

Journal Pre-proof

Exploiting discourse structure of traditional digital media to enhance automatic fake news detection

Alba Bonet-Jover, Alejandro Piad-Morffis, Estela Saquete,
Patricio Martínez-Barco, Miguel Ángel García-Cumbreras



PII: S0957-4174(20)31027-7
DOI: <https://doi.org/10.1016/j.eswa.2020.114340>
Reference: ESWA 114340

To appear in: *Expert Systems With Applications*

Received date : 5 June 2020
Revised date : 22 October 2020
Accepted date : 16 November 2020

Please cite this article as: A. Bonet-Jover, A. Piad-Morffis, E. Saquete et al., Exploiting discourse structure of traditional digital media to enhance automatic fake news detection. *Expert Systems With Applications* (2020), doi: <https://doi.org/10.1016/j.eswa.2020.114340>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier Ltd.

Exploiting Discourse Structure of Traditional Digital Media to enhance Automatic Fake News Detection

Alba Bonet-Jover^a, Alejandro Piad-Morffis^b, Estela Saquete^a, Patricio Martínez-Barco^a, Miguel Ángel García-Cumbreras^c

^a*Department of Software and Computing Systems, University of Alicante, Spain*
{alba.bonet,stela,patricio}@dlsi.ua.es

^b*School of Math and Computer Science, University of Havana, Cuba*
apiad@matcom.uh.cu

^c*CEATIC / Universidad de Jaén*
magc@ujaen.es

Abstract

This paper presents a novel architecture for dealing with Automatic Fake News detection. The architecture factors in the discourse structure of news in traditional digital media and is based on two premises. First, fake news tends to mix true and false information with the purpose of confusing readers. Second, this research is focused on fake news delivered in traditional digital media, so our approach considers the influence of the journalistic structure of news, and the way journalists tend to introduce the essential content in a news story –using 5W1H answers–. Considering both premises, this proposal deals with the news components separately because some may be true or false, instead of considering the veracity value of the news article as a unit. A two-layer architecture is proposed, Structure and Veracity layers. To demonstrate the validity of the proposal, a new dataset was created and annotated with a new fine-grained annotation scheme (FNDeepML) that considers the different elements of the news document and their veracity. Due to the severity of the COVID-19 pandemic crisis, health is the chosen domain, and Spanish is the language used to validate the architecture, given the lack of research in this language. However, the proposal can be applied to any other language or domain. The performance of the Veracity layer of our proposal, which factors in the traditional news article structure and the 5W1H annotation, is capable of delivering a result of $F_1=0.807$. This represents a strong improvement when compared to the baseline, which uses the whole document with a single veracity value, obtaining $F_1=0.605$. These findings validate the suitability and effectiveness of our approach.

Keywords: Natural Language Processing, Fake News, Automated Fact-checking, Deep Learning, Machine Learning, Human Language Technologies

1. Introduction

In the digital era, information is mostly received and accessed online and the quality of this information becomes a crucial issue. However, there is a huge world-wide problem regarding the dissemination of fake news whose aim is to create confusion and manipulate public opinions and behaviours. Fake news are structured and written in a way that makes it difficult to distinguish between what is true or false. Fake information is diffused significantly farther, faster, deeper, and more broadly than the truth in all categories of information (Vosoughi et al., 2018).

This situation is exacerbated in times of emergency such as during the 2020 global pandemic caused by COVID-19. There are several reasons that have made coronavirus hoaxes a potentially serious problem. Information on COVID-19 was scarce during the early stages of the crisis, which increased the problem of misinformation. Besides, people around the world were in lockdowns, hyperconnected and anxious, which led to exponential viralization compared to a normal situation. Finally, many hoaxes related to prevention or cures were released, albeit with the intention of protecting, but these remedies spreading unchecked can be highly damaging. One example of false information widely disseminated was the claim *“Russia released more than 500 lions to make sure that people stay inside during the COVID-19 pandemic”*. The aim was to create alarm but it was demonstrated to be false¹.

In many cases, this disinformation is delivered by digital media web pages, which present news articles following the traditional format of a news piece, but sometimes “fake” information is provided, confusing readers and, in the case of fake news related to health, putting at serious risk the well-being of these people who may follow the advice given. Detecting and tackling fake news quickly and efficiently is, therefore, crucial because once false information spreads and permeates throughout society, it becomes difficult to refute. The number of hoaxes is reaching levels that would benefit from applying automatic techniques that enable the detection of fake news before they are massively spread. This is why Artificial Intelligence and Natural Language Processing (NLP) techniques are applied, so that the process can be automated.

A common phenomenon in the context of fake news is that false infor-

¹<https://www.snopes.com/fact-check/russia-release-lions-coronavirus/>

36 mation is provided mixed with true information, to create confusion in the
37 reader, and this premise is the basis of our proposal. An example is the claim
38 “*U.S. President Donald Trump will benefit financially if hydroxychloroquine*
39 *becomes an established treatment for COVID-19*”, which was fact-checked as
40 mostly false. Furthermore, by studying the journalistic structure of news and
41 how journalists introduce the essential content in news stories, our proposal
42 considers the information as separated items, where some are true and some
43 are false, instead of considering the news article as a whole when giving it a
44 veracity value. This research proposal aims to help automatic learning sys-
45 tems to determine which parts of the structure of a news piece, or which type
46 of content is more influential in reaching a decision about the veracity of the
47 news (Conroy et al., 2015) (Pérez-Rosas et al., 2018). From hereafter, the
48 term veracity refers to the accuracy and the truthfulness of the information
49 provided in a traditional digital news document (Ciampaglia et al., 2015)
50 (Das Bhattacharjee et al., 2017) (Lewandowsky et al., 2012) (Nyhan et al.,
51 2012).

52 Considering the present context, the main contributions of this research
53 are the following:

- 54 • Firstly, the proposal of a novel architecture for automatic fake news
55 detection on traditional digital newspaper articles that can determine
56 not only the full document veracity but most importantly, the veracity
57 of the essential content elements of the news. The architecture will
58 demonstrate that it is possible to determine the veracity of the news
59 more accurately by taking advantage of the discourse structure of the
60 news, that is, the journalistic structure and the essential content of
61 the news piece, thereby reducing the noise when training automatic
62 learning systems.
- 63 • Secondly, due to the lack of resources where information is annotated as
64 independent parts, another important objective of this research is the
65 creation of a dataset using a fine-grained annotation scheme, named
66 FNDeepML. This annotation scheme is especially focused on differen-
67 tiating the structural elements and essential content of classic news ar-
68 ticles, which should respond to the 5W1H (What, When, Who, Where,
69 Why and How) questions. This approach is especially innovative be-
70 cause existing datasets tag the news as a whole, in a single veracity
71 category. The language chosen for the dataset is Spanish, because de-

72 spite being the third most spoken language in the world², there are
73 very few Spanish language resources for this task at the present time,
74 making it beneficial for the research community. Due to the alarming
75 pandemic situation, the health domain is used as a benchmark, but the
76 proposal is readily adaptable to any language and domain.

77 The rest of the paper is organized as follows: Section 2 describes the struc-
78 ture of newspaper articles and their main content as well as the background
79 of automatic fake news detection regarding NLP; Section 3 presents the def-
80 inition of a new annotation scheme and the dataset created following this
81 scheme; Section 4 shows the architecture of the automatic system proposed;
82 Section 5 describes the evaluation environment used in this research; Section
83 6 shows the evaluation results and discusses them; Section 7 presents a set
84 of experiments to compare our proposal with the state of the art (SOTA);
85 and finally, our conclusions and future work are presented in Section 8.

86 2. Background

87 The development of automatic systems for fake news detection in the con-
88 text of this proposal requires the analysis of the main features of newspaper
89 articles, such as how they are structured and how the content is presented.
90 It is important to focus on everything that can serve as a differentiating ele-
91 ment between true news and fake news. Furthermore, a revision of the most
92 relevant literature regarding computational mechanisms for automatic fake
93 news detection is presented.

94 2.1. News structure and the 5W1H method

95 News is usually presented within a specific structure to attract readers and
96 provide information in an interesting and organised way. Although there are
97 different ways of writing a news story, there are two key principles on which
98 all well-built news should be based: neutrality and the inverted pyramid
99 structure (Thomson et al., 2008). Thus, the objectivity of a news piece may
100 depend on these two factors, so the detection of unusual deviations from
101 these accepted journalistic norms could provide a clue to detect fake news.

102 In the inverted pyramid hypothesis, “certain parts of news articles carry
103 different levels of useful information” (Khan et al., 2018; Norambuena et al.,

²<https://www.cervantes.es/imagenes/File/espanol.lengua.viva.2019.pdf>

104 2020), placing the most important information first and ending with the
105 least relevant information (Zhang & Liu, 2016). The three common and
106 most important parts of the news structure are the headline, the lead and
107 the body. Other important but secondary elements of news are the subtitle or
108 the conclusion that usually appear in news articles, but they are not always
109 present (see Figure 1).

110 In a well-built article, the parts must appear the following order:

- 111 • *Headline*: This element is the title of the news article and it provides
112 the main idea of the story. Normally it summarizes, in one sentence,
113 the basic and essential information about the story. Its main objective
114 is to attract the reader's attention.
- 115 • *Subtitle*: A second title that explains the headline in a little more
116 detail. It completes the information, but it also presents the idea in a
117 very summarized way. Sometimes, it completes the information given
118 in the headline, and at times it provides other details not mentioned
119 before. Its function is to hold the reader's attention and to encourage
120 him/her to keep reading the news article.
- 121 • *Lead*: The paragraph(s) that develops the main information by fol-
122 lowing the 5W1H method and "presents the point or newsworthy el-
123 element(s) of the story and simultaneously works as a beginning of the
124 story" (Bednarek & Caple, 2012). All the main information of the news
125 article must be clearly presented in this section by answering the six
126 questions used in journalism: what, who, where, when, why and how.
127 The lead and the headline are sometimes considered as a unit because
128 the lead usually repeats the idea given by the headline, but in more
129 detail and accuracy (Thomson et al., 2008).
- 130 • *Body*: All the developed information is in this part of the news article.
131 The body presents all the background, facts, elements and reasons of
132 the story in detail. As mentioned by (Thomson et al., 2008), "the body
133 of the text does not develop new meanings but, rather, acts to refer
134 back to the headline/lead through a series of specifications." All the
135 six questions answered in the lead will be developed in the body by
136 explaining all the elements involved.
- 137 • *Conclusion/Tail*: The main idea of the story can be summarized in
138 a phrase or in a paragraph, but, even if the conclusion is a part of a

139 well-built article, it does not always appear. It does not present novel
 140 information, as it is only a summary.

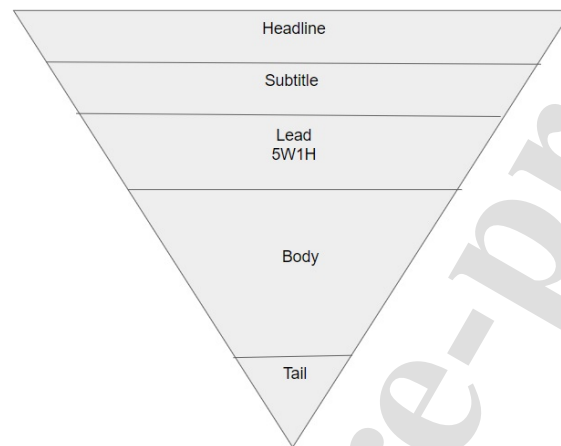


Figure 1: Inverted pyramid in newspaper articles

141 Besides the news structure, journalism purists argue that a story is not
 142 complete until the essential content is presented by answering six questions:
 143 **WHAT, WHO, WHERE, WHEN, WHY** and **HOW**. This method is known as ‘five Ws and
 144 one H’ (5W1H) (Chakma & Das, 2018a; Kim et al., 2012; Wang et al., 2010).
 145 Specifically, the six questions refer to:

- 146 ● **WHAT:** The circumstances, the event, the facts.
- 147 ● **WHO:** People involved in the events.
- 148 ● **WHERE:** The location where the events occurred.
- 149 ● **WHEN:** The time or the moment when the events occurred.
- 150 ● **WHY:** The reason or the cause of the event.
- 151 ● **HOW:** The way events have developed.

152 The 5W1H method is essential in the lead construction (Chagas, 2019).
 153 Besides, the lead is an essential part in a piece of news as it presents the
 154 main elements of an article: fact, actors, place, time, reason and manner, thus
 155 answering the six questions that are key to communicating a story accurately.

156 However, these questions are not always answered in the lead. It is sufficient
157 that only the two or three most important questions are answered and the
158 remaining questions will be answered in detail in the body of the news.

159 From a computational perspective, automatic extraction of 5W1H was
160 applied to different tasks and languages, such as English or Chinese (NIST,
161 2011) (Hamborg et al., 2018) (Chakma & Das, 2018b) (Han et al., 2013)
162 (Wang, 2012). These works demonstrated that the task is feasible, where for
163 instance, GiveMe5W1H (Hamborg et al., 2018) are obtaining a mean average
164 generalized precision 0.73 for all categories and 0.82 for ‘who’, ‘what’, ‘when’,
165 and ‘where’ for English language. Despite these results being encouraging,
166 as far as we know, those tools are not available in Spanish, and no other
167 similar resource was found at this moment.

168 *2.2. Fake News using NLP*

169 Considering that digital information is disseminated exponentially, natu-
170 ral language processing and Machine Learning (ML) approaches play a fun-
171 damental role in fake news detection (Dale, 2017). Given that assessing
172 the veracity of a news story is complex from an engineering point of view,
173 the research community is approaching this task from different perspectives
174 (Saquete et al., 2020).

175 Current fake news detection research has been conducted treating each
176 news piece as a whole to be classified with a veracity category based on:
177 lexical, syntactic and semantic content of the news as a whole (also known
178 as content-based features); or, issues related to the user or viralization of the
179 news (also known as context-based features)(Conroy et al., 2015).

180 Fake news detection currently focuses on studying linguistic aspects of
181 falsehood by identifying different types of features for fake news. (Zhou &
182 Zhang, 2008) proposed a system with the features classes, such as quan-
183 tity (amount of information), language complexity, expressiveness, message
184 content: n-grams, affect (positive or negative emotions), etc. (Pérez-Rosas
185 et al., 2018) described a similar set of features, grouped by general categories,
186 such as ngrams, punctuation, psycholinguistic features, readability and syn-
187 tax. It is very common to use the Linguistic Inquiry and Word Count(LIWC)
188 (Newman et al., 2003), which is a text analysis program that counts words in
189 psychologically meaningful categories and is available in different languages.
190 Using the Spanish language, (Almela et al., 2012) created an opinion dataset
191 consisting of 200 assessments about different topics and tested the categories
192 of LIWC.

193 Shloka Gilda (Gilda, 2017) demonstrated the relevance NLP to detect
194 fake information. They used time period frequency-inverse record frequency
195 (TFIDF) of bi-grams and probabilistic context free grammar (PCFG) detec-
196 tion. Very recent works like (Faustini & Covoos, 2020) proposed extracting
197 text features to deal with the problem at a multilingual level.

198 Stylometry is the application of the study of linguistic style generally to
199 written language. Regarding automatic Fake News detection, Potthast et al.
200 (Potthast et al., 2018) used stylometry, combining writing style features such
201 as n-grams, stop words, and parts-of-speech; and ones specific to the news
202 domain, such as 10 readability scores and dictionary features, each indicating
203 the frequency of words from a tailor-made dictionary in a document, using
204 the General Inquirer Dictionaries as a basis. The domain-specific features
205 include ratios of quoted words and external links, the number of paragraphs,
206 and their average length. Afroz et al. (Afroz et al., 2012) also used stylometry
207 to detect deception in online writing. More than 700 features were selected
208 (lexical, syntactic, content specific, grammar and vocabulary complexity, un-
209 certainty, etc). They used three feature sets to identify stylistic deception:
210 i) Writeprints feature set (lexical, syntactic, and content specific); ii) Lying-
211 detection feature set (such as q quantity, vocabulary complexity or specificity
212 and expressiveness); iii) 9-feature set (authorship-attribution features), nine
213 features that were used in the neural network experiments in Brennan’s work
214 (Brennan & Greenstadt, 2009). The main conclusion was that two kinds of
215 adversarial attacks —imitation and obfuscation— can be detected with high
216 accuracy using a large feature set. Non-content specific features have the
217 same accuracy as content-specific features, and even by ignoring the contex-
218 tual similarity of documents, it is possible to detect adversarial documents
219 with sufficient accuracy. Furthermore, previous linguistic research has shown
220 that the frequencies of common function words are content neutral and in-
221 dicative of personal writing style (Mosteller & Wallace, 1963).

222 Regarding context features, Kai et al. (Shu et al., 2019) proposed a
223 technique that exploits relationships among publishers, news pieces and users
224 to predict fake news. They employ a linear classifier and assign each user a
225 credibility score based on the user’s online behavior. A low credibility score
226 correlates to fake news.

227 Volkova et al. (Volkova et al., 2017) presented a technique that classifies
228 suspicious posts by combining content and context features via the use of
229 linguistic and network features.

230 Both Machine Learning (ML) and Deep Learning (DL) algorithms applied

231 the previously mentioned content and context features and delivered similar
 232 results when tackling the problem of classifying the text. A summary of the
 233 most commonly used detection strategies are indicated below.

- 234 • Classification approaches based on **Machine Learning**: (Gravanis
 235 et al., 2019; Conroy et al., 2015; Rubin et al., 2016; Pérez-Rosas et al.,
 236 2018; Almela et al., 2012; Afroz et al., 2012; Shu et al., 2019; Chen &
 237 Chen, 2014; Mihalcea & Strapparava, 2009)
- 238 • Classification approaches based on **Deep Learning**: (Das Bhattachar-
 239 jee et al., 2017; Volkova et al., 2017; Ren & Ji, 2017; Zhou & Zhang,
 240 2008; Monti et al., 2019; Verma et al., 2019; Rashkin et al., 2017)
- 241 • Classification approaches based on **Ensemble Learning approaches**:
 242 Very recent works are not using a single ML or DL model to tackle the
 243 problem, but an ensemble learning approach (Agarwal & Dixit, 2020).
 244 Additionally, some approaches optimize the weights of the ensemble
 245 with an external technique, such as Self-Adaptive Harmony search
 246 (Huang & Chen, 2020).
- 247 • Other approaches: (Brennan & Greenstadt, 2009)

248 Most previously cited systems use ML as a detection system, and specif-
 249 ically SVM in most cases. It is true that in recent years systems based on
 250 LSTM and DL in general have been incorporated, which use open systems
 251 such as BERT (Devlin et al., 2018). From the analysis of the literature, com-
 252 bining linguistic features with ML or DL approaches obtains some interesting
 253 results, but they seem to reach the ceiling in terms of performance. This
 254 suggests that hybrid methodologies that combine these content approaches
 255 with context information could provide a strategy to enhance performance.
 256 In addition, often, the ML or DL approximations behave like black boxes
 257 which makes it difficult to explain the generated models. The use of en-
 258 semble learning can boost performance, especially when aggregating several
 259 low-performing models, or models with different hypothesis spaces. For ex-
 260 ample, one model based on linguistic features and another model based on
 261 external knowledge. Neural models are not often explicitly ensembled, since
 262 they already offer techniques to achieve the same effect (e.g., dropout).

263 *2.2.1. Fake News datasets*

264 The purpose of this section is to present the structure currently followed
 265 by the most relevant datasets that Fake News Detection systems are using.
 266 Given that our goal is to study what type of annotation is currently being
 267 used, we have analysed the datasets presented in the literature even though
 268 their language is mainly English.

269 To the authors' knowledge, current approaches are using datasets where
 270 the news article is classified as a whole with a veracity value. (Vlachos &
 271 Riedel, 2014) are the first to release a public fake news detection and fact-
 272 checking dataset that includes 221 statements. The statements were classified
 273 by using a five-point scale: true, mostlytrue, halftrue, mostlyfalse and false.
 274 After this, (Ferreira & Vlachos, 2016) have released the Emergent dataset.
 275 In this dataset, a set of claims are classified according to their veracity and
 276 the stance of articles mentioning these claims. (Pérez-Rosas et al., 2018)
 277 introduced two new fake news datasets, one obtained through crowdsourcing
 278 and covering six news domains, and another obtained from the web cover-
 279 ing celebrities, and classified as Fake or Legitimate. BuzzFeedNews³, is a
 280 dataset comprised of a sample of news published in Facebook from 9 news
 281 agencies over a week close to the 2016 U.S. election. The LIAR dataset was
 282 presented at (Wang, 2017). They collected 12.8K manually labeled short
 283 statements from various contexts spanning a decade. This dataset is larger
 284 than the previous largest public fake news datasets of a similar type. The
 285 news articles are usually classified using a veracity scale, from true to false
 286 (pants-fire, false, barelytrue, half-true, mostly-true, and true). But, again,
 287 the whole text is annotated with a category as an atomic unit. Kaggle Fake
 288 News dataset is provided by the Kaggle competition ⁴, which is is a popu-
 289 lar platform with excellent resources for those who want to learn ML and
 290 even data science. The Kaggle dataset contains English fake and true news
 291 articles from 2015–2018. The dataset contains text and metadata from 244
 292 websites and represents 12,999 posts in total.

293 Regarding Spanish datasets, Posadas et al. (Posadas-Durán et al., 2019)
 294 presented a Spanish dataset that contains 491 true news and 480 fake news
 295 items. Almela et al. (Almela et al., 2012) presented a Spanish dataset of

³available at <https://github.com/BuzzFeedNews/2016-10-facebook-fact-check/tree/master/data>

⁴available at <https://www.kaggle.com/>

296 three different topics: opinions on homosexual adoption, opinions on bull-
 297 fighting, and feelings about one’s best friend. They collected 100 true and
 298 100 false statements for each topic, with an average of 80 words per state-
 299 ment. Arguably, there is a shortage of resources in languages other than
 300 English (Silva et al., 2020), and specifically in Spanish. Besides, although
 301 some of the datasets are in Spanish, and some are even annotated with a
 302 classification based on graded nuances for truthfulness, as was the case with
 303 the English datasets, in all cases, the annotation is of the whole textual unit
 304 rather than the parts comprising it. Given previous research on the task,
 305 the novelty of the work presented here relies on an architecture that exploits
 306 a new fine-grained annotation⁵ in a two-layer architecture. This allows the
 307 reduction of noise when training ML and DL systems. Even though the pro-
 308 posal is focused on Spanish and the health domain, it can be readily applied
 309 to different languages and domains.

310 3. A New Benchmark dataset for Spanish Fake News Detection

311 Next, the definition of the fine-grained annotation scheme, known as FN-
 312 DeepML, is presented, as well as the information about the dataset created
 313 using the said annotation scheme.

314 3.1. FNDeepML Annotation scheme

315 The annotation scheme applied to the dataset is able to distinguish the
 316 structure of the news piece, the essential parts within it and the characteristic
 317 elements that shape news. The scheme comprises two levels of representation:

- 318 1. *Newspaper article structure*: At this first level, the five elements of
 319 the newspaper article structure are annotated using one of these tags:
 320 HEADLINE, SUBTITLE, LEAD, BODY and CONCLUSION. In the case of
 321 the title and subtitle, it will almost always be the first two sentences,
 322 the lead is usually the first introductory paragraph of the news, the
 323 body usually corresponds to the remaining paragraphs of the news
 324 and the conclusion, as a rule, is the last paragraph or any concluding
 325 sentence. Furthermore, another tag has been defined at this first level:
 326 QUOTE. This tag could appear embedded in the previous elements. It

⁵The annotation is performed manually in training and automatically in testing

327 is used when an element or sentence textually quotes a message or
328 reproduces an already reported idea.

329 For each tag, there is a numerical `id` attribute to identify each element;
330 and a `type` attribute, that will indicate the value of truth or deception.
331 These values will be indicated as follows: “T” (true text), “F” (fake
332 text) or “U” (a text whose veracity is unknown). In this way, fake
333 and true elements can be detected in the same news piece. In the
334 case of the `QUOTE`, there is no `type` attribute but an attribute called
335 `author_stance` whose possible values are: “D” (the author disagrees
336 with the quote); “A” (the author agrees with the quote); and “U”
337 (Unknown, if the author’s stance is not clear). `QUOTE` is an element
338 that differs from the basic inverted pyramid structure elements because
339 it is only used to frame a set of external information –5W1H tags–
340 with a veracity value that should not be learned by the system in the
341 same way as the rest of 5W1Hs. This is due to the fact that it is
342 information reported with which the author may or may not agree
343 (depending on the `author_stance` value). So, the veracity value of the
344 5W1Hs within a `QUOTE` will be tuned by the author’s stance during the
345 training process. For that reason, the “type” attribute linked to the
346 `QUOTE` tag is not required.

347 2. *Essential news content (5W1H):*

348 In the second level of annotation, the essential content of the news
349 piece is marked by annotating the answers to “the 5 Ws and the 1 H”,
350 using the following tags for each case: `WHO`, `WHEN`, `WHERE`, `WHAT`, `WHY`
351 and `HOW`. All present 5W1H elements incorporated in the news piece
352 were annotated. These tags have two mandatory attributes and one
353 optional. All the 5W1H tags are annotated with the attributes `type`,
354 with the same description as the first level tags; and `id` to determine
355 if more than one content tag appears in the same news piece. For
356 example, if there are two `WHO` items, if they refer to different people they
357 would have a different `id` value. There is also an optional attribute,
358 termed `not_relevant`, and this term is assigned a true value when the
359 information provided by the 5W1H tag’s content is not semantically
360 relevant to determine the veracity of the news article. In order to
361 annotate the 5W1H items, first, the different facts found in the text
362 are detected, though understanding a fact from a given sentence means
363 being able to answer “Who did what to whom, when, where, why, and
364 how?”. To answer such questions of who, what, etc., it is important

365 to identify each syntactic constituent of a sentence such as predicates,
 366 subjects, objects etc. Rules already defined in the literature to identify
 367 the answer to the 5W1H questions have been followed (Voorhees, 2001)
 368 (Hamborg et al., 2018) (Chakma & Das, 2018b) (Han et al., 2013)
 369 (NIST, 2011) . Semantic role labelling tools⁶ were used to support
 370 manual annotation.

371 Moreover, metadata that are part of news content and provide information
 372 about the creation of news are the domain (**DOMAIN**), the source (**SOURCE**),
 373 the date (**DATE**), and the author (**AUTHOR**).

374 3.2. Dataset description

375 To create a Fake News dataset in Spanish (Bonet-Jover et al., 2020b),
 376 news documents in Spanish belonging to the health domain (topics such as
 377 COVID-19 which is a 50% of the dataset) were automatically collected⁷. To
 378 build the dataset in a balanced manner, fake and true news were collected
 379 from several online newspapers, blogs and fact-checking websites. The follow-
 380 ing news websites were used for collecting fake news, among others: Biosalud;
 381 Tengafe; Okdiario; Bioguia; Eje21; La Cháchara; Tudiario.net; Vidanatu-
 382 ralia; TICbeat; and, Acta sanitaria. For true news, websites such as the
 383 following were used among others: Kernpharma; Cuidateplus; Cinfasalud;
 384 Boticaria García; Comer o no comer; Julio Basulto; Nutrimedia; Vital; and,
 385 the press sections of official organizations’ sites —The World Health Organi-
 386 zation (WHO), La Asociación Española Contra el Cancer (AECC) , or the
 387 National Cancer Institute (NCI)—.

388 A total of 200 news documents were collected. More specific figures re-
 389 lating to the dataset built are presented in Table 1.

Type of News	No Docs	No tokens	Avg tokens per doc	Avg tokens Headline	Avg tokens Lead	Avg tokens Body
<i>True News</i>	105	75951	723	12	77	562
<i>Fake News</i>	95	58581	617	12	63	494
Total	200	134532	670	12	70	530

Table 1: General dataset description

⁶<http://nlp.lsi.upc.edu/freeling/node/1>

⁷Corpus download: <https://doi.org/10.5281/zenodo.4090914>

390 A manual annotation was carried out on the news collected, following the
391 FNDeepML annotation scheme described in Section 3.1.

392 To ensure the veracity of news as well as that of the different 5W1H items,
393 a manual cross-referencing information checking procedure was conducted
394 using information from official websites like WHO and the fact-checks col-
395 lected by Spanish fact-checking organizations belonging to the IFCN⁸, such
396 as Newtral⁹, Salud sin Bulos¹⁰, Maldita¹¹, Chequeado¹² or, AFP Factual¹³.
397 Fact-checking agencies verify the information delivered in the different med-
398 ia in order to determine its veracity and correctness. They publish these
399 fact-checks to make them available to the public. Furthermore, the online
400 application entitled “Google Fact Check Explorer”¹⁴ was also used to check
401 the veracity of the information.

402 This procedure verifies the veracity category of each 5W1H by searching
403 in these resources and determining if there is a previous fact-check where the
404 5W1H element is involved, whereby the corresponding category assigned to
405 the fact-check would be assumed. If information does not appear in any of
406 the sites mentioned above, we cannot determine the truthfulness or falseness
407 and hence the category of Unknown is adopted. Determining the veracity
408 category of each 5W1H element is dependent on their context and they would
409 be classified as true or false depending on their relationship with other 5W1H
410 elements, as well as the context in which the statement is included. For
411 example, it is possible to have different veracity values of the same WHO. Take
412 examples (1.a) and (1.b) where WHO=“Donald Trump” appears in different
413 newspaper documents, and after the manual cross-referencing procedure, one
414 was found True and the other was a hoax and assigned a False value.

415 (1) a. <WHO id=1 type='T'> Donald Trump </WHO> is the new candidate for US elec-
416 tions in 2020

417 b. <WHO id=1 type='F'> Donald Trump </WHO> discovers the COVID-19 vaccine

418 Table 2 presents the percentage and total items per document part clas-

⁸International Fact Checking Network (<https://www.poynter.org/ifcn/>) is a unit of the Poynter Institute dedicated to bringing together fact-checkers worldwide.

⁹<https://www.newtral.es/>

¹⁰<https://saludsinbulos.com/>

¹¹<https://maldita.es/>

¹²<https://chequeado.com/>

¹³<https://factual.afp.com/>

¹⁴<https://toolbox.google.com/factcheck/explorer>

419 sified as True, False and Unknown of the whole dataset, following the previ-
 420 ously defined annotation scheme that was carried out manually. The results
 421 indicate a very balanced dataset. Regarding the QUOTE tag, since this ele-
 422 ment does not have a veracity value, it is not included in the table, but there
 423 are 8 QUOTE in the false news part of the dataset, quoting statements that
 424 support the fake news, and 140 in the true news part of the dataset, which
 425 confirms that there is a high amount of refutations present in current true
 426 news.

Type	HEADLINE	SUBTITLE	LEAD	BODY	CONCLUSIONS
<i>True</i>	50.75%	52.22%	46.45%	53%	50.40%
<i>False</i>	45.27%	28.89%	33.88%	47%	33.60%
<i>Unknown</i>	3.98%	18.89%	19.67%	0%	16.00%
Total items	200	90	183	200	125

Table 2: Percentage and total number of items per document part classified as True, False and Unknown of the whole dataset

427 Each news piece was divided into the parts presented in Section 2.1 and
 428 the 5W1H found in the three top parts of the content (headline, subtitle and
 429 lead) were also marked. The experts were asked to mark the divided items
 430 with true, false or unknown¹⁵ based on the fact-checks of the news. Details
 431 of the figures regarding 5W1H are shown in Table 3.

Type	WHAT	WHO	WHEN	WHERE	WHY	HOW
<i>True</i>	41.64%	0.13%	30.41%	39.67%	32.26%	42.72%
<i>False</i>	35.43%	0.26%	25.77%	19.33%	45.16%	38.35%
<i>Unknown</i>	22.93%	99.61%	43.81%	41.00%	22.58%	18.93%
Total items	1112	766	194	300	62	206

Table 3: Percentage and total number of items per type of question (5W1H) classified as True, False and Unknown of the whole dataset

432 Considering the false news part of the dataset, Table 4 presents the per-
 433 centages of the different veracity values obtained for each of the news struc-
 434 ture elements as well as for the different 5W1H items. This table only includes
 435 figures extracted from false news articles of the dataset excluding true news
 436 wherein all elements are true.

¹⁵The information provided was not fact-checked as true or false

Item	False (%)	True (%)	Unknown (%)	Total items
HEADLINE	95.79	0	4.21	95
SUBTITLE	68.42	10.53	21.05	38
LEAD	75.31	6.17	18.52	81
BODY	100	0	0	95
CONCLUSION	80.38	5.88	13.73	51
WHAT	68.70	6.11	25.18	409
WHERE	24.79	22.31	52.89	121
WHEN	43.02	9.30	47.67	86
WHO	0.59	0	99.41	340
WHY	61.54	15.38	23.08	39
HOW	60.82	16.49	22.68	97

Table 4: Distribution of false, true and unknown items found in the false news part of the dataset, excluding true news

437 After a manually analyzing the dataset and the figures presented in Table
 438 4, some preliminary conclusions regarding the false part of the dataset were
 439 extracted:

- 440 • *Newspaper article structure*: The *headline* is practically always false
 441 in news documents detected as false. Obviously, the *body*, upon being
 442 annotated as a whole will be classified as false for all fake news. Fur-
 443 thermore, the *headline* and the *body* are presented in all news, but the
 444 *lead* is not always part of the false news structure.
- 445 • *The 5W1H*: *What* is the part where most false information is provided,
 446 although there is also a high degree of undefined information. The false
 447 information provided in the *Why* and *How* is also very high and close
 448 to the *What* values. In the case of *Who*, *When* and *Where* items, there
 449 is a high degree of vagueness, especially in *Who* items. Objective news
 450 provides accurate and concrete data, so detecting these inaccuracies
 451 enables us to determine if a news story is reliable. A few examples of
 452 vague *Who* tags are: “*los expertos*” (“the experts”) or “*investigadores*”
 453 (“researchers”). These *Who* terms are generic, and not specific authors
 454 because fake news usually avoid revealing a specific source that would
 455 make the information reliable. Concerning *Where* tags, some of the
 456 imprecise examples are: “*en algunas ciudades*” (“in some towns”) or
 457 “*en otros países*” (“in other countries”). In these cases, the examples
 458 indicated do not refer to a specific place and that makes the information
 459 imprecise. With regard to *When* tags, some vague expressions include:

460 “*hace unos meses*” (“some months ago”) or “*en los próximos años*”
 461 (“in the next years”). Just like places, times are also ambiguous, so
 462 information is not reliable.

463 3.3. Dataset annotation task

464 The annotation of the dataset was first applied by an annotator with
 465 linguistic training in translation and interpretation. This annotator was the
 466 person in charge of compiling the dataset and implementing the annotation
 467 schema in news. For the annotation, an annotator with journalism training
 468 was also involved. In the first phase, a simple annotation of very few
 469 documents was done in order to train the two annotators on the guidelines
 470 (Bonet-Jover et al., 2020a). Once the first annotation was done, the quality of
 471 the annotation scheme was analyzed according to the annotation agreement,
 472 including only the items where both annotators coincided. Subsequently, a
 473 meeting was held to analyse the items with different annotations with the
 474 aim of arriving at a consensus. Afterwards, the modifications required were
 475 actioned, both in terms of the guidelines and the annotations.

476 In order to measure the quality of the dataset annotation, an inter-
 477 annotator agreement between two annotators using Cohen’s *kappa* (Cohen,
 478 1960) was performed obtaining $k=0.737$ for the 1st level of annotation cate-
 479 gories and $k=0.851$ for the 2nd level of annotation categories, which validates
 480 the labeling.

481 As for the annotation time for the 200 news article dataset, 200 hours
 482 were employed (1 hour per document), 20 hours for correcting mistakes of the
 483 dataset and 30 hours for annotator training and comparison of annotations.
 484 In the case of disagreement, the annotators compared their annotations and
 485 reached a consensus, but these cases took approximately 30 extra hours to
 486 resolve, increasing the total time to complete the process to 280 hours.

487 4. Pandemic Fake News Detection system: Design and Develop- 488 ment

489 A two-layer architecture based on a pipeline is proposed. The rationale is
 490 based on the hypothesis that the structural parts and essential content of a
 491 news piece have specific veracity values, which influence the overall veracity
 492 value of the news story. This can also be inferred from the conclusions
 493 obtained in the aforementioned analysis of the dataset. The architecture

494 comprises five different phases, structured in two layers, and is graphically
 495 depicted in Figure 2.

496 The two layers and their corresponding phases of the architecture are as
 497 follows:

- 498 • **Structure Layer:** This layer is responsible for structuring the text
 499 according to the two levels of information representation. First, the
 500 news story is divided according to the journalistic structure, and then
 501 the 5W1H elements of each part of the structure are determined.
 - 502 – *Phase 1. Journalistic Structure Segmentation:* Given as input a
 503 news item from a traditional digital media, this first module is
 504 responsible for dividing the news into the parts of the structure
 505 defined for a news item. Therefore, the output of this module
 506 is the news piece divided in HEADLINE, SUBTITLE, LEAD, BODY
 507 and CONCLUSION.
 - 508 – *Phase 2. Essential content (5W1H) Extraction:* Given as input
 509 the news piece divided in parts, this module extracts the 5W1H
 510 components from each part of the news.
- 511 • **Veracity Layer:** This layer is the crucial element of this research and
 512 its purpose is to determine the veracity of each of those parts previously
 513 detected, as well as to predict the veracity of the news piece using
 514 the veracity of the different components. Determining the veracity, as
 515 explained in Section 1 implies automatically determining the accuracy
 516 and truthfulness of a piece of information within a news document.
 - 517 – *Phase 3. Essential content (5W1H) External Enrichment:* Given
 518 the 5W1H components of the news piece, this module is in charge
 519 of enriching the information of each component by using external
 520 fact-checking knowledge.
 - 521 – *Phase 4. Essential content (5W1H) Veracity Predictor:* This
 522 module, using the annotation of all the possible features (textual
 523 and fact-checking knowledge) of the 5W1H components, classifies
 524 each component in a veracity value.
 - 525 – *Phase 5. News article Veracity Predictor:* The last module, us-
 526 ing the veracity classification of each component, is in charge of
 527 predicting the veracity of the whole news item, which is the final
 528 output of the pipeline proposed.

529 The integration of the phases as a pipeline results in a prediction of
 530 the veracity of the news item. Although the Structure Layer is not the
 531 fundamental point of this research, possible approaches have been proposed
 532 for each of these phases, but without going into depth on their solution, which
 533 is a segmentation task whereas this work focuses mainly on the automatic
 534 detection of fake news. In the following sections, the development of each of
 535 the aforementioned phases is explained in more detail.

536 4.1. Journalistic Structure Segmentation

537 This phase structures the news story according to the journalistic struc-
 538 ture previously presented in Section 2.1. Given a plain news item as input,
 539 a initial preprocessing is performed to obtain HEADLINE and SUBTITLE fol-
 540 lowing a set of simple rules. These rules are variable and depend on each
 541 site’s structure. After that, the remaining text will be divided into LEAD,
 542 BODY and CONCLUSION applying a named entity recognition approach. Using
 543 Spacy library¹⁶, a tokenization of the news document is performed, and a
 544 set of features are obtained for each token (see Table 5). The features are
 545 defined in the Spacy library documentation¹⁷.

Feature	Description
text	Original text of the token.
lemma	Lemmatized version of the token.
pos	Coarse part-of-speech tag, e.g., VERB, NOUN, etc..
tags	Several fine-grained part-of-speech tags such as person, number, tense, etc.
dep	Label of the token in the dependency tree.
shape	Syntactic representation of the token shape.
ent_type	General-purpose entity label, e.g., PERSON, ORG, etc.
is_alpha	Boolean value indicating if the token is alphanumeric.
is_stop	Boolean value indicating if the token is a stopword.
index	Relative index of the token in the document, between 0 (first token) and 1 (last token).

Table 5: Token-level features, extracted with Spacy.

546 News documents are segmented at the token-level using a Conditional
 547 Random Fields (CRF) model (Sutton et al., 2012) trained on the token fea-
 548 tures described in Table 5. To introduce context, each token feature set
 549 is complemented with the features of surrounding tokens (both before and

¹⁶<https://spacy.io/>

¹⁷<https://spacy.io/api/token#attributes>

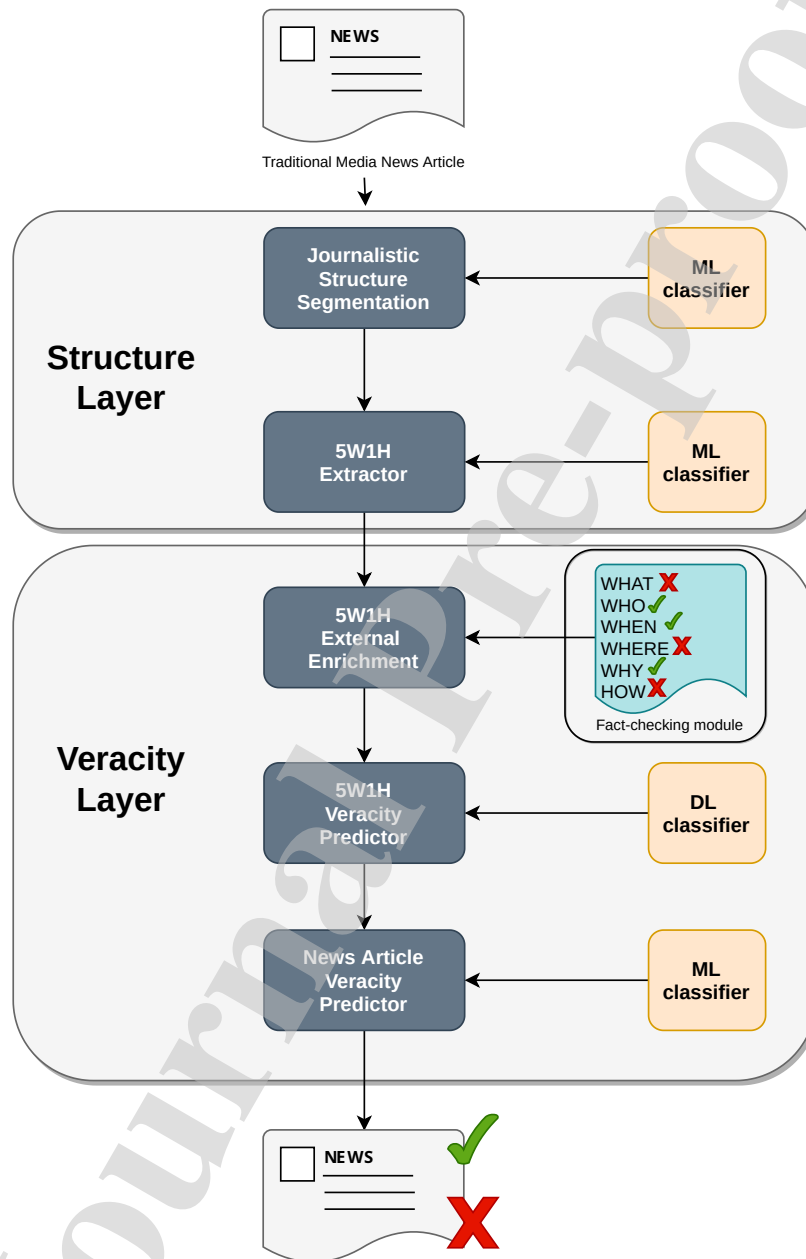


Figure 2: Pandemic Fake News Detection system's architecture

550 after) in a small window of size 0 to 3. This parameter can be adjusted
 551 to improve accuracy at the cost of a larger computational cost. The CRF
 552 model is trained using *sklearn-crfsuite*¹⁸. The segmentation problem is thus
 553 modeled as a sequence tagging problem, where each token is assigned one of
 554 these labels: LEAD, BODY and CONCLUSION.

555 After this process, the segmented news item is the output of this module,
 556 as shown in this example.

	Token	Features	Structure Part
	token1	...	=> Lead
	token2	...	=> Lead
557 (2)	token3	...	=> Body
	token4	...	=> Body
	...		
	tokenN	...	=> Conclusion

558 4.2. Essential content (5W1H) Extraction

559 Using all the features per token previously obtained, a second CRF model
 560 is used to classify each token of each part into one of the 5W1H components,
 561 or NONE. As observed in Table 3, there is a large imbalance in the labels'
 562 distribution, which provokes a poor performance of models trained to predict
 563 all classes at once. For this reason, a two-level hierarchical classification is
 564 performed, where labels are divided into two sets: the first level consists of
 565 the most common labels (NONE and WHAT) while the least common labels
 566 are grouped in a special REST class; and, the second level comprises only
 567 the least common classes (HOW, WHEN, WHERE, WHY and WHO). This
 568 allows the training of two separate models that can deal better with the
 569 unbalanced distribution of the labels, allowing each model to only focus on
 570 a smaller set of classes for which their relative numbers are similar.

571 The fact that one of the features obtained by Spacy is the Named Entities
 572 (NE) is very useful in this module since they are related to some questions
 573 such as LOCATION for WHERE, PERSON/ORGANIZATION for WHO, or
 574 TIME for WHEN. Furthermore, the same features shown in Table 5 are used
 575 to represent each token. Likewise, a window size can be adjusted to include
 576 more context at the cost of a larger feature set and increased computational
 577 cost.

578 In the case of classes with a smaller set of examples in the dataset, in
 579 addition to the features used in the first level, the semantic roles of the text

¹⁸<https://sklearn-crfsuite.readthedocs.io/en/latest/>

580 will also be used in the second level of the hierarchical model. According to
 581 (Moreda et al., 2011), the use of semantic roles can improve the detection of
 582 answers to the 5W1H, especially when dealing with questions whose answer
 583 is not itself a Named Entity. For example, in this sentence where semantic
 584 roles are annotated, the role `AM-LOC` is the answer of a Where question:

585 (3) Where was Pythagoras born? Samos Pythagoras was born [`AM-LOC` on the island
 586 of Samos].

587 In order to annotate semantic roles, Freeling(Padró & Stanilovsky, 2012) is
 588 used because this tool also annotates semantic roles in Spanish.

589 As this module performs, the different 5W1H are detected and an example
 590 of the output obtained by this module for each news document is presented
 591 next.

	Token	Features	Structure Part	5W1H
	token1	...	Lead	=> None
	token2	...	Lead	=> What
	token3	...	Body	=> What
592 (4)	token4	...	Body	=> What
	token5	...	Body	=> What
	token6	...	Body	=> Who
	...			
	tokenN	...	Conclusion	=> Where

593 As can be seen in the example, each 5W1H might span multiple tokens,
 594 as in the case of the "What" item that comprises token 3,4 and 5.

595 4.3. Essential content (5W1H) External Enrichment

596 This module is in charge of enriching each 5W1H component by using
 597 external fact-checking knowledge. As our intention is to look only for essential
 598 content, i.e. the treatment of each 5W1H element, the process is carried
 599 out using those elements rather than raw text. The first point we would
 600 like to stress is that performing a fact-checking module is not a trivial task
 601 and implies in-depth research in itself, which is beyond the scope of this
 602 work. Nevertheless, in order to add external knowledge to the proposed
 603 pipeline, a simple fact-checking module has been implemented that will be
 604 able to detect whether the 5W1H elements of a news story are part of any
 605 previous fact-check. Of course, this implies that the fact exposed in the news
 606 story has been previously refuted. The purpose of this module is not to
 607 determine the veracity of each 5W1H, but to extract external information
 608 that, in addition with the textual content, helps in the prediction of the
 609 veracity of each component performed in Phase 4.

610 More specifically, this module uses the Google Fact Check Tools API ¹⁹,
 611 which is based on ClaimReview markup²⁰. An example of a fact-check in
 Spanish is shown in Figure 3.

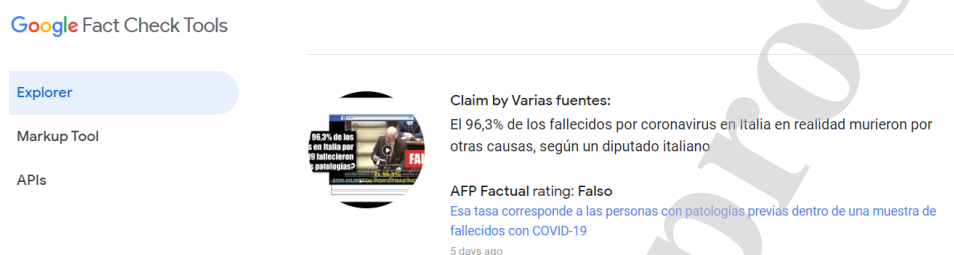


Figure 3: Screenshot of a fact-checked claim in Spanish

612 A JavaScript Client for the REST API²¹ was implemented to access this
 613 tool²². The essential content of the news (5W1H) is searched as follows.
 614 For each part of the document (title, subtitle, ...) the 5W1H items are sent
 615 separately to the API for checking their veracity. If a value is found, the
 616 label is updated to that value. If any of the 5W1H items do not receive a
 617 veracity value or receive contradictory values, a second check will be done
 618 with all 5W1H items of that part to add context information. To do that,
 619 all the items will be concatenated and sent again to check their veracity. In
 620 this case, the value obtained will serve to update the veracity value of each
 621 item. The API's textual rating is mapped to one of our True/False/Unknown
 622 categories.

624 A simple solution is proposed for this fact-checking module, but in future
 625 the fact-checking procedure should be enhanced.

626 4.4. 5W1H Veracity Predictor

627 This phase is in charge of predicting the veracity value of each 5W1H
 628 component of each news document, based on all the evidence collected in
 629 Phase 3 plus the textual content of each element. Due to the complexity

¹⁹<https://developers.google.com/fact-check/tools/api/>

²⁰<https://schema.org/ClaimReview>

²¹<https://factchecktools.googleapis.com/v1alpha1/claims:search>

²²For further information, please consult the API documentation at:
<https://developers.google.com/fact-check/tools/api/reference/rest>

630 of the task, this problem is tackled using DL, since solving the problem in
 631 this phase requires not only dealing with textual features of the components
 632 but also high level features obtained from external knowledge that enrich the
 633 components (Fact-checking in this case). In order to predict the veracity of
 634 each component, the module uses a sequential LSTM-Convolutional model
 635 with the following architecture (see Figure 4):

- 636 1. A trainable embedding layer with output dimension of 32, a maximum
 637 sequence length of 100 tokens (longer sequences are truncated) and
 638 a maximum number of 1000 vocabulary entries (built during training
 639 from the top 1000 tokens by frequency in the training set).
- 640 2. A dropout layer with a dropout rate of 0.25.
- 641 3. A 2D convolutional layer with 64 filters and kernel size of 5.
- 642 4. A max-pooling layer with a pool size of 4.
- 643 5. A second dropout layer with a dropout rate of 0.25.
- 644 6. An LSTM layer with an output dimension of 70.
- 645 7. A third dropout layer with a dropout rate of 0.25.
- 646 8. A dense layer for one-hot encoding of the label of the 5W1H compo-
 647 nent (i.e., “WHAT”, “WHERE”, “WHY”, etc.).
- 648 9. A dense layer for one-hot encoding of the label of the article part in
 649 which the 5W1H component appears in the news article (i.e., “LEAD”,
 650 “BODY”, etc.).
- 651 10. A concatenation of the previous three layers.
- 652 11. A final dense layer with 3 outputs (one for each class of *True*, *False*,
 653 *Unknown*) with a softmax activation function.

654 This model was adapted from a classic architecture for sequence classifi-
 655 cation proposed in the Keras ML library²³ and modified to fit the number of
 656 features and training examples available in this research. The exact param-
 657 eters of each layer (e.g., layer sizes, dropout rate, number of filters, etc.) were
 658 decided after a short manual tuning among a range of sensible parameters.

659 When the fact-checking information is available, a parallel two-layer dense
 660 feed-forward network (with a total of 130 trainable parameters) is added,
 661 whose output is concatenated before the final dense layer with the previous
 662 model.

²³<https://keras.io/>

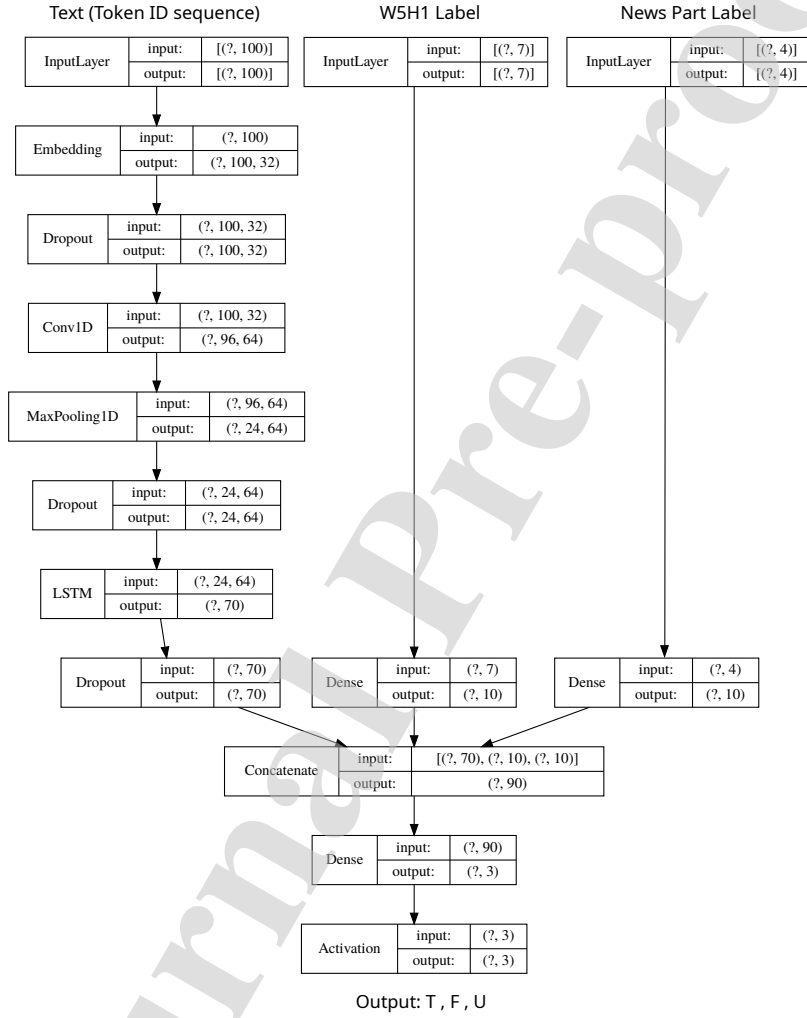


Figure 4: Graphical representation of the 5W1H Veracity Predictor DL architecture. The type of each layer and tensor shapes are reported. Shapes with size “?” indicate the batch dimension, whose size is determined at training time and does not influence the total number of parameters.

663 The overall model contains 80,377 trainable parameters (80,507 when
 664 adding the fact-checking features), and is trained with the Adam optimiza-
 665 tion scheme using categorical cross-entropy as loss function, with the recom-
 666 mended hyperparameters (Kingma & Ba, 2014). To improve performance,
 667 this model is trained with early stopping, based on the loss measured on a
 668 separate 10% of the training set, with 3 epochs of patience Prechelt (1998).
 669 The model is implemented in the Python *keras* library.

670 The DL model is trained independently on each continuous sequence of
 671 tokens that belongs to the same 5W1H part to predict their veracity value.
 672 At the end, using all the features previously extracted in the pipeline, the
 673 module is predicting the veracity of each component.

674 An example of the output of this module is:

	Token	Features	Structure Part	5W1H	Veracity
	token1	...	Lead	=> None	Null
	token2	...	Lead	=> What	T
	token3	...	Body	=> What	T
675	(5) token4	...	Body	=> What	T
	token5	...	Body	=> What	T
	token6	...	Body	=> Who	F
	...				
	tokenN	...	Conclusion	=> Where	T

676 In this example, the "What" item of the Lead is assigned a False veracity
 677 value; the "What" item of the body is assigned a True veracity value; and
 678 the "Who" is a false element. This means that the fact explained in the body
 679 happens, however the person involved in this fact was not the one indicated
 680 in the news document. The last phase will learn that certain entities are
 681 less relevant than others, which is why we consider this phase to have a
 682 regularizing effect, like an ensemble.

683 4.5. News Article Veracity Predictor

684 Finally, the last phase is in charge of giving the final prediction of the
 685 news item, using one of several classic ML models (as implemented in the
 686 *scikit-learn* package²⁴):

- 687 • Logistic Regression, with an L_2 regularization factor of 1.0 and a LBFGS
 688 optimizer.

²⁴<https://scikit-learn.org/stable/>

- 689 • Decision Trees, using GINI as the criteria for feature selection.
- 690 • Support Vector Machines, with a Radial Basis Function kernel and a
691 regularization factor of 1.0.
- 692 • Multinomial Naive Bayes, with a Laplace smoothing factor of 1.0.
- 693 • Random Baseline using a stratified random strategy.

694 In this module, to represent the documents, for each part of the structure
695 of the document, their 5W1H items are aggregated according to their verac-
696 ity value, and the number of each item within each veracity value is counted.
697 Thus, considering that there are 5 parts in the structure (HEADLINE, SUB-
698 TITLE, LEAD, BODY and CONCLUSION), and 6 possible 5W1H types of
699 items (WHAT, WHO, WHEN, WHERE, WHY and HOW) within each part,
700 and each of these 5W1H items can have one of three veracity values (TRUE,
701 FALSE, UNKNOWN) the final number of numerical features generated is
702 90.

<HEADLINE id=1 type='F'><WHO id=1 type='F'>Dr. Chen</WHO> affirmed that <WHAT id=1
type='F'>cancer is cured</WHAT> <HOW id=1 type='F'>by infusing water with a slice of lemon
</HOW> <WHEN id=1 type='F'>every day</WHEN></HEADLINE>

<LEAD id=1 type='F'><WHAT id=2 type='T'>Lemon has several properties</WHAT>, but <WHO
id=2 type='F'>medical experts</WHO> <WHAT id=3 type='F'>have used it</WHAT> <WHERE
id=1 type='F'>in Asia</WHERE> <WHEN id=2 type='F'>for millions of years</WHEN> <WHY
id=1 type='F'>because it cures cancer</WHY>. It is known that <WHAT id=4 type='T'>lemon
has health benefits</WHAT>, but <WHO id=3 type='U'>renowned oncologists </WHO> <WHEN
id=3 type='F'>now</WHEN> <WHAT id=5 type='F'>stated that it is possible to kill cancer
cells</WHAT> <HOW id=2 type='F'>by consuming hot water with citrus fruits juice</HOW>
<WHEN id=4 type='F'>every morning</WHEN> <WHY id=2 type='F'>since its vitamins are up to
100 times more effective than chemotherapy.</WHY></LEAD>

Figure 5: Graphical visualization of part of the annotation of a newspaper article using FNDeepML annotation scheme.

703 For instance, considering Figure 5 annotation of the headline and lead of
704 a specific newspaper article, the following numerical features are extracted

705 from headline and lead²⁵.

```

706 {
707     HEADLINE_WHAT_TRUE: 0,
708     HEADLINE_WHAT_FALSE: 1,
709     HEADLINE_WHAT_UNKNOWN: 0,
710     HEADLINE_WHO_TRUE: 0,
711     HEADLINE_WHO_FALSE: 1,
712     HEADLINE_WHO_UNKNOWN: 0,
713     HEADLINE_WHEN_TRUE: 0,
714     HEADLINE_WHEN_FALSE: 1,
715     HEADLINE_WHEN_UNKNOWN: 0,
716     # ...
717     LEAD_WHAT_TRUE: 2,
718     LEAD_WHAT_FALSE: 2,
719     LEAD_WHAT_UNKNOWN: 0,
720     LEAD_WHO_TRUE: 0,
721     LEAD_WHO_FALSE: 1,
722     LEAD_WHO_UNKNOWN: 1,
723     LEAD_WHEN_TRUE: 0,
724     LEAD_WHEN_FALSE: 3,
725     LEAD_WHEN_UNKNOWN: 0,
726     # ...
727 }
```

728 The same type of features will be generated from the other parts of the
729 structure of the document. Each feature indicates the number of 5W1H
730 components with a specific label and veracity that appear in each part of the
731 news. For example, `LEAD_WHAT_TRUE: 2` indicates that the LEAD contains two
732 WHAT items annotated with a TRUE veracity value. The model is trained to
733 predict the overall document veracity label based on these numerical features.

734 5. Experimental Setup and Evaluation

735 5.1. Evaluation Measures

736 In order to evaluate the proposal, the commonly used NLP measures
737 (accuracy, precision, recall and F-measure) are used.

Precision (P) is the ratio of correctly predicted positive observations to the total predicted positive observations.

$$P = \frac{\#TruePositive}{\#TruePositive + \#FalsePositive} \quad (1)$$

²⁵Only some of the features are shown to exemplify the generation of these features

Recall (R) is the ratio of correctly predicted positive observations to the all observations being actual positive.

$$R = \frac{\#TruePositive}{\#TruePositive + \#FalseNegative} \quad (2)$$

F1-Score (F_1) is the weighted average of Precision and Recall.

$$F_1 = 2 * Precision * Recall \frac{Precision * Recall}{Precision + Recall} \quad (3)$$

Accuracy (Acc) is the most intuitive performance measure and it is simply a ratio of correctly predicted observations to the total observations.

$$Acc = \frac{\#TrueP + \#TrueN}{\#TrueP + \#FalseP + \#TrueN + \#FalseN} \quad (4)$$

Furthermore, the macro and micro average of each measure is given when necessary. Macro average is the average of each of the measures, whereas micro average is an average weighted by support value—which is the number of true instances for each label—. Using these measures is also important because the macro average will be poor if any class is small, but the micro average will penalise less severely in classes with very few elements. The difference between macro and micro indicates how much damage the corpus imbalance is doing to the model.

5.2. Experiments

The main objective of the experimentation proposed in this research is to demonstrate the hypothesis that because fake news is a combination of false and true information whose aim is to create confusion among readers, an adequate approach to the problem of automatic fake news detection is a two-layer architecture.

The following set of experiments for each of the layers are performed to validate our hypothesis:

- **Structure Layer performance:** A set of experiments related to the first two phases have been carried out. The proposals made are measured to assess potential areas for improvements to increase effectiveness.

- *Phase 1 performance. Journalistic Structure Segmentation:* The performance of the module doing the segmentation into LEAD, BODY and CONCLUSION of the text is measured.

762 – *Phase 2 performance. 5W1H Extractor:* In this experiment, the
 763 performance in detecting the different segments that correspond
 764 to the answer of the 5W1H is measured.

765 • **Veracity Layer performance:** A set of experiments to measure the
 766 two phases that determine both the veracity of the components and the
 767 veracity of the news have been implemented. In addition, a final ex-
 768 periment allows the validity of this work’s hypothesis to be determined
 769 by measuring the Veracity Layer as a whole. Phase 3 does not have
 770 an individual experiment since it is an enrichment phase and its valid-
 771 ity is given by the results of phase 4, which have been measured both
 772 using the information of phase 3 and without using it to determine its
 773 benefits.

774 – *Phase 4 performance. 5W1H Veracity Predictor:* This experiment
 775 measures the performance of the module that predicts the veracity
 776 value of each element of the news piece. In order to prove this and
 777 to determine the validity of this module in isolation, the 5W1H
 778 labels of the gold standard dataset have been used and the per-
 779 formance of the module using different configurations is measured
 780 by: i) using only the textual characteristics of the content of the
 781 5W1H components; ii) using only the fact-checking characteristics
 782 and; iii) using the combination of both.

783 – *Phase 5 performance. News Article Veracity Predictor:* To mea-
 784 sure the accuracy of this phase in this experiment, the phase is
 785 measured in isolation, using as training the manually annotated
 786 gold standard news pieces with the different parts of the structure
 787 as well as the 5W1H elements with their veracity value. Thus,
 788 the errors of the previous phases are avoided, and the validity of
 789 this module alone is measured. This is one of the most important
 790 experiments since it proves the validity of the proposal.

791 – *Phase 3+4+5 performance. Veracity Layer* This experiment aims
 792 to determine the effectiveness of the Veracity Layer but avoiding
 793 segmentation errors produced by the Structure Layer. Specifically,
 794 using the gold standard segmentation of the text, Phase 3, 4 and
 795 5 together are performed and measured.

796 Finally, the performance of the full pipeline is measured and a cross-domain

797 validation is performed to explore the applicability of our proposal across
798 domains.

799 6. Results and Discussion

800 This section presents the results obtained in each of the experiments
801 described in Section 5 and a discussion of those results.

802 6.1. Phase 1 performance. Journalistic Structure Segmentation

803 Table 6 presents the performance at a token level of the Structure Seg-
804 mentation Module that corresponds to Phase 1 in the pipeline.

Features	P	R	F_1	Acc
Lead	0.851	0.772	0.810	0.851
Body	0.960	0.964	0.962	0.929
Conclusion	0.710	0.836	0.768	0.648
micro avg	0.935	0.937	0.936	0.938
macro avg	0.840	0.857	0.846	0.809

Table 6: Journalistic Structure Segmentation performance

805 Overall, this module obtains a micro F_1 score of 0.936 in an independent
806 test-set of 20% of the news items. Table 7 shows the confusion matrix over
807 the test-set. As expected from a CRF-based model, no confusion occurs
808 between classes that never overlap, i.e., Lead and Conclusion. Since Body is
809 the majority class (with a support of 23,708 tokens out of a total of 28,154
810 in the test-set, or 84.11%), it is also the class with the highest F_1 . However,
811 despite their being a significantly lower number of training instances for the
812 remaining classes, their F_1 scores are significantly higher than what can be
813 expected from a random baseline. By comparison, using only the token
814 relative index produces an overall F_1 of 0.772, which is an indication that
815 most news items (in the corpus) follow a relatively similar structure in terms
816 of the relative sizes of each segment.

817 6.2. Phase 2 performance. 5W1H Extractor

818 Table 8 presents the performance at a token level of the 5W1H Segmen-
819 tation Module that corresponds to the Phase 2 in the pipeline.

820 As explained in Section 4.2, a hierarchical model is trained on different
821 subsets of classes to deal with the imbalance of labels in the dataset. As

	Lead	Body	Conclusion
Lead	2345	692	0
Body	411	22849	418
Conclusion	0	236	1203

Table 7: Confusion matrix for the Journalistic Structure Segmentation module. For each of the 28,154 tokens in a 20% test-set, the rows indicate the real label and the columns indicate the predicted label.

822 a comparison baseline, a single linear model (logistic regression) trained on
 823 the complete set of labels obtains a micro-average $F_1 = 0.932$, but a macro-
 824 average $F_1 = 0.309$. This is because the model assigns a higher importance
 825 to the most common labels and hence performs very poorly on low-count
 826 labels such as WHY ($F_1 = 0.048$), HOW ($F_1 = 0$) and WHEN ($F_1 = 0.128$).

827 The hierarchical model is trained first only on NONE, WHAT and REST
 828 (which groups all the remaining labels), producing the results shown in Table
 829 8(top), in a test-set of 20% of the news items. Then, a second model is trained
 830 only on the subset of tokens with labels HOW, WHY, WHEN, WHERE and
 831 WHO, producing the results shown in Table 8 (bottom) in the same test-set.
 832 The first level uses only syntactic and semantic features from Spacy, while
 833 the second step includes also the semantic role features from Freeling. This
 834 configuration showed better results, presumably because semantic roles are
 835 not useful for the recognition of the WHAT class, in contrast with the rest
 836 of the 5W1H components. As can be observed, each model is significantly
 837 better (in terms of macro F_1) in the corresponding sub-problem.

838 Interestingly, the first step is able to recognize the REST class exactly,
 839 which means that we can estimate the overall performance of the model by
 840 aggregating the results of both models. The combined estimated macro-
 841 average F_1 for this two-step model is 0.661, significantly higher than the
 842 0.309 provided by a single model. Furthermore, the worst performance is
 843 obtained for the HOW and WHY labels, which have the least number of
 844 instances. If we discard these labels and only consider the remaining 5 labels
 845 (including NONE) the overall macro F_1 would be 0.774. Finally, a very good
 846 performance is obtained in the WHAT label ($F_1=0.948$), which corresponds
 847 to the most important element in terms of determining the veracity of a news
 848 item. The HOW and WHY elements are important in a fact-checking process
 849 to determine the veracity of a news item, since they add nuisance and detail
 850 and might thus change the deeper meaning of a news item. For this reason,

		First step		
		P	R	F_1
NONE		0.999901	0.983090	0.991425
WHAT		0.901966	0.999378	0.948177
REST		1.000000	1.000000	1.000000
macro avg		0.967289	0.994156	0.979867
micro avg		0.986880	0.985474	0.985784
		Second step		
		P	R	F_1
HOW		0.262500	0.750000	0.388889
WHEN		0.788732	0.629213	0.700000
WHERE		0.489583	0.566265	0.525140
WHY		0.336957	0.462687	0.389937
WHO		0.844444	0.639731	0.727969
macro avg		0.544443	0.609579	0.546387
micro avg		0.694253	0.611702	0.636717

Table 8: Results for the first level (top) and second level (bottom) of the hierarchical model trained for 5W1H extraction.

851 their failed detection in this phase is likely to cause a significant decrease in
852 the overall performance of the pipeline. In contrast, a high accuracy in the
853 extraction of the WHAT label might compensate for the performance loss.
854 However, the reliable extraction of 5W1H elements in general is a difficult
855 problem, and it is not the purpose of this research to fully address it.

856 The complexity of the phase 2 task is acknowledged and for this reason
857 the literature on Automatic extraction of 5W1H presented in Section 2.1 will
858 be taken into consideration to improve future performance of the Structure
859 Layer.

860 6.3. Phase 4 performance. 5W1H Veracity Predictor

861 As explained in Section 5, this phase is evaluated in different configura-
862 tions. Using the gold standard 5W1H elements in the dataset, the validity
863 of this module in isolation is measured.

864 Table 9 presents the performance of the 5W1H Veracity Predictor Module
865 that corresponds to the 4th phase in the pipeline, with three configurations:

866 **Deep NN (Text)** uses only textual features of the tokens within each
867 5W1H component annotated in the gold standard dataset.

868 **Deep NN (FC)** uses only fact checking features of the 5W1H components,
869 automatically obtained in Phase 3.

870 **Deep NN (Combined)** uses both textual features and fact checking fea-
 871 tures of the 5W1H components.

872 For comparison purposes, two baselines are implemented: using a strategy
 873 that always predicts the majority class (**Dummy**) and using the TF-IDF rep-
 874 resentation of the text of each 5W1H component to train a logistic regression.
 875 The values correspond to the mean precision, recall, F_1 and accuracy of each
 876 model for each veracity label (i.e., **Unknown**, **True** and **False**), averaged
 877 across 10 independent runs with 80% training and 20% testing splits.

Models	Baseline		Text	Deep Learning	
	Dummy	TF-IDF		Fact-Check	Combined
Precision (T)	0.000	0.601	0.592	0.370	0.592
Recall (T)	0.000	0.471	0.547	0.930	0.523
F_1 (T)	0.000	0.528	0.565	0.529	0.554
Precision (F)	0.000	0.476	0.512	0.000	0.507
Recall (F)	0.000	0.234	0.374	0.000	0.452
F_1 (F)	0.000	0.313	0.424	0.000	0.468
Precision (U)	0.513	0.630	0.733	0.512	0.753
Recall (U)	1.000	0.837	0.837	0.993	0.821
F_1 (U)	0.678	0.719	0.780	0.675	0.784
Accuracy	0.513	0.607	0.658	0.542	0.660
Macro- F_1	0.226	0.520	0.590	0.409	0.602

Table 9: Performance results of different configurations of 5W1H Veracity Predictor using Gold standard 5W1H segmentation

878 As can be deduced from the results obtained in Table 9, determining the
 879 veracity of each of the essential contents of a news item is not a trivial task.
 880 The figures obtained make it clear that the use of textual characteristics has
 881 a limit when it comes to improving the detection of falsehood. Also, we see
 882 that combining textual information with high-level characteristics extracted
 883 from external knowledge, such as fact-checking in this case, help to improve
 884 the prediction of the veracity of each component. Obviously, the increase is
 885 limited in our case because the tools that perform with optimal results are
 886 lacking at present, and the tool we apply at present is very limited.

887 It should be noted that although the dataset is limited in size (200 news),
 888 this phase is trained with individual 5W1H phrases; hence, there is a larger
 889 number of training examples. In total there are 2,788 different 5W1H phrases,
 890 of which 2,230 are used for training (80%) and 558 are used for validation

891 (20%). With respect to features, a maximum number of 1,000 different to-
 892 kens is allowed for the embedding layer (i.e., the 1,000 most common tokens).
 893 Similarly, a maximum of 100 tokens is allowed for any 5W1H phrase in the
 894 LSTM layer. These limits maintain a small total number of trainable param-
 895 eters, which makes it feasible to achieve a better-than-baseline performance
 896 even with such a small number of training examples.

897 To better understand the behavior of the 5W1H Veracity Predictor in
 898 different types of 5W1H components, Table 10 shows the evaluation metrics
 899 aggregated per 5W1H label.

5W1H label		HOW	WHEN	WHERE	WHY	WHAT	WHO
micro	precision	0.564	0.626	0.581	0.642	0.519	0.991
	recall	0.495	0.558	0.535	0.571	0.505	0.994
	F_1	0.481	0.526	0.501	0.557	0.502	0.992
macro	precision	0.496	0.610	0.573	0.630	0.502	0.332
	recall	0.492	0.495	0.477	0.609	0.491	0.333
	F_1	0.442	0.485	0.457	0.573	0.485	0.332
accuracy		0.495	0.558	0.535	0.571	0.505	0.197

Table 10: Evaluation metrics for the 5W1H Veracity Predictor model using combined syntactic and fact-checking features aggregated per type of 5W1H component.

900 The results obtained between the different 5W1H components are quite
 901 similar, except from WHO, that, as indicated in Table 3, has a high degree of
 902 uncertainty (U veracity), resulting in a high micro- F_1 but limited accuracy.
 903 The results indicate the need to add more complex information that implies
 904 external knowledge and context in order to improve the prediction of the
 905 veracity of each component.

906 6.4. Phase 5 performance. News Article Veracity Predictor

907 The experiments with the News Article Veracity Predictor represent the
 908 most interesting results because they demonstrate that by considering the
 909 veracity of the news structure parts and the 5W1H, a suitable solution to
 910 the problem of automatic fake news detection is provided. Therefore, to avoid
 911 problems arising from previous phases, this module is measured in isolation
 912 using the gold standard 5W1H elements and their manually-assigned veracity
 913 value. The experiment demonstrates that this information is valuable when
 914 determining the veracity of the whole news document. Table 11 presents the
 915 performance of this last phase in the pipeline. The results of the different ML

916 approaches applied are shown as well as two baselines to determine if there
 917 is an improvement when using our proposal: i) a random baseline; and ii) a
 918 baseline using the TF-IDF of the whole document annotated with a unique
 919 veracity value for the document.

Model	True News			Fake News			Acc	Macro F_1
	P	R	F_1	P	R	F_1		
Baseline (Random)	0.523	0.503	0.510	0.483	0.502	0.489	0.502	0.500
Baseline (TF-IDF)	0.609	0.868	0.715	0.726	0.381	0.494	0.637	0.605
Decision Tree	0.971	0.976	0.972	0.976	0.965	0.969	0.971	0.971
Logistic Regression	0.964	0.997	0.980	0.996	0.958	0.976	0.978	0.978
Naive Bayes	0.920	0.995	0.956	0.995	0.902	0.945	0.951	0.950
SVM	0.934	0.994	0.962	0.993	0.919	0.953	0.958	0.958

Table 11: Results of News Article Veracity Predictor performance using veracity of gold standard 5W1H components

920 As can be concluded from the results in the table, all the models proposed
 921 significantly outperform the two proposed baselines. Even so, the model that
 922 obtains the best results is Logistic Regression both for detecting false news
 923 and for determining which news stories are true, obtaining a 0.978 of macro
 924 F_1 . It is especially noteworthy that using the entire annotated document with
 925 a single truthfulness value (baseline TF-IDF) the macro F_1 is 0.605. These
 926 results validate the main hypothesis set for this research, i.e., that individual
 927 5W1H components are a good predictor of overall news story truthfulness.

928 6.5. Phase 3+4+5 performance. Veracity Layer

929 To measure the whole performance of the Veracity Layer —but avoiding
 930 the errors produced by the Structure layer, i.e. segmentation modules (phase
 931 1 and phase 2)— the gold standard elements of the dataset are used and the
 932 performance from Phase 3 to Phase 5 of the architecture is measured.

933 For this purpose, the 5W1H Veracity Predictor (phase 4) was run 10
 934 independent times in different train/test splits (80%/20%), and the results
 935 of the predicted labels (in each independent test set) were concatenated.
 936 Thus, a “new” re-sampled training set is available for training and evaluating
 937 in phase 5. This allows to train the phase 5 module directly on predicted
 938 veracity labels, instead of on the gold labels, as performed in Section 6.4.
 939 Hence, if the 5W1H Veracity Predictor makes consistent mistakes on different
 940 5W1H labels, the phase 5 module might be able to correct these mistakes in

941 the aggregated prediction by assigning less weights to those labels. Results
 942 are provided in Table 12.

Model	True News			Fake News			Acc	Macro F_1
	P	R	F_1	P	R	F_1		
Baseline (Random)	0.551	0.549	0.548	0.498	0.500	0.497	0.526	0.522
Baseline (TF-IDF)	0.609	0.868	0.715	0.726	0.381	0.494	0.637	0.605
Decision Tree	0.736	0.752	0.741	0.724	0.696	0.706	0.726	0.723
Logistic Regression	0.842	0.783	0.809	0.780	0.835	0.805	0.807	0.807
Naive Bayes (Multinomial)	0.794	0.827	0.808	0.804	0.760	0.778	0.795	0.793
SVM	0.802	0.768	0.781	0.761	0.786	0.770	0.777	0.775

Table 12: Results of News Article Veracity Predictor performance trained and evaluated on the predicted labels from phase 4.

943 As can be observed, even though the results are worse than when using
 944 gold standard annotations, they are better than what could be expected if
 945 all the errors from phase 4 were carried to phase 5. Given that phase 4 at the
 946 moment obtains a maximum of 0.660 accuracy, the fact that an average 0.805
 947 can be obtained by aggregating low-accuracy estimations for each 5W1H
 948 hints at some sort of regularizing effect. We can argue that phase 5 indeed
 949 learns to correct some of the mistakes in phase 4. This is not surprising if we
 950 consider that, in phase 5, each of the individual veracity labels for each 5W1H
 951 component in a single article can be seen as the output of a single classifier, all
 952 of which are aggregated in an ensemble fashion. Hence, even if the individual
 953 components are not very reliable (i.e., on average each 5W1H component is
 954 correct 66% of the time), the overall classifier is far more reliable. It is known
 955 that ensemble models can outperform considerably each of their components,
 956 especially when the individual components make mistakes that are mostly
 957 independent of each other (see Section 2.2). It appears that in this case, a
 958 similar effect is taking place.

959 6.6. Hyper-parameter search for full pipeline performance

960 To measure the performance of the full pipeline, we applied a hyperpa-
 961 rameter search based on the open source library AutoGOAL (Estevez-Velarde
 962 et al., 2020). The hyperparameter search enables testing a large number of
 963 parameter values for different parts of the pipeline to find the combination
 964 that produces the highest performance. A total of 24 hours of computing
 965 resources was devoted to the parameter search, which resulted in a total of

101 different pipelines tested. The best pipeline found achieved an accuracy of 0.775 on a 5-step cross-validation with a random split of 80% of the data for training and 20% for testing. The hyper-parameter space contains several different ML algorithms for each phase as well as specific configuration parameters such as window size and optimization technique for CRF taggers (Phases 2 and 3), number of filters and size of embedding vectors in Phase 4, and whether to count 5W1H components by article part (headline, body, conclusion) or aggregated in Phase 5. The best combination of parameters is summarized in Table 13.

Phase	Parameter	Value
Phase 1	Optimizer	LBFGS
Phase 1	Window size	3
Phase 2	Optimizer	Passive-Aggressive
Phase 2	Window size	3
Phase 4	Embedding vector size	32
Phase 4	CNN Kernel size	3
Phase 4	CNN filters	103
Phase 4	CNN Pooling size	4
Phase 4	LSTM Output size	75
Phase 4	Dropout	0.1
Phase 5	Algorithm	MultinomialNB
Phase 5	Separate 5W1H in parts	False

Table 13: Best combination of parameters found for the full pipeline.

After optimization, an independent test was performed on a random selection of 40 news test sets, obtaining the results summarized in Table 14. In general, the best pipeline found obtains an F_1 score of 0.74 and an accuracy score of 0.75. It obtains a larger precision on the True class and a larger recall on the Fake class, which indicates a small bias towards classifying news as Fake.

Model	True News			Fake News			Acc	Macro F_1
	P	R	F_1	P	R	F_1		
Baseline (Random)	0.551	0.549	0.548	0.498	0.500	0.497	0.526	0.522
Baseline (TF-IDF)	0.609	0.868	0.715	0.726	0.381	0.494	0.637	0.605
Full pipeline	0.920	0.550	0.790	0.680	0.950	0.690	0.750	0.740

Table 14: Full pipeline performance.

981 6.7. Cross-domain Analysis

982 In order to explore the applicability of our proposal across domains, two
 983 different experiments were performed. First, a small dataset in the political
 984 domain was created and annotated according to our annotation scheme. It
 985 contains 17 fake news and 14 true news and it was used only for testing
 986 purposes. Second, in order to be able to test the system in domains other
 987 than the political one, the two Spanish corpora available in the state of the art
 988 are studied. Since the dataset (Almela et al., 2012) is not news as such, but a
 989 dataset of opinions, the FN detection system has been tested using Posadas’
 990 Spanish dataset (Posadas-Durán et al., 2019)²⁶, which is a dataset of news
 991 websites covering different domains (Science, Sport, Economy, Education,
 992 Entertainment, Politics, Health, Security and Society). Since this corpus is
 993 only annotated with two labels (real and fake), we can not use it for training
 994 our system, only for testing it as a cross-domain experiment. Table 15 shows
 995 the results obtained given these two cross-domain scenarios.

Training	Testing	Full Pipeline				
		Acc	F_1 (True)	F_1 (Fake)	Micro F_1	Macro F_1
Health dataset	Health dataset	0.75	0.79	0.69	0.74	0.74
Health dataset	Political dataset	0.62	0.17	0.75	0.53	0.46
Health dataset	Posadas dataset	0.52	0.31	0.59	0.43	0.45

Table 15: Cross-domain analysis of the proposal

996 Not surprisingly, there is a loss in accuracy and F_1 as compared to the
 997 within-domain results shown in the first row of Table 15. Similar performance
 998 losses occurred in the literature when cross-domain is analysed (Pérez-Rosas
 999 et al., 2018) (Huang & Chen, 2020) (Hanselowski et al., 2018). Regarding
 1000 Posadas dataset, there is also a considerable loss of F_1 and accuracy in com-
 1001 parison with results obtained by the authors (Posadas-Durán et al., 2019).
 1002 One of the main causes is that Posadas’ dataset comprises documents of nine
 1003 different domains, whose vocabulary is very diverse and therefore dissimilar
 1004 to the health vocabulary on which our system is trained. Thus, it must be
 1005 considered that (Posadas-Durán et al., 2019) trained on its dataset, hence it
 1006 is expected that their results are higher.

²⁶Available at <https://github.com/jposadas/FakeNewsCorpusSpanish>

1007 The results of the cross-domain experiment show that there is still room
 1008 for improving the model in addressing the cross-domain intractability issue.
 1009 Although the system obtains reasonable accuracy results, the F_1 score drops
 1010 significantly, especially in the True class. This can be explained by consider-
 1011 ing the imbalance in terms of features in our training set, i.e., it is harder for
 1012 a news item to be classified as True, since almost any evidence of False state-
 1013 ments points to fake news. This happens because news items with both fake
 1014 and true statements are considered fake, as well as news items with only fake
 1015 statements. Hence, our model inherently learns a bias towards classifying
 1016 news items as fake, unless a sufficient number of True 5W1H components are
 1017 present. In the extreme case of having no evidence whatsoever, our model
 1018 defaults to classifying a news item as Fake. Notice that this is a sensible
 1019 default, and it is not hard-coded, but learned implicitly from the annotated
 1020 corpus. When applying our model out-of-domain, significantly less 5W1H
 1021 components are successfully extracted, since the lexical features of the other
 1022 domains differ from those where the CRF models were trained. This failure
 1023 in the earlier parts of the pipeline explains the bias towards the Fake class.

1024 7. Comparison of our proposal with the state of the art

1025 The objective of a SOTA comparison is to make a reliable comparison.
 1026 Due to the novelty and particularities of our dataset, where every essential
 1027 part of the news is detected and assigned a veracity value, and since this does
 1028 not occur in any other SOTA dataset, to the authors' knowledge, a direct
 1029 comparison of the results of the different systems published in literature on
 1030 those datasets is not possible. Nevertheless, we carried out a set of com-
 1031 parative experiments that compare our proposal with the state of the art
 1032 in three scenarios: 1) our proposal vs state-of-the-art systems, training and
 1033 testing them on our dataset; 2) our proposal's performance vs the most com-
 1034 mon method used by the SOTA approaches that use linguistic cues extracted
 1035 from LIWC for detection; and 3) our proposal configured with different state-
 1036 of-the-art fake news detection approaches that use ML or DL for each phase
 1037 of the pipeline.

1038 7.1. Our proposal vs SOTA systems

1039 To make this SOTA comparison, two outstanding works in the literature
 1040 (Pérez-Rosas et al., 2018) and (Rashkin et al., 2017) were analyzed. How-
 1041 ever, in both cases the systems were not available and have been replicated.

1042 Furthermore, another outstanding work whose code was available was also
 1043 included (Potthast et al., 2018).

1044 Regarding (Pérez-Rosas et al., 2018)’s approach, and taking into account
 1045 that their system is not available, we have replicated it considering the best
 1046 result obtained in this research. The following features have been used as
 1047 characteristics: number of characters; complex words; long words; num-
 1048 ber of syllables; word types; number of paragraphs; and readability metrics
 1049 —Flesch-Kincaid, Flesch Reading Ease, Gunning Fog, and the Automatic
 1050 Readability Index (ARI)—. Then, in line with how they describe their ex-
 1051 perimentation in their work, we have used a linear SVM classifier and five-fold
 1052 cross-validation with our English dataset.

1053 Regarding (Rashkin et al., 2017)’s approach, we replicated their DL model,
 1054 which consists of an embedding layer (using GLOVE 100-dim pre-trained
 1055 embeddings²⁷ which are fine-tuned during training, from (Pennington et al.,
 1056 2014)), an LSTM layer with 300 hidden units, and a final dense layer. The
 1057 only difference in our replication is that since our problem is binary, we apply
 1058 a sigmoid activation and binary cross-entropy loss, instead of softmax and
 1059 categorical cross-entropy, as in their original paper. Training parameters are
 1060 also replicated, i.e., 10 epochs with a batch size of 64 items. 30 independent
 1061 train/test splits were performed.

1062 The SOTA systems used in this comparison work on English datasets
 1063 where the news documents are assigned a veracity value. Therefore, to com-
 1064 pare the performance of our proposal with other SOTA fake news detection
 1065 systems, our dataset was translated into English and the three aforemen-
 1066 tioned SOTA systems were trained and tested on the translated dataset,
 1067 with a train and test configuration of 80%/20% in 30 independent evalua-
 1068 tions. The results obtained are shown in Table 16.

System	Acc	F_1 (True)	F_1 (False)	Macro- F_1
Our system	0.75	0.79	0.69	0.74
Potthast (2018)	0.66	0.63	0.69	0.66
Pérez-Rosas (2018)	0.56	0.63	0.46	0.52
Rashkin (2017)	0.53	0.46	0.55	0.51

Table 16: Comparison with SOTA systems: training and testing with our dataset

²⁷<https://nlp.stanford.edu/projects/glove/>

1069 As presented in Table 16, our proposal surpasses the other systems. Re-
 1070 garding (Potthast et al., 2018), our system obtains an improvement of 13.6%
 1071 of accuracy, 25.4% in F_1 on true news and obtains a very similar result in the
 1072 F_1 on false news. Regarding (Pérez-Rosas et al., 2018), our system obtains an
 1073 improvement of 33.45% of accuracy, 25.40% in F_1 on true news and 50.33%
 1074 in F_1 on fake news. Regarding (Rashkin et al., 2017), our system obtains an
 1075 improvement of 41.51% in accuracy, 71.73% in F_1 on true news and 25.45%
 1076 in F_1 on fake news.

1077 This SOTA comparison shows that our approach improves the results
 1078 obtained on our dataset, and that it is a robust solution. Furthermore, our
 1079 approach is more ambitious and aims to go one step further by addressing the
 1080 problem at a higher level than a simple text classification problem, whereas
 1081 these systems are acting as a black box. Hence, our goal is to give the user
 1082 the specific elements of the information that drives the system to a final
 1083 conclusion regarding the veracity of the news article.

1084 *7.2. Our proposal vs the most common method used by SOTA approaches*

1085 Considering Section 2, most of the literature’s approaches focus on study-
 1086 ing linguistic aspects of falsehood, identifying different types of linguistic fea-
 1087 tures of fake news (Zhou & Zhang, 2008) (Pérez-Rosas et al., 2018) (Almela
 1088 et al., 2012) (Afroz et al., 2012) (Shu et al., 2019) (Volkova et al., 2017).
 1089 Therefore, a comparison of our proposal with the method applied by many
 1090 different state-of-the-art approaches was performed. According to the litera-
 1091 ture, the Linguistic Inquiry and Word Count(LIWC) (Newman et al., 2003)
 1092 is widely used to extract the lexicons falling into psycholinguistic categories.
 1093 Hence, our dataset was annotated using LIWC and a set of experiments were
 1094 performed using different ML approaches with the final aim of comparing our
 1095 proposal with one of the most common state-of-the-art fake news detection
 1096 methods.

1097 For this purpose, the library AutoGOAL (Estevez-Velarde et al., 2020)
 1098 was used to search between 16 different types of ML methods (including
 1099 shallow classifiers and DL approaches) for the best algorithm and its hyper-
 1100 parameters with respect to classification accuracy. After one hour of opti-
 1101 mization, a total of 949 different variants of algorithms and parameters were
 1102 tested. Each algorithm has different optimisable parameters, such as regu-
 1103 larization factors, number of iterations, etc., which are not explicitly listed
 1104 for space reasons. In total, 72 different parameters are optimised among all
 1105 algorithms. Table 17 summarizes the results in terms of mean and standard

1106 deviation of accuracy for all of the different variants of each algorithm tested.
 1107 Each iteration consists of 30 cross-validation steps with a random 80% of the
 1108 news items for training and the remaining 20% for testing. The average ac-
 1109 curacy (micro-average) among all algorithms is 0.616, which is 15.8 percent
 1110 points below the best solution using our approach. Hence, by using LIWC
 1111 features alone, a wide range of results can be expected, ranging from 0.18 to
 1112 0.66, depending on the specific algorithms, parameters, and training used.
 1113 On average, these results do not outperform the approach presented in this
 1114 research.

Algorithm	Acc (mean)	Acc (std)	Variants
NearestCentroid	0.6657	0.1441	115
MultinomialNB	0.6394	0.1228	60
ComplementNB	0.6140	0.1241	57
NuSVC	0.5822	0.2476	88
LinearSVC	0.5646	0.2580	66
Perceptron	0.4895	0.2980	76
Neural Network	0.4833	0.0236	2
DecisionTreeClassifier	0.4571	0.0356	47
ExtraTreeClassifier	0.4490	0.0723	52
RidgeClassifier	0.4345	0.2902	69
SVM	0.4267	0.2557	55
SGDClassifier	0.4151	0.2932	54
PassiveAggressiveClassifier	0.4103	0.2571	47
KNeighborsClassifier	0.3414	0.2777	56
LogisticRegression	0.3269	0.3517	49
BernoulliNB	0.1881	0.2486	56
Our approach (full pipeline)	0.750		

Table 17: Summary of the mean and standard deviation of accuracy for the different variants of each ML algorithm trained on the LIWC characteristics.

1115 7.3. Our proposal using different SOTA fake news detection approaches for 1116 each step

1117 As presented in the background (Section 2), many SOTA systems use
 1118 different ML or DL approaches to solve the problem. In this case, to conduct
 1119 a comparison of approaches, the AutoGOAL library (see Section 6.6) is also
 1120 used to see the results of using different approaches for each of the steps in
 1121 our proposal.

1122 A total of 24 hours of computation enabled the evaluation of 101 different
 1123 combinations of algorithms applied to our proposal. For each combination,

1124 we define a set of features that correspond to specific algorithms or param-
 1125 eters used in our pipeline. Then we aggregate for each feature the accuracy
 1126 of all the pipelines in which it appeared. The same algorithm or parameter
 1127 value for a given phase appears in multiple pipelines, combined with different
 1128 options in the remaining phases. For this reason, the average accuracy of each
 1129 feature as reported is influenced by the context, i.e., by the characteristics of
 1130 the pipelines in which that feature was reported.

1131 The total number of optimisable parameters in the pipeline is 63, ranging
 1132 from numerical parameters such as the number of neurons in each layer or the
 1133 dropout rate, to categorical parameters such as which algorithm is used for
 1134 the last phase. We report only the most relevant parameters, i.e., those that
 1135 show a larger influence in the overall performance of our proposal. Figure 6
 1136 shows a graphical representation of the most relevant parameters of each
 1137 pipeline evaluated, and Table 18 summarizes the average accuracy of all
 1138 the evaluated pipelines that contained the given features. The parameters
 1139 reported are the following:

- 1140 • Optimization algorithm used in the CRF taggers (Phases 1 and 2).
- 1141 • Window size of the CRF taggers (Phases 1 and 2).
- 1142 • ML algorithm used in Phase 5.
- 1143 • Whether to aggregate 5W1H components by body part in Phase 5.

1144 As can be observed, both the algorithms used in Phase 2 & 3 as well
 1145 as those in Phase 5 have a significant impact on the overall accuracy of the
 1146 pipelines. The most consistent algorithm for the CRF components is Passive
 1147 Aggressive. For Phase 5, even if the algorithm that produces on average
 1148 the best performance is Stochastic Gradient Descent, the most consistent
 1149 option is Multinomial Naive Bayes, which is also the algorithm selected in
 1150 the best performing pipeline (see Section 6.6). The window size in the CRF
 1151 components is also a significant factor, as shown in Figure 6, since pipelines
 1152 with a larger window (size=3) consistently perform better than those with a
 1153 shorter window. This is an expected result since a larger window allows more
 1154 tokens to be considered as part of the context for a specific token. Finally,
 1155 it is interesting to note that aggregating all 5W1H components instead of
 1156 counting them within the part of the article in which they appear increases
 1157 the average accuracy by more than 3 percent points. This has the effect of

Phase	Feature	Value	Acc (mean)	Acc (std)	Variants
Phase 2 & 3	algorithm	Averaged Perceptron	0.5839	0.0923	28
Phase 2 & 3	algorithm	Adaptive Regularization	0.5455	0.0861	28
Phase 2 & 3	algorithm	Stochastic Gradient Descent	0.5875	0.0919	22
Phase 2 & 3	algorithm	L-BFGS	0.5921	0.0884	38
Phase 2 & 3	algorithm	Passive Aggressive	0.6072	0.0930	69
Phase 2 & 3	window size	0	0.5857	0.0859	70
Phase 2 & 3	window size	1	0.5960	0.0989	25
Phase 2 & 3	window size	2	0.5543	0.0964	29
Phase 2 & 3	window size	3	0.6061	0.0915	61
Phase 5	algorithm	BernoulliNB	0.5792	0.1145	6
Phase 5	algorithm	CategoricalNB	0.6031	0.0930	8
Phase 5	algorithm	ComplementNB	0.6500	0.1061	2
Phase 5	algorithm	DecisionTreeClassifier	0.5156	0.0566	8
Phase 5	algorithm	ExtraTreeClassifier	0.5528	0.0292	9
Phase 5	algorithm	GaussianNB	0.5350	0.0742	5
Phase 5	algorithm	KNeighborsClassifier	0.5788	0.0957	20
Phase 5	algorithm	MultinomialNB	0.6041	0.0953	37
Phase 5	algorithm	NuSVC	0.5000	0.0354	2
Phase 5	algorithm	SGDClassifier	0.7250	0.0707	2
Phase 5	algorithm	SVC	0.6000	0.0354	2
Phase 5	use parts	False	0.5978	0.0911	56
Phase 5	use parts	True	0.5661	0.0898	45

Table 18: Summary of the performance associated to the most relevant parameters of each pipeline evaluated in the AutoML process.

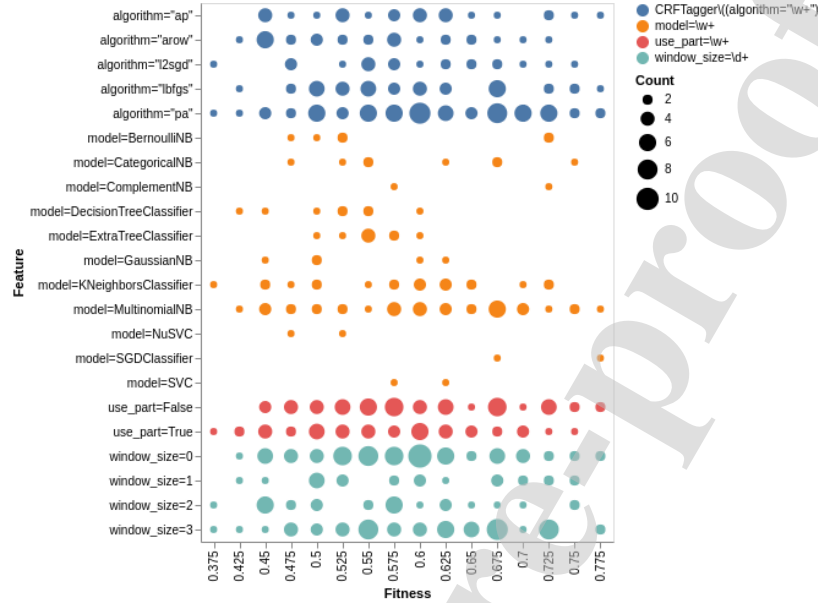


Figure 6: Graphical visualization of the most relevant parameters of each pipeline evaluated.

1158 reducing the total number of features in Phase 5, which could help alleviate
 1159 the impact of a reduced training set.

1160 8. Conclusions and further work

1161 This paper presents a novel approach to dealing with automatic fake news
 1162 detection on traditional digital media and our proposal is based on the
 1163 premise that fake news combines true and false data with the intention of
 1164 confusing readers. However, to the best of the authors' knowledge, current
 1165 datasets consider the news as a whole and assign it a single truthfulness
 1166 value, although this truthfulness value may have degrees of certainty, but
 1167 they are not determining specifically which parts within the news item are
 1168 true and which parts are false or even unverifiable.

1169 Our proposal exploits the journalistic structure of news articles and how
 1170 the content is presented, following the inverted pyramid hypothesis. More-
 1171 over, the essential content of the news is typically presented by answering
 1172 six questions that comprise the 5W1H. Based on this knowledge, a new fine-
 1173 grained annotation scheme (FNDeepML) is defined using two levels of rep-

1174 resentation: i) Newspaper article structure and ii) Essential news content
 1175 (5W1H). A new dataset in Spanish is created consisting of 200 news articles
 1176 focused on the health domain, and specifically on COVID-19 news.

1177 The proposed architecture comprises two main layers —Structure and
 1178 Veracity Layers— that predict not only the article’s veracity value but also
 1179 that of the article’s main content. The experiments have demonstrated that
 1180 the use of the veracity value of the different structural elements and that
 1181 of the 5W1H essential content within the news provides a suitable solution
 1182 to the problem. The best performance for the Veracity Layer was obtained
 1183 with a Logistic Regression model, resulting in a $F_1=0.807$, compared to a
 1184 baseline using the TF-IDF of the entire document —annotated with a unique
 1185 veracity value for the document— resulting in $F_1=0.60$. Furthermore, the
 1186 performance of the Veracity Layer using the veracity of gold standard 5W1H
 1187 components increases to $F_1=0.978$. These findings demonstrate the validity
 1188 of our proposal.

1189 Our experiments also demonstrate that determining the veracity of each
 1190 5W1H component using only textual information has a limited prediction
 1191 performance, and therefore, adding high-level features (i.e. fact-checking
 1192 information, semantic relations between components, contextual features,
 1193 among others) would be beneficial. The future goal is to predict as accurately
 1194 as possible the veracity of each essential element of the news, as this would
 1195 be a very powerful tool for readers, who would benefit from a detailed report
 1196 on the reliability of news content elements.

1197 At this stage of the research, the news elements and the news document
 1198 are classified only in True/False/Unknown categories. However, in future
 1199 developments, a weighting of elements may be used to determine an over-
 1200 all degree of veracity depending on the Fs and Ts items detected in the
 1201 same block of text. Furthermore, in future work phases 1 and 2 should be
 1202 enhanced. On the other hand, phase 3 would require an improvement in
 1203 the automatic fact-checking tool by means of determining the semantic rela-
 1204 tionship between the different 5W1H elements to provide a context to those
 1205 items. This contribution would enable the detection of contradictions in the
 1206 5W1H relations, which may be indicative of fake information.

1207 Acknowledgements

1208 This research work has been partially funded by Generalitat Valenciana
 1209 through project “SIIA: Tecnologías del lenguaje humano para una sociedad

1210 inclusiva, igualitaria, y accesible” with grant reference PROMETEU/2018/089,
 1211 by the Spanish Government through project RTI2018-094653-B-C22: “Mod-
 1212 elang: Modeling the behavior of digital entities by Human Language Tech-
 1213 nologies”, as well as being partially supported by a grant from the Fondo
 1214 Europeo de Desarrollo Regional (FEDER) and the LIVING-LANG project
 1215 (RTI2018-094653-B-C21) from the Spanish Government. Furthermore, we
 1216 would like to thank Difusión Comunicación, especially Tono Jordá and Lara
 1217 Sánchez Belda, and Newtral for their collaboration in collection and annota-
 1218 tion of datasets.

1219 References

- 1220 Afroz, S., Brennan, M., & Greenstadt, R. (2012). Detecting hoaxes, frauds,
 1221 and deception in writing style online. In *Proceedings - 2012 IEEE Sym-*
 1222 *posium on Security and Privacy, S and P 2012* Proceedings - IEEE Sym-
 1223 *posium on Security and Privacy* (pp. 461–475). Institute of Electrical and
 1224 *Electronics Engineers Inc.* doi:10.1109/SP.2012.34 33rd IEEE Symposium
 1225 *on Security and Privacy, S and P 2012* ; Conference date: 21-05-2012
 1226 *Through 23-05-2012.*
- 1227 Agarwal, A., & Dixit, A. (2020). Fake news detection: An ensemble learning
 1228 approach. In *2020 4th International Conference on Intelligent Computing*
 1229 *and Control Systems (ICICCS)* (pp. 1178–1183).
- 1230 Almela, A., Valencia-García, R., & Cantos, P. (2012). Seeing through de-
 1231 ception: A computational approach to deceit detection in written commu-
 1232 nication. In *Proceedings of the Workshop on Computational Approaches*
 1233 *to Deception Detection EACL 2012* (pp. 15–22). Stroudsburg, PA, USA:
 1234 *Association for Computational Linguistics.*
- 1235 Bednarek, M., & Caple, H. (2012). *News discourse* volume 46. A&C Black.
- 1236 Bonet-Jover, A., Saquete, E., Martínez-Barco, P., & Ángel
 1237 *García-Cumbreras, M.* (2020a). Fnddeepml annotation guide-
 1238 lines. URL: <https://doi.org/10.5281/zenodo.4091549>.
 1239 doi:10.5281/zenodo.4091549.
- 1240 Bonet-Jover, A., Saquete, E., Nieto, M., Belén, V. M., Piad-Morffis,
 1241 A., Martínez-Barco, P., & Ángel García-Cumbreras, M. (2020b).

- 1242 Fndeep dataset. URL: <https://doi.org/10.5281/zenodo.4090914>.
1243 doi:10.5281/zenodo.4090914.
- 1244 Brennan, M. R., & Greenstadt, R. (2009). Practical at-
1245 tacks against authorship recognition techniques. In K. Z.
1246 Haigh, & N. Rychtyckyj (Eds.), *IAAI*. AAAI. URL:
1247 <http://dblp.uni-trier.de/db/conf/iaai/iaai2009.html#BrennanG09>.
- 1248 Chagas, L. J. (2019). The spiral model in the text of live radio journalism.
1249 *Journal of Radio & Audio Media*, *26*, 231–246.
- 1250 Chakma, K., & Das, A. (2018a). A 5w1h based annotation scheme for seman-
1251 tic role labeling of english tweets. *Computacion y Sistemas*, *22*, 747–755.
1252 doi:10.13053/CyS-22-3-3016.
- 1253 Chakma, K., & Das, A. (2018b). A 5w1h based annotation scheme for se-
1254 mantic role labeling of english tweets. *Computación y Sistemas*, *22*.
- 1255 Chen, Y., & Chen, H. (2014). Opinion spam detection in web forum: A real
1256 case study. In *Proceedings of the 24th International Conference on World*
1257 *Wide Web* (pp. 173–183).
- 1258 Ciampaglia, G. L., Shiralkar, P., Rocha, L. M., Bollen, J., Menczer, F.,
1259 & Flammini, A. (2015). Computational fact checking from knowledge
1260 networks. *PLOS ONE*, *10*, 1–13. doi:10.1371/journal.pone.0128193.
- 1261 Cohen, J. (1960). A Coefficient of Agreement for Nominal Scales. *Educational*
1262 *and Psychological Measurement*, *20*, 37.
- 1263 Conroy, N. J., Rubin, V. L., & Chen, Y. (2015). Automatic deception detec-
1264 tion: Methods for finding fake news. In *Proceedings of the 78th ASIS&T*
1265 *Annual Meeting: Information Science with Impact: Research in and for*
1266 *the Community ASIST '15* (pp. 82:1–82:4). Silver Springs, MD, USA:
1267 American Society for Information Science.
- 1268 Dale, R. (2017). Nlp in a post-truth world. *Natural Language Engineering*,
1269 *23*, 319–324. doi:10.1017/S1351324917000018.
- 1270 Das Bhattacharjee, S., Talukder, A., & Balantrapu, B. (2017). Active learn-
1271 ing based news veracity detection with feature weighting and deep-shallow
1272 fusion. doi:10.1109/BigData.2017.8257971.

- 1273 Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-
1274 training of deep bidirectional transformers for language understanding.
- 1275 Estevez-Velarde, S., Piad-Morffis, A., Gutiérrez, Y., Montoyo, A., Munoz,
1276 R., & Almeida-Cruz, Y. (2020). Solving Heterogeneous AutoML Problems
1277 with AutoGOAL. In *Proceedings of the 7th ICML Workshop on Automated
1278 Machine Learning*. ICML.
- 1279 Faustini, P., & Covoos, T. (2020). Fake news detection in multiple plat-
1280 forms and languages. *Expert Systems with Applications*, (p. 113503).
1281 doi:10.1016/j.eswa.2020.113503.
- 1282 Ferreira, W., & Vlachos, A. (2016). Emergent: a novel data-set for stance
1283 classification. In *Proceedings of the Conference of the North American
1284 Chapter of the Association for Computational Linguistics* (pp. 1163–1168).
1285 Association for Computational Linguistics. doi:10.18653/v1/N16-1138.
- 1286 Gilda, S. (2017). Evaluating machine learning algorithms for fake news detec-
1287 tion. In *2017 IEEE 15th Student Conference on Research and Development
1288 (SCORED)* (pp. 110–115). doi:10.1109/SCORED.2017.8305411.
- 1289 Gravanis, G., Vakali, A., Diamantaras, K., & Karadais, P. (2019). Behind
1290 the cues: A benchmarking study for fake news detection. *Expert Systems
1291 with Applications*, 128. doi:10.1016/j.eswa.2019.03.036.
- 1292 Hamborg, F., Breiterger, C., Schubotz, M., Lachnit, S., & Gipp, B.
1293 (2018). Extraction of main event descriptors from news articles by an-
1294 swering the journalistic five w and one h questions. In *Proceedings of the
1295 18th ACM/IEEE on Joint Conference on Digital Libraries JCDL '18* (p.
1296 339–340). New York, NY, USA: Association for Computing Machinery.
1297 doi:10.1145/3197026.3203899.
- 1298 Han, S., Lee, K., Lee, D., & Lee, G. G. (2013). Counseling dialog system
1299 with 5W1H extraction. In *Proceedings of the SIGDIAL 2013 Conference*
1300 (pp. 349–353). Metz, France: Association for Computational Linguistics.
- 1301 Hanselowski, A., PVS, A., Schiller, B., Caspelherr, F., Chaudhuri,
1302 D., Meyer, C. M., & Gurevykh, I. (2018). A retrospective anal-
1303 ysis of the fake news challenge stance-detection task. In *Proceed-
1304 ings of the 27th International Conference on Computational Linguis-*

- 1305 *tics* (pp. 1859–1874). Association for Computational Linguistics. URL:
1306 <https://www.aclweb.org/anthology/C18-1158>.
- 1307 Huang, Y.-F., & Chen, P.-H. (2020). Fake news detection using an ensemble
1308 learning model based on self-adaptive harmony search algorithms. *Expert*
1309 *Systems with Applications*, *159*, 113584. doi:10.1016/j.eswa.2020.113584.
- 1310 Khan, S. U. R., Islam, M. A., Aleem, M., Iqbal, M. A., & Ahmed, U. (2018).
1311 Section-based focus time estimation of news articles. *IEEE Access*, *6*,
1312 75452–75460.
- 1313 Kim, J.-D., Son, J., & Baik, D.-K. (2012). Ca 5w1h onto: Ontological
1314 context-aware model based on 5w1h. *International Journal of Distributed*
1315 *Sensor Networks*, *2012*. doi:10.1155/2012/247346.
- 1316 Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic opti-
1317 mization. URL: <http://arxiv.org/abs/1412.6980> cite arxiv:1412.6980
1318 Comment: Published as a conference paper at the 3rd International Con-
1319 ference for Learning Representations, San Diego, 2015.
- 1320 Lewandowsky, S., Ecker, U. K. H., Seifert, C. M., Schwarz, N., & Cook,
1321 J. (2012). Misinformation and its correction: Continued influence and
1322 successful debiasing. *Psychological Science in the Public Interest*, *13*, 106–
1323 131. doi:10.1177/1529100612451018.
- 1324 Mihalcea, R., & Strapparava, C. (2009). The lie detector: explorations in the
1325 automatic recognition of deceptive language. In *Proceedings of the ACL-*
1326 *IJCNLP 2009 Conference* (pp. 309–312). Association for Computational
1327 Linguistics.
- 1328 Monti, F., Frasca, F., Eynard, D., Mannion, D., & Bronstein, M. M. (2019).
1329 Fake news detection on social media using geometric deep learning. *CoRR*,
1330 *abs/1902.06673*. URL: <http://arxiv.org/abs/1902.06673>.
- 1331 Moreda, P., Llorens, H., Saquete, E., & Palomar, M. (2011). Combining
1332 semantic information in question answering systems. *Inf. Process. Manag.*,
1333 *47*, 870–885. doi:10.1016/j.ipm.2010.03.008.
- 1334 Mosteller, F., & Wallace, D. L. (1963). Inference in an authorship prob-
1335 lem. *Journal of the American Statistical Association*, *58*, 275–309.
1336 doi:10.1080/01621459.1963.10500849.

- 1337 Newman, M. L., Pennebaker, J. W., Berry, D. S., & Richards, J. M. (2003).
1338 Lying words: Predicting deception from linguistic styles. *Personality and*
1339 *Social Psychology Bulletin*, *29*, 665 – 675.
- 1340 NIST (2011). Tac 2011 guided summarization task guidelines. URL:
1341 [https://tac.nist.gov/2011/Summarization/Guided-Summ.2011-](https://tac.nist.gov/2011/Summarization/Guided-Summ.2011-.guidelines.html)
1342 [.guidelines.html](https://tac.nist.gov/2011/Summarization/Guided-Summ.2011-.guidelines.html).
- 1343 Norambuena, B., Horning, M., & Mitra, T. (2020). Evaluating
1344 the inverted pyramid structure through automatic 5w1h extraction
1345 and summarization. *Computational Journalism Symposium*, . URL:
1346 <http://par.nsf.gov/biblio/10168974>.
- 1347 Nyhan, B., Reifler, J., & Ubel, P. (2012). The hazards of
1348 correcting myths about health care reform. *Medical care*, *51*.
1349 doi:10.1097/MLR.0b013e318279486b.
- 1350 Padró, L., & Stanilovsky, E. (2012). Freeling 3.0: Towards wider multilingual-
1351 ity. In *Proceedings of the Language Resources and Evaluation Conference*
1352 *(LREC 2012)*. Istanbul, Turkey: ELRA.
- 1353 Pennington, J., Socher, R., & Manning, C. (2014). Glove: Global Vectors
1354 for Word Representation. In *Proc. 2014 Conf. Empir. Methods Nat. Lang.*
1355 *Process.* (pp. 1532–1543). doi:10.3115/v1/D14-1162.
- 1356 Pérez-Rosas, V., Kleinberg, B., Lefevre, A., & Mihalcea, R. (2018). Auto-
1357 matic detection of fake news. In *Proceedings of the 27th International*
1358 *Conference on Computational Linguistics* (pp. 3391–3401). Santa Fe,
1359 New Mexico, USA: Association for Computational Linguistics. URL:
1360 <https://www.aclweb.org/anthology/C18-1287>.
- 1361 Posadas-Durán, J., Gomez-Adorno, H., Sidorov, G., & Escobar, J. (2019).
1362 Detection of fake news in a new corpus for the spanish language. *Journal of*
1363 *Intelligent and Fuzzy Systems*, *36*, 4868–4876. doi:10.3233/JIFS-179034.
- 1364 Potthast, M., Kiesel, J., Reinartz, K., Bevendorff, J., & Stein, B.
1365 (2018). A stylometric inquiry into hyperpartisan and fake news. In
1366 *Proceedings of the 56th Annual Meeting of the Association for Com-*
1367 *putational Linguistics (Volume 1: Long Papers)* (pp. 231–240). Mel-
1368 bourne, Australia: Association for Computational Linguistics. URL:

- 1369 <https://www.aclweb.org/anthology/P18-1022>. doi:10.18653/v1/P18-
1370 1022.
- 1371 Prechelt, L. (1998). Early stopping-but when? In *Neural Networks: Tricks*
1372 *of the trade* (pp. 55–69). Springer.
- 1373 Rashkin, H., Choi, E., Jang, J. Y., Volkova, S., & Choi, Y. (2017).
1374 Truth of varying shades: Analyzing language in fake news and polit-
1375 ical fact-checking. In *Proceedings of the 2017 Conference on Empiri-*
1376 *cal Methods in Natural Language Processing* (pp. 2931–2937). Copen-
1377 hagen, Denmark: Association for Computational Linguistics. URL:
1378 <https://www.aclweb.org/anthology/D17-1317>. doi:10.18653/v1/D17-
1379 1317.
- 1380 Ren, Y., & Ji, D. (2017). Neural networks for deceptive opinion spam de-
1381 tection: An empirical study. *Information Sciences*, 385-386, 213 – 224.
1382 doi:<https://doi.org/10.1016/j.ins.2017.01.015>.
- 1383 Rubin, V., Conroy, N., Chen, Y., & Cornwell, S. (2016). Fake news or truth?
1384 using satirical cues to detect potentially misleading news. In *Proceedings of*
1385 *the Second Workshop on Computational Approaches to Deception Detec-*
1386 *tion* (pp. 7–17). San Diego, California: Association for Computational
1387 Linguistics. URL: <https://www.aclweb.org/anthology/W16-0802>.
1388 doi:10.18653/v1/W16-0802.
- 1389 Saquete, E., Tomás, D., Moreda, P., Martínez-Barco, P., & Palomar, M.
1390 (2020). Fighting post-truth using natural language processing: A review
1391 and open challenges. *Expert Syst. Appl.*, 141.
- 1392 Shu, K., Wang, S., & Liu, H. (2019). Beyond news contents: The role of
1393 social context for fake news detection. In *WSDM 2019 - Proceedings of*
1394 *the 12th ACM International Conference on Web Search and Data Mining*
1395 *WSDM 2019 - Proceedings of the 12th ACM International Conference on*
1396 *Web Search and Data Mining* (pp. 312–320). Association for Comput-
1397 ing Machinery, Inc. doi:10.1145/3289600.3290994 12th ACM International
1398 Conference on Web Search and Data Mining, WSDM 2019 ; Conference
1399 date: 11-02-2019 Through 15-02-2019.
- 1400 Silva, R. M., Santos, R. L. S., Almeida, T. A., & Pardo, T. A. S. (2020). To-
1401 wards automatically filtering fake news in portuguese. *Expert Syst. Appl.*,
1402 146, 113199.

- 1403 Sutton, C., McCallum, A. et al. (2012). An introduction to conditional
1404 random fields. *Foundations and Trends® in Machine Learning*, 4, 267–
1405 373.
- 1406 Thomson, E. A., White, P. R., & Kitley, P. (2008). “objectivity” and “hard
1407 news” reporting across cultures: Comparing the news report in english,
1408 french, japanese and indonesian journalism. *Journalism studies*, 9, 212–
1409 228.
- 1410 Verma, A., Mittal, V., & Dawn, S. (2019). Find: Fake information and news
1411 detections using deep learning. In *2019 Twelfth International Conference
1412 on Contemporary Computing (IC3)* (pp. 1–7).
- 1413 Vlachos, A., & Riedel, S. (2014). Fact checking: Task definition and dataset
1414 construction. In *Proceedings of the ACL 2014 Workshop on Language Tech-
1415 nologies and Computational Social Science* (pp. 18–22). Baltimore, MD,
1416 USA: Association for Computational Linguistics. doi:10.3115/v1/W14-
1417 2508.
- 1418 Volkova, S., Shaffer, K., Jang, J. Y., & Hodas, N. (2017). Separating
1419 facts from fiction: Linguistic models to classify suspicious and trusted
1420 news posts on twitter. In *Proceedings of the 55th Annual Meeting of
1421 the Association for Computational Linguistics (Volume 2: Short Pa-
1422 pers)* (pp. 647–653). Vancouver, Canada: Association for Computational
1423 Linguistics. URL: <https://www.aclweb.org/anthology/P17-2102>.
1424 doi:10.18653/v1/P17-2102.
- 1425 Voorhees, E. (2001). The trec question answering track. *Nat. Lang. Eng.*, 7,
1426 361–378. doi:10.1017/S1351324901002789.
- 1427 Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news
1428 online. *Science*, 359, 1146–1151. doi:10.1126/science.aap9559.
- 1429 Wang, W. (2012). Chinese news event 5w1h semantic elements extrac-
1430 tion for event ontology population. In *Proceedings of the 21st Inter-
1431 national Conference on World Wide Web WWW '12 Companion* (p.
1432 197–202). New York, NY, USA: Association for Computing Machinery.
1433 doi:10.1145/2187980.2188008.
- 1434 Wang, W., Zhao, D., Zou, L., Wang, D., & Zheng, W. (2010). Extract-
1435 ing 5w1h event semantic elements from chinese online news. In L. Chen,

- 1436 C. Tang, J. Yang, & Y. Gao (Eds.), *Web-Age Information Management*
1437 (pp. 644–655). Berlin, Heidelberg: Springer Berlin Heidelberg.
- 1438 Wang, W. Y. (2017). "liar, liar pants on fire": A new benchmark dataset for
1439 fake news detection. *CoRR*, *abs/1705.00648*.
- 1440 Zhang, H., & Liu, H. (2016). Visualizing structural "inverted pyramids" in
1441 english news discourse across levels. *Text & Talk*, *36*, 89–110.
- 1442 Zhou, L., & Zhang, D. (2008). Following linguistic footprints: Automatic de-
1443 ception detection in online communication. *Communications of the ACM*,
1444 *51*, 119–122. doi:10.1145/1378727.1389972.

- A novel Automatic Fake News detection proposal based on determining the veracity of the essential content of news articles.
- A new benchmark Spanish Fake News dataset focused on health news is presented
- A new Fake News Detection architecture comprising two layers (Structure Layer and Veracity Layer) is presented
- Each layer of the architecture involves a set of phases and each phase is thoroughly described
- Performance of each layer of the architecture is measured and analysed

Conceptualization	Alba Bonet-Jover, Estela Saquete, Patricio Martínez-Barco
Methodology	Estela Saquete, Patricio Martínez-Barco, Alejandro Piad-Morffis, Miguel Angel García-Cumbreras
Software	Alejandro Piad-Morffis, Miguel Angel García-Cumbreras
Validation	Alejandro Piad-Morffis, Miguel Angel García-Cumbreras
Formal analysis	Alba Bonet-Jover, Estela Saquete, Patricio Martínez-Barco
Investigation	Estela Saquete, Patricio Martínez-Barco, Alejandro Piad-Morffis, Miguel Angel García-Cumbreras, Alba Bonet-Jover
Resources	Alba Bonet-Jover, Miguel Angel García-Cumbreras
Data Curation	Alba Bonet-Jover
Writing - Original Draft	Alba Bonet-Jover, Alejandro Piad-Morffis, Estela Saquete
Writing - Review & Editing	Patricio Martínez-Barco, Miguel Angel García-Cumbreras
Visualization	Estela Saquete, Alejandro Piad-Morffis
Supervision	Patricio Martínez-Barco, Estela Saquete
Project administration	Patricio Martínez-Barco, Estela Saquete
Funding acquisition	Patricio Martínez-Barco, Estela Saquete, Miguel Angel García-Cumbreras

Declaration of interests

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

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proof