
Unbiased Watermark for Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The recent advancements in large language models (LLMs) have sparked a growing
2 apprehension regarding the potential misuse. One approach to mitigating this risk
3 is to incorporate watermarking techniques into LLMs, allowing for the tracking and
4 attribution of model outputs. This study examines a crucial aspect of watermark-
5 ing: how significantly watermarks impact the quality of model-generated outputs.
6 Previous studies have suggested a trade-off between watermark strength and out-
7 put quality. However, our research demonstrates that it is possible to integrate
8 watermarks without affecting the output probability distribution with appropriate
9 implementation. We refer to this type of watermark as an **unbiased watermark**.
10 This has significant implications for the use of LLMs, as it becomes impossible
11 for users to discern whether a service provider has incorporated watermarks or not.
12 Furthermore, the presence of watermarks does not compromise the performance
13 of the model in downstream tasks, ensuring that the overall utility of the language
14 model is preserved. Our findings contribute to the ongoing discussion around
15 responsible AI development, suggesting that unbiased watermarks can serve as
16 an effective means of tracking and attributing model outputs without sacrificing
17 output quality.

18 1 Introduction

19 In recent years, large language models (LLMs) [19, 39, 40] have become an indispensable tool for a
20 wide range of tasks, including text generation [27, 10], translation [7, 5], summarization [36], etc.
21 With the escalating misuse of LLMs, such as plagiarism, tracking the usage of text generated by
22 machines has become increasingly important. One viable method to monitor the usage of LLMs
23 is watermarking [20, 32, 59], which embeds imperceptible information within the generated text,
24 thereby allowing for efficient detection and tracking of the model’s potential abuse.

25 Watermarking techniques can serve multiple purposes, such as embedding ownership information
26 within the generated text to protect the intellectual property rights of the model. It can also help
27 mitigate potential harm caused by LLMs by monitoring where the model is being used and whether it
28 is being misused or abused.

29 A good watermarking method should not adversely affect the normal usage of the language model or
30 degrade the quality of the generated text. However, a prevailing belief holds that there is an inevitable
31 trade-off between the strength of the watermark and the quality of the output text. For instance,
32 recent work by Kirchenbauer et al. [32] introduced a method that augmented the logits of a randomly
33 selected set of "green" tokens. By tuning the "magnitude of logits adjustment", they demonstrated a
34 trade-off between watermark strength and text quality.

35 Our primary contribution is to challenge this conventional wisdom. We show that with the right
36 implementation, watermarking can be accomplished without affecting the output quality. We refer to
37 this particular type of watermark as an **unbiased watermark**. We approach the problem of output
38 quality degradation from the perspective of watermark detection. We posit that if the watermark

39 causes a decline in output quality, there should be a method to guess the presence of the watermark
40 based on the quality. Conversely, if the watermark cannot be detected, it implies that the output
41 quality remains unaffected. Specifically, we provide a proof that with a suitable implementation,
42 watermarking does not affect the output probability distribution. This has significant implications,
43 as users who do not have the private key are unable to discern whether a service provider has
44 applied watermarking to the model. Furthermore, the addition of watermarking does not affect
45 the performance of the generated text in any downstream tasks. **Our main contributions can be**
46 **summarized as follows:**

- 47 • We introduce *unbiased watermark*, an innovative family of watermark methods that guarantee the
48 non-degradation of text quality. In addition, we offer a comprehensive framework that facilitates
49 the design and detection of unbiased watermarks.
- 50 • We propose two innovative and practical watermarking techniques known as δ -reweight and
51 γ -reweight. Through extensive experimentation, we demonstrate that these techniques preserve
52 output quality in machine translation and text summarization tasks.
- 53 • We develop an advanced maximin variant of the original log-likelihood ratio test for watermark
54 detection. This novel detection method comes with theoretical guarantees, specifically an upper
55 bound on type I error, thus enhancing the reliability of watermark detection in language models.

56 2 Preliminary

57 In this section, we delve into the problem of watermarking in the context of LLMs. We begin by
58 setting up the problem and defining essential concepts.

59 **Problem Modeling:** We first introduce several notations to formalize the problem. Let Σ denote the
60 vocabulary set, which is the set of all possible tokens an LLM can generate in a single step. We then
61 define the set Σ^* as the collection of all possible strings of any length, including those of length zero.

62 An LLM generates a sequence of tokens conditioned on a given context. In a single step, the
63 probability of generating the next token $x_{n+1} \in \Sigma$ given the current context, x_1, x_2, \dots, x_n , can be
64 denoted as $P_M(x_{n+1} \mid x_1, x_2, \dots, x_n)$. The LLM operates in an autoregressive fashion, which means
65 the joint probability of generating multiple tokens x_{n+1}, \dots, x_{n+m} can be written as:

$$P_M(x_{n+1}, \dots, x_{n+m} \mid x_1, x_2, \dots, x_n) = \prod_{i=1}^m P_M(x_{n+i} \mid x_1, x_2, \dots, x_n, x_{n+1}, \dots, x_{n+i-1}).$$

66 For simplicity, we use the following notation: $P_M(\mathbf{x}_{n+1:n+m} \mid \mathbf{x}_{1:n})$, where $\mathbf{x}_{n+1:n+m} =$
67 $(x_{n+1}, \dots, x_{n+m}) \in \Sigma^*$.

68 In the context of watermarking, we introduce a service provider that holds a private key k from the key
69 space K . The key $k \in K$ is chosen at random from the prior distribution $P_K(k)$. The watermarked
70 output of the LLM follows distribution $P_{M,w}(x_{n+1} \mid x_1, x_2, \dots, x_n; k)$, which is conditioned on both
71 the key k and the context $\mathbf{x}_{1:n}$. Similarly, we use the notation $P_{M,w}(\mathbf{x}_{n+1:n+m} \mid \mathbf{x}_{1:n}; k)$ for the
72 probability of generating a sequence of tokens in a watermarked model.

73 **Objective.** Our goal is to devise a watermarking scheme that: a) is efficiently detectable by the
74 service provider; b) can't be detected by users and does not negatively impact the quality of the
75 output.

76 The reason we focus on the detection of watermarks by users is that it is closely related to the output
77 quality. If the watermark causes a degradation in the output quality, there should exist a method
78 to infer the presence of the watermark by examining the quality. Conversely, if the watermark is
79 undetectable, it implies that it does not impact the output quality.

80 From a statistical testing perspective, a watermark is considered strictly undetectable if the probability
81 distributions of the watermarked and non-watermarked outputs are identical. To capture this notion,
82 we define several desirable properties of watermarking schemes.

83 **Definition 1** (*n*-shot-undetectable). *For a fixed input sequence $\mathbf{a} \in \Sigma^*$, we say that watermarked*
84 *LLM and key prior pair $(P_{M,w}, P_K)$ is *n*-shot-undetectable compared to original LLM P_M if*

$$\prod_{i=1}^n P_M(\mathbf{x}^i \mid \mathbf{a}) = \sum_{k \in K} P_K(k) \prod_{i=1}^n P_{M,w}(\mathbf{x}^i \mid \mathbf{a}; k), \quad \text{for any } n \text{ number of strings } \mathbf{x}^i \in \Sigma^*.$$

85 **Definition 2** (downstream-invariant). We say the watermarked LLM and key prior pair $(P_{M,w}, P_K)$
 86 are invariant compared to original LLM P_M on downstream tasks iff

$$\mathbb{E}_{\mathbf{x} \sim P_{M,w}(\cdot | \mathbf{a}; k), k \sim P_K} [f(\mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim P_M(\cdot | \mathbf{a})} [f(\mathbf{x})],$$

87 for any strings $\mathbf{x}, \mathbf{a} \in \Sigma^*$, and for any metric $f : \Sigma^* \rightarrow \mathbb{R}$.

88 Note that the one-shot-undetectable property implies the downstream invariant property. Interestingly,
 89 this implication does not require the n -shot-undetectable property for $n > 1$, which means a water-
 90 marking scheme that is one-shot-undetectable can still maintain the output quality for downstream
 91 tasks even if the user might discern the existence of the watermark through multiple generation
 92 requests.

93 In summary, we have outlined the preliminary concepts and objectives for developing a watermarking
 94 scheme for LLMs. We highlight the desired properties of n -shot-undetectability and downstream
 95 invariance, as they provide a rigorous theoretical guarantee of quality preservation and integrity in
 96 the deployment of watermark schema. In Section 4, we will present a watermark framework that is
 97 provably n -shot-undetectable for any given integer $n \geq 1$.

98 3 Warm up: undetectability in a simplified toy environment

99 In this subsection, we aim to prove the feasibility of undetectability in a highly simplified toy
 100 environment. This preliminary analysis serves as a foundation for understanding the more complex
 101 scenarios that follow.

102 **Settings.** Consider a service provider that offers a random number generation service. The service
 103 outputs a uniformly distributed random number in the set $\{0, 1\}$. The clean generation process can
 104 be represented as $P_M(x) = 1/2, \forall x \in \{0, 1\}$. We assume that the key k belongs to the set $\{0, 1\}$
 105 and is selected with equal probability. With the watermark added, the probability of the new output
 106 can be expressed as: $P_{M,w}(x | k) = \delta_k(x)$.

107 Recall that the one-shot-undetectable property can be represented as $P_M(x) = \sum_{k \in K} P_{M,w}(x |$
 108 $k)P_K(k)$. Suppose that a user can only make a single request to the service. If the user is unaware
 109 of the key, the user will be unable to distinguish whether the received result is watermarked or not.
 110 Therefore, in this simplified scenario, the undetectability of the watermark is achieved.

111 However, there is a considerable gap between this toy example and the practical implementation of
 112 watermarking in LLMs. Firstly, the symbol set Σ in LLMs is far more complex than the binary set
 113 $\{0, 1\}$, and the probability distribution is not uniform. Besides, the generation process in LLMs is
 114 autoregressive, which means that more than one symbol are generated iteratively. Furthermore, the
 115 toy example does not satisfy the n -shot-undetectable property for $n > 1$.

116 Despite these differences, this simple example provides essential insights that help in understanding
 117 the following sections where we address these challenges. The underlying principles of undetectability
 118 remain constant, while their application becomes more intricate in a more complex environment.

119 4 Watermarking with unbiased reweighting

120 In this section, we build upon the intuition from the previous section and extend the approach to
 121 LLMs' generation. The section is structured as follows: Section 4.1 introduces a fundamental
 122 mathematical tool for addressing the reweighting problem in general discrete probability distributions.
 123 Section 4.2 applies the reweighting technique to LLMs. Section 4.3 presents the final framework.

124 4.1 Distribution reweighting

125 In its most general form, we consider a random watermark code E and a reweight function $R_E :$
 126 $\Delta_\Sigma \rightarrow \Delta_\Sigma$, which depends on the random watermark code E . The set of all possible probability
 127 distributions on the symbol set Σ is denoted as Δ_Σ , which forms a simplex.

128 **Definition 3.** A **reweighting function** is a tuple (\mathcal{E}, P_E, R) where \mathcal{E} is called the watermark code
 129 space, P_E is a probability distribution on space \mathcal{E} , and R is a function $R : \mathcal{E} \times \Delta_\Sigma \rightarrow \Delta_\Sigma$.
 130 For a specific watermark code $E \in \mathcal{E}$, we denote the partially evaluated reweighting function as
 131 $R_E : \Delta_\Sigma \rightarrow \Delta_\Sigma$.

132 **Definition 4.** Given a random watermark code E and a reweighting function $R_E : \Delta_\Sigma \rightarrow \Delta_\Sigma$, we
 133 say that R is an **unbiased reweighting function** if and only if for all $P \in \Delta_\Sigma$, $\mathbb{E}_E[R_E(P)] = P$.

134 **4.1.1 Existing reweighting methods**

135 Kirchenbauer et al. [32] essentially comprise two reweighting methods in their work, but neither of
 136 them satisfies the unbiased property.

137 Both methods have \mathcal{E} as the set of mappings $f : \Sigma \rightarrow \{\text{red}, \text{green}\}$, such that f maps half of the
 138 tokens in Σ to ‘red’ and the other half to ‘green’, and P_E as a uniform distribution. Therefore, the
 139 random watermark code E assigns each symbol to either *red* or *green*. The ‘‘Hard Red List’’ method
 140 sets the probability of all red symbols to zero and renormalizes the probabilities of the remaining
 141 vocabulary. The second method is ‘‘Soft Red List’’ blocking, where they randomly select the same
 142 ‘‘Red List’’ as the first method and decrease the corresponding probability for red symbols by adding a
 143 constant δ to the logits of the green symbols, then apply softmax to obtain the final probabilities.

144 **4.1.2 Unbiased reweighting methods**

145 In this section, we present two reweighting methods that satisfy the unbiased property.

146 **δ -reweight:** Let the watermark code space \mathcal{E} be the interval $[0, 1]$, and let P_E be the uniform
 147 probability on \mathcal{E} . Leveraging *Inverse Transform Sampling*¹ [14], we can sample from distribution
 148 $P \in \Delta_\Sigma$ using a uniformly distributed random number in $[0, 1]$. Therefore, we have a mapping
 149 $\text{sampling}_P : \mathcal{E} \rightarrow \Sigma$. The δ -reweight just returns a delta distribution $R_E(P) = \delta_{\text{sampling}_P(E)}$.

150 It is important to note that while the reweighted distribution for each individual random event E
 151 is a delta distribution, the mean output token probabilities remain the original distribution P when
 152 considering the randomness of E .

153 **γ -reweight:** Let the watermark code space \mathcal{E} be the set of all bijective function between vocabularies
 154 set Σ and a set of indices $[\Sigma] = \{1, \dots, |\Sigma|\}$, where $|\Sigma|$ is the size of vocabularies set Σ . Essentially,
 155 any watermark code E is an indexing function for vocabularies set Σ , and is also equivalent to a total
 156 order on Σ . Let P_E be the uniform probability on \mathcal{E} , it is easy to sample a watermark code E by
 157 randomly shuffling the symbol list.

158 Assume the original distribution is $P_T(t) \in \Delta_\Sigma, \forall t \in \Sigma$. Given the watermark code $E : \Sigma \rightarrow [\Sigma]$,
 159 we construct auxiliary functions $F_I(i) = \sum_{t \in \Sigma} \mathbf{1}(E(t) \leq i) P_T(t)$, $F_S(s) = \max(2s - 1, 0)$,
 160 $F_{I'}(i) = F_S(F_I(i))$. The γ -reweight yields new distribution $P_{T'}(t) = F_{I'}(E(t)) - F_{I'}(E(t) - 1)$.

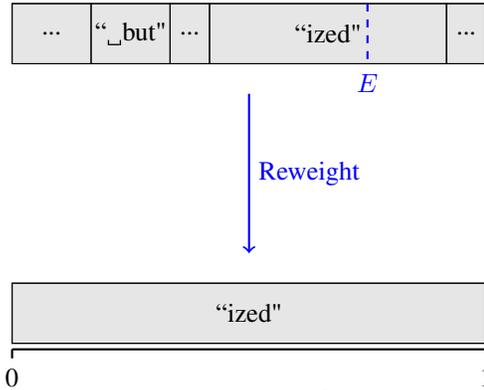


Figure 1: Illustration of δ -reweight.

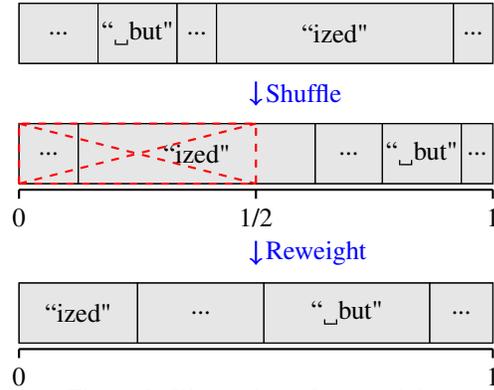


Figure 2: Illustration of γ -reweight.

161 We provide illustrations of the δ -reweight and γ -reweight methods in Figures 1 and 2. Each block
 162 represents a token, and the width represents the probability of that token, so the total length is 1. The
 163 left panel shows the δ -reweight method, where each individual random watermark code $E \in [0, 1]$
 164 uniformly sampled from interval $[0, 1]$ corresponds to a specific token according to the horizontal axis,
 165 and the reweighted distribution is just a δ distribution on that token, such that the selected token has 1
 166 probability, and all other vocabulary tokens have a probability of 0. The right panel demonstrates the
 167 γ -reweight method. First, the symbol set is shuffled. Then, the left half of the regions are rejected,
 168 and the remaining regions are amplified with a factor of 2.

169 Both methods are unbiased¹ when considering the randomness of the watermark code E . For δ -
 170 reweight, we can see that by noticing that the probability of returning a δ distribution on a token is

¹Detailed definition and rigorous proof can be found in Appendix D

171 just the original probability on that token, therefore the weighted average of all delta distributions is
 172 still the original probability. In the case of γ -reweight, although certain regions are rejected and the
 173 other regions are amplified, every token has the same probability to be in the rejected or amplified
 174 region, thus ensuring the unbiased property.

175 4.2 Reweighting for autoregressive model

176 The reweighting methods presented in the previous section can be applied to single token-generation
 177 directly. Given a prefix $\mathbf{x}_{1:n}$, the probability distribution for generating a new token without a
 178 watermark is denoted as $P_M(\cdot|\mathbf{x}_{1:n}) \in \Delta_\Sigma$. For a random watermark code E , we sample from a
 179 new distribution $P_{M,w}(\cdot|\mathbf{x}_{1:n}) = R_E(P_M(\cdot|\mathbf{x}_{1:n})) \in \Delta_\Sigma$. If the reweighting function is unbiased,
 180 we have $\mathbb{E}_E[R_E(P_M(\cdot|\mathbf{x}_{1:n}))] = P_M(\cdot|\mathbf{x}_{1:n})$. This ensures that, for an individual unaware of
 181 the watermark code, it is impossible to determine whether a new token is sampled directly from
 182 $P_M(\cdot|\mathbf{x}_{1:n})$ or from $P_{M,w}(\cdot|\mathbf{x}_{1:n}; E)$ for a random watermark E . However, if the watermark code is
 183 known, one can perform statistical hypothesis testing to determine the likelihood of a token being
 184 sampled from either distribution.

185 The main challenge now is constructing the watermark code E . Since the LLM generation task is
 186 autoregressive, multiple reweighting steps are required, with each step needing a watermark code E_i
 187 for reweighting the distribution of token x_i .

188 4.2.1 Independence of watermark codes

189 It is crucial that E_i values are independent to ensure the unbiased nature of the entire sequence, rather
 190 than just the single-token generation process.

191 **Theorem 5.** *Given an unbiased reweighting function (\mathcal{E}, P_E, R) , if E_i values are i.i.d. with the*
 192 *distribution P_E , we have: $\mathbb{E}_{E_1, \dots, E_n}[P_{M,w}(\mathbf{x}_{1:n}|\mathbf{a}_{1:m})] = P_M(\mathbf{x}_{1:n}|\mathbf{a}_{1:m})$.*

193 If the E_i values are not independent, we cannot guarantee that the generation probability of the entire
 194 sequence remains unbiased. As an extreme example, consider a case where all E_i values are identical.
 195 Referring to the random bit example in the previous section, assume that the correct distribution is
 196 a sequence where each token is a random 0 or 1 with equal probability. Identical E_i values would
 197 result in identical token outputs, ultimately producing sequences consisting solely of 0's or 1's, which
 198 is clearly biased.

199 4.2.2 Context code

200 To construct a large number of independent watermark codes E_i during watermarking and to know
 201 the used E_i values during watermark detection, we follow an approach similar to Kirchenbauer et al.
 202 [32] by combining the information from the prefix and a secret key to construct E_i .

203 For a single token generation process, given a prefix x_1, x_2, \dots, x_n , we consider an abstract context
 204 code space C and an abstract context code generation function $cc : \Sigma^* \rightarrow C$. Based on the prefix,
 205 we construct the context code $c_{n+1} = cc(x_1, x_2, \dots, x_n)$. Specific examples include using the entire
 206 prefix $c_{n+1} = (x_1, x_2, \dots, x_n)$, and using the m most recent prefixes $c_{n+1} = (x_{n-m+1}, \dots, x_n)$. Our
 207 comprehensive framework accommodates diverse context code generation approaches, particularly
 208 those that integrate error-correcting mechanisms to augment watermark resilience in the face of text
 209 manipulation attacks. Nevertheless, we refrain from delving into these strategies within the confines
 210 of this paper and consider it a subject for subsequent investigation.

211 The final watermark code is defined as $E_i = \hat{E}(c_i, k)$, using a watermark code generation function
 212 $\hat{E} : C \times K \rightarrow \mathcal{E}$.

213 **Definition 6.** *Given an unbiased reweighting function (\mathcal{E}, P_E, R) and a context code space C , an*
 214 *unbiased watermark code generation function is a tuple $(\mathcal{E}, P_E, R, C, K, P_K, \hat{E})$ that satisfies:*

- 215 1. *Unbiasedness:* $\mathbb{E}_{k \sim P_K}[R_{\hat{E}(c,k)}(P)] = P, \forall P \in \Delta_\Sigma, \forall c \in C$.
- 216 2. *Independence:* For any n distinct $c_1, \dots, c_n \in C$, the values $R_{\hat{E}(c_i,k)}(P)$ are mutually
 217 independent.

218 **Theorem 7.** *For any unbiased reweighting function and context code space, an unbiased watermark*
 219 *code generation function always exists.*

220 In practice, pseudorandom numbers can be used to implement the unbiased watermark code generation
 221 function in the above theorem. Specifically, the hash value $\text{hash}(c, k)$ can be used as a random seed

222 to sample E from P_E as an implementation of $E = \hat{E}(c, k)$. In this paper, we employ SHA-256 for
 223 hash function and a 1024-bit random bitstring as the key k .

224 An unbiased watermark code generation function ensures that watermark codes E_i are independent
 225 with each other if only their context codes are different. During the generation of a sequence,
 226 context codes may be repeated, although this is a rare event in practice. If c_i and c_j are equal,
 227 then E_i and E_j are also equal, violating the independence of E_i . A simple workaround is to skip
 228 reweighting for a token when encountering a previously used context code. In other words, we set
 229 $P_{M,w}(\cdot | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1}) = P_M(\cdot | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1})$ if the context code has appeared before.

230 4.3 Framework

Algorithm 1 Watermarking framework

```

1: Input: key for watermark  $k \in K$ , prompt  $\mathbf{a}_{1:m} \in \Sigma^*$ , generate length  $n \in \mathbb{N}$ , initial code
   history  $cch \in 2^C$ , context code function  $cc : \Sigma^* \rightarrow C$ , watermark code generation function
    $\hat{E} : C \times K \rightarrow \mathcal{E}$ , and reweighting function  $R : \mathcal{E} \times \Delta_\Sigma \rightarrow \Delta_\Sigma$ .
2: for  $t = 1, \dots, n$  do
3:    $P_i \leftarrow P_M(\cdot | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1})$  ▷ original distribution
4:    $c_i \leftarrow cc(\cdot | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1})$  ▷ context code
5:   if  $c_i \in cch$  then
6:      $Q_i \leftarrow P_i$  ▷ skip the reweighting
7:   else
8:      $cch \leftarrow cch \cup \{c_i\}$  ▷ record history
9:      $E_i \leftarrow \hat{E}(c_i, k)$  ▷ watermark code
10:     $Q_i \leftarrow R_{E_i}(P_i)$  ▷ reweighted distribution
11:    Sample the next token  $x_i$  using distribution  $Q_i$ 
12: return  $\mathbf{x}_{1:n}$ 

```

231 Integrating the tools discussed earlier, we present a general framework for watermarking here. The
 232 algorithm for this framework is outlined in Algorithm 1.

233 We note that our abstract framework requires the specification of two key components in order to be
 234 practically implemented: the unbiased reweight function R_E and the context code function cc .

235 5 Statistical hypothesis testing for watermark detection

236 In the previous section, we discussed the process of adding a watermark to a text based on a secret
 237 key k and a given prompt $\mathbf{a}_{1:m}$. The watermark-embedded text can be sampled from the distribution
 238 $P_{M,w}(\mathbf{x}_{1:n} | \mathbf{a}_{1:m}; k)$. In this section, we focus on the watermark detection task, which is the inverse
 239 problem of watermark embedding.

240 Given a text $\mathbf{x}_{1:n}$, the goal of watermark detection is to infer whether it is more likely to be generated
 241 from the unmarked distribution $P_M(\mathbf{x}_{1:n} | \mathbf{a}_{1:m})$ or the marked distribution $P_{M,w}(\mathbf{x}_{1:n} | \mathbf{a}_{1:m}; k)$.
 242 This problem can be formulated as a statistical hypothesis test between two competing hypotheses:
 243 H_0 , which posits that $\mathbf{x}_{1:n}$ follows the unmarked distribution, and H_1 , which posits that $\mathbf{x}_{1:n}$ follows
 244 the marked distribution.

245 5.1 Score-based testing

246 We focus on a particular kind of score-based testing, which assigns a score to each token in the text.
 247 The score can be interpreted as the confidence that the token was generated by the watermark model
 248 rather than the original model. Scores s_i can be computed based on $\mathbf{x}_{1:i}$, in accordance with the
 249 autoregressive manner of the generation process.

250 The total score S is given by $S = \sum_{i=1}^n s_i$. A threshold \hat{S} is set such that if $S < \hat{S}$, the null
 251 hypothesis H_0 is accepted, indicating insufficient evidence to conclude that the text contains a
 252 watermark. Otherwise, the null hypothesis is rejected. There are two types of error probabilities
 253 associated with this decision process: Type I error, which is the probability of incorrectly rejecting

254 the null hypothesis under H_0 , denoted as $P_{H_0}(S \geq \hat{S})$, and Type II error, which is the probability of
 255 incorrectly accepting the null hypothesis under H_1 , denoted as $P_{H_1}(S < \hat{S})$.

256 To derive theoretical results, we require the scores to have a specific property: under the null
 257 hypothesis H_0 , the exponential momentum of s_i is bounded, conditioned on the preceding context
 258 $\mathbf{x}_{1,i-1}$. This requirement leads to an upper bound on α , the Type I error probability.

259 To derive theoretical results, we require that the scores have a particular property: the exponential
 260 moment of s_i under H_0 should be bounded, conditioned on the previous text $\mathbf{x}_{1,i-1}$. This requirement
 261 leads to an upper bound on the Type I error rate.

262 **Theorem 8.** *Given a probability space (Ω, \mathcal{A}, P) and a Σ -valued stochastic process $x_i : 1 \leq i \leq n$,
 263 as well as an \mathbb{R} -valued stochastic process $s_i : 1 \leq i \leq n$, let $\mathcal{F}_i^x := \sigma(x_j \mid 1 \leq j \leq i)$ and
 264 $\mathcal{F}_i^s := \sigma(s_j \mid 1 \leq j \leq i)$ be the corresponding filtrations, where $\sigma(\cdot)$ denotes the σ -algebra
 265 generated by random variables. If $\mathcal{F}_i^s \subseteq \mathcal{F}_i^x$ and $\mathbb{E}[\exp(s_i) | \mathcal{F}_{i-1}^x] \leq 1$, then $P(\sum_{i=1}^n s_i \geq t) \leq e^{-t}$.*

266 Therefore, to ensure that the Type I error probability has an upper bound α , we can set the threshold
 267 \hat{S} as $\hat{S} = -\log(\alpha)$. In the following, we discuss two special scores.

268 5.2 Log likelihood ratio (LLR) score

269 According to the Neyman-Pearson lemma, the likelihood ratio test is the most powerful test among
 270 all tests with the same Type I error rate. Specifically, the log-likelihood ratio (LLR) score is defined
 271 as $s_i = \log \frac{P_{M,w}(x_i | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1}; k)}{P_M(x_i | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1})}$, and the total score becomes $S = \log \frac{P_{M,w}(\mathbf{x}_{1:n} | \mathbf{a}_{1:m}; k)}{P_M(\mathbf{x}_{1:n} | \mathbf{a}_{1:m})}$.

272 We now provide an optimization derivation of the above s_i to gain intuition and set the foundation
 273 for the maximin variant of the LLR score in the next section. Let $P_i = P_M(\cdot | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1})$,
 274 $Q_i = P_{M,w}(\cdot | \mathbf{a}_{1:m}, \mathbf{x}_{1:i-1}; k)$, and let $s_i = S_i(x_i) \in \mathbb{R}$ denote the score corresponding to different
 275 x_i . Note that P_i , Q_i , and S_i are all functions with signature $\Sigma \rightarrow \mathbb{R}$, therefore equivalent to vectors
 276 of dimension $|\Sigma|$. We can define the inner product as $\langle P_i, S_i \rangle = \sum_{x \in \Sigma} P_i(x) S_i(x)$.

277 The requirement $\mathbb{E}[\exp(s_i) | \mathcal{F}_{i-1}^x] \leq 1$ can be reformulated as $\langle P_i, \exp(S_i) \rangle \leq 1$, where the expo-
 278 nential function is applied element-wise. Instead of minimizing the Type II error directly, we aim to
 279 maximize the average score under H_1 , i.e., $\langle Q_i, S_i \rangle$.

280 The optimization problem becomes $\max_{S_i} \langle Q_i, S_i \rangle$, s.t. $\langle P_i, \exp(S_i) \rangle \leq 1$. The optimal solution is
 281 given by $S_i(x) = \log \frac{Q_i(x)}{P_i(x)}$, which recovers the optimal log likelihood ratio score.

282 5.3 Maximin variant of LLR score

283 One major limitation of the LLR score described in the previous section is that when $Q_i(x) = 0$,
 284 $S_i(x) = -\infty$. This means that as long as a single token does not come from the watermark model
 285 $P_{M,w}$, the score becomes negative infinity, making it impossible to reject the null hypothesis H_0 .

286 A more general reason for this issue is that the watermark model $P_{M,w}$ used in the detection process
 287 may not exactly match the true distribution of the watermarked text. In practice, potential sources of
 288 discrepancy include editing (e.g., a text sampled from $P_{M,w}$ may undergo some degree of editing
 289 before being watermark detection) and imperfect estimation of the generation process (e.g., due to
 290 lack of knowledge of the exact prompt and temperature used during generation).

291 To address this problem, we consider a perturbed generation distribution. Instead of the original
 292 hypothesis H_1 , where $\mathbf{x}_{1:n}$ follows the watermark distribution $P_{M,w}$, we now assume that $\mathbf{x}_{1:n}$
 293 follows a distribution $P'_{M,w}$, which is similar to but not identical to $P_{M,w}$. Specifically, during the
 294 generation of each token, the total variation (TV) distance between Q'_i and Q_i is bounded by d .

295 The corresponding new optimization problem is

$$\max_{S_i} \min_{Q'_i \in \Delta_{\Sigma}, TV(Q'_i, Q_i) \leq d} \langle Q'_i, S_i \rangle, \quad \text{s.t. } \langle P_i, \exp(S_i) \rangle \leq 1.$$

296 Intuitively, the optimal solution for Q'_i in the inner optimization decreases $Q'_i(x)$ when $S_i(x)$ is large
 297 and increases $Q'_i(x)$ when $S_i(x)$ is small.

298 The computation of the maximin solution can be done efficiently in $\tilde{O}(|\Sigma|)$ time and the specific
 299 algorithm is shown in Appendix C.

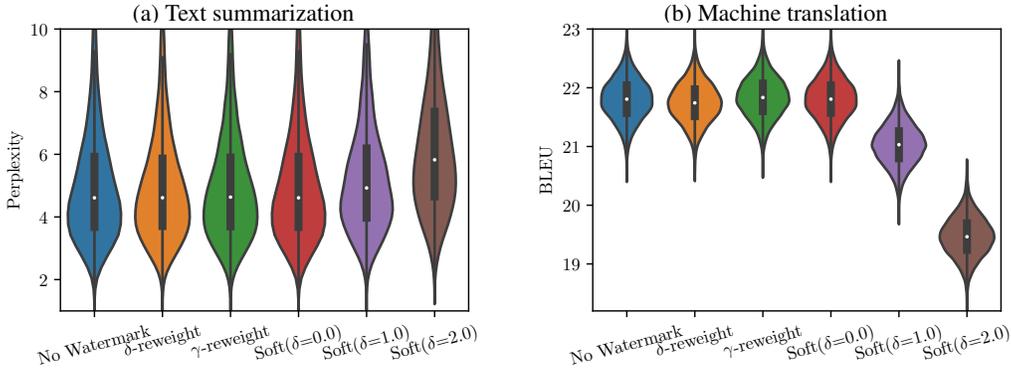


Figure 3: Distribution of perplexity of output for TS and BLEU score for MT.

300 It is important to note that the maximin variant of the LLR score is more robust than the standard
 301 LLR score, as it yields higher scores when the text has undergone some degree of editing. However,
 302 it is not specifically designed to defend against any attacks.

303 A hyperparameter $d \in [0, 1]$ that represent the perturbation strength is introduced in the score.
 304 Intuitively, if the text to be detected has undergone more editing and deviates further from the
 305 distribution $P_{M,w}$, d should be larger. In practice, we recommend using grid search to select the best
 306 value of d . Assuming there are A candidate values for d , corresponding to A different scores $s_i^{(a)}$
 307 ($1 \leq a \leq A$), we can modify Theorem 8 as follows.

308 **Theorem 9.** *Under the same conditions as Theorem 8, but with multiple scores $s_i^{(a)}$, we have*

$$P \left(\max_{1 \leq a \leq A} \left(\sum_{i=1}^n s_i^{(a)} \right) \geq t \right) \leq Ae^{-t}.$$

309 Thus, when using grid search, the final threshold should be adjusted as $\hat{S} = -\log(\alpha) + \log(A)$. This
 310 ensures that the upper bound of the type I error is still α .

311 6 Experiments

312 We evaluate the performance of our Unbiased Watermarks on two important applications of seq2seq
 313 models: text summarization (TS) and machine translation (MT). For the TS task, we use the BART-
 314 large model [37] and the CNN-DM [25] corpus as our training dataset. The MT task involves
 315 translating English to Romanian, for which we employ the Multilingual BART (MBart) [37] model
 316 on the WMT’14 En-Ro corpus. For further details on the experiment setup, please refer to Appendix E.

Table 1: Performance of different watermarking methods on TS and MT. We use F1 scores of BERTScore and scale BERTScore and ROUGE-1 with a factor of 100.

	Text summarization			Machine translation	
	BERTScore \uparrow	ROUGE-1 \uparrow	Perplexity \downarrow	BERTScore \uparrow	BLEU \uparrow
No Watermark	32.70 ± 0.08	38.56 ± 0.09	5.024 ± 0.018	55.9 ± 0.3	21.8 ± 0.3
δ -reweight	32.71 ± 0.08	38.57 ± 0.09	5.022 ± 0.018	56.3 ± 0.3	21.7 ± 0.3
γ -reweight	32.69 ± 0.08	38.60 ± 0.09	5.019 ± 0.018	56.2 ± 0.3	21.8 ± 0.3
Soft($\delta=0.0$)	32.70 ± 0.08	38.56 ± 0.09	5.024 ± 0.018	55.9 ± 0.3	21.8 ± 0.3
Soft($\delta=1.0$)	32.35 ± 0.08	38.20 ± 0.09	5.313 ± 0.018	55.1 ± 0.3	21.0 ± 0.3
Soft($\delta=2.0$)	31.21 ± 0.08	37.17 ± 0.08	6.253 ± 0.022	53.8 ± 0.3	19.5 ± 0.3

317 Our primary focus is to compare the performance of our proposed unbiased watermarking methods
 318 including the δ -reweight and γ -reweight, with the soft-red-list method presented by Kirchenbauer
 319 et al. [32]. The strength of the watermark in the soft-red-list approach is controlled by a parameter δ .

320 The quality of output post-watermarking is presented in Table 1. We observed that the output quality
 321 remains unaffected by our unbiased watermark methods, both for the δ -reweight and γ -reweight,

Table 3: Text sampled from OPT-6.7B, with and without watermarks. For "No watermark" (NW), the score is computed based on δ -reweight. When watermarks are included, the corresponding reweighting function is used for computing score. The optimal perturbation strengths d obtained by grid search are 0.9, 0.0, 0.0 for three outputs respectively.

Prompt	What is a watermark? What's the purpose of it?	score	p-value ²
NW	Why don't you want it on there? I'm confused. I think he meant to say logo, since he wrote "watermark", so the first word must be a typo. Yes! Exactly typo. Sorry, English is not my native language. Thanks for the explanation!	0.30	8.14
δ -RW	It is supposed to be watermarking the pictures that you took with your phone. I think. So, so you can share your pictures and not take credit for them.	75.9	1.2e-32
γ -RW	A watermark is a small image or logo (often in square pixels) that is placed over the larger, original image. It serves primarily to distinguish copyright or ownership of large images (such as banners and logos) and, on rare occasion, to identify small images (such as thumbnail images for blog posts and pictures).	32.9	5.7e-14

irrespective of the task and metric. Conversely, the soft-red-list method, when $\delta = 0$, does not introduce any watermark and hence does not affect output quality. However, for $\delta > 0$, it significantly deteriorate the quality of output.

Figure 3 provides a more intuitive depiction of the score distributions. It is evident that our unbiased watermark methods not only ensure that the mean performance remains unaffected but also that the performance distribution is stable. Conversely, the soft-red-list method shows a notable performance decrease.

In terms of watermark detection, we compute score associated with each token. The mean and variance of score per token for TS and MT are presented in Table 2. As a heuristic, if the sum of the scores for all tokens in a sentence reaches 10, a p-value of less than 0.0005 is ensured. If the sum score hits 20, the p-value must be less than $3e-8$.

Table 2: Mean and variance of score per token for different reweighting methods and different tasks.

	Text summarization	Machine translation
δ -RW	0.8784 ± 1.4354	0.4192 ± 1.1361
γ -RW	0.2207 ± 0.3678	0.1056 ± 0.2916

Additionally, we provide an example of watermarking applied to a completion task in Table 3. It visually demonstrates the score distribution across tokens: positive scores are represented in green, and negative ones in red. The intensity of the color corresponds to the magnitude of the score, with darker shades representing larger absolute values.

7 Related work

The idea of watermarking text has been widely explored by many researchers [11, 31, 44, 45, 4, 28, 49, 43], even before the advent of large language models. Several techniques involve editing existing text to add a watermark, such as changing synonyms [54, 57, 9, 59, 66] or visually indistinguishable words [46], altering sentence structures [56, 55, 38], and employing neural networks [22, 23, 67].

Recent advancements in generative models have opened new possibilities for directly generating watermarked results. Two relevant works in this domain are by Kirchenbauer et al. [32] and Aaronson [1]. Due to space constraints, we moved the in-depth analysis and other related work to Section B.

8 Conclusion

Overall, this paper provides a novel framework of watermarking for language models, demonstrating that it is possible to use watermark to protect intellectual property and monitor potential misuse without compromising the quality of the generated text. This research serves as a valuable foundation for future work in the field of watermarking for large language models.

²This is an upper bound computed based on Theorem 9. The upper bound could be larger than 1, but this does not necessarily imply that the p-value exceeds 1.

354 References

- 355 [1] Scott Aaronson. My ai safety lecture for ut effective altruism. November 2022. URL [https://](https://scottaaronson.blog/?p=6823)
356 scottaaronson.blog/?p=6823.
- 357 [2] Sahar Abdelnabi and Mario Fritz. Adversarial watermarking transformer: Towards tracing text provenance
358 with data hiding. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 121–140. IEEE, 2021.
- 359 [3] Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. Turning your weakness
360 into a strength: Watermarking deep neural networks by backdooring. In *27th USENIX Security Symposium*,
361 pages 1615–1631, 2018.
- 362 [4] Mikhail J Atallah, Victor Raskin, Michael Crogan, Christian Hempelmann, Florian Kerschbaum, Dina
363 Mohamed, and Sanket Naik. Natural language watermarking: Design, analysis, and a proof-of-concept
364 implementation. In *Information Hiding: 4th International Workshop, IH 2001 Pittsburgh, PA, USA, April*
365 *25–27, 2001 Proceedings 4*, pages 185–200. Springer, 2001.
- 366 [5] Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham,
367 Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu
368 Pal, Matt Post, and Marcos Zampieri. Findings of the 2019 conference on machine translation (WMT19).
369 In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*,
370 pages 1–61, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/
371 W19-5301.
- 372 [6] Franziska Boenisch. A systematic review on model watermarking for neural networks. *Frontiers in big*
373 *Data*, 4:729663, 2021.
- 374 [7] Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang,
375 Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post,
376 Raphael Rubino, Lucia Specia, and Marco Turchi. Findings of the 2017 conference on machine translation
377 (WMT17). In *Proceedings of the Second Conference on Machine Translation*, pages 169–214. Copenhagen,
378 Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4717.
- 379 [8] Nicholas Boucher, Iliia Shumailov, Ross Anderson, and Nicolas Papernot. Bad characters: Imperceptible
380 nlp attacks. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1987–2004. IEEE, 2022.
- 381 [9] Yuei-Lin Chiang, Lu-Ping Chang, Wen-Tai Hsieh, and Wen-Chih Chen. Natural language watermarking
382 using semantic substitution for chinese text. In *Digital Watermarking: Second International Workshop,*
383 *IWDW 2003, Seoul, Korea, October 20-22, 2003. Revised Papers 2*, pages 129–140. Springer, 2004.
- 384 [10] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi
385 Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv*
386 *preprint arXiv:2210.11416*, 2022.
- 387 [11] Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. *Digital watermarking and*
388 *steganography*. Morgan kaufmann, 2007.
- 389 [12] Evan Crothers, Nathalie Japkowicz, and Herna Viktor. Machine generated text: A comprehensive survey
390 of threat models and detection methods. *arXiv preprint arXiv:2210.07321*, 2022.
- 391 [13] Falcon Z Dai and Zheng Cai. Towards near-imperceptible steganographic text. *arXiv preprint*
392 *arXiv:1907.06679*, 2019.
- 393 [14] Luc Devroye. *Non-Uniform Random Variate Generation*. Springer New York, 1986.
- 394 [15] Tina Fang, Martin Jaggi, and Katerina Argyraki. Generating steganographic text with lstms. *arXiv preprint*
395 *arXiv:1705.10742*, 2017.
- 396 [16] Evgeniy Gabrilovich and Alex Gontmakher. The homograph attack. *Communications of the ACM*, 45(2):
397 128, 2002.
- 398 [17] Margherita Gambini, Tiziano Fagni, Fabrizio Falchi, and Maurizio Tesconi. On pushing deepfake tweet
399 detection capabilities to the limits. In *14th ACM Web Science Conference 2022*, pages 154–163, 2022.
- 400 [18] Riley Goodside. There are adversarial attacks for that proposal as well — in particular, generating
401 with emojis after words and then removing them before submitting defeats it,.. January 2023. URL
402 <https://twitter.com/goodside/status/1610682909647671306>.
- 403 [19] Google. Palm-2-llm. <https://blog.google/technology/ai/google-palm-2-ai-large-language-model/>, 2023.

- 404 [20] Chenxi Gu, Chengsong Huang, Xiaoqing Zheng, Kai-Wei Chang, and Cho-Jui Hsieh. Watermarking
405 pre-trained language models with backdooring. *arXiv preprint arXiv:2210.07543*, 2022.
- 406 [21] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine
407 learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- 408 [22] Xuanli He, Qionikai Xu, Lingjuan Lyu, Fangzhao Wu, and Chenguang Wang. Protecting intellectual
409 property of language generation apis with lexical watermark. In *Proceedings of the AAAI Conference on*
410 *Artificial Intelligence*, volume 36, pages 10758–10766, 2022.
- 411 [23] Xuanli He, Qionikai Xu, Yi Zeng, Lingjuan Lyu, Fangzhao Wu, Jiwei Li, and Ruoxi Jia. Cater: Intellectual
412 property protection on text generation apis via conditional watermarks. *arXiv preprint arXiv:2209.08773*,
413 2022.
- 414 [24] James N Helfrich and Rick Neff. Dual canonicalization: An answer to the homograph attack. In *2012*
415 *eCrime Researchers Summit*, pages 1–10. IEEE, 2012.
- 416 [25] Karl Moritz Hermann, Tomás Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman,
417 and Phil Blunsom. Teaching machines to read and comprehend. In Corinna Cortes, Neil D. Lawrence,
418 Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information Processing*
419 *Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015,*
420 *Montreal, Quebec, Canada*, pages 1693–1701, 2015.
- 421 [26] Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. Automatic detection of
422 generated text is easiest when humans are fooled. *arXiv preprint arXiv:1911.00650*, 2019.
- 423 [27] Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Dániel Simig, Ping Yu, Kurt
424 Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. Opt-impl: Scaling language model instruction
425 meta learning through the lens of generalization. *arXiv preprint arXiv:2212.12017*, 2022.
- 426 [28] Zunera Jalil and Anwar M Mirza. A review of digital watermarking techniques for text documents. In
427 *2009 International Conference on Information and Multimedia Technology*, pages 230–234. IEEE, 2009.
- 428 [29] Ganesh Jawahar, Muhammad Abdul-Mageed, and Laks VS Lakshmanan. Automatic detection of machine
429 generated text: A critical survey. *arXiv preprint arXiv:2011.01314*, 2020.
- 430 [30] Hengrui Jia, Christopher A Choquette-Choo, Varun Chandrasekaran, and Nicolas Papernot. Entangled
431 watermarks as a defense against model extraction. In *USENIX Security Symposium*, pages 1937–1954,
432 2021.
- 433 [31] Nurul Shamimi Kamaruddin, Amirrudin Kamsin, Lip Yee Por, and Hameedur Rahman. A review of text
434 watermarking: theory, methods, and applications. *IEEE Access*, 6:8011–8028, 2018.
- 435 [32] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A watermark
436 for large language models. *arXiv preprint arXiv:2301.10226*, 2023.
- 437 [33] Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. Paraphrasing evades
438 detectors of ai-generated text, but retrieval is an effective defense. *arXiv preprint arXiv:2303.13408*, 2023.
- 439 [34] Zheng Li, Chengyu Hu, Yang Zhang, and Shanqing Guo. How to prove your model belongs to you: A
440 blind-watermark based framework to protect intellectual property of dnn. In *Proceedings of the 35th*
441 *Annual Computer Security Applications Conference*, pages 126–137, 2019.
- 442 [35] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches*
443 *out*, pages 74–81, 2004.
- 444 [36] Yang Liu and Mirella Lapata. Text summarization with pretrained encoders. In *Proceedings of the*
445 *2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint*
446 *Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, Hong Kong, China,
447 November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1387.
- 448 [37] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and
449 Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. *Transactions of the*
450 *Association for Computational Linguistics*, 8:726–742, 2020.
- 451 [38] Hasan Mesut Meral, Bülent Sankur, A Sumru Özsoy, Tunga Güngör, and Emre Sevinç. Natural language
452 watermarking via morphosyntactic alterations. *Computer Speech & Language*, 23(1):107–125, 2009.
- 453 [39] OpenAI. Chatgpt. <https://openai.com/blog/chatgpt>, 2023a.

- 454 [40] OpenAI. Gpt-4 technical report. *arXiv*, 2023b.
- 455 [41] Luca Pajola and Mauro Conti. Fall of giants: How popular text-based mlaas fall against a simple evasion
456 attack. In *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 198–211. IEEE,
457 2021.
- 458 [42] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation
459 of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational*
460 *Linguistics*, pages 311–318, 2002.
- 461 [43] Fabien AP Petitcolas, Ross J Anderson, and Markus G Kuhn. Information hiding-a survey. *Proceedings of*
462 *the IEEE*, 87(7):1062–1078, 1999.
- 463 [44] Christine I Podilchuk and Edward J Delp. Digital watermarking: algorithms and applications. *IEEE signal*
464 *processing Magazine*, 18(4):33–46, 2001.
- 465 [45] Vidyasagar M Potdar, Song Han, and Elizabeth Chang. A survey of digital image watermarking techniques.
466 In *INDIN'05. 2005 3rd IEEE International Conference on Industrial Informatics, 2005.*, pages 709–716.
467 IEEE, 2005.
- 468 [46] Stefano Giovanni Rizzo, Flavio Bertini, and Danilo Montesi. Fine-grain watermarking for intellectual
469 property protection. *EURASIP Journal on Information Security*, 2019:1–20, 2019.
- 470 [47] Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi. Can
471 ai-generated text be reliably detected? *arXiv preprint arXiv:2303.11156*, 2023.
- 472 [48] M Hassan Shirali-Shahreza and Mohammad Shirali-Shahreza. A new synonym text steganography. In
473 *2008 international conference on intelligent information hiding and multimedia signal processing*, pages
474 1524–1526. IEEE, 2008.
- 475 [49] Katzenbeisser Stefan, A Petitcolas Fabien, et al. Information hiding techniques for steganography and
476 digital watermarking, 2000.
- 477 [50] Yuchen Sun, Tianpeng Liu, Panhe Hu, Qing Liao, Shouling Ji, Nenghai Yu, Deke Guo, and Li Liu. Deep
478 intellectual property: A survey. *arXiv preprint arXiv:2304.14613*, 2023.
- 479 [51] Reuben Tan, Bryan A Plummer, and Kate Saenko. Detecting cross-modal inconsistency to defend against
480 neural fake news. *arXiv preprint arXiv:2009.07698*, 2020.
- 481 [52] Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. The science of detecting llm-generated texts. *arXiv preprint*
482 *arXiv:2303.07205*, 2023.
- 483 [53] Yi Tay, Dara Bahri, Che Zheng, Clifford Brunk, Donald Metzler, and Andrew Tomkins. Reverse engineering
484 configurations of neural text generation models. *arXiv preprint arXiv:2004.06201*, 2020.
- 485 [54] Mercan Topkara, Cuneyt M Taskiran, and Edward J Delp III. Natural language watermarking. In *Security,*
486 *Steganography, and Watermarking of Multimedia Contents VII*, volume 5681, pages 441–452. SPIE, 2005.
- 487 [55] Mercan Topkara, Giuseppe Riccardi, Dilek Hakkani-Tür, and Mikhail J Atallah. Natural language
488 watermarking: Challenges in building a practical system. In *Security, Steganography, and Watermarking of*
489 *Multimedia Contents VIII*, volume 6072, pages 106–117. SPIE, 2006.
- 490 [56] Mercan Topkara, Umut Topkara, and Mikhail J Atallah. Words are not enough: sentence level natural
491 language watermarking. In *Proceedings of the 4th ACM international workshop on Contents protection*
492 *and security*, pages 37–46, 2006.
- 493 [57] Umut Topkara, Mercan Topkara, and Mikhail J Atallah. The hiding virtues of ambiguity: quantifiably
494 resilient watermarking of natural language text through synonym substitutions. In *Proceedings of the 8th*
495 *workshop on Multimedia and security*, pages 164–174, 2006.
- 496 [58] Honai Ueoka, Yugo Murawaki, and Sadao Kurohashi. Frustratingly easy edit-based linguistic steganography
497 with a masked language model. *arXiv preprint arXiv:2104.09833*, 2021.
- 498 [59] Ashish Venugopal, Jakob Uszkoreit, David Talbot, Franz Josef Och, and Juri Ganitkevitch. Watermarking
499 the outputs of structured prediction with an application in statistical machine translation. In *Proceedings of*
500 *the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 1363–1372, 2011.
- 501 [60] Hong Wang, Xuan Luo, Weizhi Wang, and Xifeng Yan. Bot or human? detecting chatgpt imposters with a
502 single question. *arXiv preprint arXiv:2305.06424*, 2023.

- 503 [61] Alex Wilson and Andrew D Ker. Avoiding detection on twitter: embedding strategies for linguistic
504 steganography. Society of Photo-optical Instrumentation Engineers, 2016.
- 505 [62] Alex Wilson, Phil Blunsom, and Andrew D Ker. Linguistic steganography on twitter: hierarchical language
506 modeling with manual interaction. In *Media Watermarking, Security, and Forensics 2014*, volume 9028,
507 pages 9–25. SPIE, 2014.
- 508 [63] Alex Wilson, Phil Blunsom, and Andrew Ker. Detection of steganographic techniques on twitter. In
509 *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages
510 2564–2569, 2015.
- 511 [64] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric
512 Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art
513 natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- 514 [65] Max Wolff and Stuart Wolff. Attacking neural text detectors. *arXiv preprint arXiv:2002.11768*, 2020.
- 515 [66] Xi Yang, Jie Zhang, Kejiang Chen, Weiming Zhang, Zehua Ma, Feng Wang, and Nenghai Yu. Tracing text
516 provenance via context-aware lexical substitution. In *Proceedings of the AAAI Conference on Artificial
517 Intelligence*, volume 36, pages 11613–11621, 2022.
- 518 [67] KiYoon Yoo, Wonhyuk Ahn, Jiho Jang, and Nojun Kwak. Robust natural language watermarking through
519 invariant features. *arXiv preprint arXiv:2305.01904*, 2023.
- 520 [68] Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin
521 Choi. Defending against neural fake news. *Advances in neural information processing systems*, 32, 2019.
- 522 [69] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy.
523 Protecting intellectual property of deep neural networks with watermarking. In *Proceedings of the 2018 on
524 Asia Conference on Computer and Communications Security*, pages 159–172, 2018.
- 525 [70] Xuandong Zhao, Yu-Xiang Wang, and Lei Li. Protecting language generation models via invisible
526 watermarking. *arXiv preprint arXiv:2302.03162*, 2023.
- 527 [71] Zachary M Ziegler, Yuntian Deng, and Alexander M Rush. Neural linguistic steganography. *arXiv preprint
528 arXiv:1909.01496*, 2019.