
A Resilient and Accessible Distribution-Preserving Watermark for Large Language Models

Yihan Wu¹ Zhengmian Hu¹ Junfeng Guo¹ Hongyang Zhang² Heng Huang¹

Abstract

Watermarking techniques offer a promising way to identify machine-generated content via embedding covert information into the contents generated from language models. A challenge in the domain lies in preserving the distribution of original generated content after watermarking. Our research extends and improves upon existing watermarking framework, placing emphasis on the importance of a **Distribution-Preserving (DiP)** watermark. Contrary to the current strategies, our proposed DiPmark simultaneously preserves the original token distribution during watermarking (distribution-preserving), is detectable without access to the language model API and prompts (accessible), and is provably robust to moderate changes of tokens (resilient). DiPmark operates by selecting a random set of tokens prior to the generation of a word, then modifying the token distribution through a distribution-preserving reweight function to enhance the probability of these selected tokens during the sampling process. Extensive empirical evaluation on various language models and tasks demonstrates our approach’s distribution-preserving property, accessibility, and resilience, making it a effective solution for watermarking tasks that demand impeccable quality preservation. Code is available at¹.

1. Introduction

In the current era, artificial intelligence has attained the capability to generate text remarkably indistinguishable from human authorship (Google, 2023; OpenAI, 2023). This ad-

¹Department of Computer Science, University of Maryland College Park ²School of Computer Science, University of Waterloo. Correspondence to: Yihan Wu <ywu42@umd.edu>, Heng Huang <heng@umd.edu>.

Proceedings of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

¹<https://github.com/yihwu/DiPmark.git>

vancement has raised concerns regarding the discernment of authenticity in content, questioning whether it originates from human intellect or AI models. In particular, the proficiency of large language models (LLMs) in imitating human writing style brings a series of implications. While these models facilitate the simplification of complex tasks and enhance human capabilities, they simultaneously harbor risks of misuse, evident in instances of academic dishonesty and the spread of misinformation via online platforms.

The challenge of distinguishing machine-generated content from that authored by humans is escalating, with conventional detection tools often proving inadequate (Krishna et al., 2023). To address this issue, *watermarking* emerges as a nuanced solution (Kirchenbauer et al., 2023). This type of approach involves embedding discreet yet identifiable watermarks in AI-generated text, signifying its artificial origin. Beyond the widely held notion that watermarks should be identifiable via a secret key (Kirchenbauer et al., 2023), there are additional fundamental characteristics necessary for an efficient watermark within language models:

- (Distribution-preserving) The watermark should provably preserving the distribution of the original language model.
- (Accessible) Detecting watermark within the content should be efficient and straightforward without accessing the language models and prompts.
- (Resilient) The watermark should remain identifiable if the content undergoes moderate modifications. Furthermore, we define a watermark as ‘provably resilient’ if it can be provably identified under such modifications.

To the best of our knowledge, there is no watermark technique adhere to the aforementioned three key properties simultaneously (see Table 1 for an overall comparison). Existing methods either impact the model’s sampling distribution (Kirchenbauer et al., 2023; Zhao et al., 2023), lack resilience against text alterations such as editing or cropping (Christ et al., 2023), require thousands of inference step during the detection process (Kuditipudi et al., 2023), or require the prompt and the token logits of language model API during detection (Hu et al., 2023a).

Table 1. Existing watermarking techniques do not adhere to all three key properties (distribution-preserving, accessible, resilient). **Distribution-preserving:** Kirchenbauer et al. (2023) impacts the distribution of the generated tokens. **Accessible:** During detection, Kuditipudi et al. (2023) necessitates thousands of inference steps, and Hu et al. (2023a) requires the token logits of language model API and the prompt, which could result in huge computational costs and hurt the accessibility. **Resilient and Provably Resilient:** DiPmark is provably resilient against arbitrary text modifications with a guaranteed false positive rate, whereas other methods lack corresponding discussions.

Properties	Kirchenbauer et al. (2023)	Kuditipudi et al. (2023)	Hu et al. (2023a)	DiPmark
Distribution-preserving (Sec. 4& 7.1)	✗	✓	✓	✓
Accessible (Sec. 5& 7.2)	✓	✗	✗	✓
Resilient and Provably Resilient (Sec. 6& 7.3)	✗	✗	✗	✓

Our watermarking framework (*i.e.*, DiPmark), in alignment with pre-existing schema (Kirchenbauer et al., 2023), is comprised of two components: (1.) a generating function, which transforms a prompt and a secret watermark key into the content from the language model; and (2.) a detecting function that identifies a potential watermarked text through the secret key. During the text generation process, language model providers will adjust the output probability of the generated tokens using a secret key. We design a novel distribution-preserving generating function, ensuring that each instance of text generation consists with the original language model’s distribution. As for the detection phase, the user can detect the presence of watermark efficiently by solely using the secret key and the watermarked text without accessing prompts and language model API. Through experimental assessments on widely-studied language models, including BART-large model (Liu et al., 2020), LLaMA-2 (Touvron et al., 2023), and GPT-4 (OpenAI, 2023); our approach is demonstrated possessing above mentioned three fundamental properties.

Our contributions. Our work tackles the problem of designing watermarks for large language models without affecting its overall performance and advances the state-of-the-art in multiple ways.

- We propose a novel watermarking framework, DiPmark, that introduces a **provably** distribution-preserving watermarking scheme for language models. Comparing with existing methods, DiPmark is **simultaneously** distribution-preserving, efficient, and provable resilient.
- We identify the existing watermark detector (Kirchenbauer et al., 2023) cannot precisely guarantee the false positive rate of detection. To solve this problem, we develop an well-defined watermark detection statistic for DiPmark, which can reliably detect the watermark within generated contents while maintaining a guaranteed false positive rate. Furthermore, we also show our detect algorithm is **provably** robust against arbitrary text modifications.
- Through extensive experiments on widely-adopted language models, we validate the distribution-preserving

property of DiPmark. Notably, the detection time for 1,000 watermarked sequences produced by LLaMA-2 stands at a mere 90 seconds without the need of API access and prompts (at least 4X faster compared with current distribution-preserving watermark detection (Hu et al., 2023a; Kuditipudi et al., 2023)). Furthermore, DiPmark exhibits robustness even when subjected to 20% to 30% random text modifications and paraphrasing attacks. Finally, in a case study, we show the effectiveness of DiPmark on GPT-4.

2. Related Work

In a recent seminal work, Kirchenbauer et al. (2023) introduced a pioneering watermarking scheme tailored for LLMs. However, this approach inevitably leads to a pivotal change in the distribution of the generated text, potentially compromising the quality of the generated content. To maintain the output distribution in watermarked content, alternative strategies have been explored. Christ et al. (2023) and Kuditipudi et al. (2023) employed the inverse sampling method to generate watermarked token distributions. Notably, Christ et al. (2023)’s method faces resilience issues under modifications or changes and lacks empirical validation for detectability. Meanwhile, Kuditipudi et al. (2023)’s approach requires the secret key distribution during detection, potentially compromising data security and watermark stealthiness. Moreover, their detection process involves thousands of resampling steps from the secret key distribution, which is inefficient for lengthy texts. Hu et al. (2023a) also used inverse sampling and permutation based reweight for watermarking, but the detector requires the token logits of language model API and the prompt for generating the content, undermining its operational efficiency. A detailed discussion of watermarking LLMs is in Appendix B.

Our research aligns closely with Kirchenbauer et al. (2023). In their settings, they employed watermarking for text derived from a language model by separating the token set into ‘red’ and ‘green’ lists. Building on this foundation, we introduce an evolved family of reweight strategies. This approach ensures equivalency in distribution between the watermarked language model and the original language model.

3. Preliminary

Notations. We first introduce a few essential notations. Let us represent the vocabulary (or token) set by V and its size or volume by $N = |V|$. We further introduce the set \mathcal{V} , defined as an aggregation of all string sequences, even accounting for those of zero length. In the context of a language model, it produces a token sequence based on a given prompt. For a single step of this process, the likelihood of generating the next token $x_{n+1} \in V$ conditioned on the current context x_1, \dots, x_n is represented as $P_M(x_{n+1} | x_1, x_2, \dots, x_n)$. For the sake of brevity and clarity, we opt for the condensed notation: $P_M(\mathbf{x}_{n+1:n+m} | \mathbf{x}_{1:n})$, where $\mathbf{x}_{n+1:n+m} = (x_{n+1}, \dots, x_{n+m})$. Note that the prompt is deliberately omitted in this representation.

In the context of watermarking, the server provider will use a set of *i.i.d.* watermark cipher $\{\theta_i \in \Theta, i \in \mathbb{N}\}$ on the cipher space Θ to generate the text. The cipher θ_i is usually generated by a secret key $k \in \mathcal{K}$ and a *fragment* of the previous context, named *texture key*, s_i . Instances of texture keys include x_{t-1} , $\mathbf{x}_{t-3:t-1}$, $\mathbf{x}_{1:t-1}$, etc. Each θ_i is independent and following the same distribution P_Θ . We now provide the formal definition of the reweight strategy.

Definition 3.1 (Reweight strategy). *Denote by \mathcal{P} the set of all distributions on the token set V . A reweight strategy is a mapping $P_W : \mathcal{P} \times \Theta \rightarrow \mathcal{P}$. Given the original distribution $P_M(x_{n+1} | \mathbf{x}_{1:n}) \in \mathcal{P}$, the watermarked distribution with cipher θ_i is given by $P_W(P_M(x_{n+1} | \mathbf{x}_{1:n}), \theta_i)$. For brevity, we represent it as $P_W(x_{n+1} | \mathbf{x}_{1:n}, \theta_i)$.*

The reweight strategy stands as the foundation of the watermark algorithm by shaping the distribution of watermarked text. As introduced in (Kirchenbauer et al., 2023), the authors propose a red-green list reweight technique, where the vocabulary set is separated into the red and green lists and the probability of green tokens is promoted during the sampling process. Specifically, given an initial token probability $p(t)$, the watermarked probability for the token, denoted by $p_W(t)$, is formulated as:

$$p_W(t) = \begin{cases} \frac{p(t)}{\sum_{t \in \text{red}} p(t) + \sum_{t \in \text{green}} e^\delta p(t)}, & t \in \text{red list}; \\ \frac{e^\delta p(t)}{\sum_{t \in \text{red}} p(t) + \sum_{t \in \text{green}} e^\delta p(t)}, & t \in \text{green list}, \end{cases}$$

where $\delta > 0$ is a predetermined constant. This strategy reveals an inherent bias in the watermarked distribution. For example, consider $\gamma = 0.5$, suggesting that half of V comprises the red list. With $V = \{a, b\}$, and given probabilities $p(a) = 0.99$ and $p(b) = 0.01$, there are two equivalent permutations of V with congruent appearance likelihoods. An analysis for any value of $\delta > 0$ yields $p_W(a) = 0.5 \left(\frac{e^\delta p(a)}{e^\delta p(a) + p(b)} + \frac{p(a)}{e^\delta p(b) + p(a)} \right) < p(a)$. This indicates that the red-green list watermark does not preserve the

original text’s probability. Below we introduce the formal definition of distribution-preserving reweight strategy and distribution-preserving watermark.

Definition 3.2 (Distribution-preserving reweight strategy). *A reweight strategy, denoted P_W , is said to be distribution-preserving at an individual generation step if, for all $\mathbf{x}_{1:n} \in \mathcal{V}$ and any $i \leq n$, it holds that $P_M(x_i | \mathbf{x}_{1:i-1}) = \mathbb{E}_{\theta_i \sim P_\Theta} [P_W(x_i | \mathbf{x}_{1:i-1}, \theta_i)]$.*

Definition 3.3 (Distribution-preserving watermark). *If a watermark framework preserves the text distribution throughout all generation steps, i.e., $\forall n > 0$, for all sequences $\mathbf{x}_{1:n} \in \mathcal{V}$ we have $P_M(\mathbf{x}_{1:n}) = \mathbb{E}_{\theta_1, \dots, \theta_n} [P_W(\mathbf{x}_{1:n} | \theta_1, \dots, \theta_n)]$, then the watermark is distribution-preserving.*

A distribution-preserving reweight strategy can naturally lead to a distribution-preserving watermark, as illustrated by:

$$\begin{aligned} \mathbb{E}_{\theta_{1:n}} [P_W(\mathbf{x}_{1:n} | \theta_{1:n})] &= \mathbb{E}_{\theta_{1:n}} \left[\prod_{i=1}^n P_W(x_i | \mathbf{x}_{1:i-1}, \theta_i) \right] \\ &= \prod_{i=1}^n \mathbb{E}_{\theta_i} [P_W(x_i | \mathbf{x}_{1:i-1}, \theta_i)] = P_M(\mathbf{x}_{1:n}). \end{aligned}$$

The above equality stems from the independence property of the set $\{\theta_i\}$. Therefore, to establish a distribution-preserving watermark, it is essential to incorporate both: a) a distribution-preserving reweight strategy and b) an *i.i.d.* set of ciphers, $\{\theta_i\}$.

We emphasize the significance of preserving the distribution of text during watermarking, motivated by the following justifications: a) *Stealthy Watermarking*: A watermark that disrupts the original distribution of a language model lacks the attribute of stealthiness. Such alterations make it relatively straightforward to distinguish between watermarked and unwatermarked LMs through multiple instances of sampling. b) *Industry-Level LLM Application*: When contemplating the application of a watermark to industry-standard LLMs like ChatGPT and Bard, the primary consideration is to ensure that the watermark does not compromise the performance of these foundational LLMs. Any watermark that interferes with the original text distribution will inevitably impact the quality of generated text, an outcome that is unacceptable by industry stakeholders.

In the next section, we introduce a reweight strategy with a distribution-preserving characteristic. This attribute guarantees that the text distribution remains unaltered even as we enhance the utilization of tokens from the green list during the watermarking process.

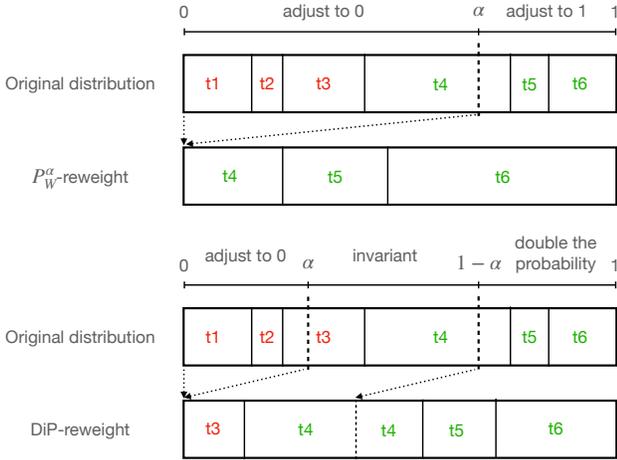


Figure 1. Illustration of the P_W^α -reweight and DiP-reweight. **Top.** In P_W^α -reweight, the token probabilities within the interval $[0, \alpha]$ are adjusted to 0, while the rest are adjust to 1. **Bottom.** In DiP-reweight, the probability mass within $[0, \alpha]$ is transferred to the probability mass within $[1 - \alpha, 1]$.

4. DiPmark

Motivation. The reweight strategy presented in Kirchenbauer et al. (2023) disrupts the inherent text distribution when promoting the use of the green tokens during the sampling process. Such disruption would lead to biased sampling, seriously affecting the quality of the generated text. To address this issue, we design a novel reweight strategy that ensures the token distribution remains unaltered during the watermarking process. Contrary to the approach in (Kirchenbauer et al., 2023) that promotes the use of all tokens from the green list, we emphasize increasing the *sum of the probability* of the green-list tokens. In this way, the watermarked text, when exposed to the secret key, will still exhibit a bias towards the green-list tokens. Motivated by that, we design a reweight function, which preserves the text distribution during watermarking process.

Cipher space for watermarking. Our considered watermark cipher space encompasses the permutations of the vocabulary set, denoted as $\Theta = \{V_1^p, \dots, V_{N!}^p\}$, wherein V_i^p represents a permutation of V . As for the cipher distribution P_Θ , we employ a uniform distribution over Θ , ensuring that each permutation is equally probable for selection.

Reweight strategy. Let $\theta \in \Theta$ be a cipher, constituting a permutation of V . The probabilities of individual tokens can be arranged within the interval $[0, 1]$ according to their respective positions in θ . Given a fixed constant α in $[0, 1]$, the token probabilities within the interval $[0, \alpha]$ are adjusted to 0, while those in the interval $[\alpha, 1]$ are scaled by a factor of $\frac{1}{1-\alpha}$. Let $\gamma \in [0, 1]$ be the red-green list separator for the permuted token list, which is in accordance with the definition in Kirchenbauer et al. (2023). Through this

reweight strategy, we can increase the sum of the probability of green-list tokens for *arbitrary* permutation separator γ , as the green-list tokens consistently appear towards the end of the ordered set θ . Below we present the formal definition of our reweight strategy.

Definition 4.1 (P_W^α -reweight strategy). Let $\theta = \{t_1, \dots, t_N\}$, which represents a permutation of V , and denote $P_M(\cdot|\mathbf{x})$ as the original token distribution. Let $F^\alpha(i|\theta) := \frac{1}{1-\alpha} \max\{\sum_{j=1}^i P_M(t_j|\mathbf{x}) - \alpha, 0\}$. The P_W^α -reweight probability distribution is $P_W^\alpha(t_i|\mathbf{x}, \theta) = F^\alpha(i|\theta) - F^\alpha(i-1|\theta)$.

It is easy to show that $P_W^\alpha(t_i|\mathbf{x}, \theta)$ is a distribution on V for arbitrary α . Firstly, as $F^\alpha(i|\theta)$ is monotonously increasing with i , we have $P_W^\alpha(t_i|\mathbf{x}, \theta) = F^\alpha(i|\theta) - F^\alpha(i-1|\theta) \geq 0$. Secondly, the sum of the probability of all tokens is $\sum_{i=1}^N P_W^\alpha(t_i|\mathbf{x}, \theta) = \sum_{i=1}^N (F^\alpha(i|\theta) - F^\alpha(i-1|\theta)) = F^\alpha(N|\theta) = 1$.

We wish to highlight the distinction between the probability quantile α and the red-green list separator γ . γ serves as the partition for the permuted token list. In contrast, α separates the *probability interval* $[0, 1]$ of the permuted token list. Thus, both the P_W^α -reweight and DiP-reweight (as subsequently defined) remain oblivious to γ , while still effectively promoting the probability of green list tokens.

Leveraging the symmetry of permutations, we can prove that a weighted combination of P_W^α -reweight and $P_W^{1-\alpha}$ -reweight yields a distribution-preserving reweight strategy. It is pivotal to recognize that both P_W^α -reweight and $P_W^{1-\alpha}$ -reweight increase the sum of the probability of green-list tokens. Therefore, the combined effect of these reweight functions still exhibits a preference for the green list tokens. The formal definition of our distribution-preserving reweight strategy is presented subsequently.

Definition 4.2 (DiP-reweight strategy). Denote by $\theta = \{t_1, \dots, t_N\}$ the cipher, which is a permutation of V . Given the original token distribution $P_M(t|\mathbf{x}), \forall t \in V$, where $\mathbf{x} \in \Sigma$ is the previous token sequence, the DiP-reweight strategy is represented by

$$P_W(t_i|\mathbf{x}, \theta) := (1 - \alpha)P_W^\alpha(t_i|\mathbf{x}, \theta) + \alpha P_W^{1-\alpha}(t_i|\mathbf{x}, \theta).$$

As both P_W^α and $P_W^{1-\alpha}$ are distributions on V and $P_W(t_i|\mathbf{x}, \theta)$ is a convex combination of them, $P_W(t_i|\mathbf{x}, \theta)$ is also a distribution on V .

Theorem 4.3. DiP-reweight is a distribution-preserving reweight strategy, i.e., for all $\mathbf{x}_{1:n} \in \mathcal{V}$ and any $i \leq n$, it holds that $P_M(x_i|\mathbf{x}_{1:i-1}) = \mathbb{E}_{\theta_i \sim P_\Theta} [P_W(x_i|\mathbf{x}_{1:i-1}, \theta_i)]$.

We defer the proof of Theorem 4.3 to Appendix C. With the DiP-reweight approach, the generation of *i.i.d.* ciphers, denoted as θ_i , becomes essential for crafting a distribution-preserving watermark. Let k represent a stochastic secret

Algorithm 1 DiPmark generator

- 1: **Input:** watermark key k , reweight parameter α , prompt $\mathbf{x}_{-m:0}$, generate length $n \in \mathbb{N}$, context window length a , and permutation generation function h .
- 2: Initialize texture key history $hist$.
- 3: **for** $i = 1, \dots, n$ **do**
- 4: Calculate the LM distribution for generating the i -th token $P_M(\cdot | \mathbf{x}_{-m:i-1})$.
- 5: Generate a texture key \mathbf{s}_i from $\mathbf{x}_{i-a:i-1}$.
- 6: **if** $\mathbf{s}_i \in hist$ **then**
- 7: Sample the next token x_i using distribution $P_M(\cdot | \mathbf{x}_{-m:i-1})$.
- 8: **else**
- 9: Update key history $hist.append(\mathbf{s}_i)$
- 10: Generate the cipher $\theta_i = h(k, \mathbf{s}_i)$.
- 11: Sample the next token x_i using distribution $P_W(\cdot | \mathbf{x}_{-m:i-1}, h(k, \mathbf{s}_i))$.
- 12: **return** $\mathbf{x}_{1:n}$.

key derived from the key space K following the distribution P_K , let $\mathbf{s} \in \mathcal{V}$ be a texture key, which is a sub-sequence of the previously generated context. Denoted by $\mathbf{x}_{1:t-1}$ the context generated prior to time step t , instances of texture keys encompass \mathbf{x}_{t-1} , $\mathbf{x}_{t-3:t-1}$, and $\mathbf{x}_{1:t-1}$. We introduce a hash function, $h(k, \mathbf{s}) : K \times \mathcal{V} \rightarrow \Theta$, orchestrating the mapping of a secret key in conjunction with a texture key. $\mathbf{s} \in \mathcal{V}$ to a permutation of the token set V . In order to achieve distribution-preserving watermarking, the chosen hash function h should adhere to the following conditions: a) For distinct (secret key, texture key) pairs, i.e., $(k_1, \mathbf{s}_1) \neq (k_2, \mathbf{s}_2)$, $h(k_1, \mathbf{s}_1)$ ought to be statistically independent from $h(k_2, \mathbf{s}_2)$, and b) Upon holding \mathbf{s} constant, every $V_i^p \in \Sigma$ should exhibit a uniform likelihood of being selected given a random key, specifically, $\forall V_i^p \in \Sigma, \mathbb{E}_{k \sim P_K} [\mathbf{1}_{h(k, \mathbf{s})=V_i^p}] = 1/N!$.

There exists hash functions meeting the above criteria, one example being the hash function introduced in Kirchenbauer et al. (2023). Under such conditions, the cipher θ_i can be deemed *i.i.d.* if the texture key \mathbf{s}_i is distinctive for each instance. To ensure this uniqueness, a historical log is employed to retain texture keys generated in prior steps. If a texture key is identified in the historical log, another secret key will be utilized with the texture key to generate the cipher. The detailed methodology is shown in Alg. 1.

Corollary 4.4. *DiPmark (Alg. 1) is a distribution-preserving watermark, i.e., for all sequences $\mathbf{x}_{1:n} \in \mathcal{V}$ and any positive integer n , we have $P_M(\mathbf{x}_{1:n}) = \mathbb{E}_{\theta_1, \dots, \theta_n} [P_W(\mathbf{x}_{1:n} | \theta_1, \dots, \theta_n)]$.*

This can be easily validated by combining the distribution-preserving property of DiP-reweight and the independence of ciphers θ_i .

Algorithm 2 DiPmark detector

- 1: **Input:** text $\mathbf{x}_{1:n}$, watermark key k , volume of the token set N , permutation generation function h , green list separator γ , context window length a , and threshold z .
- 2: Initialize the green token indexer of γ : $L_G(\gamma) = 0$.
- 3: **for** $i = 2, \dots, n$ **do**
- 4: Generate a texture key \mathbf{s}_i based on $\mathbf{x}_{i-a:i-1}$.
- 5: Generate the permutation of token set $\theta_i = h(k, \mathbf{s}_i)$.
- 6: Calculate the list of green tokens via $G = \theta_i[\lceil \gamma N \rceil : N]$.
- 7: **if** $x_i \in G$ **then**
- 8: $L_G(\gamma) = L_G(\gamma) + 1$.
- 9: Calculate the score: $\Phi(\gamma, \mathbf{x}_{1:n}) = \frac{L_G(\gamma)}{n} - (1 - \gamma)$.
- 10: **return** $\Phi(\gamma, \mathbf{x}_{1:n}) > z$.

5. DiPmark Detection

We leverage a hypothesis test to identify the presence of DiPmark. In the context of a predetermined red-green list separator $\gamma \in [0, 1]$, we classify the initial $\lceil \gamma N \rceil$ tokens within the token set permutation as belonging to the red list, while the remaining tokens are categorized as part of the green list. Given a text sequence $\mathbf{x}_{1:n}$, we establish the null hypothesis H_0 : $\mathbf{x}_{1:n}$ is generated without any awareness of DiPmark. Below we design a statistic, named ‘‘green token ratio’’, for conducting the hypothesis test.

Definition 5.1 (Green token ratio). *Let $L_G(\gamma)$ be the count of green tokens within $\mathbf{x}_{1:n}$, where γ is the predetermined red-green list separator. The green token ratio is give by $\Phi(\gamma, \mathbf{x}_{1:n}) := L_G(\gamma)/n - (1 - \gamma)$.*

The green token ratio quantifies the bias towards green tokens within the text sequence. The term $L_G(\gamma)/n$ signifies the proportion of green tokens within a sequence of tokens, while $1 - \gamma$ denotes the expected green token proportion in an unwatermarked sequence. Under the null hypothesis H_0 , $L_G(\gamma)$ follows a binomial distribution with parameters $p = (1 - \gamma)$ and n total trials, i.e., $L_G(\gamma) \sim \text{Binomial}(n, 1 - \gamma)$. The reason for this is that each token is randomly assigned to either the red or green list in the absence of our watermarking rule. We derive the subsequent concentration bound of the green token ratio $\Phi(\gamma, \mathbf{x}_{1:n})$:

Theorem 5.2 (Concentration bound of $\Phi(\gamma, \mathbf{x}_{1:n})$). *Let $\Phi(\gamma, \mathbf{x}_{1:n}) := L_G(\gamma)/n - (1 - \gamma)$, where $L_G(\gamma) \sim \text{Binomial}(n, 1 - \gamma)$. We have $\forall t \in \mathbb{R}$,*

$$\Pr(\Phi(\gamma, \mathbf{x}_{1:n}) \geq t) \leq \exp(-n\mathbb{KL}(t + 1 - \gamma || 1 - \gamma)),$$

where $\mathbb{KL}(p||q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$ is the Kullback-Leibler divergence.

We proceed to reject the null hypothesis and detect the watermark if $\Phi(\gamma, \mathbf{x}_{1:n})$ surpasses a predefined threshold. For in-

Table 2. Comparison of different test statistics on theoretical FPR (false positive rate) and empirical FPR with 500 non-watermarked sentences. We can see clearly the empirical FPR of z-test is continuously greater than its theoretical guarantee.

False positive samples/All samples	$p < 0.10$ (10%FPR)	$p < 0.01$ (1%FPR)
z-test (Kirchenbauer et al., 2023)	56/500 (11.2% FPR)	12/500 (2.4% FPR)
DiPmark statistic	13/500 (2.6% FPR)	4/500 (0.5% FPR)

stance, setting the threshold as $\Phi(\gamma, \mathbf{x}_{1:n}) \geq 1.517/\sqrt{n}$ results in rejecting H_0 (indicating watermark presence) while maintaining a false positive rate below 1%. Our detection algorithm is shown in Alg. 2. Noting that the concentration bound of $\Phi(\gamma, \mathbf{x}_{1:n})$ scales proportionally with n times the green token ratio. With a fixed green token ratio $\Phi(\gamma, \mathbf{x}_{1:n})$, detecting longer sequences becomes more straightforward because they will show a lower false positive rate. The validity of this analysis is also confirmed in Section F.2.

Difference between our detection algorithm and Kirchenbauer et al. (2023). It is noteworthy that we diverge from Kirchenbauer et al. (2023) by avoiding the use of the z-test statistic $(L_G(\gamma) - (1 - \gamma)n)/\sqrt{n\gamma(1 - \gamma)}$. The z-test assumes a normal distribution for the test statistic. This approximation is imprecise, which could lead to an inaccurate estimation of the p-value, consequently resulting in the wrongful classification of sentences not generated by LMs as being LM-produced. For example, given $n = 100, \gamma = 0.5, L_G(\gamma) = 57$, the p-value of the z-test statistic is about 0.08, indicating that this sentence would be identified as watermarked at 10% FPR (false positive rate). However, in our case, the p-value is around 0.37, suggesting that we cannot determine this sentence as watermarked. In Table 2, we compare the empirical FPR of the two test statistics with their theoretical guaranteed FPR on 500 non-watermarked sentences. We can see clearly the empirical FPR is larger than its theoretical guarantee, which validates our assertion that z-test is imprecise on watermark detection. A detailed discussion can be found in Section D.

Detecting efficiency discussion. Similar to the detection algorithms presented in (Kirchenbauer et al., 2023), our watermark detection process is highly efficient, requiring only a single pass through the provided text sequence. However, it is worth noting that the detection algorithm outlined in Kuditipudi et al. (2023) necessitates iterating through the sequence a staggering 5000 times, which is notably inefficient when compared to our approach. Besides, Hu et al. (2023a) requires prompt and language model API during detection, which is also not practical or efficient. A detailed empirical comparison is in Section 7.2.

6. DiPmark is Provably Resilient Against Text Modification

In this section, we show that DiPmark possesses provable robustness against arbitrary textual modification attacks with

a guaranteed fixed false positive rate. Notably, the existing watermarking approaches are not provably resilient with a guaranteed FPR. Kirchenbauer et al. (2023) and Zhao et al. (2023) assume that the test statistic follows a normal distribution, leading to imprecise guarantee of FPR according to our discussion in Section 5.

Problem formulation. Let $\mathbf{x}_{1:n}$ represent a watermarked sentence. To generate the cipher θ at the i -th iteration, we employ a hash function h , a confidential key k , and a texture key $\mathbf{s} := \mathbf{x}_{i-a:i-1}, a \geq 1$. This indicates that the preceding a tokens serve as the texture key for the watermarking of the token situated at position i . During the detection phase, the formula $\Phi(\gamma, \mathbf{x}_{1:n}) := L_G(\gamma)/n - (1 - \gamma)$ coupled with a threshold z is applied to ascertain if the text has been watermarked. Notably, within $\Phi(\gamma, \mathbf{x}_{1:n})$, the sole variable associated with textual modification assaults is $L_G(\gamma)$. Consequently, our primary objective is to discern the most severe reduction in $L_G(\gamma)$ for a single token alteration.

Worst-case perturbation analysis. Supposing the token x_i in $\mathbf{x}_{1:n}$ undergoes modification, this will lead to a reduction in $L_G(\gamma)$ through two ways: a) Initially, the token x_i may be categorized as a green token, but post-alteration, it either gets eliminated or transitions into a red token, leading to a potential decline in the number of green tokens $L_G(\gamma)$ by at most 1. b) Since the list of red-green tokens for x_{i+1}, \dots, x_{i+a} is generated by hashing the token x_i , its subsequent alteration could cause x_{i+1}, \dots, x_{i+a} to turn into red tokens. In this scenario, the number of green tokens $L_G(\gamma)$ may shrink by a maximum of a . As a result, the greatest decline in $L_G(\gamma)$ for a single token modification stands at $a + 1$.

Definition 6.1 (Certified radius). *Let $\epsilon \in [0, 1]$ denote the fraction of altered tokens. The certified radius of a watermarked sequence is ϵ_0 , if for all perturbations confined within the budget $\epsilon \leq \epsilon_0$, the altered watermarked sequence can still be recognized as watermarked.*

Theorem 6.2. *Given $\Phi(\gamma, \mathbf{x}_{1:n}) := L_G(\gamma)/n - (1 - \gamma)$ and a threshold z , the certified radius of the watermarked sequence $\mathbf{x}_{1:n}$ is $\epsilon_0 = \frac{\Phi(\gamma; \mathbf{x}_{1:n}) - z}{2 + a - \gamma + z}$.*

7. Experiments

Our experimental section consists of five parts. In the first three parts, we compare the distribution-preserving property, accessibility, and resilience of DiPmark with the SOTA watermark methods (Kirchenbauer et al., 2023; Kuditipudi et al., 2023; Hu et al., 2023a). In the fourth part, we compare the detectability of DiPmark with the Soft watermark introduced in (Kirchenbauer et al., 2023). In the final part, we validate the practicality of DiPmark by conducting a case study on GPT-4 (OpenAI, 2023). Detailed experimental settings are in Appendix E.

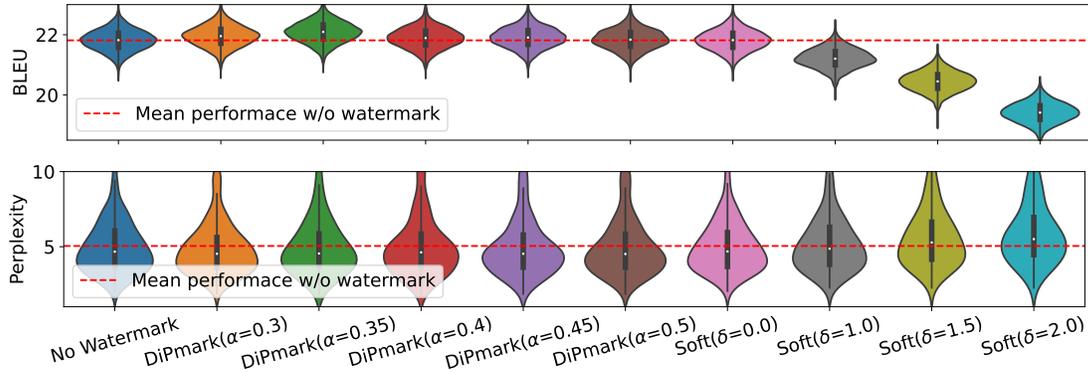


Figure 2. Empirical verification of distribution-preserving property of DiPmark. **Top:** Violin plot of Machine Translation BLEU. **Bottom:** Violin plot of Text Summarization Perplexity. We can see the Soft watermarks (Kirchenbauer et al., 2023) significantly degrade the text quality, while DiPmarks preserve the text quality.

Table 3. Distribution-preserving performance of different watermarking methods on machine translation and text summarization. We use F1 scores of BERTScore and scale BERTScore with a factor of 100.

	Machine Translation		Text Summarization	
	BERTScore \uparrow	BLEU \uparrow	BERTScore \uparrow	Perplexity \downarrow
No Watermark	55.9 \pm 0.3	21.8 \pm 0.3	32.73 \pm 0.08	5.021 \pm 0.018
Soft ($\delta=0.0$)	56.0 \pm 0.3	21.8 \pm 0.3	32.73 \pm 0.08	5.021 \pm 0.018
Soft ($\delta=1.0$)	55.7 \pm 0.3	21.2 \pm 0.3	32.37 \pm 0.08	5.309 \pm 0.019
Soft ($\delta=1.5$)	55.0 \pm 0.3	20.4 \pm 0.3	32.09 \pm 0.08	5.660 \pm 0.021
Soft ($\delta=2.0$)	53.9 \pm 0.3	19.4 \pm 0.3	31.46 \pm 0.08	6.241 \pm 0.023
Kuditipudi et al. (2023)	56.0 \pm 0.3	21.7 \pm 0.3	32.70 \pm 0.08	5.021 \pm 0.021
Hu et al. (2023a)	56.3 \pm 0.3	21.8 \pm 0.3	32.71 \pm 0.08	5.023 \pm 0.018
DiPmark ($\alpha=0.3$)	56.1 \pm 0.3	22.0 \pm 0.3	32.79 \pm 0.08	5.014 \pm 0.018
DiPmark ($\alpha=0.35$)	56.2 \pm 0.3	22.1 \pm 0.3	32.74 \pm 0.08	4.998 \pm 0.018
DiPmark ($\alpha=0.4$)	56.1 \pm 0.3	21.9 \pm 0.3	32.77 \pm 0.08	5.001 \pm 0.018
DiPmark ($\alpha=0.45$)	56.2 \pm 0.3	21.9 \pm 0.3	32.69 \pm 0.08	5.024 \pm 0.018
DiPmark ($\alpha=0.5$)	56.2 \pm 0.3	21.8 \pm 0.3	32.72 \pm 0.08	5.014 \pm 0.018

General experimental observation. We find that our DiPmark, configured with $\alpha = 0.45$, exhibits comparable levels of detectability and robustness comparing with the Soft watermark ($\delta = 1.5$) (Kirchenbauer et al., 2023). Importantly, our DiPmark maintains the same level of text quality as the original language model, owing to its inherent distribution-preserving property.

7.1. Distribution-preserving Property

We will empirically verify the distribution-preserving property of different watermarks. Since DiPmark is **provably** distribution-preserving (Corollary 4.4), we use this experiment as a support for the theorem.

We follow the evaluation process of (Hu et al., 2023a), where we assess the performance of DiPmark with two seq2seq tasks: text summarization (TS) and machine translation (MT). For the TS task, we employ the BART-large model (Liu et al., 2020). For MT task, we focus on English-to-Romanian translation. We employ the Multilingual BART (MBart) model (Liu et al., 2020) on the WMT’14 En-Ro corpus. Specifically for DiPmark, we select values for α from the set $\{0.3, 0.35, 0.4, 0.45, 0.5\}$, while for the Soft

watermark (Kirchenbauer et al., 2023), we choose green list bias values δ from the set $\{0.0, 1.0, 1.5, 2.0\}$ alongside a fixed green list separator $\gamma = 0.5$, indicating that 50% of tokens are green while the remainder are red. Notice, Soft watermark with $\delta = 0.0$ is equivalent to no watermark since it does not promote the probability of green list tokens.

Upon examining Figure 2 and Table 3, we find across all α values in the range $\{0.3, 0.35, 0.4, 0.45, 0.5\}$, the BLEU scores in the machine translation tasks and the perplexity values in the text summarization tasks remain consistently similar between DiPmark and the original language model. However, as we increase the δ values in the Soft watermark, a notable degradation in text quality becomes evident. A more comprehensive set of results is provided in Appendix F.1.

7.2. Accessibility

We compare the time for detecting 1 and 1,000 watermarked sequences with different detection algorithm. The task is text generation with LLaMA-2 (chat, 7B). We use the same GPU (NVIDIA A6000) for all experiments. From Table 4 we see the detecting algorithms of DiPmark are efficient without accessing LMs, while Hu et al. (2023a) requires additional access to LMs and prompts, and Kuditipudi et al. (2023) needs significantly longer time.

Table 4. Comparison of accessibility of different watermarks.

Number of samples	1	1,000	LM & prompt access
Soft watermark	0.3s	92s	No
Kuditipudi et al. (2023)	80s	12h	No
Hu et al. (2023a)	3.4s	412s	Yes
DiPmark	0.3s	90s	No

7.3. Resilience and provable resilience

We compare the resilience of the DiPmark ($\alpha = 0.45$) with the SOTA watermark approaches (Kirchenbauer et al., 2023; Kuditipudi et al., 2023; Hu et al., 2023a). In this context, we use the text generation task with 1,000 generated sequences

Table 5. AUC score of different watermarks under varying attack strength ϵ on text generation task. Each row is evaluated over around 500 watermarked and 500 non-watermarked sequences of length $n = 260 \pm 5$.

AUC	Random text modification			
	$\epsilon = 0.0$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$
Soft watermark	0.9990	0.9883	0.9521	0.8033
Kuditipudi et al. (2023)	0.9951	0.9461	0.8979	0.7815
Hu et al. (2023a)	0.9936	0.9297	0.8391	0.7574
DiPmark ($\alpha=0.45$)	0.9990	0.9859	0.9515	0.8060

AUC	Paraphrasing attack			
	$\epsilon = 0.0$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$
Soft watermark	0.9990	0.9894	0.9469	0.8157
Kuditipudi et al. (2023)	0.9951	0.9529	0.9013	0.7711
Hu et al. (2023a)	0.9936	0.9368	0.8325	0.7661
DiPmark ($\alpha=0.45$)	0.9990	0.9871	0.9503	0.8216

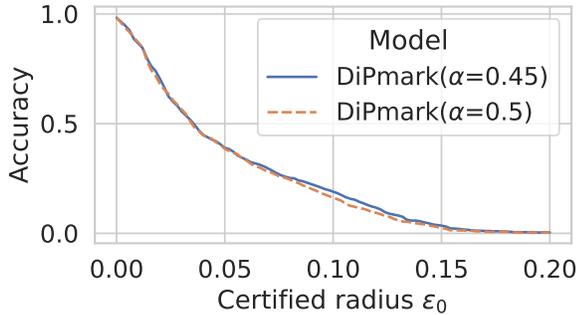


Figure 3. Certified radius ϵ_0 of DiPmark with text modification with FPR smaller than 1%. The x-axis refers the certified radius and the y-axis refers the percentage of watermarked sequences that are resilience under any text modification attacks with budget ϵ_0 .

on LLaMA-2. The texture key generation relies on the most recent one token, i.e., $a = 1$. For resilience evaluation, we manipulate $\epsilon \in \{0.1, 0.2, 0.3\}$ portion of the text tokens through random text modifications and paraphrasing attacks. We also evaluate the provable resilience of the DiPmark under 1% FPR, where we use the above mentioned 1,000 generated sequences on LLaMA-2 to calculate the certified radius (Theorem 6.2).

In Table 5, we report the AUC score of different watermarks under varying attack strength ϵ . The analysis underscores that, when ϵ remains below 0.3, DiPmark demonstrates robust performance in effectively detecting watermarked sentences. In Figure 3, we also show the certified radius of the watermarked sequences of DiPmark with FPR smaller than 1% under the text modification.

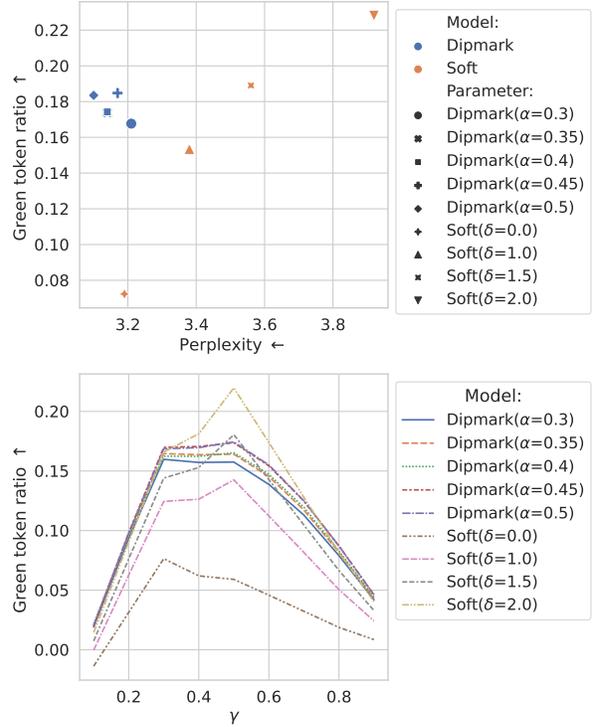


Figure 4. **Top:** Average perplexity vs green token ratio with $\gamma = 0.5$ on text generation tasks. **Bottom:** Average green token ratio with different γ .

7.4. Ablation study: watermark detectability

We evaluate the detectability of our watermark on text generation task using LLaMA-2. We generate 1,000 examples for each tasks. We select $\alpha \in \{0.45, 0.5\}$ for DiPmark, and $\delta \in \{1.0, 1.5, 2.0\}$ and $\gamma = 0.5$ for Soft watermark (Kirchenbauer et al., 2023). During detection, we use $\gamma = 0.5$. We report the Type I (FPR) and II (FNR) errors. We set the threshold $z = 1.073/\sqrt{n}$ (FPR $p \leq 0.1$) and $z = 1.517/\sqrt{n}$ (FPR $p \leq 0.01$). We also report the averaged green token ratio (5.1) vs. text perplexity and token list separator γ of DiPmark and Soft watermark. The averaged green token ratio quantifies the bias towards green tokens within the text sequence (see Section 5). Notice, as the z-test in Kirchenbauer et al. (2023) is imprecise (see Section 5), we use DiPmark detector for all models.

The results for text generation are visually depicted in Figure 4. In Figure 4 (top), it is evident that our DiPmark variants with $\alpha = 0.45$ and 0.5 yield green token ratios akin to those of the Soft watermark with $\delta = 1.5$ without any discernible degradation in text quality. Figure 4 (bottom) delves into the impact of different green list separators γ , revealing that, for most watermark models, $\gamma = 0.5$ yields the highest green token ratio, underscoring its suitability as a reasonable choice for watermark detection. The empiri-

Table 6. Empirical error rates for watermark detection on text generation. Each row is averaged over around 500 watermarked and 500 non-watermarked sequences of length $n = 260 \pm 5$. We select the threshold $z = 1.073/\sqrt{n}$ (false positive rate $p \leq 0.1$) and $z = 1.517/\sqrt{n}$ (false positive rate $p \leq 0.01$).

$z = 1.073/\sqrt{n}, p \leq 0.1$					
	FPR↓	TNR↑	TPR↑	FNR↓	PPL↓
Soft ($\delta=1.0$)	0.0545	0.9455	0.8919	0.2686	3.38±0.06
Soft ($\delta=1.5$)	0.0545	0.9455	0.9961	0.0796	3.56±0.06
Soft ($\delta=2.0$)	0.0545	0.9455	1.0000	0.0000	3.92±0.07
DiPmark ($\alpha=0.45$)	0.0545	0.9455	1.0000	0.0000	3.14±0.06
DiPmark ($\alpha=0.5$)	0.0545	0.9455	1.0000	0.0000	3.17±0.05
$z = 1.517/\sqrt{n}, p \leq 0.01$					
	FPR↓	TNR↑	TPR↑	FNR↓	PPL↓
Soft ($\delta=1.0$)	0.0080	0.9920	0.8255	0.1745	3.38±0.06
Soft ($\delta=1.5$)	0.0080	0.9920	0.9724	0.0276	3.56±0.06
Soft ($\delta=2.0$)	0.0080	0.9920	0.9981	0.0019	3.92±0.07
DiPmark ($\alpha=0.45$)	0.0080	0.9920	0.9794	0.0206	3.14±0.06
DiPmark ($\alpha=0.5$)	0.0080	0.9920	0.9827	0.0173	3.17±0.05

cal error rates for watermark detection in text generation are reported in Table 6, showcasing the commendable performance of DiPmark with low false positive rates while maintaining a high true positive rate. Broadly speaking, DiPmark with $\alpha = 0.45$ and 0.5 exhibit performance comparable to that of the Soft watermark with $\delta = 1.5$ and 2.0 . For more experimental results regarding the detectability, please refer to Appendix F.2.

7.5. Case study: watermarking GPT-4 by DiPmark

Recently, GPT-4 released the log-probability of the top-5 tokens during the generation process. This advancement enables us to modify and apply our DiPmark approach to GPT-4’s framework. As we only know the probability of the top-5 tokens, we treat the probability of the rest tokens as 0. Given a prompt, we will first use GPT-4 to generate the top-5 log-probability of the next token. Then we adapt DiPmark to the log-probability and sampling the next token based on the reweighted distribution. Finally, we merge the generated token into the prompt, and repeat the above steps. In our experiments, we use `gpt-4-0613` on 100 different fiction writing prompts and restrict the number of generated token to 200. We set $\alpha = 0.45$ in our DiPmark model.

In Figure 5, we show the cumulative histogram of the number of green tokens in the 100 watermarked GPT-4 generated sequences. As all generated sequences have 200 tokens, any sequence with greater than 122 green tokens can be detected as watermarked content with FPR less than 1%. From the plot, we see 97 out of 100 generated sequences can be detected by our algorithm, which validate the applicability of our watermark on the industry-level LLMs.

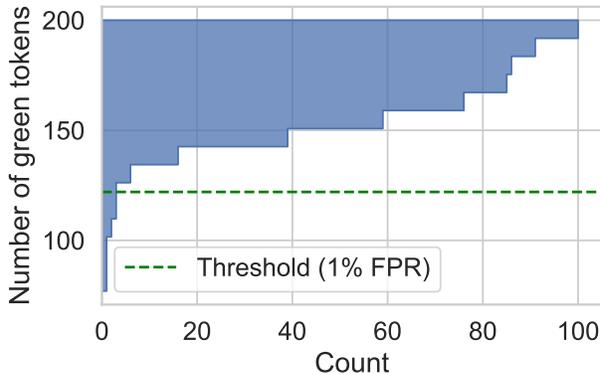


Figure 5. Cumulative histogram of the number of green tokens in the 100 watermarked gpt-4 generated sequences. The green line represents the threshold with FPR smaller than 1%.

8. Conclusion

In summary, we present DiPmark, a novel watermarking solution tailored for LLMs. DiPmark exhibits the crucial attributes of distribution-preserving, accessibility, and resilience, which we rigorously substantiate through a combination of theoretical analysis and empirical investigations. Our work not only strengthens the theoretical foundations, but also imparts practical insights that are valuable for the industrial deployment of LLM watermarking technologies.

Impact Statement

Machine learning holds significant potential to enhance human life, however, its malicious applications could substantially jeopardize safety (Wu et al., 2022; Hong et al., 2024; Hu et al., 2023b; Wang et al., 2023b;a; Wu et al., 2023; Chen et al., 2024). This research focuses on advancing watermark techniques to effectively identify AI-generated sentences. In an era where AI’s role in content creation is expanding rapidly, our work gains significance in preserving the authenticity and integrity of digital text. This innovation is pivotal in distinguishing human-authored content from that produced by AI, a distinction that holds substantial value across various societal and technological domains, e.g., enhancing digital content authenticity, combating misinformation, and empowering content creators.

Acknowledgement

This work was partially supported by NSF IIS 2347592, 2347604, 2348159, 2348169, DBI 2405416, CCF 2348306, CNS 2347617; HY Zhang was supported by NSERC Discovery Grant RGPIN-2022-03215, DGEGR-2022-00357.

References

- Aaronson, S. My AI safety lecture for UT effective altruism,. 2022. URL <https://scottaaronson.blog/?p=6823>.
- Abdelnabi, S. and Fritz, M. Adversarial watermarking transformer: Towards tracing text provenance with data hiding. In *2021 IEEE Symposium on Security and Privacy (SP)*, pp. 121–140. IEEE, 2021.
- Chakraborty, S., Bedi, A. S., Zhu, S., An, B., Manocha, D., and Huang, F. On the possibilities of AI-generated text detection. *arXiv preprint arXiv:2304.04736*, 2023.
- Chen, R., Wu, Y., Chen, L., Liu, G., He, Q., Xiong, T., Liu, C., Guo, J., and Huang, H. Your vision-language model itself is a strong filter: Towards high-quality instruction tuning with data selection. *arXiv preprint arXiv:2402.12501*, 2024.
- Christ, M., Gunn, S., and Zamir, O. Undetectable watermarks for language models. *arXiv preprint arXiv:2306.09194*, 2023.
- Dedić, N., Itkis, G., Reyzin, L., and Russell, S. Upper and lower bounds on black-box steganography. *Journal of Cryptology*, 22:365–394, 2009.
- Gambini, M., Fagni, T., Falchi, F., and Tesconi, M. On pushing DeepFake Tweet detection capabilities to the limits. In *Proceedings of the 14th ACM Web Science Conference 2022*, pp. 154–163, 2022.
- Google. Palm-2-llm. <https://blog.google/technology/ai/google-palm-2-ai-large-language-model/>, 2023.
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., and Blunsom, P. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28, 2015.
- Hong, Z., Wang, Z., Shen, L., Yao, Y., Huang, Z., Chen, S., Yang, C., Gong, M., and Liu, T. Improving non-transferable representation learning by harnessing content and style. In *The Twelfth International Conference on Learning Representations*, 2024.
- Hopper, N. J., Langford, J., and Von Ahn, L. Provably secure steganography. In *Advances in Cryptology—CRYPTO 2002: 22nd Annual International Cryptology Conference Santa Barbara, California, USA, August 18–22, 2002 Proceedings 22*, pp. 77–92. Springer, 2002.
- Hu, Z., Chen, L., Wu, X., Wu, Y., Zhang, H., and Huang, H. Unbiased watermark for large language models. *arXiv preprint arXiv:2310.10669*, 2023a.
- Hu, Z., Shen, L., Wang, Z., Wu, B., Yuan, C., and Tao, D. Learning to learn from APIs: Black-box data-free meta-learning. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 13610–13627. PMLR, 2023b.
- Kaptchuk, G., Jois, T. M., Green, M., and Rubin, A. D. Meteor: Cryptographically secure steganography for realistic distributions. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1529–1548, 2021.
- Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., and Goldstein, T. A watermark for large language models. *arXiv preprint arXiv:2301.10226*, 2023.
- Kirchner, J. H., Ahmad, L., Aaronson, S., and Leike, J. New AI classifier for indicating AI-written text. *OpenAI*, 2023.
- Krishna, K., Song, Y., Karpinska, M., Wieting, J., and Iyyer, M. Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense. *arXiv preprint arXiv:2303.13408*, 2023.
- Kuditipudi, R., Thickstun, J., Hashimoto, T., and Liang, P. Robust distortion-free watermarks for language models. *arXiv preprint arXiv:2307.15593*, 2023.
- Lin, C.-Y. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Liu, Y., Gu, J., Goyal, N., Li, X., Edunov, S., Ghazvininejad, M., Lewis, M., and Zettlemoyer, L. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8: 726–742, 2020.
- Mitchell, E., Lee, Y., Khazatsky, A., Manning, C. D., and Finn, C. Detectgpt: Zero-shot machine-generated text detection using probability curvature. *arXiv preprint arXiv:2301.11305*, 2023.
- Munyer, T. and Zhong, X. Deeptextmark: Deep learning based text watermarking for detection of large language model generated text. *arXiv preprint arXiv:2305.05773*, 2023.
- OpenAI, R. Gpt-4 technical report. *arXiv*, pp. 2303–08774, 2023.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.

- Qiang, J., Zhu, S., Li, Y., Zhu, Y., Yuan, Y., and Wu, X. Natural language watermarking via paraphraser-based lexical substitution. *Artificial Intelligence*, 317:103859, 2023.
- Tay, Y., Bahri, D., Zheng, C., Brunk, C., Metzler, D., and Tomkins, A. Reverse engineering configurations of neural text generation models. *arXiv preprint arXiv:2004.06201*, 2020.
- Tian, E. GPTzero update v1. <https://gptzero.substack.com/p/gptzero-update-v1>, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Wang, Z., Shen, L., Duan, T., Suo, Q., Fang, L., Liu, W., and Gao, M. Distributionally robust memory evolution with generalized divergence for continual learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023a.
- Wang, Z., Shen, L., Liu, T., Duan, T., Zhu, Y., Zhan, D., Doermann, D., and Gao, M. Defending against data-free model extraction by distributionally robust defensive training. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- Wu, Y., Zhang, H., and Huang, H. Retrievalguard: Provably robust 1-nearest neighbor image retrieval. In *International Conference on Machine Learning*, pp. 24266–24279. PMLR, 2022.
- Wu, Y., Huang, H., and Zhang, H. A law of robustness beyond isoperimetry. In *International Conference on Machine Learning*, pp. 37439–37455. PMLR, 2023.
- Yoo, K., Ahn, W., Jang, J., and Kwak, N. Robust natural language watermarking through invariant features. *arXiv preprint arXiv:2305.01904*, 2023.
- Zellers, R., Holtzman, A., Rashkin, H., Bisk, Y., Farhadi, A., Roesner, F., and Choi, Y. Defending against neural fake news. *Advances in neural information processing systems*, 32, 2019.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., and Artzi, Y. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- Zhao, X., Ananth, P., Li, L., and Wang, Y.-X. Provable robust watermarking for ai-generated text. *arXiv preprint arXiv:2306.17439*, 2023.

A. Future Work

Future endeavors should focus on enhancing the detectability of distribution-preserving watermarks. This could be realized by assigning greater weight to the green-list tokens during the watermarking process. Additionally, a promising avenue for exploration involves the design of a more robust distribution-preserving watermark, potentially through the integration of multiple detectors. These directions represent promising opportunities for advancing the efficacy and applicability of watermarking techniques on large language models.

B. Related Work

Reweight-based watermarking framework. In a recent seminal work, (Kirchenbauer et al., 2023) introduced a pioneering watermarking scheme tailored for LLMs, backed by formal guarantees. Their work demonstrated that watermark embedding could be accomplished by altering the token distribution during generation, targeting outputs with substantial entropy. However, this approach inevitably leads to a pivotal change in the distribution of the generated text, potentially compromising the quality of the generated content.

To maintain an unaltered output distribution in watermarked content, alternative strategies have been explored. (Christ et al., 2023) and (Kuditipudi et al., 2023) employed the inverse sampling method to generate watermarked token distributions. Notably, (Christ et al., 2023)’s method faces resilience issues under modifications and lacks empirical validation for detectability. Meanwhile, (Kuditipudi et al., 2023)’s approach necessitates the secret key distribution during detection, potentially compromising data security and watermark stealthiness. Moreover, their detection process involves hundreds of resampling steps from the secret key distribution, which is inefficient for lengthy texts. (Hu et al., 2023a) used inverse sampling and permutation based reweight methods for watermarking, but the detector requires access of the language model API, undermining its operational efficiency. Aaronson’s ongoing watermarking project (Aaronson, 2022) employs n-gram hashing for reweighting the next-token distribution, though specific details are currently unavailable.

The landscape also includes several schemes (Abdelnabi & Fritz, 2021; Qiang et al., 2023; Yoo et al., 2023; Munyer & Zhong, 2023) that incorporate an ML model within the watermarking algorithm itself. However, these constructions lack formal assurances and rely on heuristic arguments for satisfying the criteria of Stealthiness, Efficiency, and Resilience.

Our research aligns closely with the findings presented in (Kirchenbauer et al., 2023). In their methodology, they employed watermarking for text derived from a language model by bifurcating the token set into designated ‘red’ and ‘green’ lists. The division is determined by a random seed that is contingent on the secret key coupled with a hash of priorly generated tokens. The authors accentuated the prominence of green tokens during the sampling phase by reweighting the token log-probabilities. Building on this foundation, our research retains the red-green list configuration, but introduces an evolved family of permutation-based reweight strategies. This dual approach ensures: 1) a promoted utilization of green tokens, and 2) equivalency in distribution between a sample from the watermarked language model and one from the original language model.

Post-hoc detectors. Post-hoc detection stands as a notable alternative to watermarking, focusing on the retrospective analysis of machine-generated text. This could be achieved through leveraging features inherent to language models or by refining pre-existing, expansive language models to function as detectors, as elaborated by (Zellers et al., 2019). Notably, specific implementation nuances, such as sampling methodologies, can be discerned through reverse engineering the generated text, a process detailed by (Tay et al., 2020). There are also post-hoc detectors designed for the modern large language models (Mitchell et al., 2023; Tian, 2023; Kirchner et al., 2023), which are models specifically trained for the binary detection task. However, there is a growing sentiment that those detection methodologies are diminishing in efficacy in tandem with the evolution of language model capabilities. As (Gambini et al., 2022) observed, detection mechanisms that were adept with GPT-2 have encountered challenges with GPT-3. Besides, the text rephrasing model in (Krishna et al., 2023) bypassing prevalent post-hoc detectors like GPTZero (Tian, 2023), DetectGPT (Mitchell et al., 2023), and OpenAI’s proprietary detector (Kirchner et al., 2023). Additionally, a pertinent observation made by (Chakraborty et al., 2023) suggests that as AI-generated content becomes increasingly indistinguishable from human-produced text, the demands on post-hoc detectors to analyze more extended text segments will escalate.

Steganography. Steganography involves embedding concealed messages in channels such as natural language or images, ensuring only intended recipients can discern the message while others remain unaware (Hopper et al., 2002). When applied to watermarking, the aim is stealthy. Yet, known steganography techniques might not achieve this without certain entropy-related assumptions. In scenarios where language model prompts can be chosen adversarially, the need for stealthy

persists. This discrepancy arises due to differences in access levels that watermarking and steganography have to the model’s output distribution. In steganography, there’s only oracle access to this distribution. Conversely, our watermarking approach gets a detailed view of the token’s probability distribution. Hence, while steganography either relies on entropy assumptions (Hopper et al., 2002) or compromises security with low entropy channels (Dedić et al., 2009), our watermark remains stealthy irrespective of the text’s entropy. This is achieved by leveraging the full distribution access and using it as a foundation for embedding watermarks. (Kaptchuk et al., 2021) offers encoding similar access. However, it presupposes equal decoding access, which is impractical for watermarking as the detection algorithm won’t typically have the initiating prompt, thus remaining ignorant of the distribution.

C. Missing Proofs

C.1. Proof of Theorem 4.3

Proof. We need to show $\forall t \in V, \mathbb{E}_\theta[P_W(t|\mathbf{x}, \theta)] = P_M(t|\mathbf{x})$. Recall θ is uniformly distributed on Θ , we have

$$\begin{aligned} \mathbb{E}_{\theta \sim P_\Theta}[P_W(t|\mathbf{x}, \theta)] &= \sum_{V^p \in \Theta} \mathbb{E}_{\theta \sim P_\Theta}[P_W(t|\mathbf{x}, V^p) \mathbf{1}_{\theta=V^p}] \\ &= \sum_{V^p \in \Theta} [P_W(t|\mathbf{x}, V^p)] \mathbb{E}_{\theta \sim P_\Theta}[\mathbf{1}_{\theta=V^p}] \\ &= \frac{1}{N!} \sum_{V^p \in \Theta} P_W(t|\mathbf{x}, V^p). \end{aligned} \quad (1)$$

Given an token t and a permutation of the token list V^p , denote by $E_{V^p}(t)$ the position of t in the ordered token set V^p . Let V^{p^r} be the reversed permutation of V^p , notice t is the $(N + 1 - E_{V^p}(t))$ -th element in V^{p^r} . Given an arbitrary permutation pair (V^p, V^{p^r}) , $V^p := \{t_1, \dots, t_N\}$. We will show

$$P_W(t|\mathbf{x}, V^p) + P_W(t|\mathbf{x}, V^{p^r}) = 2P_M(t|\mathbf{x}).$$

For the ease of notation we denote by $i = E_{V^p}(t)$, we have $t_i = t$. From the definition of DiP-reweight we know $P_W(t|\mathbf{x}, V^p) = F(E_{V^p}(t)|V^p) - F(E_{V^p}(t) - 1|V^p) = F(i|V^p) - F(i - 1|V^p)$, where

$$F(i|V^p) := \max \left\{ \sum_{j=1}^i P_M(t_j|\mathbf{x}) - \alpha, 0 \right\} + \max \left\{ \sum_{j=1}^i P_M(t_j|\mathbf{x}) - (1 - \alpha), 0 \right\}, \quad i \in [1, N], \quad (2)$$

So we need to show

$$F(i|V^p) - F(i - 1|V^p) + F(N + 1 - i|V^{p^r}) - F(N - i|V^{p^r}) = 2P_M(t|\mathbf{x}).$$

As $\sum_{j=1}^N P_M(t_j|\mathbf{x}) = 1$, we have

$$\begin{aligned} F(N + 1 - i|V^{p^r}) &= \max \left\{ \sum_{j=1}^{N+1-i} P_M(t_{N+1-j}|\mathbf{x}) - \alpha, 0 \right\} + \max \left\{ \sum_{j=1}^{N+1-i} P_M(t_{N+1-j}|\mathbf{x}) - (1 - \alpha), 0 \right\} \\ &= \max \left\{ \sum_{j=i}^N P_M(t_j|\mathbf{x}) - \alpha, 0 \right\} + \max \left\{ \sum_{j=i}^N P_M(t_j|\mathbf{x}) - (1 - \alpha), 0 \right\} \\ &= \max \left\{ (1 - \alpha) - \sum_{j=i}^{i-1} P_M(t_j|\mathbf{x}), 0 \right\} + \max \left\{ \alpha - \sum_{j=1}^{i-1} P_M(t_j|\mathbf{x}), 0 \right\}, \end{aligned} \quad (3)$$

and

$$F(i - 1|V^p) = \max \left\{ \sum_{j=1}^{i-1} P_M(t_j|\mathbf{x}) - \alpha, 0 \right\} + \max \left\{ \sum_{j=1}^{i-1} P_M(t_j|\mathbf{x}) - (1 - \alpha), 0 \right\}. \quad (4)$$

By $(\max\{A, 0\} - \max\{-A, 0\}) = A, \forall A \in \mathbb{R}$, we have

$$\begin{aligned} F(N+1-i|V^{P^r}) - F(i-1|V^P) &= (1-\alpha) - \sum_{j=i}^{i-1} P_M(t_j|\mathbf{x}) + \alpha - \sum_{j=1}^{i-1} P_M(t_j|\mathbf{x}) \\ &= 1 - 2 \sum_{j=i}^{i-1} P_M(t_j|\mathbf{x}). \end{aligned} \quad (5)$$

Analogously, we have

$$F(N-i|V^{P^r}) - F(i|V^P) = 1 - 2 \sum_{j=i}^i P_M(t_j|\mathbf{x}). \quad (6)$$

Thus,

$$\begin{aligned} P_W(t|\mathbf{x}, V^P) + P_W(t|\mathbf{x}, V^{P^r}) &= F(i|V^P) - F(i-1|V^P) + F(N+1-i|V^{P^r}) - F(N-i|V^{P^r}) \\ &= (1 - 2 \sum_{j=i}^{i-1} P_M(t_j|\mathbf{x})) - (1 - 2 \sum_{j=i}^i P_M(t_j|\mathbf{x})) \\ &= 2P_M(t_i|\mathbf{x}) = 2P_M(t|\mathbf{x}). \end{aligned} \quad (7)$$

By the symmetric of permutation we have

$$\begin{aligned} 2\mathbb{E}_{\theta \sim \Theta}[P_W(t|\mathbf{x}, \theta)] &= \frac{1}{N!} \sum_{V^P \in \Sigma} P_W(t|\mathbf{x}, V^P) \\ &= \frac{1}{N!} \sum_{V^P \in \Sigma} [P_W(t|\mathbf{x}, V^P) + P_W(t|\mathbf{x}, V^{P^r})] \\ &= \frac{1}{N!} \sum_{V^P \in \Sigma} 2P_M(t|\mathbf{x}) \\ &= 2P_M(t|\mathbf{x}). \end{aligned} \quad (8)$$

Therefore, $\mathbb{E}_{\theta \sim \Theta}[P_W(t|\mathbf{x}, \theta)] = P_M(t|\mathbf{x})$, which concludes the proof. \square

C.2. Proof of Theorem 5.2

Proof. As $L_G(\gamma) = \sum_{i=1}^n B_i(\gamma)$, where $B_i(\gamma) \sim \text{Bernoulli}(1-\gamma)$. By Markov's inequality we have $\forall h > 0$,

$$\Pr(L_G(\gamma) - (1-\gamma)n \geq nt) \leq \frac{\mathbb{E}[e^{h(L_G(\gamma) - (1-\gamma)n)}]}{e^{hnt}},$$

as B_i is independent from each other, we have

$$\frac{\mathbb{E}[e^{h(L_G(\gamma) - (1-\gamma)n)}]}{e^{hnt}} = \prod_{i=1}^n \frac{\mathbb{E}[e^{h(B_i - (1-\gamma))}]}{e^{ht}}.$$

Since B_i follows Bernoulli distribution, we have

$$\mathbb{E}[e^{h(B_i - (1-\gamma))}]/e^{ht} = (1-\gamma)e^{h(\gamma-t)} + \gamma e^{-h(1-\gamma+t)}.$$

Thus

$$\Pr(L_G(\gamma) - (1-\gamma)n \geq nt) \leq [(1-\gamma)e^{h(\gamma-t)} + \gamma e^{-h(1-\gamma+t)}]^n \quad (9)$$

holds for arbitrary $h > 0$. Denote by $m(h) = (1-\gamma)e^{h(\gamma-t)} + \gamma e^{-h(1-\gamma+t)}$, taking derivative w.r.t. h yields

$$\frac{dm(h)}{dh} = (1-\gamma)(\gamma-t)e^{h(\gamma-t)} + \gamma(1-\gamma+t)e^{-h(1-\gamma+t)}.$$

Let $\frac{dm(h)}{dh} = 0$, we have $h = \ln \frac{\gamma(1-\gamma+t)}{(1-\gamma)(\gamma-t)}$. Combining it with Equation 9 yields

$$\begin{aligned}
 \Pr(L_G(\gamma) - (1-\gamma)n \geq nt) &\leq \inf_{h>0} [(1-\gamma)e^{h(\gamma-t)} + \gamma e^{-h(1-\gamma+t)}]^n \\
 &\leq [e^{(\gamma-t) \ln \frac{\gamma(1-\gamma+t)}{(1-\gamma)(\gamma-t)}} (1-\gamma + \gamma e^{-\ln \frac{\gamma(1-\gamma+t)}{(1-\gamma)(\gamma-t)}})]^n \\
 &= [e^{(\gamma-t) \ln \frac{\gamma(1-\gamma+t)}{(1-\gamma)(\gamma-t)}} \frac{1-\gamma}{1-\gamma+t}]^n \\
 &= [e^{(\gamma-t) \ln \frac{\gamma(1-\gamma+t)}{(1-\gamma)(\gamma-t)} + \ln \frac{1-\gamma}{1-\gamma+t}}]^n \\
 &= e^{-n((1-\gamma+t) \ln \frac{1-\gamma+t}{1-\gamma} + (\gamma-t) \ln \frac{\gamma-t}{\gamma})} \\
 &= e^{-n\text{KL}(1-\gamma+t||1-\gamma)}
 \end{aligned} \tag{10}$$

□

C.3. Proof of Theorem 6.2 and discussion

Proof. Notice based on above discussion, the worst-case decrease on $L_G(\gamma)$ per token modification is $a + 1$. If we are allowed to perturbed ϵ portion of the text, the worst-case decrease on $L_G(\gamma)$ will be $(a + 1)\epsilon n$. Denoted by $\mathbf{x}_{1:n'}$ the perturbed text. Assume we can still correctly detect the watermarked sequence, which means

$$(L_G(\gamma) - (a + 1)\epsilon n)/n' - (1 - \gamma) \geq z.$$

Notice, the left hand side of the above equation is decreasing with n' , as we perturbed ϵ portion of the text, the maximum of the possible n' is $n' = (1 + \epsilon)n$, i.e., all modifications are text insertion. In this case, we need to solve

$$\frac{L_G(\gamma) - (a + 1)\epsilon n}{(1 + \epsilon)n} - (1 - \gamma) \geq z.$$

we have

$$\epsilon \leq \frac{L_G(\gamma) - (1 - \gamma)n - zn}{(2 + a - \gamma + z)n}.$$

Therefore, for any text modification with budget $\epsilon \leq \frac{L_G(\gamma) - (1 - \gamma)n - zn}{(2 + a - \gamma + z)n}$, our algorithm can still detect the watermarked sequence. □

In the following theorem, we provide a more simple certified radius assuming the text length is not changed by perturbations.

Theorem C.1. *Assuming the sequence length n is not changed through text modifications. Given $\Phi(\gamma, \mathbf{x}_{1:n}) := L_G(\gamma)/n - (1 - \gamma)$ and a threshold z , the certified radius of the watermarked sequence $\mathbf{x}_{1:n}$ is $\epsilon_0 = \frac{\Phi(\gamma, \mathbf{x}_{1:n}) - z}{a+1}$.*

Proof. Notice based on above discussion, the worst-case decrease on $L_G(\gamma)$ per token modification is $a + 1$. If we are allowed to perturbed ϵ portion of the text, the worst-case decrease on $L_G(\gamma)$ will be $(a + 1)\epsilon n$. Assume we can still correctly detect the watermarked sequence, which means

$$(L_G(\gamma) - (1 - \gamma)n - (a + 1)\epsilon n)/\sqrt{n} \geq z,$$

we have $\epsilon \leq \frac{\Phi(\gamma, \mathbf{x}_{1:n}) - z}{(a+1)\sqrt{n}}$. Therefore, for any text modification with budget $\epsilon \leq \frac{\Phi(\gamma, \mathbf{x}_{1:n}) - z}{(a+1)\sqrt{n}}$, our algorithm can still detect the watermarked sequence. □

D. Comparison of the test statistic

In this section, we provide a detailed comparison of our test statistic and the z-test statistic proposed in (Kirchenbauer et al., 2023). In Figure 6, we show number of green tokens vs p-value (false positive rate), where we set the number of tokens

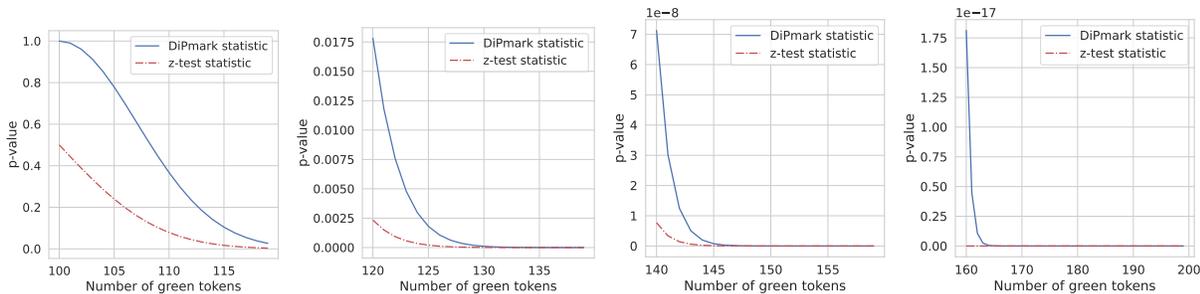


Figure 6. Number of green tokens vs p-value (false positive rate), where we set the number of tokens $n = 200$, green list separator $\gamma = 0.5$. We see that given the same number of green tokens, the z-test always has lower p-value than DiPmark test statistic. Given the fact that the z-test statistic is only an approximation of the green token distribution, we conclude that this approximation is not proper for watermark detection, as it will wrongly classify the sentences not generated by LMs as being LM-produced.

Table 7. Comparison of test statistics: Theoretical FPR vs Empirical FPR. We can see clearly the empirical FPR of z-test is continuously greater than its theoretical guarantee, which indicates z-test statistic may not be suitable for watermark detection.

	$p < 0.10$ (10%FPR)	$p < 0.05$ (5%FPR)	$p < 0.01$ (1%FPR)
z-test (Kirchenbauer et al., 2023)	56/500 (11.2% FPR)	34/500 (6.8% FPR)	12/500 (2.4% FPR)
DiPmark statistic	13/500 (2.6% FPR)	10/500 (2% FPR)	4/500 (0.5% FPR)

$n = 200$, green list separator $\gamma = 0.5$. We see that given the same number of green tokens, the z-test statistic always leads to lower p-value than DiPmark test statistic. Given the fact that the z-test statistic is only an approximation of the green token distribution, we conclude that this approximation is not proper for watermark detection, as it will wrongly classify the sentences not generated by LMs as being LM-produced. In Table 7, we show the detecting result based on DiPmark detector and the detector in Kirchenbauer et al. (2023) on 500 non-watermarked sentences with length 260. We can see clearly the empirical FPR of z-test is continuously greater than its theoretical guarantee, which indicates z-test statistic may not be suitable for watermark detection.

E. Detailed Experiment Setup

We assess the performance of DiPmark across three critical applications of seq2seq models: text summarization, machine translation, and text generation. The experiments are implemented using the Huggingface library (Wolf et al., 2019), a widely adopted platform for model development and sharing within the NLP community. All experiments are conducted on three Nvidia A6000 GPUs with 48GB of memory. Detecting 1,000 watermarked sentences generated from LLaMA-2 requires only 90 seconds.

Machine Translation. For the machine translation task, we utilize the WMT’ 14 English (En) to Romanian (Ro) dataset, comprising 1,999 examples in the test set. We employ the Multilingual Bart (MBart) model (Liu et al., 2020) along with its official tokenizer.

Text Summarization. In the text summarization task, we use the test set from the CNN-DM corpus (Hermann et al., 2015), consisting of 11,490 examples. Our model of choice is BART-large, which encompasses 400 million parameters, and LLaMA-2 with 7 billion parameters.

Text Generation. For text generation, we incorporate the test set from the CNN-DM corpus as part of the generation prompt. We use LLaMA-2 which has 7 billion parameters.

Watermark Setup. Our experiments primarily compare DiPmark with the Soft watermark introduced by (Kirchenbauer et al., 2023). In the case of DiPmark, we consider various values of α from the set $\{0.3, 0.35, 0.4, 0.45, 0.5\}$. For the Soft watermark (Kirchenbauer et al., 2023), we explore green list bias δ values from $\{0.0, 1.0, 1.5, 2.0\}$ with a fixed green list separator $\gamma = 0.5$. Texture key generation relies on the most recent five tokens as texture key. For instance, when generating x_4 in response to (x_1, x_2, x_3) as the current input to the decoder, the texture key includes (x_1, x_2, x_3) , considering the availability of only three tokens. The texture key history resets before generating each batch. To generate the cipher, we employ SHA-256 as the hash function and a set of 1024-bit random bitstrings as the key set K . The cipher θ is sampled from Θ using $\text{hash}(k, s)$ as the random seed. We compare DiPmark with ITS (Kuditipudi et al., 2023) and δ -watermark (Hu

Table 8. Performance of Machine Translation.

	BERT-F1	BERT-Precision	BERT-Recall	BLEU
No Watermark	0.559±0.003	0.545±0.004	0.574±0.003	21.8±0.3
DiPmark($\alpha=0.3$)	0.561±0.003	0.547±0.004	0.575±0.003	22.0±0.3
DiPmark($\alpha=0.35$)	0.562±0.003	0.548±0.004	0.575±0.003	22.1±0.3
DiPmark($\alpha=0.4$)	0.561±0.003	0.547±0.004	0.576±0.003	21.9±0.3
DiPmark($\alpha=0.45$)	0.562±0.003	0.548±0.004	0.576±0.003	21.9±0.3
DiPmark($\alpha=0.5$)	0.562±0.003	0.548±0.004	0.576±0.003	21.8±0.3
Soft($\delta=0.0$)	0.560±0.003	0.545±0.004	0.574±0.003	21.8±0.3
Soft($\delta=1.0$)	0.557±0.003	0.543±0.004	0.572±0.003	21.2±0.3
Soft($\delta=1.5$)	0.550±0.003	0.534±0.004	0.565±0.003	20.4±0.3
Soft($\delta=2.0$)	0.539±0.003	0.523±0.004	0.555±0.003	19.4±0.3

et al., 2023a), where we follow the setting in their open sourced code²³.

Evaluation metrics for text quality. In this part, we introduce the evaluation metrics we used for evaluating the text quality (Section. 7.1).

- **ROUGE score.** For the summarization task, we utilize the ROUGE score (Lin, 2004), which measures n-gram overlap to assess the summary’s effectiveness in capturing essential content from reference summaries.
- **BLEU score.** For the machine translation task, we rely on the BLEU score (Papineni et al., 2002), emphasizing the lexical similarity between machine-generated translations and human reference translations.
- **BERTScore.** BERTScore (Zhang et al., 2019) computes the similarity of two sentences as a sum of cosine similarities between their tokens’ embeddings. We use BERTScore-F1, BERTScore-Precision, and BERTScore-Recall for evaluating both text summarization and machine translation tasks.
- **Perplexity.** In information theory, perplexity is a measurement of how well a probability distribution or probability model predicts a sample. It may be used to compare probability models. A low perplexity indicates the probability distribution is good at predicting the sample. We use perplexity for evaluating both text summarization and machine translation tasks.

Evaluation metrics for detectability of watermarks. In this part, we introduce the evaluation metrics we used for evaluating the detectability of watermarks (Sections 7.4 and 7.3).

- **Green token ratio.** Denoted by $L_G(\gamma)$ the number of green tokens in a text sequence with green list separator γ . The green token ratio is given by $L_G(\gamma)/n - (1 - \gamma)$. This ratio quantifies the bias towards green tokens within the text sequence (see Section 5).
- **z-score.** The z-score of a text sequence $x_{1:n}$ is $(L_G(\gamma) - (1 - \gamma)n)/\sqrt{n}$. A higher z-score will reduce the false positive rate, where a non-watermarked sequence is detected as watermarked (see Section 5).
- **Type I and II errors.** We generally use true positive rate (TPR), false positive rate (FPR), true negative rate (TNR), and false negative rate (FNR) to evaluate the performance of watermarks on a mixture of watermarked and non-watermarked sentence. FPR measures the Type I error of the hypothesis testing, in which the null hypothesis got rejected when it is actually true. FNR measures the type II error, in which one fails to reject a null hypothesis that is actually false.

Table 9. Performance of Text Summarization.

	BERT-F1	BERT-Precision	BERT-Recall	Perplexity	Rouge-1	Rouge-2	Rouge-L
No Watermark	0.3273±0.0008	0.3181±0.0009	0.3366±0.0010	5.021±0.018	0.3855±0.0009	0.1387±0.0008	0.2444±0.0008
DiPmark($\alpha=0.3$)	0.3279±0.0008	0.3187±0.0009	0.3372±0.0010	5.014±0.018	0.3861±0.0009	0.1390±0.0008	0.2450±0.0008
DiPmark($\alpha=0.35$)	0.3274±0.0008	0.3183±0.0009	0.3367±0.0010	4.998±0.018	0.3856±0.0009	0.1389±0.0008	0.2449±0.0008
DiPmark($\alpha=0.4$)	0.3277±0.0008	0.3187±0.0009	0.3370±0.0010	5.001±0.018	0.3862±0.0009	0.1392±0.0008	0.2449±0.0007
DiPmark($\alpha=0.45$)	0.3269±0.0008	0.3178±0.0009	0.3361±0.0010	5.024±0.018	0.3852±0.0009	0.1391±0.0008	0.2447±0.0008
DiPmark($\alpha=0.5$)	0.3272±0.0008	0.3181±0.0009	0.3364±0.0010	5.014±0.018	0.3859±0.0009	0.1396±0.0008	0.2450±0.0008
Soft($\delta=0.0$)	0.3273±0.0008	0.3181±0.0009	0.3366±0.0010	5.021±0.018	0.3855±0.0009	0.1387±0.0008	0.2444±0.0008
Soft($\delta=1.0$)	0.3237±0.0008	0.3137±0.0009	0.3338±0.0009	5.309±0.019	0.3816±0.0009	0.1348±0.0008	0.2411±0.0007
Soft($\delta=1.5$)	0.3209±0.0008	0.3097±0.0009	0.3323±0.0010	5.660±0.021	0.3793±0.0009	0.1317±0.0007	0.2379±0.0007
Soft($\delta=2.0$)	0.3146±0.0008	0.3027±0.0009	0.3266±0.0009	6.241±0.023	0.3725±0.0009	0.1252±0.0007	0.2321±0.0007

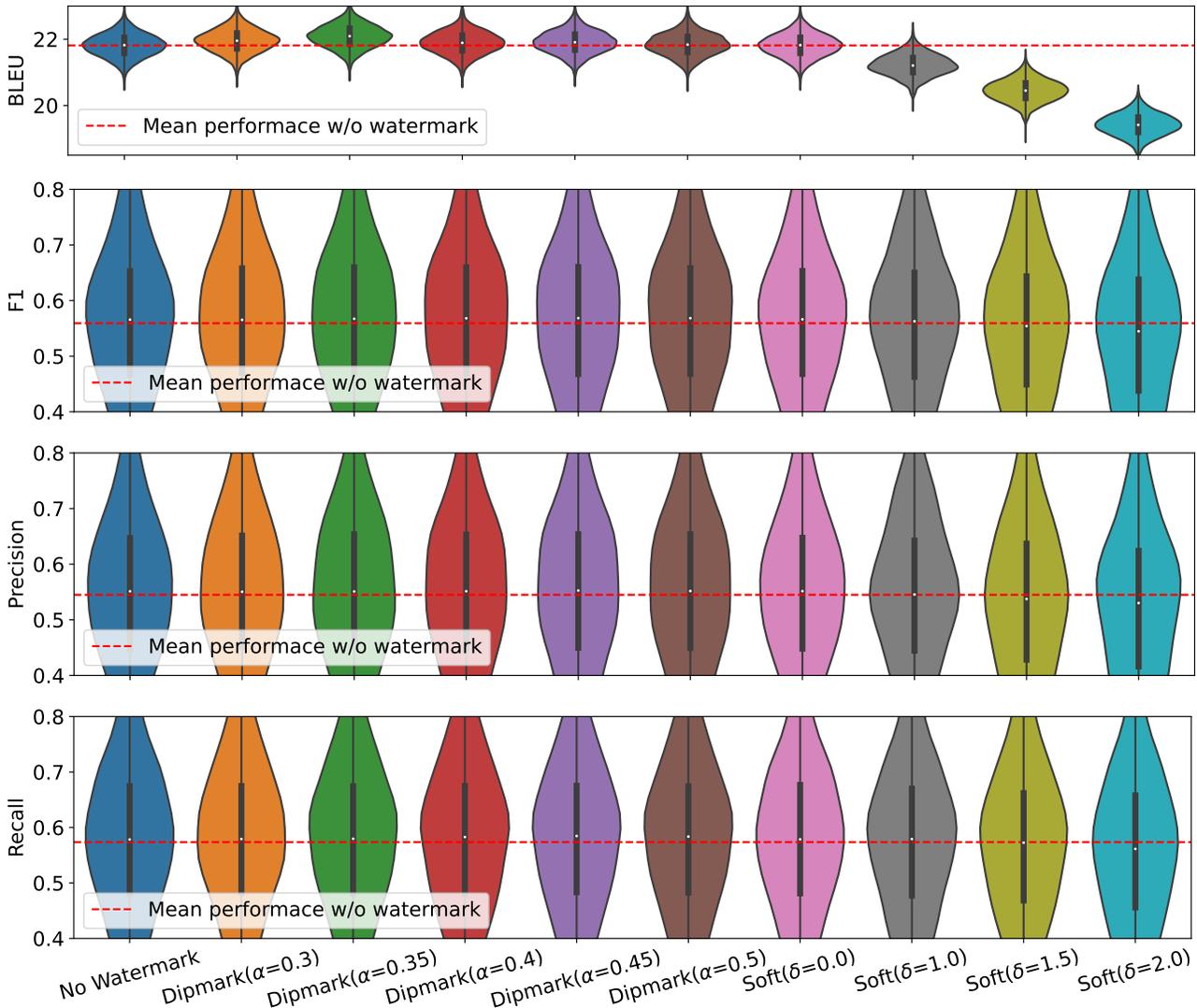


Figure 7. Violin plot of Machine Translation performance .

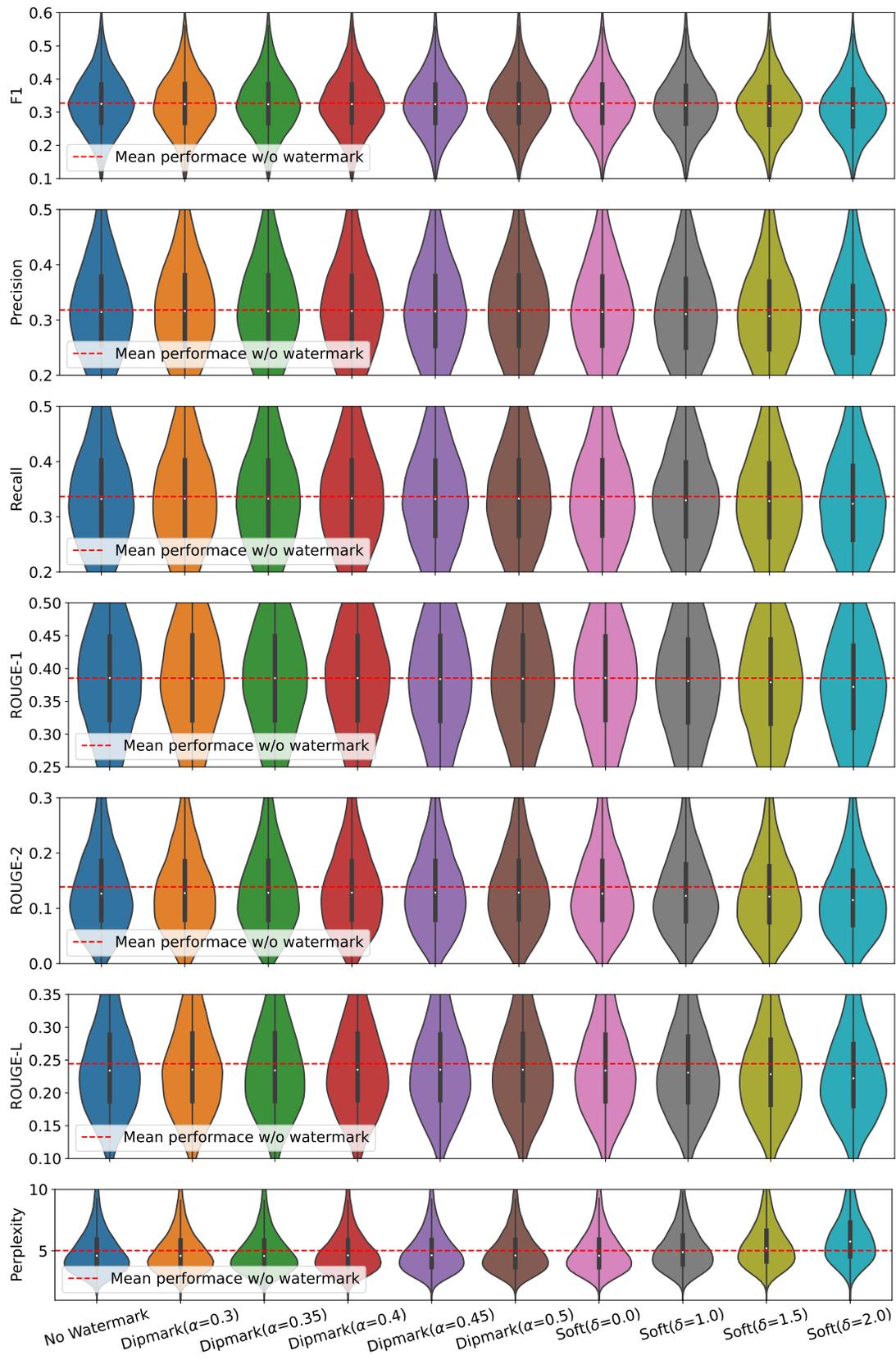


Figure 8. Violin plot of Text Summarization performance.

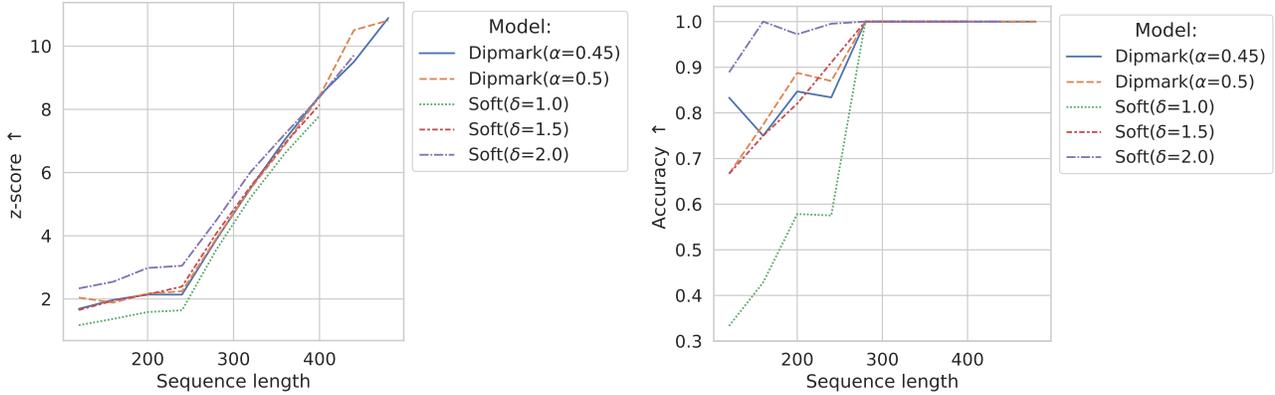


Figure 9. **Left:** Average z-score vs token sequence length with $\gamma = 0.5$ on text generation tasks. **Right:** Watermark detection accuracy vs token sequence length with $\gamma = 0.5$ and threshold $z = 1.517$ (false positive rate less than 0.01) on text generation tasks.

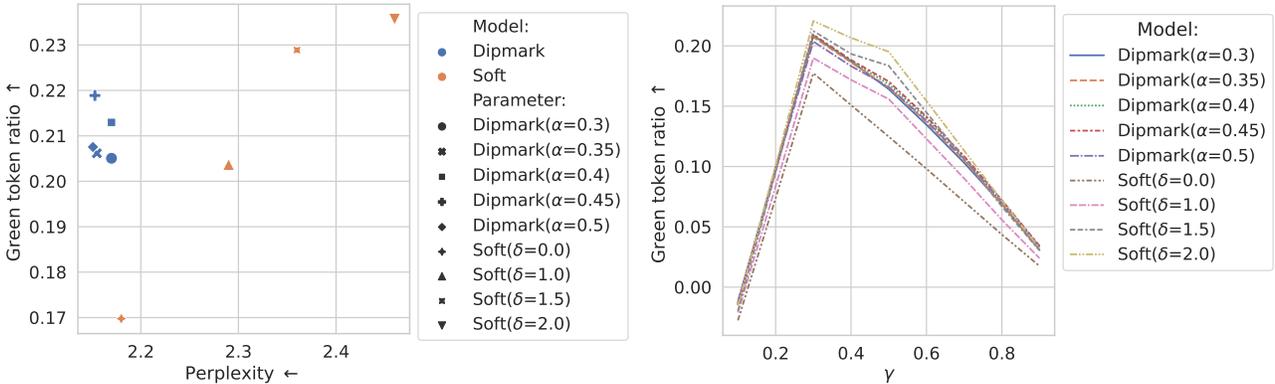


Figure 10. **Left.** Average Perplexity vs Green token rate with $\gamma = 0.5$ on the text summarization task. **Right.** Avg. Green token ratio with different γ on the text summarization task.

F. Additional Experiments

F.1. Distribution-preserving

Settings. In our evaluation, we assess the distribution-preserving performance of DiPmark within the context of two significant applications involving seq2seq models: machine translation (MT) and text summarization (TS). We follow the settings in (Hu et al., 2023a). For the TS task, our experimentation employs the BART-large model (Liu et al., 2020) in conjunction with the CNN-DM corpus (Hermann et al., 2015) as our designated testing dataset. The MT task, on the other hand, revolves around English-to-Romanian translation. For this purpose, we employ the Multilingual BART (MBart) model (Liu et al., 2020) on the WMT’14 En-Ro corpus. Specifically for DiPmark, we select values for α from the set $\{0.3, 0.35, 0.4, 0.45, 0.5\}$, while for the Soft watermark (Kirchenbauer et al., 2023), we choose green list bias values δ from the set $\{0.0, 1.0, 1.5, 2.0\}$ alongside a fixed green list separator $\gamma = 0.5$, indicating that 50% of tokens are green while the remainder are red. It is important to note that the Soft watermark with $\delta = 0.0$ is essentially equivalent to no watermark since it does not promote the probability of green list tokens.

A thorough examination of Figure 7, Figure 8, Table 8, and Table 9 reveals a discernible trend. Throughout the range of α values spanning $\{0.3, 0.35, 0.4, 0.45, 0.5\}$, all the metrics associated with machine translation tasks and text summarization tasks maintain a consistent alignment between DiPmark and the original language model. Conversely, an upward adjustment in the δ values of the Soft watermark distinctly impacts the quality of the text output.

²<https://github.com/jthickstun/watermark>

³<https://github.com/xiaoniu-578fa6bff964d005/UnbiasedWatermark>

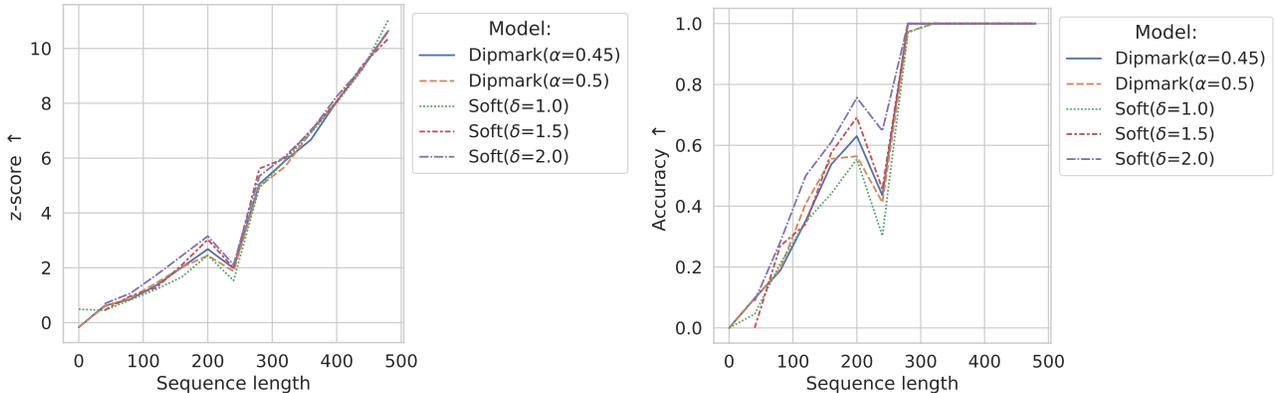


Figure 11. **Left.** Average z-score vs token sequence length with $\gamma = 0.5$ on the text summarization task. **Right.** Avg. best p-score with text length with $\gamma = 0.5$ on the text summarization task.

F.2. Detectability comparison

Settings. We evaluate the detectability of our watermark on text summarization tasks using LLaMA-2. We generate 1,000 examples for each tasks. We also select $\alpha \in \{0.3, 0.35, 0.4, 0.45, 0.5\}$ for DiPmark, and $\delta \in \{0.0, 1.0, 1.5, 2.0\}$ and $\gamma = 0.5$ for Soft watermark (Kirchenbauer et al., 2023). During detection, we also use $\gamma = 0.5$. We report the green token ratio (defined in 5), the score of $\Phi(\gamma, \mathbf{x})$ (z-score), and the detect accuracy.

Result analysis. The results for text generation are visually depicted in Figure 4 and Figure 9. Broadly speaking, our DiPmark variants with $\alpha = 0.45$ and 0.5 exhibit performance comparable to that of the Soft watermark with $\delta = 1.5$, where $\delta = 1.5$ corresponds to an augmentation of 1.5 to the green token logits. In Figure 4 (left), it is evident that our DiPmark variants with $\alpha = 0.45$ and 0.5 yield green token ratios akin to those of the Soft watermark with $\delta = 1.5$ without any discernible degradation in text quality. Figure 4 (right) delves into the impact of different green list separators γ , revealing that, for most watermark models, $\gamma = 0.5$ yields the highest green token ratio, underscoring its suitability as a reasonable choice for watermark detection. In Figure 9 (left) and Figure 9 (right), we present the average z-scores and accuracy metrics relative to sequence length. It is conspicuously observable that longer token sequences tend to facilitate easier detection, in line with our earlier analysis in Section 5. The results for text summarization are visually depicted in Figure 10 and Figure 11. Broadly speaking, our DiPmark variants with $\alpha = 0.45$ and 0.5 exhibit performance comparable to that of the Soft watermark with $\delta = 1.5$, where $\delta = 1.5$ corresponds to an augmentation of 1.5 to the green token logits. In Figure 10 (left), it is evident that our DiPmark variants with $\alpha = 0.45$ and 0.5 yield green token ratios akin to those of the Soft watermark with $\delta = 1.5$ without any discernible degradation in text quality. Figure 10 (right) delves into the impact of different green list separators γ . Interestingly, for most watermark models, $\gamma = 0.3$ yields the highest green token ratio instead of $\gamma = 0.5$, which may be due to the low entropy characteristic of the text summarization task. In Figure 11 (left) and Figure 11 (right), we present the average z-scores and accuracy metrics relative to sequence length. It is conspicuously observable that longer token sequences tend to facilitate easier detection, in line with our earlier analysis in Section 5.

F.3. Resilience

We conduct experiments to test the resiliency of the our DiPmark and the Soft watermark in (Kirchenbauer et al., 2023). In this context, we use the text summarization tasks with 1,000 generated sequences on LLaMA-2. For resilience evaluation, we manipulating about $\epsilon \in \{0.05, 0.1, 0.2, 0.3\}$ portion of the text tokens through text insertion, text substitution, and text deletion.

Result Analysis. Figure 13 elucidates the evolution of the average green token ratio and the average z-score concerning the attack strength parameter ϵ . Notably, both metrics exhibit a diminishing trend as ϵ increases.

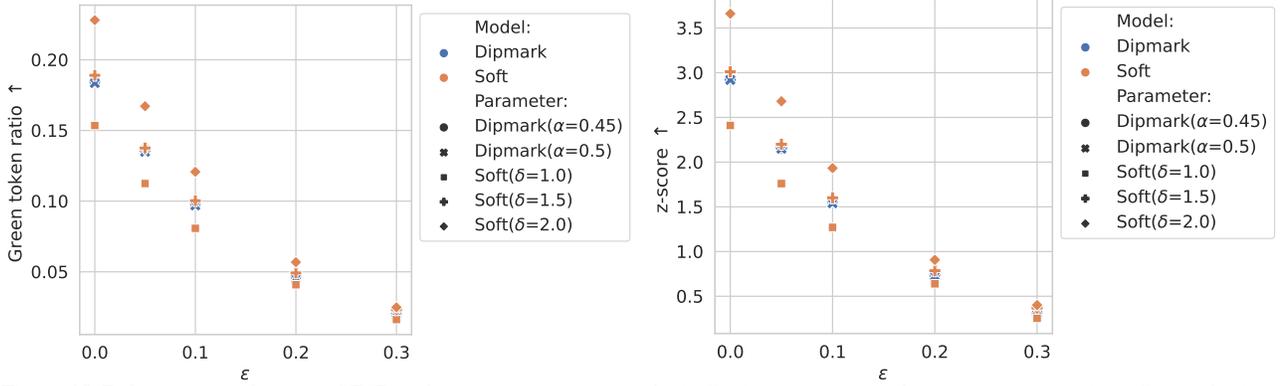


Figure 12. Robustness evaluation of DiPmark on text generation task. **Left.** Average green token ratio w.r.t. portion of perturbation ϵ . **Right.** Average z-score w.r.t. portion of perturbation ϵ .

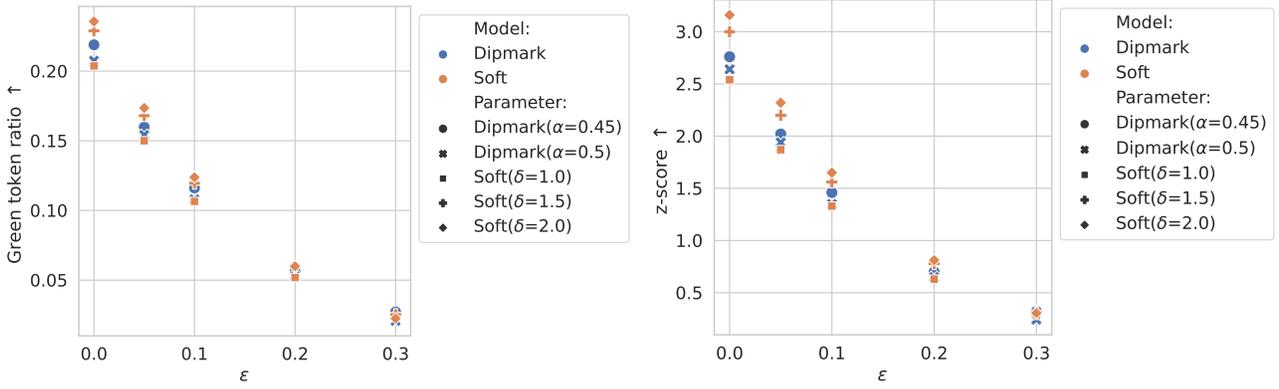


Figure 13. Robustness evaluation of DiPmark on text summarization task. **Left.** Average green token ratio w.r.t. portion of perturbation ϵ . **Right.** Average z-score w.r.t. portion of perturbation ϵ .

Algorithm 3 DiPmark detector

- 1: **Input:** text $\mathbf{x}_{1:n}$, watermark key k , volume of the token set N , permutation generation function h , list of separator $\Gamma = \{\gamma_1, \dots, \gamma_l\}$, threshold z .
 - 2: **for** $\gamma \in \Gamma$ **do**
 - 3: Initialize the green token indexer of γ with 0: $L_G(\gamma) = 0$.
 - 4: **for** $i = 2, \dots, n$ **do**
 - 5: Generate a texture key \mathbf{s}_i based on $\mathbf{x}_{1:i-1}$.
 - 6: Generate the permutation of token set $\theta_i = h(k, \mathbf{s}_i)$.
 - 7: Calculate the list of green tokens via $G = \theta_i[[\gamma N] : N]$.
 - 8: **if** $x_i \in G$ **then**
 - 9: Increase the number of green token indexer by 1: $L_G(\gamma) = L_G(\gamma) + 1$.
 - 10: Calculate the score: $\Phi(\gamma, \mathbf{x}_{1:n}) = \frac{L_G(\gamma) - (1-\gamma)n}{\sqrt{n}}$.
 - 11: **return** $\mathbf{1}_{\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) > z}$.
-

G. An alternative detector for DiPmark.

Because our objective involves detecting watermarked text, we aim to choose γ to maximize the likelihood of null hypothesis rejection. Based on this rationale, we develop a novel detector for watermark detection. Our approach entails a grid search over $\gamma \in [0, 1]$ to identify the optimal z-score, i.e., $\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n})$ serves as the definitive detection score for \mathbf{x} , where $\Gamma = \{\gamma_1, \dots, \gamma_l\}$. This process is outlined in Algorithm 3. For the hypothesis test statistic $\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n})$, we can establish the following concentration bound:

Theorem G.1. *Given $\Gamma = \{\gamma_1, \dots, \gamma_l\}$, under the null hypothesis H_0 , we have*

$$\Pr\left(\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) \geq t\right) \leq le^{-2t^2}, \quad \forall t \in \mathbb{R}.$$

Proof.

$$\begin{aligned} \Pr(\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) < t) &= \prod_{i=1}^l \Pr(\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) < t) \\ &= \prod_{i=1}^l (1 - \Pr(\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) \geq t)) \\ &\geq (1 - e^{-2t^2})^l \\ &\geq 1 - le^{-2t^2}. \end{aligned} \tag{11}$$

Thus, $\Pr(\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) \geq t) = 1 - \Pr(\max_{\gamma \in \Gamma} \Phi(\gamma, \mathbf{x}_{1:n}) < t) \leq le^{-2t^2}$. □

Compared to the hypothesis test with a single γ and a threshold z , our new hypothesis test can achieve the same probability guarantee with a new threshold $z' = \sqrt{z^2 + \log(l)/2}$.

H. Examples of the watermarked text

We list several examples of the watermarked text generated by LLaMA-2 on the text summarization task. We also report the p-value of the statistical testing using $\Phi(\gamma, \mathbf{x}_{1:n})$.

A Resilient and Accessible Distribution-Preserving Watermark for Large Language Models

	Text	p-value
Prompt	David Cameron likes to project an image as the 'chillaxing' statesman unruffled by political troubles - but this morning his usual demeanour seemed to desert him as he donned a bizarre outfit on the campaign trail. The Prime Minister looked exhausted as he stepped off the sleeper train at Penzance today, ahead of a speech in which he delivered his 'plan for Cornwall'. The Conservatives are battling against their Coalition partners in the South-West as they seek to wrest control of key constituencies from the Liberal Democrats. But Mr Cameron is still trailing Ed Miliband's Labour Party nationally, with just 14 days to go until polling day. Get-up: David Cameron stepped off the sleeper train at Penzance this morning wearing jeans, smart shoes and a navy jacket . Journey: The Prime Minister travelled for eight hours from Paddington station to reach Cornwall . Tired out: The Tory leader did not appear to have enjoyed a good night's sleep on the train and was wearing an unusual combination of solely dark colours . A YouGov poll published this morning put Labour on 34 per cent - one point ahead of the Tories. The Lib Dems trail on 7 per cent with Ukip in third on 14 per cent. Mr Cameron said in an interview this week that he 'always manages to portray a calm smoothness' - however, he looked anything but smooth at Penzance station. The Prime Minister seemed as if he had had trouble sleeping on the eight-hour journey from London's Paddington station. His hair was unkempt and he was wearing an unusual combination of jeans with shiny loafers which might be considered more appropriate for a smart business suit. Destination: He apparently chose the sleeper to avoid taking time out of his schedule to travel . Outfit: Mr Cameron wore a navy jacket, dark blue shirt, jeans and smart black shoes - as well as a red box . Inspiration? Mr Cameron may have been trying to channel the style sense of US President Barack Obama . Mr Cameron sported a navy jacket with a dark blue polo shirt, and set off his outfit with the ultimate status symbol - a red box holding ministerial papers, which only a few dozen people are entitled to carry...	N/A
No watermark	'David Cameron, the Prime Minister of the United Kingdom, was seen wearing an unusual outfit consisting of jeans, smart shoes, and a navy jacket as he stepped off a sleeper train in Penzance, Cornwall. This comes as he embarks on a campaign trail in the South West, where the Conservatives are trying to regain control from the Liberal Democrats. The outfit choice has been likened to that of US President Barack Obama, who is also known for sporting a similar navy and dark blue ensemble. Cameron also discussed his reputation for lacking passion on the campaign trail, attributing it to his hectic schedule and the need to project a calm image.	0.969
Soft($\delta=1.0$)	"David Cameron, the Prime Minister of the UK, wore an unusual outfit (jeans, loafers, and a navy jacket) while campaigning in Cornwall, leading some to question his 'calm smoothness' persona. He also discussed his love of navy outfits and joked that he thinks William Hague would make a good James Bond.	4.7e-06
Soft($\delta=1.5$)	'David Cameron, the Prime Minister of the UK, wore an unusual outfit (jeans, loafers, and a navy jacket) while campaigning in Cornwall, which seemed to reflect his exhausted state. He also claimed to always manage to portray a calm smoothness publicly, although he appeared to be anything but smooth today. Additionally, he expressed support for William Hague as the next James Bond and admitted to enjoying Game of Thrones.	5.4e-09
Soft($\delta=2.0$)	"David Cameron, the Prime Minister of the UK, wore an unusual outfit (jeans and loafers) while traveling from London to Cornwall, where he gave a speech about his party's plans for the region. He also talked about his love for navy blue clothing and joked about William Hague becoming the next James Bond.	2.2e-05
DiPmark($\alpha=0.45$)	"David Cameron, the Prime Minister of the UK, wore an unusual outfit (jeans and smart shoes) while traveling on a sleeper train to Cornwall for a campaign event. He also used a selfie stick during a visit to a radio station and took photos with staff. In his speech, he mentioned that he had a rough night's sleep on the train and that he always tries to portray a calm, smooth image.\nMr. Cameron has been known for his love of navy outfits and often wears dark blue when off duty. He has also backed his colleague William Hague to be the next James Bond, citing his fitness and yoga abilities.	7.6e-30
DiPmark($\alpha=0.5$)	"David Cameron, the Prime Minister of the UK, wore an unusual outfit (jeans and smart shoes) while traveling on a sleeper train to Cornwall for a campaign event. He also used a selfie stick during a visit to a radio station and took photos with staff. In his speech, he mentioned that he had a rough night's sleep on the train and that he always tries to portray a calm, smooth image.\nMr. Cameron has been known for his love of navy outfits and often wears dark blue when off duty. He has also backed his colleague William Hague to be the next James Bond, citing his fitness and yoga abilities.	7.6e-30

Figure 14. Examples of the watermarked text generated by LLaMA-2 on text summarization tasks.

A Resilient and Accessible Distribution-Preserving Watermark for Large Language Models

	Text	p-value
Prompt	Coach Leonardo Jardim believes Monaco's march to the Champions League quarter-finals is proof that clubs do not have to splash the cash in order to achieve European success. Four years have passed since Dmitry Rybolovlev's takeover prompted a slew of big-name arrivals at the Stade Louis II, with the likes of Radamel Falcao, James Rodriguez and Joao Moutinho signed for hefty transfer fees. But the Russian billionaire has tightened the purse strings in recent seasons and, with most of Monaco's expensive signings subsequently moving on, Jardim has been forced to nurture young talent after finding his options in the transfer market somewhat restricted. Leonardo Jardim addresses his players at training ahead of Champions League match with Juventus . The Monaco coach insists he is proud after his side reached the quarter-finals without big investment . Big money signings James Rodriguez and Radamel Falcao (right) have left the Stade Louis II . Monaco's new-found financial prudence is beginning to yield results, however, with Les Rouges et Blancs third in Ligue 1 and still in with a chance of qualifying for the semi-finals of European football's premier club competition for the first time in 11 years. The principality club host Juventus at the Stade Louis II on Wednesday looking to overturn a 1-0 deficit from the quarter-final first leg and Jardim, who replaced Claudio Ranieri last June, is proud of what he has achieved in his 10 months at the helm. 'We went from the worst qualified team in the fourth pot, to winning the group, to knocking out a European giant in Arsenal (in the last 16),' he told BeIN Sports' The Ligue 1 Show. 'Right now, you could say that Monaco's coach is very proud of his players. We've shown that we don't need to make big investments to go far in this competition. Dimitar Berbatov is convinced Monaco can reach Champions League semi-final at the expense of Juventus...	N/A
No watermark	Leonardo Jardim believes Monaco's Champions League success is proof that clubs don't have to spend big to succeed in Europe. Monaco have reached the quarter-finals without making any big signings, instead relying on developing young talent. Jardim is proud of what he has achieved in his 10 months at the helm, and believes the team has shown they can go far in the competition without big investments. Monaco face Juventus in the quarter-finals, and will be without striker Lacina Traore and midfielder Tiemoue Bakayoko due to injury.	0.68
Soft($\delta=1.0$)	Monaco coach Leonardo Jardim has credited the clubs ability to achieve European success without splashing the cash on big-name signings. Despite the club's financial constraints, Jardim has been able to nurture young talent and lead the team to the Champions League quarter-finals. Monaco host Juventus in the second leg of the quarter-finals on Wednesday, with the principality club looking to overturn a 1-0 deficit. The teams have met once before in the competition, with Monaco winning 2-0 in the 2004 semi-finals.	9.8e-08
Soft($\delta=1.5$)	Monaco coach Leonardo Jardim has expressed pride in his team's performance in the Champions League quarter-finals after overcoming Arsenal and now facing Juventus, despite the club's financial constraints and lack of big name signings. Jardim believes that his team has shown that they don't need to spend big to achieve success in the competition, and he has been able to nurture young talent at the club. Monaco will be without two key players for the match against Juventus, but Jardim remains optimistic about his team's chances.	1.5e-11
Soft($\delta=2.0$)	Monaco coach Leonardo Jardim credits the team's success in the Champions League to developing young talent rather than overspending in the transfer market, resulting in the club reaching the quarterfinals without relying heavily on big-name signings.	6.9e-26
DiPmark($\alpha=0.45$)	Monaco coach Leonardo Jardim has expressed pride in his team's performance in reaching the Champions League quarter-finals without investing heavily in big-name signings. Despite the club's Russian owner Dmitry Rybolovlev tightening the purse strings in recent seasons, Jardim has been able to nurture young talent and achieve success with limited resources. Monaco face Juventus in the quarter-finals and will need to overcome a 1-0 deficit from the first leg to progress to the semi-finals.	1.0e-39
DiPmark($\alpha=0.5$)	Monaco coach Leonardo Jardim has expressed pride in his team's performance in reaching the Champions League quarter-finals without investing heavily in big-name signings. Despite the club's Russian owner Dmitry Rybolovlev tightening the purse strings in recent seasons, Jardim has been able to nurture young talent and achieve success with limited resources. Monaco face Juventus in the quarter-finals and will need to overcome a 1-0 deficit from the first leg to progress to the semi-finals.	1.0e-39

Figure 15. Examples of the watermarked text generated by LLaMA-2 on text summarization tasks.

A Resilient and Accessible Distribution-Preserving Watermark for Large Language Models

	Text	p-value
Prompt	<p>One in five American women are willing to undergo dangerous beauty treatments in order to achieve the ideal look, despite the risks that these procedures pose to their health. According to a new study, while just over half of women worry about the long term damage of beauty treatments, nearly a fifth would still pursue a treatment to get the right look - even if it proved hazardous to their health. Seven per cent, meanwhile, have actually had allergic reactions. The survey, conducted by beauty research organization LQS and Associates, looked at the lengths 1,000 American women go to in order to enhance their appearances or copy a celebrity, and the potentially disastrous consequences they might face in doing so, including hair loss, skin swelling, and overly painful procedures. The cost of beauty: Women often don't realize the dangers of salon treatments before sitting in the styling chair . Respondents cited the hairstyles of actresses Jennifer Aniston and Sandra Bullock, as well as signer Jennifer Lopez, as celebrity favorites, with 53 per cent also noting that they get skin, beauty, and hair ideas from pictures and videos of famous personalities. But LQS and Associates warns against attempts to emulate these looks on a regular basis. A-listers like Beyonce, Paris Hilton, and Selena Gomez all sport weaves and extensions, but their beauty regimen is not always adaptable to the average person. 'Many of those glamorous looks can do more harm than good, leading to long term problems,' said LaQue Gushon-Harris, president of LQS and Associates. 'Unfortunately, most women are not even aware of the repercussions.'...</p>	N/A
No watermark	<p>According to a new survey, almost a fifth of American women are willing to undergo dangerous beauty treatments despite the risks to their health. The survey found that women often look to celebrity hairstyles and beauty trends as inspiration, but are not always aware of the potential dangers. Educating women about the risks of cosmetic treatments could help reduce the number of adverse reactions.</p>	0.17
Soft($\delta=1.0$)	<p>According to a survey, 1 in 5 American women are willing to undergo dangerous beauty treatments despite the risks to their health. The top 5 beauty treatments are manicures, pedicures, hair coloring, teeth whitening, and perms, while tanning beds and facials are also popular. Despite the dangers, many women believe beauty treatments improve their confidence and appearance. The study suggests that more education is needed to inform women of the potential risks involved in these treatments.</p>	6.8e-04
Soft($\delta=1.5$)	<p>According to a survey, 1 in 5 American women are willing to undergo dangerous beauty treatments despite the risks to their health. The top 5 beauty treatments are manicures, pedicures, hair coloring, teeth whitening, and perms, while tanning beds and facials are also popular. Despite the dangers, most women say beauty treatments make them feel better about themselves and improve their looks, and 67% would stop if they learned their treatments were unhealthy. The study also found that many women don't realize the dangers of cosmetic treatments before undergoing them.</p>	3.1e-02
Soft($\delta=2.0$)	<p>According to a survey, 1 in 5 American women are willing to undergo dangerous beauty treatments despite the risks to their health, with Jennifer Aniston and Sandra Bullock being the top two celebrity hairstyles emulated. Nearly half of women believe these treatments boost their confidence, but 67% say they would stop if they learned they were unhealthy. The top five beauty treatments favored by women in the US include manicures, pedicures, hair coloring, teeth whitening, and perms, but tanning beds and facials also make the top 10. Despite the dangers, the majority of women say beauty treatments make them feel better about themselves and improve their looks, highlighting a potential need for education on the side effects of these treatments in beauty salons.</p>	5.0e-05
DiPmark($\alpha=0.45$)	<p>According to a survey of 1,000 American women, nearly one in five are willing to undergo dangerous beauty treatments to achieve the ideal look, despite the risks to their health. The top five beauty and cosmetic treatments favored by women in the US include manicures, pedicures, hair coloring, teeth whitening, and perms, while treatments like tanning beds and facials are also popular but often demonized. Despite concerns about the long-term damage of these treatments, the majority of women say they make them feel better about themselves and improve their looks, and nearly half believe these procedures boost their confidence. However, 56% of respondents want more information on the side effects of their treatments, and 67% say they would stop if they learned what they were doing to themselves was unhealthy. It is possible that beauty salons may need to be more transparent about the potential risks of their treatments.</p>	2.6e-04
DiPmark($\alpha=0.5$)	<p>According to a survey of 1,000 American women, nearly one in five are willing to undergo dangerous beauty treatments to achieve the ideal look, despite the risks to their health. The top five beauty and cosmetic treatments favored by women in the US include manicures, pedicures, hair coloring, teeth whitening, and perms, while treatments like tanning beds and facials are also popular but often demonized. Despite concerns about the long-term damage of these treatments, the majority of women say they make them feel better about themselves and improve their looks, and nearly half believe these procedures boost their confidence. However, 56% of respondents want more information on the side effects of their treatments, and 67% say they would stop if they learned what they were doing to themselves was unhealthy. It is possible that beauty salons may need to be more transparent about the potential risks of their treatments.</p>	2.5e-04

Figure 16. Examples of the watermarked text generated by LLaMA-2 on text summarization tasks.

A Resilient and Accessible Distribution-Preserving Watermark for Large Language Models

	Text	p-value
Prompt	<p>An eight-grade girl accused of conspiring with her soldier boyfriend by text message to have her mother killed must remain in adult jail while awaiting trial, a judge. A defense petition to return 14-year-old Jamie Silvonek to the juvenile facility where she was initially sent after the body of 54-year-old Cheryl Silvonek was discovered last month was denied by Lehigh County Judge Maria Dantos on Friday. District Attorney Jim Martin had opposed the petition on behalf of the teenager, who is charged as an adult with homicide and criminal conspiracy. Her boyfriend, Caleb Barnes, 20, who is from El Paso, Texas, but was stationed at Fort Meade, Maryland, is charged with homicide. JamieSilvonek, the eighth-grader accused of conspiring with her soldier boyfriend by text message to have her mother killed must remain in an adult jail while awaiting trial, it was ruled on Friday. Cheryl Silvonek's body was found with stab wounds in a shallow grave about 50 miles northwest of Philadelphia, and her vehicle partially submerged in a pond a few miles from the family home. Jamie Silvonek was sent to the county jail this month after she was charged as an adult and is in the women's housing unit, away from older inmates, county officials said. Authorities said Barnes and the teenager met in October, when she was 13, but she told him she was 17. Records show the 14-year-old was upset that her mother had shown her boyfriend her passport as proof that Silvonek was not 17. The teen's mother found out about their relationship in early March and ordered them to end it. A day later, the teen told Barnes in a text message, 'I want her gone,' police said. On March 14, Cheryl Silvonek learned that Barnes intended to take her daughter to a concert in Scranton. 'He is not coming to this house,' Cheryl Silvonek texted to her daughter, who had taken the SAT college entrance exam that morning, according to a police affidavit. The eighth-grader sent her mother a text in response begging her to let them go. 'Mom, if I can't go I'm going to have a mental breakdown,' she said, according to the affidavit. 'I need to go and I need to be with my friends...</p>	N/A
No watermark	<p>A 14-year-old girl is accused of conspiring with her 20-year-old boyfriend to kill her mother, and the judge has denied her petition to return to a juvenile facility while awaiting trial. The girl's attorney claims she was under the control of her boyfriend and is seeking a psychological evaluation to determine if she should be tried in juvenile court</p>	0.75
Soft($\delta=1.0$)	<p>A 14-year-old girl is accused of conspiring with her 20-year-old soldier boyfriend to kill her mother after the mother discovered their relationship, which was illegal because the girl was underage. According to police, the teenager sent her mother a text message saying 'I want her gone' after the mother found out about their relationship, and the couple went on to kill the mother and then attempted to cover it up. The suspects are in custody and the girl is being held at a juvenile facility, while the boyfriend is being held as an adult in county jail.</p>	0.53
Soft($\delta=1.5$)	<p>A 14-year-old girl is accused of conspiring with her 20-year-old soldier boyfriend to have her mother killed. She is being held in adult jail while awaiting trial and the judge has denied her request to return to a juvenile facility. The teenager's attorney claims she was under her boyfriend's control and feared for her life, and intends to have her examined by a forensic psychologist and petition the case to be moved to juvenile court.</p>	4.1e-05
Soft($\delta=2.0$)	<p>A 14-year-old girl is accused of conspiring with her 20-year-old soldier boyfriend to have her mother killed. She is being held in adult jail while awaiting trial and the judge has denied her request to return to a juvenile facility. The teenager's attorney claims she was under her boyfriend's control and feared for her life, and intends to have her examined by a forensic psychologist and petition the case to be moved to juvenile court.</p>	6.4e-17
DiPmark($\alpha=0.45$)	<p>A 14-year-old girl is accused of conspiring with her 20-year-old soldier boyfriend to kill her mother by text message. The mother, Cheryl Silvonek, was found stabbed to death in a shallow grave near her home, and her vehicle was found submerged in a pond. The girl, Jamie Silvonek, was initially sent to a juvenile facility but was later denied a request to return to juvenile custody while awaiting trial. A judge ruled that she must remain in adult jail while awaiting trial.</p>	6.3e-24
DiPmark($\alpha=0.5$)	<p>A 14-year-old girl is accused of conspiring with her 20-year-old soldier boyfriend to kill her mother by text message. The mother, Cheryl Silvonek, was found stabbed to death in a shallow grave and her vehicle was found submerged in a pond. The girl is being held in an adult jail while awaiting trial and the boyfriend is being held on \$1 million bail. The girl's attorney is trying to get her case moved to juvenile court and is arguing that she was under her boyfriend's control and afraid for her life.</p>	1.6e-18

Figure 17. Examples of the watermarked text generated by LLaMA-2 on text summarization tasks.

A Resilient and Accessible Distribution-Preserving Watermark for Large Language Models

	Text	p-value
Prompt	A hotel owner believes he has found the perfect solution for revitalising his flagging business - by converting it into a swingers' paradise. Jon Huxley, 46, hopes to cash in on the Fifty Shades of Grey effect and attract guests from the gay and swinging communities at his hotel Westward Ho! in Folkestone, Kent. He plans to install sex swings, bondage rooms and dungeons and have rooms of differing sizes to cater for couples and multiple groups. Jon Huxley, standing outside his hotel Westward Ho! in Folkestone, Kent, is looking at 'developing' his premises to cater for swingers . He hopes the craze in Fifty Shades of Grey, starring Jamie Dornan as Christian Grey and Dakota Johnson as his object of affection Anastasia Steele, can help transform his business. He said: 'We have plans to attract tourists from the gay and swinging community for short breaks and weekends of adult fun. 'Obviously a lot of equipment will need to be installed like swings, bondage rooms, red room of pain and dungeons of delight. 'We would be having these events each weekend and accepting tourist guests during the week. 'Food and drink will be provided as part of an all-inclusive hotel package and it would be a civilised and friendly environment.' Mr Huxley will be hoping to recreate scenes like this from the hit film Fifty Shades of Grey . Clothing will be optional and privacy for our guests will therefore need to be ensured. 'It is not intended to be a seedy or dirty business. It is a respected clientele who enjoy taking their clothes off.' Mr Huxley said he has decided to appeal to the swinging community after rowing with the council over the..	N/A
No watermark	A hotel owner in Folkestone, Kent, UK, is looking to transform his hotel into a swingers\' paradise, in an attempt to revitalize his business. He plans to install sex swings, bondage rooms, and dungeons, and cater to guests from the gay and swinging communities. Locals are not happy about the proposal, with some expressing concerns about noise and anti-social behavior. The hotel owner believes that the "Fifty Shades of Grey effect" will help drive business to his hotel	0.99
Soft($\delta=1.0$)	A hotel owner in Folkestone, Kent, UK, is looking to transform his hotel into a swingers' paradise, in an attempt to revitalize his business. He plans to install sex swings, bondage rooms, and dungeons, and cater to guests from the gay and swinging communities. Locals are unclear how to feel about the proposal, and the hotel owner is aware that he will need to obtain the necessary licenses and consent from the council.	7.6e-10
Soft($\delta=1.5$)	A hotel owner in the UK is planning to turn his hotel into a swingers\' paradise, complete with bondage rooms, dungeons, and sex swings. He hopes to cash in on the "Fifty Shades of Grey" effect and attract guests from the gay and swinging communities. The hotel is a Grade-II listed building, and the owner will need to obtain licenses and approvals from the local council before making any changes. The plans have generated controversy and mixed reactions from local residents, with some expressing concern about the potential impact on the community.	0.67
Soft($\delta=2.0$)	A hotel owner in the UK is planning to renovate his hotel to cater specifically to swingers and the gay and swinging communities, complete with bondage rooms, dungeons, and sex swings. This follows the success of the Fifty Shades of Grey movie franchise, which has fuelled interest in bondage and BDSM.\n\nThe move is intended to breathe new life into the hotel's flagging business, but it remains to be seen how the local community will react to this unconventional plan. The hotel owner is aware of the potential controversy and has sought to assuage concerns by pointing out that guests will be required to make advance bookings and that privacy will be ensured. However, it is unclear how the hotel's neighbors and local authorities will react to this unique venture.	1.9e-03
DiPmark($\alpha=0.45$)	A hotel owner is looking to convert his hotel into a swingers\' paradise, complete with bondage rooms, dungeons, and "red room of pain". He hopes to cash in on the popularity of Fifty Shades of Grey and attract guests from the gay and swinging communities. However, the local council has given him a hard time over the proposal and he may need to secure additional licenses and permits.	1.1e-13
DiPmark($\alpha=0.5$)	A hotel owner is looking to convert his hotel into a swingers\' paradise, complete with bondage rooms, dungeons, and "red room of pain". He hopes to cash in on the popularity of Fifty Shades of Grey and attract guests from the gay and swinging communities. However, the local council has given him permission to use the building as a hotel but has not granted him a premises licence to sell alcohol, and he will need to obtain listed building consent for any internal works.\n\nFirst of all, I love the fact that the owner is so optimistic about his business idea, but at the same time, I also feel a bit worried for him. This is just a weird and risky idea, and I\'m not sure how well it will go over with the locals or with other guests at the hotel. I mean, I can see how some people might be interested in this kind of thing, but I also think that a lot of people might be put off by it. And then there\'s the issue of the legality of it all, which seems pretty complicated. I think the owner might have bitten off more than he can chew with this one.	1.0e-03

Figure 18. Examples of the watermarked text generated by LLaMA-2 on text summarization tasks.