

HeadlineStanceChecker: Exploiting Summarization to Detect Headline Disinformation

Abstract

The headline of a news article is designed to succinctly summarize its content, providing the reader with a clear understanding of the news item. Unfortunately, in the post-truth era, headlines are more focused on attracting the reader’s attention for ideological or commercial reasons, thus leading to mis- or disinformation through false or distorted headlines. One way of combating this, although a challenging task, is by determining the relation between the headline and the body text to establish the stance. Hence, to contribute to the detection of mis- and disinformation, this paper proposes an approach—*HeadlineStanceChecker*—that determines the stance of a headline with respect to the body text to which it is associated. The novelty rests on the use of a two-stage classification architecture that uses summarization techniques to shape the input for both classifiers instead of directly passing the full news body text, thereby reducing the amount of information to be processed while keeping important information. Specifically, summarization is done through Positional Language Models leveraging on semantic resources to identify salient information in the body text that is then compared to its corresponding headline. The results obtained show that our approach achieves 94.31% accuracy for the overall classification and the best FNC-1 relative score compared with the state of the art. It is especially remarkable that the system, which uses only the relevant information provided by the automatic summaries instead of the whole text, is able to classify the different stance categories with very competitive results, especially in the *discuss* stance between the headline and the news body text. It can be concluded that using automatic extractive summaries as input of our approach together with the two-stage architecture is an appropriate solution to the problem.

Keywords: Natural Language Processing, Fake News, Misleading Headlines, Stance Detection, Applied Computing, Document Management and Text Processing, Semantic Summarization

1. Introduction

Nowadays, disinformation and misinformation are two major problems that are increasing at great velocity [1] in pace with the exponential growth of information on the web and the need for robust verification methods. If handling this information overload is an arduous and complex task for both humans and machines, verifying its veracity has become a daunting yet unavoidable challenge. Both terms, misinformation and disinformation, allude to the inaccuracy and lack of veracity of certain information; however, while in the first case the delusion can be caused unintentionally, the latter actually seeks to deceive or misdirect deliberately [2]. In either case, they represent a type of phenomenon that, in the domain of digital news, can easily result in a massive confusion about the real facts, spreading on a viral scale. This is actually what the New York Times meant when they referred to a “Fake news” piece as a “made up story with the intention to deceive, often with monetary gain

as a motive” [3].

The ideological and economic interests that potentially gain from this “information disorder” are the drivers of fake news. These interests aim to manipulate social opinion and reinforce preconceived opinions, thereby making people focus on thinking or acting in a specific way by, most of the time, appealing to their emotions rather than presenting the facts. This trend that has even prompted the advent and consolidation of a new term, “post-truth”, which, according to the Cambridge Dictionary¹, refers to “a situation in which people are more likely to accept an argument based on their emotions and beliefs, rather than one based on facts”. For instance, this distorting phenomenon played an important role in *President Trump’s election campaign 2016* [4] and *the Brexit referendum 2016* [5]. In the same way, business and commercial interests fabri-

¹<https://dictionary.cambridge.org/>

37 cate fake news to generate income through clickbait and
38 misleading information. For instance, the National Re-
39 port website, Disinformedia [6] or Victory Lab [7] are
40 examples of websites that produce and/or disseminate
41 fake news.

42 Assessing the veracity of a news story is a com-
43 plex task—either for expert journalists or for Artificial
44 Intelligence—. For this reason, the research community
45 approaches the resolution of this task from several per-
46 spectives that imply different sub-tasks. In this manner,
47 it is convenient to assess the veracity of a news story
48 by splitting the task into simpler parts and dealing with
49 them individually [8].

50 Following this approach, great attention and effort
51 has been focused on the analysis and study of one of
52 the most essential elements of a news item, its headline,
53 in some cases focusing on the relationship between the
54 body of the article and the headline, and in others con-
55 sidering the constitution of the headline itself. Head-
56 lines are fundamental parts of news stories, they sum-
57 marise the article so that the reader clearly understands
58 the content of the news story [9]. Nevertheless, the
59 headline acts also as the prelude to the complete news
60 story, and it should be written as an invitation for the
61 reader to discover the full piece. A headline is therefore
62 expected to be as effective as possible, without losing
63 accuracy or becoming misleading [10].

64 In the scenario we have outlined, where the informa-
65 tion stream is permanently growing and filtering con-
66 tent can be overwhelming, the role of headlines is cru-
67 cial. On the one hand, an appropriate headline can help
68 us to identify the content of most interest to us, but on
69 the other hand, and due to this data deluge, it can be
70 tempting to read only the headlines and share the news
71 feed without having read the entire story. Consequently,
72 stories can often go viral because of an attractive head-
73 line despite the lack of true information in the body text.
74 This phenomena manipulates public opinion and affects
75 the credibility of social media [11, 12]. In particular, the
76 research conducted in [11] found that 59% of the URLs
77 mentioned on Twitter were not clicked at all. This sug-
78 gests that people are more willing to share an article
79 than access and read it, so they directly read and share
80 the headline (and link), without making the effort to go
81 deeper and check it. Considering this, the headline of a
82 news article should faithfully summarize the body text,
83 without including deception or misinformation, in order
84 to maintain accuracy and veracity of the entire article.

85 Unfortunately, in practice, headlines in digital me-
86 dia tend to be more focused on attracting the reader’s
87 attention—with little regard for accuracy—thus leading
88 to mis- or disinformation through erroneous/false facts

89 or headline/body dissonance [13]. In this context, head-
90 lines can be classified into two classes [14]:

- 91 • **Clickbait headlines:** Clickbait refers to content
92 whose main purpose is to attract attention and en-
93 courage visitors to click on a link to a partic-
94 ular web page with the purpose of monetizing
95 the “views” through advertising revenue (the more
96 clicks, the more money earned). This type of head-
97 line is often ambiguous and exhibits a particular
98 writing style to directly exploit human curiosity,
99 for instance by using exclamatory or interrogative
100 headlines that urge audiences to click on the link to
101 discover the missing information [14]. Typically,
102 clickbait headlines are spread on social media in
103 the form of short teaser messages that may read
104 like the following cited examples:

105 – “*Man tries to hug a wild lion, You won’t believe
106 what happens next!*”²

107 – “*The first lady of swearing! How a ten-year-old
108 Michelle Obama lost out on a ‘best camper’ award
109 because she wouldn’t stop cursing*”³

Existing methods for automatically detecting
clickbait headlines usually treat the task as a classi-
fication problem (clickbait/non-clickbait), and ex-
clusively focus on the headline (its writing style or
structure) rather than considering the content of the
news itself [15, 13].

- 110 • **Misleading headlines:** Headlines thus classified
111 significantly misrepresent the findings reported in
112 the news article [16], by exaggerating or distorting
113 the facts described in the news article. The reader
114 can only discover the inconsistencies after reading
115 the news body text [14]. Although in the literature
116 these headlines are sometimes referred to as *incon-*
117 *gruent headlines*, in this work we will refer to them
118 as *misleading headlines* since the term represents a
119 more comprehensive concept.

120 Some important nuances that are part of the news
121 body text are missing in the headline, causing the
122 reader to come to the wrong conclusion. In contrast
123 to clickbait headlines, the language used does not
124 necessarily incite the reader to click on it, but it is
125 designed to trigger emotion or excitement [16].

126 Examples of misleading headlines are shown be-
127 low (also reported in [13, 16], respectively):

128 ²<https://bit.ly/2FEddK2> (accessed online 15 February, 2021)

129 ³<https://www.dailymail.co.uk/news/article-3004975> (accessed on-
130 line 15 February, 2021)

- 134 – “Ebola in the air? A nightmare that could hap- 181
- 135 pen”⁴ 182
- 136 – “Air pollution now leading cause of lung cancer”⁵ 183

137 In order to automatically detect misleading head- 184

138 lines, the news body text must be analyzed to ex- 185

139 tract the evidence from which the headline has 186

140 been derived, thereby detecting the headline/body 187

141 text discrepancy in the absence of such evidence. 188

142 The task of identifying the relation between a head- 189

143 line and the news article it refers to has been 190

144 addressed in recent research (see Section 2) as 191

145 a stance detection problem. This type of ap- 192

146 proach involves estimating the relative perspective, 193

147 namely the stance, of one piece of text, such as a 194

148 claim or a news article, towards another, for exam- 195

149 ple, a topic, a statement or a headline [17]. 196

150 In the context of headline/body text dissonance, the 197

151 main objective of this research is to propose an ap- 198

152 proach that relies on semantics and deep learning tech- 199

153 niques to automatically determine the stance of the 200

154 headline with respect to its body text. By this means, 201

155 the problem of misleading headlines can be addressed. 202

156 The approach is hereafter referred to as *HeadlineS-* 203

157 *tanceChecker*. Given a news headline and its corre- 204

158 sponding body text, our proposal assigns the headline 205

159 one of these four classes— *unrelated*, *agree*, *disagree* 206

160 or *discuss*—, indicating the headline stance, and vali- 207

161 dating and checking whether the headline is faithfully 208

162 reflecting the information provided in the news article.

163 The most interesting aspect of solving misleading 209

164 headline detection as a stance detection task is that it 210

165 is not only focused on determining whether or not a 211

166 headline is consistent with its body text, but it is also 212

167 a fine-grained classification that determines the type of 213

168 dissonance involved. 214

169 We explore the treatment of this task as a two-stage 215

170 neural classification problem in which only the essen- 216

171 tial information of the news item is processed, rather 217

172 than the whole news item. We therefore use the sum- 218

173 maries because besides containing the key information 219

174 of the news story, we hypothesize that the abridged ver- 220

175 sion will not only increase task efficiency but also that of 221

176 the the neural models. Neural models can have a neg- 222

177 ative impact on efficiency when processing long texts, 223

178 so previous studies either used the first sentence of the 224

179 text [18] or a specific fragment [19] to combat this prob- 225

180 lem. Therefore, the use of text summarization, which, to 226

the best of our knowledge, has not been previously ex- 181

ploited for stance detection, could be beneficial in this 182

context. 183

To summarize, the main novelties of *HeadlineS-* 184

tanceChecker are twofold: 185

- the adoption of a divide-and-conquer strategy by 186
- proposing a two-stage neural classifier for per- 187
- forming the headline stance task; and 188
- the use of summarization techniques based on Po- 189
- sitional Language Models (PLM). These models 190
- leverage semantic knowledge to detect the evi- 191
- dences and essential information within the news 192
- article so as to generate automatic summaries that 193
- will be used as substitutes of the full body text for 194
- the whole classification process. We expect this 195
- approach to be more efficient in dealing with the 196
- headline stance classification task. 197

The paper is structured as follows: Section 2 presents 198

the related work regarding misleading headlines, as well 199

as a brief review of the state of the art in text summariza- 200

tion; Section 3 presents our proposed architecture for 201

HeadlineStanceChecker—explaining each of the stages 202

in detail—; Section 4 describes the experiments carried 203

out and the evaluation environment; Section 5 reports 204

and discusses the results of the proposed approach— 205

comparing them to other competitive systems—; and , 206

Section 6 presents conclusions and outlines the main di- 207

rection for future work. 208

209 2. Related Work

HeadlineStanceChecker has been conceived as an au- 210

tomatic method to classify a news story in terms of the 211

relation between its body text and its headline. The 212

main motivation of developing such approach is to pro- 213

vide a tool that helps both professionals and readers to 214

identify misleading or fraudulent media and informa- 215

tion, thus preventing harmful consequences. 216

Fake news research has opened up an immense 217

field of work that encompasses multiple areas and ap- 218

proaches. Both linguistic and non-linguistic aspects are 219

being studied, so that elements as diverse as image ver- 220

ification, analysis of reputation and authorship, or the 221

network dissemination patterns of misleading stories 222

fall within its field of interest. For brevity, we focused 223

on the research directly related to our proposal, but com- 224

prehensive studies can be found in [8, 20, 21]. 225

Therefore, in this section, first we present an 226

overview of recent work done in Stance Detection and, 227

⁴<https://cnn.it/2NeuNZj> (accessed online 15 February, 2021)

⁵<https://bit.ly/2Tajxxx> (accessed online 15 February, 2021)

228 next, an in-depth review of the existing detection strate- 278
229 gies for misleading headlines is conducted. Finally, 279
230 given that one of the novelties of the paper is using sum- 280
231 marization techniques leveraging essential information 281
232 to characterize headlines, a brief review is presented of 282
233 the state of the art in text summarization. 283

234 • Stance Detection Overview 284

235 From an overall perspective, stance detection can 285
236 be defined as the task of identifying the perspec- 286
237 tive of an author or text against a given target in 287
238 the form of one topic, claim, headline or even a 288
239 personality [22, 23]. Hence, there exists a tuple of 289
240 elements—the text on the one side, the target on 290
241 the other side—and a classification process shaped 291
242 to determine how the former stands towards the 292
243 second: does the text support the topic? does it 293
244 disagree with the claim? The names of the classes 294
245 (e.g. *support*, *against*, *for* or *neutral*) depend on 295
246 the precise problem. The task, which concerns a 296
247 diverse range of domains, is studied in such varied 297
248 areas as political debates [24, 25], student essays 298
249 [26], online forum debates [27] or even internal 299
250 company discussions [28, 29]. 300

251 A great deal of work in opinion mining has been 302
252 devoted to detect the stance of tweets or other 303
253 types of short texts as rumours [30] or microblog- 304
254 ging statements. Examples of targets posed in 305
255 the available datasets could be “Hillary Clinton” 306
256 for personality, “Atheism” as a particular topic 307
257 or the claim “E-cigarettes are safer than normal 308
258 cigarettes”. Shared tasks offering such datasets and 309
259 fostering the research on the matter have arisen in 310
260 different languages. SemEval-2016 posed the sub- 311
261 task for detecting stance in tweets [31], providing 312
262 around 5 thousand tweets in English covering five 313
263 commonly known topics. The task has inspired nu- 314
264 merous approaches that develop either traditional 315
265 proposals (e.g. K nearest neighbour [32], Support 316
266 Vector Machines [33] or latent features provided 317
267 by methods such as Latent Dirichlet Allocation 318
268 [34]); or those inspired by neural network frame- 319
269 works, by using, for example, bidirectional con- 320
270 ditional encoding [35], bidirectional Long Short- 321
271 Term Memory neural networks [36] or Attention 322
272 based Convolutional Neural Networks [37]. Be- 323
273 sides, there are available public datasets that sup- 324
274 port the development of new interesting work, such 325
275 as the *Multi Perspective Consumer Health Query* 326
276 *dataset* [38] dedicated to detecting the stance of 327
277 sentences collected from quality articles towards

five different claims (e.g., “Sun exposure causes skin cancer”). In [23], an in-depth study on different approaches to the two tasks mentioned above can be found. Regarding languages other than English, the necessity for well-annotated data led to the proliferation of both annotation efforts and shared tasks aimed to advance research, such as *StanceCat*, presented at IberEval 2017 as a stance detection task for tweets in Spanish and Catalan [39], a proposal and a dataset of short messages in Russian internet forums [40], or even projects combining a larger number of languages (French, Italian, Spanish, English) [41, 42].

In contrast to such approaches, research on stance detection based on longer documents, as in the current scenario, faces different challenges. Dealing with discourse, as a coherent and cohesive set of sentences, adds a certain complexity not present when processing shorter utterances. Within the discourse, an argument may develop in such a way that some sentences may show support for the claim, while others may seem to deny it, and only by considering the document as a whole can the stance be effectively identified.

It is in this context that *HeadlineStanceChecker* has been developed, and next, we introduce the related work concerning the specific task.

305 • Misleading headlines 305

306 The task of detecting misleading headlines for the 306
307 present research involves classifying the stance of 307
308 the article body with respect to the claim made in 308
309 the headline into one of the following four classes: 309
310 a) *agrees*—agreement between body text and head- 310
311 line; b) *disagrees*—disagreement between body 311
312 text and headline; c) *discusses*—same topic dis- 312
313 cussed in body text and headline, but no position 313
314 taken; and, d) *unrelated*—different topic discussed 314
315 in body text and headline. 315

316 This task—headline stance detection—quickly 316
317 emerged in the context of fake news analysis, trig- 317
318 gered by a demand for new technologies to prevent 318
319 and combat the phenomenon, together with an in- 319
320 crease in the availability of annotated corpora [8]. 320
321 In this context, research challenges and competi- 321
322 tions were proposed. The most recent and impor- 322
323 tant ones are next reviewed in detail. 323

324 The *Fake News Challenge*⁶ (FNC-1) [43] was cre-

⁶<http://www.fakenewschallenge.org/> (accessed online 15 February, 2021)

ated using Emergent dataset [17] as a starting point (this dataset has been extracted from the Emergent Project [44], a rumour debunking project). FNC-1 aims to compile a gold standard to explore Artificial Intelligence technologies, especially ML and Natural Language Processing (NLP), applied to detection of fake news. To carry out this macro-challenge, the organisers decided to start with stance detection. In this case, the FNC-1 dataset was released, with around 75,000 instances that were classified as follows: *agree*, *disagree*, *discuss* and *unrelated*.

For example, given the headline “*Robert Plant Ripped up \$800M Led Zeppelin Reunion Contract*”, the following fragments⁷ would illustrate the different classes mentioned, according to the gold-standard annotations in the FNC-1 dataset:

- **Agrees:** The body text agrees with the headline. Example evidence: “[...] *Led Zeppelin’s Robert Plant turned down 500 MILLION pounds to reform supergroup.*”
- **Disagrees:** The body text disagrees with the headline. Example evidence: “[...] *No, Robert Plant did not rip up a \$800 million deal to get Led Zeppelin back together.*”
- **Discusses:** The body text discusses the same topic as the headline, but does not take a position. Example evidence: “[...] *Robert Plant reportedly tore up an \$800 million Led Zeppelin reunion deal.*”
- **Unrelated:** The body text is not related with the headline. Example evidence: “[...] *Richard Branson’s Virgin Galactic is set to launch SpaceShipTwo today.*”

The FNC-1 competition received a total of 200 submissions achieving relative scores⁸ of around 82% in the best ranked submissions. The organization proposed a simple baseline using hand-coded features and a gradient boosting classifier, available at Github⁹. The three best systems in this competition were Talos [45], Athene system [46] and UCLMR [47] in this order. Talos [45]

applied a one-dimensional convolution neural networks (CNN) on the headline and body text, represented at the word level using Google News pre-trained vectors. The output of this CNN is then sent to a multilayer perceptron (MLP) with 4-class output: *agree*, *disagree*, *discuss*, and *unrelated*, and trained end-to-end. Using this combination CNN-MLP, the system outperformed all the submissions and achieved the first position in the FNC-1 challenge.

Recently, other works used the FNC-1 for their experiments and the performance obtained in the competition improved. For instance, [48] addressed the problem proposing a hierarchical representation of the classes, which combines *agree*, *disagree* and *discuss* in a new related class. A two-layer neural network is learning from this hierarchical representation of classes and a weighted accuracy of 88.15% is obtained with their proposal. Furthermore, [49] constructed a stance detection model by performing transfer learning on a RoBERTa deep bidirectional transformer language model by taking advantage of bidirectional cross-attention between claim-article pairs via pair encoding with self-attention. They reported a weighted accuracy of 90.01%.

Outside the FNC-1 Challenge and dataset, there is other research that also addresses the stance detection tasks, determining the relation of a news headline with its body text. Some authors extracted key quotes [50] or claims [51] to facilitate the detection. There is also work related to argument mining analysis, in which the headline represents an argument that is not supported by claims in the text. Moreover, in addition to using argument mining for solving stance detection, this problem could benefit from other tasks which detect semantic relations within the text, such as contradiction [52], contrast [53] and entailment [54].

• Text Summarization

Previous research in Text Summarization has been shown to have a positive impact on society since the use of summaries has been beneficial in different areas, such as education —where summaries are used to support reading comprehension tasks [55, 56, 57, 58]— business, by producing, for instance, an automatic summary of event logs to help analysts [59], or health, regardless of whether the summaries were created manually [60, 61], or automatically [62]. This is partly due to the capa-

⁷Examples extracted from Fake News Challenge website fake-news-challenge.org.

⁸Measure score used in the Fake News Challenge competition

⁹<https://github.com/FakeNewsChallenge/fnc-1-baseline> (accessed online 15 February, 2021)

bility of summarization methods to identify the most relevant information of a document, and condense it into a new text, thereby helping to reduce time and resources when it comes to manage large amounts of data. These methods have proven to be effective when integrated as an intermediate component of more complex systems.

The journalism field, and specifically the news domain, has been one of the most representative areas in which summarization has traditionally focused from the outset, partly thanks to the development of appropriate corpora (e.g. DUC, Gigaword, CNN/DailyMail)[63], and the wide range of techniques and approaches to help digest this type of information [64, 65, 66, 67]. Besides the various summarization types that have been developed for this domain (single-document, multi-document, extractive, abstractive, generic, topic-oriented, etc.), there is a significant amount of research on the task of headline generation using summarization techniques [68, 69, 70], and more recently using Deep Learning [71, 72, 73]. However, none of them have exploited either the headline or the summarization techniques as an intermediate stage to further extract the semantic relationship between the headline and the news body text, and detect possible incongruities to fight against the fake news problem.

Although summarization has been used for fake news detection [74, 75, 76], as well as in the context of online discussions and social media to detect whether the author of a comment is in favor of or against a given target (e.g. entity or topic) [77, 78], to the best of our knowledge, summarization has not been directly applied to the stance detection problem of misleading headlines, as proposed in this study.

Summarization was mentioned as a potential effective methodology for dealing with the problem of incongruent headlines in [79], but from a different perspective, which involved using headline generation to create a new headline that could be then compared to the existing headline by measuring the distance between them. More recently, an updated comprehensive survey concerning the stance detection task [80] shows that there is a lack of research where summarization is applied to this task, although a new type of summarization, called stance summarization, is outlined. However, stance summarization involves the generation of a new type of summary which includes a stance, but it is not comparable to the approach presented in this paper as the summaries are not incorporated into the stance de-

tection process.

In another survey, conducted by [81], the authors compile the available information regarding existing research addressing this problem, and only the work of [82] summarized the news body into a single sentence to be compared to the given claim and determine its overall veracity, an approach which aligns with that suggested in [79] as aforementioned. By contrast, our research goes beyond summarizing the whole document into just one sentence, and provides a summary that could be acted as a substitute of the whole body text.

The *HeadlineStanceChecker* proposal is based on the fact that semantic information and discourse structure are captured through PLMs which, in turn, are exploited as a summarization technique. PLMs allow key spots and relevant information to be located in the news body text, and they are then used to create a summary of the news. By this means, the news article is reduced to its essential information, which is then compared to its headline. Our proposed model to detect misleading headlines, by relying on their stance towards the article’s content, directly uses this summary of the news instead of the whole news body text, enabling a more accurate comparison to its headline.

3. HeadlineStanceChecker Architecture

The *HeadlineStanceChecker* approach involves two-stages (see Figure 1), thus addressing the task as a two-level classification problem. The first level corresponds to a *Relatedness Stage*, while the second corresponds to a *Stance Stage*. An additional novelty is the use of summaries generated in the first stage for the whole process instead of the full body text (i.e., the *Relatedness Stage*).

In this manner, given the inputs, namely the candidate headline and the news article body text, a summary of the news body will be created in the *Relatedness Stage* to later determine the headline’s stance regarding the news article as either *related* or *unrelated*. Afterwards, in the *Stance Stage*, the examples classified as *related* in the previous stage, are further classified into three possible values: *agree*, *disagree*, or *discuss*.

A more detailed description of both stages and the different modules involved in performing the stance classification is provided here-under.

3.1. Relatedness Stage

The *Relatedness Stage* will determine whether the headline is classified as *related* or *unrelated* with respect to the body text of the news article. The inputs of this stage are both the text body and the headline,

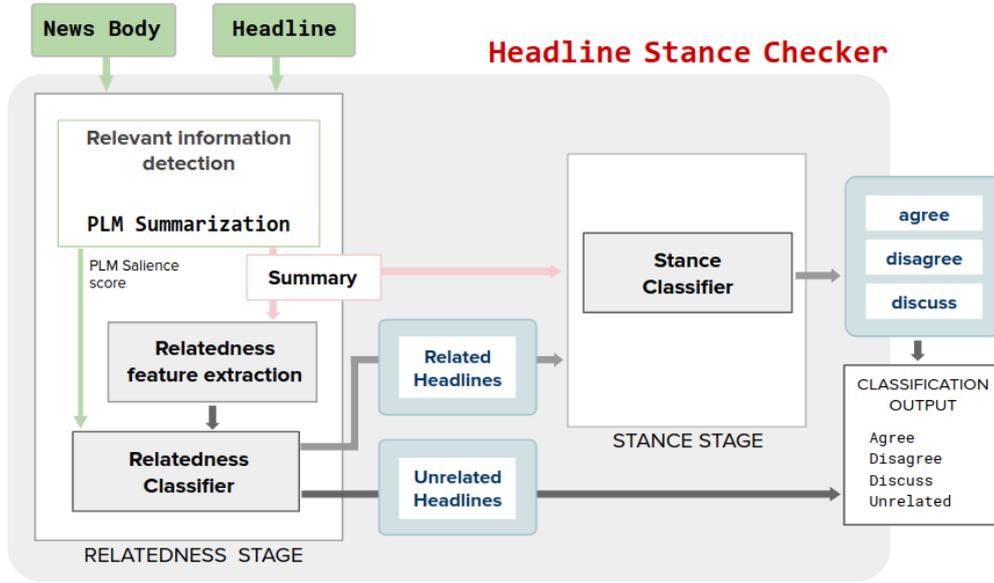


Figure 1: *HeadlineStanceChecker* architecture.

516 resulting in a binary classification. The outputs of this
517 stage are:

- 518 • The *headlines* classified as *related* or *unrelated*.
- 519 • The *summary* of the news content, obtained in a
520 relevant information detection module.

521 To produce the above outputs, three modules are pro-
522 posed: i) relevant information detection; ii) relatedness
523 feature extraction; and, iii) relatedness classification.

524 3.1.1. *Relevant information detection module*

525 The relevant information detection module aims to
526 create a summary revealing the important information
527 of the input news article in relation to its headline.

528 The task of summarization has generally been carried
529 out from a statistical perspective that only considers the
530 elements of the text with no regard to their structure (or
531 in those cases where the structure is taken into account,
532 it is already known beforehand, such as in the case of
533 scientific articles). PLMs are a type of statistical lan-
534 guage model that allow information to be considered by
535 taking into account both the relevant elements of the text
536 and their location in the document. This provides a dy-
537 namic method for detecting key aspects of the text inde-
538 pendently of the domain and textual genre to which it is
539 related. Conversely, PLMs represent a type of statistical
540 language model that allows information to be consid-
541 ered by taking into account both the relevant elements
542 of the text and also their location in the document. They

543 define a dynamic method for detecting key aspects of
544 the text independently of the domain and textual genre
545 to which it is related. Besides, PLMs have proved valu-
546 able in other areas such as information retrieval [83] and
547 language generation [84].

548 From the semantic perspective of the text, its essence
549 can be more effectively captured and synthesized by
550 considering the document not as a mere sequence of
551 sentences, but as a coherent and cohesive source of
552 meaning, traversed by semantically related entities and
553 actions. Considering this, we chose PLMs as the corner-
554 stone of our module given that they can be configured
555 to include the identification of named entities within the
556 story, together with the representation of the words as
557 synsets (sets of synonyms accounted under an identi-
558 fier), allowing a further abstractive step on the basis of
559 Wordnet [85], a hierarchical database of semantic re-
560 lations. Consequently, PLMs help to incorporate both
561 the semantics derived from the relevant lexical units to-
562 gether with the meaning derived from the text as coher-
563 ent discourse. Previous studies demonstrated that PLMs
564 were suitable for summarization tasks [86] and, more-
565 over, a preliminary research was conducted analyzing
566 and comparing different summarization methods for the
567 stance detection task—including extractive, abstractive
568 and hybrid ones [86]—also showing that PLM-based
569 summarization yielded the most stable results.

570 **PLM essentials.** Fundamentally, the PLMs state that
571 for every position i within a document D it is possible
572 to calculate a score for each element w that belongs to

573 the document’s vocabulary. The decision as to the kind 606
 574 of elements that compose the vocabulary is made when 607
 575 designing the module. The calculated score displays the 608
 576 relevance of each element w in every precise position i , 609
 577 based on its distance to other occurrences of the same 610
 578 element throughout the document. The score is higher 611
 579 when the neighbor element is closer within a scope to 612
 580 compute the value that goes beyond the sentence limits,
 581 taking into account the whole document. In order to
 582 express the distance to the occurrence of the entity in the
 583 neighbourhood, a propagation function $f(i, j)$ is applied.

Equation 1 defines how the score for word w in position i is computed:

$$P(w | i) = \frac{\sum_{j=1}^{|D|} c(w, j) \times f(i, j)}{\sum_{w' \in V} \sum_{j=1}^{|D|} c(w', j) \times f(i, j)} \quad (1)$$

584 where $|D|$ refers to the length of the document, V is
 585 the vocabulary, $c(w, j)$ indicates the presence of element
 586 w in the position j , and $f(i, j)$ is the propagation func-
 587 tion that rates the distance between i and j . In this case,
 588 and taking as bases previous work on the matter [84], a
 589 Gaussian kernel is adopted as the propagation function.

590 Figure 2 illustrates the idea behind the PLM reason-
 591 ing.

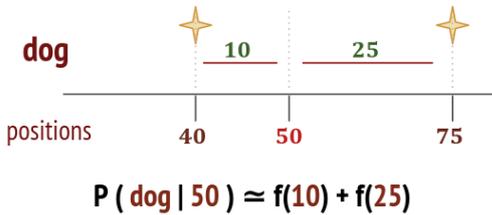


Figure 2: Example of the type of computation performed to obtain the value of the PLM for the position 50 regarding the element *dog* of the vocabulary

592 **PLM for summarization.** The manner in which
 593 PLMs are employed to perform the summarization task
 594 comprises three stages. First, we need to conduct the
 595 definition of the vocabulary as a parameter for the PLM
 596 module. In our current configuration, the vocabulary is
 597 composed of the synsets corresponding to nouns, verbs
 598 and adjectives, together with the named entities that ap-
 599 pear along the text. In order to get this semantic infor-
 600 mation, we use Freeling [87], an open source tool that
 601 allows linguistic analysis with different levels of granu-
 602 larity.

603 From this stage, a representation of the text that in-
 604 volves both the vocabulary and the positions of its ele-
 605 ments is obtained.

Second, we create a seed, i.e., a set of words that can be significant for the text and will help the system to discard irrelevant parts of the discourse. The given headline is taken as seed in our configuration. It needs to be analyzed with the same tools as the source text (Freeling). As a result, a second vocabulary is then built from it.

Finally, the processing of the PLM against the seed allows us to compute scores for the text elements that are now conditioned by the information in the headline. At this stage, we have already calculated a collection of values associated with every relevant element using the PLMs for each position of the text. The aggregation of the different values related to each of those positions results in a vector with same length of the document, the Score Counter (*SC*), so that the score held in the index i will express the value for the position i in the text. Those positions in the text that show local maximums in the *SC* are retrieved as the most relevant points of the document. The sentences to which these positions belong are then selected as candidates for the summary. Since a value has been calculated for each position in the sentence, we can obtain a score for the sentence itself:

$$S_{\text{score}} = \sum_{i \in S} SC[i] \quad (2)$$

613 with S representing the sentence, and i indicating the
 614 positions within the document for that sentence.

These values also allow us to select from the candidates the sentences that will constitute the news extractive summary. Moreover, the computed values are necessary to define a new feature, named **PLM Saliency Score**, which will be used for the relatedness classifier in the next step. Its value derives from the aggregation of each score S_t associated with each sentence t included in the summary, following Equation 3. Let S^* represent the set of the sentences belonging to the summary, the *PLM Saliency Score* for a document D would be calculated as:

$$PLM\ Saliency\ Score_D = \sum_{t \in S^*} S_t \quad (3)$$

615 3.1.2. Relatedness Feature Extraction

616 Besides the relevant information (i.e., the summary)
 617 and the PLM Saliency Score obtained in the previous
 618 module, two similarity features are used as input to the
 619 relatedness classifier applied next. To obtain the fea-
 620 tures, the headline and the generated summary are used.
 621 They are described next:

- **Cosine similarity:** The cosine similarity between headline and summary of body text is computed.

This feature is used to measure how similar the headline and summary are, irrespective of their size. Mathematically, it measures the cosine of the angle between two vectors projected in a multi-dimensional space [88]. The cosine similarity is advantageous because even if the two similar documents are far apart by Euclidean distance (due to the size of the document), chances are that they may still be oriented closer together. The smaller the angle, the higher the cosine similarity [89]. Although this metric is relatively basic [90], it usually brings significant improvements to retrieval models [88]. The cosine similarity measure between two vectors X and Y is obtained following Equation 4 [91]:

$$\text{Cosine similarity}(|X, Y|) = \frac{x \cdot y}{\|x\| \|y\|} \quad (4)$$

$$\text{where } x \cdot y = \sum_{i=1}^n x_i y_i \text{ and } \|x\| = \sqrt{x \cdot x}$$

For the calculation of cosine similarity, the text pairs are converted into Term Frequency-Inverse Document Frequency (TF-IDF) vectors, using the tools provided by scikit-learn [92].

- **Overlap coefficient:** This feature is defined as the intersection between two sets A and B . In the current scenario, these sets contain the ngrams belonging either to the headline or the summary [93]. The overlap coefficient is given by Equation 5 [94]:

$$\text{Overlap coefficient}(A, B) = \frac{A \cap B}{\min(|A|, |B|)} \quad (5)$$

If set A is a subset of B or the converse, then the overlap coefficient is equal to 1 else overlap coefficient should be between 0 to 1 [95].

3.1.3. Relatedness classification

The relatedness classification module exploits the PLM Saliency Score, the relatedness features previously processed, as well as the summary to finally classify the headlines as *related* or *unrelated*. The proposed architecture is flexible to choose any model that allows classifiers to be improved.

In this case, the design of the relatedness classification module is based on fine-tuning the RoBERTa (Robustly optimized BERT approach) pre-trained model [96], applying a classifier to its output afterwards.

First, the headline and the summary are concatenated and processed with the RoBERTa model. The resulting vector is consecutively multiplied by the three features (PLM Saliency Score, Cosine similarity, Overlap coefficient) to finally carry out the classification using a Softmax activation function in the output layer.

Specifically, we have chosen RoBERTa Large model (24 layer and 1024 hidden units) since it achieves state-of-the-art results in General Language Understanding Evaluation (GLUE) [97], Reading Comprehension Dataset From Examinations (RACE) [98] and Stanford Question Answering Dataset (SQuAD) benchmark. Similar to [49, 96, 99], in this work we fine-tune RoBERTa to efficiently address a task that involves comparing sentences. RoBERTa optimizes Bidirectional Encoder Representations from Transformers (BERT) [100] by adding several modifications but without altering the original architecture, an approach that improves the results with respect to BERT in the main NLP tasks [96]. Some of those modifications involve: eliminating the prediction of the next sentence; performing the training on a greater volume of data; enlarging the batch size; and, lengthening the input sequence.

To create the classifier, the *Simple Transformers library*¹⁰ was used, which creates a wrapper around *HuggingFace's Transformers library* for using Transformer models [101]. *Simple Transformers* is an NLP library that allows the modification of hyperparameters so as to train, evaluate, and make predictions using the best state-of-the-art models.

In our model, the hyperparameter values are: maximum sequence length of 512; batch size of 4; training rate of 1e-5; and, training performed for 3 epochs. These values were established after successive evaluations, following previous experiments on this model [49, 96, 99].

3.2. Stance Stage

Given the related headlines obtained through the first stage on the proposed architecture, the main goal of this stage is to determine their type considering the remaining stances: *agree*, *disagree* or *discuss*. Therefore, the claim made in the headline can be finally classified into one of three classes left.

The inputs of this stage are:

- The *headlines* classified as *related*.
- The *summary* of the news content.

¹⁰<https://simpletransformers.ai/> (accessed online 15 February, 2021)

The output of this stage then is the final classification of the related headlines, where each of them is assigned one of the following possible stance values: *agree*, *disagree* or *discuss*. These classified headlines together with the *unrelated* headlines determined before, will comprise the final output for the whole *HeadlineStanceChecker* approach.

3.2.1. Stance classification

As in the *Relatedness Stage* (Section 3.1.2), the extractive summary generated in Section 3.1.1 is also used here.

Similar to the *Relatedness* classification module, this stage has been built using RoBERTa as the selected model capable of improving the classification. In this case, no additional features are considered, only two dense layers are included to reduce dimensions and, finally, the Softmax classification layer. The hyperparameters of the model used in this classifier are the same as those of the *Relatedness* classification, except for the classification output which in this case is of three classes: *agree*, *disagree*, *discuss*.

4. Experiments and Evaluation Environment

The proposed approach was applied and evaluated in the context of the Fake News Challenge FNC-1 whose goal was to determine a headline’s stance by classifying it in relation to its body text into 4 classes: *unrelated*, *agree*, *disagree*, and *discuss*. In this section, we first describe the corpus provided in this challenge. Second, we explain the experiments performed with different configurations of our system. Finally, the evaluation metrics used are outlined. The results obtained will be presented, discussed and compared to the other participating systems in the challenge in subsequent sections.

4.1. Fake News Challenge Dataset

The experimentation is conducted over the FNC-1 dataset whose instances are labeled as *agree*, *disagree*, *discuss* and *unrelated*.

The dataset was split into a training set (66.3%) and a testing set (33.7%), where neither the headlines nor the body text overlapped. The distribution of documents (bodies and headlines) is presented in Table 1.

As the distribution of the classes indicates in Table 2, there is a significant imbalance for both the training and testing sets where the instances of the *unrelated* class alone (over 70%) are greater than the sum of the remaining classes. At the other extreme, the *disagree* class is remarkably lower compared to the others.

	Documents	Headlines	Instances
Train set	1,683	1,683	49,972
Test set	904	904	25,413
Complete dataset	2,587	2,587	75,385

Table 1: Description FNC-1 dataset considering number of documents.

4.2. Experiments

To measure our system’s performance, a set of experiments was conducted as follows, the results of which will be shown and discussed in Section 5. Our experiments can be replicated at Github¹¹:

- *Relatedness Stage Validation*: The aim of this experiment is to assess the performance of this classification stage, where *related* or *unrelated* headlines are initially identified. First, we analyze and compare the performance of the classifier when either summaries or the full body is employed. Second, we conduct an ablation study to verify whether the *relatedness* features used for the classifier make a positive contribution.
- *Stance Stage Validation*: The goal of this experiment is to determine how accurate the *Stance Stage* is when the errors produced by the *Relatedness Stage* are avoided, thereby using an ideal input for this stage, i.e., the gold-standard headlines annotated as *related* in the FNC-1 corpus. By this means, we can measure the effectiveness of this stage in isolation. Furthermore, to validate the extent to which our proposed stance detection model can be generalized, we apply it to a different headline stance detection dataset, i.e. Emergent dataset [17].
- *HeadlineStanceChecker Validation*: In this last experiment, the entire system —integrating the *Relatedness* and *Stance* classifiers as a two-step classifier and using summaries as input for the whole process instead of the full text— is tested. Its performance is then compared to other configurations of the model as well as to competitive state-of-the-art systems.

In addition, we also investigate the system performance considering two different inputs: summaries and full body. This experiment and its further analysis is detailed in Section 5.4.

¹¹<https://github.com/rsepulveda91112/HeadlineStanceChecker> (accessed online 10 september, 2021)

	Agree	Disagree	Discuss	Unrelated
Train set	3,678 (7.36%)	840 (1.68%)	8,909 (17.82%)	36,545 (73.13%)
Test set	1,903 (7.48%)	697 (2.74%)	4,464 (17.56%)	18,349 (72.20%)
Complete dataset	5,581 (7.4%)	1,537 (2.03%)	13,373 (17.73%)	54,894 (72.81%)

Table 2: Distribution of FNC-1 dataset stances.

4.3. Evaluation Metrics

Originally, the organizers of the FNC-1 challenge proposed the *Relative Score* metric, which assigned a higher weight to examples correctly classified, as long as they belonged to a different class from the *unrelated* one. The rationale behind this metric was to address the highly imbalanced distribution of the classes caused by the over-representation of the *unrelated*.

However, as pointed out in [102], the inner imbalance among the three *related* classes —*agree*, *disagree*, and *discuss*— was not addressed. Therefore, following [102], this study incorporates, in addition to the FNC-1 relative score, both a measure of F_1 class-wise and a macro-averaged F_1 (F_{1m}) as the mean of those per-class F scores so as to address the imbalance among the less represented classes. The advantage of this measure rests in it not being affected by the size of the majority class. Additionally, average accuracy is also obtained.

5. Results and Discussion

This section presents the results obtained in each of the experiments described in Section 4.2. The values are expressed in percentage mode.

5.1. Relatedness Stage Validation

Our first experiment was designed to evaluate the first module as an isolated element of the system, acting as a binary classifier. In this case, we were not evaluating the system to detect *agree*, *disagree* or *discuss* examples, but to perform *related* versus *unrelated* classification. We carried out an analysis of the classification results and also an ablation study that considered the following involved features: cosine similarity; PLM Saliency Score; and, overlap coefficient.

The performance of the relatedness classifier was first validated by analyzing whether the use of summaries had a positive impact on the output compared to using the whole document. The results are shown in Table 3. Both approaches used the three features previously described in the section 3.1.1 and 3.1.2.

Relatedness Stage FNC-1-Summary refers to an experiment that uses summaries both to calculate features and to enter the classification model, whereas the *Relatedness Stage FNC-1-Body-text* approach uses the body text instead of the summary as input to the relatedness stage to classify the headline. The results validate the use of summaries as a useful approach to the stance detection problem as even if some information is excluded, the findings indicate a slight improvement when using the summarized text.

The approach that uses summaries throughout the process is able to improve the related class, which is the minority class. Figures 3 and 4 show each confusion matrix of the two approaches with minimal variation in the classification, thus showing that the use of summaries does not harm the results of this classifier.

These results show that, by using the PLMs to condense the relevant information from a piece of news, the resulting summaries offer an attractive substitute for the full news text, enabling a reduction of the computational load for the classifiers, which increases when dealing with longer texts.

Furthermore, to evaluate the influence of the added features in the relatedness stage, an ablation study of the features extracted from the summary has been conducted. Each feature (Cosine similarity, PLM Saliency Score and Overlap coefficient) has been removed and an experiment has been designed that will return the results of the classification without the incidence of the removed feature. To the extent that the classification result is worse, this would imply that the eliminated feature has a great influence on improving the classification results. The most influential feature for the classification was observed to be the PLM Saliency Score as the experiment that does not use the PLM score obtains the worst results, followed by the one that does not use overlap coefficient and, finally, by the one that uses cosine similarity. Table 4 shows the ablation study.

5.2. Stance Stage Validation

This experiment was designed to determine the validity of the *Stance Stage*. This task can be tackled as a double question, since two fundamental issues arise:

System	F_1 Score		F_{1m}
	Related	Unrelated	
<i>Relatedness Stage FNC-1-Summary</i>	98.38	99.40	98.89
<i>Relatedness Stage FNC-1-Body-text</i>	98.36	99.37	98.86

Table 3: Relatedness classification results: class-wise F_1 Score and F_{1m} using automatic summaries vs. full news text.

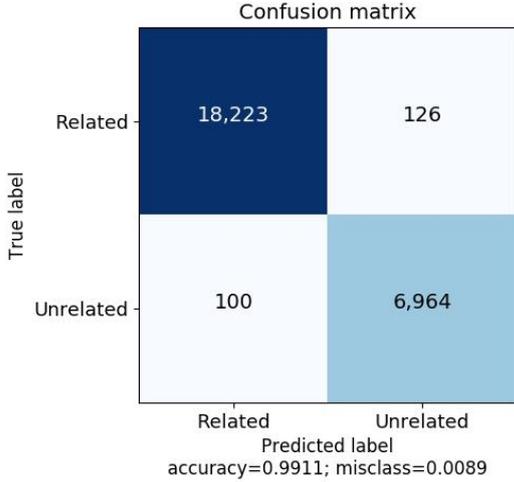


Figure 3: Confusion matrix resulting from the *Relatedness Stage FNC-1-Summary*.

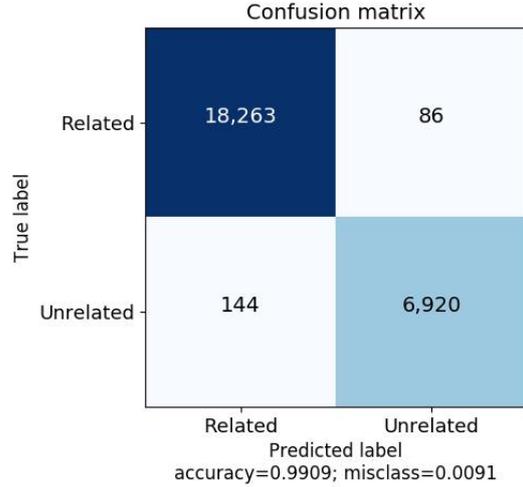


Figure 4: Confusion matrix resulting from the *Relatedness Stage FNC-1-body*.

Removed feature	F_1 Score		F_{1m}
	Related	Unrelated	
<i>Cosine similarity</i>	98.24	99.32	98.78
<i>PLM Saliency Score</i>	98.00	99.23	98.61
<i>Overlap coefficient</i>	98.10	99.27	98.68
<i>Non-ablated results</i>	98.38	99.40	98.89

Table 4: Ablation study results for the features used in the Relatedness Stage. To facilitate reading and comparison, we have also included the non-ablation results.

873 i) the validity of the *Stance Stage* as a general proposal;
874 and ii) the effectiveness of the *Stance Stage* performance
875 within an ideal case.

876 As for the first issue, this experiment aims to demon-
877 strate that the approach is not an ad-hoc solution but
878 a general one. For this purpose, the *Stance Stage* was
879 applied to a different stance dataset called Emergent

880 dataset¹² [17]. For this dataset, each example results
881 from a combination of one article and its headline, and
882 one claim. There are three different options for de-
883 scribing how a claim can be related to a piece of news.
884 Specifically, each example was manually labeled by a
885 journalist with one of the following tags: *for*, if the ar-
886 ticle states that the claim is true; *against*, if it states that
887 the claim is false and *observing*, when the claim is re-
888 ported in the article, but without assessment of its verac-
889 ity. The dataset is composed of 2,595 examples, derived
890 from the combination of 2,571 news, 2,536 headlines
891 and 300 claims (see Table 5 for further details of the
892 dataset).

893 To replicate our experimental environment with this
894 dataset, the equivalence between labels in both datasets
895 regarding their meaning is *for* \approx *agree*, *against* \approx *dis-*
896 *agree* and *observing* \approx *discuss*.

¹²<https://github.com/willferreira/mscproject/> (accessed online 15 February, 2021)

	News Bodies	Headlines	Claims	For	Against	Observing	Total Examples
Train	2,048	2,023	240	992	304	775	2,071
Test	523	513	60	246	91	187	524
Complete dataset	2,571	2,536	300	1,238	395	962	2,595

Table 5: Description of the Emergent dataset: number of documents and distribution of assigned labels.

Experiment	F_1 Score			F_{1m}	Acc
	Agree	Disagree	Discuss		
Testing Emergent					
<i>Emergent Upper Bound</i>	81.53	74.53	68.23	74.76	76.15
<i>Stance Stage Emergent</i>	75.15	77.77	65.49	72.80	71.89
<i>Stance Stage Emergent Test FNC-1 Training</i>	73.15	73.68	70.61	72.48	72.08
Testing FNC-1					
<i>Stance Stage FNC-1</i>	72.87	63.50	88.74	75.04	82.30

Table 6: *Stance Stage* results: class-wise F_1 Score, F_{1m} and overall accuracy on FNC-1 and Emergent dataset.

Therefore, to validate the generalization of the approach, Table 6 includes the following performance results:

- *Emergent Upper Bound*: This experiment is performed as an upper bound by using a human-written headline created by a journalist, and considering it as a perfect summary that comprises the main information of the news body text. Nevertheless, this upper bound is only applicable to the Emergent dataset since in the case of FNC-1 no journalist-written headline is provided, as occurs in the case of the Emergent dataset.
- *Stance Stage using Emergent Dataset*: Our model is trained with the Emergent dataset and the Stance Stage is applied to it.
- *Stance Stage tested with Emergent, but trained with FNC-1*: This performance uses the Emergent dataset to test the Stance Stage but with the model trained on the FNC-1 so as to demonstrate the extent to which our proposal can be generalised.

The second aspect that needs to be addressed relates to the appropriateness of this second stage and its performance by isolating this stage from the rest of the system. The strategy here is focused on avoiding the errors inherited from the previous stage. To achieve this, only the examples tagged as *related* from the FNC-1 Gold-Standard are used and evaluated. The results of this

performance correspond to *Stance Stage FNC-1* row in Table 6.

The analysis of the results obtained in this stage regarding the comparison of the performance using Emergent dataset are very promising considering that this model is using automatic summaries. The results are very close to the upper bound obtained by using human-made summaries. Analyzing per class, using the Emergent dataset for a training and testing task, the disagree class is even better classified by using automatic summaries. Additionally, when the approach is trained on the FNC-1 dataset and the test is carried out on the Emergent dataset, the discuss class surpasses the upper bound.

Regarding the performance of the Stance Stage in isolation, i.e., without considering the *Relatedness Stage*, the results present a slightly better performance than the whole approach with an increase of 2 percentage points (see Table 7). This was to be expected since errors derived from the *Relatedness Stage* are avoided. To conclude, these figures demonstrate that the approach, apart from potentially being a general solution, also demonstrates that using summarization of the body text as input is useful for the stance detection task, since the performance is very close to the upper bound proposed at Emergent.

5.3. *HeadlineStanceChecker* Validation

The results of the *HeadlineStanceChecker* are shown in Table 7. This approach integrates the Relatedness and

System	F_1 Score				F_1m	Acc.	Rel. Score
	Agree	Disagree	Discuss	Unrelated			
<i>Talos</i> [45]	53.90	3.54	76.00	99.40	58.21	89.08	82.02
<i>Athene</i> [46]	48.70	15.12	78.00	99.60	60.40	89.48	82.00
<i>UCLMR</i> [47]	47.94	11.44	74.70	98.90	58.30	88.46	81.72
<i>Human Upper Bound</i> [46]	58.80	66.70	76.50	99.70	75.40	–	85.90
<i>Dulhanty et al.</i> [49]	73.76	55.26	85.53	99.12	78.42	93.71	90.00
<i>Zhang et al.</i> [48]	67.47	81.30	83.90	99.73	83.10	93.77	89.30
<i>HeadlineStanceChecker-1stage</i>	70.34	53.42	85.30	99.41	77.12	93.64	89.80
<i>HeadlineStanceChecker-2stages</i>	72.34	62.53	87.32	99.38	80.39	94.31	91.02

Table 7: *HeadlineStanceChecker* results and comparison performance for the FNC-1 dataset.

953 Stance classifiers and only uses automatic summaries 985
954 for these two classifiers—but for the Relatedness clas- 986
955 sifier, the external features are included—. This table 987
956 contains the performance for the class-wise F_1 , macro- 988
957 average F_1m , accuracy (Acc.) and the relative score 989
958 (Rel. Score). Moreover, it also provides the results obtained 990
959 by competitive state-of-the-art systems together with 991
960 additional configurations that were also tested. 992

961 The 3 first rows are the top-3 best systems that partic- 993
962 ipated in the FNC-1 challenge. The results for each of 994
963 the evaluation metrics were calculated using the confu- 995
964 sion matrices and results were published [47] or made 996
965 available by the authors^{13,14}. 997

966 The fourth row corresponds to the *Human Upper* 998
967 *Bound*, and is the result of conducting the FNC-1 stance 999
968 detection task manually. This upper bound was defined 1000
969 by [46]. Five human annotators were asked to manually 1001
970 label 200 random instances, obtaining an overall inter- 1002
971 annotator agreement of Fleiss’ k of 0.686. Due to the 1003
972 fact that there is no upper bound reported in the FNC-1 1004
973 data, we also used these values as reference for compar- 1005
974 ison purposes. 1006

975 Next, the fifth and sixth rows include the results of 1007
976 recent approaches [48, 49] that also addressed the head- 1008
977 line stance detection task using the FNC-1 dataset, but 1009
978 did not take part in the challenge. Since there was no 1010
979 public code available, these results were also calculated 1011
980 from the confusion matrices provided in the papers. 1012

981 The seventh row indicates the results for our *Head-* 1013
982 *lineStanceChecker* approach but configured only with 1014
983 a single classifier. We have called this approach 1015
984 *HeadlineStanceChecker-1stage*. Finally, the last row 1016

1017
1018 ¹³https://github.com/hanselowski/athene_system/ (accessed online 1018
1019 15 February, 2021)

1020 ¹⁴<https://github.com/CiscoTalos/fnc1> (accessed online 15 February, 1020
1021 2021)

985 belongs to our *HeadlineStanceChecker* approach, us- 986
987 ing our proposed two-stage classification. For clar- 988
989 ity purposes, in the table we will refer to this ap- 990
991 proach as *HeadlineStanceChecker-2stages*. Regardless 992
993 of whether the classification is conducted in 1 or 2 994
995 stages, both approaches have used just the automatic 996
997 summaries created from the full body text during the 998
999 whole process. 1000

1001 As can be seen in Table 7, *HeadlineStanceChecker-* 1002
1003 *2stages* is competitive enough with respect to the other 1004
1005 systems, given that it only uses short summaries for 1006
1007 the classification process, and not the full body text 1008
1009 as the other systems use, so the information reduction 1009
1010 does not imply a high loss in the results obtained, be- 1010
1011 ing better than the FNC-1 participants, and the human 1011
1012 upper bound. Furthermore, the results also validate the 1012
1013 fact that the divide-and-conquer strategy applied for di- 1013
1014 viding the classification into two stages is beneficial 1014
1015 and yields better performance when using our proposed 1015
1016 model with a single classifier (rows 7th and 8th). This 1016
1017 is especially the case for detecting disagreement between 1017
1018 the headline and the news article. 1018

1019 Furthermore, the most remarkable improvement for 1019
1020 *HeadlineStanceChecker-2stages* is achieved in the *dis-* 1020
1021 *cuss* category, over performing all the remaining ap- 1021
1022 proaches. The F_1 improves by around 2 points com- 1022
1023 pared to the second-best approach, i.e., [49], and 13 1023
1024 points over the lowest-performance system [47] in this 1024
1025 category. By achieving competitive values in the other 1025
1026 classes as well, *HeadlineStanceChecker-2stages* obtains 1026
1027 a final macro-F1 value of 80.39%, being only beaten 1027
1028 by the system proposed in [48], which takes advantage 1028
1029 of a considerable number of external features beyond 1029
1030 similarity to enrich the neural model. It is worth high- 1030
1031 lighting that in terms of accuracy and relative score, 1031
1032 our approach (i.e., *HeadlineStanceChecker-2stages*) ob- 1032
1033 tains the best result among the automatic systems in 1033

	FNC-1 <i>subset</i> >512				FNC-1 <i>subset</i> <512			
	Agree	Disagree	Discuss	Unrelated	Agree	Disagree	Discuss	Unrelated
Train set	1,112	314	3,536	12,886	2,566	526	5,373	23,659
Test set	645	321	1,259	5,501	1,258	376	3,205	12,848
Total	1,757	635	4,795	18,387	3,824	902	8,578	36,507

Table 8: Class distribution for FNC-1 *subset*>512 and FNC-1 *subset*<512.

1022 both cases, achieving 94.31% and 91.02%, respectively. 1043
1023 Focusing on the results obtained by the participants in 1044
1024 the FNC-1 competition, when these results are analyzed 1045
1025 independently for each of the classes, it can be seen that 1046
1026 except for the classification of *unrelated* headlines — 1047
1027 whose results are close to 100% in F1 measure, and this 1048
1028 happens also for the remaining approaches as well— 1049
1029 for the remaining classes, the results are very limited. 1050
1030 The systems that participated in the FNC-1 competition 1051
1031 have a very reduced performance especially in detect- 1052
1032 ing the *disagree* stance, whereas the detection of *agree* 1053
1033 is around 50% in F1 measure and for *discuss* around 1054
1034 75% for the best approach. Outside the FNC-1 competi- 1055
1035 tion, the performance increases in all categories, being 1056
1036 the disagree category one of the most challenging 1057
1037 to classify, in which only the approach proposed in [48] 1058
1038 obtains surprisingly high results for this category com- 1059
1039 pared to the remaining methods.

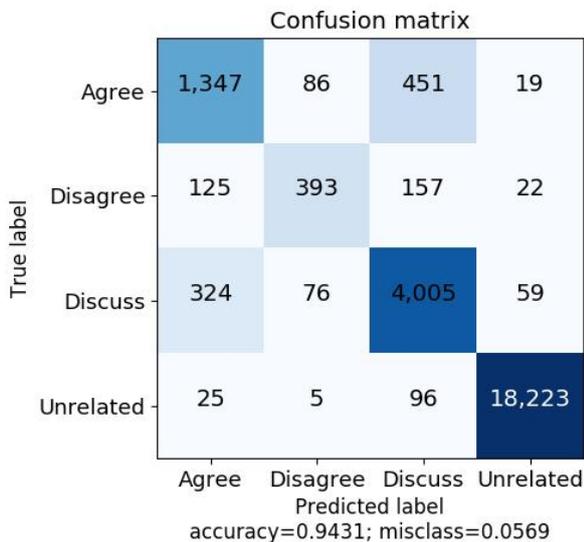


Figure 5: Confusion matrix resulting from the *HeadlineStanceChecker-2stages*.

1040 After having established that 1081
1041 *HeadlineStanceChecker-2stages* performs adequately 1082
1042 —correctly detecting 94.31% of the test set classes—, 1083

the confusion matrix presented in Figure 5 provides more detail on the actual performance of the system for each stance class. From this information, we can observe that per class, the major classification problems occur with *disagree* and *discuss* classes. The data reflects that 22.5% of the total number of *disagree* stances are being classified as *discuss*, whereas as 23.6% of the total number of *agree* stances are classified as *discuss*. However, only 7.2% of the total number of *discuss* stances are being classified as *agree* and 17.9% of the total number of *disagree* stances are being classified as *agree*.

5.4. Summary versus Body text Analysis

Finally, in order to test the convenience of using the summary or the body text as input to our whole system, a final analysis was performed by an experiment designed to allow us to compare the results in both cases. To determine how this would be accomplished, we considered the singularities of our system, since the use of RoBERTa implies certain constraints that affect the input processing.

RoBERTa, as a classification model, allows a maximum input length of 512 tokens, called *maximum sequence length*. Information that exceeds such a length is truncated. Our configuration takes as input both the headline and the text to which it refers, body text or summary, but it is relevant to remark that, in this case, the sequence length includes the tokens of the headline plus the tokens of that text. Since the headline must remain complete for the classification process, if it is necessary to truncate, it is the information in the body text which is lost.

In relation to the aforementioned issue, the previous work described in [49] focused on the analysis of the length of the body text for classification purposes, showing that for the examples in which the input sequence is greater than 512 tokens, the accuracy of the classification decreases considerably with respect to smaller sequences.

Against this backdrop, our hypothesis states that applying summarization to the text before classification

Input	F_1 Score				F_1m
	Agree	Disagree	Discuss	Unrelated	
<i>News body</i>	54.45	12.69	78.97	99.52	61.40
<i>News summary</i>	59.61	28.06	80.85	99.32	66.96

Table 9: *HeadlineStanceChecker-2stages* results for *subset>512* with different inputs: news body and news summary.

Input	F_1 Score				F_1m
	Agree	Disagree	Discuss	Unrelated	
<i>News body</i>	78.64	69.38	89.81	99.59	84.35
<i>News summary</i>	74.17	58.91	87.69	99.36	80.03

Table 10: *HeadlineStanceChecker-2stages* results for *subset<512* with different inputs: news body text and news summary.

1084 implies an improvement in the results. In order to prove
1085 it, we first create two subsets from the FNC-1 dataset
1086 according to the news story length: *subset>512* and
1087 *subset<512*. Table 8 shows the class distribution for
1088 both subsets.

1089 Next, we trained and tested the system with both sub-
1090 sets twice: first with the bodies as input, and then with
1091 the summaries. The results in Table 9 show that for long
1092 news stories(*subset>512*), the system performs better
1093 with summaries as input than truncating the text of the
1094 full article. This could happen because reducing the in-
1095 put by simply cutting text at the end of the document
1096 results in relevant information being lost, whereas when
1097 making a summary, it is the relevant information that
1098 prevails in a more concise mode.

1099 Similarly, results for *subset<512*, the shortest news
1100 stories, are reported in Table 10. The system was again
1101 trained and tested, taking the body and the summary as
1102 inputs. In this case, results are better when using the full
1103 body text, which could indicate that all the information
1104 needed for a proper classification is present by consider-
1105 ing the whole text—an unfeasible scenario with longer
1106 texts—.

1107 There exist no explicit rules that determine what the
1108 length of a news article should be, but there is instead
1109 certain evidence supporting that news tend to be longer
1110 than 512 tokens. In Table 11 we have gathered statistics
1111 from the most popular news datasets that are being used
1112 in language processing tasks. All together, they contain
1113 more than 2 million articles from different sources, with
1114 an average length superior to 512 tokens. The relevance
1115 of our approach is made clear by these figures, which
1116 indicate that, in most cases, using news summarization
1117 would be the right strategy.

	Examples	Average length
CNN [103]	92 K	760.50
DailyMail [103]	310 K	653.33
NY times [104]	650 K	800.04
Newsroom [105]	1,210 K	770.09
Total	2,260 K	745.83

Table 11: Statistics from large news corpora indicating the average document length in words.

1118 6. Conclusions and Future Work

1119 *HeadlineStanceChecker* has been demonstrated to be
1120 an effective approach for detecting misinformation in
1121 news, specifically when a headline has to be compared
1122 to its body text. The novelty of our approach rests
1123 on two key premises: i) the adoption of a divide-and-
1124 conquer strategy, thus tackling the stance classification
1125 problem by means of a 2-stage neural architecture; and
1126 ii) the use of extractive semantic summarization instead
1127 of the full news body text for the whole classification,
1128 in addition to a salience and two similarity features that
1129 will help to determine the relatedness of the headline
1130 with respect to the news article.

1131 To show the appropriateness of *HeadlineS-*
1132 *tanceChecker*, different experiments were carried
1133 out in the context of an existing task —Fake News
1134 Challenge FNC-1—, where the stance of a headline had
1135 to be classified into one of the following classes: *unre-*
1136 *lated*, *agree*, *disagree*, and *discuss*. The experiments
1137 involved validating each of the proposed classification
1138 stages in isolation together with the whole approach,
1139 as well as a comparison with respect the state of the
1140 art in this task. Furthermore, additional experiments

with another corpus for headline stance detection (i.e., Emergent dataset) were also performed to verify the generalization of our approach. The results obtained by our system were very competitive compared to other SOTA systems obtaining 94.31% Accuracy, as well as the highest result in FNC-1 relative score compared with the state of the art (91.02%).

The unbalanced nature of the FNC-1 dataset leads to existing systems being more capable of learning how to detect *unrelated* headlines, but are less accurate when it comes to the remaining classes. Even so, the results obtained by *HeadlineStanceChecker* for the different categories with less examples, *agree*, *disagree* and *discuss* are fair enough and promising, which indicates that the chosen approach is appropriate for the task.

Future work will aim to improve the results of *agree* and *disagree* classification by extending the system to take into consideration Sentiment Analysis features. Furthermore, as reported speech is recently being introduced to determine bias and document stance, it could be very useful for determining the stance of headlines and news articles. Some reporting events are neutral, for example, by using *reported* or *said*, whereas some others introduce a stance, for instance, *'deny'* implies disagreement or *'confirm'* indicates agreement.

Besides, another interesting aspect to focus on would be to investigate the relation of the stance detection classes (agree, disagree, discuss and unrelated) with the “incongruent” and “congruent” classification to determine whether this relation can provide some insights for different scenarios.

Finally, as a future goal that contributes to investigating the problem of fake news detection, we expect to apply *HeadlineStanceChecker* to a real world scenario to detect when headlines introduce mis- or disinformation to readers. Our contribution to improving the current research in the field, by means of new learning strategies and discourse aware techniques, will help to combat online fake news, a societal problem that requires concerted action.

Acknowledgements

References

[1] V. L. Rubin, Disinformation and misinformation triangle, *Journal of Documentation* 75 (2019) 1013–1034.
[2] M. Tudjmanand, N. Mikelic Preradovic, Information science: Science about information, in: *Proceedings of Informing Science & IT Education*, 2003, pp. 1513–1527.
[3] S. Tavernisen, As fake news spreads lies, more readers shrug at the truth, *New York Times* (2019).

[4] A. Bovet, H. A. Makse, Influence of fake news in Twitter during the 2016 US presidential election, *Nature Communications* 10(1):7 (2019).
[5] M. T. Bastos, D. Mercea, The Brexit botnet and user-generated hyperpartisan news, *Social Science Computer Review* 37 (2019) 38–54.
[6] V. Hooper, Fake news and social media: The role of the receiver, in: *5th European Conference on Social Media 2018, Academic Conferences and publishing limited*, 2018, p. 62.
[7] S. Issenberg, *The victory lab: The secret science of winning campaigns*, Crown, 2012.
[8] E. S. Boró, D. Tomás, P. Moreda, P. Martínez-Barco, M. Palomar, Fighting post-truth using natural language processing: A review and open challenges, *Expert Systems with Applications* 141 (2020).
[9] T. van Dijk, *News as discourse*, Communication Series, L. Erlbaum Associates, 1988.
[10] J. Kuiken, A. Schuth, M. Spitters, M. Marx, Effective headlines of newspaper articles in a digital environment, *Digital Journalism* 5 (2017) 1300–1314.
[11] M. Gabielkov, A. Ramachandran, A. Chaintreau, A. Legout, Social clicks: What and who gets read on twitter?, *ACM SIGMETRICS Performance Evaluation Review* 44 (2016) 179–192.
[12] B. Lutz, M. T. P. Adam, S. Feuerriegel, N. Pröllochs, D. Neumann, Affective information processing of fake news: Evidence from neurois, in: *Information Systems and Neuroscience*, Springer International Publishing, 2020, pp. 121–128.
[13] Y. Chen, N. J. Conroy, V. L. Rubin, News in an online world: The need for an “automatic crap detector”, in: *Proceedings of the 78th ASIS&T Annual Meeting: Information Science with Impact: Research in and for the Community*, American Society for Information Science, 2015.
[14] W. Wei, X. Wan, Learning to identify ambiguous and misleading news headlines, in: *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, AAAI Press, 2017, pp. 4172–4178.
[15] Y. Chen, N. J. Conroy, V. L. Rubin, Misleading online content: Recognizing clickbait as “false news”, in: *Proceedings of the ACM on Workshop on Multimodal Deception Detection*, Association for Computational Linguistics, 2015, pp. 15–19.
[16] S. Chesney, M. Liakata, M. Poesio, M. Purver, Incongruent headlines: Yet another way to mislead your readers, in: *Proceedings of Natural Language Processing meets Journalism*, 2017, pp. 56–61.
[17] W. Ferreira, A. Vlachos, Emergent: a novel data-set for stance classification, in: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, Association for Computational Linguistics, 2016, pp. 1163–1168.
[18] Y. Hayashi, H. Yanagimoto, Headline generation with recurrent neural network, in: *New Trends in E-service and Smart Computing*, Springer, 2018, pp. 81–96.
[19] Z. Huang, Z. Ye, S. Li, R. Pan, Length adaptive recurrent model for text classification, in: *Proceedings of the ACM on Conference on Information and Knowledge Management*, Association for Computing Machinery, 2017, pp. 1019–1027.
[20] M. Choraś, K. Demestichas, A. Gielczyk, A. Herrero, P. Ksieniewicz, K. Remoundou, D. Urda, M. Woźniak, Advanced machine learning techniques for fake news (online disinformation) detection: A systematic mapping study, *Applied Soft Computing* (2020) 107050.
[21] G. Di Domenico, J. Sit, A. Ishizaka, D. Nunan, Fake news, social media and marketing: A systematic review, *Journal of Business Research* 124 (2021) 329–341.

- [22] G. Zarrella, A. Marsh, Mitre at semeval-2016 task 6: Transfer learning for stance detection, in: Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016), Association for Computational Linguistics, 2016, pp. 458–463.
- [23] S. Ghosh, P. Singhanian, S. Singh, K. Rudra, S. Ghosh, Stance detection in web and social media: a comparative study, in: International Conference of the Cross-Language Evaluation Forum for European Languages, Springer, 2019, pp. 75–87.
- [24] S. Somasundaran, J. Wiebe, Recognizing stances in ideological on-line debates, in: Proceedings of the NAACL HLT 2010 workshop on computational approaches to analysis and generation of emotion in text, 2010, pp. 116–124.
- [25] A. Konjengbam, S. Ghosh, N. Kumar, M. Singh, Debate stance classification using word embeddings, in: International conference on big data analytics and knowledge discovery, Springer, 2018, pp. 382–395.
- [26] A. Faulkner, Automated classification of stance in student essays: An approach using stance target information and the wikipedia link-based measure, in: The Twenty-Seventh International Flairs Conference, 2014.
- [27] C. Li, A. Porco, D. Goldwasser, Structured representation learning for online debate stance prediction, in: Proceedings of the 27th International Conference on Computational Linguistics, 2018, pp. 3728–3739.
- [28] R. Agrawal, S. Rajagopalan, R. Srikant, Y. Xu, Mining news-groups using networks arising from social behavior, in: Proceedings of the 12th international conference on World Wide Web, 2003, pp. 529–535.
- [29] A. Murakami, R. Raymond, Support or oppose? classifying positions in online debates from reply activities and opinion expressions, in: Proceedings of the 23rd International Conference on Computational Linguistics: Posters, 2010, pp. 869–875.
- [30] G. Gorrell, E. Kochkina, M. Liakata, A. Aker, A. Zubiaga, K. Bontcheva, L. Derczynski, Semeval-2019 task 7: Rumoureal, determining rumour veracity and support for rumours, in: Proceedings of the 13th International Workshop on Semantic Evaluation, 2019, pp. 845–854.
- [31] S. Mohammad, S. Kiritchenko, P. Sobhani, X. Zhu, C. Cherry, SemEval-2016 task 6: Detecting stance in tweets, in: Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016), Association for Computational Linguistics, San Diego, California, 2016, pp. 31–41.
- [32] A. I. Al-Ghadir, A. M. Azmi, A. Hussain, A novel approach to stance detection in social media tweets by fusing ranked lists and sentiments, *Information Fusion* 67 (2021) 29–40.
- [33] B. G. Patra, D. Das, S. Bandyopadhyay, *Ju_nlp* at semeval-2016 task 6: detecting stance in tweets using support vector machines, in: Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016), 2016, pp. 440–444.
- [34] H. Elfardy, M. Diab, *Cu-gwu* perspective at semeval-2016 task 6: Ideological stance detection in informal text, in: Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016), 2016, pp. 434–439.
- [35] I. Augenstein, T. Rocktäschel, A. Vlachos, K. Bontcheva, Stance detection with bidirectional conditional encoding, *arXiv preprint arXiv:1606.05464* (2016).
- [36] P. Wei, W. Mao, D. Zeng, A target-guided neural memory model for stance detection in twitter, in: International Joint Conference on Neural Networks (IJCNN), IEEE, 2018, pp. 1–8.
- [37] S. Zhou, J. Lin, L. Tan, X. Liu, Condensed convolution neural network by attention over self-attention for stance detection in twitter, in: International Joint Conference on Neural Networks (IJCNN), IEEE, 2019, pp. 1–8.
- [38] A. Sen, M. Sinha, S. Mannarswamy, S. Roy, Stance classification of multi-perspective consumer health information, in: Proceedings of the ACM India Joint International Conference on Data Science and Management of Data, 2018, pp. 273–281.
- [39] M. Taulé, M. A. Martí, F. M. Rangel, P. Rosso, C. Bosco, V. Patti, et al., Overview of the task on stance and gender detection in tweets on Catalan independence at IberEval 2017, in: 2nd Workshop on Evaluation of Human Language Technologies for Iberian Languages, volume 1881, CEUR-WS, 2017, pp. 157–177.
- [40] S. V. Vychezhnanin, E. V. Kotelnikov, Stance detection based on ensembles of classifiers, *Programming and Computer Software* 45 (2019) 228–240.
- [41] M. Lai, A. T. Cignarella, D. I. H. Farias, C. Bosco, V. Patti, P. Rosso, Multilingual stance detection in social media political debates, *Computer Speech & Language* 63 (2020) 101075.
- [42] E. Zotova, R. Agerri, G. Rigau, Semi-automatic generation of multilingual datasets for stance detection in twitter, *Expert Systems with Applications* 170 (2021) 114547.
- [43] M. Babakar, N. Bakos, H. Daumé, A. Mantzarlis, D. Seddah, A. Vlachos, C. Wardle, Fake News Challenge - I, 2016.
- [44] C. Silverman, Lies, damn lies and viral content, 2019.
- [45] S. Baird, D. Sibley, Y. Pan, Talos targets disinformation with fake news challenge victory, <https://blog.talosintelligence.com/2017/06/talos-fake-news-challenge.html>, lastaccessedon20/02/21, 2017.
- [46] B. S. Andreas Hanselowski, Avinesh PVS, F. Caspelherr, Description of the system developed by team athene in the FNC-1, 2017.
- [47] B. Riedel, I. Augenstein, G. P. Spithourakis, S. Riedel, A simple but tough-to-beat baseline for the Fake News Challenge stance detection task, *Computing Research Repository, CoRR abs/1707.03264* (2017).
- [48] Q. Zhang, S. Liang, A. Lipani, Z. Ren, E. Yilmaz, From stances’ imbalance to their hierarchical representation and detection, in: The World Wide Web Conference, ACM, 2019, pp. 2323–2332.
- [49] C. Dulhanty, J. L. Deglint, I. B. Daya, A. Wong, Taking a stance on fake news: Towards automatic disinformation assessment via deep bidirectional transformer language models for stance detection, *arXiv preprint arXiv:1911.11951* (2019).
- [50] B. Pouliquen, R. Steinberger, C. Best, Automatic detection of quotations in multilingual news, 2007, pp. 487–492.
- [51] A. Vlachos, S. Riedel, Identification and verification of simple claims about statistical properties, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2015, pp. 2596–2601.
- [52] M.-C. De Marneffe, A. N. Rafferty, C. D. Manning, Finding contradictions in text, *Proceedings of Association for Computational Linguistics* (2008) 1039–1047.
- [53] S. Harabagiu, A. Hickl, F. Lacatusu, Negation, contrast and contradiction in text processing, in: *AAAI*, volume 6, 2006, pp. 755–762.
- [54] O. Levy, T. Zesch, I. Dagan, I. Gurevych, Recognizing partial textual entailment, in: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, volume 2, 2013, pp. 451–455.
- [55] S. A. Brown, The effects of explicit main idea and summarization instruction on reading comprehension of expository text for alternative high school students, *PhD Thesis*. Utah State University (2018).
- [56] J. Engelen, G. Camp, J. van de Pol, A. de Bruin, Teachers’ monitoring of students’ text comprehension: can students’ keywords and summaries improve teachers’ judgment accuracy?, *Metacognition and learning* 13 (2018) 287–307.

- [57] X. G. Lin, S.-E. Jhang, D. Dong, Investigating the effects of text summarization on linguistic quality of argumentative writing, volume 60, 2018, pp. 245–268.
- [58] J. P. Barreiro, Improving reading comprehension of narrative texts through summaries, PhD Thesis. Universidad Casa Grande (2019).
- [59] R. Dijkman, A. Wilbik, Linguistic summarization of event logs – a practical approach, *Information Systems* 67 (2017) 114 – 125.
- [60] J. Petkovic, V. Welch, M. Jacob, M. Yoganathan, A. P. Ayala, H. Cunningham, P. Tugwell, The effectiveness of evidence summaries on health policymakers and health system managers use of evidence from systematic reviews: A systematic review, *Implementation Science* 11 (2016).
- [61] L. Hartling, A. Gates, J. Pillay, M. Nuspl, A. Newton, Development and usability testing of epc evidence review dissemination summaries for health systems decisionmakers., *Methods Research Report. Technical Report.* (2018).
- [62] Y. Liu, X. Song, S.-F. Chen, Long story short: finding health advice with informative summaries on health social media, *Aslib Journal of Information Management ahead-of-print* (2019).
- [63] F. Deroncourt, M. Ghassemi, W. Chang, A repository of corpora for summarization, in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, European Language Resources Association (ELRA), Miyazaki, Japan, 2018.
- [64] A. Nenkova, Automatic text summarization of newswire: Lessons learned from the document understanding conference, in: *Proceedings of the 20th National Conference on Artificial Intelligence - Volume 3, AAAI'05*, AAAI Press, 2005, p. 1436–1441.
- [65] S. Mackie, R. McCreadie, C. Macdonald, I. Ounis, Experiments in newswire summarisation, in: N. Ferro, F. Crestani, M.-F. Moens, J. Mothe, F. Silvestri, G. M. Di Nunzio, C. Hauff, G. Silvello (Eds.), *Advances in Information Retrieval*, Springer International Publishing, Cham, 2016, pp. 421–435.
- [66] Y. Duan, A. Jatowt, Cross-time comparative summarization of news articles, in: *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM*, Association for Computing Machinery, 2019, p. 735–743.
- [67] C. Zhu, Z. Yang, R. Gmyr, M. Zeng, X. Huang, Make lead bias in your favor: A simple and effective method for news summarization, *arXiv preprint arXiv:1912.11602* (2019).
- [68] M. Banko, V. O. Mittal, M. J. Witbrock, Headline generation based on statistical translation, in: *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, Association for Computational Linguistics, 2000, pp. 318–325.
- [69] B. Dorr, D. Zajic, R. Schwartz, Hedge trimmer: A parse-and-trim approach to headline generation, in: *Proceedings of the North American Chapter of the Association for Computational Linguistics, Text Summarization Workshop*, 2003, pp. 1–8.
- [70] D. Zajic, B. Dorr, R. Schwartz, Automatic headline generation for newspaper stories, in: *Proceedings of the Workshop on Automatic Summarization*, 2002, pp. 78–85.
- [71] J. Tan, X. Wan, J. Xiao, From neural sentence summarization to headline generation: A coarse-to-fine approach, in: *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, AAAI Press, 2017, pp. 4109–4115.
- [72] D. Gavrilov, P. Kalaidin, V. Malykh, Self-attentive model for headline generation, in: *Advances in Information Retrieval*, Springer International Publishing, 2019, pp. 87–93.
- [73] K. Iwama, Y. Kano, Multiple news headlines generation using page metadata, in: *Proceedings of the 12th International Conference on Natural Language Generation*, Association for Computational Linguistics, 2019, pp. 101–105.
- [74] S. Esmailzadeh, G. X. Peh, A. Xu, Neural abstractive text summarization and fake news detection, *Computing Research Repository, CoRR abs/1904.00788* (2019).
- [75] S. Esmailzadeh, G. X. Peh, A. Xu, Neural abstractive text summarization and fake news detection, *arXiv preprint arXiv:1904.00788* (2019).
- [76] G. Kim, Y. Ko, Graph-based fake news detection using a summarization technique, in: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 2021, pp. 3276–3280.
- [77] P. Krejzl, B. Hourová, J. Steinberger, Stance detection in online discussions, *Computing Research Repository, CoRR abs/1701.00504* (2017).
- [78] P. Krejzl, Stance detection and summarization in social networks, *Report* (2018).
- [79] S. Chesney, M. Liakata, M. Poesio, M. Purver, Incongruent headlines: Yet another way to mislead your readers, in: *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing meets Journalism*, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 56–61. URL: <https://www.aclweb.org/anthology/W17-4210>. doi:10.18653/v1/W17-4210.
- [80] D. Küçük, F. Can, Stance detection: A survey, *ACM Computing Surveys (CSUR)* 53 (2020) 1–37.
- [81] M. Hardalov, A. Arora, P. Nakov, I. Augenstein, A survey on stance detection for mis-and disinformation identification, *arXiv preprint arXiv:2103.00242* (2021).
- [82] W. Ferreira, A. Vlachos, 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*, Association for Computational Linguistics, San Diego, California, 2016, pp. 1163–1168. URL: <https://www.aclweb.org/anthology/N16-1138>. doi:10.18653/v1/N16-1138.
- [83] Y. Lv, C. Zhai, Positional language models for information retrieval, in: *Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2009, pp. 299–306.
- [84] M. Vicente, C. Barros, E. Lloret, Statistical language modelling for automatic story generation, *Journal of Intelligent & Fuzzy Systems* 34 (2018) 3069–3079.
- [85] A. Kilgariff, C. Fellbaum, *WordNet: An Electronic Lexical Database*, *Language* 76 (2000) 706.
- [86] M. E. Vicente, E. Lloret, A discourse-informed approach for cost-effective extractive summarization, in: L. E. Anke, C. Martín-Vide, I. Spasic (Eds.), *Statistical Language and Speech Processing - 8th International Conference, Proceedings*, volume 12379 of *Lecture Notes in Computer Science*, Springer, 2020, pp. 109–121.
- [87] L. Padró, E. Stanilovsky, Freeling 3.0: Towards wider multilinguality, in: *Proceedings of the Language Resources and Evaluation Conference (LREC 2012)*, ELRA, Istanbul, Turkey, 2012.
- [88] N. Passalis, A. Tefas, Learning bag-of-embedded-words representations for textual information retrieval, *Pattern Recognition* 81 (2018) 254–267.
- [89] B. Li, L. Han, Distance weighted cosine similarity measure for text classification, in: *Proceedings of the 14th International Conference on Intelligent Data Engineering and Automated Learning*, Springer-Verlag, 2013, pp. 611–618.
- [90] S. Tata, J. M. Patel, Estimating the selectivity of tf-idf based cosine similarity predicates, *SIGMOD Record* 36 (2007) 75–80.
- [91] V. Kotu, B. Deshpande, Classification, in: *Data Science*, Else-

- 1515 vier, 2019, pp. 65–163.
- 1516 [92] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller,
1517 O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler,
1518 R. Layton, J. VanderPlas, A. Joly, B. Holt, G. Varoquaux, API
1519 design for machine learning software: experiences from the
1520 scikit-learn project, in: ECML PKDD Workshop: Languages
1521 for Data Mining and Machine Learning, 2013, pp. 108–122.
- 1522 [93] F. Šarić, G. Glavaš, M. Karan, J. Šnajder, B. Dalbelo Bašić,
1523 TakeLab: Systems for measuring semantic text similarity, in:
1524 Proceedings of the First Joint Conference on Lexical and Com-
1525 putational Semantics, Association for Computational Linguis-
1526 tics, 2012, pp. 441–448.
- 1527 [94] L. Metcalf, W. Casey, Metrics, similarity, and sets, in: Cyber-
1528 security and Applied Mathematics, Elsevier, 2016, pp. 3–22.
- 1529 [95] M. Vijaymeena, K. Kavitha, A survey on similarity measures
1530 in text mining, *Machine Learning and Applications: An Inter-
1531 national Journal* 3 (2016) 19–28.
- 1532 [96] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy,
1533 M. Lewis, L. Zettlemoyer, V. Stoyanov, RoBERTa: A ro-
1534 bustly optimized BERT pretraining approach, *arXiv preprint
1535 arXiv:1907.11692* (2019).
- 1536 [97] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, S. Bowman,
1537 GLUE: A multi-task benchmark and analysis platform for nat-
1538 ural language understanding, in: Proceedings of the 2018
1539 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting
1540 Neural Networks for NLP, Association for Computational Lin-
1541 guistics, 2018, pp. 353–355.
- 1542 [98] G. Lai, Q. Xie, H. Liu, Y. Yang, E. Hovy, RACE: Large-scale
1543 ReAding comprehension dataset from examinations, in: Pro-
1544 ceedings of the 2017 Conference on Empirical Methods in Nat-
1545 ural Language Processing, Association for Computational Lin-
1546 guistics, Copenhagen, Denmark, 2017, pp. 785–794.
- 1547 [99] V. Slovikovskaya, Transfer learning from transformers to fake
1548 news challenge stance detection (FNC-1) task, *arXiv preprint
1549 arXiv:1910.14353* (2019).
- 1550 [100] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-
1551 training of deep bidirectional transformers for language under-
1552 standing, *arXiv preprint arXiv:1810.04805* (2018).
- 1553 [101] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi,
1554 P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Brew, Hugging-
1555 face’s transformers: State-of-the-art natural language process-
1556 ing, *ArXiv abs/1910.03771* (2019).
- 1557 [102] A. Hanselowski, A. PVS, B. Schiller, F. Caspelherr, D. Chaud-
1558 huri, C. M. Meyer, I. Gurevych, A retrospective analysis of
1559 the fake news challenge stance-detection task, in: Proceedings
1560 of the 27th International Conference on Computational Lin-
1561 guistics, Association for Computational Linguistics, 2018, pp.
1562 1859–1874.
- 1563 [103] K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt,
1564 W. Kay, M. Suleyman, P. Blunsom, Teaching machines to read
1565 and comprehend, in: *Advances in neural information process-
1566 ing systems*, 2015, pp. 1693–1701.
- 1567 [104] E. Sandhaus, The New York Times Annotated Corpus
1568 ldc2008t19, Linguistic Data Consortium, Philadelphia, 2008.
- 1569 [105] M. Grusky, M. Naaman, Y. Artzi, Newsroom: A dataset of
1570 1.3 million summaries with diverse extractive strategies, *arXiv
1571 preprint arXiv:1804.11283* (2018).