Linguistics.	973 974 975 976 977 978 979 980
924			
925			
926	Kareem Darwish, Dimitar Alexandrov, Preslav Nakov, and Yelena Mejova. Seminar users in the Arabic Twitter sphere. In <i>SocInfo</i> .	H. Finberg et al. 2002. Digital journalism credibility study. <i>Online News Association</i> . Retrieved November, 3:2003.	981 982 983
927			
928			
929	Kareem Darwish, Peter Stefanov, Michaël J. Aupetit, and Preslav Nakov. 2020. Unsupervised user stance detection on twitter. In <i>AAAI ICWSM</i> , pages 141–152.	Soumen Ganguly, Juhi Kulshrestha, Jisun An, and Hae-woon Kwak. 2020. Empirical evaluation of three common assumptions in building political media bias datasets. In <i>AAAI ICWSM</i> .	984 985 986 987
930			
931			
932			
933	Sohan De Sarkar, Fan Yang, and Arjun Mukherjee. 2018. Attending sentences to detect satirical fake news . In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 3371–3380, Santa Fe, New Mexico, USA.	Matthew Gentzkow and Jesse M. Shapiro. 2010. What drives media slant? evidence from u.s. daily newspapers. <i>Econometrica</i> , 78(1):35–71.	988 989 990
934			
935			
936			
937			
938	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for</i>	Martin Gilens and Craig Hertzman. 2000. Corporate ownership and news bias: Newspaper coverage of the 1996 telecommunications act. <i>The Journal of Politics</i> , 62(2):369–386.	991 992 993 994
939			
940			
941			
942			

995	Jesse Graham, Jonathan Haidt, and Brian A Nosek.	Joan B Hooper. 1975. <i>On assertive predicates</i> , volume 4.	1041
996	2009. Liberals and conservatives rely on different	Academic Press, New York.	1042
997	sets of moral foundations. <i>Journal of personality and</i>		
998	<i>social psychology</i> , 96(5):1029.	Benjamin Horne and Sibel Adali. 2017. This just in: Fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news.	1043
999	Tim Groeling. 2013. Media bias by the numbers: Chal-		1044
1000	lenges and opportunities in the empirical study of		1045
1001	partisan news. <i>Annual Review of Political Science</i> ,	<i>CoRR</i> , abs/1703.09398.	1046
1002	16.	Benjamin D. Horne, Sibel Adali, and Sujoy Sikdar. 2017. Identifying the social signals that drive online discussions: A case study of Reddit communities. In <i>IEEE ICCCN</i> , pages 1–9.	1047
1003	Tim Groseclose and Jeffrey Milyo. 2005. A measure		1048
1004	of media bias. <i>The Quarterly Journal of Economics</i> ,		1049
1005	120(4):1191–1237.		1050
1006	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vla-	Benjamin D. Horne, Sara Khedr, and Sibel Adali. 2018. Sampling the news producers: A large news and feature data set for the study of the complex media landscape. In <i>Proceedings of the Twelfth International Conference on Web and Social Media</i> , ICWSM '18, pages 518–527.	1051
1007	chos. 2022a. A survey on automated fact-checking.		1052
1008	<i>Transactions of the Association for Computational</i>		1053
1009	<i>Linguistics</i> , 10:178–206.		1054
1010	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vla-	Austin Hounsel, Jordan Holland, Ben Kaiser, Kevin Bor-	1055
1011	chos. 2022b. A survey on automated fact-checking.	golte, Nick Feamster, and Jonathan Mayer. 2020. Identifying disinformation websites using infrastructure features. In <i>Proceedings of the 10th USENIX Workshop on Free and Open Communications on the Internet</i> , FOCI '20.	1056
1012	<i>Transactions of the Association for Computational</i>		1057
1013	<i>Linguistics</i> , 10:178–206.		1058
1014	Benjamin Hadwiger, Daniele Baracchi, Alessandro Piva,	Minyoung Huh, Andrew Liu, Andrew Owens, and	1059
1015	and Christian Riess. 2019. Towards learned color representations for image splicing detection. In <i>IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019</i> , pages 8281–8285. IEEE.	Alexei A. Efros. 2018. Fighting fake news: Image splice detection via learned self-consistency. In <i>ECCV</i> .	1060
1016			1061
1017		Clayton J. Hutto and Eric Gilbert. 2014. VADER: A parsimonious rule-based model for sentiment analysis of social media text. In <i>Proceedings of the 8th International Conference on Weblogs and Social Media</i> , ICWSM '14.	1062
1018			1063
1019			1064
1020	S. Hao, A. Kantchelian, B. Miller, V. Paxson, and	Ken Hyland. 2005. <i>Metadiscourse: Exploring Interaction in Writing</i> . Bloomsbury Publishing.	1065
1021	N. Feamster. 2016. Predator: proactive recognition		1066
1022	and elimination of domain abuse at time-of-		1067
1023	registration. In <i>ACM CCS</i> , pages 1568–1579.		1068
1024	Shuang Hao, Nadeem Ahmed Syed, Nick Feamster,	Luca Invernizzi, Stanislav Miskovic, Ruben Torres,	1069
1025	Alexander G Gray, and Sven Krasser. 2009. Detect-	Christopher Kruegel, Sabyasachi Saha, Giovanni Vigna, Sung-Ju Lee, and Marco Mellia. 2014. Nazca: Detecting malware distribution in large-scale networks. In <i>NDSS</i> , volume 14, pages 23–26.	1070
1026	ing spammers with snare: Spatio-temporal network-		1071
1027	level automatic reputation engine. In <i>USENIX Security</i> ,		1072
1028	volume 9.		1073
1029	Momchil Hardalov, Arnav Arora, Preslav Nakov, and	Md. Rafiqul Islam, Shaowu Liu, Xianzhi Wang, and Guan-	1074
1030	Isabelle Augenstein. 2022a. A survey on stance detection for mis- and disinformation identification. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 1259–1277, Seattle, United States. Association for Computational Linguistics.	dong Xu. 2020. Deep learning for misinformation detection on online social networks: a survey and new perspectives. <i>SNAM</i> , 10(1):1–20.	1075
1031			1076
1032		Massimo Iuliani, Marco Fanfani, Carlo Colombo, and	1077
1033		Alessandro Piva. 2017. Reliability assessment of principal point estimates for forensic applications. <i>J. Vis. Comm. Img. Repr.</i> , 42:65–77.	1078
1034			1079
1035		Massimo Iuliani, Dasara Shullani, Marco Fontani, Saverio Meucci, and Alessandro Piva. 2018. A video forensic framework for the unsupervised analysis of mp4-like file container. <i>IEEE TIFS</i> , 14(3):635–645.	1080
1036	Momchil Hardalov, Anton Chernyavskiy, Ivan Koychev,		1081
1037	Dmitry Ilvovsky, and Preslav Nakov. 2022b. Crowd-	J. James, L. Sandhya, and C. Thomas. 2013. Detection of phishing urls using machine learning techniques. In <i>IEEE ICC</i> , pages 304–309.	1082
1038	Checked: Detecting previously fact-checked claims		1083
1039	in social media. In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing, AACL_IJCNLP '22, online.</i>		1084
1040		Zhiwei Jin, Juan Cao, Yongdong Zhang, Jianshe Zhou, and Qi Tian. 2016. Novel visual and statistical image features for microblogs news verification. <i>IEEE TMM</i> , 19(3):598–608.	1085
			1086
			1087
			1088
			1089
			1090
			1091
			1092
			1093
			1094

1201	Jing Ma, Wei Gao, Zhongyu Wei, Yueming Lu, and Kam-Fai Wong. 2015. Detect rumors using time series of social context information on microblogging websites . In <i>Proceedings of the 24th ACM International Conference on Information and Knowledge Management, CIKM 2015, Melbourne, VIC, Australia, October 19 - 23, 2015</i> , pages 1751–1754.	Subhabrata Mukherjee and Gerhard Weikum. 2015. Leveraging joint interactions for credibility analysis in news communities . In <i>Proceedings of the 24th ACM International Conference on Information and Knowledge Management, CIKM 2015, Melbourne, VIC, Australia, October 19 - 23, 2015</i> , pages 353–362.	1254 1255 1256 1257 1258 1259
1208	Jing Ma, Wei Gao, and Kam-Fai Wong. 2017. Detect rumors in microblog posts using propagation structure via kernel learning . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 708–717, Vancouver, Canada.	Preslav Nakov, David Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021. Automated fact-checking for assisting human fact-checkers. In <i>Proceedings of the 30th International Joint Conference on Artificial Intelligence, IJCAI '21</i> , pages 4551–4558.	1260 1261 1262 1263 1264 1265 1266
1214	Suman Kalyan Maity, Aishik Chakraborty, Pawan Goyal, and Animesh Mukherjee. 2017. Detection of sockpuppets in social media. In <i>CSCW, CSCW '17</i> .	An T. Nguyen, Aditya Kharosekar, Matthew Lease, and Byron C. Wallace. 2018. An interpretable joint graphical model for fact-checking from crowds . In <i>Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, AAAI '18</i> , pages 1511–1518.	1267 1268 1269 1270 1271
1217	Giovanni Da San Martino, Stefano Cresci, Alberto Barrón-Cedeño, Seunghak Yu, Roberto Di Pietro, and Preslav Nakov. 2020. A survey on computational propaganda detection . In <i>Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020</i> , pages 4826–4832. ijcai.org.	Van-Hoang Nguyen, Kazunari Sugiyama, Preslav Nakov, and Min-Yen Kan. 2022. Fang: Leveraging social context for fake news detection using graph representation . <i>Commun. ACM</i> , 65(4):124–132.	1272 1273 1274 1275
1223	Falko Matern, Christian Riess, and Marc Stamminger. 2020. Gradient-based illumination description for image forgery detection. <i>IEEE TIFS</i> , 15:1303–1317.	O'Brien et al. 2012. Exposing photo manipulation with inconsistent reflections. <i>ACM Trans. Graph.</i> , 31(1):4:1–4:11.	1276 1277 1278
1226	Todor Mihaylov and Preslav Nakov. 2016. Hunting for troll comments in news community forums . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 399–405, Berlin, Germany.	Jeremy Padgett, Johanna L Dunaway, and Joshua P Darr. 2019. As seen on tv? how gatekeeping makes the us house seem more extreme. <i>Journal of Communication</i> , 69(6):696–719.	1279 1280 1281 1282
1231	Tsvetomila Mihaylova, Preslav Nakov, Lluís Màrquez, Alberto Barrón-Cedeño, Mitra Mohtarami, Georgi Karadzhov, and James R. Glass. 2018. Fact checking in community forums . In <i>Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18)</i> , pages 5309–5316, New Orleans, Louisiana, USA.	Jeff Z. Pan, Siyana Pavlova, Chenxi Li, Ningxi Li, Yangmei Li, and Jinshuo Liu. 2018. Content based fake news detection using knowledge graphs. In <i>Proceedings of the International Semantic Web Conference, ISWC '18</i> .	1283 1284 1285 1286
1238	Lewis Mitchell, Kameron Decker Harris, Morgan R. Frank, Peter Sheridan Dodds, and Christopher M. Danforth. 2013. The geography of happiness: Connecting Twitter sentiment and expression, demographics, and objective characteristics of place. <i>PLoS one</i> , 8(5):e64417.	Panayot Panayotov, Utsav Shukla, Husrev Taha Sen-car, Mohamed Nabeel, and Preslav Nakov. 2022. GREENER: Graph neural networks for news media profiling. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP '22, Abu Dhabi, UAE</i> .	1287 1288 1289 1290 1291 1292
1243	R. M. Mohammad, F. Thabtah, and L. McCluskey. 2012. An assessment of features related to phishing websites using an automated technique. In <i>IEEE ICITST</i> , pages 492–497.	Kunwoo Park, Haewoon Kwak, Jisun An, and Sanjay Chawla. 2021. Understanding effects of editing tweets for news sharing by media accounts through a causal inference framework. In <i>AAAI ICWSM</i> .	1293 1294 1295 1296
1247	R. M. Mohammad, F. Thabtah, and L. McCluskey. 2014. Predicting phishing websites based on self-structuring neural network. <i>Neural Computing and Applications</i> , 25(2).	Bo Peng, Wei Wang, Jing Dong, and Tieniu Tan. 2017. Optimized 3d lighting environment estimation for image forgery detection. <i>IEEE TIFS</i> , 12(2):479–494.	1297 1298 1299
1251	Morstatter et al. 2018. Identifying framing bias in online news. <i>ACM Transactions on Social Computing</i> , 1(2):1–18.	Verónica Pérez-Rosas, Bennett Kleinberg, Alexandra Lefevre, and Rada Mihalcea. 2018. Automatic detection of fake news . In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 3391–3401, Santa Fe, New Mexico, USA.	1300 1301 1302 1303 1304
1252		Kashyap Papat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2016. Credibility assessment of textual claims on the web. In <i>CIKM</i> .	1305 1306 1307

1308	Kashyap Papat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2017. Where the truth lies: Explaining the credibility of emerging claims on the Web and social media . In <i>Proceedings of the 26th International Conference on World Wide Web Companion, WWW '17 Companion</i> , page 1003–1012, Perth, Australia.	Sencar et al., editor. 2022. <i>Multimedia Forensics</i> . Springer.	1362 1363
1310		Shivam Sharma, Firoj Alam, Md. Shad Akhtar, Dimitar Dimitrov, Giovanni Da San Martino, Hamed Firooz, Alon Halevy, Fabrizio Silvestri, Preslav Nakov, and Tanmoy Chakraborty. 2022. Detecting and understanding harmful memes: A survey. In <i>Proceedings of the 31st International Joint Conference on Artificial Intelligence, IJCAI-ECAI '22</i> , Vienna, Austria.	1364 1365 1366 1367 1368 1369 1370
1311		Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. <i>SIGKDD Explor. Newsl.</i> , 19(1).	1371 1372 1373
1312		Shivangi Singhal, Rajiv Ratn Shah, Tanmoy Chakraborty, Ponnurangam Kumaraguru, and Shin'ichi Satoh. 2019. Spottfake: A multi-modal framework for fake news detection. In <i>IEEE BigMM</i> , pages 39–47.	1374 1375 1376 1377
1313		Stuart N Soroka. 2012. The gatekeeping function: Distributions of information in media and the real world. <i>The Journal of Politics</i> , 74(2):514–528.	1378 1379 1380
1314		Peter Stefanov, Kareem Darwish, Atanas Atanasov, and Preslav Nakov. 2020. Predicting the topical stance and political leaning of media using tweets . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 527–537, Online. Association for Computational Linguistics.	1381 1382 1383 1384 1385 1386
1315	Kashyap Papat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2018. CredEye: A credibility lens for analyzing and explaining misinformation . In <i>Companion Proceedings of the The Web Conference 2018, WWW '18</i> , page 155–158, Lyon, France.	Robert L. Stevenson, Richard A. Eisinger, Barry M. Feinberg, and Alan B. Kotok. 1973. Untwisting the news twisters: A replication of efron's study . <i>Journalism Quarterly</i> , 50(2):211–219.	1387 1388 1389 1390
1316		Edson C Tandoc Jr. 2014. Journalism is twerking? how web analytics is changing the process of gatekeeping. <i>New media & society</i> , 16(4):559–575.	1391 1392 1393
1317		James Thorne and Andreas Vlachos. 2018a. Automated fact checking: Task formulations, methods and future directions . In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 3346–3359, Santa Fe, New Mexico, USA.	1394 1395 1396 1397 1398
1318		James Thorne and Andreas Vlachos. 2018b. Automated fact checking: Task formulations, methods and future directions. <i>arXiv preprint arXiv:1806.07687</i> .	1399 1400 1401
1319		Luisa Verdoliva. 2020. Media forensics and deepfakes: an overview. <i>IEEE J. of STSP</i> , 14(5):910–932.	1402 1403
1320	Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek Bevendorff, and Benno Stein. 2018. A stylometric inquiry into hyperpartisan and fake news . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 231–240, Melbourne, Australia.	S. Vosoughi et al. 2018. The spread of true and false news online. <i>Science</i> , 359(6380).	1404 1405
1321		Paul Waldman and James Devitt. 1998. Newspaper photographs and the 1996 presidential election: The question of bias. <i>Journal Mass Commun Q</i> , 75(2):302–311.	1406 1407 1408
1322		Robert Walecki, Ognjen Rudovic, Vladimir Pavlovic, and Maja Pantic. 2016. Copula ordinal regression for joint estimation of facial action unit intensity. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 4902–4910.	1409 1410 1411 1412 1413
1323			
1324			
1325			
1326	Rizka Purwanto, Arindam Pal, Alan Blair, and Sanjay Jha. 2020. Phishzip: A new compression-based algorithm for detecting phishing websites. In <i>IEEE CNS</i> , pages 1–9.		
1327			
1328			
1329			
1330	Peng Qi, Juan Cao, Tianyun Yang, Junbo Guo, and Jintao Li. 2019. Exploiting multi-domain visual information for fake news detection. In <i>IEEE ICDM</i> , pages 518–527.		
1331			
1332			
1333			
1334	Marta Recasens, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky. 2013. Linguistic models for analyzing and detecting biased language . In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1650–1659, Sofia, Bulgaria.		
1335			
1336			
1337			
1338			
1339			
1340	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.		
1341			
1342			
1343			
1344			
1345			
1346			
1347			
1348	Ribeiro et al. 2018. Media bias monitor: Quantifying biases of social media news outlets at large-scale. In <i>AAAI ICWSM</i> , pages 290–299.		
1349			
1350			
1351	Christian Riess, Mathias Unberath, Farzad Naderi, Sven Pfaller, Marc Stamminger, and Elli Angelopoulou. 2017. Handling multiple materials for exposure of digital forgeries using 2-d lighting environments. <i>MM. Tools Appl.</i> , 76(4):4747–4764.		
1352			
1353			
1354			
1355			
1356	Ellen Riloff and Janyce Wiebe. 2003. Learning extraction patterns for subjective expressions . In <i>Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing</i> , pages 105–112.		
1357			
1358			
1359			
1360	Husrev Taha Sencar and Nasir Memon, editors. 2013. <i>Digital Image Forensics</i> . Springer.		
1361			

1414 G. Wang, J. W. Stokes, C. Herley, and D. Felstead. 2013.
1415 Detecting malicious landing pages in malware distribu-
1416 tion networks. In *IEEE/IFIP DSN*, pages 1–11.

1417 David A Weaver and Bruce Bimber. 2008. Finding news
1418 stories: a comparison of searches using LexisNexis and
1419 Google News. *Journal Mass Commun Q*, 85(3):515–
1420 530.

1421 Felix Ming Fai Wong, Chee Wei Tan, Soumya Sen, and
1422 Mung Chiang. 2013. Quantifying political leaning from
1423 tweets and retweets. In *AAAI ICWSM*, pages 640–649.

1424 Pengpeng Yang, Daniele Baracchi, Massimo Iuliani,
1425 Dasara Shullani, Rongrong Ni, Yao Zhao, and Alessan-
1426 dro Piva. 2020. Efficient video integrity analysis
1427 through container characterization. *IEEE J. of STSP*,
1428 14(5):947–954.

1429 Zaman et al. 2014. A bayesian approach for predicting the
1430 popularity of tweets. *Ann. Appl. Stat.*, 8(3).

1431 Yifan Zhang, Giovanni Da San Martino, Alberto Barrón-
1432 Cedeño, Salvatore Romeo, Jisun An, Haewoon Kwak,
1433 Todor Staykovski, Israa Jaradat, Georgi Karadzhov,
1434 Ramy Baly, Kareem Darwish, James Glass, and Preslav
1435 Nakov. 2019. [Tanbih: Get to know what you are reading](#).
1436 In *Proceedings of the 2019 Conference on Empirical
1437 Methods in Natural Language Processing and the 9th
1438 International Joint Conference on Natural Language
1439 Processing (EMNLP-IJCNLP): System Demonstrations*,
1440 pages 223–228, Hong Kong, China.

1441 Xinyi Zhou and Reza Zafarani. 2020. A survey of fake
1442 news: Fundamental theories, detection methods, and
1443 opportunities. *CSUR*.

1444 Dimitrina Zlatkova, Preslav Nakov, and Ivan Koychev.
1445 2019. [Fact-checking meets fauxtography: Verifying
1446 claims about images](#). In *Proceedings of the 2019 Con-
1447 ference on Empirical Methods in Natural Language
1448 Processing and the 9th International Joint Conference
1449 on Natural Language Processing (EMNLP-IJCNLP)*,
1450 pages 2099–2108, Hong Kong, China. Association for
1451 Computational Linguistics.

1452 Arkaitz Zubiaga, Geraldine Wong Sak Hoi, Maria Liakata,
1453 Rob Procter, and Peter Tolmie. 2016. Analysing how
1454 people orient to and spread rumours in social media by
1455 looking at conversational threads. *PLOS ONE*, 11(3):1–
1456 29.