

Journal Pre-proof

Ensuring News Integrity against Online Information Disorder through Text Watermarking and Blockchain

Flavio Bertini, Alessandro Benetton and Danilo Montesi

PII: S2096-7209(25)00141-1
DOI: <https://doi.org/10.1016/j.bcra.2025.100414>
Reference: BCRA 100414

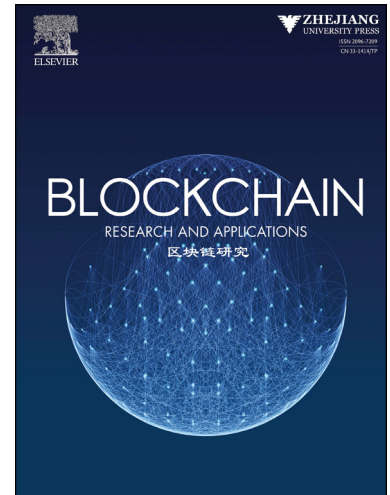
To appear in: *Blockchain: Research and Applications*

Received date: 1 February 2025
Revised date: 18 September 2025
Accepted date: 15 October 2025

Please cite this article as: F. Bertini, A. Benetton and D. Montesi, Ensuring News Integrity against Online Information Disorder through Text Watermarking and Blockchain, *Blockchain: Research and Applications*, 100414, doi: <https://doi.org/10.1016/j.bcra.2025.100414>.

This is a PDF of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability. This version will undergo additional copyediting, typesetting and review before it is published in its final form. As such, this version is no longer the Accepted Manuscript, but it is not yet the definitive Version of Record; we are providing this early version to give early visibility of the article. Please note that Elsevier's sharing policy for the Published Journal Article applies to this version, see: <https://www.elsevier.com/about/policies-and-standards/sharing#4-published-journal-article>. Please also 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.

© 2025 Published by Elsevier.



Highlights

- A platform for ensuring the integrity of textual news, combining text watermarking with blockchain technology to track provenance.
- Fine-grained watermarking to enable the identification of sources, even from small text excerpts.
- Blockchain-based storage of verification information in a decentralised manner, guaranteeing the authenticity of textual news.
- A browser extension for real-time verification of textual news, flagging any detected manipulations.

Ensuring News Integrity against Online Information Disorder through Text Watermarking and Blockchain

Flavio Bertini ^{a,*}, Alessandro Benetton^b, Danilo Montesi ^b



^a*Department of Mathematical, Physical and Computer Sciences, University of Parma, Parco Area delle Scienze 53/A, Parma, 43124, Italy*

^b*Department of Computer Science and Engineering, University of Bologna, Mura Anteo Zamboni 7, Bologna, 40126, Italy*

Abstract

The rapid proliferation of online information disorder poses a significant societal challenge. This phenomenon has been further exacerbated by the pervasive influence of social media, affecting a broad range of domains. Addressing the spread of information disorder through manual approaches (*e.g.*, human fact-checking) is impractical due to the vast volume of textual content daily generated, often facilitated by disruptive generative technologies. Similarly, the implementation of automated tools presents considerable obstacles, primarily due to the inherent ambiguity and complexity of natural language. This paper introduces the CERVANTES platform: an innovative application of blockchain technology integrated with text watermarking techniques. It is conceived to support news producers and enhance awareness among readers engaging with content across various social media platforms. CERVANTES is a multiplatform and language-independent solution that addresses online

*Corresponding author

Email addresses: flavio.bertini@unipr.it (Flavio Bertini ) ,
alessandro.benetton@studio.unibo.it (Alessandro Benetton),
daniilo.montesi@unibo.it (Danilo Montesi )

information disorder from an innovative perspective. It allows news producers to automatically embed a unique watermark within the content they create, ensuring the integrity and authenticity of the text, and safeguarding it against manipulation and misattribution. Furthermore, the association between the watermark and the corresponding news item is securely recorded on the blockchain, mitigating the risk of manipulation that might result from centralised management. We conduct an extended evaluation on twelve different social media platforms with a cohort of twenty users, observing ease of use and a high degree of satisfaction.

Keywords: Text watermarking, Text news sealing, Text provenance checking, Blockchain application, Online information disorder

1. Introduction

Over the past two decades, the proliferation of user-generated content on websites and social media platforms has significantly enhanced information exchange across ethnic, political, and geographical boundaries. This development has fostered the emergence of communities built around shared interests, perspectives, and narratives. However, in this landscape of disintermediation, it has become increasingly evident that information disorders frequently undermine human rights and the foundational principles of democracy [1]. A report by the Council of Europe identified the proliferation of information disorders, marked by the emergence of three distinct forms of misleading or false content: mis-, dis-, and mal-information [2]. These phenomena have had a significant impact across a wide range of domains, including democratic elections, public health communication, climate change

14 discourse, and geopolitical conflicts. Digital platforms and social media no-
15 tably amplify this trend. Such false, inaccurate, or misleading content has
16 the potential to manipulate public opinion, erode trust in institutions, and
17 fuel scepticism about critical issues such as vaccination [3, 4].

18 Online information disorder represents a multifaceted research challenge,
19 as its accurate classification necessitates addressing both malicious activities
20 and entirely lawful behaviours. Tandoc Jr. et al. have identified a range of
21 online behaviours that can be categorised within the framework of informa-
22 tion disorders [5]. Information disorder manifests in various forms, including
23 *news satire*, characterised by mock news with an explicitly humorous intent;
24 *news parody*, which employs non-factual information for comedic purposes
25 but, unlike news satire, fails to clearly indicate its non-journalistic nature;
26 *news fabrication* and *photo manipulation*, where articles and images devoid
27 of factual accuracy are intentionally altered to construct false narratives; and
28 *propaganda*, referring to news content created by political actors to influence
29 public opinion. Importantly, not all of these forms necessarily constitute
30 malicious behaviour.

31 The global expansion of online information disorder across websites and
32 social media platforms underscores the urgent need for effective counter-
33 measures. This issue is particularly complex, as it spans various forms of
34 online content, with text-based content posing distinct challenges due to its
35 nuanced nature [6]. The literature identifies various methodologies for ad-
36 dressing information disorder, broadly classified into language-based, topic-
37 agnostic, machine learning, knowledge-based, and hybrid approaches [7]. De-
38 spite their potential, automated methods exhibit significant limitations. In

39 particular, the inherent ambiguity of natural language poses substantial chal-
40 lenges to the development and application of content-based approaches, often
41 rendering them impractical. This issue has been further exacerbated by the
42 widespread adoption of Large Language Models (LLMs), which are capable
43 of generating human-like text often embedded with unverified or inaccurate
44 information [8]. Conversely, manual fact-checking approaches are equally
45 unfeasible due to the sheer volume of text disseminated daily, which makes
46 human intervention at scale impractical.

47 Moreover, one of the most common activities among online users is sharing
48 excerpts from original news, a practice in which both automatic and manual
49 methods fail, as this copy-and-paste behaviour removes the full context and
50 obscures the source information. Reader deception becomes even more evi-
51 dent when text content is presented as legitimate news and is accompanied
52 by seemingly authoritative markers of the user profile's integrity, such as the
53 blue checkmark on X (formerly Twitter). Although this verification badge
54 signifies that an account has undergone a verification process, it does not
55 guarantee the accuracy or authenticity of the content disseminated by that
56 account.

57 This paper addresses the challenge of online information disorder in tex-
58 tual news by introducing the CERVANTES platform designed to secure text
59 content and trace its source provenance. The proposed platform integrates
60 text watermarking techniques with blockchain technology to ensure the in-
61 tegrity of text news, even when partially re-shared, and to enable the iden-
62 tification of its source. Specifically, the watermarking is embedded in a fine-
63 grained manner, enabling even small excerpts of news content to be copy-and-

64 pasted and shared while maintaining traceability. The embedded watermark,
65 based on a previous text watermarking technique [9], is generated from both
66 the data and metadata of the original textual content, seamlessly integrated
67 into the text and securely stored on a blockchain for verification purposes.
68 In particular, the blockchain records critical information in a decentralised
69 and distributed manner, providing a consensus-driven framework for authen-
70 ticating textual news and enabling content traceability across diverse social
71 media platforms. A browser extension further enhances usability by allowing
72 readers to verify in real-time whether the content they are accessing corre-
73 sponds to the original one. Even when only a fragment of the original text
74 is displayed, the extension retrieves supplementary information from the dis-
75 tributed ledger. Additionally, any detected manipulations are immediately
76 flagged and highlighted to ensure transparency for the reader.

77 The CERVANTES platform is designed to empower readers to make
78 more informed decisions about the content they encounter, thereby reduc-
79 ing reliance on third-party fact-checking organisations. It is designed to be
80 both multiplatform and language-independent, capable of scaling to meet the
81 needs of diverse users. It operates independently of the social media plat-
82 forms utilised, supports the watermarking of texts written in the Latin al-
83 phabet irrespective of the language, and accommodates high-volume textual
84 content production without requiring human intervention. These features
85 enable news organizations of all sizes, as well as independent journalists,
86 to effectively track and disseminate text content across multiple platforms.
87 Moreover, the platform's capacity to manage the provenance, trustworthi-
88 ness, and verification of online news lies in its ability to strike a balance

89 between protecting freedom of expression and maintaining information qual-
90 ity. The proposed CERVANTES platform overcomes the content-based lim-
91 itations of traditional automated methods and addresses the deficiencies of
92 manual approaches, which are often influenced by censorship dynamics and
93 are impractical or prone to inaccuracies. To assess the effectiveness of the
94 CERVANTES platform, we performed an extensive evaluation across twelve
95 social media platforms. Additionally, a preliminary A/B testing experiment
96 was conducted involving twenty participants, who were asked to evaluate the
97 veracity of both original and manipulated news items on various topics, with
98 and without the assistance of the CERVANTES platform. Despite the small
99 size of the testing cohort, the experiment enabled the collection of qualita-
100 tive insights into users' responses to the introduction of a new interaction
101 element.

102 The paper is structured as follows. Section 2 offers a comprehensive
103 review of the literature relevant to this study. Section 3 presents a brief
104 overview of text watermarking attacks, to clarify the rationale behind their
105 integration with blockchain technology. Section 4 provides a detailed overview
106 of the CERVANTES architecture. Section 5 discusses the implementation
107 choices and presents the evaluation results. Section 6 discusses the limi-
108 tations of our work. Finally, Section 7 concludes the manuscript with key
109 insights and outlines potential directions for future research.

110 **2. Related Work**

111 In this section, we review relevant literature spanning multiple domains
112 that converge in the scope of our research. Specifically, we examine studies on

113 online information disorder, text watermarking techniques, and blockchain
114 technologies, highlighting their potential to verify the provenance of textual
115 content and ensure its integrity.

116 In recent years, the phenomenon of information disorder has witnessed
117 a significant escalation. While the internet has enabled the amplification of
118 diverse voices and democratised access to information, it has also introduced
119 a series of technological vulnerabilities that exacerbate the spread of informa-
120 tion disorder. Notably, a study conducted by MIT demonstrated that false in-
121 formation spreads more rapidly, extensively, and deeply compared to truthful
122 content, with falsehoods being 70% more likely to be retweeted than accurate
123 information [10]. Wardle et al. [2] proposed a tripartite classification of in-
124 formation disorder, distinguishing between misinformation, disinformation,
125 and malinformation. However, their framework does not encompass other
126 forms of expression, such as satire and parody, which introduce additional
127 complexities in efforts to address the dissemination of fake news [5]. Var-
128 ious methodologies address information disorder, including language-based,
129 topic-agnostic, machine learning, knowledge-based, and hybrid approaches
130 [7]. Language-based methods detect textual inconsistencies. Topic-agnostic
131 strategies rely on meta-information, such as advertisement density or sensa-
132 tionalist language. Machine learning uses annotated datasets to identify fake
133 news patterns. Knowledge-based methods combine machine learning with
134 knowledge engineering but struggle to keep pace with fabricated content.
135 Hybrid approaches merge human analysis and machine learning, focusing on
136 text, audience engagement, and source credibility. The multi-dimensional na-
137 ture of this issue, compounded by the inherent ambiguity of natural language

138 and the sheer volume of information shared online daily, renders both man-
139 ual and automated approaches to combating information disorder frequently
140 impractical and often imprecise [7]. Many existing EU projects require exten-
141 sive end-user involvement to evaluate individual posts and assess their trust-
142 worthiness [11], incorporate overly complex content evaluation procedures
143 [12], or lack the capability to address fine-grained textual content [13]. More
144 recently, a collaborative effort by Microsoft Research, BBC, CBC/Radio-
145 Canada, The New York Times, and Truepic introduced a platform aimed at
146 curbing the dissemination of fraudulent images and videos online [14].
147 The CERVANTES platform offers a fully automated solution for tracking and
148 verifying textual content. By integrating text watermarking techniques with
149 blockchain technology, the platform ensures both the integrity and traceabil-
150 ity of information, while simultaneously fostering greater awareness among
151 online readers regarding content authenticity.

152 Watermarking embeds a mark in digital content to prove ownership, au-
153 thorship, and verify content [15]. Text watermarking, due to its low embed-
154 ding capacity and limited syntactic and semantic variations, presents signif-
155 icant challenges. These have led to the development of various approaches,
156 broadly categorised as follows.

157 *Zero-watermarking* - In these methods, no direct watermark is embedded
158 in the text itself; instead, defining characteristics of the text are stored on
159 a third-party server [16]. A key limitation of zero-watermarking lies in its
160 centralised nature and the associated lack of transparency.

161 *Image-based approaches* - These methods necessitate the conversion of the
162 text into an image to apply the watermark [17, 18]. This approach alters

163 the original nature of the document and it is often impractical in contexts
164 such as online social media, where image-based methods are incongruent and
165 significantly constrain sharing practices.

166 *Syntactic methods* - These techniques operate on the syntactic structure of
167 natural language text by modifying the syntactic tree of a sentence to embed
168 a watermark. Examples of such syntactic operations include clefting, pas-
169 sivation, or activation [19, 20, 21]. A notable limitation of these methods
170 is their language dependency and low embedding capacity. Additionally, the
171 assumption that different syntactic forms convey identical meanings is not
172 always accurate, further constraining their applicability.

173 *Semantic methods* - These approaches leverage the semantic similarity be-
174 tween words by replacing them with their synonyms [22]. Other techniques
175 operate at the sentence level, capitalizing on the implicit presuppositions em-
176 bedded in each sentence [23, 24]. However, even when combined with syn-
177 tactic approaches to enhance the embedding capacity [25], semantic methods
178 share the same limitations as syntactic methods, particularly their low em-
179 bedding capacity and heavy dependence on language.

180 *Structural methods* - These methods do not modify the textual content di-
181 rectly; instead, they alter its structure, such as by inserting empty lines [26]
182 or utilising different Unicode whitespace characters [27, 28]. In recent work,
183 the authors employed an alphabet comprising five visually indistinguishable
184 Unicode whitespace characters to encode and embed a secret message within
185 a cover text [29]. Because the watermark is embedded within the underlying
186 representation of the text, these techniques offer the notable advantage of
187 preserving the original content without the need for an external database.

188 CERVANTES integrates a structural text watermarking technique for sealing
189 text, ensuring the preservation of content and length [9].

190 Blockchain technology enables the storage of transactional data in an
191 open, decentralised ledger, ensuring both verifiability and immutability. To
192 enhance efficiency, transactional data (*i.e.*, tokens) are organized into time-
193 stamped blocks, each linked to its predecessor [30]. The decision to add a
194 new block is made collectively by peers using a consensus protocol. Once
195 added, a block becomes immutable and cannot be removed. Transaction va-
196 lidity can be publicly verified by any client connecting to a peer that holds
197 the full blockchain. Originally designed for decentralized digital currencies,
198 this technology has since found applications in contexts where ensuring the
199 authenticity and verification of data is critical, without the reliance on a cen-
200 tralized, trusted authority [31, 32].

201 Several pioneering works have proposed blockchain-based solutions to address
202 the issue of online misinformation [33, 34, 35]. In [36], Arquam et al. intro-
203 duced a blockchain-based framework for verifying information propagation.
204 A blockchain-based model incorporating a secure voting system is proposed
205 in [37]. News reviewers provide feedback on news items, and a probabilistic
206 mathematical model is then used to estimate the truthfulness of each item
207 based on the feedback received. Alexandrescu et al. propose a decentralised
208 blockchain-based architecture for news retrieval and aggregation, separat-
209 ing crawling and scraping phases and ensuring information reliability via a
210 majority-based mechanism [38]. A conceptual framework for fake news de-
211 tection based on machine learning and blockchain is proposed in [39], where
212 both experts and users rate content to assess its authenticity. To counter the

213 spread of fake news on social media, [40] proposes a multilevel model based on
214 blockchain and deep learning techniques, designed to detect and prevent the
215 propagation of rumours. More recently, several approaches combining wa-
216 termarking and blockchain have been developed specifically for multimedia
217 content to combat fake news. In [41], the authors combine basic blockchain
218 and watermarking techniques to trace the origin of fake news and limit its
219 spread on social media. A theoretical framework combining blockchain and
220 watermarking is proposed in [42]. The authors claim that the framework
221 ensures the integrity of posted content and enables accountability of the
222 post’s owner or user. In [43], the authors present a blockchain-based system
223 for tracking photographs, which accommodates image transformations while
224 ensuring provenance and integrity, and demonstrates advantages over tra-
225 ditional watermarking approaches. A recent study proposes a semi-fragile,
226 blind watermarking system within a blockchain-based framework, designed
227 to detect fake images on social media, limit their spread, and support trace-
228 ability, manipulation tracking, and blacklisting [44]. The authors in [45] also
229 propose a blockchain-enabled watermarking technique to address deepfake-
230 related challenges, integrating content tracking, decentralised identity, and
231 cryptographic verification.

232 These solutions are primarily aimed at audio and video content, and the reli-
233 ability depends on social network parameters, whereas our approach does not
234 require any network structure to validate textual content. Moreover, in con-
235 trast to these systems, our solution leverages blockchain technology to track
236 textual content directly at the paragraph level, without relying on content
237 evaluation methods. These features enhance the versatility of CERVANTES,

238 making it readily applicable across multiple platforms and websites and fa-
239 cilitating its integration into various online information production systems.

240 **3. Background on Text Watermarking Attacks**

241 As this paper brings together expertise from distinct disciplines, we deemed
242 it necessary to include a section outlining some fundamental concepts related
243 to text watermarking attacks. This should later help clarify the rationale be-
244 hind the integration of blockchain technology within the proposed platform.

245 Focusing on structural text watermarking techniques - considered more
246 flexible and promising, as discussed in the Section 2 - several attacks can be
247 identified that aim to remove or compromise the watermark embedded in the
248 text. Based on a careful analysis of the literature [46, 47, 48], we identified
249 the following types of attack:

- 250 • *Deletion* - this type of attack consists in the removal of segments of text
251 to compromise the integrity of the watermark. This strategy can be
252 employed to reshape the message, aligning it with a specific narrative or
253 obscuring particular information, while preserving the overall coherence
254 of the text.
- 255 • *Insertion* - this attack is complementary to the previous one, as it
256 involves adding new words or characters to alter its meaning. A simple
257 example would be the insertion of the word “not” before a statement,
258 thereby reversing its intended message.
- 259 • *Replacement* - this attack involves replacing certain words or characters
260 in the text with different ones. It may be viewed as a combination of

261 deletion and insertion at the same location. A typical example is the
262 modification of a pronoun to alter its gender expression.

263 • *Copy&Paste* - this type of attack involves copying and reusing portions
264 of text, and is particularly insidious as it often relies on extracting
265 only small segments from a watermarked source. Such attacks can be
266 employed, for instance, to misrepresent the statements of a political
267 opponent by selectively quoting only those fragments that support the
268 attacker's narrative.

269 • *Retyping* - this attack occurs when a malicious user manually rewrites
270 the text in a different file or platform. As all structural methods are
271 inherently vulnerable to such attacks, this case is excluded and instead
272 informs the design of the three evaluation states of the CERVANTES
273 browser extension (see Section 4.2 for details).

274 • *Reformatting* - this attack involves altering the formatting character-
275 istics of a text. However, within the context of our work and the
276 intended application of the CERVANTES platform, such attacks can
277 be excluded without loss of generality. This is because social plat-
278 forms typically do not allow stylistic modifications that would affect
279 the watermarking method employed, unlike several earlier approaches
280 (for further details, see [49]).

281 Section 5.1 will illustrate how these attacks were adapted to reflect realistic
282 user behaviour, thereby emulating practical usage scenarios of the CER-
283 VANTES platform.

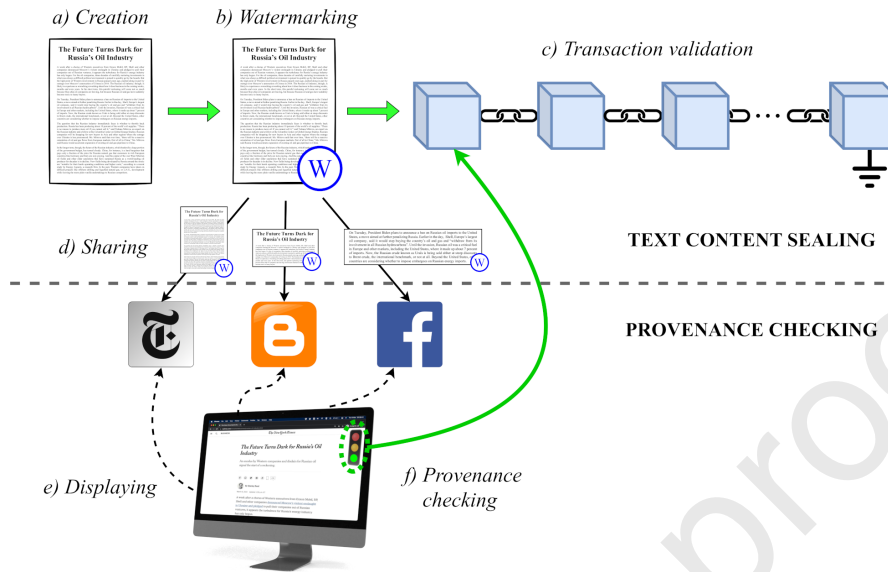


Figure 1: Overview of the architecture of the CERVANTES platform, illustrating its key components and functionalities.

284 4. The CERVANTES Architecture

285 This section outlines the architecture of the proposed platform, high-
 286 lighting its two key components. First, we examine the sealing of text news
 287 through text watermarking, as shown above the dotted line in Figure 1. Sec-
 288 ond, we provide a detailed explanation of the verification process, depicted
 289 below the dotted line in Figure 1.

290 4.1. Text Content Sealing

291 The content sealing process through structural text watermarking em-
 292 ploys homoglyph-based character substitution. The technique has already
 293 been evaluated in a previous publication [9]. In the present work, it is sub-
 294 ject only to minor modifications, with negligible impact on its performance
 295 and robustness, in order to meet the requirements of the proposed platform.

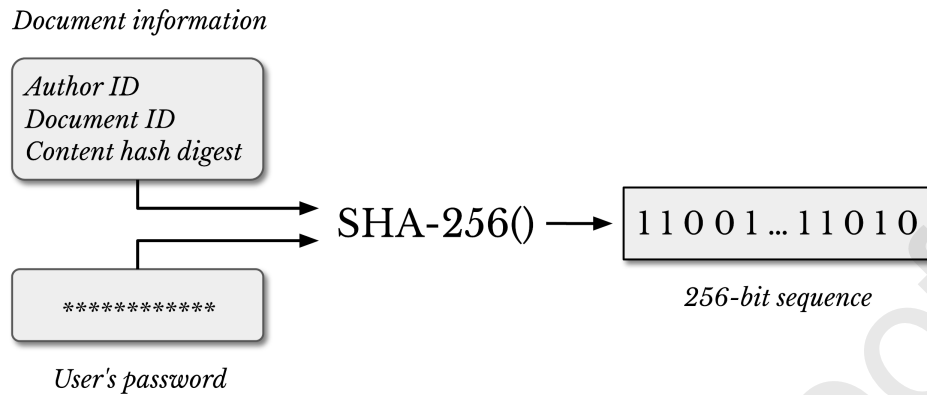


Figure 2: Generation of the watermark: given the document information and a secret user password, the hash function generates a bit sequence representing the watermark to be embedded.

296 Homoglyphs¹ are Unicode characters that visually resemble each other but
 297 differ in their underlying representations. This visual similarity is leveraged
 298 to embed a bit sequence watermark, which is generated using a hash function
 299 that combines the document's information with the user's password, thereby
 300 enhancing both security and resilience. Figure 2 shows the use of the SHA-
 301 256 hash function; however, the function is parametric, and a more secure
 302 variant can be selected. This would simply require a greater number of char-
 303 acters to embed the resulting, longer bit sequence.

304 Continuous and fine-grained watermarking, applied at the level of a few dozen
 305 words, is accomplished by concatenating the watermark throughout the doc-
 306 ument. This approach enables the validation of any sufficiently long copied
 307 portion that contains a complete watermark sequence. This requirement does
 308 not pose a limitation, even when considering character restrictions imposed

¹<https://www.unicode.org/reports/tr36>

309 by certain social media platforms, as the watermarking technique operates
310 effectively with between 46 and 101 characters [49]. To maintain content in-
311 tegrity, the CERVANTES platform stores hashes for each sub-portion, which
312 are explicitly identified to contain the full payload (*i.e.*, the watermark). An
313 invisible separator character is introduced before and after each sub-portion
314 during the watermark embedding process to ensure consistent separation be-
315 tween the sealing and validation phases. The embedded payload to generate
316 the watermark comprises the *author ID*, *document ID*, and a *seal* (*i.e.*, the
317 hash digest), thereby safeguarding the integrity of the text and confirming
318 that it has not been altered.

319 The lists of characters and whitespace symbols, along with their respec-
320 tive homoglyphs - which ensure that the final watermarked text remains
321 visually indistinguishable across all evaluated social media platforms [49] -
322 are shown in Table 1 and Table 2, respectively (further details are provided
323 in Section 5.3). The embedding process will employ the homoglyphs listed in
324 Tables 1 and 2 in order to enhance embedding capacity. Furthermore, Table
325 3 illustrates the frequency distribution of the original symbols and letters
326 across six Latin-based languages (*i.e.*, English, Spanish, French, Portuguese,
327 German, and Italian). For the English language, the frequencies were derived
328 from the *New York Times* corpus, comprising approximately 14 million words
329 [50]. For other languages, frequency calculations were based on the *Full-text*
330 *Corpus Data* website for Spanish and Portuguese², the *88milSMS* corpus for
331 French [51], the *DeReWo* corpus for German [52], and the *Il Post* corpus

²<https://www.corpusdata.org>

Table 1: The set of Latin letters and punctuation homoglyphs utilised during the watermarking process.

Character	Homoglyph	Original Code	Homoglyph Code
-	-	U+002D	U+2010
C	C	U+0043	U+216D
D	D	U+0044	U+216E
L	L	U+004C	U+216C
M	M	U+004D	U+216F
V	V	U+0056	U+2164
X	X	U+0058	U+2169
c	c	U+0063	U+217D
d	d	U+0064	U+217E
i	i	U+0069	U+2170
j	j	U+006A	U+0458
l	l	U+006C	U+217C
v	v	U+0076	U+2174
x	x	U+0078	U+2179

332 for Italian [53]. In addition to the homoglyphs listed in Table 1, the invis-
 333 ible character U+200B, commonly referred to as the *Zero Width Space*, is
 334 employed as a separator during the multiple stages of watermark embedding.

335 As shown in Figure 3, the embedding algorithm operates through the
 336 following sequence of steps:

- 337 1. the invisible character U+200B (*i.e.*, *Zero Width Space*) character is
 338 inserted to delineate the start of the embedding process (highlighted
 339 in red in Figure 3);

Table 2: Name and Unicode code of the eight whitespace characters, including the encodable bit sequences.

Whitespace	Unicode Code	Sequence Encodable
Space	U+0020	000
En quad	U+2000	001
Three-per-em space	U+2004	010
Four-per-em space	U+2005	011
Punctuation space	U+2008	100
Thin space	U+2009	101
Narrow no-break space	U+202f	110
Medium mathematical space	U+205f	111

- 340 2. the text is scanned sequentially until a character or whitespace listed
341 in Table 1 and Table 2 is identified (all highlighted in green in Figure
342 3);
- 343 3. if the subsequent bit in the watermark sequence is ‘1’, the identified
344 character is substituted with its equivalent homoglyph; otherwise, the
345 original character is retained. If the identified character in the previ-
346 ous step is a whitespace, this step remains unchanged. However, the
347 availability of eight distinct homoglyphs for whitespace enables the con-
348 sumption of three bits of the watermark per character (as illustrated
349 in Table 2 third column, and in Figure 3);
- 350 4. steps 2 and 3 are iterated until the entire sequence of watermark bits
351 has been embedded. Upon reaching the final bit of the watermark, an-
352 other *Zero Width Space* character is introduced, and the watermark bit
353 sequence is reset to its initial state to restart a new round of embedding.

Table 3: Frequency distribution of original symbols and letters across major Latin-based languages.

Character	English	Spanish	French	Portuguese	German	Italian
-	4.34%	0.11%	0.18%	0.79%	< 0.00%	0.04%
C	0.35%	0.23%	0.26%	0.34%	0.06%	0.24%
D	0.20%	0.14%	0.22%	0.21%	0.20%	0.12%
L	0.16%	0.20%	0.16%	0.20%	0.24%	0.25%
M	0.40%	0.16%	0.33%	0.28%	0.35%	0.19%
V	0.05%	0.05%	0.05%	0.08%	0.25%	0.06%
X	0.01%	0.01%	0.02%	0.01%	< 0.00%	0.01%
c	3.01%	4.35%	2.95%	3.44%	2.79%	4.13%
d	3.63%	4.86%	2.95%	5.22%	2.12%	3.80%
i	6.94%	6.32%	7.55%	6.39%	6.78%	11.33%
j	0.10%	0.38%	1.28%	0.28%	0.09%	0.01%
l	3.92%	5.21%	4.20%	2.75%	4.64%	6.11%
v	1.00%	0.95%	1.77%	1.28%	0.64%	1.44%
x	0.19%	0.18%	0.36%	0.24%	0.10%	0.03%

354 The four steps are iterated continuously until the entire document has been
355 processed. The extraction procedure mirrors these steps: it begins by lo-
356 cating the *Zero Width Space* character, subsequently identifies homoglyph-
357 compatible characters listed in Table 1 and Table 2, and reconstructs the
358 binary watermark sequence based on the presence or absence of homoglyph.
359 It is worth noting that the use of the invisible character U+200B as a wa-
360 termark delimiter can be further secured. Currently, malicious users familiar
361 with the watermarking scheme and the role of this special character could

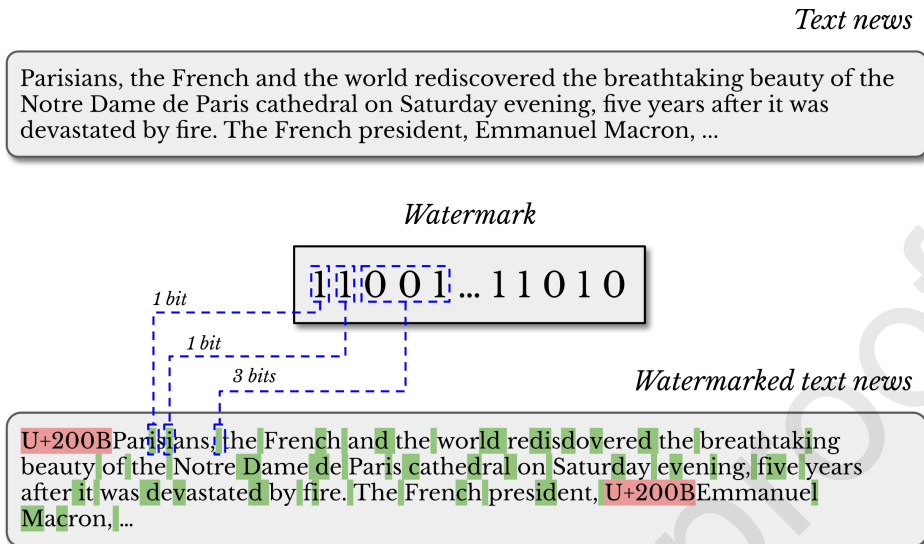


Figure 3: Embedding watermark bits by replacing confusable symbols. Only those characters with a corresponding duplicate in Tables 1 and 2 (highlighted in green) are used to embed the watermark bits according to the bit sequence. One bit can be embedded using Latin letters, and three bits using space characters. The U+200B character is used to delimit successive rounds of embedding. For the sake of simplicity, a shortened bit sequence is shown.

362 easily remove it, potentially resulting in a false negative (*i.e.*, the news item
 363 is reliable but cannot be verified). This aspect could be strengthened by
 364 obfuscating the generation of the character or sequence of characters used as
 365 delimiters, for example by incorporating the Innamark method proposed in
 366 [29].

367 The proposed watermarking approach enables content creators to secure
 368 their textual content, as described in phases *a* and *b* of Figure 1, and of-
 369 fers two key advantages. Firstly, by repeatedly applying the watermark and
 370 distributing uniquely signed versions across multiple social media platforms

371 (as illustrated in phase *d* of Figure 1), authors can enhance traceability and
372 pinpoint the origins of specific quotes. Secondly, the system enables the
373 immediate detection of alterations, even in brief excerpts of the original doc-
374 ument, while leveraging the pointer to the blockchain to provide readers with
375 supplementary metadata - such as the author’s identity, timestamps, and the
376 original source URL - thereby enhancing transparency and fostering trust in
377 the authenticity of the content.

378 *4.2. Provenance Checking*

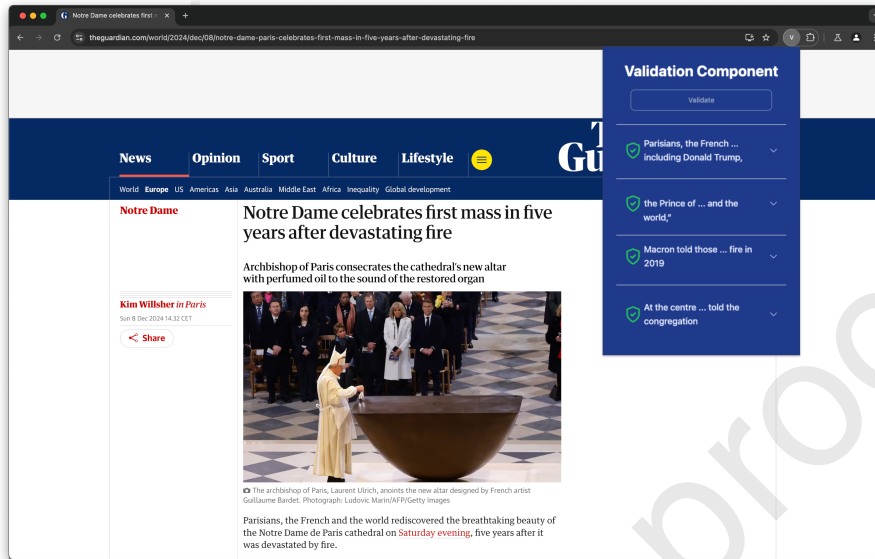
379 Relying exclusively on embedded information for text validation presents
380 several limitations. First, it is relatively straightforward to create a document
381 containing a valid watermark without utilizing the designated sealing tool.
382 Additionally, embedding extensive metadata directly into the text via water-
383 marking is inefficient, as minimizing the payload size is essential to reduce
384 the required embedding length. Larger payloads necessitate more text for
385 successful embedding, which becomes impractical when it is necessary to val-
386 idate excerpts of the original document. To address this, the CERVANTES
387 platform employs a minimal payload design, ensuring efficient and scalable
388 verification. Furthermore, the platform incorporates a provenance-checking
389 protocol that tracks watermarks via a blockchain, enabling the validation of
390 text fragments without requiring duplication of the entire document on the
391 ledger.

392 The CERVANTES platform integrates blockchain technology as a second
393 key component (depicted as phase *c* in Figure 1). This integration ensures a
394 consensus-driven approach for recording supplementary metadata and pro-
395 vides a decentralized framework for verifying the integrity of text content

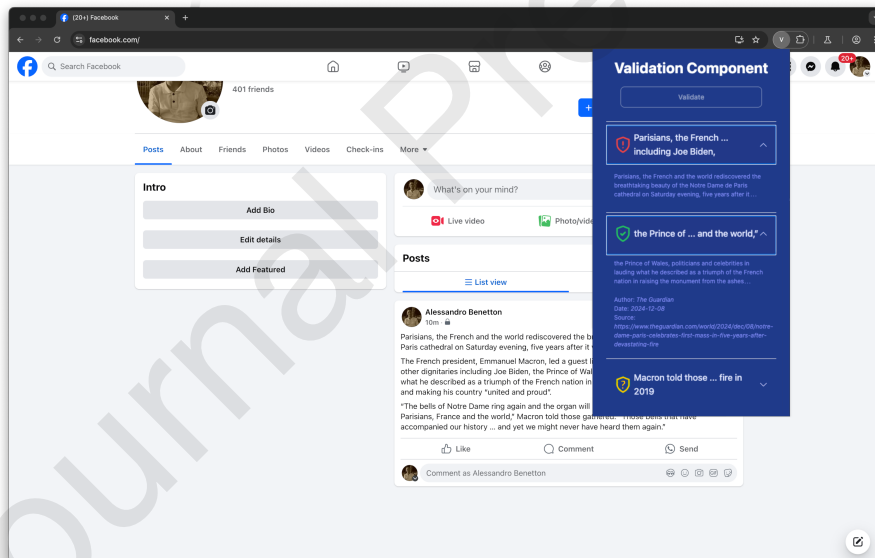
396 (illustrated as phase f in Figure 1). The distributed ledger functions as the
397 authoritative source of truth within the system and is designed to meet three
398 critical criteria: reliability, ensuring that all recorded data is accurate and
399 trustworthy; independence, guaranteeing that no single entity has control
400 over the stored information; and immutability, ensuring that once data is
401 entered into the ledger, it remains permanently unalterable.

402 An integral feature of the platform is its ability to automatically en-
403 able real-time content validation for readers. This functionality is facilitated
404 through a browser extension that actively scans the content of the displayed
405 web page for embedded watermarks. Upon detecting a watermark, the ex-
406 tension initiates a blockchain-based verification process to authenticate the
407 content. The results of this process are then presented to the user in the
408 form of a reliability indicator and any additional metadata retrieved (refer
409 to phases e and f in Figure 1). Notably, the validation mechanism remains
410 effective regardless of the length of the text excerpt being evaluated. Further-
411 more, any detected discrepancies between the watermark and the blockchain
412 record are flagged, promptly alerting the reader to potential issues with the
413 content’s integrity.

414 Figure 4 illustrates the result of the validation process for the displayed text
415 content, as observed both on the original news website (*i.e.*, The Guardian,
416 Figure 4a) and on the social media platform where the news was shared (*i.e.*,
417 Facebook, Figure 4b). Upon completion of the validation, the extension pro-
418 vides the user with a detailed report, including supplementary information
419 alongside the identified watermarked content. Specifically, each paragraph
420 is assigned one of three status labels: *valid*, *unknown*, or *invalid*. A *valid*



(a)



(b)

Figure 4: The browser extension for provenance verification: enhancing reader awareness on websites (*i.e.*, The Guardian) (a) and social media platforms (*i.e.*, Facebook) (b).

421 status indicates that the paragraph has been copied from a known source
 422 without modification, as the extracted watermark matches an entry in the
 423 blockchain (the green shields in Figure 4b). The *unknown* status signifies
 424 an incomplete or corrupted watermark, which may occur when the copied
 425 portion fails to capture the full watermark or when part of the seal is miss-
 426 ing (the yellow shield in Figure 4b). In contrast, an *invalid* status denotes
 427 that the watermark is absent or does not correspond to any entry in the
 428 blockchain, suggesting that the original seal has been deliberately altered or
 429 replaced (the red shield in Figure 4b).

430 4.3. Content Sharing and Verification Protocol

431 This section provides an overview of the processes involved in content
 432 sealing, dissemination, and verification, building upon the foundational com-
 433 ponents outlined in the preceding sections.



Figure 5: The procedural flow of text content sharing and verification.

434 The protocol for content sharing and verification consists of four primary
 435 steps, as illustrated in Figure 5:

- 436 1. The initial stage, referred to as *document writing & watermarking*,
 437 involves the creation of textual content by the news producer. During
 438 this process, the system automatically applies a unique watermark to
 439 the content through the functionalities provided by the *text content*
 440 *sealing* module.

- 441 2. In the subsequent stage, referred to as *manifest creation & transac-*
442 *tion storing*, the manifest associated with the newly created document
443 is constructed and securely stored on a blockchain. The correspond-
444 ing address is then registered on a distributed ledger. This manifest
445 encapsulates key metadata, including details about the author, the
446 document, and the hash digest of the original text, thereby enabling a
447 robust framework for validation and provenance verification of all the
448 excerpts of the original document.
- 449 3. In *document/short excerpts sharing* stage, the text - either in its en-
450 tirety or as selected portions - can be distributed across websites and
451 social media platforms.
- 452 4. In the final stage, referred to as the *visualisation & validation* phase,
453 any sufficiently lengthy segment of the original document that encom-
454 passes a complete watermark sequence can undergo validation through
455 the use of the *provenance checking* components.

456 The subsequent section will discuss the various implementation decisions
457 involved in the development of the platform prototype, as well as its evalua-
458 tion across twelve social media platforms with a cohort of twenty users.

459 **5. Prototyping and Evaluation**

460 The following section outlines the rationale behind key technical choices
461 to ensure the replicability of the platform. We begin by presenting results
462 on the robustness of the text watermarking technique against the attacks
463 outlined in Section 3, in order to complement the prior evaluation in [9].
464 Additionally, we conducted an evaluation of CERVANTES on twelve of the

465 most widely utilised social media platforms. Lastly, we report the results of
 466 the user assessment.

467 5.1. Text Watermarking Robustness

468 Section 3 outlined several types of attacks; however, an important aspect
 469 to consider is the proportion of text that must be affected to simulate a
 470 realistic and credible scenario. In [54], the authors suggest that altering 10%
 471 of a text constitutes a significant attack size for any watermarking technique.
 472 In the case of *copy & paste* attack, this percentage should be interpreted in
 473 reverse: as the proportion of copied text increases, watermark reconstruction
 474 becomes easier. In our tests, we consider portions of text copied from the
 475 original up to a minimum of 5% of its length. A set of 1,000 articles from
 476 the New York Times Corpus was used to carry out the tests, and the results
 477 are presented in Table 4.

Table 4: Success rate of watermark reconstruction following each specific attack involving 10% of the original text (*i.e.*, 5% in the case of the *copy&paste* attack).

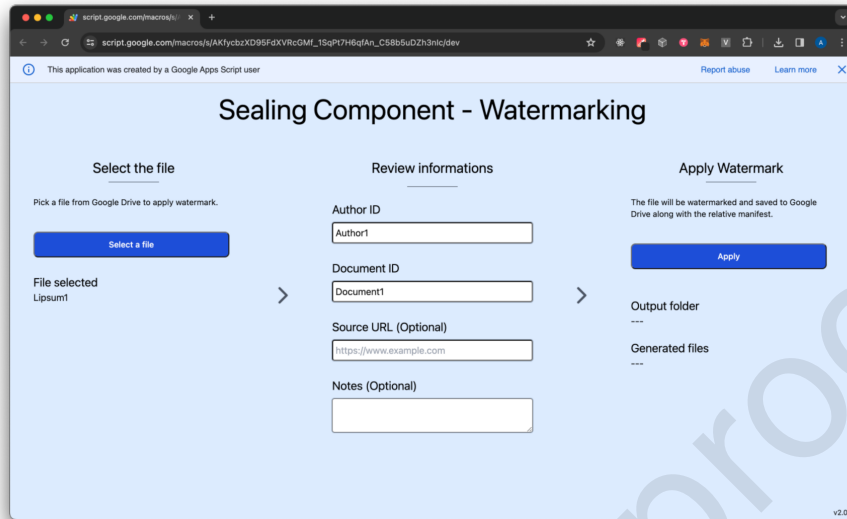
Deletion	Insertion	Replacement	Copy&Paste
98.20%	96.40%	98.30%	99.92%

478 The retyping and reformatting attacks have been excluded from the eval-
 479 uation. The former renders structural watermarking methods ineffective by
 480 definition, while the latter is not applicable in the context of social media
 481 platforms, which do not support formatting changes. The results obtained
 482 indicate that the watermarking technique integrated into the CERVANTES
 483 platform demonstrates strong robustness, largely attributable to its ability
 484 to watermark the original text repeatedly.

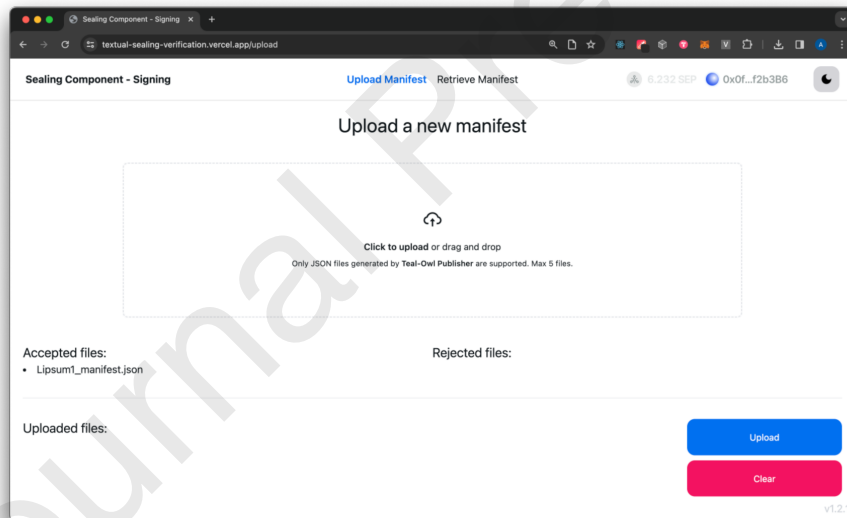
485 5.2. Platform Prototyping

486 CERVANTES platform comprises two primary components: *watermark-*
487 *ing & sealing* and *storing & validation*. The watermarking process embeds a
488 unique identifier into the text document with minimal computational over-
489 head, generating a specific file, referred to as the manifest. The manifest
490 contains critical metadata designed for secure storage on the blockchain,
491 facilitating subsequent validation and provenance verification. The Google
492 Apps Script environment was selected as the scripting language for imple-
493 menting the watermarking process within Google Docs. As illustrated in
494 Figure 6a, the script seamlessly and automatically embeds a watermark into
495 the text document and generates a JSON file, *i.e.*, the manifest. The mani-
496 fest includes essential metadata such as the *author ID*, *document ID*, *sealing*
497 *timestamp*, *source URL*, and a *list of hash digests* corresponding to each sub-
498 section of the text. These hash digests are crucial for verifying the provenance
499 of each of the short excerpts from the original document.

500 The second component consists of a web application designed to upload
501 the generated JSON file (*i.e.*, the manifest) to the InterPlanetary File Sys-
502 tem (IPFS) and record its corresponding address on a smart contract. This
503 process is facilitated through the user-friendly interface shown in Figure 6b.
504 Future developments aim to streamline the workflow by automating both the
505 watermarking and upload operations through an integrated interface. The
506 IPFS functions as a decentralized and distributed storage network, leverag-
507 ing a peer-to-peer architecture to improve content availability and fault toler-
508 ance. By generating addresses derived from the content's hash, IPFS ensures
509 document immutability while providing an efficient mechanism for identify-



(a)



(b)

Figure 6: The Google Apps Script interface facilitates the watermarking of text documents within Google Docs (a) and enables the subsequent upload of the generated manifest (b).

510 ing duplicates. The smart contract, deployed on Ethereum-based *testnets* for
511 testing purposes, securely stores the IPFS address, facilitating the identifi-
512 cation and retrieval of the associated manifest. A Chromium-based browser
513 extension has been developed to validate text documents accessed by users.
514 When activated, the extension analyses the webpage, detects paragraphs
515 containing watermarked text, and cross-references them with data stored on
516 IPFS and entries recorded in the smart contract, ensuring authenticity and
517 provenance.

518 5.3. Platform Evaluation Across Social Media

519 The prototype of CERVANTES was tested in a controlled environment
520 using basic HTML pages to conduct a preliminary evaluation. Subsequently,
521 we assessed the platform's performance in real-world web environments through
522 the following protocol: *i*) posting a text containing all the special characters
523 listed in Tables 1 and 2 on a social media platform to verify their compatibil-
524 ity, and *ii*) posting multiple encrypted contents to examine the effectiveness
525 of the watermark extraction and validation process.

526 The proposed platform was evaluated across the following twelve social
527 media platforms: Bluesky, Discord, Facebook, LinkedIn, Mastodon, Quora,
528 Reddit, Telegram web, Whatsapp web, Wordpress, X (formerly Twitter), and
529 Youtube.

530 All the evaluated platforms demonstrated full compatibility with the special
531 characters outlined in Table 1, enabling the successful extraction and valida-
532 tion of the watermarked text. Minor adjustments were required for certain
533 platforms; for instance, Reddit required the implementation of dedicated
534 code to remove predictable formatting characters (*i.e.*, additional lines and

spaces at the start and end of HTML section containing text) that would otherwise interfere with the watermark extraction process. Similarly, on Facebook, the watermark could only be fully detected if the entire text was displayed (*i.e.*, the “show more” is clicked).

Whitespace support varies significantly across platforms. Discord, LinkedIn, WordPress, and X do not support the “En quad” whitespace (code U+2000). In contrast, Facebook, Quora, and Telegram Web do not support any of the whitespace characters listed in Table 2. For the first group of social media, the remaining whitespace still allows for the embedding of two bits of the watermark. While this limitation increases the length of the text required for watermark embedding, it does not significantly impact overall performance and the verification process.

It is worth noting the following characteristics of Quora and Mastodon. Quora imposes a 250-character limit on question text, although no such restriction applies to answers. Mastodon, on the other hand, features a multi-server architecture in which each server is highly customizable, which may require tailored solutions, as was the case with Reddit.

5.4. Users Evaluation

To preliminarily assess the effectiveness of the CERVANTES platform in enhancing user awareness, we conducted an A/B test involving 20 participants. The cohort was structured to ensure a balanced representation in terms of age (ranging from 18 to 65 years), education level (including both secondary school and university graduates), and gender. It is worth noting that the sample size is limited: achieving statistical significance in theoretical terms would require several thousand participants. Nevertheless, the cohort

560 enabled us to gather qualitative insights into users' responses to the intro-
 561 duction of a new interaction element (*i.e.*, the browser plug-in) within the
 562 content evaluation process.

563 Participants were then divided into two groups, with only one granted access
 564 to CERVANTES. Both groups were asked to evaluate ten news articles cover-
 565 ing five different topics (*i.e.*, politics, economy, culture, science, and sports).

566 For each topic, one article was slightly altered to reverse its original meaning.
 567 Participants were instructed to assess the reliability of the entire article and,
 568 if deemed unreliable, identify the specific paragraphs that had been modi-
 569 fied. The CERVANTES group had the option to use the platform at their
 570 discretion, ensuring a realistic scenario without imposing any constraints.

571 To assess and compare the performance of the two groups, we employed a set
 572 of well-established metrics commonly used for classifier evaluation, namely
 573 sensitivity, specificity, F1-score, and accuracy. The results are presented in
 574 Table 5.

Table 5: A/B test results for the two participant groups.

Group	Sensitivity	Specificity	F1-score	Accuracy
Without CERVANTES	0.778	0.311	0.631	0.544
Using CERVANTES	0.818	0.655	0.756	0.736

575 In our context, sensitivity represents the ability to correctly identify reli-
 576 able news, while specificity denotes the ability to detect manipulated articles.
 577 These metrics assess the participants' ability to distinguish reliable news and

578 evaluate the effectiveness of the CERVANTES platform in supporting this
579 process. In particular, the participants who had access to CERVANTES
580 chose to use it in 68.2% of cases. The improvement of the group that au-
581 tonomously chose to use CERVANTES is clearly reflected in the summary
582 metrics, with F1-score and accuracy increasing by 19.81% and 35.29%, re-
583 spectively. These results highlight the platform’s effectiveness in enhancing
584 users’ ability to assess news reliability. Also, a significant improvement was
585 observed in the ability to identify altered paragraphs, with the group that did
586 not have access to CERVANTES achieving only 20.0%, compared to 61.8%
587 for the group that used the platform.

588 **6. Discussion**

589 In this section, we critically analyse the crucial aspects of the CER-
590 VANTES platform, assessing the potential limitations and outlining possible
591 strategies for enhancement.

592 The text watermarking technique is designed for languages that use the
593 Latin script, which may limit its applicability in non-Western contexts. How-
594 ever, according to Britannica³, the Latin alphabet is the most widely adopted
595 writing system globally, used by nearly 70% of the world’s population. Fur-
596 thermore, many Western democracies, which are frequent targets of online
597 information disorder campaigns [55, 56], predominantly employ the Latin
598 alphabet. Therefore, this constraint does not represent a significant limita-
599 tion. Currently, the Latin alphabet serves as the primary writing system for

³“The World’s 5 Most Commonly Used Writing Systems”, <https://www.britannica.com/list/the-worlds-5-most-commonly-used-writing-systems>.

600 numerous widely spoken languages, including English, Spanish, French, Por-
601 tuguese, German, and Italian, with approximately 1.5 billion, 595 million,
602 321 million, 300 million, 200 million, and 81 million speakers worldwide, re-
603 spectively. As such, it encompasses a substantial proportion of global online
604 users.

605 It is worth noting that the proposed watermarking method, which relies
606 on the use of specific Unicode characters with well-defined visual properties,
607 is font-dependent and therefore susceptible to reformatting attacks. It con-
608 stitutes only a partial limitation of the CERVANTES platform. Firstly, the
609 platform is designed to interface with twelve social media services, offering
610 broad coverage in terms of user reach. The list of homoglyphs presented in
611 Tables 1 and 2 has been verified to be supported by all twelve platforms. It
612 is important to emphasise that, on these platforms, users are not given the
613 option to select or modify the font used for publishing content. Furthermore,
614 in the event of a reformatting attack involving a change of font, the water-
615 mark is not removed; rather, some characters may appear slightly different.
616 As a result, the provenance verification process would remain feasible. Only
617 the complete removal of the watermark would compromise the verification
618 process, a situation that, by design, would prompt the browser extension
619 to report the content as unverifiable, even if the content itself remains un-
620 changed.

621 Another critical aspect concerns the time required to register the man-
622 ifest, the costs associated with maintaining the blockchain, and scalability.
623 According to the protocol, news content can only be circulated after es-
624 tablishing the necessary groundwork for validation (*i.e.*, manifest storage).

625 While this process may introduce a slight delay, it represents a reasonable
626 trade-off, as CERVANTES is not intended for real-time applications, such
627 as stock trading, where immediate execution is critical. Regarding mainte-
628 nance costs, these primarily consist of a one-time expense for uploading a
629 manifest to the smart contract and recurring costs associated with maintain-
630 ing storage on IPFS. The initial cost is determined by the gas fees required
631 to execute an operation on an Ethereum smart contract. For ongoing stor-
632 age, several platforms provide “pinning” services to ensure the persistence of
633 IPFS documents, typically at an average rate of \$0.1 per GB. For instance,
634 assuming an average JSON manifest containing 20 paragraphs with a size of
635 2KB, the monthly storage cost would approximate \$0.1 for 20,000 manifests.
636 Moreover, the modular architecture of the CERVANTES platform enables
637 the substitution of this component with alternative blockchain technologies
638 that may offer enhanced performance or cost-efficiency. Finally, with the
639 increasing volume of news content sealing and provenance checking requests,
640 the CERVANTES platform - with minimal adjustments - could also address
641 the scalability challenge. Specifically, in the event of a substantial influx of
642 data (*i.e.*, new texts to be sealed) and transactions (*i.e.*, provenance ver-
643 ification requests), the platform could adopt modular and complementary
644 strategies to maintain efficiency and cost-effectiveness. These might include
645 the use of multiple blockchains to distribute certification requests, as well as a
646 multi-tiered verification process designed to optimise provenance verification
647 requests, which are computationally less intensive but expected to be more
648 numerous. For example, the sealing phase could be managed through differ-
649 ent blockchain technologies offering diverse service characteristics, while the

650 browser extension may interface with multiple blockchains operating at vary-
651 ing levels of granularity, such as national systems, journalist associations, or
652 individual newspapers, in a coordinated manner. This architecture would en-
653 hance both performance and operational efficiency. Therefore, high volumes
654 of content and frequent transactions do not pose a scalability constraint.

655 Although this lies outside the scope of the present work, as it concerns
656 a purely blockchain-related security issue, we note that the platform is not
657 entirely immune to poisoning attacks on the blockchain. In such attacks, a
658 malicious actor introduces manipulated data with the aim of compromising
659 the integrity of the ledger. In its current configuration, the platform would
660 produce false positives in the presence of such attacks. Nevertheless, the
661 adoption of multiple blockchains with cross-verification mechanisms, com-
662 bined with a decentralised consensus-based system for verifying the identity
663 and trustworthiness of authors prior to publication, may prove effective in
664 mitigating these risks.

665 CERVANTES addresses the challenge of online information disorder from
666 an orthogonal perspective, offering a user-friendly solution tailored to the
667 behaviour of online readers who often skim headlines, exhibit a limited in-
668 clination to delve deeper into topics or face difficulties in distinguishing fake
669 news [57]. The platform aims to maintain independence from social media
670 platforms, which have limited incentives to curb the spread of partially true
671 information that enhances user engagement [58]. To this end, while the se-
672 lection of platforms analysed may initially appear restrictive, it encompasses
673 333 widely used platforms, including Facebook, X, Reddit, and Threads,
674 which collectively account for approximately 2.9 billion, 611 million, 267.5

675 million, and 275 million monthly active users, respectively. Furthermore,
676 given the orthogonal nature of our approach to online information disorder,
677 a quantitative and qualitative comparison with machine learning methods
678 focused on direct content evaluation is scarcely feasible. Nonetheless, we do
679 not consider the CERVANTES platform a definitive solution to the problem,
680 but rather an additional tool that can be made available to users. We be-
681 lieve that future developments may well integrate multiple approaches, both
682 provenance verification and content evaluation, for instance, by triggering a
683 semantic content analysis only when the provenance verification platform is
684 unable to resolve the pointer to the blockchain. This would have the dual
685 advantage of reducing the computational workload and relying on a set of
686 previously certified reference news items for any given topic.

687 Finally, we briefly clarify the role of LLMs in online information disorder
688 and the potential contribution of CERVANTES in this context. Previous
689 research has explored the idea of embedding watermarks into textual content
690 generated by service providers, such as LLMs. Major industry players are
691 actively seeking solutions along these lines^{4,5}. The main reason is that text
692 generation services based on LLMs are widely used by malicious actors to
693 produce information disorder and create information pollution on a global
694 scale [59]. CERVANTES would enable the use of such generative tools by

⁴“Google Is Paying Publishers to Test an Unreleased Gen AI Platform” <https://www.adweek.com/media/google-paying-publishers-unreleased-gen-ai>.

⁵“Meta has created a way to watermark AI-generated speech” <https://www.technologyreview.com/2024/06/18/1094009/meta-has-created-a-way-to-watermark-ai-generated-speech>.

695 legitimate news producers while enhancing reader awareness in situations
696 where content verification is not feasible.

697 In the future, the platform may be enhanced with truth and consensus-
698 based mechanisms to mitigate the reliance on centralised evaluations of the
699 reputation of news creators, an aspect that lies beyond the scope of this pa-
700 per. An alternative approach could involve collaboration with established
701 organizations, such as professional journalist associations or international
702 fact-checking networks. This proposed platform aligns with the perspective
703 outlined by historian Y. N. Harari, who emphasizes that the fundamental
704 distinction between dictatorships and democracies lies in their approach to
705 information management [60]. While dictatorships prioritize control and re-
706 striction, democracies focus on the sharing, consensus-driven evaluation, and
707 dissemination of accurate information.

708 **7. Conclusion**

709 The widespread dissemination of online information disorder has emerged
710 as a critical societal challenge, exerting a profound influence on public dis-
711 course. This phenomenon increasingly undermines societal values and demo-
712 cratic processes, fostering opinion polarisation on critical issues and reshaping
713 perceptions of facts, truths, and beliefs without substantiated foundations.
714 While notable progress has been made in the development of automated
715 systems for detecting fake news, research on the human factors driving indi-
716 viduals to believe and share such information disorder remains in its infancy.

717 Instead of tackling the inherently complex task of reliably detecting on-
718 line information disorder based solely on content, this study introduces the

719 CERVANTES platform, designed to assist end users in identifying trustwor-
720 thy news sources. The platform employs an innovative integration of text
721 watermarking techniques with blockchain technology to achieve fine-grained
722 marking of textual news. This approach ensures that any sufficiently long
723 excerpt of the original news, shared across web pages or social media plat-
724 forms and containing a complete watermark sequence, can be authenticated
725 for its provenance.

726 The primary objective of the CERVANTES platform is to empower online
727 readers to form informed opinions about textual news independently, with-
728 out reliance on third-party fact-checking organisations. It effectively secures
729 even short excerpts of text, addressing the prevalent trend of users engag-
730 ing only with news headlines. Its robustness in managing the provenance,
731 trustworthiness, and verification of online news lies in its capacity to strike a
732 balance between freedom of expression and the assurance of information qual-
733 ity. As future work, CERVANTES can be further enhanced with mechanisms
734 grounded in truth and consensus, mitigating potential biases associated with
735 centralised assessments of the reputation of news creators.

736 **CRedit authorship contribution statement**

737 **Flavio Bertini:** Conceptualization, Methodology, Validation, Writing –
738 original draft, Writing – review & editing. **Alessandro Benetton:** Concep-
739 tualization, Software, Investigation. **Danilo Montesi:** Conceptualization,
740 Writing – review & editing, Supervision.

741 **Declaration of competing interest**

742 The authors declare that they have no known competing financial inter-
743 ests or personal relationships that could have appeared to influence the work
744 reported in this paper.

745 **Funding**

746 This work is partially supported by the INCA project, funded under the
747 EU Horizon Europe Programme, grant agreement n° 101061653.

748 **References**

- 749 [1] L. Monsees, Information disorder, fake news and the future of democ-
750 racy, *Globalizations* 20 (1) (2023) 153–168.
- 751 [2] C. Wardle, H. Derakhshan, Information disorder: Toward an interdis-
752 ciplinary framework for research and policymaking, Vol. 27, Council of
753 Europe Strasbourg, F-67075 Strasbourg Cedex, 2017.
- 754 [3] K. M. d. Treen, H. T. Williams, S. J. O’Neill, Online misinformation
755 about climate change, *Wiley Interdisciplinary Reviews: Climate Change*
756 11 (5) (2020) e665.
- 757 [4] Y. Wang, M. McKee, A. Torbica, D. Stuckler, Systematic literature
758 review on the spread of health-related misinformation on social media,
759 *Social science & medicine* 240 (2019) 112552.
- 760 [5] E. C. Tandoc Jr, Z. W. Lim, R. Ling, Defining “fake news” a typology
761 of scholarly definitions, *Digital journalism* 6 (2) (2018) 137–153.

- 762 [6] D. De Beer, M. Matthee, Approaches to identify fake news: a systematic
763 literature review, *Integrated science in digital age 2020* 136 (2021) 13–22.
- 764 [7] X. Zhou, R. Zafarani, A survey of fake news: Fundamental theories, de-
765 tection methods, and opportunities, *ACM Computing Surveys (CSUR)*
766 53 (5) (2020) 1–40.
- 767 [8] C. Chen, K. Shu, Can llm-generated misinformation be detected?
768 (2024). [arXiv:2309.13788](https://arxiv.org/abs/2309.13788).
- 769 [9] S. Branchetti, Testing attacks on structural text watermarking tech-
770 niques, Master’s thesis, Department of Computer Science and Engineer-
771 ing, University of Bologna (2023).
- 772 [10] S. Vosoughi, D. Roy, S. Aral, The spread of true and false news online,
773 *science* 359 (6380) (2018) 1146–1151.
- 774 [11] L. Toumanidis, R. Heartfield, P. Kasnesis, G. Loukas, C. Patrikakis,
775 A prototype framework for assessing information provenance in decen-
776 tralised social media: The eunomia concept, in: *International Confer-*
777 *ence on e-Democracy*, Springer, Springer, Cham, 2019, pp. 196–208.
- 778 [12] M. Choraś, M. Pawlicki, R. Kozik, K. Demestichas, P. Kosmides,
779 M. Gupta, Socialtruth project approach to online disinformation (fake
780 news) detection and mitigation, in: *Proceedings of the 14th Interna-*
781 *tional Conference on Availability, Reliability and Security*, Association
782 for Computing Machinery, New York, NY, USA, 2019, pp. 1–10.
- 783 [13] B. Yousuf, M. A. Qureshi, B. Spillane, G. Munnely, O. Carroll,
784 M. Runswick, K. Park, E. Culloty, O. Conlan, J. Suiter, Provenance:

- 785 An intermediary-free solution for digital content verification (2021).
786 arXiv:2111.08791.
- 787 [14] E. Strickland, This election year, look for content credentials: Media or-
788 ganizations combat deepfakes and disinformation with digital manifests,
789 IEEE Spectrum 61 (01) (2024) 24–27.
- 790 [15] I. Cox, M. Miller, J. Bloom, J. Fridrich, T. Kalker, Digital watermark-
791 ing and steganography, Morgan kaufmann, Burlington, Massachusetts,
792 United States, 2007.
- 793 [16] S. Kaur, G. Babbar, A zero-watermarking algorithm on multiple occur-
794 rences of letters for text tampering detection, International Journal on
795 Computer Science and Engineering 5 (5) (2013) 294.
- 796 [17] Y.-W. Kim, I.-S. Oh, Watermarking text document images using
797 edge direction histograms, Pattern Recognition Letters 25 (11) (2004)
798 1243–1251.
- 799 [18] D. Huang, H. Yan, Interword distance changes represented by sine waves
800 for watermarking text images, IEEE Transactions on Circuits and Sys-
801 tems for Video Technology 11 (12) (2001) 1237–1245.
- 802 [19] M. J. Atallah, V. Raskin, M. Crogan, C. Hempelmann, F. Kerschbaum,
803 D. Mohamed, S. Naik, Natural language watermarking: Design, analy-
804 sis, and a proof-of-concept implementation, in: Information Hiding: 4th
805 International Workshop, IH 2001 Pittsburgh, PA, USA, April 25–27,
806 2001 Proceedings 4, Springer, Springer, Berlin, Heidelberg, Berlin, Hei-
807 delberg, 2001, pp. 185–200.

- 808 [20] K. Lambrecht, A framework for the analysis of cleft constructions, *Linguistics* 39 (3) (2001) 463–516.
809
- 810 [21] H. M. Meral, B. Sankur, A. S. Özsoy, T. Güngör, Sevinç, Emre, Natural language watermarking via morphosyntactic alterations, *Computer
811 Speech & Language* 23 (1) (2009) 107–125.
812
- 813 [22] U. Topkara, M. Topkara, M. J. Atallah, The hiding virtues of ambiguity: quantifiably resilient watermarking of natural language text through
814 synonym substitutions, in: *Proceedings of the 8th workshop on Multimedia and security*, Association for Computing Machinery, New York,
815 NY, USA, 2006, pp. 164–174.
816
- 818 [23] O. Vybornova, B. Macq, A method of text watermarking using presuppositions, in: *Security, Steganography, and Watermarking of Multimedia Contents IX*, Vol. 6505, SPIE, SPIE, San Jose, CA, United States,
819 2007, pp. 613–622.
820
- 822 [24] O. Vybornova, B. Macq, Natural language watermarking and robust hashing based on presuppositional analysis, in: *2007 IEEE International
823 Conference on Information Reuse and Integration*, IEEE, IEEE, Las Vegas, NV, USA, 2007, pp. 177–182.
824
- 826 [25] M. Topkara, C. M. Taskiran, E. J. Delp III, Natural language watermarking, in: *Security, Steganography, and Watermarking of Multimedia
827 Contents VII*, Vol. 5681, SPIE, SPIE, San Jose, California, United States, 2005, pp. 441–452.
828
829

- 830 [26] L. Y. Por, K. Wong, K. O. Chee, Unispach: A text-based data hiding
831 method using unicode space characters, *Journal of Systems and Software*
832 85 (5) (2012) 1075–1082.
- 833 [27] R. A. Alotaibi, L. A. Elrefaei, Improved capacity arabic text watermark-
834 ing methods based on open word space, *Journal of King Saud University-*
835 *Computer and Information Sciences* 30 (2) (2018) 236–248.
- 836 [28] N. Mir, Copyright for web content using invisible text watermarking,
837 *Computers in Human Behavior* 30 (2014) 648–653.
- 838 [29] M. Hellmeier, H. Norkowski, E.-C. Schrewe, H. Qarawlus, F. Howar,
839 Innamark: A whitespace replacement information-hiding method, arXiv
840 preprint arXiv:2502.12710 (2025).
- 841 [30] A. Narayanan, J. Bonneau, E. Felten, A. Miller, S. Goldfeder, *Bitcoin*
842 *and cryptocurrency technologies: a comprehensive introduction*, Princeton
843 University Press, Princeton, New Jersey, United States, 2016.
- 844 [31] C. Di Ciccio, G. Meroni, P. Plebani, On the adoption of blockchain
845 for business process monitoring, *Software and Systems Modeling* 21 (3)
846 (2022) 915–937.
- 847 [32] F. Donini, A. Marcelletti, A. Morichetta, A. Polini, Coordinating rest
848 interactions in service choreographies using blockchain, *Blockchain: Re-*
849 *search and Applications* (2024) 100241.
- 850 [33] S. Paul, J. I. Joy, S. Sarker, S. Ahmed, A. K. Das, et al., Fake news de-
851 tection in social media using blockchain, in: 2019 7th international Con-

- 852 ference on smart computing & communications (ICSCC), IEEE, IEEE,
853 Sarawak, Malaysia, 2019, pp. 1–5.
- 854 [34] Z. Shae, J. Tsai, Ai blockchain platform for trusting news, in: 2019
855 IEEE 39th International Conference on Distributed Computing Systems
856 (ICDCS), IEEE, IEEE, Dallas, TX, USA, 2019, pp. 1610–1619.
- 857 [35] W. Shang, M. Liu, W. Lin, M. Jia, Tracing the source of news based on
858 blockchain, in: 2018 IEEE/ACIS 17th International Conference on Com-
859 puter and Information Science (ICIS), IEEE, IEEE, Singapore, 2018, pp.
860 377–381.
- 861 [36] M. Arquam, A. Singh, R. Sharma, A blockchain-based secured and
862 trusted framework for information propagation on online social net-
863 works, *Social Network Analysis and Mining* 11 (1) (2021) 49.
- 864 [37] E. Sengupta, R. Nagpal, D. Mehrotra, G. Srivastava, Problock: a
865 novel approach for fake news detection, *Cluster Computing* 24 (2021)
866 3779–3795.
- 867 [38] A. Alexandrescu, C. N. Butincu, Decentralized news-retrieval architec-
868 ture using blockchain technology, *Mathematics* 11 (21) (2023) 4542.
- 869 [39] A. R. Faridi, R. Singh, F. Masood, M. Y. Salmony, Machine learn-
870 ing based novel framework for fake news detection and prevention us-
871 ing blockchain, in: 2023 10th International Conference on Comput-
872 ing for Sustainable Global Development (INDIACom), IEEE, 2023, pp.
873 751–755.

- 874 [40] P. Rani, V. Jain, J. Shokeen, A. Balyan, Blockchain-based rumor de-
875 tection approach for covid-19, *Journal of Ambient Intelligence and Hu-*
876 *manized Computing* 15 (1) (2024) 435–449.
- 877 [41] A. D. Dwivedi, R. Singh, S. Dhall, G. Srivastava, S. K. Pal, Tracing the
878 source of fake news using a scalable blockchain distributed network, in:
879 2020 IEEE 17th international conference on mobile ad hoc and sensor
880 systems (MASS), IEEE, 2020, pp. 38–43.
- 881 [42] S. Dhall, A. D. Dwivedi, S. K. Pal, G. Srivastava, Blockchain-based
882 framework for reducing fake or vicious news spread on social media/mes-
883 saging platforms, *Transactions on Asian and Low-Resource Language*
884 *Information Processing* 21 (1) (2021) 1–33.
- 885 [43] X. Li, L. Wei, L. Wang, Y. Ma, C. Zhang, M. Sohail, A blockchain-based
886 privacy-preserving authentication system for ensuring multimedia con-
887 tent integrity, *International journal of intelligent systems* 37 (5) (2022)
888 3050–3071.
- 889 [44] T. K. Araghi, D. Megías, V. Garcia-Font, M. Kuribayashi, W. Mazur-
890 czyk, Disinformation detection and source tracking using semi-fragile
891 watermarking and blockchain, in: *Proceedings of the 2024 European*
892 *Interdisciplinary Cybersecurity Conference*, 2024, pp. 136–143.
- 893 [45] Q.-u.-A. Mastoi, M. F. Memon, S. Jan, A. Jamil, M. Faique, Z. Ali,
894 A. Lakhan, T. A. Syed, Enhancing deepfake content detection through
895 blockchain technology, *Int. J. Adv. Comput. Sci. Appl.* 16 (6) (2025).

- 896 [46] Z. Jiang, H. Wang, S. Han, A robust pdf watermarking scheme with
897 versatility and compatibility, *Multimedia Tools and Applications* 83 (24)
898 (2024) 64341–64367.
- 899 [47] S. N. KL, B. K. R, Text steganography: enhanced character-level embed-
900 ding algorithm using font attribute with increased resilience to statistical
901 attacks, *Multimedia Tools and Applications* (2024) 1–26.
- 902 [48] F. N. Al-Wesabi, F. Alrowais, H. G. Mohamed, M. Al Duhayyim, A. M.
903 Hilal, A. Motwakel, Heuristic optimization algorithm based watermark-
904 ing on content authentication and tampering detection for english text,
905 *IEEE Access* 11 (2023) 86104–86111.
- 906 [49] C. Stomeo, Text watermarking e social network: uno studio sperimen-
907 tale, Master’s thesis, Department of Computer Science and Engineering,
908 University of Bologna (2016).
- 909 [50] M. N. Jones, D. J. Mewhort, Case-sensitive letter and bigram frequency
910 counts from large-scale english corpora, *Behavior research methods, in-
911 struments, & computers* 36 (3) (2004) 388–396.
- 912 [51] Panckhurst, Rachel and Détrie, Catherine and Lopez, Cédric and Moïse,
913 Claudine and Roche, Mathieu and Verine, Bertrand, 88milSMS. A cor-
914 pus of authentic text messages in French, *Banque de corpus CoMeRe*.
915 Chanier T.(éd)-Ortolang: Nancy (2016).
- 916 [52] Kupietz, Marc and Lungen, Harald and Kamocki, Paweł and Witt, An-
917 dreas, The German reference corpus DeReKo: New developments–new

- 918 opportunities, in: Proceedings of the eleventh international conference
919 on language resources and evaluation (LREC 2018), 2018.
- 920 [53] Landro, Nicola and Gallo, Ignazio and La Grassa, Riccardo and Fed-
921 erici, Edoardo, Two new datasets for italian-language abstractive text
922 summarization, *Information* 13 (5) (2022) 228.
- 923 [54] M. Bashardoost, M. S. Mohd Rahim, T. Saba, A. Rehman, Replacement
924 attack: A new zero text watermarking attack, *3D Research* 8 (1) (2017)
925 8.
- 926 [55] M. Reglitz, Fake news and democracy, *J. Ethics & Soc. Phil.* 22 (2022)
927 162.
- 928 [56] D. Gordon, Targeted systems and democracy: Russia, iran, and china's
929 cyber threats and disinformation campaigns to weaken and undermine
930 western democracies, Master's thesis, Utica College (2020).
- 931 [57] N. Newman, R. Fletcher, C. T. Robertson, A. R. Arguedas, R. K.
932 Nielsen, Reuters institute digital news report 2024, Reuters Institute
933 for the study of Journalism (2024).
- 934 [58] N. Marchal, B. Kollanyi, L.-M. Neudert, P. N. Howard, Junk news dur-
935 ing the eu parliamentary elections: Lessons from a seven-language study
936 of twitter and facebook, University of Oxford (2019).
- 937 [59] J. Zhou, Y. Zhang, Q. Luo, A. G. Parker, M. De Choudhury, Syn-
938 thetic lies: Understanding ai-generated misinformation and evaluating
939 algorithmic and human solutions, in: Proceedings of the 2023 CHI Con-
940 ference on Human Factors in Computing Systems, 2023, pp. 1–20.

- ⁹⁴¹ [60] Y. N. Harari, Nexus: A Brief History of Information Networks from the
⁹⁴² Stone Age to AI, Random House, 1745 Broadway New York, NY 10019,
⁹⁴³ 2024.

Journal Pre-proof

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.

Journal Pre-proof

Flavio Bertini: Conceptualization, Methodology, Validation, Writing – original draft, Writing – review & editing.

Alessandro Benetton: Conceptualization, Software, Investigation.

Danilo Montesi: Conceptualization, Writing – review & editing, Supervision.

Journal Pre-proof