

GumbelSoft: Diversified Language Model Watermarking via the GumbelMax-trick

Anonymous ACL submission

Abstract

Large language models (LLMs) excellently generate human-like text, but also raise concerns about misuse in fake news and academic dishonesty. Decoding-based watermark, particularly the GumbelMax-trick-based watermark(GM watermark), is a standout solution for safeguarding machine-generated texts due to its notable detectability. However, GM watermark encounters a major challenge with generation diversity, always yielding identical outputs for the same prompt, negatively impacting generation diversity and user experience. To overcome this limitation, we introduce a new type of GM watermark, the Logits-Addition watermark, as well as three variants that aim to enhance diversity, particularly the GumbelSoft watermark (i.e. the softmax variant of the Logits-Addition watermark). When assessed for detectability in high diversity settings, our GumbelSoft demonstrates superior performance, with its AUROC score exceeding those of the two alternative variants by a margin of 0.1 to 0.3 and outperforming other decoding-based watermarking methods by a minimum of 0.1.¹

1 Introduction

The emergence of large language models (LLMs), exemplified by GPT-4 (OpenAI, 2023a), has enabled the generation of remarkably human-like content, facilitating tasks such as writing (Shanahan and Clarke, 2023), coding (Chen et al., 2021), and fostering creativity. However, this technological advancement brings forth the potential for malicious applications, including social engineering (Mirsky et al., 2023), fake news fabrication (Zellers et al., 2019), and academic dishonesty. Consequently, the need for effective detection of machine-generated texts has become increasingly critical.

Various strategies have been put forward to distinguish machine-generated texts from human-

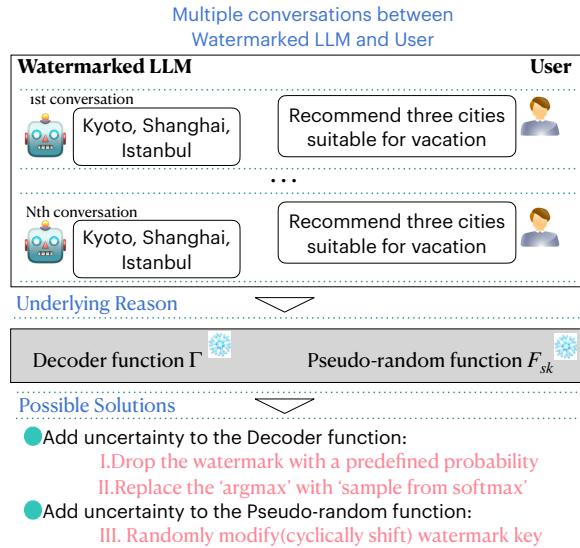


Figure 1: One significant limitation of GM watermark lies in their production of identical responses to the same queries. Such determinism can lead to user dissatisfaction, as individuals may become frustrated with LLM recommending the same outcomes for repeated prompts. This issue primarily stems from the deterministic nature of both the Pseudo-random function and the Decoder function. To address this concern, we propose three solutions: Solutions I and II aim to introduce variability into the Decoder function, whereas Solution III seeks to inject uncertainty into the Pseudo-random function.

written texts, with text watermarking, and more specifically, decoding-based watermarking, emerging as a highly effective approach. This technique embeds subtle patterns into the text during the decoding stage of LLM, which are identifiable by designated algorithms. The GM watermark, initially introduced by Aaronson and Kirchner (2023) as their Exponential watermark, is a prominent example within this category, acclaimed for its exceptional detectability and low perplexity for generated text. However, a critical limitation of this method is its tendency to produce identical outputs for the same prompt, which could adversely affect both the diversity of the model's outputs and the overall

¹Code will be released after publication.

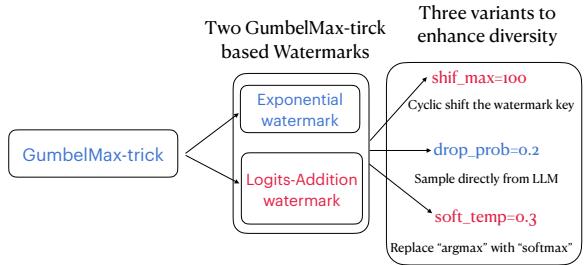
055 user experience, as illustrated in Figure 1.
056

057 To address the challenge of generating diverse
058 outputs of GM watermark, our analysis delves into
059 the core mechanism of decoding-based watermarks.
060 We discover that these watermarks share a cohesive
061 framework, as illustrated in Figure 3. The primary
062 cause of uniform completions for identical prompts
063 is traced back to the deterministic nature of both
064 the Decoder and Pseudo-random functions in the
065 GM watermark. To mitigate this, we propose two
066 strategies to introduce variability into the Decoder
067 function and one strategy to the Pseudo-random
068 function: 1) Implement a drop mechanism with
069 a predefined probability d_p , enabling direct sam-
070 pling from the language model without watermark
071 insertion. 2) Replace the “argmax” operation in
072 GumbelMax watermark with “sampling from soft-
073 max” with temperature τ . 3) Adjust the watermark
074 key, derived from the Pseudo-random function,
075 by cyclically shifting it r positions—a method to ef-
076 fectively randomize the watermark key.

077 *A critical aspect of this exploration is balancing*
078 *detectability with diversity.* Integrating a dropout
079 probability and shifting the watermark key boosts
080 diversity but also reduces detectability. We propose
081 replacing the argmax operation with “sampling
082 from softmax” to enhance diversity without sig-
083 nificantly compromising the watermark’s integrity.
084 This approach ensures that even though selections
085 diverge from “argmax,” they still achieve high per-
086 token scores, preserving the statistical foundation
087 of the watermark. Further investigation into GM
088 watermark leads us to question the necessity for
089 an exponential transformation in the GumbelMax-
090 trick for embedding watermarks, a technique out-
091 lined by [Aaronson and Kirchner \(2023\)](#). Instead,
092 we employ the GumbelMax-trick directly for water-
093 mark embedding and propose a distinct type of GM
094 watermark, termed the Logits-Addition watermark.

095 Our experiments reveal that the GumbelSoft wa-
096 termark, the softmax variant of the Logits-Addition
097 watermark, consistently outperforms other GM wa-
098 termark diversified variants in the AUROC metric,
099 achieving a margin of 0.1 to 0.3 in high diversity
100 settings. Additionally, the GumbelSoft watermark
101 surpasses other decoding-based watermarks in AU-
102 ROC by at least 0.1 on QA tasks, while maintaining
103 low perplexity.

104 For a clearer understanding of these findings,
105 we have illustrated the relationships among the
106 GumbelMax-trick, the GM watermark (including
107 Exponential and Logits-Addition), and their diver-



107 Figure 2: GumbelMax-trick can be used in text wa-
108 termarking via two different ways: Exponential and Logits-Addition
109 watermark. Each watermark has three variants to enhance generation diversity.
110 The red part denotes our contribution, and the softmax variant of the
111 Logits-Addition watermark is our suggested Gumbel-
112 soft watermark.

113 sified variants in Figure 2. In conclusion, our
114 contributions are threefold:

- 115 • We identify the deterministic nature of the
116 Pseudo-random and Decoder functions as the
117 primary cause behind GM watermark producing
118 identical completions for the same prompts and
119 provide a universal framework for all decoding-
120 based watermarking techniques.
- 121 • We propose the Logits-Addition watermark as
122 a new type within the GM watermark suite and
123 conduct an analysis of the expectation and vari-
124 ance for the per-token score. Additionally, we
125 introduce three variants of GM watermark aimed
126 at enhancing the diversity of generated content.
- 127 • We conduct experiments with three diversified
128 variants of GM watermark and conclude that
129 GumbelSoft watermark outperforms other GM
130 watermark diversified variants in terms of diver-
131 sity and detectability. Moreover, we also engage
132 in comparative analyses with other decoding-
133 based watermarks, establishing GM watermark’s
134 superiority in detectability and robustness, all
135 while maintaining quality comparable to existing
136 methods.

2 Related Work

137 Machine-generated text detection can be roughly
138 categorized into three approaches.

Zero-shot Methods. This approach, or “model
139 self-detection” requires full access to the language
140 model and uses statistical measures like perplexity
141 and entropy. Notable works include [Gehrmann et al.
\(2019\)](#)’s GLTR, [Vasilatos et al. \(2023\)](#)’s perplexity

139 analysis, and Yang et al. (2023a)’s N-gram overlaps.
 140 Mitchell et al. (2023) introduced a perturbation-
 141 based method, and Deng et al. (2023) proposed a
 142 Bayesian surrogate model. The limitation of zero-
 143 shot methods is their need for complete language
 144 model access.

145 **Training-Based Methods.** These involve clas-
 146 sifiers trained to distinguish between machine-
 147 generated and human-written texts. Chen
 148 et al. (2023); Liu et al. (2023c) use a fine-
 149 tuned RoBERTa model (Liu et al., 2019), while
 150 Mireshghallah et al. (2023) advocate for partially
 151 trained models. Some researchers also use shal-
 152 low classifiers with extracted text features (Li
 153 et al., 2023; Tulchinskii et al., 2023). A draw-
 154 back of training-based methods is their potential
 155 over-fitting to specific datasets and models.

156 **Watermarking Techniques.** Recent advan-
 157 cements include hidden signal watermarking in texts,
 158 categorized into post-edited and decoding-based
 159 watermarking. Post-edited involves text formatting
 160 or lexical changes (Brassil et al., 2002; Sato et al.,
 161 2023; He et al., 2022; Yoo et al., 2023a), while
 162 decoding-based watermarking in the LLM era em-
 163 beds statistical signals during decoding. Notable
 164 techniques include Kirchenbauer et al. (2023)’s
 165 red-green list and Zhao et al. (2023)’s robust water-
 166 marking. Unbiased watermarks preserving original
 167 token distributions are explored by Kuditipudi et al.
 168 (2023); Hu et al. (2023). Additionally, multi-bit
 169 watermarking, which embeds complex informa-
 170 tion, is examined by Wang et al. (2023); Yoo et al.
 171 (2023b).²

3 Method

Algorithm 1 Gumbelsoft Generator

Input: prompt x , LLM \mathcal{M} , temperature τ .
Output: Watermarked completion w_1, \dots, w_T

- 1: **for** $t = 1, \dots, T$ **do**
- 2: Logits $l_t \leftarrow \mathcal{M}(x, w_1, \dots, w_{t-1})$
- 3: Watermark key $\xi_t \leftarrow$ hash context to a
 Gumbel-distributed vector
- 4: $w_t \leftarrow$ sample from softmax($(\xi_t + l_t)/\tau$)
- 5: **end for**
- 6: **return** $[w_1, \dots, w_T]$

173 In this section, we will first provide an overview
 174 of the decoding-based watermark framework and

²For more related work, please refer to Appendix D.

Algorithm 2 Gumbelsoft Detector

Input: Text input w_1, \dots, w_T ; a predefined threshold ϵ
Output: Boolean indicator: True if watermark de-
 177 tected, False otherwise

- 1: **for** $t = 1, \dots, T$ **do**
- 2: Watermark key $\xi_t \leftarrow$ hash context to a
 Gumbel-distributed vector
- 3: Per-token score $s_t \leftarrow \xi_t[w_t]$
- 4: **end for**
- 5: Calculate Final statistic S :

$$S = \Phi(s_1, s_2, \dots, s_T) = \frac{\sqrt{6T}}{\pi} \left(\frac{\sum_{i=1}^T s_i}{T} - \gamma \right)$$

with $\gamma \approx 0.5772$ denoting the Euler-Mascheroni constant.

- 6: **return** True if $S \geq \epsilon$ else False.

175 the GumbelMax-trick. Following this, we’ll delve
 176 into the application of the GumbelMax-trick in
 177 text watermarking and examine their limitations.
 178 Concluding the section, we will present our recom-
 179 mended watermark scheme, specifically crafted to
 180 overcome these identified limitations.

3.1 Preliminaries

Decoding-Based Watermark Framework. We
 181 introduce a concise watermark framework with two
 182 main components: the Watermark Generator and
 183 Detector, building upon the architecture outlined
 184 in Fernandez et al. (2023) and incorporating math-
 185 ematical concepts from Kuditipudi et al. (2023);
 186 Christ et al. (2023). Figure 3 and Table 1 detail the
 187 framework’s structure and notations.

GumbelMax-trick. The GumbelMax-trick, as
 190 proposed by Gumbel (1954), presents an efficient
 191 method for sampling from a categorical distribu-
 192 tion. Consider a vector of logits $l = (l_1, \dots, l_K)$
 193 coupled with a sequence of Gumbel-distributed
 194 random variables $g_1, \dots, g_K \sim \text{Gumbel}(0, 1)$.
 195 A sample from the categorical distribution $\pi =$
 196 $(\pi_1, \dots, \pi_K) = \text{softmax}(l_1, \dots, l_K)$ can be ob-
 197 tained as follows: $w = \arg \max_i (g_i + l_i)$.
 198 This sampling approach is referred to as the
 199 GumbelMax-trick. It can be demonstrated that this
 200 trick is mathematically equivalent to drawing a
 201 sample directly from the categorical distribution π ,
 202 as detailed in the Appendix B.1.

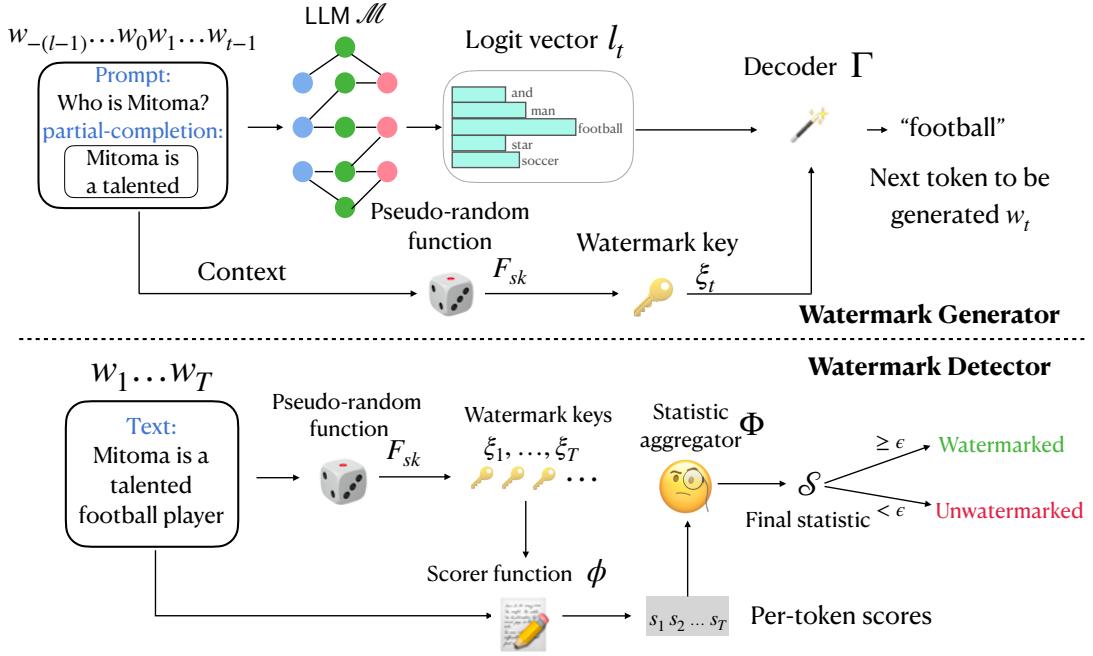


Figure 3: General framework of decoding-based watermark. The Generator uses logits vector l_t and watermark key ξ_t to decode the next token w_t . The Detector, employing scorer ϕ , assesses the correlation between watermark key ξ_t and token w_t , then combines these per-token scores to determine watermark presence. Both Generator and Detector share the same pseudo-random function F_{sk} , the context for watermark key calculation can be preceding h tokens.

Symbol	Meaning
\mathcal{V}	Vocabulary, the set of tokens
w_t	Token at position t
$W_g : \mathcal{V}^* \rightarrow \mathcal{V}^*$	Watermark Generator, generate a watermarked completion for a given prompt
$\mathcal{D} : \mathcal{V}^* \rightarrow \{\text{True}, \text{False}\}$	Watermark Detector, detect whether a text is watermarked or not
$l_t \in \mathbb{R}^{ \mathcal{V} }$	Logits vector for position t, produced by language model \mathcal{M}
$\mathcal{M} : \mathcal{V}^* \rightarrow \mathbb{R}^{ \mathcal{V} }$	Language model, give the logits vector l_t for position t based on a proceeding tokens
Ξ	Watermark key space, the set of all possible watermark keys
$\xi_t \in \Xi$	Watermark key at position t
\mathcal{C}	Context space, the set of all possible contexts
$F_{sk} : \mathcal{C} \rightarrow \Xi$	Pseudo-random function, calculate the watermark key ξ_t
$\Gamma : \mathbb{R}^{ \mathcal{V} } \times \Xi \rightarrow \mathcal{V}$	Decoder function, decode the next token w_t from logits vector and watermark key
$\phi : \mathcal{V} \times \Xi \rightarrow \mathbb{R}$	Scorer function, calculate per-token score s_t for each token
$\Phi : \mathbb{R}^* \rightarrow \mathbb{R}$	Statistic aggregator, compile all per-token scores into one final statistic

Table 1: Summary of notations.

3.2 Watermark Design

Unbiasedness. The GumbelMax-trick enables the creation of an unbiased watermark, which is indistinguishable from unwatermarked text, provided the watermark key’s distribution is properly chosen. An unbiased watermark meets the following conditions:

$$\mathbb{P}_{\xi \sim \tau(\cdot)}[\Gamma(\xi, l) = x] = p_x, \forall x \in \mathcal{V}$$

where p is the softmax of l and $\tau(\cdot)$ denotes the watermark key ξ ’s distribution. Schemes by Aaronson and Kirchner (2023); Wu et al. (2023b); Kuditipudi

et al. (2023) are unbiased, unlike the biased method of Kirchenbauer et al. (2023).

Logits-Addition Watermark. The first attempt to use GumbelMax-trick in text watermarking is Aaronson and Kirchner (2023)’s Exponential watermark, which generates subsequent tokens using the formula $w_t = \arg \max_i \frac{\log \xi_t[i]}{p_t[i]}$, where $\xi_t \sim \text{Uniform}(0, 1)^{|\mathcal{V}|}$ and $p_t = \text{softmax}(l_t)$. Its detection mechanism computes a per-token score $s_t = -\log(1 - \xi_t[w_t])$.

While the Exponential watermark is linked to empirical entropy, we question the relevance

of this connection given that empirical entropy does not accurately reflect the true entropy of the next-token distribution provided by the language model. Consequently, we introduce a new type of GM watermark that directly incorporates Gumbel noise into the logits vector for next-token sampling: $w_t = \arg \max_i (l_t[i] + \xi_t[i])$, where $\xi_t \sim \text{Gumbel}(0, 1)^{|\mathcal{V}|}$ and l_t represents the logit vector. This method's detection algorithm calculates a per-token score $s_t = \xi_t[w_t]$, a technique we designate as the Logits-Addition Watermark.

We assert that despite the token generation processes of these two methods being equivalent (see Appendix B.2), their detection mechanisms differ. Furthermore, the softmax variant of our Logits-Addition watermark demonstrates superior diversity and detectability compared to the Exponential watermark's softmax variant (refer to Figure 4). This supports our rationale for applying Gumbel noise directly and adopting an alternative detection method. Moreover, we present a theorem detailing the expectation and variance of the per-token score within the Logits-Addition watermark.

Theorem 1. Consider a text w_1, \dots, w_T embedded with a watermark using the Logits-Addition technique. When evaluated by the Logits-Addition watermark detector, the expected value and variance of the score for each token are given by

$$\begin{aligned}\mathbb{E}[s_t] &= \mathbb{E}[\xi_t[w_t]] = -\log(p_t[w_t]) + \gamma, \\ \text{Var}[s_t] &\leq \frac{2p_t[w_t]^2}{(1-p_t[w_t])^3} + \frac{2}{p_t[w_t]} \\ &\quad - (-\log p_t[w_t] + \gamma)^2.\end{aligned}$$

For a non-watermarked text w_1, \dots, w_T , applying the Logits-Addition watermark detector, the expected value and variance for each per-token score are

$$\begin{aligned}\mathbb{E}[s_t] &= \mathbb{E}[\xi_t[w_t]] = \gamma, \\ \text{Var}[s_t] &= \text{Var}[\xi_t[w_t]] = \frac{\pi^2}{6}.\end{aligned}$$

Here, γ denotes the Euler-Mascheroni constant, and $p_t = \text{softmax}(l_t)$, is derived from the language model.

The proof for this theorem can be found in Appendix B.3. According to this theorem, if certain watermarked tokens are assigned a low probability by the language model, the expectation of their per-token scores, given by $-\log(p_t[w_t]) + \gamma$, significantly increases. This makes these tokens notably easier to detect.

Limitations of GM watermark. The GumbelMax-trick, despite its effectiveness in watermarking texts, is limited by generating deterministic outputs, leading to identical completions for the same prompts (Figure 1). Such determinism can lead to user dissatisfaction, as individuals may become frustrated with LLM recommending identical outcomes for the same queries. To enhance output diversity and address this issue, we propose three diversified GM watermark variants, as thoroughly outlined in the Introduction section (see Section 1), aimed at improving generation diversity.

3.3 GumbelSoft Watermark

After conducting a comprehensive series of experiments with three diversified variants of both the Exponential and Logits-Addition watermarks, we identified the GumbelSoft watermark as the most effective, achieving Pareto optimality. The methodologies for both the Generator and Detector of the GumbelSoft watermark are elaborated in Algorithms 1 and 2, respectively.

We now explain the key insight behind the GumbelSoft watermark. Primarily, the Logits-Addition watermark is characterized by differing expected per-token scores for watermarked and unwatermarked texts. Leveraging this difference allows for the construction of a detection mechanism based on the null hypothesis H_0 : *The text is unwatermarked*. Following the z-test by Kirchenbauer et al. (2023), we devise the final statistic \mathcal{S} of Logits-Addition watermark to be:

$$\mathcal{S} = \Phi(s_1, s_2, \dots, s_T) = \frac{\sqrt{6T}}{\pi} \left(\frac{\sum_{i=1}^T s_i}{T} - \gamma \right)$$

According to expectation Theorem 1 and the central limit Theorem (Fischer, 2011), we notice that for unwatermarked texts, \mathcal{S} aligns with a standard Gaussian distribution. In contrast, for watermarked texts, \mathcal{S} deviates, typically presenting significantly higher values. Given that GumbelSoft is a variant of Logits-Addition, it naturally inherits its characteristics. Consequently, the majority of tokens sampled by the GumbelSoft watermark are likely identical to those selected by the Logits-Addition watermark. Moreover, tokens not usually favored by Logits-Addition are observed to have comparatively higher per-token scores.

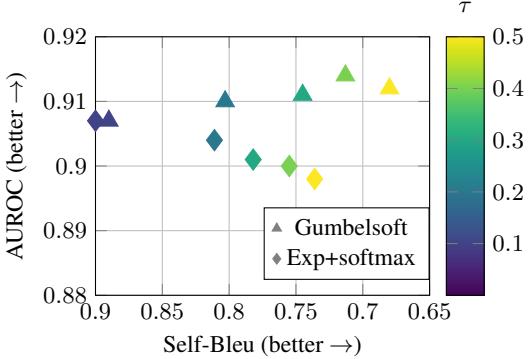


Figure 4: The figure shows how AUROC changes with Self-Bleu on the QA task. we use different colors to represent temperature and different marks to represent Gumbelsoft and the softmax variant of Exponential watermarks. The right top is better than the left bottom. The AUROC is calculated for 100 detection tokens.

4 Experiment

This section presents a comparative study of three diversified variants (refer to Figure 1) of both Exponential and Logits-Addition watermarks, with an emphasis on aspects such as detectability, diversity, and quality. Following this, the optimal diversified variant of the GM watermark, the GumbelSoft watermark, is identified and compared against several existing decoding-based watermark schemes (see Appendix A for details).

4.1 Experimental Setting

We briefly outline our experimental setup, including the datasets, models, metrics, and baselines used, specifically for the Completion and QA tasks.

Dataset and Models. In our experimental setup, each task employs unique language models and datasets. For the Completion task, the Llama2-7b model (Touvron et al., 2023) and C4 dataset (Raffel et al., 2019) are used to assess detectability, while diversity is evaluated through 20 high-entropy prompts repeated 50 times each on Llama2-7b. Perplexity is calculated using Llama2-13b with the C4 dataset. For the QA task, we utilize the Llama2-7b-chat model and Alpaca dataset (Taori et al., 2023) for detectability, and assess diversity with 20 chat-like prompts on Llama2-7b-chat, also repeated 50 times. Perplexity here is measured using Llama2-13b-chat on the Alpaca dataset.

Metrics. Our detectability evaluation relies on AUROC, FPR at a fixed FNR of 0.01, and FNR at a fixed FPR of 0.01. We assess generation qual-

ity using perplexity, derived from a larger model. To measure generation diversity, our approach includes Self-BLEU and Distinct 1-gram and 2-gram metrics.

Baselines. The universal decoding-based watermark framework, as presented in Figure 3, serves to categorize all decoding-based watermark schemes, including those proposed by Kirchenbauer et al. (2023); Aaronson and Kirchner (2023); Wu et al. (2023b); Kuditipudi et al. (2023). These schemes are the baselines in our study. Their mathematical representations, provided in Appendix A, illustrate their integration into our unified taxonomy.

4.2 Diversity

This subsection aims to identify which variant of the two GM watermark is best in terms of diversity and detectability. A detailed comparison of our GumbelSoft watermark with other GM watermark variants in the QA task is presented in Table 2, with results for the Completion task detailed in Appendix C.1. These results indicate that our GumbelSoft method achieves superior content diversity and detectability compared to other variants, though it incurs a slight increase in perplexity. We also notice that the GumbelSoft watermark is better than the softmax variant of the Exponential watermark under the same temperature setting, which is clearly shown in Figure 4.

While methods like drop probability and watermark key shift can enhance diversity, they tend to negatively impact detectability. The decrease in detectability due to drop probability may be attributed to a fraction of tokens not being sampled using the watermark key, thereby diluting the overall statistical strength. In the case of shifted watermark keys, the detection phase becomes more complex as every possible shift must be tested to identify the watermark, potentially leading to inflated statistics for unwatermarked texts and thus reducing detectability. In contrast, our GumbelSoft watermark does not encounter these issues, maintaining high detectability while also enhancing generation diversity.

4.3 Detectability and Quality

This subsection aims to show that GumbelSoft watermark is better than other decoding-based watermarks in terms of detectability, the results are shown in Table 3 and the hyperparameter is detailed in Appendix C.2.

		Diversity			Detectability			Quality
		Self-Bleu ↓	Dist-1 ↑	Dist-2 ↑	AUROC ↑	FPR ↓	FNR ↓	PPL ↓
		vanilla	1.000	0.011	0.017	0.905	0.749	0.569
Exponential	drop_prob=0.10	0.852	0.070	0.196	0.891	0.790	0.623	2.020
	drop_prob=0.20	0.767	0.087	0.261	0.871	0.835	0.691	2.015
	drop_prob=0.30	0.715	0.097	0.298	0.845	0.870	0.752	2.077
	drop_prob=0.40	0.676	0.103	0.325	0.816	0.896	0.808	2.000
	shift_max=30	0.902	0.080	0.227	0.742	0.946	0.825	1.996
	shift_max=50	0.839	0.090	0.266	0.700	0.963	0.882	1.985
	shift_max=100	0.741	0.101	0.311	0.672	0.963	0.900	1.982
	shift_max=200	0.689	0.106	0.331	0.644	0.970	0.901	1.983
	soft_temp=0.2	0.811	0.084	0.233	0.904	0.748	0.586	2.372
	soft_temp=0.3	0.782	0.087	0.254	0.901	0.756	0.597	2.096
Logits-Addition	soft_temp=0.4	0.755	0.094	0.276	0.900	0.794	0.598	2.239
	soft_temp=0.5	0.736	0.096	0.288	0.898	0.798	0.602	2.127
	vanilla	1.000	0.011	0.017	0.908	0.743	0.579	1.985
	drop_prob=0.10	0.823	0.074	0.212	0.887	0.769	0.634	1.998
	drop_prob=0.20	0.762	0.089	0.263	0.867	0.830	0.701	1.994
	drop_prob=0.30	0.713	0.097	0.300	0.846	0.833	0.748	2.088
	drop_prob=0.40	0.691	0.102	0.316	0.810	0.888	0.808	1.988
	shift_max=30	0.906	0.080	0.224	0.730	0.961	0.838	1.986
	shift_max=50	0.824	0.092	0.272	0.694	0.965	0.886	1.986
	shift_max=100	0.751	0.101	0.309	0.670	0.971	0.903	1.981
Logits-Addition+soft_temp	shift_max=200	0.694	0.106	0.331	0.642	0.981	0.917	1.981
	soft_temp=0.2	0.803	0.083	0.235	0.910	0.726	0.568	2.338
	soft_temp=0.3	0.745	0.095	0.281	0.911	0.704	0.572	2.027
	soft_temp=0.4	0.713	0.098	0.300	0.914	0.713	0.570	2.169
	soft_temp=0.5	0.680	0.105	0.326	0.912	0.742	0.571	2.221

Table 2: Comparison of three diversified variants of both Exponential and Logits-Addition watermarks in the QA task. These variants include **drop_prob=0.2**, sampling from the language model directly at a 0.2 probability; **shift_max=100**, where the watermark key is cyclically shifted within a 0-100 range; and **soft_temp=0.3**, which uses a softmax sampling with a temperature of 0.3 to balance randomness. **Vanilla** is the original GM watermark(Exponential and Logits-Addition) without any technique to enhance diversity. The detectability is measured by 100 detection tokens. Note that Logits-Addition+soft_temp is the Gumbelsoft watermark.

GumbelSoft watermark exhibits the highest detectability, likely explained by the expectation and variation theory in Theorem 1. Increased detection token amounts also improve detectability, aligning with findings from Chakraborty et al. (2023). The high-entropy Llama2-7b model in Completion tasks shows greater detectability than the lower entropy Llama2-7b-chat in QA tasks, as high entropy facilitates easier watermark embedding. Regarding generation quality (perplexity), GumbelSoft shows relatively low perplexity. In contrast, the KGW watermark’s biased logits modification leads to high perplexity, while Dipmark’s strategy of amplifying high-probability tokens results in the lowest perplexity in the Completion task. For the QA task, the low perplexity across all methods, attributed to the low entropy of Llama2-7b-chat, diminishes the value of comparative perplexity analysis.

4.4 Robustness

In this section, we assess the robustness of various decoding-based watermarking schemes, with results for the Completion task in Figures 5 and for the QA task in Appendix C.3. All texts, both watermarked and unwatermarked, were tested under the T5-span attack (explained in Appendix C.3).

Our key finding reveals that the Exponential and GumbelSoft watermarks are particularly robust against the T5-span attack, in contrast to other watermarks. Their AUROC values, as well as FPR and FNR metrics, remained stable post-attack, while other schemes experienced significant declines. This robustness can be attributed to the effective embedding of watermarks by the GumbelMax-trick, ensuring significant final statistics despite per-token score alterations. Further analysis, involving a comparative study of final statistic distributions between KGW and Gumbel-

	Completion	# tokens=40			# tokens=60			# tokens=100			
		AUROC ↑	FPR ↓	FNR ↓	AUROC ↑	FPR ↓	FNR ↓	AUROC ↑	FPR ↓	FNR ↓	PPL ↓
QA	Unwatermarked	-	-	-	-	-	-	-	-	-	11.576
	KGW	0.970	0.616	0.361	0.988	0.329	0.164	0.997	0.078	0.041	14.217
	Exponential	0.997	0.012	0.012	0.999	0.000	0.006	1.000	0.000	0.000	10.953
	Dipmark	0.935	0.693	0.565	0.968	0.483	0.362	0.988	0.274	0.153	8.664
	ITS	0.961	0.073	1.000	0.978	0.040	1.000	0.994	0.010	0.402	11.843
	Gumbelsoft	0.998	0.011	0.010	1.000	0.000	0.005	1.000	0.000	0.001	11.820
QA	Unwatermarked	-	-	-	-	-	-	-	-	-	1.980
	KGW	0.657	0.985	0.969	0.701	0.978	0.945	0.754	0.949	0.901	2.081
	Exponential	0.780	0.892	0.813	0.840	0.852	0.738	0.905	0.749	0.569	1.985
	Dipmark	0.588	0.988	0.982	0.615	0.981	0.984	0.646	0.979	0.970	1.792
	ITS	0.583	1.000	1.000	0.618	0.963	1.000	0.665	0.954	1.000	2.011
	Gumbelsoft	0.788	0.866	0.812	0.848	0.837	0.722	0.911	0.704	0.572	2.027

Table 3: A comparative analysis of the detectability across various decoding-based watermarking schemes. Detectability is assessed for varying token counts: 40, 60, and 100. The temperature for Gumbelsoft is set to 0.3.

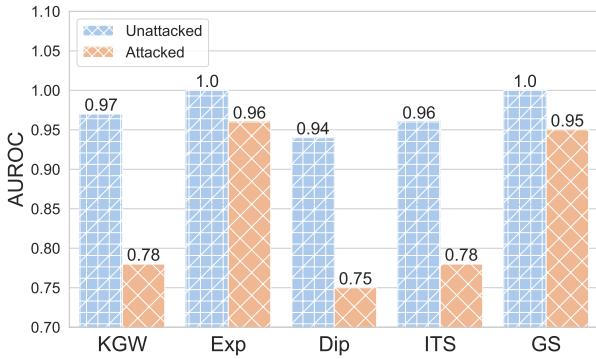


Figure 5: Comparison of the robustness of decoding-based watermark on Completion task. Blue histograms indicate unattacked conditions and red histograms show attacked scenarios. The AUROC is calculated for 40 detection tokens, with Gumbel-softmax set at a 0.3 temperature. **Exp**, **Dip**, and **GS** refer to Exponential, Dipmark, and Gumbelsoft, respectively.

Soft watermarks, is shown in Figure 6. The results demonstrate that while attacked watermarked texts under KGW show considerable overlap with unwatermarked texts, our GumbelSoft watermark displays less overlap, indicating its greater robustness.

5 Conclusion

We observed that the GumbelMax-trick-based watermark(GM watermark) produces identical responses to identical queries due to the deterministic nature of both the Decoder and Pseudo-random functions. To address this, we introduce three diversified variants aimed at enhancing GM watermark diversity. Furthermore, we question the need for an Exponential transformation (Aaronson and Kirchner, 2023) in watermark embedding and propose a new approach named Logits-Addition watermark. Our experiments across these variants for both Ex-

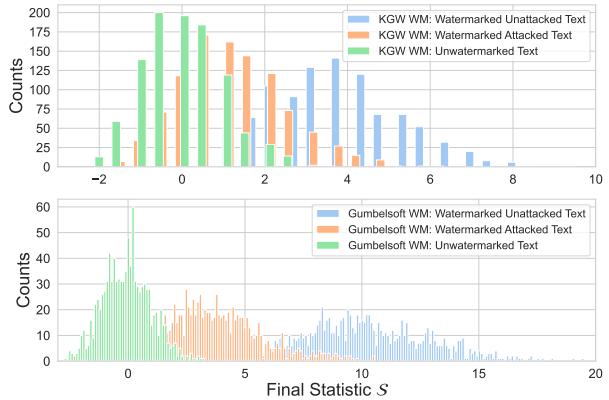


Figure 6: Comparison of final statistic for KGW and Gumbelsoft watermark on Completion task. The final statistic is calculated for 40 detection tokens, with Gumbelsoft set at the temperature of 0.3.

ponential and Logits-Addition watermarks identified GumbelSoft, a softmax-based Logits-Addition variant, as the optimal choice. Comparative analysis with other decoding-based watermarks demonstrated that GumbelSoft surpasses in detectability, maintains lower perplexity, and ensures higher robustness.

Limitations

GumbelSoft watermark’s Ngram pseudo-random function is susceptible to paraphrase attacks due to its dependence on the previous h tokens for key determination. In terms of quality assessment, we solely rely on perplexity, whereas some studies utilize downstream tasks for evaluation. Our mathematical analysis is focused solely on the Logits-Addition watermarking technique, we do not provide a comprehensive mathematical analysis of the GumbelSoft watermark.

451 Ethical Considerations

452 As advanced language models increasingly demonstrate remarkable capabilities, concerns regarding
453 their misuse have escalated. Consequently, the development of effective methods for detecting
454 machine-generated text has become crucial. The GM watermark has emerged as a highly effective
455 technique for differentiating between machine-generated and natural text. Nevertheless, the GM
456 watermark is limited by issues of diversity, which may hinder its practical application. The Gumbel-
457 Soft watermark represents a straightforward yet effective strategy to address this limitation. This
458 approach maintains the watermark's detectability while significantly enhancing its generative diversity.
459 We believe that our method will facilitate the broader implementation of the GM watermark in
460 various applications.

469 References

470 S. Aaronson and H. Kirchner. 2023. [Watermarking gpt outputs](#). Technical report, openai.

472 Mikhail J. Atallah, Victor Raskin, Michael Crogan,
473 Christian Hempelmann, Florian Kerschbaum, Dina
474 Mohamed, and Sanket Naik. 2001. Natural lan-
475 guage watermarking: Design, analysis, and a proof-
476 of-concept implementation. In *Proceedings of the 4th*
477 *International Workshop on Information Hiding*, IHW
478 '01, page 185–199, Berlin, Heidelberg. Springer-
479 Verlag.

480 Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi
481 Yang, and Yue Zhang. 2023. [Fast-detectgpt: Efficient](#)
482 [zero-shot detection of machine-generated text via](#)
483 [conditional probability curvature](#).

484 J. Brassil, S. Low, N. Maxemchuk, and L. O'Gorman.
485 2002. [Electronic marking and identification tech-](#)
486 [niques to discourage document copying](#). In *Pro-*
487 *ceedings of INFOCOM '94 Conference on Computer*
488 *Communications*.

489 Souradip Chakraborty, AmritSingh Bedi, Sicheng Zhu,
490 Bang An, Dinesh Manocha, and Furong Huang. 2023.
491 On the possibilities of ai-generated text detection.
492 *arXiv: Learning*.

493 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming
494 Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-
495 plan, Harri Edwards, Yuri Burda, Nicholas Joseph,
496 Greg Brockman, Alex Ray, Raul Puri, Gretchen
497 Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-
498 try, Pamela Mishkin, Brooke Chan, Scott Gray,
499 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz
500 Kaiser, Mohammad Bavarian, Clemens Winter,
501 Philippe Tillet, Felipe Petroski Such, Dave Cum-
502 mings, Matthias Plappert, Fotios Chantzis, Eliza-
503 beth Barnes, Ariel Herbert-Voss, William Heben

504 Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie
505 Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,
506 William Saunders, Christopher Hesse, Andrew N.
507 Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan
508 Morikawa, Alec Radford, Matthew Knight, Miles
509 Brundage, Mira Murati, Katie Mayer, Peter Welinder,
510 Bob McGrew, Dario Amodei, Sam McCandlish, Ilya
511 Sutskever, and Wojciech Zaremba. 2021. [Evaluating](#)
512 [large language models trained on code](#).

513 Yutian Chen, Hao Kang, Vivian Zhai, Liangze Li, Rita
514 Singh, and Bhiksha Raj. 2023. [Gpt-sentinel: Disting-](#)
515 [guishing human and chatgpt generated content](#).

516 Miranda Christ, Sam Gunn, and Or Zamir. 2023. [Unde-](#)
517 [etectable watermarks for language models](#).

518 Zhijie Deng, Hongcheng Gao, Yibo Miao, and Hao
519 Zhang. 2023. [Efficient detection of llm-generated](#)
520 [texts with a bayesian surrogate model](#).

521 Jaiden Fairoze, Sanjam Garg, Somesh Jha, Saeed
522 Mahloujifar, Mohammad Mahmoody, and Mingyuan
523 Wang. 2023. [Publicly detectable watermarking for](#)
524 [language models](#). Cryptology ePrint Archive, Paper
525 2023/1661. <https://eprint.iacr.org/2023/1661>.

526 Pierre Fernandez, Antoine Chaffin, Karim Tit, Vivien
527 Chappelier, and Teddy Furon. 2023. [Three bricks to](#)
528 [consolidate watermarks for large language models](#).

529 Hans Fischer. 2011. [A History of the Central Limit](#)
530 [Theorem](#). New York: Springer.

531 Sebastian Gehrmann, Hendrik Strobelt, and Alexander
532 Rush. 2019. [Gltr: Statistical detection and visual-](#)
533 [ization of generated text](#). In *Proceedings of the 57th*
534 *Annual Meeting of the Association for Computational*
535 *Linguistics: System Demonstrations*.

536 E.J. Gumbel. 1954. Statistical theory of extreme values
537 and some practical applications: A series of lectures.
538 *U.S. Government Printing Office eBooks*, U.S. Gov-
539 *ernment Printing Office eBooks*.

540 Xuanli He, Qiongkai Xu, Yi Zeng, Lingjuan Lyu,
541 Fangzhao Wu, Jiwei Li, and Ruoxi Jia. 2022. [Cater:](#)
542 [Intellectual property protection on text generation](#)
543 [apis via conditional watermarks](#).

544 Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu,
545 Hongyang Zhang, and Heng Huang. 2023. [Unbiased](#)
546 [watermark for large language models](#).

547 John Kirchenbauer, Jonas Geiping, Yuxin Wen,
548 Jonathan Katz, Ian Miers, and Tom Goldstein. 2023.
549 [A watermark for large language models](#). In *Proceed-*
550 *ings of the 40th International Conference on Machine*
551 *Learning*, volume 202 of *Proceedings of Machine*
552 *Learning Research*, pages 17061–17084. PMLR.

553 Kalpesh Krishna, Yixiao Song, Marzena Karpinska,
554 John Wieting, and Mohit Iyyer. 2023. [Paraphras-](#)
555 [ing evades detectors of ai-generated text, but retrieval](#)
556 [is an effective defense](#).

557	Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. 2023. Robust distortion-free watermarks for language models.	607
558		608
559		609
560	Linyang Li, Pengyu Wang, Ke Ren, Tianxiang Sun, and Xipeng Qiu. 2023. Origin tracing and detecting of llms.	610
561		611
562		
563	Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and Lijie Wen. 2023a. A semantic invariant robust watermark for large language models.	612
564		613
565		614
566	Aiwei Liu, Leyi Pan, Yijian Lu, Jingjing Li, Xuming Hu, Xi Zhang, Lijie Wen, Irwin King, Hui Xiong, and Philip S. Yu. 2024. A survey of text watermarking in the era of large language models.	615
567		
568		
569		
570	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>Cornell University - arXiv, Cornell University - arXiv</i> .	616
571		617
572		618
573		619
574		620
575		621
576	Yixin Liu, Hongsheng Hu, Xun Chen, Xuyun Zhang, and Lichao Sun. 2023b. Watermarking classification dataset for copyright protection.	622
577		623
578		624
579	Zeyan Liu, Zijun Yao, Fengjun Li, and Bo Luo. 2023c. Check me if you can: Detecting chatgpt-generated academic writing using checkgpt.	625
580		626
581		627
582	Hasan Mesut Meral, Bülent Sankur, A. Sumru Özsoy, Tunga Güngör, and Emre Sevinç. 2009. Natural language watermarking via morphosyntactic alterations. <i>Computer Speech & Language</i> , page 107–125.	628
583		629
584		630
585		631
586	Fatemehsadat Mireshghallah, Justus Mattern, Sicun Gao, Reza Shokri, and Taylor Berg-Kirkpatrick. 2023. Smaller language models are better black-box machine-generated text detectors.	632
587		633
588		634
589		635
590	Yisroel Mirsky, Ambra Demontis, Jaidip Kotak, Ram Shankar, Deng Gelei, Liu Yang, Xiangyu Zhang, Maura Pintor, Wenke Lee, Yuval Elovici, and Battista Biggio. 2023. The threat of offensive ai to organizations. <i>Computers & Security</i> , page 103006.	636
591		
592		
593		
594		
595	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature.	637
596		638
597		639
598		640
599	OpenAI. 2023a. Gpt-4 technical report. <i>ArXiv</i> , abs/2303.08774.	641
600		642
601	OpenAI. 2023b. New ai classifier for indicating ai-written text. Technical report, openai.	643
602		644
603	Lip Yee Por, KokSheik Wong, and Kok Onn Chee. 2012. Unispach: A text-based data hiding method using unicode space characters. <i>Journal of Systems and Software</i> , 85(5):1075–1082.	645
604		646
605		647
606		648
557	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and PeterJ. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>arXiv: Learning</i> .	649
558		650
559		651
560	Jie Ren, Han Xu, Yiding Liu, Yingqian Cui, Shuaiqiang Wang, Dawei Yin, and Jiliang Tang. 2023. A robust semantics-based watermark for large language model against paraphrasing.	652
561		653
562		654
563	Stefano Giovanni Rizzo, Flavio Bertini, and Danilo Montesi. 2016. Content-preserving text watermarking through unicode homoglyph substitution. In <i>Proceedings of the 20th International Database Engineering & Applications Symposium on - IDEAS '16</i> .	655
564		656
565		657
566	Ryoma Sato, Yuki Takezawa, Han Bao, Kenta Niwa, and Makoto Yamada. 2023. Embarrassingly simple text watermarks.	658
567		659
568		660
569	Murray Shanahan and Catherine Clarke. 2023. Evaluating large language model creativity from a literary perspective.	661
570	Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, Miles McCain, Alex Newhouse, Jason Blazakis, Kris McGuffie, and Jasmine Wang. 2019. Release strategies and the social impacts of language models.	662
571		663
572		664
573		665
574		666
575		667
576	Jinyan Su, Terry Yue Zhuo, Di Wang, and Preslav Nakov. 2023. Detectllm: Leveraging log rank information for zero-shot detection of machine-generated text.	668
577		669
578		670
579	Ruixiang Tang, Qizhang Feng, Ninghao Liu, Fan Yang, and Xia Hu. 2023. Did you train on my dataset? towards public dataset protection with clean-label backdoor watermarking.	671
580		672
581		673
582	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	674
583		675
584		676
585		677
586	Edward Tian. 2023. More than an ai detector preserve what's human. Technical report, Princeton University.	678
587		679
588		680
589		681
590	Yuchuan Tian, Hanting Chen, Xutao Wang, Zheyuan Bai, Qinghua Zhang, Ruifeng Li, Chao Xu, and Yunhe Wang. 2023. Multiscale positive-unlabeled detection of ai-generated texts.	682
591		683
592		684
593		685
594		686
595	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	687
596		688
597		689
598		690
599		691
600		692
601		693
602		694
603		695
604		696
605		697
606		698

662	Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.	716
663		717
664		718
665		
666		
667		
668		
669		
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686	Eduard Tulchinskii, Kristian Kuznetsov, Laida Kushnareva, Daniil Cherniavskii, Serguei Baranikov, Irina Piontkovskaya, Sergey Nikolenko, and Evgeny Burnaev. 2023. Intrinsic dimension estimation for robust detection of ai-generated texts.	722
687		723
688		724
689		725
690		726
691		
692		
693	Christoforos Vasilatos, Manaar Alam, Talal Rahwan, Yasir Zaki, and Michail Maniatakos. 2023. Howkgpt: Investigating the detection of chatgpt-generated university student homework through context-aware perplexity analysis.	727
694		728
695		729
696		
697	Saranya Venkatraman, Adaku Uchendu, and Dongwon Lee. 2023. Gpt-who: An information density-based machine-generated text detector.	
698		
699		
700	Lean Wang, Wenkai Yang, Deli Chen, Hao Zhou, Yankai Lin, Fandong Meng, Jie Zhou, and Xu Sun. 2023. Towards codable watermarking for injecting multi-bit information to llm.	
701		
702		
703		
704	Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Derek F. Wong, and Lidia S. Chao. 2023a. A survey on llm-generated text detection: Necessity, methods, and future directions.	
705		
706		
707		
708	Yihan Wu, Zhengmian Hu, Hongyang Zhang, and Heng Huang. 2023b. Dipmark: A stealthy, efficient and resilient watermark for large language models.	
709		
710		
711		
712	Xi Yang, Jie Zhang, Kejiang Chen, Weiming Zhang, Zehua Ma, Feng Wang, and Nenghai Yu. 2021. Tracing text provenance via context-aware lexical substitution.	
713		
714		
715		
716	Xianjun Yang, Wei Cheng, Yue Wu, Linda Petzold, William Yang Wang, and Haifeng Chen. 2023a. Dna-gpt: Divergent n-gram analysis for training-free detection of gpt-generated text.	
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		

730 A Baselines

731 Here, we present a consolidated mathematical
732 representation within a unified taxonomy for the base-
733 line watermark schemes. For the KGW scheme, as
734 proposed by Kirchenbauer et al. (2023):

- 735 • **Context:** The previous h tokens.
- 736 • **Pseudo-random Function:** $F_{sk}(\text{context})$
737 hashes the context to seed, then uses this seed
738 to generate a random vector in $\{0, 1\}^{|\mathcal{V}|}$, the
739 vector has $\gamma|\mathcal{V}|$ 1's and $(1 - \gamma)|\mathcal{V}|$ 0's.
- 740 • **Decoder:** $\Gamma(\xi_t, l_t)$ samples a token from
741 $\text{softmax}(\delta * \xi_t + l_t)$.
- 742 • **Scorer:** $\phi(\xi_t, w_t) = \xi_t[w_t]$.
- 743 • **Statistic Aggregator:**

$$\Phi(s_1, \dots, s_T) = \frac{\sum_{t=1}^T s_t - \gamma T}{\sqrt{T}\gamma(1 - \gamma)}$$

743 For the Exponential scheme, as proposed by
744 Aaronson and Kirchner (2023):

- 745 • **Context:** The previous h tokens.
- 746 • **Pseudo-random Function:** $F_{sk}(\text{context})$
747 hashes the context to a seed, then uses this
748 seed to generate a random vector in $(0, 1)^{|\mathcal{V}|}$,
749 each element is uniformly sample from $(0, 1)$.
- 750 • **Decoder:** $\Gamma(\xi_t, l_t) = \arg \max_{1 \leq i \leq |\mathcal{V}|} \frac{\log \xi_t[i]}{p_t[i]}$,
751 where $p_t = \text{softmax}(l_t)$.
- 752 • **Scorer:** $\phi(\xi_t, w_t) = -\log(1 - \xi_t[w_t])$.
- 753 • **Statistic Aggregator:**

$$\Phi(s_1, \dots, s_T) = \frac{1}{\sqrt{T}} \sum_{t=1}^T s_t - \sqrt{T}$$

753 For the Dipmark scheme, as proposed by Wu
754 et al. (2023b):

- 755 • **Context:** The previous h tokens.
- 756 • **Pseudo-random Function:** $F_{sk}(\text{context})$
757 hashes the context to a seed, then uses this
758 seed to generate a random permutation on the
759 vocabulary \mathcal{V} .

- 760 • **Decoder:** $\Gamma(\xi_t, l_t)$ samples token $\xi_t[i]$
761 with probability $\lambda(i) = \lambda(i - 1)$, where
762 $\lambda(i) = \max\{\sum_{j=1}^i p_t(\xi_t[j]) - \alpha, 0\} +$
763 $\max\{\sum_{j=1}^i p_t(\xi_t[j]) - (1 - \alpha), 0\}$, where
764 $p_t = \text{softmax}(l_t)$.

- 765 • **Scorer:** $\phi(\xi_t, w_t) = \mathbf{1}_{\{w_t \in \xi_t[\gamma|\mathcal{V}|:|\mathcal{V}|]\}}$.

- 766 • **Statistic Aggregator:**

$$\Phi(s_1, s_2, \dots, s_T) = \frac{\sum_{t=1}^T s_t - (1 - \gamma)T}{\sqrt{T}}$$

766 For the ITS scheme, as proposed by Kuditipudi
767 et al. (2023):

- 768 • **Context:** A global watermark key sequence ξ -
769 list and the position index t . Each watermark
770 key ξ -list[i] consists of a permutation π on the
771 vocabulary and a random number μ in $(0, 1)$.
- 772 • **Pseudo-random Function:** $F_{sk}(\text{context}) =$
773 ξ -list[t].
- 774 • **Decoder:** $\Gamma(\xi_t, l_t) = \pi^{-1}(\min\{\pi(i) : p_t(\{j : \pi(j) \leq \pi(i)\}) \geq \mu\})$, where $\xi_t = (\mu, \pi)$ and
775 $p_t = \text{softmax}(l_t)$.
- 776 • **Scorer:** $\phi(\xi_t, w_t) = |\mu - \eta(\pi(w_t))|$, where
777 $\eta(k) = \frac{k-1}{|\mathcal{V}|-1}$.
- 778 • **Statistic Aggregator:**

$$\Phi(s_1, s_2, \dots, s_T) = -\frac{1}{T} \sum_{t=1}^T s_t$$

779 We now explore the design principles underlying
780 these fundamental components.

781 **Context** For Watermark generators and detec-
782 tors, it is essential to recognize that they share
783 the same context, which is constrained to the pre-
784 vious tokens of w_t . This limitation arises from
785 the auto-regressive nature of the Watermark gen-
786 erator, which sequentially generates tokens from
787 left to right. A conventional approach for con-
788 text selection is to use the previous h tokens:
789 w_{t-h}, \dots, w_{t-1} . However, this design is vuln-
790 erable to paraphrase attacks, as such attacks can alter
791 these preceding tokens, subsequently modifying
792 the watermark key ξ_t , and ultimately affecting the
793 per-token score s_t . A more robust approach in-
794 volves considering the semantic meaning of previ-
795 ous tokens, based on the rationale that paraphrasing

796 maintains the semantics despite changing the tokens (Liu et al., 2023a). Kuditipudi et al. (2023)
 797 suggest utilizing a global list for storing all watermark keys and retrieving a specific watermark key
 798 using the position index t

801 **Pseudo-random Function.** The pseudo-random
 802 function’s role is to determine the watermark key
 803 ξ_t based on the given context. This function could
 804 be as basic as a hash function of the context or
 805 might involve leveraging an embedding model to
 806 extract the context’s semantic content. An alterna-
 807 tive method is to use the position index t to retrieve
 808 a watermark key from a global list. It is crucial to
 809 note that both the Watermark generator and detec-
 810 tor share the same pseudo-random function.

811 **Decoder.** The decoder is integral to the Water-
 812 mark generator, utilizing the watermark key ξ_t and
 813 the logits vector l_t to determine the subsequent
 814 token w_t . Implementation methods for this compo-
 815 nent vary among different watermarks.

816 **Scorer.** The scorer is to establish a correlation
 817 between the watermark key ξ_t and the token w_t .
 818 Nevertheless, using a global watermark key list and
 819 the position index t for key retrieval can result in
 820 a significant alignment shift issue. This problem
 821 manifests as a misalignment between the water-
 822 mark key ξ_t and the token w_t in texts subjected to
 823 insertion or deletion attacks. To address this, Ku-
 824 ditipudi et al. (2023) recommend using alignment
 825 cost or edit distance for computing the sequence
 826 score, as opposed to the per-token score.

827 **Statistic Aggregator.** Finally, the statistic aggre-
 828 gator compiles all per-token scores or employs a
 829 single sequence score to ascertain the presence of a
 830 watermark. A typical method involves calculating
 831 the z-score and p-value of collected scores. Alter-
 832 natively, one could use the empirical cumulative
 833 distribution function (Kuditipudi et al., 2023) for
 834 final statistical analysis.

835 B Mathematical Proofs

836 B.1 Unbiasedness for GumbelMax

837 We demonstrate that the GumbelMax-trick is math-
 838 ematically equivalent to directly sampling from the
 839 categorical distribution π , thereby establishing its
 840 unbiased nature when utilized in text watermarking
 841 applications.

842 Denote the vocabulary as \mathcal{V} , the vector of log-
 843 its as $l = (l_1, \dots, l_{|\mathcal{V}|})$, and a sequence of inde-

pendent Gumbel-distributed random variables as
 $\xi_1, \dots, \xi_{|\mathcal{V}|} \sim \text{Gumbel}(0, 1)$.

$$\mathbb{P}_{\xi \sim \text{Gumbel}(0,1)^{|\mathcal{V}|}} \left[\arg \max_{1 \leq i \leq |\mathcal{V}|} \{\xi_i + l_i\} = x \right] \quad (1)$$

$$= \mathbb{P}_{\eta \sim Q(\cdot)} \left[\arg \max_{1 \leq i \leq |\mathcal{V}|} \eta_i = x \right] \quad (2)$$

$$= \mathbb{P}_{\eta \sim Q(\cdot)} [\eta_x \geq \eta_i, \forall i \neq x] \quad (3)$$

$$= \mathbb{E}_{Y \sim q(\cdot)} \left[\prod_{i \neq x} \mathbb{P}[Y \geq \eta_i] \right] \quad (4)$$

$$= \int_{-\infty}^{+\infty} f(y - l_x) \prod_{i \neq x} F(y - l_i) dy \quad (5)$$

$$= \int_{-\infty}^{+\infty} e^{-((y-l_x)+e^{-(y-l_x)})} \prod_{i \neq x} e^{-e^{-(y-l_i)}} dy \quad (6)$$

$$= \int_{-\infty}^{+\infty} e^{\sum_{i=1}^{|\mathcal{V}|} -e^{l_i-y}} e^{l_x-y} dy \quad (7)$$

$$= \int_{-\infty}^{+\infty} e^{-e^{-y} \sum_i e^{l_i}} e^{-y} e^{l_x} dy \quad (8)$$

$$= Z p_x \int_{-\infty}^{+\infty} e^{-Z e^{-y}} e^{-y} dy \quad (9)$$

$$= Z p_x \frac{1}{Z} \quad (10)$$

$$= p_x \quad (11)$$

In equation (2), we introduce a variable substitu-
 857 tion $\eta_i = \xi_i + l_i$ for simplification. Moving to
 858 equation (4), we define the random variable Y to
 859 represent η_x for enhanced clarity, and we assume
 860 that Y follows the distribution $q(\cdot)$. Furthermore,
 861 we utilize the independence of $\eta_i, i = 1, \dots, |\mathcal{V}|$.
 862 In equation (5), $f(\cdot)$ denotes the probability den-
 863 sity function of the Gumbel(0,1) distribution, while
 864 $F(\cdot)$ represents its cumulative distribution function.
 865 Finally, in equation (9), we introduce $Z = \sum_i e^{l_i}$
 866 as a notation simplification.

867 B.2 Equivalence of Two Representations

868 We contend that the token generation processes
 869 for the Logits-Addition watermark and the Expon-
 870 ential watermark are mathematically equivalent,
 871 though their respective per-token scoring mech-
 872 anisms differ. To illustrate this equivalence, we
 873 present the following equations, equation (12) de-
 874 fines the Logits-Addition watermark, while equa-
 875 tion (15) corresponds to the definition of the Expo-

877 nential watermark.

$$w = \arg \max_{1 \leq i \leq |\mathcal{V}|} \{\eta_i + l_i\} \quad (12)$$

$$= \arg \max_{1 \leq i \leq |\mathcal{V}|} e^{\eta_i + l_i} \quad (13)$$

$$= \arg \max_{1 \leq i \leq |\mathcal{V}|} \frac{-p_i}{\log \xi_i} \quad (14)$$

$$= \arg \max_{1 \leq i \leq |\mathcal{V}|} \frac{\log \xi_i}{p_i} \quad (15)$$

882 Here, we utilize the relationship $p_i =$
883 $\text{softmax}(l_i) \propto e^{l_i}$ and $\eta_i = -\log(-\log \xi_i)$.
884 In these notations, we omit the position index t for
885 simplicity, and we assume $\eta_i \sim \text{Gumbel}(0, 1)$ and
886 $\xi_i \sim \text{Uniform}(0, 1)$.

887 While the token generation processes for the
888 Logits-Addition watermark and the Exponential
889 watermark are equivalent, their scoring methods
890 are distinct:

$$w = \eta[w] \quad (16)$$

$$= -\log(-\log \xi[w]) \quad (17)$$

$$\neq -\log(1 - \xi[w]) \quad (18)$$

894 Equation (16) specifies the per-token scoring for
895 the Logits-Addition watermark while equation (18)
896 is the scoring method for the Exponential water-
897 mark.

898 B.3 Expectation and Variance for Per-token 899 Score

900 We now establish the expected per-token score for
901 texts, distinguishing between those with and without
902 the Logits-Addition watermark. In the case of
903 unwatermarked texts, the watermark token, w_t , ex-
904 hibits no correlation with ξ_t . Consequently, $\xi_t[w_t]$
905 adheres to a Gumbel(0,1) distribution. This leads
906 to

$$907 \mathbb{E}[s_t] = \mathbb{E}[\xi_t[w_t]] = \gamma,$$

$$\text{Var}[s_t] = \text{Var}[\xi_t[w_t]] = \frac{\pi^2}{6}.$$

908 Conversely, for watermarked texts, a correlation
909 exists between w_t and ξ_t . This correlation alters
910 the distribution of $\xi_t[w_t]$, diverging it from the stan-
911 dard Gumbel(0,1) form. To compute its expected
912 value, we define $\xi_t[w_t]$ as a random variable X.
913 We then deduce its cumulative distribution func-
914 tion (CDF), $F(x)$, and probability density function
915 (PDF), $f(x)$. Utilizing this PDF, we calculate the
916 expectation of X. For simplification, we exclude
917 the position index ‘t’ in the subsequent equations.

918 Here, ξ_i represents $\xi[i]$, implying ξ_w is equivalent
919 to $\xi_t[w_t]$, and similar conventions apply to other
920 notations.

$$F(x) = \mathbb{P}[X \leq x] \quad (19)$$

$$= \mathbb{P}[\xi_w \leq x] \quad (20)$$

$$= \mathbb{P}[\xi_i + l_i - l_w \leq x, \forall i] \quad (21)$$

$$= \mathbb{P}\left[e^{\xi_i + l_i - l_w} \leq e^x, \forall i\right] \quad (22)$$

$$= \mathbb{P}\left[\frac{-1}{\log h_i} \frac{p_i}{p_w} \leq e^x, \forall i\right] \quad (23)$$

$$= \mathbb{P}\left[\frac{p_i}{p_w} e^{-x} \leq -\log h_i, \forall i\right] \quad (24)$$

$$= \prod_{i=1}^{|\mathcal{V}|} \mathbb{P}\left[\frac{p_i}{p_w} e^{-x} \leq -\log h_i\right] \quad (25)$$

$$= \prod_{i=1}^{|\mathcal{V}|} 1 - \mathbb{P}\left[-\log h_i \leq \frac{p_i}{p_w} e^{-x}\right] \quad (26)$$

$$= \prod_{i=1}^{|\mathcal{V}|} 1 - (1 - e^{-\frac{p_i}{p_w} e^{-x}}) \quad (27)$$

$$= \prod_{i=1}^{|\mathcal{V}|} e^{-\frac{p_i}{p_w} e^{-x}} \quad (28)$$

$$= e^{\sum_{i=1}^{|\mathcal{V}|} -\frac{p_i}{p_w} e^{-x}} \quad (29)$$

$$= e^{\frac{-e^{-x}}{p_w}} \quad (30)$$

In equation (22), we utilize the fact that $p_i \propto e^{l_i}$ and $\xi_i = -\log(-\log h_i)$. Equation (24) leverages the independence of h_i . Equation (26) uses the fact that $-\log h_i \sim \text{Exp}(1)$. Finally, equation (29) employs the fact that $\sum_{i=1}^{|\mathcal{V}|} p_i = 1$. Hence, the density function:

$$f(x) = F'(x) = \frac{e^{-x}}{p_w} e^{\frac{-e^{-x}}{p_w}}$$

The expectation:

$$\mathbb{E}[\xi_t[w_t]] = \mathbb{E}[X] \quad (31)$$

$$= \int_{-\infty}^{+\infty} x f(x) dx \quad (32)$$

$$= - \int_{-\infty}^{+\infty} x \frac{e^x}{p_w} e^{\frac{-e^x}{p_w}} dx \quad (33)$$

$$= -\frac{1}{p_w} [p_w \log p_w - \gamma p_w] \quad (34)$$

$$= -\log p_w + \gamma \quad (35)$$

The equation (32) use the fact that:

$$\int_{-\infty}^{+\infty} x e^x e^{\frac{-e^x}{t}} dx = t \log t - \gamma t$$

We now prove the fact. This is not a standard integral that can be solved by elementary functions. However, we can attempt to solve it using the substitution method and some properties of the Gamma and incomplete Gamma functions, which are commonly used to handle integrals involving exponentials of exponentials.

$$\int_{-\infty}^{+\infty} xe^x e^{-\frac{e^x}{t}} dx \quad (36)$$

$$= \int_0^{+\infty} \log(u) e^{-u/t} du \quad (37)$$

$$= \int_0^{+\infty} \log(vt) e^{-v} dv \quad (38)$$

$$= t \int_0^{+\infty} (\log(v) + \log(t)) e^{-v} dv \quad (39)$$

$$= t \log(t) \int_0^{+\infty} e^{-v} dv + t \int_0^{+\infty} \log(v) e^{-v} dv \quad (40)$$

$$= t \log(t) + t \int_0^{+\infty} \log(v) e^{-v} dv \quad (41)$$

$$= t \log(t) - \gamma t \quad (42)$$

In Equation(35), we use variable substitution $u = e^x$, In Equation(36), we use variable substitution $v = \frac{u}{t}$, In Equation (39), we use the definition of Euler-Mascheroni constant: $\gamma = -\int_0^{+\infty} \log(v) e^{-v} dv$.

As for the variance of the per-token score for watermarked text, we can also derive it via the probability density function $f(x)$.

$$\text{Var}[s_t] \quad (43)$$

$$= \text{Var}[\xi_t[w_t]] \quad (44)$$

$$= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \quad (45)$$

$$= \int_{-\infty}^{+\infty} x^2 f(x) dx - (-\log p_w + \gamma)^2 \quad (46)$$

$$= \int_{-\infty}^{+\infty} x^2 \frac{e^{-x}}{p_w} e^{-\frac{e^{-x}}{p_w}} dx - (-\log p_w + \gamma)^2 \quad (47)$$

$$\leq \frac{2p_w^2}{(1-p_w)^3} + \frac{2}{p_w} - (-\log p_w + \gamma)^2 \quad (48)$$

In equation(48), we use the fact that

$$\int_{-\infty}^{+\infty} x^2 \frac{e^{-x}}{p_w} e^{-\frac{e^{-x}}{p_w}} dx \leq \frac{2p_w^2}{(1-p_w)^3} + \frac{2}{p_w}$$

We now prove it:

$$\int_0^{+\infty} x^2 \frac{e^{-x}}{p_w} e^{-\frac{e^{-x}}{p_w}} dx \quad (49)$$

$$\leq \int_0^{+\infty} x^2 \frac{e^{-x}}{p_w} dx \quad (50)$$

$$= \frac{2}{p_w} \quad (51)$$

$$\int_{-\infty}^0 x^2 \frac{e^{-x}}{p_w} e^{-\frac{e^{-x}}{p_w}} dx \quad (52)$$

$$= \int_0^{+\infty} x^2 \frac{e^x}{p_w} e^{-\frac{e^x}{p_w}} dx \quad (53)$$

$$= \int_0^{+\infty} \frac{x^2}{p_w} e^{(x-\frac{e^x}{p_w})} dx \quad (54)$$

$$\leq \int_0^{+\infty} \frac{x^2}{p_w} e^{(1-\frac{1}{p_w})x} dx \quad (55)$$

$$= \frac{2p_w^2}{(1-p_w)^3} \quad (56)$$

A similar theorem also holds for the Exponential watermark. For unwatermarked text:

$$\mathbb{E}[s_t] = \mathbb{E}[-\log(1 - \xi_t[w_t])] = 1$$

$$\text{Var}[s_t] = \text{Var}[-\log(1 - \xi_t[w_t])] = 1$$

For watermarked text,

$$\begin{aligned} \mathbb{E}[s_t] &= \mathbb{E}[-\log(1 - \xi_t[w_t])] \\ &\geq 1 + \left(\frac{\pi^2}{6} - 1\right)(-p_w \log p_w), \end{aligned}$$

$$\begin{aligned} \text{Var}[s_t] &= \text{Var}[-\log(1 - \xi_t[w_t])] \\ &= \psi_1(1) - \psi_1(1 + \frac{1}{p_w}), \end{aligned}$$

where ψ_1 is the trigamma function. The proof can be found in [Fernandez et al. \(2023\)](#)

C Experiment details

We describe all experiment details here.

C.1 Diversity

We began by carefully selecting 40 high-entropy prompts to elicit a wide range of completions. These prompts were split evenly into two categories: 20 prompts followed a Completion format tailored for Llama2-7b, while the remaining 20 were structured in a QA format, specifically designed for Llama2-7b-chat. Each prompt was queried 50 times, and we assessed the resulting

		Diversity			Detectability			Quality
		Self-Bleu ↓	Dist-1 ↑	Dist-2 ↑	AUROC ↑	FPR ↓	FNR ↓	PPL ↓
		vanilla	1.000	0.010	0.014	1.000	0.000	0.000
Exponential	drop_prob=0.05	0.367	0.222	0.529	1.000	0.000	0.000	11.450
	drop_prob=0.10	0.227	0.254	0.633	1.000	0.000	0.001	11.423
	drop_prob=0.20	0.146	0.300	0.733	1.000	0.000	0.001	11.839
	drop_prob=0.30	0.113	0.307	0.766	1.000	0.000	0.002	11.911
	drop_prob=0.40	0.087	0.317	0.788	1.000	0.001	0.005	11.964
	shift_max=10	0.991	0.079	0.146	0.999	0.000	0.003	11.307
	shift_max=30	0.798	0.158	0.333	0.999	0.000	0.002	11.084
	shift_max=50	0.645	0.184	0.403	1.000	0.000	0.003	11.222
	shift_max=100	0.414	0.221	0.496	0.999	0.000	0.003	11.102
	shift_max=200	0.247	0.235	0.546	0.999	0.000	0.004	11.068
Logits-Addition	soft_temp=0.1	0.388	0.210	0.490	1.000	0.000	0.000	11.218
	soft_temp=0.2	0.244	0.238	0.565	1.000	0.000	0.001	11.353
	soft_temp=0.3	0.202	0.265	0.630	1.000	0.000	0.001	11.610
	soft_temp=0.4	0.169	0.275	0.669	1.000	0.000	0.001	11.874
	soft_temp=0.5	0.146	0.285	0.686	1.000	0.000	0.001	12.222
	vanilla	1.000	0.010	0.014	1.000	0.000	0.000	10.953
	drop_prob=0.05	0.421	0.205	0.493	0.999	0.000	0.003	11.561
	drop_prob=0.10	0.209	0.268	0.652	0.999	0.000	0.003	11.754
	drop_prob=0.20	0.143	0.292	0.729	0.999	0.000	0.005	11.997
	drop_prob=0.30	0.097	0.309	0.774	0.999	0.001	0.005	11.890
Logits-Addition+soft_temp	drop_prob=0.40	0.093	0.319	0.790	0.999	0.003	0.006	12.156
	shift_max=10	0.991	0.078	0.144	0.999	0.000	0.006	11.228
	shift_max=30	0.806	0.159	0.335	0.998	0.002	0.006	11.243
	shift_max=50	0.627	0.188	0.412	0.998	0.000	0.006	11.250
	shift_max=100	0.417	0.220	0.497	0.998	0.000	0.006	11.536
	shift_max=200	0.246	0.242	0.559	0.998	0.001	0.007	11.243
	soft_temp=0.1	0.370	0.213	0.497	1.000	0.000	0.001	11.159
	soft_temp=0.2	0.227	0.243	0.581	1.000	0.000	0.002	11.309
	soft_temp=0.3	0.158	0.254	0.608	1.000	0.000	0.001	11.820
	soft_temp=0.4	0.121	0.276	0.661	1.000	0.000	0.001	12.831
	soft_temp=0.5	0.100	0.298	0.699	1.000	0.000	0.001	14.140

Table 4: Comparison of three variants of both Exponential and Logits-Addition watermarks in the Completion task. The variants include **drop_prob=0.2**, sampling from the language model directly at a 0.2 probability; **shift_max=100**, where the watermark key is cyclically shifted within a 0-100 range; and **soft_temp=0.3**, which uses a softmax sampling with a temperature of 0.3 to balance randomness. **Vanilla** is the original two GumbelMax watermarks without any technique to enhance diversity. The detectability is measured by 100 detection tokens. Note that Logits-Addition+soft_temp is our Gumbel-softmax watermark.

completions using metrics such as Self-Bleu, Distinct 1-gram, and Distinct 2-gram. The average values of these metrics were then computed for the 20 prompts in each category. We control the max generation length for each prompt to be 256 tokens.

For the soft_temp parameter, we tested five different temperature settings: 0.1, 0.2, 0.3, 0.4, and 0.5. In the case of the shifted watermark key, we experimented with five maximum shift values: 10, 30, 50, 100, and 200. For drop probability, the tested probabilities were 5%, 10%, 15%, 20%, and 40%. We evaluated detectability and quality using a sample of 100 generated tokens, while diversity assessments were conducted with a sample size of 256 tokens.

C.2 Detectability

The objective of text watermarking is to embed a concealed pattern into generated texts and subsequently detect this pattern to ascertain if the text is watermarked. We gathered 1,000 lengthy texts from the news-like subset of the C4 dataset, dividing each text into two parts: the first 50 words as prompt and the remaining as gold-completion. For each watermarking scheme, we utilized Llama2-7b to create both watermarked and unwatermarked completions for these 1,000 prompts. The effectiveness of each scheme was then assessed using the corresponding detector to evaluate 2,000 completions. Key metrics reported include AUROC (Area Under the Receiver Operating Characteristic), FPR

(False Positive Rate) at a fixed FNR (False Negative Rate) of 0.01, and FNR at a fixed FPR of 0.01.

Additionally, we compiled 1,000 lengthy texts from the alpaca dataset. Unlike the C4 dataset, here we used only the question as a prompt to query Llama2-7b-chat, with the subsequent detection process mirroring that of the C4 dataset.

In line with the detectability theorem by Chakraborty et al. (2023), we anticipate higher detectability in longer texts. Therefore, we report detection metrics for generated token lengths of 40, 60, 80, and 100. However, for quality assessment, we calculate perplexity only for texts with 100 generated tokens, as fewer tokens would inadequately represent quality measures. We use llama2-13b and llama2-13b-chat to evaluate ppl for the texts generated by llama2-7b and llama2-7b-chat respectively.

The hyper-parameters employed for each watermarking scheme are specified as follows: All experiments are conducted at a temperature setting of 1, except the Gumbel-softmax, which utilized a temperature setting of 0.3 to achieve an equilibrium between detectability and generation diversity. For KGW, we adopt $\delta = 2$ and $\gamma = 0.1$, following the recommendations by Kirchenbauer et al. (2023). For Dipmark, the parameters are set to $\alpha = 0.45$ and $\gamma = 0.5$, by Wu et al. (2023b). Regarding ITS, we utilize a sample of 500 texts from the C4 subset for the computation of reference scores.

We repeat the experiment 5 times to calculate the average value for each metric.

C.3 Robustness



Figure 7: Comparison of robustness of decoding-based watermark on QA task. Blue histograms indicate unattacked conditions and red histograms show attacked scenarios. The AUROC is calculated for 40 detection tokens, with Gumbel-softmax set at a 0.3 temperature. **Exp**, **Dip**, and **GS** refer to Exponential, Dipmark, and Gumbelsoft, respectively.

We employ the T5-span attack (Kirchenbauer et al., 2023) on both watermarked and unwatermarked texts. Each word in a text undergoes a potential attack with a probability of 0.5. For attacked words, we use their immediate five-word context (preceding and following) and apply t5-large (Raffel et al., 2019) for context-based word prediction, replacing the original word with the predicted one. This process may occasionally retain the original word; however, we opt not to enforce unique substitutions to avoid excessive time consumption.

D More Related Work

Zero-shot Methods. More Zero-shot methods are listed as follows: Recent studies include Krishna et al. (2023) advocating retrieval against paraphrase attacks, Su et al. (2023) leveraging log-rank ratios, Solaiman et al. (2019) using log probability, Bao et al. (2023) focusing on conditional probability curvature, and Venkatraman et al. (2023) employing uniform information density for improved detection.

Training-based methods. More training-based methods are listed as follows: OpenAI (2023b); Tian (2023) training classifiers from mixed sources, Yin et al. (2023) using graph structures and contrastive learning, and Tian et al. (2023) applying positive unlabeled training for classifier development.

Watermarking Techniques. More watermarking techniques are listed as follows: Techniques include text formatting for embedding watermarks by Por et al. (2012); Rizzo et al. (2016), context-aware lexical substitution by Yang et al. (2021), syntactic modifications by Atallah et al. (2001); Meral et al. (2009), training data watermarking by Liu et al. (2023b); Tang et al. (2023), a publicly detectable watermark proposed by Fairoze et al. (2023), and leveraging semantic meaning for robustness by Ren et al. (2023).

There are also some surveys on machine-generated content detection (Wu et al., 2023a; Yang et al., 2023b) and text watermarking (Liu et al., 2024).