

DisTrack: a new Tool for Semi-automatic Misinformation Tracking in Online Social Networks

Guillermo Villar-Rodríguez^a (guillermo.villar@upm.es), Álvaro
Huertas-García^a (alvaro.huertas@upm.es), Alejandro Martín^a
(alejandro.martin@upm.es), Javier Huertas-Tato^a
(javier.huertas.tato@upm.es) and David Camacho^a
(david.camacho@upm.es)

^a Departamento de Sistemas Informáticos. Universidad Politécnica de Madrid,
España.

Corresponding Author:

Alejandro Martín

Departamento de Sistemas Informáticos, Universidad Politécnica de Madrid,
Spain. Campus Sur de la UPM, C. de Alan Turing, s/n, 28031 Madrid

Email: alejandro.martin@upm.es

DisTrack: a new Tool for Semi-automatic Misinformation Tracking in Online Social Networks

Guillermo Villar-Rodríguez^a, Álvaro Huertas-García^a, Alejandro Martín^{a,*},
Javier Huertas-Tato^a, David Camacho^a

^a*Departamento de Sistemas Informáticos. Universidad Politécnica de Madrid, Spain*

Abstract

Introduction: This article introduces DisTrack, a methodology and a tool developed for tracking and analyzing misinformation within Online Social Networks (OSNs). DisTrack is designed to combat the spread of misinformation through a combination of Natural Language Processing (NLP) Social Network Analysis (SNA) and graph visualization. The primary goal is to detect misinformation, track its propagation, identify its sources, and assess the influence of various actors within the network.

Methods: DisTrack’s architecture incorporates a variety of methodologies including keyword search, semantic similarity assessments, and graph generation techniques. These methods collectively facilitate the monitoring of misinformation, the categorization of content based on alignment with known false claims, and the visualization of dissemination cascades through detailed graphs. The tool is tailored to capture and analyze the dynamic nature of misinformation spread in digital environments.

Results: The effectiveness of DisTrack is demonstrated through three case studies focused on different themes: discredit/hate speech, anti-vaccine misinformation, and false narratives about the Russia-Ukraine conflict. These studies show DisTrack’s capabilities in distinguishing posts that propagate falsehoods from those that counteract them, and tracing the evolution of misinformation from its inception.

*Corresponding author

Email addresses: guillermo.villar@upm.es (Guillermo Villar-Rodríguez), alvaro.huertas@upm.es (Álvaro Huertas-García), alejandro.martin@upm.es (Alejandro Martín), javier.huertas.tato@upm.es (Javier Huertas-Tato), david.camacho@upm.es (David Camacho)

Conclusions: The research confirms that DisTrack is a valuable tool in the field of misinformation analysis. It effectively distinguishes between different types of misinformation and traces their development over time. By providing a comprehensive approach to understanding and combating misinformation in digital spaces, DisTrack proves to be an essential asset for researchers and practitioners working to mitigate the impact of false information in online social environments.

Keywords: Misinformation, Transformers, COVID-19, Hoax, Natural Language Inference, Semantic Similarity

1. Introduction

At present, the impact of misinformation on our societies is beyond question. Governments, companies, researchers, among many others, are putting a great deal of effort into providing solutions and limiting the damage caused by this problem. The challenge, however, lies in the multifaceted nature of misinformation, making a comprehensive approach to combat it elusive. Addressing misinformation is contingent upon numerous factors, notably the underlying intent. Thus, when we use the word “misinformation”, we are referring to information that is false by definition or unintentionally false information [1, 2], whereas “disinformation” is used when there is a deliberate intention to misinform.[3], following the differences underlined by Karlova and Fisher in 2013 [4]. In addition to these terms, there is the concept of “malinformation”, which groups every content created as a weapon [5, 6].

In order to understand the phenomenon of misinformation in the current era, and beyond intentionality, it is essential to include the means of transmission in the equation. As such, social media, particularly Online Social Networks (OSNs), has risen as society’s primary channel for accessing information. Social media has transformed information dissemination from direct peer-to-peer exchanges to expansive many-to-many propagations, altering societal content consumption habits and thereby easing the spread of false information [7]. Literature has demonstrated with specific cases how social media fosters falsehoods easily [8].

The combined effects of OSNs’ dominance as information sources and the growing disenchantment with traditional news media amplify concerns about misinformation’s influence on the public. According to the Reuters Institute Digital News Report, there’s a notable decline in the practice of

27 consulting a variety of news outlets and a general waning interest in news
28 over the years. This trend includes the avoidance of topics like climate change
29 and the Ukraine invasion by specific segments of the audience. This study
30 underlines the contrast between worrying about misinformation and then
31 depending more on social media to receive information [9].

32 Despite the lack of expert consensus on the recent escalation of misin-
33 formation [7], academia has witnessed a significant surge in research articles
34 tackling this issue, reflecting its growing impact [2, 10]. In the years before
35 the pandemic, these academic works had already increased exponentially [10].
36 Subsequent studies associate the surge in academic interest regarding misin-
37 formation, particularly from 2017, with events like the 2016 US elections [11],
38 suggesting these events catalyzed but did not solely trigger the focus [12].
39 After these contributions to research, the emergence of more relevant interna-
40 tional events and their originated misinformation highlights the importance
41 of tackling this problem.

42 In the fight against this phenomenon, fact-checkers represent the foremost
43 defense. They are responsible for verifying thousands of pieces of information
44 daily, cross-checking them and issuing statements on social media to counter
45 the spread of false information. The International Fact-Checking Network
46 (IFCN) built a “Code of Principles” to establish the criteria for the task
47 of debunking misinformation. During the coronavirus crisis, the then Asso-
48 ciate Director Cristina Tardáguila highlighted an unprecedented volume of
49 misinformation, presenting a novel challenge for fact-checkers compared to
50 prior years [13]. Recent research on COVID-19 misinformation has shown
51 the emergence of misinformation waves, with a plethora of hoaxes spreading
52 rapidly, complicating efforts to address them all at once [14]. Evidence like
53 this demands new automated mechanisms not only at the level of verifica-
54 tion but also in the following steps of mitigation on OSNs for a coordinated
55 response against emerging falsehoods.

56 Mitigating the problem of misinformation therefore requires a complex
57 approach where all factors of the problem are considered. In addition to the
58 verification of the content itself, it is necessary to take into account how it
59 is spread, the actors involved in its distribution and, in general, any element
60 that participates in the cascade of misinformation propagation. In Fig. 1
61 the described cascade is represented, where several actors interact over time
62 about a given falsehood. We want to find and appropriately characterize
63 this misinformation cascade, allowing for a better understanding of online
64 discourse. As a result of considering all these factors, we draw the following

65 research question for our work: “*Is it possible to track conversations around*
 66 *specific hoaxes on Twitter (X)?*” which we sub-divide into four subquestions
 67 to answer it:

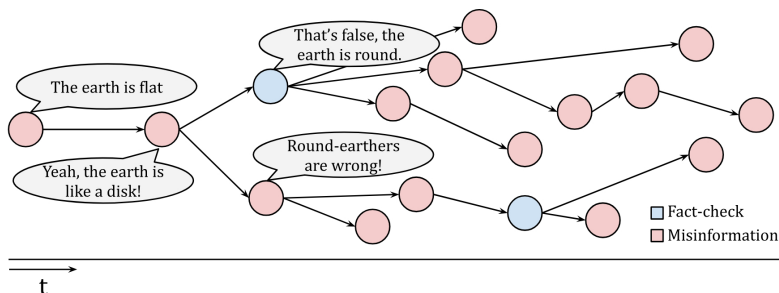


Figure 1: Visualization of the misinformation cascade. The t axis represents time from left to right, vertices are claims made by actors in any OSN, while edges represent a relation whether implicit (semantic similarity) or explicit (a repost of another piece of content).

- 68 1. Can we extract the conversation about a hoax on Twitter (X)?
 69 2. Can we separate tweets related to the hoax in the extracted conversa-
 70 tion from tweets not related to it?
 71 3. Can we distinguish between hoaxes that propagate a hoax from those
 72 that contradict it?
 73 4. Can we trace the movement of tweets related to a hoax with the users
 74 that spread them from beginning to end?

75 We integrate cutting-edge Language Models and graph generation tech-
 76 niques to explore content propagation within Online Social Networks (OSNs)
 77 through the lens of Social Network Analysis (SNA). The social platform cho-
 78 sen for this objective is Twitter (X), given its dominant use for information:
 79 in proportion, it is more used than other social media for consuming main-
 80 stream news but also smaller or alternative news sources, and politicians and
 81 political activists, too [9]. Furthermore, Twitter surpassed the other networks
 82 in this survey for the consumption of news about politics in the country and
 83 the Russia-Ukraine war, among others. This suggests the impact of false
 84 information and polarization could be disorders of these informative and po-
 85 litical interests in this ecosystem and, thus, there would be a need to avoid
 86 them.

87 This research enables not only the detection of misinformation but also
88 the tracking of its journey from inception to conclusion within a Twitter (X)
89 conversation, distinguishing between posts that propagate the falsehood and
90 those that refute it. Our architecture is designed with two primary objectives:
91 first, to **combat misinformation** by doing more than just verifying the
92 veracity of a claim—we aim to identify every tweet that supports it; second,
93 to **monitor and trace the sources of misinformation**. In the context
94 of X, these sources of misinformation will be the tweets stating a false claim,
95 not only the first tweets but all tweets that contribute to this, since all of
96 them harvest this misinformation in the circles they belong to, whether the
97 users inside them continue sharing it or not. This monitoring enables the
98 analysis of the influence of various tweets and the identities of the individuals
99 disseminating them.

100 Our methodology leverages a semi-automated fact-checking architecture [14]
101 that avoids the pitfalls of relying solely on a fully automated system. The
102 traditional automated method, which categorizes information as true or false
103 based on a pre-defined database, suffers from two major issues: it relies on a
104 static collection of data that may not keep pace with emerging hoaxes, and
105 it tends to categorize new, unverified information based on existing patterns,
106 possibly mislabeling them. In contrast, our semi-automated strategy miti-
107 gates these issues by utilizing a dynamic database curated by fact-checkers
108 that specifically targets false claims and can be regularly updated with new
109 misinformation. Additionally, this method prioritizes the evaluation of claims
110 in the database that align with the content of new information being assessed,
111 ensuring a more accurate and timely verification process.

112 This integration of fact-checking organizations with AI techniques enables
113 the automated identification of posts that are part of a hoax, as flagged in
114 the fact-checkers' database. Distrack's semi-automated method advances
115 Social Network Analysis by innovatively creating graphs that categorize the
116 elements that create a cascade of information on social media according to
117 their relationship with misinformation or with its debunking. This approach
118 refines the analysis of viral content flows by filtering out irrelevant data,
119 a significant departure from prior methodologies that did not distinguish
120 between relevant and unrelated content.

121 This article presents the following contributions:

- 122 • **An integration of NLP and SNA to combat misinformation.**

123 The study integrates cutting-edge language models and graph gener-

124 ation techniques to explore content propagation within online social
125 networks (OSNs), specifically focusing on Twitter.

- 126 • **A graph generation architecture for tracking misinformation.**
127 The tool creates graphs that categorize the elements of information cas-
128 cades on social media into three types: supporting a specific database
129 claim, contradicting it, or being unrelated. This approach refines the
130 analysis of viral content flows by filtering out irrelevant data and pro-
131 vides a comprehensive view of the spread of misinformation.
- 132 • **Case studies evaluation.** The paper presents three case studies that
133 illustrate the application of DisTrack in tracking different types of false
134 information related to topics such as COVID-19 vaccines, political is-
135 sues, and more. These case studies demonstrate the tool’s versatility
136 and effectiveness in misinformation tracking across various subjects.

137 The product of this research is conceived as a cognitive computation sys-
138 tem, according to its purpose of replicating human problem-solving [15]. In
139 this paper, it entails the human tasks of detecting false information, discover-
140 ing all the posts related to it and analyzing the evolution of them, towards the
141 final tracking of falsehoods. Through the computational methods presented,
142 these manual processes are automated, assisted through AI.

143 More specifically, Natural Language Processing is cited as one of the
144 areas in charge of cognitive systems and analyzing content in OSNs as one of
145 the examples of computerized cognitive abilities. Likewise, the extraction of
146 knowledge with Social Network Analysis is one of the techniques mentioned in
147 Decision Support Systems (DSS), being these a combination of the outcomes
148 of complex data analyses and machine learning [15].

149 This paper continues as follows: Background section reviews the key con-
150 cepts and models used in this research; the next section describes the three
151 modules in the methodology, the tweets extraction procedure, Natural Lan-
152 guage Inference and graph generation; after this, the following section dis-
153 entangles three examples in which this described methodology has been set
154 into practice. The answers to the research questions are derived from these
155 evaluated use cases, composing the Results and Discussion section that sub-
156 stantiates the final conclusions in the last part and the future lines of work.

157 2. Background

158 In this section, earlier contributions to the visualization of misinformation
159 are explored. Discovering the veracity of a claim is crucial to visualizing the
160 dissemination of misinformation because there are several lines of discourse
161 surrounding any false-information claim, with fact-checkers on one side and
162 misinforming actors on the other. It is expected that any information cascade
163 about a hoax is surrounded by two different and opposing narratives.

164 Our research is motivated by earlier works performing semi-automated
165 fact-checking using transformer-based language models, able to detect whether
166 pieces of content (or claims) are either factually fake or have undetermined
167 veracity. Following the previous rationale we motivate our techniques by ex-
168 ploring language models applied to fact-checking, as well as Social Network
169 Analysis (or SNA) to accurately portray the dissemination of misinformation
170 across online social networks (OSNs).

171 2.1. Language Models

172 The development of machine learning and deep learning models in the
173 field of Natural Language Processing has made it possible to deal with com-
174 plex tasks related to Natural Language Understanding (NLU). One of the
175 most important steps was the emergence of language models with the intro-
176 duction of the attention mechanism, leading to the development of Trans-
177 former models like BERT [16], RoBERTa [17], and XML [18]. Unlike earlier
178 embeddings (like word2vec [19] or Glove [20]), Transformer models gener-
179 ate vectorial representations using contextual information from neighboring
180 words in the surrounding text, where each word is semantically informed by
181 the sentence.

182 These advances opened the era of Language Models (LMs), architectures
183 trained for tasks such as predicting the next word, but designed for multiple
184 NLP problems. Among these, there are many scenarios where LMs can be
185 deployed to fuel fact-checking processes, providing significant improvements
186 over traditional machine learning methods [21]. For example, these models
187 facilitate the automation of fact-checking by employing binary classification
188 to discern false facts within the input. The state-of-the-art literature shows
189 promising results in this line of work [21, 22]. Furthermore, this approach can
190 be refined beyond simple binary outcomes by using varied labels to provide
191 more nuanced distinctions between types of information [23].

192 *2.2. Automated fact-checking*

193 Misinformation is an ever-shifting issue. New pieces of misinformation
194 may emerge as time passes, new narratives may become misinformation
195 whereas older known hoaxes become real after an unexpected world event
196 happens. Automated models trained without information retrieval tech-
197 niques will inevitably become obsolete within the span of a few months.
198 Allowing a model to retrieve information from trusted sources allows for
199 proper decision making [23]. In contrast to this, an alternative approach
200 arises where the dataset is conceived as a knowledge base [24]. In this source,
201 the data consists of textual statements containing verified falsehoods. Us-
202 ing these falsehoods, a model can compare an unverified claim against any
203 verified falsehood and if there is any match, we can determine that the un-
204 verified claim also contains a falsehood. In this structure fact-checkers have a
205 double role: they are responsible for the curation of the database, as well as
206 the interpretation of the model output, giving complete control of the semi-
207 automated model to fact-checkers and responsibility of its application [14].

208 The process has two steps: information retrieval (IR) and Natural Lan-
209 guage Inference (NLI). For IR, one of the most commonly used methods is the
210 calculation of the semantic distance between semantic embeddings [25, 26].
211 This approach does not depend on a preliminary dataset of posts on the so-
212 cial network chosen to assess the veracity of an unseen post on that platform,
213 allowing the classification of texts that belong to messaging environments in
214 which full datasets are rarely obtained, such as WhatsApp [27]. The ad-
215 vances in this type of pipelines in the era of coronavirus encouraged research
216 to focus on refining that knowledge base for that specific context [28]. Nowa-
217 days, tools in the fact-checking process based on the cosine similarity of texts
218 have already been implemented in newsrooms and show their success over
219 other methods [29].

220 The second step in this fact-checking process is determining the align-
221 ment between retrieved falsehoods and the original unverified content [24].
222 This alignment can be applied through the use of Natural Language Inference
223 (NLI) as a subset of Natural Language. The task of NLI consists of checking
224 if Sentence A (hypothesis) is inferred from Sentence B (premise) [30]. In
225 essence, this involves demonstrating that both sentences make the same as-
226 sertion. The relationships obtained from NLI are the following: *entailment*,
227 when A and B refer to the same statement; *neutral*, when A and B are not
228 related to each other, and *contradiction*, when A and B refer to the opposite
229 statement [31].

230 The classification of statements in this task is performed by feeding the
231 algorithms with datasets specifically designed for it. Stanford Natural Lan-
232 guage Inference corpus (SNLI), with 570,000 pairs of sentences annotated
233 with one of the three categories mentioned above [32], stands as the main
234 reference for NLI, but further datasets solve the drawbacks of SNLI: for ex-
235 ample, MultiNLI, with texts extracted from images cutlines, presents a more
236 enriched text [33], or XNLI, with its cross-lingual approach, does not re-
237 strict NLI to one language [34]. The use of NLI datasets enhanced with
238 Transformers facilitates the comparison among languages [35].

239 *2.3. Social Network Analysis*

240 Studying and mitigating the problem of misinformation involves detecting
241 the misinformation itself, but also tackling the pathways by which it spreads.
242 Thus, understanding how a piece of misinformation is disseminated on a
243 social network is a vital tool.

244 The flow of social media posts on Online Social Networks (OSNs) can
245 be effectively represented as a graph. This graph-based structure organizes
246 data into vertices linked by edges, providing a clear visualization of complex
247 relationships [36]. Within this framework, vertices represent either users or
248 their posts, while edges illustrate the myriad interactions or relationships
249 between them. These connections encompass the explicit social interactions
250 derived from the network’s metadata but, additionally, the more subtle, la-
251 tent properties that link posts together.

252 Social Network Analysis (SNA) is the area in charge of studying these
253 graphs from social platforms. The directions in this discipline comprise both
254 the extraction of the common features of networks and the identification of
255 aspects from the users from the graph [37]. Algorithms can be trained with
256 this information, using, for example, the dynamics of likes [38], to distinguish
257 between types of posts. Misinformation arises as one of the emerging domains
258 of application of SNA, together with politics and multimedia, being fields
259 such as marketing, tourism, healthcare or cybersecurity more settled in this
260 sort of studies [36].

261 Surveys in the field of misinformation have highlighted that the charac-
262 teristics inside the post are one of the indicators to detect falsehoods, but
263 the properties of the OSN itself play an important role. Sharma et al. [39]
264 enumerated key elements such as the source/promoters, user responses and
265 the information content parts. On the other hand, Parikh et al. [40] specified
266 non-text cues-based methods in the fight against false information, with user

267 behaviour analysis as one of the subareas covered. In 2018, different studies
268 demonstrated how false pieces of information were disseminated much fur-
269 ther on Twitter than those that were true, by looking at properties such as
270 the depth of the cascades generated by the post, the accounts contributing
271 to spreading them and the duration of the propagation [41]. However, in the
272 COVID-19 context, graph analyses revealed this expansion was only in terms
273 of vastness: both false and legitimate information have the same influence,
274 but actors spreading false information post more than those publishing the
275 true one [42].

276 However, regardless of the temporal context, repetitive patterns can be
277 found. Before the coronavirus, SNA showed that individuals' decisions re-
278 lated to vaccines could be influenced inside circles debating about vacci-
279 nation [43], indicating that the connections among users on these platforms
280 matter to the extent of dangerous implications if the communities approached
281 are anti-vaccine. These contagions from one group to another bring the issue
282 of virality, taken from the propagation of viruses. The creation of a cascade
283 responds to one of the models of infection, the viral model where infected
284 vertices by others can exchange the virus too, in contrast to the broadcast
285 models in which contagions derive from a main vertex [44].

286 The role of the main spreaders in these cascades of misinformation among
287 these circles has also been the focus of social-media-driven analyses. The
288 change of information from one community to another through 'super-spreaders'
289 and their characteristics were also disseminated with graphs [45]. Verified
290 Twitter accounts (in the era before Elon Musk ownership and X) were shown
291 to be 50 times more powerful in terms of propagating content about vaccines
292 in comparison to non-verified profiles [46].

293 To conclude, the use of NLP models without considering the dynamics
294 of the social network allows us to visualize only part of the problem, leaving
295 out key details [25, 26, 35]. Every instance of misinformation is surrounded
296 by a community of users who interact with, share, comment on, support, or
297 dispute it. Relying solely on content analysis limits our capabilities, over-
298 looking the crucial task of unraveling the impact of false claims on Online
299 Social Networks (OSNs), which extends beyond mere verification to include
300 users' responses. The combination of NLP and SNA leads to a more realistic
301 picture of the whole conversation of each falsehood, a technique yet to be
302 sufficiently explored in the literature as we have seen in this section.

303 3. Tracking misinformation in OSNs

304 NLP and SNA can represent an alternative map of misinformation in
305 OSNs through all the posts about a claim spreading a falsehood. In it, false
306 information appears from viralization, but also from messages of different
307 shapes from a wide to a short range of interactions and from a variety of
308 users, not necessarily with the same impact in terms of their popularity
309 in the social network. As an example, this approach allows to model the
310 contagion from one vertex to the rest [44] in two senses: on the one hand,
311 NLP-driven research has demonstrated that there is not a unique message
312 repeated in the propagation of misinformation, but many of them expressed
313 differently; on the other hand, SNA-oriented studies show unconnected users
314 outside the cascades also distribute falsehoods on social media.

315 Nevertheless, these studies contradicting the only focus on broadcast
316 models [47, 48] in their combination of NLP and SNA are limited to the
317 analysis of the data of the properties from the social network chosen in their
318 final goal. Graph generation is ignored and the possibility of representing
319 these ecosystems of falsehoods is missed. This results in an obstacle between
320 the theory that confirms how misinformation is not just a cascade and the
321 practice of representing what it is instead. This mentioned practice, in con-
322 trast to the previous approaches, would go further than the demonstration of
323 a different model of diffusion of misinformation to reveal a more realistic flux
324 of the posts causing it, examining its origin, evolution and end, if applicable.

325 3.1. The DisTrack architecture

326 The main goal of DisTrack is to create a complete representation of the
327 propagation cascade of a piece of misinformation. To do so, it integrates
328 different language models and SNA techniques that allow building a graphical
329 representation of this cascade, providing information about the content and
330 the actors in the social network that have played an essential role in the
331 dissemination. DisTrack consists of three main sequential steps:

- 332 1. **Information retrieval from OSNs:** this module comprises the ex-
333 traction of relevant keywords, the generation of queries through their
334 possible combinations and the use of Twitter API to download all the
335 tweets.
- 336 2. **Semantic and Natural Language Inference:** this module refers
337 to the conversion of tweets into Transformer-based embeddings that

338 capture their meaning and context and the extraction of metrics based
339 on their inference (if a tweet supports a false claim, contradicts it or
340 is unrelated. For this second step, we make use of FacTeR-Check [14],
341 which implements a semantic similarity filtering process followed by
342 Natural Language Inference.

343 3. **Graph generation:** this module ends with the hydration of the tweets
344 downloaded to extract the insights to be used in the graph and the cre-
345 ation of it. Its vertices will be the tweets extracted and whose edges
346 will correspond to the interactions among them. After this, the partic-
347 ularity of DisTrack is the use of the NLP-related and Twitter-related
348 metrics to show the rest of the properties (position, size or colour).

349 DisTrack leverages the concept of tracking by modeling a graph based
350 on a set of tweets downloaded from the OSN and labeled according to the
351 alignment with a specific claim. This modeling process includes information
352 extracted from the OSN such as following between users, retweets, or replies.
353 As a result, the mechanisms of DisTrack can be distributed in three modules
354 (see Fig. 2).

355 The final output of these three modules is a visualization that acts as
356 the operation center for the supervision of each piece of misinformation from
357 the beginning to the end, assessing about the flux of a certain claim, the
358 most influential spreaders and tweets with their connections or the periods
359 in which this false information has had more impact.

360 This architecture allows us to characterize the discourse surrounding any
361 piece of information. The visualization aims to highlight specific phenom-
362 ena in social media discourse, as well as provide answers to questions for
363 policymakers and content moderators. We provide some examples of these
364 properties:

- 365 • Proliferation (and decay) over time of falsehoods on social media. How
366 does a piece of content propagate across time in an OSN?
- 367 • Impact of fact-checkers on the discourse surrounding a falsehood. Do
368 they contribute to the proliferation or the decay of the surrounding
369 discourse?
- 370 • Relationship between the influence of a social network actor and their
371 impact on the discourse. How much do influential accounts dominate
372 online discourse? Do smaller accounts have any impact on the evolution
373 of the cascade?

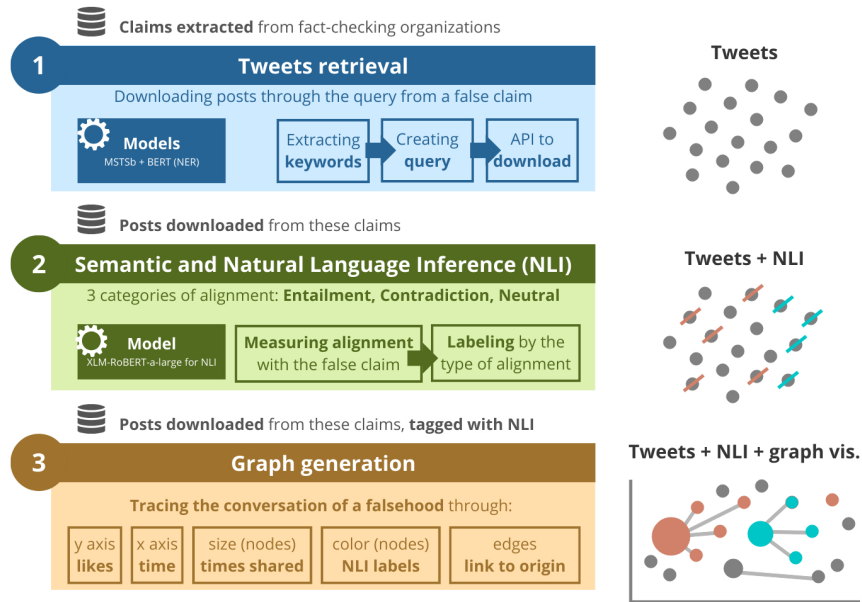


Figure 2: DisTrack modules, top to bottom: 1) Information Retrieval, 2) Natural Language Inference, 3) OSN Tracking Visualization.

- 374 • Detection of possible astroturfing campaigns. Has the falsehood spawned
 375 from several disconnected accounts? Is a coordinated attack on the so-
 376 cial network being performed?

377 This division into modules gives them independence in performing a task.
 378 For example, the extraction of Twitter- and NLI-related insights without the
 379 need to follow the trajectory of false information has already been applied,
 380 without any further implementation, to discover the type of tweets, according
 381 to their interactions, that represent the true volume of misinformation beyond
 382 viral posts. These modules are explained in the subsections below.

383 To enable our research we begin exploring Twitter ¹. This will be useful to
 384 evaluate our method on a heterogeneous set of use cases, but it is important
 385 to note that our methods can be applied to any social network that enables
 386 API search of textual content.

¹now X, however we will call it Twitter for the sake of simplicity and understandability.

387 *3.2. Retrieving Twitter Content*

388 The process of extracting information from Twitter requires the execution
389 of a series of searches for certain keywords. Since the Twitter API restricts
390 the search to the exact keywords that are given as input, in order to retrieve
391 a large representative sample of the relevant tweets and interactions related
392 to one specific claim, it is necessary to build a set of multiple queries. By
393 using different words and expressions in each of these queries it is possible
394 to cover a large part of the tweets referring to the claim.

395 *3.2.1. Keyword search*

396 The query is the input to search and download tweets through Twitter
397 API. These downloaded posts will contain the keywords inside that query.
398 Through logical operators, queries can request the API tweets with every
399 keyword inside it or for every tweet that at least has one of the keywords.
400 However, this process is manual and the downloaded tweets are the result
401 of the subjective decision of typing a query with a set of keywords instead
402 of another. How these keywords are distributed with the logical operators
403 to optimize the search is also an arbitrary decision. Furthermore, thinking
404 of all these steps of converting a false claim into a query to download as
405 many tweets as possible requires time, slowing down the computer-based
406 fight against misinformation.

407 To solve these drawbacks of the manual creation of a query, an infor-
408 mation retrieval module to generate queries and automate this process is
409 needed. This research follows FacTer-CheckKey [14], a designed automated
410 search to extract the most relevant keywords from a claim and concatenate
411 them through the logical operator “AND” to generate the final query. For
412 example, the sentence “Massive protest in France against the mandatory im-
413 plementation of the COVID passport in public spaces” results in the concate-
414 nation of keywords “(protest AND france AND passport AND covid AND
415 public)”. This development takes inspiration from the model KeyBERT [49]
416 and it uses multilingual Transformer-based models to take the context from
417 the meaning of a claim to extract the keywords from it, in addition to Name
418 Entity Recognition (NER) to improve this extraction for those cases in Span-
419 ish, when needed.

420 This automated extraction faces the problem of the subjective decision of
421 building the query manually, and they also deal with the subjective wording
422 of the claim itself, from which the automated query is generated. In practice,
423 including all the important keywords in these generated queries may not lead

424 to the expected results because there can be tweets with some of the keywords
425 from the claim but not necessarily all of them. Synonyms, abbreviations or
426 word camouflage are examples of figures of speech contributing to a difficult
427 keyword search.

428 We present a refined keyword extraction that tries to widen the search
429 space and minimize the obstacles of this process, introducing three major
430 add-ons:

- 431 1. The introduction of the logical operator *OR* in the query, to allow
432 multiple combinations of the keywords for the automated search
- 433 2. The inclusion of a parameter that indicates the number of keywords
434 that will be discarded from each combination of words, being the result
435 of the query a concatenation of all the possible mixtures excluding two
436 different keywords in each of them. In the example previously seen, the
437 query would be optimized in this way: “((protest AND france AND
438 passport AND covid) OR (protest AND france AND passport AND
439 public) OR (protest AND france AND covid AND public) OR (protest
440 AND passport AND covid AND public) OR (france AND passport
441 AND covid AND public))”.
- 442 3. The conversion of numerical numbers into all the possible forms, also
443 textual, to avoid ignoring pieces of misinformation that include figures
444 cited in a different way (e.g., “10000 OR 10,000 OR 10.000 OR ‘10
445 thousand’ OR ‘ten thousand’”)

446 In this case, `MSTsb-paraphrase-multilingual-mpnet-base-v2` model
447 is used in the query generator in parallel with `FacTer-Check`, to encode the
448 meaning together with the context of the words, and includes the model
449 `bert-spanish-cased-finetuned-ner` for Name Entity Recognition (NER).
450 For this implementation, adverbs, conjunctions, adpositions and other stop
451 words have been removed for the final composition with the multilingual
452 Flair tagger and Spacy models.

453 Although the default settings could capture how a certain hoax has been
454 diffused massively in a short time-lapse. For the sake of a better understand-
455 ing of social media discourse, we are interested in extracting the maximum
456 possible number of texts, actors and connections involved in the early steps
457 and also the evolution of every falsehood, given the dimensions of misinform-
458 ation on social media nowadays.

459 *3.2.2. Technical details*

460 The query generated through the claim of each piece of misinformation
461 constitutes the input of Twitter API, whose access is offered through a de-
462 veloper account on Twitter.

463 The reduced limitations granted by the academic API access soften two
464 restrictions: the maximum number of tweets extracted and the time of the
465 start of the final download. This last aspect can be modified to obtain
466 information from a certain timestamp. This becomes crucial in the field of
467 misinformation because the search must be constrained after the birth of the
468 topic that is referenced in each hoax (for example, claims about COVID will
469 not be found prior to the emergence of coronavirus).

470 Overall, the created query, the credential to validate the permissions (the
471 token found on the developer account) and the selected time-lapse and max-
472 imum number of posts as parameters are included for the automated request
473 for the needed tweets about misinformation. The data of each downloaded
474 tweet consists of a JSON with Twitter-based information structured as meta-
475 data in different fields.

476 However, since this cannot be exploited enough, NLP-based features in
477 the following step will contribute to filtering non-related hoaxes and devel-
478 oping the final tool. These technical details describe a fixed frame of the
479 use of the API at the moment of the experiment, but the current and future
480 restrictions of this developers' system do not distort the steps of this research
481 if the API is not available. Twitter versions of the API are a hands-on au-
482 tomated solution based on X advanced search, available through the search
483 bar on the interface. What the API returns in terms of content is the same
484 as the results that the interface gives to the user, also with the same query
485 as input. For this reason, parameters such as the chosen time-lapse can be
486 selected on X advanced search too.

487 The difference involves the codification of the output: whereas the API
488 offers the massive download of tweets in JSON format, the advanced search
489 gives the visual integration of those posts in the interface, which would need
490 further processing to transform them as raw data and to structure their
491 metadata.

492 *3.3. Automated Verification*

493 Once we have retrieved a sample of information from the OSN related
494 to the input claim, each tweet is labeled according to the alignment with
495 the original claim. We define alignment as the result of performing Natural

496 Language Inference over a pair of content (the retrieved tweet from the last
497 step) and falsehood (the piece of misinformation that is being tracked). After
498 evaluating each retrieved Tweet we assign their labels.

499 *3.3.1. Natural Language Inference*

500 Natural language inference (NLI) is crucial to distinguish if a tweet is
501 aligned or contradicts a false claim. Again, Transformer-based architectures
502 are applied to measure this alignment regardless of how different each tweet
503 is formulated in comparison to the false claim selected. The NLI task consists
504 of discovering if a hypothesis h can be inferred from the premise p in a pair
505 of texts (p, h) . In the misinformation and fact-checking domains, p will be
506 each tweet from the conversation downloaded from Twitter and h will be
507 the piece of misinformation debunked by fact-checkers. Thus: h is h_f when
508 stating that falsehood, and h is h_u when that factuality is undetermined.

509 The classification of posts according to their content through NLI is:

- 510 • **Entailment** (when the falsehood is enunciated): a post that entails a
511 piece of misinformation is a post that supports it and spreads it.
- 512 • **Contradiction** (when the negation of the falsehood is enunciated): a
513 post that contradicts a piece of misinformation is a post that denies it
514 and may act as a protective shield in the circles where it arrives.
- 515 • **Neutral** (when the falsehood or its contradiction is not enunciated):
516 a post that is neutral bears no relevance to the discourse whatsoever a
517 secondary effect of widening the spread of the keyword search.

518 *3.3.2. Technical details*

519 The semantic search uses the methodology proposed in FacTer-Check step
520 by step without ad-dons or modifications. For the NLI task, a fine-tuned
521 XLM-RoBERTa-large [50] was used for this module. For the NLI task, the
522 Machine Translated MultiNLI (MNLI-MT) [33] and XNLI [34] datasets were
523 used. Additional datasets such as ANLI [51], SNLI [32] and FEVER [52] for
524 English have been also included, using two training processes (inspired by
525 FacTer-Check): one only with MNLI (for the cross-lingual texts) and one
526 with all the mentioned datasets. The hyperparameters chosen are: 1024 as
527 batch size, $2e-5$ as the learning rate, with Adam [53] as the optimizer. Same

528 for warmup and linear decay. The validation data in XNLI determines the
529 optimal selected network after a manual hyper-parameter finetuning².

530 3.4. Graph visualization

531 Our main contribution lies in the graph visualization module where we
532 conceive how to translate an interconnected graph of tweets into a useful
533 visualization of online discourse around a topic. Our DisTrack architecture
534 makes use of all the properties stored in each tweet to make a readable
535 composition, clearly showing how any piece of misinformation has evolved.

536 3.4.1. Cascade graph building

537 The information stored for each tweet includes the preceding author and
538 the preceding tweet from which it derives. This information is enough to
539 generate a directed graph $G = (V, E)$, with V the vertex set containing tweets
540 published by an author. In this graph, the E edges represent a connection
541 between vertices due to being either a reply, quote or retweet from another
542 tweet, which will be the parent vertex.

543 Each vertex has additional information contained by the contextual meta-
544 data of the author, likes among other details. Conveying this metadata to
545 make it understandable is a non-trivial challenge that we address in the fol-
546 lowing subsections. In particular:

- 547 • *Time* is contained within the metadata, but the information is so un-
548 evenly spaced over time that the interpretation of the cascade is com-
549 promised.
- 550 • The cascade of misinformation is perpetrated mainly by actors that
551 influence online discourse. Representing their *influence* accurately in
552 the visualization aside from the vertices is crucial to understanding the
553 actual impact of actors.
- 554 • The information cascade of misinformation is twofold, containing usu-
555 ally two opposing narratives, the misinformation itself and the fact-
556 checkers combating it. Greater insight can be achieved by using the
557 results of NLI and integrating them into the visualization to evaluate
558 the *veracity* of the evaluated claims.

²Model available at [AIDA-UPM/xlm-roberta-large-snli_mnli_xnli_fever_r1_r2_r3](https://huggingface.co/AIDA-UPM/xlm-roberta-large-snli_mnli_xnli_fever_r1_r2_r3)

559 *3.4.2. Non-linear time representation*

560 The X-axis is defined by time. Each vertex will be placed according to
561 the moment it has been published. However, time itself is not represented
562 linearly in this axis. Falsehoods act in waves, the flow of information may be
563 quick and sudden during a short time-lapse generating concentrated content,
564 preventing the user from reading the evolution of that piece of false informa-
565 tion properly when it explodes in terms of impact. Furthermore, if there is
566 much distance in time from one vertex to the next one, again the space in the
567 graph between these two points will create empty areas in the plot whereas
568 the areas filled would be crowded with vertices (during the propagation of
569 a hoax, there are moments of inactivity or a lower number of publications
570 about it). Our approach just sorts posts chronologically, regardless of the
571 time passed between one vertex and the next one, and sets different time
572 stamps on the x-axis that improve the understanding of the propagation
573 cascade.

574 *3.4.3. Author influence representation*

575 Regarding the y-axis, it measures the degree of influence of the author
576 of a vertex/tweet. The number of followers allows us to distinguish between
577 the influence of the different actors responsible for the spread of falsehood or
578 their contradiction.

579 However, the achieved impressions by the tweet are also affected by the
580 number of likes, in addition to retweets and quotes, which are already rep-
581 resented through the children vertices. For this reason, likes will affect the
582 size of each vertex proportionately in the final visualization. Obtaining this
583 would result in a visualization that shows the evolution of a conversation of
584 the tweets extracted from a query but also the explanation of it through the
585 impact that users, retweets and likes generate.

586 *3.4.4. Veracity representation*

587 Finally, the outputs of NLI will serve as the colours to differentiate posts
588 containing misinformation from those that contradict or are unrelated. This
589 is the advance that allows us to check the beginning, transformation and
590 current status of a hoax rather than just all the tweets downloaded through
591 the main keywords appearing in a specific hoax, and that would culminate
592 in the goal of tracing the conversations about misinforming posts, if the
593 hypothesis is confirmed.

594 **4. Case studies**

595 Three use cases illustrate the application of DisTrack. They contain the
596 beginning and evolution of a piece of misinformation on Twitter, all of them
597 with a different topic to show the versatility of this tool. Firstly, an ex-
598 ploratory data analysis of the downloaded tweets is made to disentangle the
599 types of tweets and of their authors according to their NLI- and Twitter-
600 based metrics. After this, as the proof of the variety of tweets contributing
601 to the expansion of content related to misinformation, the representation of
602 the final graphs is made thanks to the final module of graph generation from
603 DisTrack.

604 These three cases represent three different topics:

- 605 • **Case 1: “The 80 percent of Muslims living in Europe live from**
606 **social welfare and they refuse to work”**. The first case is linked
607 to the disbelief in institutions and hate against Muslims. It constitutes
608 32 original tweets and a positive balance of tweets involving *entailment*
609 in contrast to *contradiction* (i.e., denying the hoax). The weight of
610 *entailment* increases much more when every post is shared. Thus, this
611 case contains a total of 84 representative posts, with all the retweets
612 included.
- 613 • **Case 2: “RNA vaccines against coronavirus includes graphene**
614 **oxid”**. It focuses on COVID-19-related antivaccine statements. It also
615 includes 32 original tweets and a balanced number between *entailment*
616 and *contradiction*. This second case has a total of 128 representative
617 posts, including retweets.
- 618 • **Case 3: “Zelensky sold 17 million hectares of land to Mon-**
619 **santo, Dupont and Cargill”**. It is an example of misinformation
620 around the Russia-Ukraine war. It departs from 26 original tweets and
621 most of the total tweets represent *entailment* with the hoax (80%).
622 This third case involves a total of 916 tweets.

623 The specific distribution of the number of tweets supporting or denying
624 each false claim is represented in Figure 3. Cases 1 and 3 are examples of the
625 general trend that can be observed in the dissemination of false information.
626 The presence of posts that support misinformation is much higher than those
627 that contradict the news, mostly due to the activity of fact-checker accounts.
628 In contrast, during the pandemic, we could see how a large part of the social

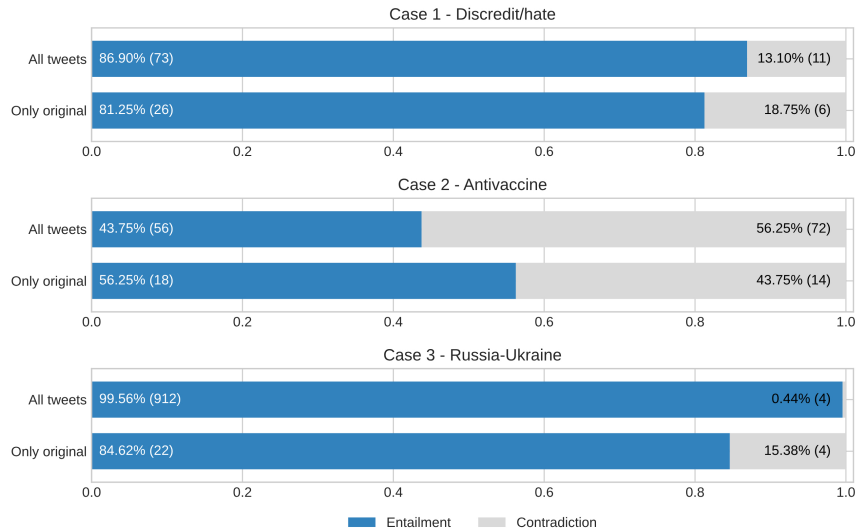


Figure 3: Distribution of the groups of posts according to their type of NLI-based alignment (*entailment* or *contradiction*), in each of the three cases.

629 media community actively participated in countering false information, as in
 630 the example of case 2.

631 In order to evaluate the impact of a post on a social network, the number of
 632 retweets or likes are two of the most commonly used indicators. As can
 633 be seen in Fig. 4, in terms of retweets, tweets from one to ten reposts are
 634 the main group, followed by those with zero retweets. Regarding likes, most
 635 of the posts in the flux of misinformation do not contain any interaction,
 636 in contrast with the lower percentages of tweets between one and ten likes.
 637 Tweets between one and ten likes are a small proportion except in the case
 638 of the hoax related to the vaccines, and only the propagation of the misin-
 639 formation cited about the war between Russia and Ukraine has a tweet with
 640 more than 100 likes.

641 We also analyzed the number of followers of the user accounts involved
 642 in these three case studies (see Fig. 5, showing different types of users. The
 643 original tweets in the three fluxes about misinformation are mainly shared by
 644 users from 1000 or more followers. Whereas the propagation related to the
 645 hoax about discredit/hate has more users between 1001 and 10.000 followers
 646 than those that are below or above these numbers, the other two examples
 647 indicate that authors between 101 and 1000 correspond to the largest pro-
 648 portion.

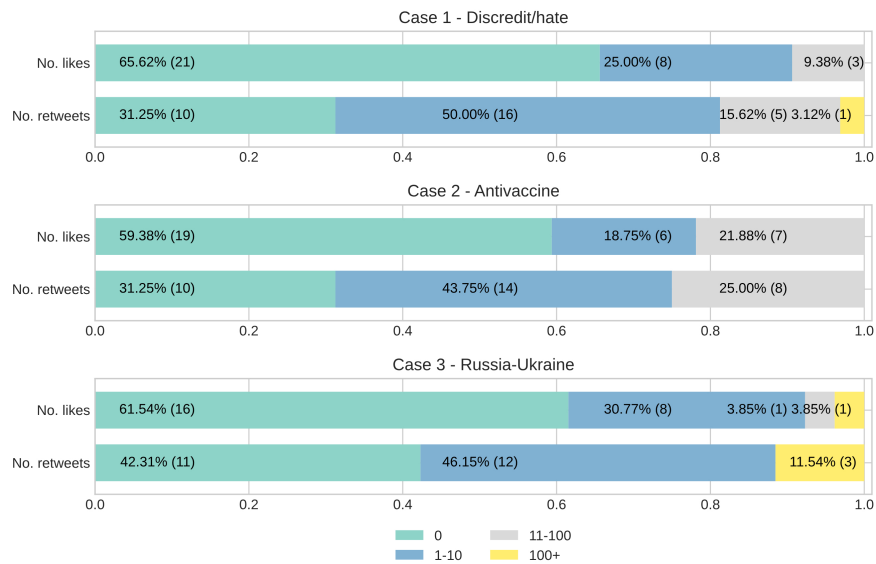


Figure 4: Distribution of the groups of posts according to their number of retweets and number of likes, in each of the three cases, with only original tweets into account.

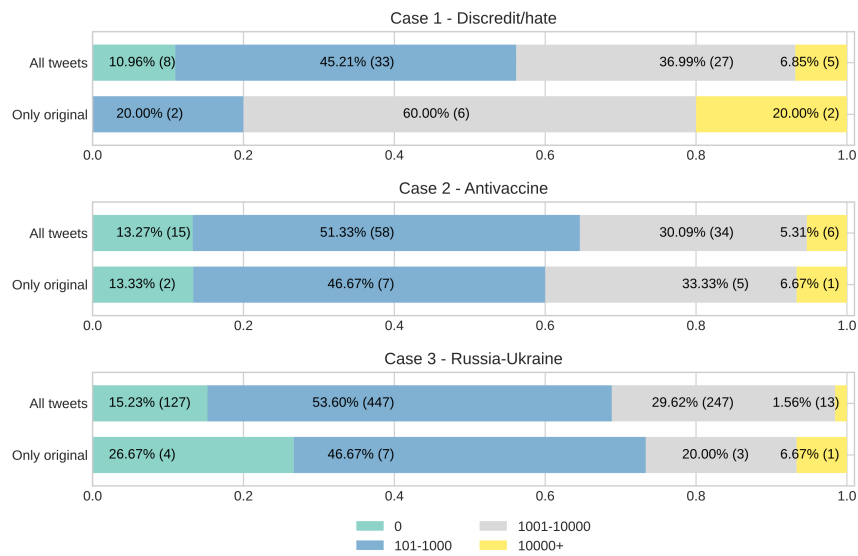


Figure 5: Distribution of the groups of posts according to their number of followers, in each of the three cases.

Author	No. Followers	No. interactions	Max No. Retweets	Max No. Likes	No. Tweets
0	25464	0	0	0	3
1	12537	7	7	5	4
2	8795	4	1	2	1
3	3881	1	1	0	1
4	3856	0	0	1	1
5	2669	0	0	1	1
6	2641	27	22	72	2
7	1141	1	1	2	1
8	417	0	0	0	1
9	146	38	32	60	5

Table 1: Ranking of the most active accounts in the spread of the tweets of the case 1.

649 In the following subsections we show and describe in detail each of the
650 three case studies. In the graphs, posts are rounded vertices, sized by the
651 number of likes, whereas reposts are vertices with the shape of a rhom-
652 bus. Posts and reposts tagged as *entailment* are orange, those tagged as
653 *contradiction* are blue and *neutral* ones are grey. The x-axis represents the
654 chronological order (with annotated dates each 100 days from the first post
655 in the visualization as a guide) and the y-axis (log scale) shows the number
656 of followers of the authors of each post and repost.

657 *4.1. Case 1: Discredit/hate*

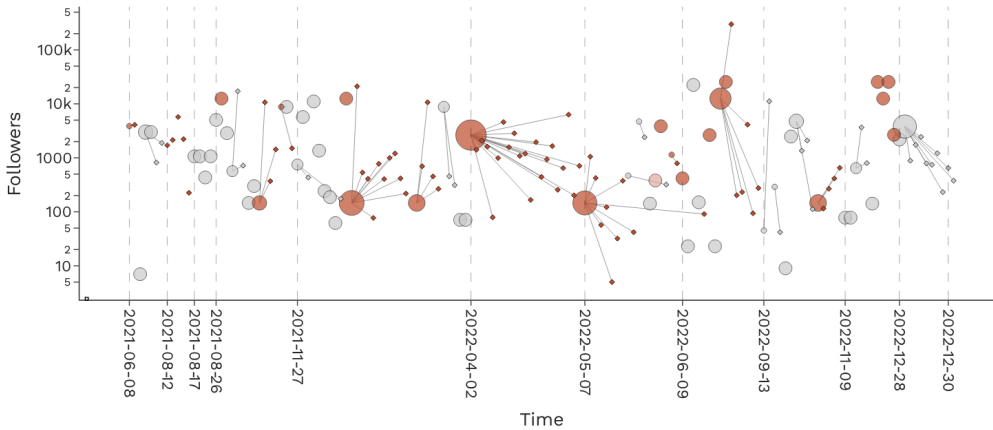


Figure 6: Visualization of the graph derived from the claim “80 percent of Muslims in Europe live from social welfare and they do not want to work”. The colour of the vertices mean blue for *contradiction*, red for *entailment* and gray for *neutral*.

658 For this first case (see Fig. 6), we focus on a hoax that circulated on Twit-
659 ter asserting that “80 percent of Muslims living in Europe live from social
660 welfare and they refuse to work”. For this hoax, we retrieved a large pool
661 of information from the social network, including mentions to other users,
662 popular or not, fake attributions to researchers to make the text trustworthy
663 and personal views and opinions reinforcing that hoax.

664 Likewise, fact-checks announce the same sentence in the negative form
665 to combat it, as seen in the graphs generated. The variations in this case
666 mention the name of the accounts of fact-checking and/or add links to the
667 news debunking the hoax. It is important to note that the same statement
668 to negate the false information is shared repeatedly at different moments in
669 the life of the hoax.

670 The most retweeted post affirming the hoax took place on April 2nd,
671 2022, with 22 retweets, however, we can see how this false information was
672 originally posted around 10 months before. On June 8th, 2021, an account
673 with 3,881 followers (at the time of the data extraction) published it and
674 received a unique retweet. It took exactly five months up to that date to see
675 the second most retweeted post supporting this claim, with 11 retweets, and
676 this type of misinformation resists at least until October.

677 The most interesting fact about the two users with more than 10.000 fol-
678 lowers who actively wrote the contradiction of the hoax is that both of them
679 reproduce the same action several times: one of them (25,464 followers) does
680 it three times; the other one (12,537 followers), four. Active spreaders posting
681 misinformation by themselves have a lower number of followers. However,
682 in the case of retweets, there are actors with more than 10,000 followers
683 that spread this falsehood. Maldito Bulo, from the fact-checker Maldita,
684 emerges as the user that exceeds that number of followers as a repost of its
685 fact-checking through another user.

686 This confirms that our system shows the whole cascade of misinformation
687 and not just the static picture of the most propagated tweet in April 2022.
688 Furthermore, this visualization also demonstrates that any virality of a post,
689 regardless of its impact, can be preceded by posts with zero or a few active
690 interactions (retweet or like), as shown in this case in 2021. Finally, it also
691 reveals how after the most retweeted post supporting the claim, the flux of
692 false information continues with different users and in different ways up to
693 the ones shown in the last trimester of 2022.

694 Regarding the involvement of the specific actors that spread misinforma-
695 tion (see Table 1), several accounts with an important number of followers

Author	No. Followers	No. interactions	Max No. Retweets	Max No. Likes	No. Tweets
0	9547	26	23	23	1
1	1903	6	6	13	1
2	1854	8	1	4	3
3	632	8	6	10	1
4	445	9	9	11	1
5	248	2	1	1	1
6	158	11	9	11	1
7	116	2	1	2	3
8	81	0	0	0	1
9	1	0	0	1	1

Table 2: Ranking of the most active accounts in the spread of the tweets of the case 2.

696 posted the hoax. It is interesting to note, however, that the ones that had
 697 the most interactions, and therefore caused the most discussion around the
 698 falsehood, were some with a smaller number of followers.

699 4.2. Case 2 - Antivaccine

700 The propagation of how falsely RNA vaccines against coronavirus include
 701 graphene oxide has also had an important impact on Twitter and other social
 702 media. This hoax has succeeded through paraphrasing in different periods.
 703 Whereas the previous case study showed a fixed structure and shape, this
 704 one discovers a range of posts expressing the same falsehood with different
 705 words and tones, from more declarative sentences to more aggressive ones.

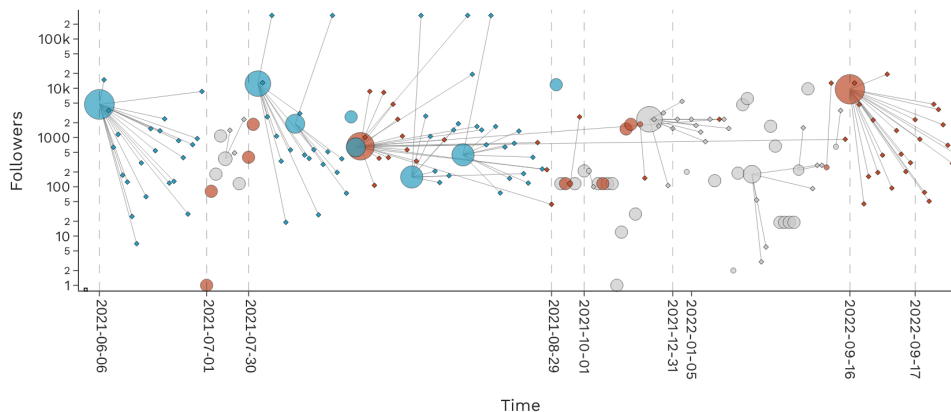


Figure 7: Visualization of the graph derived from the claim “Messenger RNA vaccines against COVID-19 contain graphene oxide”.

706 Again, fact-checks are mostly composed of the same sentence, including
 707 small variations mentioning the fact-checking source and/or the link. In some

708 specific cases, mainly in the first interactions of fact-checkers to contradict
709 false information, we can observe longer explanations, substantiating the
710 falsity of the fact.

711 In this particular case study, we can see that the first tweets composing
712 the propagation cascade deny the hoax. This surprising effect is present in
713 those cascades where the origin is in other environments (i.e., chain messages
714 on WhatsApp). The chain of propagation originates on 6 June 2021, and
715 about a month later we can observe some tweets disseminating the hoax. It
716 is noteworthy that the tweet that has the greatest effect on the dissemination
717 of the hoax (16 September 2022) occurs after months in which the only
718 references to the hoax are *neutral*.

719 Analyzing in more detail the accounts involved (see Table 4.2), the two
720 accounts that typed the contradiction of the claim have more than 10,000
721 followers, but none with this influence actively expresses the hoax itself (the
722 most influential one has 9,547 followers). Close to these numbers, the last
723 original tweet sharing this piece of misinformation had 9,547 followers. Re-
724 garding retweets, Maldito Bulo, from the fact-checker Maldita (300,081 fol-
725 lowers at that moment) played an important role in this hoax.

726 This case study reveals the limitations of traditional methods, which
727 might have analyzed this publication in August or its final iteration in Septem-
728 ber, without considering the origins or evolution of the misinformation in-
729 volved. By adopting this new approach, we gain insight into the broader
730 context of how misinformation spreads. It shows that the groups opposing a
731 piece of false information represent just a segment of the overall dissemina-
732 tion of posts linked to misinformation, as it continues to spread among other
733 individuals who either refute or challenge the assertion. The spread of the
734 latest hoax on September 16, which occurred without any users debunking
735 it, underscores this point. This example highlights the complexity of misin-
736 formation propagation and the necessity of considering its full lifecycle for
737 effective analysis.

738 4.3. Case 3: Russia-Ukraine

739 This third case, shown in Fig. 8, considers the hoax “Zelensky sold 17 mil-
740 lion hectares of land to Monsanto, Dupont and Cargill”. Unlike the earlier
741 cases, the visualizations here highlight variations in impact, illustrating the
742 hoax’s dissemination through a multitude of messages. Echoing the patterns
743 seen in the previous example, there is not a single form of post; instead, a
744 diversity of presentations emerge, ranging from the use of hashtags and user

Author	No. Followers	No. Interactions	Max No. Retweets	Max No. Likes	No. Tuits
0	43502	774	663	1000	1
1	2933	8	7	11	1
2	2211	0	0	1	1
3	1925	3	1	1	1
4	843	2	2	4	1
5	248	2	1	4	1
6	223	0	0	1	1
7	147	0	0	0	1
8	113	1	0	1	1
9	8	0	0	0	1

Table 3: Ranking of the most active accounts in the spread of the tweets of the case 3.

745 mentions to varying styles and tactics to engage the audience. The spread of
746 this misinformation is primarily driven by reposts, yet these varied expres-
747 sions of the same false claim also play a significant role, particularly at the
748 initial stages of its spread. This multiplicity of formats and channels under-
749 scores the complex nature of misinformation propagation and the challenges
750 in tracing and countering it.

751 This example showcases a viral post from September 19th, 2022, that
752 stood out with 663 retweets and 1,000 likes. Remarkably, its retweets oc-
753 curred not just on the day of the original tweet but continued sporadically
754 until December 30th of the same year, although diminishing impact. This
755 trend is visualized along the x-axis and highlights both the immediate impact
756 of the hoax and its prolonged presence in the digital discourse.

757 Thanks to this visualization, abnormal activity is shown on the same user
758 who writes the viral post. This suggests that this person retweets it several
759 times, according to the download of Twitter API data. This continuous ac-
760 tivity preserves the propagation of this false information and prevents it from
761 dying on Twitter (X). In addition, there was another post with even more
762 retweets whose initial user has been deleted or suspended (1.597 retweets),
763 as there are no edges linking them to the original post.

764 In this case, the virality of the post on September 19th already shows
765 the vast propagation of this misinformation, how it arrives to other users
766 throughout time, how it can be the continuation of other viral posts in the
767 past, and even the strategies led by users to keep the impact of a tweet. This
768 viral post has the same publication date as others with their own content
769 that also spread content about the hoax. In Table 3, it can be seen how one
770 user is the one leading the whole propagation cascade, spreading the hoax
771 while being the center of a high number of interactions.

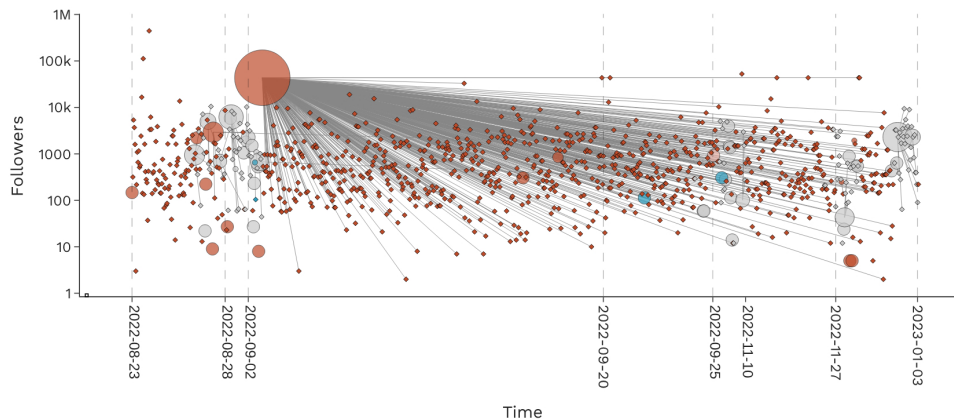


Figure 8: Visualization of the graph derived from the claim “Zelensky sold 17 million hectares of land to Monsanto, Dupont and Cargill”.

772 **5. Results and discussion**

773 In every analyzed example, DisTrack answers positively to the question
 774 “Can we extract the conversation about a hoax on Twitter?”. The conversa-
 775 tions about the three chosen falsehoods have been visualized after download-
 776 ing and processing and refining the data and metadata of the posts related
 777 to them. This demonstrates the success in generating automated queries to
 778 extract as many tweets as possible about false claims to later filter them
 779 depending on their type of alignment (*entailment*, *contradiction* or *neutral*).

780 This leads to the affirmative answer to the second subquestion “Can we
 781 separate tweets related to the hoax in the extracted conversation from tweets
 782 not related to it?”) and to the third one (“Can we distinguish between hoaxes
 783 that propagate a hoax from those that contradict it?”). The application of
 784 NLI leveraged with Transformers has been satisfactory at two levels: not
 785 only by separating tweets that are not related to each false claim, but also
 786 by separating the *entailment* posts (those that state misinformation) from
 787 *contradiction* posts (those that deny them).

788 In particular, the use of Transformers in the application of NLI to detect
 789 misinformation in this case has allowed DisTrack to identify falsehoods on X,
 790 regardless of the way they were exposed in the platform. In case 1, the posts
 791 were similar to the false claim taken as the reference, with subtle changes, but
 792 in cases 2 and 3, the posts paraphrased the content of the false claim together
 793 with hashtags, links or other elements, and their *entailment* was also guessed

794 by the NLI model. In these downloaded conversations about misinformation,
795 The existence of accounts reproducing exactly the same original content and
796 of those that paraphrase it instead encourages research to unveil the dynamics
797 of falsehoods and the existence of bots in contexts of crises [54].

798 Thanks to the whole process of searching tweets, aligning them with the
799 validated claim and visual representation, what we present is a complete
800 architecture that allows the automated and detailed analysis of the spread
801 of disinformation on a social network. From the first references to the hoax,
802 to the nodes that most influence its dissemination (thanks to their number
803 of followers, for example) or the patterns observed in the cascade, DisTrack
804 allows us to understand first-hand the dynamics of disinformation and how
805 it spreads and generates a certain impact.

806 The fourth research question “Can we identify the users involved in the
807 conversation of a hoax from beginning to end?” is also answered affirmatively.
808 DisTrack shows a closer image of the whole life of a piece of misinformation:
809 viral tweets matter, but all the participants in the conversation are relevant.
810 In the three case studies analyzed, falsehoods appear before and after the
811 spread of the most retweeted posts, in the shape of less viral tweets. This
812 research reveals the characteristics of every actor inside the ecosystem of
813 false information: the case studies unveil users with many followers as actors
814 in the propagation of a hoax, those conceived as ‘super-spreaders’ [45, 46],
815 but also accounts with different weights in terms of followers, not always the
816 most followed ones [54].

817 The results of the four subquestions allow us to answer the main re-
818 search question: “Is it possible to track conversations around specific hoaxes
819 on Twitter (X)?”. The three use cases confirm that misinformation can be
820 traced on Twitter by modeling it, using the three presented modules: firstly,
821 the tweets of the conversation about falsehoods were extracted through au-
822 tomated queries (subquestion 1); secondly, they were separated from each
823 other according to their alignment to that false information (subquestions 2
824 and 3), and, finally, their authors and their characteristics were also iden-
825 tified (subquestion 4) through the final generation of graphs that show the
826 evolution of that falsehood.

827 These outcomes follow recent research relying on tweet extraction and
828 NLI-based classification in the field of misinformation. These studies show
829 the orchestration of posts with different ranges of influence contributing to a
830 more complex propagation, and also the different nature of users in the con-
831 versation about specific falsehoods, in accordance with the results provided

832 by DisTrack. The exploratory data analysis of the three presented use cases
833 was a demonstration of this variety of posts and users.

834 Nevertheless, DisTrack adds a layer of research through the generation
835 of the graphs to exploit more the synergies between NLP and SNA as the
836 fields in charge of combating misinformation [23]: a falsehood does not die
837 progressively after its more viral publication and is not always born directly
838 from it, and fact-checks are also repeated in the successive lives of each
839 falsehood. In the cases analyzed, the conversation is reshaped and hoaxes
840 arise again. Furthermore, the evolution of tweets from fact-checkers and of
841 those that refuse verbally the selected hoaxes does not only occur with the
842 virality of that type of misinformation at its best, but also in other periods,
843 as observed in the final visualizations of this research.

844 The combination of falsehoods and the posts that debunk them in a final
845 visualization is also worth mentioning. DisTrack proposes an additional step
846 to the precursory work: whereas previous experiments contribute to a better
847 knowledge of the nature of the posts about misinformation, they do not offer
848 any surveillance to monitor and mitigate it. Nevertheless, DisTrack brings
849 back those steps of the tweets extraction and NLI filtering to put them at
850 the service of graph generation for that desired supervision. This enables the
851 coordinated response of fact-checkers and the control of their effects on false
852 information.

853 Overall, DisTrack reinforces the need to study misinformation as a viral
854 model [44] in a chain of infections by posts with different ranges of influence.
855 Its opposed model covered in research, the broadcast model, would only
856 have shown a part of misinformation because it only conceives contagion as
857 a unique primary vertex infecting the rest. On the contrary, the examples
858 made by Distrack print several versions of the vertex, various infections and,
859 as a consequence, many propagations instead of one.

860 This is also important in the combination of *entailment* and *contradiction*
861 posts. Not only do broadcast models [44] isolate a cascade of misinforma-
862 tion from the rest of the generated false information, but they also put it
863 apart from the posts that contradict it and from fact-checkers, preventing
864 researchers from depicting their appearances in the ecosystem of that cas-
865 cade. Likewise, a broadcast-model cascade of fact-checks without showing
866 the rest of the conversation through the steps of DisTrack would remove them
867 from the falsehoods they counteract.

868 6. Conclusions

869 All in all, our work releases a line of action through the shape of Dis-
870 Track, as the beginning of a tool able to merge Transformer-leveraged tweet
871 extraction, NLI-driven tagging of misinformation from the posts retrieved
872 and graph generation of Twitter-based and user-based properties in an out-
873 put that shows chronologically the evolution of a conversation motivated by
874 misinformation (spreaders, fact-checkers and other users) across the different
875 publications and actors involved.

876 With this proposed line of action, future experiments can study how
877 DisTrack modules can be modified. With NLP as one of the core parts of this
878 research, there is room for improvement given the advance of new Language
879 Models (LM): on the one hand, with their application for new models that
880 capture better the topics and, thus, the keywords for the query that enable
881 the download of posts; on the other hand, with their use for NLI to increase
882 the level of accuracy in the classification of posts as *entailment*, *contradiction*
883 of *neutral*, the three tags used to colour the vertices in the final graphs to
884 build the representative picture of misinformation.

885 Future work includes the application of DisTrack beyond misinformation.
886 For instance, previous research shows how sentiment analysis has become
887 relevant in the studies about aggressive discourse in the context of govern-
888 ment elections [55] and has stated Twitter as a “sentiment thermometer”
889 through VADER [56]. Although the first steps of these experiments evoke
890 the implementation of DisTrack, involving the download of posts from X and
891 the extraction of features for the study of specific behaviours, the part of
892 monitoring polarization through the combination of NLP and SNA through
893 graph generation is missing, unlike DisTrack. For this reason, future exper-
894 iments could be oriented to develop the same modules as DisTrack but by
895 substituting the task of detecting false information in favour of analyzing
896 sentiment in this platform, a space with more proportion of politicians and
897 political activists than others [9].

898 Accordingly, this research opens the door to the development of the part
899 of tracking in other scenarios, where NLI does not have to be necessarily
900 excluded and can enrich the information in the graph (in the previous case,
901 for example, to know the users stating the same political information and,
902 thus, enunciating the same message to the audience for the elections). These
903 other scenarios can include: a more general approach, beyond politics, to
904 take more advantage of the advances of Transformers for hate speech de-

905 tection [57, 58] and follow the trajectory of harmful posts on social media;
906 an alternative approach in the domain of author profiling, where graphs can
907 leverage recent research in LLM for this task [59] by tracing information and
908 linking the vertices that are likely to have the same attribution, or the area
909 of topic modeling, fueled by models such as BERTopic [60], in which the con-
910 catenation of the modules proposed by DisTrack could result in edges tying
911 the posts about the same issue, among other fields of study.

912 DisTrack also serves as an initiative to explore platforms beyond Twitter
913 (X). The international survey developed by Reuters Institute Digital News
914 Report in 2022 already revealed the prevalence of Facebook, YouTube, Insta-
915 gram or TikTok as well as WhatsApp, Telegram or Facebook Messenger as
916 ways of consuming news [61], which has continued in 2023 [9]. Each of them
917 has a different structure but DisTrack arises as a proposal to be adapted to
918 other scenarios of misinformation on these ecosystems.

919 Although OSNs mutate or change, with Twitter, now X, as an exam-
920 ple of that, users continue searching platforms to be in the circles of their
921 ecosystems. Research about migrations from Twitter to Bluesky, Mastodon
922 and Threads has already been covered [62], showing the interest in the cur-
923 rent microblogging social media and, thus, the necessity of fighting against
924 their information disorders, motivated by tools and methodologies such as
925 DisTrack.

926 **Declarations**

927 *Competing Interests*

928 The authors declare no competing interests.

929 *Funding*

930 This work has been funded by the project PCI2022-134990-2 (MARTINI)
931 of the CHISTERA IV Cofund 2021 program, funded by MCIN/AEI/10.13039/
932 501100011033 and by the “European Union NextGenerationEU/PRTR”; by
933 the research project DisTrack: Tracking disinformation in Online Social Net-
934 works through Deep Natural Language Processing, granted by Mobile World
935 Capital Foundation; by the Spanish Ministry of Science and Innovation under
936 FightDIS (PID2020-117263GB-I00); by MCIN/AEI/10.13039/501100011033/
937 and European Union NextGenerationEU/PRTR for XAI-Disinfodemics (PLEC
938 2021-007681) grant, by European Commission under IBERIFIER Plus - Iberian
939 Digital Media Observatory (DIGITAL-2023-DEPLOY- 04-EDMO-HUBS 101158511);

940 and by EMIF managed by the Calouste Gulbenkian Foundation, in the
941 project MuseAI.

942 **References**

- 943 [1] R. Salaverría, N. Buslón, F. López-Pan, B. León, I. López-Goñi, M.-C.
944 Erviti, Desinformación en tiempos de pandemia: tipología de los bulos
945 sobre la covid-19, *Profesional de la Información* 29 (3) (2020).
- 946 [2] E. M. Said-Hung, M. A. Merino-Arribas, J. Martínez-Torres, Evolución
947 del debate académico en la web of science y scopus sobre unfaking news
948 (2014-2019), *Estudios Sobre el Mensaje Periodístico* 27 (3) (2021) 961.
- 949 [3] A. M. Guess, B. A. Lyons, Misinformation, disinformation, and on-
950 line propaganda, *Social media and democracy: The state of the field,*
951 *prospects for reform* 10 (2020).
- 952 [4] N. A. Karlova, K. E. Fisher, A social diffusion model of misinformation
953 and disinformation for understanding human information behaviour, *In-*
954 *formation Research* (2013).
- 955 [5] C. Wardle, H. Derakhshan, *Information disorder: Toward an interdis-*
956 *ciplinary framework for research and policymaking*, Vol. 27, Council of
957 Europe Strasbourg, 2017.
- 958 [6] C. Ireton, J. Posetti, *Journalism, fake news & disinformation: handbook*
959 *for journalism education and training*, Unesco Publishing, 2018.
- 960 [7] J. Posetti, A. Matthews, A short guide to the history of ‘fake news’
961 and disinformation, *International Center for Journalists* 7 (2018) (2018)
962 2018–07.
- 963 [8] R. Kouzy, J. Abi Jaoude, A. Kraitem, M. B. El Alam, B. Karam,
964 E. Adib, J. Zarka, C. Traboulsi, E. W. Akl, K. Baddour, Coronavirus
965 goes viral: quantifying the covid-19 misinformation epidemic on twitter,
966 *Cureus* 12 (3) (2020).
- 967 [9] N. Newman, R. Fletcher, K. Eddy, C. T. Robertson, R. K. Nielsen, *Digi-*
968 *tal news report 2023*, RISJ: Reuters Institute for the Study of Journalism
969 (2023).

- 970 [10] M. Choraś, K. Demestichas, A. Gielczyk, Á. Herrero, P. Ksieniewicz,
971 K. Remoundou, D. Urda, M. Woźniak, Advanced machine learning tech-
972 niques for fake news (online disinformation) detection: A systematic
973 mapping study, *Applied Soft Computing* 101 (2021) 107050.
- 974 [11] D. Freelon, C. Wells, Disinformation as political communication, *Polit-
975 ical communication* 37 (2) (2020) 145–156.
- 976 [12] S. Altay, M. Berriche, H. Heuer, J. Farkas, S. Rathje, A survey of expert
977 views on misinformation: Definitions, determinants, solutions, and fu-
978 ture of the field, *Harvard Kennedy School Misinformation Review* 4 (4)
979 (2023) 1–34.
- 980 [13] J. S. Brennan, F. M. Simon, P. N. Howard, R. K. Nielsen, Types, sources,
981 and claims of covid-19 misinformation, *Reuters Institute for the Study
982 of Journalism* (2020).
- 983 [14] A. Martín, J. Huertas-Tato, Á. Huertas-García, G. Villar-Rodríguez,
984 D. Camacho, Facter-check: Semi-automated fact-checking through se-
985 mantic similarity and natural language inference, *Knowledge-Based Sys-
986 tems* (2022) 109265.
- 987 [15] I. Hasan, S. Rizvi, Review of ai techniques and cognitive computing
988 framework for intelligent decision support, in: *2021 8th International
989 Conference on Computing for Sustainable Global Development (INDI-
990 ACom)*, IEEE, 2021, pp. 891–898.
- 991 [16] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training
992 of deep bidirectional transformers for language understanding, *arXiv
993 preprint arXiv:1810.04805* (2018).
- 994 [17] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis,
995 L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert pre-
996 training approach, *arXiv preprint arXiv:1907.11692* (2019).
- 997 [18] G. Lample, A. Conneau, Cross-lingual language model pretraining,
998 *arXiv preprint arXiv:1901.07291* (2019).
- 999 [19] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed
1000 representations of words and phrases and their compositionality, *Ad-
1001 vances in neural information processing systems* 26 (2013).

- 1002 [20] J. Pennington, R. Socher, C. D. Manning, Glove: Global vectors for
1003 word representation, in: Proceedings of the 2014 conference on empirical
1004 methods in natural language processing (EMNLP), 2014, pp. 1532–1543.
- 1005 [21] A. Tretiakov, A. Martín, D. Camacho, Detection of false information
1006 in spanish using machine learning techniques, in: International Confer-
1007 ence on Intelligent Data Engineering and Automated Learning, Springer,
1008 2022, pp. 42–53.
- 1009 [22] H. Jwa, D. Oh, K. Park, J. M. Kang, H. Lim, exbake: Automatic fake
1010 news detection model based on bidirectional encoder representations
1011 from transformers (bert), Applied Sciences 9 (19) (2019) 4062.
- 1012 [23] A. Montoro-Montarroso, J. Cantón-Correa, P. Rosso, B. Chulvi,
1013 Á. Panizo-Lledot, J. Huertas-Tato, B. Calvo-Figueras, M. J. Rementeria,
1014 J. Gómez-Romero, Fighting disinformation with artificial intelligence:
1015 fundamentals, advances and challenges, Profesional de la información
1016 32 (3) (2023).
- 1017 [24] R. Vijjali, P. Potluri, S. Kumar, S. Teki, Two stage transformer model
1018 for covid-19 fake news detection and fact checking, arXiv preprint
1019 arXiv:2011.13253 (2020).
- 1020 [25] Á. Huertas-García, J. Huertas-Tato, A. Martín, D. Camacho, Civic-
1021 upm at checkthat! 2021: Integration of transformers in misinformation
1022 detection and topic classification., in: CLEF (Working Notes), 2021, pp.
1023 520–530.
- 1024 [26] Á. Huertas-García, J. Huertas-Tato, A. Martín, D. Camacho, Coun-
1025 tering misinformation through semantic-aware multilingual models, in:
1026 International conference on intelligent data engineering and automated
1027 learning, Springer, 2021, pp. 312–323.
- 1028 [27] J. Gaglani, Y. Gandhi, S. Gogate, A. Halbe, Unsupervised whatsapp
1029 fake news detection using semantic search, in: 2020 4th International
1030 Conference on Intelligent Computing and Control Systems (ICICCS),
1031 IEEE, 2020, pp. 285–289.
- 1032 [28] X. Guo, H. Mirzaalian, E. Sabir, A. Jaiswal, W. Abd-Almageed,
1033 Cord19sts: Covid-19 semantic textual similarity dataset, arXiv preprint
1034 arXiv:2007.02461 (2020).

- 1035 [29] I. Larraz, F. Sallicati, et al., Semantic similarity models for automated
1036 fact-checking: Claimcheck as a claim matching tool, *Profesional de la*
1037 *información* 32 (3) (2023).
- 1038 [30] B. MacCartney, *Natural language inference*, Stanford University, 2009.
- 1039 [31] S. Gururangan, S. Swayamdipta, O. Levy, R. Schwartz, S. R. Bowman,
1040 N. A. Smith, Annotation artifacts in natural language inference data,
1041 arXiv preprint arXiv:1803.02324 (2018).
- 1042 [32] S. R. Bowman, G. Angeli, C. Potts, C. D. Manning, A large anno-
1043 tated corpus for learning natural language inference, arXiv preprint
1044 arXiv:1508.05326 (2015).
- 1045 [33] A. Williams, N. Nangia, S. R. Bowman, A broad-coverage challenge
1046 corpus for sentence understanding through inference, arXiv preprint
1047 arXiv:1704.05426 (2017).
- 1048 [34] A. Conneau, G. Lample, R. Rinott, A. Williams, S. R. Bowman,
1049 H. Schwenk, V. Stoyanov, Xnli: Evaluating cross-lingual sentence rep-
1050 resentations, arXiv preprint arXiv:1809.05053 (2018).
- 1051 [35] J. Huertas-Tato, A. Martín, D. Camacho, Silt: Efficient transformer
1052 training for inter-lingual inference, *Expert Systems with Applications*
1053 200 (2022) 116923.
- 1054 [36] D. Camacho, Á. Panizo-LLedot, G. Bello-Orgaz, A. Gonzalez-Pardo,
1055 E. Cambria, The four dimensions of social network analysis: An
1056 overview of research methods, applications, and software tools, *Informa-*
1057 *tion Fusion* 63 (2020) 88–120.
- 1058 [37] A. Panizo-LLedot, J. Torregrosa, G. Bello-Orgaz, J. Thorburn, D. Ca-
1059 macho, Describing alt-right communities and their discourse on twitter
1060 during the 2018 us mid-term elections, in: *International conference on*
1061 *complex networks and their applications*, Springer, 2019, pp. 427–439.
- 1062 [38] E. Tacchini, G. Ballarin, M. L. Della Vedova, S. Moret, L. De Alfaro,
1063 Some like it hoax: Automated fake news detection in social networks,
1064 arXiv preprint arXiv:1704.07506 (2017).

- 1065 [39] K. Sharma, F. Qian, H. Jiang, N. Ruchansky, M. Zhang, Y. Liu, Com-
1066 bating fake news: A survey on identification and mitigation techniques,
1067 ACM Transactions on Intelligent Systems and Technology (TIST) 10 (3)
1068 (2019) 1–42.
- 1069 [40] S. B. Parikh, P. K. Atrey, Media-rich fake news detection: A survey,
1070 in: 2018 IEEE conference on multimedia information processing and
1071 retrieval (MIPR), IEEE, 2018, pp. 436–441.
- 1072 [41] S. Vosoughi, D. Roy, S. Aral, The spread of true and false news online,
1073 science 359 (6380) (2018) 1146–1151.
- 1074 [42] D. Saby, O. Philippe, N. Buslón, J. del Valle, O. Puig, R. Salaverría,
1075 M. J. Rementeria, Twitter analysis of covid-19 misinformation in spain,
1076 in: Computational Data and Social Networks: 10th International Confer-
1077 ence, CSoNet 2021, Virtual Event, November 15–17, 2021, Proceed-
1078 ings 10, Springer, 2021, pp. 267–278.
- 1079 [43] G. Bello-Orgaz, J. Hernandez-Castro, D. Camacho, Detecting discussion
1080 communities on vaccination in twitter, Future Generation Computer
1081 Systems 66 (2017) 125–136.
- 1082 [44] S. Goel, A. Anderson, J. Hofman, D. J. Watts, The structural virality
1083 of online diffusion, Management Science 62 (1) (2016) 180–196.
- 1084 [45] A. Bodaghi, J. Oliveira, The theater of fake news spreading, who plays
1085 which role? a study on real graphs of spreading on twitter, Expert
1086 Systems with Applications 189 (2022) 116110.
- 1087 [46] R. Carrasco Polaino, M. Á. Martín Cárdbaba, E. Villar Cirujano, Par-
1088 ticipación ciudadana en twitter. polémicas anti-vacunas en tiempos de
1089 covid-19, Comunicar: Revista científica iberoamericana de comunicación
1090 y educación.(Ejemplar dedicado a: Participación ciudadana en la esfera
1091 digital) 29 (69) (2021) 21–31.
- 1092 [47] G. Villar-Rodríguez, M. Souto-Rico, A. Martín, Virality, only the tip of
1093 the iceberg: ways of spread and interaction around covid-19 misinfor-
1094 mation in twitter, Communication & Society (2022) 239–256.

- 1095 [48] J. M. Noguera-Vivo, M. del Mar Grandío-Pérez, G. Villar-Rodríguez,
1096 A. Martín, D. Camacho, Disinformation and vaccines on social net-
1097 works: Behavior of hoaxes on twitter, *Revista Latina de Comunicación*
1098 *Social* (81) (2023) 44–62.
- 1099 [49] M. Grootendorst, [Keybert: Minimal keyword extraction with bert.](#)
1100 (2020). doi:10.5281/zenodo.4461265.
1101 URL <https://doi.org/10.5281/zenodo.4461265>
- 1102 [50] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek,
1103 F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, V. Stoyanov, Unsu-
1104 pervised cross-lingual representation learning at scale, arXiv preprint
1105 arXiv:1911.02116 (2019).
- 1106 [51] Y. Nie, A. Williams, E. Dinan, M. Bansal, J. Weston, D. Kiela, Adver-
1107 sarial nli: A new benchmark for natural language understanding, arXiv
1108 preprint arXiv:1910.14599 (2019).
- 1109 [52] J. Thorne, A. Vlachos, C. Christodoulopoulos, A. Mittal, Fever: a
1110 large-scale dataset for fact extraction and verification, arXiv preprint
1111 arXiv:1803.05355 (2018).
- 1112 [53] D. Kinga, J. B. Adam, et al., A method for stochastic optimization, in:
1113 *International conference on learning representations (ICLR)*, Vol. 5, San
1114 Diego, California;, 2015, p. 6.
- 1115 [54] M. Himelein-Wachowiak, S. Giorgi, A. Devoto, M. Rahman, L. Ungar,
1116 H. A. Schwartz, D. H. Epstein, L. Leggio, B. Curtis, Bots and misin-
1117 formation spread on social media: Implications for covid-19, *Journal of*
1118 *medical Internet research* 23 (5) (2021) e26933.
- 1119 [55] J. Torregrosa, S. D’Antonio-Maceiras, G. Villar-Rodríguez, A. Hussain,
1120 E. Cambria, D. Camacho, A mixed approach for aggressive political
1121 discourse analysis on twitter, *Cognitive computation* 15 (2) (2023) 440–
1122 465.
- 1123 [56] C. Hutto, E. Gilbert, Vader: A parsimonious rule-based model for senti-
1124 ment analysis of social media text, in: *Proceedings of the international*
1125 *AAAI conference on web and social media*, Vol. 8, 2014, pp. 216–225.

- 1126 [57] R. T. Mutanga, N. Naicker, O. O. Olugbara, Hate speech detection in
1127 twitter using transformer methods, *International Journal of Advanced*
1128 *Computer Science and Applications* 11 (9) (2020).
- 1129 [58] S. G. Roy, U. Narayan, T. Raha, Z. Abid, V. Varma, Leveraging
1130 multilingual transformers for hate speech detection, arXiv preprint
1131 arXiv:2101.03207 (2021).
- 1132 [59] J. Huertas-Tato, A. Martín, D. Camacho, Understanding writing style
1133 in social media with a supervised contrastively pre-trained transformer,
1134 *Knowledge-Based Systems* 296 (2024) 111867.
- 1135 [60] M. Grootendorst, Bertopic: Neural topic modeling with a class-based
1136 tf-idf procedure, arXiv preprint arXiv:2203.05794 (2022).
- 1137 [61] S. Hölig, J. Behre, W. Schulz, Reuters institute digital news report 2022:
1138 *Ergebnisse für deutschland*, Reuters Institute (2022).
- 1139 [62] U. Jeong, A. Nirmal, K. Jha, S. X. Tang, H. R. Bernard, H. Liu, User
1140 migration across multiple social media platforms, in: *Proceedings of the*
1141 *2024 SIAM International Conference on Data Mining (SDM)*, SIAM,
1142 2024, pp. 436–444.