PSEUDO- VS. TRUE-RANDOMNESS: RETHINKING DISTORTION-FREE WATERMARKS OF LANGUAGE MOD ELS UNDER WATERMARK KEY COLLISIONS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

022

024

025

026

027

028 029

031

Paper under double-blind review

ABSTRACT

Language model (LM) watermarking techniques inject a statistical signal into LM-generated content by substituting the random sampling process with pseudorandom sampling, using watermark keys as the random seed. Among these statistical watermarking approaches, distortion-free watermarks are particularly crucial because they embed watermarks into LM-generated content without compromising generation quality. However, one notable limitation of pseudo-random sampling compared to true-random sampling is that, under the same watermark keys (i.e., key collision), the results of pseudo-random sampling exhibit correlations. This limitation could potentially undermine the distortion-free property. Our studies reveal that key collisions are inevitable due to the limited availability of watermark keys, and existing distortion-free watermarks exhibit a significant distribution bias toward the original LM distribution in the presence of key collisions. Moreover, we go beyond the key collision condition and prove that achieving a perfect distortion-free watermark is impossible. To study the trade-off between watermark strength and its distribution bias, we introduce a new family of distortion-free watermarks-beta-watermark. Experimental results support that the beta-watermark can effectively reduce the distribution bias under key collisions.

1 INTRODUCTION

In an era where artificial intelligence surpasses human capabilities in generating text, the authenticity
and origin of such AI-generated content have become paramount concerns. Language model watermarking (Aaronson, 2022; Kirchenbauer et al., 2023; Christ et al., 2023; Kuditipudi et al., 2023; Hu
et al., 2023) provides a promising solution for distinguishing between human and machine-generated
text. This technique secretly embeds a statistical signal into the generated text using a pseudo-random
generator seeded with watermark keys. The embedded signal is then detected through a statistical hypothesis test, ensuring the traceability and verification of the text's origin.

Distortion-free watermarks (Aaronson, 2022; Christ et al., 2023; Kuditipudi et al., 2023; Hu et al., 2023) represent one of the most compelling techniques in language model watermarking. These watermarks are particularly valuable because they provably preserve the output distribution of the original language model. Specifically, the expected watermarked distribution with respect to the watermark keys remains identical to the original language model distribution, thus offering significant practical application potential.

However, the pseudo-random nature of the watermark generator may lead to correlations between
generated content when the watermark keys are identical (i.e., key collision). In extreme cases,
such as when the prompt remains the same, key collisions can result in identical generated content,
significantly limiting its application scenarios. For instance, when using GPT-4 to generate content,
if the initial output is unsatisfactory, a request to regenerate would typically yield a different result.
However, under a distortion-free watermarking scheme, the output may remain unchanged due to the
consistent application of the same watermark key. This limitation highlights a critical challenge in
the practical deployment of such watermarking techniques.

053 In our research, we comprehensively analyze the existing distortion-free watermarks and demonstrate, through both theoretical and empirical evidence, that *no distortion-free watermark can fully preserve* the original LM distribution under key collisions. Specifically, we categorize the level of distortion-free capability into three types: a) Step-wise distortion-free—the watermark preserves the LM distribution at a single token generation step; b) Weakly distortion-free—the watermark preserves the LM distribution for a one-time sentence generation; c) Strongly distortion-free—the watermark preserves the LM distribution across multiple sentence generations. Our findings indicate that all existing distortion-free watermarks are weakly distortion-free but not strongly distortion-free due to key collisions. Under the key collisions, In particular, we theoretically prove that there does not exist any detectable strongly distortion-free watermark. We also show that key collisions are inevitable given the limited number of watermark keys available in current schemes.

To mitigate the distribution bias caused by key collisions, we introduce the beta-watermark and
 develop a model-agnostic detector that can identify watermarks without requiring access to prompts
 or language models. Additionally, we design empirical metrics to measure the distribution bias
 resulting from key collisions. Through rigorous testing on widely-studied language models, including
 BART-large model (Liu et al., 2020) and LLaMA-2 (Touvron et al., 2023), our beta-watermark has
 demonstrated effectiveness in significantly reducing the distribution bias induced by key collisions.

- 069 Our main contributions are summarized as follows:
 - We identify three levels of distortion-free capabilities in watermarks—Step-wise, Weakly, and Strongly Distortion-free—revealing that existing watermarks are not strongly distortion-free and cannot preserve the original language model distribution under multiple generations

due to the inevitability of key collisions.

- Under watermark key collisions, we theoretically demonstrate a trade-off between watermark strength and its distribution bias to the original LM distribution—a smaller distribution bias results in weaker watermark strength. Based on our discussion on the distribution bias, we proposed a black-box distortion-free watermark detection approach, which can effectively check if an LM is watermarked given the original LM. Furthermore, we go beyond the key collision assumption and prove that strongly distortion-free watermarks does not exist.
 - We introduce beta-watermark, a new family of weakly distortion-free watermarks that can provably reduce the distribution bias by trading the watermarking strength. Through experiments on popular language models like BART-large and LLaMA-2, we demonstrate our theoretical findings that existing watermarks are not strongly distortion-free and beta-watermark can effectively reduce the distribution bias.
- 084 085

087

071

073

075

076

077

078

079

081

082

2 RELATED WORK

Statistical watermarks. Kirchenbauer et al. (2023) enhanced the statistical watermark framework originally introduced by Aaronson (2022), demonstrating the effectiveness of statistical watermarking 090 through extensive experiments on large language models. They splited the LM tokens into red 091 and green list, then promoted the use of green tokens by adding a fixed parameter δ to their logits. 092 Zhao et al. (2023) proposed the unigram watermark, which enhances the robustness of the statistical watermark by using one-gram hashing to produce watermark keys. Liu et al. (2023b) also improved 094 the robustness of statistical watermarking by leveraging the semantics of generated content as 095 watermark keys. Additionally, Liu et al. (2023a) proposed an unforgeable watermark scheme that 096 employs neural networks to modify token distributions instead of using traditional watermark keys. 097 However, these approaches may lead to significant changes in the distribution of generated text, 098 potentially compromising content quality.

099 **Distortion-free watermarks.** To preserve the original output distribution in watermarked content, 100 researchers have explored alternative strategies to modify the token distribution. Aaronson (2022) 101 introduced the first distortion-free watermarking strategy, which utilized Gumbel-reparametrization 102 to alter token distribution and the prefix n-gram content as the watermark keys. Christ et al. (2023) 103 and Kuditipudi et al. (2023) adopted the inverse-sampling and Gumbel-reparametrization to modify 104 the watermarked token distributions, where the watermark keys are based on the token position or a 105 fixed key list respectively. Notice Christ et al. (2023)'s method encounters resilience challenges under modifications and lacks empirical evidence regarding its detectability. Meanwhile, Kuditipudi et al. 106 (2023)'s detection process involves hundreds of resampling steps from the secret key distribution, 107 proving inefficient for processing lengthy texts. Hu et al. (2023) employed inverse-sampling and

permute-reweight methods for watermarking. But their detector is not model-agnostic, which requires access to the language model API and prompts. Wu et al. (2023) improved the permute-reweight methods and designed a model-agnostic detector. A detailed related work section is in Appendix B.

112 113 3 PRELIMINARY

114 **Notations.** Denote by $V := \{t_1, ..., t_N\}$ the vocabulary (or token) set of a language model, and 115 by N = |V| its size. Let V represent the set of all conceivable string sequences, including those 116 of zero length. A language model generates a token sequence based on a predetermined prompt. 117 For a single step in this process, the probability of generating the next token $x_{n+1} \in V$, given the 118 current context from x_1 to x_n , is represented as $P_M(x_{n+1} | x_1, x_2, \ldots, x_n)$. For brevity, we adopt 119 the condensed notation: $P_M(\boldsymbol{x}_{n+1:n+m} \mid \boldsymbol{x}_{1:n})$, where $\boldsymbol{x}_{n+1:n+m} = (x_{n+1}, \dots, x_{n+m})$. Note that 120 the prompt is deliberately omitted in this representation. Inherent to its design, the language model operates in an autoregressive mode. This implies that the combined probability of generating several 121 tokens, specifically from x_{n+1} to x_{n+m} , takes the form $P_M(x_{n+1:n+m} \mid x_{1:n}) = \prod_{i=1}^m P_M(x_{n+i} \mid x_{n+i})$ 122 $x_{1:n+i-1}$). 123

Watermarking problem definition. A language model (LM) service provider aims to watermark the generated content such that all other users can verify if the content is generated by the LM without needing access to the LM or the original prompt. A watermark framework primarily consists of two components: a *watermark generator* and a *watermark detector*. The watermark generator embeds a watermark into the text through a *Pseudo-random Distribution Adjustment rule* (PDA-rule), which is seeded by watermark keys. The watermark detector, on the other hand, detects the presence of the watermark within the content using a statistical hypothesis test.

Definition 3.1 (PDA-rule). Let \mathcal{P} represent the space of token distributions and let K denote the space of watermark keys. A Pseudo-random Distribution Adjustment rule (**PDA-rule**), defined as $F : \mathcal{P} \times K \to \mathcal{P}$, adjusts the token distribution based on a given watermark key.

134 Watermark generator. During the watermark generation process, the service provider modifies 135 the original language model distribution P_M using a *watermark key* $k \in K$ and a PDA-rule. Here, 136 the watermark key acts as a random seed to modify the distribution, after which the next token is 137 sampled from this modified distribution. A watermark key usually consists of a secret key sk and a context key (e.g., n-gram (Aaronson, 2022) or token position (Christ et al., 2023)). Let $\mathcal{F} := \{F : f : f \in \mathbb{C}\}$ 138 $\mathcal{P} \times K \to \mathcal{P}$ denote the set of PDA-rules. Specifically, let P_W denote the distribution of the LM 139 after watermarking, and k the watermark key, $P_W(t \mid \boldsymbol{x}_{1:n-1}) := F(P_M(\cdot \mid \boldsymbol{x}_{1:n-1}), k)(t), \forall t \in V$, 140 where $P_M(\cdot \mid x_{1:n-1})$ is the LM token distribution for sampling the *n*-th token. When sampling the 141 next token x_n , the language model samples from $P_W(\cdot | \boldsymbol{x}_{1:n-1})$ instead of $P_M(\cdot | \boldsymbol{x}_{1:n-1})$. This 142 mechanism allows the service provider to inject a statistical signal into the generated content. 143

The PDA-rule is the core of the watermark generator. A PDA-rule is considered distortion-free if and
only if it preserves the token distribution during watermark generation. To the best of our knowledge,
there are three types of distortion-free PDA-rules: inverse-sampling (Christ et al., 2023; Kuditipudi
et al., 2023; Hu et al., 2023), Gumbel-reparametrization (Aaronson, 2022; Kuditipudi et al., 2023; Fu
et al., 2024), and permute-reweight (Hu et al., 2023). A detailed introduction to these methods can be
found in Section 4.1. The formal definition of a distortion-free PDA-rule is presented below.

Definition 3.2 (Distortion-free PDA-rule). A PDA-rule F, is a distortion-free PDA-rule, if and only if for an arbitrary $LM P_M$, $\forall \mathbf{x}_{1:n} \in \mathcal{V}$, and $\forall i \leq n$, it holds that $P_M(x_i|\mathbf{x}_{1:i-1}) = \mathbb{E}_{k_i}[F(P_M(\cdot|\mathbf{x}_{1:i-1}), k_i)(x_i)].$

153 Watermark Detector. During the process of watermark detection, the user will have access only 154 to the watermark key and the PDA-rule F. The detector employs a hypothesis testing approach to identify the presence of the watermark signal. The hypothesis test is defined as: H_0 : The content is generated without the presence of watermarks, and H_1 : The content is generated with 156 the presence of watermarks. For the purposes of the statistical test, a score function s(x, k, F): $V \times K \times F \rightarrow \mathbb{R}$ is employed. Under H_0 , the score function is a random variable S_{H_0} where 157 158 $\Pr(S_{H_0} = s(t,k,F)|k,F) = \sum_{s(t',k,F)=s(t,k,F)} P_M(t'), \forall t \in V$, while under H_1 , the random 159 variable S_{H_1} becomes $\Pr(S_{H_1} = s(t,k,F)|k,F) = \sum_{s(t',k,F)=s(t,k,F)} P_W(t')$. Thus, we can use 160 the discrepancy between S_{H_0} and S_{H_1} to detect the watermark content. Given an observation (text 161 sequence) $x_{1:n}$, we define the test statistic $S(x_{1:n}) = \sum_{i=1}^{n} s(x_i, k, F)$ as the measure for the test.

The decision to reject the null hypothesis is based on the difference between $S(x_{1:n})$ and the expected value $\mathbb{E}_{H_0}[S(x_{1:n})]$.

Watermark Key. For each generating step, we will use a watermark key to seed the PDA-rule. There are generally three key sampling methods:

- (n-gram hashing) Aaronson (2022), Christ et al. (2023) and Hu et al. (2023) use a fixed secret key Sk_0 and the prefix n-gram s (e.g., $s = \mathbf{x}_{l-n:l-1}$ for generating x_l) to form the watermark keys, i.e., $K = \{(\mathsf{sk}_0, s) \mid s \in \mathcal{V}_n\}$, where \mathcal{V}_n represents the set of all n-grams with token set V. A history list is kept during one generation to ensure the watermark keys are unique. If the length of previously generated tokens is less than n, all preceding tokens are used as s.
 - (position hashing) Christ et al. (2023) uses a fixed secret key sk_0 and the token position are used as watermark keys, i.e., $K = \{(\mathsf{sk}_0, i) \mid i \in \mathbb{N}\}.$
- (fixed key set) Kuditipudi et al. (2023) uses a fixed secret key sk_0 generates a set of watermark keys, $K = \{k_1, \ldots, k_{n_0}\}$. During token generation at step *i*, a random integer *r* is sampled, and $k_{(i+r) \mod n_0}$ is used as the seed for the PDA-rule. If the token length exceeds n_0 , we will sample from the original LM distribution instead.

Definition 3.3 (Key collision). Key collision refers to scenarios where the same watermark keys are used to seed the PDA-rule.

All three watermark key sampling methods mentioned previously have a limited number of keys given the fixed secret key \mathbf{sk}_0 . The maximum key volume is $|V|^n$ for n-gram hashing, l_0 for position hashing, and n_0 for the fixed key set. Here, l_0 represents the maximum token length for the language model, typically ranging from 10^4 to 10^6 . Therefore, if we only have one secret key, key collisions will occur when the number of queries and the generated tokens exceeds the key volume.

187

167

168

170

171

172 173

174

175

176

177

178

188 189

4 CURSE OF KEY COLLISION ON DISTORTION-FREE WATERMARKS

190 We start with showing key collision is inevitable. In the previous section, we show that given a fixed 191 secret key sk_0 , the watermark key space is finite. Consequently, key collisions will occur with a 192 sufficient number of queries to the language model. One might naturally question whether using an 193 infinite number of secret keys (e.g., a unique key for each generation) could expand the watermark key space to infinity, thereby reducing the likelihood of collisions. However, this approach is impractical 194 because it would substantially reduce detection efficiency. When analyzing a watermarked sequence, 195 the detection algorithm would need to be applied to all possible secret keys, even though only one key 196 corresponds to the watermark. Thus, the watermark information becomes obscured by the numerous 197 other keys. All missing proofs can be found in Appendix D.

Lemma 4.1 (Detection efficiency with multiple secret keys). Denote by $S(\cdot|\mathsf{sk})$ the test statistic. Under the null hypothesis H_0 , given a random text $\mathbf{x}_{1:n}$, we have $\Pr(S(\mathbf{x}_{1:n}|\mathsf{sk}_0) - \mathbb{E}_{H_0}[S] \ge t|H_0) = p_0(t)$, i.e., $p_0(t)$ is the false positive rate of threshold t under single secret key detection. Given M different secret keys, if we use the maximum of the score as the test statistic, we have

203 204

$$\Pr\left(\max_{i\in[M]} (S(\boldsymbol{x}_{1:n}|\mathbf{S}\mathbf{k}_i) - \mathbb{E}_{H_0}[S]) \ge t|H_0\right) = 1 - (1 - p_0(t))^M, \quad \forall t \in \mathbb{R}$$

205 206

Corollary 4.2. Under the existing watermark key sampling schemes, key collision is inevitable.

Lemma 4.1 states that, given the same threshold t, the false positive rate increases with the number of secret keys. Especially, when $M \to \infty$, the false positive rate will tend to 1, which indicates every sentence will be detected as watermarked. Thus, the number of secret keys should be finite, and key collision is inevitable.

We then provide the definition of the three levels of distortion-free capabilities in watermarks: 1) distortion-free within a single token generation, 2) distortion-free in one entire generation, 3) distortion-free across multiple generations.

Definition 4.3 (Step-wise distortion-free watermark). *If a watermark framework adopts a distortion-free PDA-rule, then it is a step-wise distortion-free watermark.*

Definition 4.4 (Weakly distortion-free watermark). A step-wise distortion-free watermark P_W is weakly distortion-free, if $\forall n \in \mathbb{N}_+, \forall \mathbf{x}_{1:n} \in \mathcal{V}$, we have $P_M(\mathbf{x}_{1:n}) = \mathbb{E}_{\mathbf{k}_{1:n}}[P_W(\mathbf{x}_{1:n}|\mathbf{k}_{1:n})]$.

Definition 4.5 (Strongly distortion-free watermark). A step-wise distortion-free watermark P_W is strongly distortion-free if for arbitrary number of generation N_0 and $\forall \boldsymbol{x}_{1:n}^{(i)} \in \mathcal{V}, i \in [N_0]$, it holds that $\prod_{i=1}^{N_0} P_M(\boldsymbol{x}_{1:n}^{(i)}) = \mathbb{E}_{\boldsymbol{k}_{1:n}^{(1)}, \dots, \boldsymbol{k}_{1:n}^{(N_0)}} [\prod_{i=1}^{N_0} P_W(\boldsymbol{x}_{1:n}^{(i)} | \boldsymbol{k}_{1:n}^{(i)})].$

In the next theorem, we show the sufficient conditions for achieving a weakly/strongly distortion-free watermark.

Theorem 4.6. A watermark framework is a weakly/strongly distortion-free watermark if a) it adopts
 a distortion-free PDA-rule and b) there is no key collision during watermark generation.

Corollary 4.7. A watermark that consists of a distortion-free PDA-rule with n-gram hashing, position hashing or fixed key set is a weakly distortion-free watermark.

The proof of this corollary is straightforward because all these watermark key samplers guarantee the uniqueness of each watermark key in a single generation. However, across multiple generations, key collisions become inevitable as the number of generated tokens can surpass the volume of available keys. In the rest of this section, we will explain how key collisions can impact the generation quality and lead to a biased watermarked distribution compared to the original language model distribution.

235

237

222

236 4.1 EXISTING DISTORTION-FREE PDA-RULES

To analyze the influence of key collision on the distortion-free watermarks, we begin with introducing the existing PDA-rules. We also provide a detailed illustration of the existing PDA-rules in Figure 1.

Gumbel-reparametrization. In the Gumbel-reparametrization rule, when sampling x_i with the watermark key k_i , we first sample Gumbel pseudo-random variables $g_1(k_i), ..., g_N(k_i) \sim Gumbel(0, 1)$ with the watermark key k_i . These N independent Gumbel random variables are added to the log-probability of tokens $\log P_M(t_1|\mathbf{x}_{1:i-1}), ..., \log P_M(t_N|\mathbf{x}_{1:i-1})$. The token that achieves the maximum value is then selected as the next token x_i . This process can be formulated through the following equation: $F_{GR}(P_M(\cdot|\mathbf{x}_{1:i-1}), k_i) = \delta_{t_m*}$, where $m^* = \arg \max_{m \in [N]} (g_m(k_i) + \log P_M(t_m|\mathbf{x}_{1:i-1}))$ and δ is the Dirac function.

1nverse-sampling. In the inverse-sampling rule, when sampling x_i with the watermark key k_i , we first organize the LM token probability $P_M(t_1|\mathbf{x}_{1:i-1}), ..., P_M(t_N|\mathbf{x}_{1:i-1})$ within the interval [0, 1]. Then we will sample a pseudo-random variable $r(k_i) \in U(0, 1)$, where U(0, 1) is the uniform distribution on [0, 1]. The next token is selected based on the location of $r(k_i)$ within the cumulative probability intervals on [0, 1]. This process can be formulated through the following equation: $F_{IS}(P_M(\cdot|\mathbf{x}_{1:i-1}), k_i) = \delta_{t_m^* \in V}$, where $r(k_i) \in [\sum_{j=1}^{m^*-1} P_M(t_j|\mathbf{x}_{1:i-1}), \sum_{j=1}^{m^*} P_M(t_j|\mathbf{x}_{1:i-1})]$ and δ is the Dirac function.

254 **Permute-reweight.** In the permute-reweight rule, when sampling x_i with the watermark key k_i , 255 we first generate a pseudo-random token permutation $\pi(\cdot|k_i): V \to [N]$, which is a bijection 256 between token set V and [N]. The token permutations are uniformly distributed with the watermark 257 keys. The LM token probabilities are then rearranged within the interval [0, 1] according to the 258 permutation $\pi(\cdot|k_i)$. The token probability within [0, 1/2] will be scaled to 0, and the rest half will be scaled to 1. Subsequently, x_i is randomly sampled following this adjusted distribution. We can 259 formulate the permute-reweight rule through the following formula: $F_{PR}(P_M(\cdot|\boldsymbol{x}_{1:i-1}), k_i)(t) = \max\{2\sum_{t', \pi(t'|k_i) \le \pi(t|k_i)} P_M(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} - \max\{2\sum_{t', \pi(t'|k_i) \le \pi(t|k_i) - 1} P_M(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\}$ 260 261 1, 0. 262

Pseudo- vs True- Randomness. Based on the above discussion, it is clear that token sampling using Gumbel-reparametrization or inverse-sampling relies entirely on pseudo-randomness, as the watermark distribution for these methods is deterministic given the watermark key. Consequently, for the same token distribution, key collisions result in identical token generation. For instance, when
 generating multiple responses with the same prompt, the first token will always be identical. In contrast, token sampling with the permute-reweight rule does not fully depend on pseudo-randomness. The permute-reweight PDA-rule only scales the first half of the distribution to zero, preserving the rest of the token probabilities. True-random sampling is then applied to the remaining tokens.



Figure 1: Pseudo-randomness in a token sampling step for watermarked LMs. "Before" refers the original LM token distribution and "After" refers the watermarked token distribution. Given a fixed watermark key, both inverse-sampling and Gumbel reparametrization methods become deterministic. In contrast, the permute-reweight method retains elements of randomness.

4.2 DISTRIBUTION BIAS OF DISTORTION-FREE WATERMARKS UNDER KEY COLLISIONS

In this subsection, we explore the distribution bias introduced by the watermark. Given that the distribution overlap between two distributions $P_1, P_2 \in \mathcal{P}$ is represented by $\sum_{t \in V} \min\{P_1(t), P_2(t)\}$, we use $1 - \sum_{t \in V} \min\{P_1(t), P_2(t)\}$, i.e., the total variation, to measure the distribution bias between P_1 and P_2 . Under the key collisions, the bias introduced by a PDA-rule F on a token distribution $P \in \mathcal{P}$ is $1 - \sum_{t \in V} \min\{P(t), F(P|k)(t)\}$. Thus, we introduce the *expected total variation* as a metric for measuring distribution bias.

Definition 4.8 (Expected total variation). *Given a token distributions* $P \in \mathcal{P}$ *and a PDA-rule* F, *the expected total variation between them is given by* $\mathbb{D}(P, F) := 1 - \mathbb{E}_k[\sum_{t \in V} \min\{P(t), F(P|k)(t)\}].$

Trade-off between watermark strength and distribution bias under key collisions. Interestingly, the expected total variation also reflects the watermark's strength. In statistical watermarking, where the goal is to embed a statistical signal into generated content, a larger total variation enhances the strength of this signal and improve the detection efficiency. However, under key collisions, it is desirable for the expected total variation to be minimized to better preserve the original LM distribution. Therefore, a trade-off exists between watermark strength and distribution bias under key collisions.

We compute the expected distribution bias of the existing distortion-free PDA-rules: Gumbelreparametrization F_{GR} , inverse-sampling F_{IS} , and permute-reweight F_{PR} .

Theorem 4.9. Given an arbitrary token distribution $P \in \mathcal{P}$, we have

$$\mathbb{D}(P, F_{GR}) = \mathbb{D}(P, F_{IS}) = 1 - \sum_{t \in V} P(t)^2,$$

and

313 314 315

316 317

$$0.5(1 - \max_{t \in V} P(t)) \le \mathbb{D}(P, F_{PR}) \le 0.5 - \max\{\max_{t \in V} P(t) - 0.5, 0\}.$$

318 Moreover,
$$\mathbb{D}(P, F_{PR}) \leq \mathbb{D}(P, F_{IS}) = \mathbb{D}(P, F_{GR}).$$

From this theorem, we find that the permute-reweight watermark exhibits a smaller distribution
 bias compared to the Gumbel-reparametrization and inverse-sampling watermarks. This finding
 aligns with our analysis in Section 4.1, where we assert that Gumbel-reparametrization and inverse sampling become deterministic with a fixed watermark key, while permute-reweight maintains an
 element of randomness, resulting in a smaller distribution bias. In the next theorem, we will show

289 290 291

292

293

287

324 that under key collisions, a watermark with a PDA-rule F is strongly distortion-free if and only if 325 $\mathbb{D}(P, F) = 0, \forall P \in \mathcal{P}$, which indicates that no signal can be embedded into the generated content. 326

Theorem 4.10. Under key collisions, a watermark with a distortion-free PDA-rule F is strongly 327 distortion-free if and only if $\forall P \in \mathcal{P}$, $\mathbb{D}(P, F) = 0$. 328

By integrating Theorem 4.10 with Theorem 4.9, we find that F_{GR} , F_{IS} , and F_{PR} are unable to yield 330 a strongly distortion-free watermark when key collisions occur. Thus, all existing distortion-free 331 watermarks (Aaronson, 2022; Christ et al., 2023; Kuditipudi et al., 2023; Hu et al., 2023) are not 332 strongly distortion-free. Following the above discussion, we summarize the characteristics of existing distortion-free watermarks in Table 1. 333

334 **Corollary 4.11.** Under key collisions, a strongly distortion-free watermark does not exist.

If $\forall P \in \mathcal{P}, \mathbb{D}(P, F) = 0$, the watermarked LM shows no distribution bias towards the original LM 336 under the watermark key, i.e., $\forall k \in K, F(P|k) = P$. In this case, no watermark is added to the 337 generated content. As key collision is inevitable, we can conclude that with the current watermark 338 key sampling approaches, a strongly distortion-free watermark does not exist. 339

4.3 A BLACK-BOX DISTORTION-FREE WATERMARK DETECTION APPROACH 341

342 As all watermarking approaches present distribution bias towards the original LM under key collisions, 343 we can naturally design a watermark detection approach for the distortion-free watermarks based on 344 the performance difference between the watermarked and the original LMs. 345

We define a new metric Δ , which measures the performance gap between the watermarked model and 346 the original LM. For n random prompts $p_1, ..., p_n$ with m responses for each $g_1^{p_i}, ..., g_m^{p_i}$, denoted 347 by Met an arbitrary performance metric (e.g., perplexity), P_M the original LM, P_T the test LM, we define $\Delta Met(P_M, P_T) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{m} |\sum_{j=1}^{m} Met(g_j^{p_i}(P_M)) - \sum_{j=1}^{m} Met(g_j^{p_i}(P_T))|$. Our watermark detection statistic is given by DetWmk $(P_M, P_T) := \Delta Met(P_M, P_T) - \Delta Met(P_M, P_{M'})$, where 348 349 350 $P_{M'}$ is identically distributed with P_M 351

Theorem 4.12. If P_T is identically distributed with P_M , and $\forall g, |Met(g)| \leq A$, we have $\forall t > 0$,

$$\Pr(|\text{DetWmk}(P_M, P_T)| \ge t) \le \exp(-\frac{m^2 n t^2}{12A^4})$$
(1)

356 With this concentration bound, we can efficiently detect whether a language model has been wa-357 termarked by examining the performance gap between the test model and the original model. 358 Theorem 4.12 provides a statistical guarantee that if the test model P_T is identically distributed with the original model P_M (i.e., unwatermarked), the probability that our detection statistic 359 $DetWmk(P_M, P_T)$ exceeds a threshold t diminishes exponentially with the number of prompts 360 n and responses m. Specifically, the bound ensures that false positives are highly unlikely when 361 the performance metric is bounded by A. This allows us to confidently and efficiently identify 362 watermarked language models by detecting significant deviations in performance metrics. 363

364 4.4 BEYOND KEY COLLISION - STRONGLY DISTORTION-FREE WATERMARK DOES NOT EXIST

365 366 We extend our analysis on strongly distortion-free watermarks and prove that a detectable, strongly 367 distortion-free watermark does not exist. From the above analysis we know that the independence of 368 PDA-rule is a necessary condition for strongly distortion-free watermark, and the independence of 369 PDA-rule stems from the independence of hashed watermark keys h(k). Thus, we can divide the 370 proof into two parts: 1) A strongly distortion-free watermark must use independent hashed watermark 371 keys, denoted as h(k), where h is the hash function employed by the PDA-rule. 2) A watermark 372 using a distortion-free and independent PDA-rule is undetectable by arbitrary detector. Combining 1) and 2) we have the following theorem: 373

374 **Theorem 4.13.** A detectable strongly distortion-free watermark does not exist.

375

335

340

352 353

354 355

Theorem 4.13 establishes a fundamental limitation in the design of watermarking schemes by stating 376 that a detectable, strongly distortion-free watermark cannot exist. This theorem also highlights the 377 trade-off in watermarking systems between distribution bias and watermark strength (see Sec. 4.2). If a watermark is designed to be unbiased to the original data distribution (strongly distortion-free), it
 cannot be reliably detected using standard detection methods. Conversely, introducing detectability
 requires some form of alteration or pattern that can be recognized, which compromises the strongly
 distortion-free property.

382

384

5 REDUCING DISTRIBUTION BIAS VIA BETA-WATERMARK

In this section, we introduce a new family of watermarks, beta-watermark, which trades watermark strength for low distribution bias. The betawatermark is based on a distortion-free beta PDArule and n-gram hashing. Additionally, we present a model-agnostic detection method for it. In Appendix A Alg. 1 and 2 we show the algorithms of the generator and detector of beta-watermark.

The beta PDA-rule is a variation of the permutereweight PDA-rule (another example is DiPmark (Wu et al., 2023)) that introduces greater true randomness during sampling. Similar to permute-reweight watermark, When sampling x_i with the watermark key k_i ,



Figure 2: Illustration of Beta PDA-rule.

we first generate a pseudo-random token permutation $\pi(\cdot|k_i): V \to [N]$. Then we arrange the LM token probability within the interval [0, 1] following the permutation $\pi(\cdot|k_i)$. The first half of token probability (token probability within [0, 1/2]) will be scaled to β , and the rest half probability will be scaled to $1 - \beta$ (See Figure 2 for a detailed illustration). The next token is randomly sampled from the new distribution. Notice, when $\beta = 0$, the permute-reweight PDA-rule is applied and when $\beta = 0.5$, the original LM distribution is used.

Definition 5.1 (Beta PDA-rule). *Beta PDA-rule* F_{β} *is defined by:* $F_{\beta}(P_M(\cdot|\boldsymbol{x}_{1:i-1}), k_i)(t) = (1 - \beta)F_{PR}(P_M(\cdot|\boldsymbol{x}_{1:i-1}), k_i)(t) + \beta[\max\{2\sum_{t', \pi(t'|k_i) \ge \pi(t|k_i)} P_M(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\}] - \max\{2\sum_{t', \pi(t'|k_i) \ge \pi(t|k_i) + 1} P_M(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\}].$ Notice, the range of β is from 0 to 0.5.

407 **Theorem 5.2.** Beta PDA-rule is a distortion-free PDA-rule, i.e., $\forall \mathbf{x}_{1:n} \in \mathcal{V}, \forall i \leq n$, 408 $P_M(x_i|\mathbf{x}_{1:i-1}) = \mathbb{E}_{k_i}[F(P_M(\cdot|\mathbf{x}_{1:i-1}), k_i)(x_i)].$

409 410 **Corollary 5.3.** *Beta-watermark is a weakly distortion-free watermark.*

The proof is straightforward, as the beta-watermark consists of a distortion-free PDA-rule and the
 n-gram hashing. In the subsequent theorem, we theoretically demonstrate that the beta PDA-rule
 introduces a smaller distribution bias compared to the permute-reweight watermark.

414 **Theorem 5.4.** Given an arbitrary token distribution $P \in \mathcal{P}$, $\mathbb{D}(P, F_{\beta}) \leq \mathbb{D}(P, F_{PR}) - \beta(1 - \max_{t \in V} P(t))$. Besides, if $\beta_1 < \beta_2$, $\mathbb{D}(P, F_{\beta_1}) > \mathbb{D}(P, F_{\beta_2})$. 416

As the detector of the permute-reweight watermark (Hu et al., 2023) is dependent on the logits from the original LM, we design a new model-agnostic detection algorithm for the beta-watermark. As shown in Figure 2, beta-reweighting tends to enhance the token probability towards the end of the permutation. During detection, given an input token, we can determine its position within the permutation using $\pi(x|k)$. Thus, a higher score should be assigned to larger values of $\pi(x|k)$. We use a sigmoid function: sigmoid($C(\pi(x|k)/|V| - 0.5)$), where C is a scaling parameter, to appropriately scale the scores.

Definition 5.5 (Model-agnostic beta-reweight detection). We use score function $s(x,k,F) = sigmoid(C(\pi(x|k)/|V| - 0.5))$ to conduct detection. Given a random observation $\mathbf{x}_{1:n}$, under the null hypothesis, we have $\Pr(S(\mathbf{x}_{1:n}) - \mathbb{E}_{H_0}[S(\mathbf{x}_{1:n})] > t\sqrt{n}|H_0) \le \exp(-2t^2)$.

427

6 EXPERIMENTS

428 429

Our experimental section consists of two parts. In the first part, we compare the weakly and strongly
 distortion-free nature of the beta watermark with that of existing watermarks, and validate the
 trade-off between the watermark strength and distribution bias. In the second part, we evaluate



Figure 3: Performance of different watermarks under one-time generation. **Top:** Violin plot of Text Summarization Perplexity. **Bottom:** Violin plot of Machine Translation BLEU. We can see the weakly distortion-free watermarks preserve the generation quality.

Table 2: Performance of different watermarks under multi-time generations. We randomly selected 1000 prompts and generated 100 responses for each. The baseline is the Δ metrics between two no-watermarked models.

	Te	xt Summarization	l	Machine Trai	nslation
	Δ BERT Score \downarrow	Δ ROUGE-1 \downarrow	Δ Perplexity \downarrow	Δ BERT Score \downarrow	Δ BLEU \downarrow
Baseline	0.0062	0.0070	0.3028	0.0180	0.7716
Beta-Reweight ($\beta = 0$)	0.0090	0.0093	0.3753	0.0267	1.2373
Beta-Reweight ($\beta = 0.05$)	0.0084	0.0085	0.3549	0.0248	1.1806
Beta-Reweight ($\beta = 0.1$)	0.0079	0.0081	0.3453	0.0230	1.0316
Beta-Reweight ($\beta = 0.2$)	0.0070	0.0077	0.3368	0.0203	0.9475
Beta-Reweight ($\beta = 0.3$)	0.0066	0.0073	0.3144	0.0195	0.8638
Inverse-sampling	0.0446	0.0494	1.7846	0.1316	5.5354
Gumbel-reparametrization	0.0428	0.0488	1.8892	0.1341	5.6438
$\text{Soft}(\delta = 1.0)$	0.0091	0.0099	0.5473	0.0428	1.4660
$\text{Soft}(\delta = 1.5)$	0.0128	0.0136	1.1237	0.0808	2.5310
$\text{Soft}(\delta = 2.0)$	0.0195	0.0194	2.0817	0.1274	3.7758

the detection efficiency of the beta watermark against existing watermarks. In the third part, we assess the robustness of the beta watermark when subjected to random paraphrasing attacks. We focus on three seq2seq tasks in our experiments: machine translation, text summarization and text generation. Detailed experimental settings are provided in Appendix E and additional experimental results, including the detectability and the robustness of beta-watermark, are in Appendix F.

6.1 WEAKLY AND STRONGLY DISTORTION-FREENESS

In this section, we conduct experiments to validate our theoretical analysis. We evaluate the weakly and strongly distortion-free properties of existing watermark strategies as defined in Def-initions 4.4 and 4.5. We validate the weakly distortion-free property by assessing the qual-ity of the watermarked text generated once for each prompt. For the strongly distortion-free property, we examine the quality of the watermarked text for 1000 prompts, where for each prompt we have 100 generations. We define a new metric Δ , which measures the performance gap between the watermarked model and the original LM. For n prompts $p_1, ..., p_n$ with m re-sponses for each $g_{1}^{p_i}, ..., g_m^{p_i}$, denoted by Met an arbitrary performance metric (e.g., perplexity), $\Delta \text{Met} = \frac{1}{n} \sum_{i=1}^n \frac{1}{m} |\sum_{j=1}^m \text{Met}(g_j^{p_i}(\text{No watermark})) - \sum_{j=1}^m \text{Met}(g_j^{p_i}(\text{Watermarked}))|$

Weakly Distortion-Free. The results are presented in Figure 3. This figure shows that compared to the model without watermarks, all weakly distortion-free watermarks exhibit no significant performance bias in text summarization and text generation tasks. However, for the Soft-watermark (Kirchenbauer et al., 2023), a significant performance bias is observable as δ increases.

Strongly Distortion-Free. The results are displayed in Table 2. From this table, it is evident that compared to the baseline, which is the Δ metrics between two non-watermarked models, all weakly distortion-free watermarks demonstrate performance bias across all tasks. In contrast, the Beta-watermark exhibits less bias compared to other weakly distortion-free watermarks. Additionally, as β increases, the distribution bias is further reduced, consistent with our theoretical analysis.



Figure 4: Trade-off between distribution bias and watermark strength under key collision. The TPR is measured under 1% (Left), 0.1% (Right) FPR. We can see Δ Perplexity (distribution bias) increase with the TPR.

Table 3: p-value of our black-box distortion-free watermark detection algorithm on text summarization and machine translation tasks. TS: Text Summarization; MT: Machine Translation; IS: Inversesampling; GR: Gumble-reparametrization; PR: Permute-reweight. The definition of DetWmk is shown in Sec. 4.3.

		IS	GR	PR	Beta watermark (β)		hark (β)	
					0.05	0.1	0.2	0.3
TS	DetWmk	1.4818	1.5864	0.0725	0.0521	0.0425	0.0340	0.0116
	p-value	0	0	0.0125	0.1041	0.2219	0.3816	0.8939
MT	DetWmk	4.7638	4.8722	0.4597	0.4090	0.2600	0.1795	0.0992
	p-value	0	0	1.6594e-05	1.6450e-04	0.0295	0.1867	0.5989

Trade-off between watermark strength and distribution bias. We use the beta-watermark to empirically verify the trade-off between watermark strength and distribution bias. As shown in Figure 4, with increasing values of β , the distribution bias decreases, but there is also a corresponding decrease in the true positive rate of watermark detection. This indicates that reducing the distribution bias of the watermark compromises its detectability.

6.2 BLACK-BOX DISTORTION-FREE WATERMARK DETECTION

In this section, we present experimental results to validate our proposed black-box distortion-free watermark detection method (Theorem 4.12). We evaluate the performance on text summarization and machine translation tasks using perplexity and BLEU as the metrics, respectively. To ensure these metrics are bounded, we clip the perplexity to the interval [0,10] and the BLEU score to the interval [0,20]. We report the p-value calculated according to Eq. (1).

From Table 3, we observe that, under a 5% false positive rate (FPR), our detection method successfully
 identifies inverse-sampling, Gumbel-reparametrization, and permute-reweight watermarks for both
 text summarization and machine translation tasks. However, the beta-watermark is able to significantly
 reduce the detection accuracy.

7 CONCLUSION

In conclusion, this work identifies three levels of distortion-free capabilities in watermarks—Step-wise, Weakly, and Strongly Distortion-free—and demonstrates that existing watermarks are not strongly distortion-free due to key collisions, which disrupt the original language model distribution across multiple generations. We theoretically establish a trade-off between watermark strength and distribution bias, and introduce a black-box detection approach for identifying watermarked models. Additionally, we prove that strongly distortion-free watermarks are theoretically unattainable. As a practical solution, we propose beta-watermark, a new weakly distortion-free watermark that effectively reduces distribution bias at the cost of watermarking strength. Future research direction includes 1) exploring further details of the trade-off between the distribution bias and the watermarking strengh and 2) developing more efficient watermark detection methods for weakly distortion-free watermarks.

540 REFERENCES

559

560

561

562

566

567

568

569

578

579

580 581

582

583

- Scott Aaronson. My AI safety lecture for UT effective altruism, 2022. URL https://
 scottaaronson.blog/?p=6823.
- Souradip Chakraborty, Amrit Singh Bedi, Sicheng Zhu, Bang An, Dinesh Manocha, and Furong Huang. On the possibilities of AI-generated text detection. *arXiv preprint arXiv:2304.04736*, 2023.
- 548 Miranda Christ and Sam Gunn. Pseudorandom error-correcting codes. In *Annual International* 549 *Cryptology Conference*, pp. 325–347. Springer, 2024.
- Miranda Christ, Sam Gunn, and Or Zamir. Undetectable watermarks for language models. *arXiv preprint arXiv:2306.09194*, 2023.
- Nenad Dedić, Gene Itkis, Leonid Reyzin, and Scott Russell. Upper and lower bounds on black-box
 steganography. *Journal of Cryptology*, 22:365–394, 2009.
- Jiayi Fu, Xuandong Zhao, Ruihan Yang, Yuansen Zhang, Jiangjie Chen, and Yanghua Xiao. Gumbelsoft: Diversified language model watermarking via the gumbelmax-trick. *arXiv preprint arXiv:2402.12948*, 2024.
 - Margherita Gambini, Tiziano Fagni, Fabrizio Falchi, and Maurizio Tesconi. On pushing DeepFake Tweet detection capabilities to the limits. In *Proceedings of the 14th ACM Web Science Conference 2022*, pp. 154–163, 2022.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa
 Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28, 2015.
 - Nicholas J Hopper, John Langford, and Luis Von Ahn. 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.
- Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu, Hongyang Zhang, and Heng Huang. Unbiased
 watermark for large language models. *preprint*, 2023.
- Gabriel Kaptchuk, Tushar M Jois, Matthew Green, and Aviel D Rubin. Meteor: Cryptographically secure steganography for realistic distributions. In *Proceedings of the 2021 ACM SIGSAC Conference* on Computer and Communications Security, pp. 1529–1548, 2021.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A
 watermark for large language models. *arXiv preprint arXiv:2301.10226*, 2023.
 - Jan Hendrik Kirchner, Lama Ahmad, Scott Aaronson, and Jan Leike. New AI classifier for indicating AI-written text. *OpenAI*, 2023.
 - Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense. *arXiv preprint arXiv:2303.13408*, 2023.
- Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. Robust distortion-free
 watermarks for language models. *arXiv preprint arXiv:2307.15593*, 2023.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Aiwei Liu, Leyi Pan, Xuming Hu, Shu'ang Li, Lijie Wen, Irwin King, and Philip S Yu. An unforgeable
 publicly verifiable watermark for large language models. *arXiv preprint arXiv:2307.16230*, 2023a.
- 593 Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and Lijie Wen. A semantic invariant robust watermark for large language models. *arXiv preprint arXiv:2310.06356*, 2023b.

594 595 596	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742, 2020.
597 598 599 600	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. Detectgpt: Zero-shot machine-generated text detection using probability curvature. <i>arXiv preprint arXiv:2301.11305</i> , 2023.
601 602 603	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pp. 311–318, 2002.
604 605 606 607	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67, 2020.
608 609 610	Yi Tay, Dara Bahri, Che Zheng, Clifford Brunk, Donald Metzler, and Andrew Tomkins. Reverse engineering configurations of neural text generation models. <i>arXiv preprint arXiv:2004.06201</i> , 2020.
611 612	Edward Tian. GPTzero update v1. https://gptzero.substack.com/p/ gptzero-update-v1, 2023.
613 614 615	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
616 617 618	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. <i>arXiv preprint arXiv:1910.03771</i> , 2019.
620 621	Yihan Wu, Zhengmian Hu, Hongyang Zhang, and Heng Huang. Dipmark: A stealthy, efficient and resilient watermark for large language models. <i>arXiv preprint arXiv:2310.07710</i> , 2023.
622 623 624	Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending against neural fake news. <i>Advances in neural information processing systems</i> , 32, 2019.
625 626 627	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> , 2019.
628 629 630	Xuandong Zhao, Prabhanjan Ananth, Lei Li, and Yu-Xiang Wang. Provable robust watermarking for ai-generated text. <i>arXiv preprint arXiv:2306.17439</i> , 2023.
631 632	
633	
634	
635	
636	
637	
638	
639	
640	
04 I 640	
6/3	
644	
645	
646	
647	

A ALGORITHMS OF BETA-WATERMARK

1: Input: secret key SK, parameter β , prompt $\boldsymbol{x}_{-m:0}$, generate length $n \in \mathbb{N}$, texture key histor hist, n-gram parameter a , and permutation generation function h . 2: for $i = 1,, n$ do 3: Calculate the LM distribution for generating the i -th token $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 4: Generate a watermark key $k_i = (Sk, \boldsymbol{x}_{i-a,i-1})$. 5: if $k_i \in hist$ then 6: Sample the next token x_i using original LM distribution $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_{\beta}(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key Sk, volume of the token set N , score function s , n-gram parameter a , threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (Sk, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the score function via $S = S + s(\pi(x_i k_i), k_i, F_{\beta})$. 7: return $S > z\sqrt{n}$.		orithm I Beta-watermark generator
<i>hist</i> , <i>n</i> -grain parameter <i>a</i> , and permutation generation function <i>h</i> . 2: for <i>i</i> = 1,, <i>n</i> do 3: Calculate the LM distribution for generating the <i>i</i> -th token $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 4: Generate a watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: if $k_i \in hist$ then 6: Sample the next token x_i using original LM distribution $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_{\beta}(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key Sk, volume of the token set N , score function s , n-gram parameter a , threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_{\beta})$. 7: return $S > z\sqrt{n}$.	1:	Input: secret key SK, parameter β , prompt $x_{-m:0}$, generate length $n \in \mathbb{N}$, texture key history
2: for $i = 1,, n$ do 3: Calculate the LM distribution for generating the <i>i</i> -th token $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 4: Generate a watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: if $k_i \in hist$ then 6: Sample the next token x_i using original LM distribution $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_\beta(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key sk , volume of the token set <i>N</i> , score function <i>s</i> , n-gram parameter <i>a</i> , threshold <i>z</i> . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z_\sqrt{n}$.	2.	nist, n-gram parameter a , and permutation generation function h .
Since Calculate the LM distribution for generating the <i>i</i> -th token $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 4: Generate a watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: if $k_i \in hist$ then 6: Sample the next token x_i using original LM distribution $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_\beta(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key sk , volume of the token set <i>N</i> , score function <i>s</i> , n-gram parameter <i>a</i> , threshold <i>z</i> . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z\sqrt{n}$.	2:	10 $i = 1,, n$ do
4: Generate a Watermark key $k_i = (SK, x_{i-a,i-1})$. 5: if $k_i \in hist$ then 6: Sample the next token x_i using original LM distribution $P_M(\cdot \boldsymbol{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_\beta(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key Sk, volume of the token set N , score function s , n-gram parameter a , threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (SK, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z\sqrt{n}$.	3:	Calculate the LNI distribution for generating the <i>i</i> -th token $P_M(\cdot \mid \boldsymbol{x}_{-m:i-1})$.
5: If $\mathbf{k}_i \in mst$ then 6: Sample the next token x_i using original LM distribution $P_M(\cdot \mathbf{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \mathbf{x}_{-m:i-1}) = F_\beta(P_M(\cdot \mathbf{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \mathbf{x}_{-m:i-1})$. 11: return $\mathbf{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\mathbf{x}_{1:n}$, secret key Sk, volume of the token set N , score function s , n-gram parameter a, threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (Sk, \mathbf{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z_\sqrt{n}$.	4:	Generate a watermark key $\kappa_i = (SK, x_{i-a,i-1})$.
6: Sample the next token x_i using original LW distribution $F_M(\cdot \boldsymbol{x}_{-m:i-1})$. 7: else 8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_\beta(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key Sk, volume of the token set N , score function s , n-gram parameter a , threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (Sk, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z_{\sqrt{n}}$.	5:	If $\kappa_i \in Hist$ then Some lot the part taken κ_i using original I M distribution D_i ($ m_i \rangle$)
8: Generate the permutation of token set $\pi(\cdot k_i)$. 9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_\beta(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i)$. 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1})$. 11: return $\boldsymbol{x}_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key Sk, volume of the token set N , score function s , n-gram parameter a , threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (\text{Sk}, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z_{\sqrt{n}}$.	0: 7.	sample the next token x_i using original LW distribution $F_M(\cdot \mid x_{-m:i-1})$.
9: Calculate watermarked distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = F_\beta(P_M(\cdot \boldsymbol{x}_{-m:i-1}), k_i).$ 10: Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}).$ 11: return $\boldsymbol{x}_{1:n}.$ Algorithm 2 Beta-watermark detector 1: Input: text $\boldsymbol{x}_{1:n}$, secret key sk, volume of the token set N , score function s , n-gram parameter a , threshold z . 2: Initialize the score function: $S = 0.$ 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (\text{sk}, \boldsymbol{x}_{i-a,i-1}).$ 5: Generate the permutation of token set $\pi(\cdot k_i).$ 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta).$ 7: return $S > z_{\sqrt{n}}.$	/. o.	Concrete the permutation of taken set $\pi(k)$
 9. Calculate watermarked distribution T_W(¬ x=m:i=1) = T_β(T_M(¬ x=m:i=1), k_i). 10: Sample the next token x_i using distribution P_W(¬ x=m:i=1). 11: return x_{1:n}. Algorithm 2 Beta-watermark detector Input: text x_{1:n}, secret key sk, volume of the token set N, score function s, n-gram parameter a, threshold z. Initialize the score function: S = 0. 3: for i = 2,, n do 4: Generate the watermark key k_i = (sk, x_{i=a,i=1}). 5: Generate the permutation of token set π(· k_i). 6: Update the score function via S = S + s(π(x_i k_i), k_i, F_β). 	0. 0.	Calculate watermarked distribution $P_{\rm ext}(\mathbf{x}_i , \cdot, \cdot) = F_0(P_{\rm ext}(\mathbf{x}_i , \cdot, \cdot)) k_i)$
11: return $x_{1:n}$. Algorithm 2 Beta-watermark detector 1: Input: text $x_{1:n}$, secret key sk, volume of the token set N , score function s , n-gram parameter a, threshold z . 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (Sk, x_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_{\beta})$. 7: return $S > z\sqrt{n}$.	9. 10:	Sample the next token x_i using distribution $P_W(\cdot \boldsymbol{x}_{-m:i-1}) = P_\beta(T_M(\cdot \boldsymbol{x}_{-m:i-1}),\kappa_i)$.
 Algorithm 2 Beta-watermark detector 1: Input: text x_{1:n}, secret key Sk, volume of the token set N, score function s, n-gram paramete a, threshold z. 2: Initialize the score function: S = 0. 3: for i = 2,, n do 4: Generate the watermark key k_i = (Sk, x_{i-a,i-1}). 5: Generate the permutation of token set π(· k_i). 6: Update the score function via S = S + s(π(x_i k_i), k_i, F_β). 7: return S > z√n. 	11:	return $x_{1:n}$.
 Algorithm 2 Beta-watermark detector 1: Input: text x_{1:n}, secret key \$\mathbf{s}\$, volume of the token set N, score function s, n-gram parameter a, threshold z. 2: Initialize the score function: S = 0. 3: for i = 2,, n do 4: Generate the watermark key k_i = (\$\mathbf{s}\$, x_{i-a,i-1}). 5: Generate the permutation of token set π(· k_i). 6: Update the score function via S = S + s(π(x_i k_i), k_i, F_β). 7: return S > z√n. 		
 Input: text x_{1:n}, secret key Sk, volume of the token set N, score function s, n-gram parameter a, threshold z. Initialize the score function: S = 0. for i = 2,, n do Generate the watermark key k_i = (Sk, x_{i-a,i-1}). Generate the permutation of token set π(· k_i). Update the score function via S = S + s(π(x_i k_i), k_i, F_β). return S > z√n. 	Alg	orithm 2 Beta-watermark detector
a, threshold z. 2: Initialize the score function: $S = 0$. 3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (Sk, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z\sqrt{n}$.	1:	Input: text $x_{1:n}$, secret key Sk, volume of the token set N, score function s, n-gram parameter
 2: Initialize the score function: S = 0. 3: for i = 2,, n do 4: Generate the watermark key k_i = (sk, x_{i-a,i-1}). 5: Generate the permutation of token set π(· k_i). 6: Update the score function via S = S + s(π(x_i k_i), k_i, F_β). 7: return S > z√n. 		a, threshold z.
3: for $i = 2,, n$ do 4: Generate the watermark key $k_i = (\mathbf{sk}, \boldsymbol{x}_{i-a,i-1})$. 5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_{\beta})$. 7: return $S > z\sqrt{n}$.	2:	Initialize the score function: $S = 0$.
 Generate the watermark key k_i = (sk, x_{i-a,i-1}). Generate the permutation of token set π(· k_i). Update the score function via S = S + s(π(x_i k_i), k_i, F_β). return S > z√n. 	3:	for $i = 2,, n$ do
5: Generate the permutation of token set $\pi(\cdot k_i)$. 6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z\sqrt{n}$.	4:	Generate the watermark key $k_i = (sk, \boldsymbol{x}_{i-a,i-1})$.
6: Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$. 7: return $S > z\sqrt{n}$.	5:	Generate the permutation of token set $\pi(\cdot k_i)$.
7: return $S > z\sqrt{n}$.	6:	Update the score function via $S = S + s(\pi(x_i k_i), k_i, F_\beta)$.

674 675

676 677

648

649

B RELATED WORK

Statistical watermarks. Kirchenbauer et al. (2023) enhanced the statistical watermark framework 678 originally introduced by Aaronson (2022), demonstrating the effectiveness of statistical watermarking 679 through extensive experiments on large language models. They splited the LM tokens into red 680 and green list, then promoted the use of green tokens by adding a fixed parameter δ to their logits. 681 Zhao et al. (2023) proposed the unigram watermark, which enhances the robustness of the statistical 682 watermark by using one-gram hashing to produce watermark keys. Liu et al. (2023b) also improved 683 the robustness of statistical watermarking by leveraging the semantics of generated content as 684 watermark keys. Additionally, Liu et al. (2023a) proposed an unforgeable watermark scheme that 685 employs neural networks to modify token distributions instead of using traditional watermark keys. 686 However, these approaches may lead to significant changes in the distribution of generated text, potentially compromising content quality. 687

688 **Distortion-free watermarks.** To preserve the original output distribution in watermarked content, 689 researchers have explored alternative strategies to modify the token distribution. Aaronson (2022) 690 introduced the first distortion-free watermarking strategy, which utilized Gumbel-reparametrization 691 to alter token distribution and the prefix n-gram content as the watermark keys. Christ et al. (2023) 692 and Kuditipudi et al. (2023) adopted the inverse-sampling and Gumbel-reparametrization to modify 693 the watermarked token distributions, where the watermark keys is based on the token position or a fixed key list respectively. Notice Christ et al. (2023)'s method encounters resilience challenges under 694 modifications and lacks empirical evidence regarding its detectability. Meanwhile, Kuditipudi et al. 695 (2023)'s detection process involves hundreds of resampling steps from the secret key distribution, 696 proving inefficient for processing lengthy texts. Hu et al. (2023) employed inverse-sampling and 697 permute-reweight methods for watermarking. But their detector is not model-agnostic, which requires 698 access to the language model API and prompts, which compromises its operational efficiency. 699

Post-hoc Detectors. Post-hoc detection serves as a significant alternative to watermarking, focusing
 on the retrospective analysis of machine-generated text. This can be achieved by leveraging inherent
 features of language models or by enhancing pre-existing expansive models to function as detectors,

702 as detailed by (Zellers et al., 2019). Specific implementation nuances, such as sampling methods, 703 can be uncovered through reverse engineering the generated text, a process described by (Tay 704 et al., 2020). Additionally, there are post-hoc detectors designed for modern large language models 705 (Mitchell et al., 2023; Tian, 2023; Kirchner et al., 2023), specifically trained for binary detection 706 tasks. However, there is a growing consensus that these detection methods are becoming less effective as language models evolve. As observed by Gambini et al. (2022), detection mechanisms effective against GPT-2 have struggled with GPT-3. Moreover, text rephrasing models like those in (Krishna 708 et al., 2023) are bypassing prevalent post-hoc detectors such as GPTZero (Tian, 2023), DetectGPT 709 (Mitchell et al., 2023), and OpenAI's proprietary detector (Kirchner et al., 2023). Additionally, 710 Chakraborty et al. (2023) notes that as AI-generated content becomes increasingly indistinguishable 711 from human-produced text, the demand on post-hoc detectors to analyze longer text segments will 712 likely increase. 713

Steganography. Steganography involves embedding hidden messages in media such as natural 714 language or images, ensuring only intended recipients can detect the message while it remains 715 concealed from others (Hopper et al., 2002). When applied to watermarking, the goal is to maintain 716 stealth. However, established steganography techniques may not meet this goal without certain 717 entropy-related assumptions. In scenarios where language model prompts can be adversarially 718 chosen, the need for stealth remains. This discrepancy arises due to the different levels of access that 719 watermarking and steganography have to the model's output distribution. In steganography, there 720 is only oracle access to this distribution, whereas our watermarking approach provides a detailed 721 view of the token's probability distribution. Hence, while steganography either depends on entropy 722 assumptions (Hopper et al., 2002) or risks security with low entropy channels (Dedić et al., 2009), 723 our watermark remains stealthy regardless of the text's entropy. This is achieved by leveraging full distribution access and using it as a foundation for embedding watermarks. Kaptchuk et al. (2021) 724 discusses encoding with similar access but presupposes equal decoding access, which is impractical 725 for watermarking as the detection algorithm typically lacks the initiating prompt, thus remaining 726 unaware of the distribution. 727

727 728

C DISCUSSION

729 730 731

In this section, we provide detailed discussion of two "undetectable" scheme (Christ et al., 2023;
Christ & Gunn, 2024). We claim neither of them can achieve strongly distortion-free.

733 For the undetectable scheme proposed by Christ & Gunn (2024), it is important to note that strongly 734 distortion-free watermarks require the independence of $F(P_M(x_i \mid x_{1:i-1}), k_i)$ at every generation 735 step i. Existing distortion-free watermarks achieve this by using distinct watermark keys k_i and a hash 736 function to ensure independence, which requires the 'key collision' not occuring. In contrast, Christ 737 & Gunn (2024) achieves the independence of $F(P_M(x_i \mid x_{1:i-1}), k_i)$ by developing a key sampling 738 method (termed PRC), which aims for i.i.d. sampling of watermark keys with true randomness, i.e., 739 replacing the pseudorandomness in the PDA-rule with true randomness by randomly sampling the watermark keys. However, despite these efforts, their method still does not achieve strongly distortion-740 free watermarks, as PRC methods are only close to, but cannot fully achieve, i.i.d. sampling of true 741 randomness (see Lemma 9 in their paper). Therefore, their method can still not achieve strongly 742 distortion-free watermarks. 743

744 The undetectable watermark (Christ et al., 2023) is another example of trading detectability for reducing the distribution bias. In Christ et al. (2023), watermarked tokens are produced only if the 745 hash window has an entropy larger than a given threshold λ , i.e., they skip watermarking the first 746 several tokens to accumulate enough true randomness. However, there is also a trade-off between 747 the watermark strength and the distribution bias under their scheme. This trade-off is controlled by 748 the entropy threshold λ . When λ increases, the number of watermarked tokens decreases, and it will 749 become more difficult to detect the watermark, but the key collision is less likely to occur, and the 750 distribution bias decreases. 751

For example, if we use Hoeffding's concentration bound as the p-value estimator, i.e., $P(S_n - E[S_n] > s) \le \exp(-2\frac{s^2}{n})$, when the generated sequence length does not change, the p-value upper bound exponentially increases with the number of non-watermarked tokens (because $s = S_n^{watermarked} - E[S_n]$ is linearly related to the number of watermarked tokens). Thus, although the order of hash windows is 2^{λ} , the detectability could be the order $O(e^{-\lambda})$ (assuming λ is linearly

related to the number of watermarked tokens), which show a trade-off between the distortion bias and
 the watermark detectability.

Besides, the key sampling methods Alg. 3 and 5 of Christ et al. (2023) are also not resilient to even single text modification. The watermark key space of undetectable watermark is (sk, texture key, position key) (see line 12 of Alg.3, line 7 of Alg.5 and the discussion at the second last paragraph of Section 4.2). The texture key is similar to the definition of n-gram. Since they also use position keys in their watermark key, a single deletion will remove the watermark.

Notice, the texture key in Christ et al. (2023) is similar but not equal to the n-gram hashing. In Alg.3,
they use the same texture key during one generation. From my perspective, the texture keys in Christ
et al. (2023) are more likely to serve as increasing the diversity of the secret key sk to reduce the
distribution-bias in "multiple generation with the same prompt". The diversity of watermark keys in
one generation is ensured by the position key.

D MISSING PROOFS

D.1 PROOF OF THEOREM 4.1

Proof.

$$\Pr\left(\max_{i\in[M]} (S(\boldsymbol{x}_{1:n}|\mathbf{S}\mathbf{k}_i) - \mathbb{E}_{H_0}[S]) \le t|H_0\right) = \prod_{i=1}^M \Pr\left(S(\boldsymbol{x}_{1:n}|\mathbf{S}\mathbf{k}_i) - \mathbb{E}_{H_0}[S] \le t|H_0\right)$$
$$= \prod_{i=1}^M (1 - \Pr\left(S(\boldsymbol{x}_{1:n}|\mathbf{S}\mathbf{k}_i) - \mathbb{E}_{H_0}[S] \ge t|H_0\right)\right) \quad (2)$$
$$= (1 - p_0(t))^M.$$

Thus,

$$\Pr\left(\max_{i\in[M]}(S(\boldsymbol{x}_{1:n}|\mathbf{s}\mathbf{k}_i) - \mathbb{E}_{H_0}[S]) \ge t|H_0\right) = 1 - \Pr\left(\max_{i\in[M]}(S(\boldsymbol{x}_{1:n}|\mathbf{s}\mathbf{k}_i) - \mathbb{E}_{H_0}[S]) \le t|H_0\right)$$
$$= 1 - (1 - p_0(t))^M.$$
(3)

D.2 PROOF OF THEOREM 4.6

Proof. We first show the weakly distortion-free case: firstly, if key collision does not occur, we have

$$\mathbb{E}_{\boldsymbol{k}_{1:n}}[P_{W}(\boldsymbol{x}_{1:n}|\boldsymbol{k}_{1:n})] = \mathbb{E}_{\boldsymbol{k}_{1:n}}\left[\prod_{i=1}^{n} F(P_{M}(x_{i}|\boldsymbol{x}_{1:i-1}), k_{i})\right]$$

$$= \prod_{i=1}^{n} \mathbb{E}_{k_{i}}[F(P_{M}(x_{i}|\boldsymbol{x}_{1:i-1}), k_{i})].$$
(4)

The above equality stems from the independence property of the PDA-rule $F(P_M(x_i|\boldsymbol{x}_{1:i-1}), k_i)$. Christ et al. (2023) and Hu et al. (2023) show that if there is no repeating keys in $\boldsymbol{k}_{1:n}$, the independence property can be satisfied with hash functions.

Since F is a distortion-free PDA-rule, we have $\mathbb{E}_{k_i}[F(P_M(x_i|\boldsymbol{x}_{1:i-1}),k_i)] = P_M(x_i|\boldsymbol{x}_{1:i-1})$. Thus,

$$\mathbb{E}_{\boldsymbol{k}_{1:n}}[P_W(\boldsymbol{x}_{1:n}|\boldsymbol{k}_{1:n})] = \prod_{i=1}^n \mathbb{E}_{k_i}[F(P_M(x_i|\boldsymbol{x}_{1:i-1}), k_i)] = \prod_{i=1}^n P_M(x_i|\boldsymbol{x}_{1:i-1}) = P_M(\boldsymbol{x}_{1:n}).$$
 (5)

Analogously, for the strongly distortion-free case, if key collision does not occur, we will have distinct $k_{1:n}^{(i)}$. By the independence property of the PDA-rule, we have

$$\mathbb{E}_{\boldsymbol{k}_{1:n}^{(1)},\dots,\boldsymbol{k}_{1:n}^{(N_0)}}[\prod_{i=1}^{N_0} P_W(\boldsymbol{x}_{1:n}^{(i)}|\boldsymbol{k}_{1:n}^{(i)})] = \prod_{i=1}^{N_0} \mathbb{E}_{\boldsymbol{k}_{1:n}^{(i)}}[P_W(\boldsymbol{x}_{1:n}^{(i)}|\boldsymbol{k}_{1:n}^{(i)})] \\ = \prod_{i=1}^{N_0} \prod_{j=1}^{n} \mathbb{E}_{\boldsymbol{k}_j^{(i)}}[P_W(\boldsymbol{x}_j^{(i)}|\boldsymbol{x}_{1:j-1}^{(i)},\boldsymbol{k}_j^{(i)})] \\ = \prod_{i=1}^{N_0} \prod_{j=1}^{n} \mathbb{E}_{\boldsymbol{k}_j^{(i)}}[F(P_M(\boldsymbol{x}_j^{(i)}|\boldsymbol{x}_{1:j-1}^{(i)}),\boldsymbol{k}_j^{(i)})]$$
(6)
$$= \prod_{i=1}^{N_0} \prod_{j=1}^{n} P_M(\boldsymbol{x}_j^{(i)}|\boldsymbol{x}_{0:j-1}^{(i)}) \\ = \prod_{i=1}^{N_0} P_M(\boldsymbol{x}_{1:n}^{(i)}).$$

D.3 PROOF OF THEOREM 4.9

Proof. Part 1. We start from proving $\mathbb{D}(P, F_{GR}) = \mathbb{D}(P, F_{IS}) = 1 - \sum_{t \in V} P(t)^2$. Since both F_{GR} and F_{IS} are distortion-free PDA-rule, $P(t) = \mathbb{E}_k[F_{GR}(P|k)(t)] = \mathbb{E}_k[F_{IS}(P|k)(t)]$. Since $F_{GR}(P|k)$ and $F_{IS}(P|k)$ are Dirac distribution, when $F_{GR}(P|k)(t) > 0$, $F_{GR}(P|k)(t) = 1$, and $\mathbb{E}_{k}[F_{GR}(P|k)(t)] = \mathbb{E}_{k}[\mathbf{1}_{F_{GR}(P|k)(t)>0}] = \Pr(F_{GR}(P|k)(t)>0), \forall t \in V.$ Thus,

$$\mathbb{E}_{k} [\sum_{t \in V} \min\{P(t), F_{GR}(P|k)(t)\}] = \sum_{t \in V} \mathbb{E}_{k} [P(t) \mathbf{1}_{F_{GR}(P|k)(t)>0}]$$

$$= \sum_{t \in V} \mathbb{E}_{k} [P(t) \mathbf{1}_{F_{GR}(P|k)(t)>0}] \Pr(F_{GR}(P|k)(t)>0)$$

$$= \sum_{t \in V} \mathbb{E}_{k} [P(t)|\mathbf{1}_{F_{GR}(P|k)(t)>0}] \Pr(F_{GR}(P|k)(t)>0)$$

$$= \sum_{t \in V} P(t)^{2}.$$
(7)

Analogously, $\mathbb{E}_k[\sum_{t \in V} \min\{P(t), F_{IS}(P|k)(t)\}] = \sum_{t \in V} P(t)^2$. Therefore, we have

$$\mathbb{D}(P, F_{GR}) = \mathbb{D}(P, F_{IS}) = 1 - \sum_{t \in V} P(t)^2.$$
(8)

Part 2. Next, we show $0.5(1 - \max_{t \in V} P(t)) \le \mathbb{D}(P, F_{PR}) \le 0.5 - \max\{\max_{t \in V} P(t) - 0.5, 0\}$. Given a permutation on the token list, assume w.l.o.g. the permutation is of order $\{t_1, ..., t_N\}$, in F_{PR} we will arrange the token probabilities on the interval [0, 1] following the permutation order. Denote by i_0 the index of the token such that 0.5 lies in its probability region, then the token probabilities of $\{t_{i_0+1}, ..., t_N\}$ will be doubled, while the token probabilities of $\{t_1, ..., t_{i_0-1}\}$ will be scaled to 0. Thus, under this permutation,

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} = \sum_{i=i_0+1}^{N} P(t_i) + \min\{P(t_{i_0}), 2\xi_{i_0}\},$$

where ξ_{i_0} is the probability mass of t_{i_0} that is in the interval [0.5, 1], max $\{P(t_{i_0}) - 0.5, 0\} \le \xi_{i_0} \le$ $\min\{0.5, P(t_{i_0})\}$. Next, we consider the reverse permutation $\{t_N, \dots, t_1\}$, following the similar discussion, we have

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k^r)(t)\} = \sum_{i=1}^{i_0-1} P(t_i) + \min\{P(t_{i_0}), 2(P(t_{i_0}) - \xi_{i_0})\}$$

where k^r refers the key that lead to the reserved permutation. Thus,

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^r)(t)\}$$
(9)

$$=1+\min\{P(t_{i_0}),2\xi_{i_0}\}+\min\{P(t_{i_0}),2(P(t_{i_0})-\xi_{i_0})\}-P(t_{i_0}).$$

870 Next, we show $P(t_{i_0}) \ge \min\{P(t_{i_0}), 2\xi_{i_0}\} + \min\{P(t_{i_0}), 2(P(t_{i_0}) - \xi_{i_0})\} - P(t_{i_0}) \ge \max\{\max_{t \in V} P(t) - 0.5, 0\}$. The left hand side inequality is trivial, as $\min\{P(t_{i_0}), 2\xi_{i_0}\} + \min\{P(t_{i_0}), 2(P(t_{i_0}) - \xi_{i_0})\} \le 2P(t_{i_0})$.

For the right hand side inequality, given $\min\{A, 2x\} + \min\{A, 2A - 2x\} = A + \min\{2A - 2x, 2x\}$, we have

$$\min\{P(t_{i_0}), 2\xi_{i_0}\} + \min\{P(t_{i_0}), 2(P(t_{i_0}) - \xi_{i_0})\} - P(t_{i_0}) = 2\min\{P(t_{i_0}) - \xi_{i_0}, \xi_{i_0}\}.$$
 (10)

877 Since $0 \le \max\{P(t_{i_0}) - 0.5, 0\} \le \xi_{i_0} \le \min\{0.5, P(t_{i_0})\} \le P(t_{i_0})$, the minimum value of 878 $\min\{P(t_{i_0}) - \xi_{i_0}, \xi_{i_0}\}$ when ξ_{i_0} take either $\max\{P(t_{i_0}) - 0.5, 0\}$ or $\min\{0.5, P(t_{i_0})\}$, thus

$$\min\{P(t_{i_0}) - \xi_{i_0}, \xi_{i_0}\} \ge \max\{P(t_{i_0}) - 0.5, 0\}.$$
(11)

If $P(t_{i_0}) - 0.5 > 0$, it is obvious that $\max_{t \in V} P(t) = P(t_{i_0})$. So

$$\min\{P(t_{i_0}) - \xi_{i_0}, \xi_{i_0}\} \ge \max\{\max_{t \in V} P(t) - 0.5, 0\}.$$
(12)

Combining it with Equation 9, we have

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^{r})(t)\}$$

=1 + min{ $P(t_{i_{0}}), 2\xi_{i_{0}}\}$ + min{ $P(t_{i_{0}}), 2(P(t_{i_{0}}) - \xi_{i_{0}})\}$ - $P(t_{i_{0}})$.
 \leq 1 + $P(t_{i_{0}}) \leq$ 1 + max $P(t)$, (13)

and

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^{r})(t)\}$$

=1 + min{ $P(t_{i_{0}}), 2\xi_{i_{0}}\}$ + min{ $P(t_{i_{0}}), 2(P(t_{i_{0}}) - \xi_{i_{0}})\}$ - $P(t_{i_{0}}).$
 $\geq 1 + 2 \max\{\max_{t \in V} P(t) - 0.5, 0\}.$ (14)

Since the permutation over V is uniformly seeded with the watermark keys,

$$\mathbb{D}(P, F_{PR}) = 1 - \mathbb{E}_{k} [\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\}]$$

$$= 1 - \frac{1}{2} \mathbb{E}_{k} [\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^{r})(t)\}].$$
(15)

Combining it with Equation 13 and Equation 14, we have

$$0.5(1 - \max_{t \in V} P(t)) \le \mathbb{D}(P, F_{PR}) \le 0.5 - \max\{\max_{t \in V} P(t) - 0.5, 0\}.$$
(16)

Part 3. Finally, we show $\mathbb{D}(P, F_{PR}) \leq \mathbb{D}(P, F_{IS}) = \mathbb{D}(P, F_{GR})$. We only need to prove $0.5 - \max\{\max_{t \in V} P(t) - 0.5, 0\} \leq 1 - \sum_{t \in V} P(t)^2$. We have two steps for Part 3.

912 **Lemma D.1.** Given
$$0 \le x_1, x_2 \le r_0 \le r_1, x_1 + x_2 = r_1 \le 1$$
, we have $x_1^2 + x_2^2 \le r_0^2 + (r_1 - r_0)^2$
913

914 **Proof.**
$$x_1^2 + x_2^2 = x_1^2 + (r_1 - x_1)^2 = 2x_1^2 - 2x_1r_1 + r_1^2 = 2(x_1 - r_1/2)^2 + r_1^2/2 \le \min_{x_1} 2(x_1 - r_1/2)^2 + r_1^2/2 = r_0^2 + (r_1 - r_0)^2.$$

917 Thus, by inductive we have $1 - \sum_{t \in V} P(t)^2 \ge 1 - \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor (\max_{t \in V} P(t))^2 - (1 - \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor \max_{t \in V} P(t))^2$. Now we continue the proof of the main theorem. Step 1. When $\max_{t \in V} P(t) \ge 0.5$, $1 - \sum_{t \in V} P(t)^2 \ge 1 - (\max_{t \in V} P(t))^2 - (1 - \max_{t \in V} P(t))^2$ $= 2 \max_{t \in V} P(t) - 2(\max_{t \in V} P(t))^2$ $= 0.5 - 2(\max_{t \in V} P(t) - 0.5)^2$ (17) $\geq 0.5 - (\max_{t \in V} P(t) - 0.5)$ $= 0.5 - \max\{\max_{t \in V} P(t) - 0.5, 0\}.$

Step 2. When $\max_{t \in V} P(t) \leq 0.5$,

97[.]

$$1 - \sum_{t \in V} P(t)^{2} \ge 1 - \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor (\max_{t \in V} P(t))^{2} - (1 - \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor \max_{t \in V} P(t))^{2}$$

= $2 \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor \max_{t \in V} P(t) - (\lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor + \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor^{2}) (\max_{t \in V} P(t))^{2},$
(18)

denote by $\epsilon = \frac{1}{\max_{t \in V} P(t)} - \lfloor \frac{1}{\max_{t \in V} P(t)} \rfloor$, we have $0 \le \epsilon < 1$ and $1 - \sum_{t \in V} P(t)^2 = 2(\frac{1}{\max_{t \in V} P(t)} - \epsilon) \max_{t \in V} P(t) - ((\frac{1}{\max_{t \in V} P(t)} - \epsilon) + (\frac{1}{\max_{t \in V} P(t)} - \epsilon)^2)(\max_{t \in V} P(t))^2,$

$$\frac{1}{V} \qquad \max_{t \in V} P(t) \qquad t \in V \qquad \max_{t \in V} P(t) \qquad t \in V$$

$$= 2 - 2\epsilon \max_{t \in V} P(t) - \left(\max_{t \in V} P(t) - \epsilon \max_{t \in V} P(t)^2 + 1 - 2\epsilon \max_{t \in V} P(t) + \epsilon^2 \max_{t \in V} P(t)^2 \right)$$

$$= 1 - \max_{t \in V} P(t) + (\epsilon - \epsilon^2) \max_{t \in V} P(t)^2$$

$$\geq 1 - \max_{t \in V} P(t) \ge 0.5 = 0.5 - \max\{\max_{t \in V} P(t) - 0.5, 0\}.$$
(19)

By Step 1 and Step 2, we have $\mathbb{D}(P, F_{PR}) \leq \mathbb{D}(P, F_{IS}) = \mathbb{D}(P, F_{GR})$.

D.4 PROOF OF THEOREM 4.10

Proof. Consider the scenario of generating multiple responses with the **same-prompt single-token**-generation task. According to Definition 4.5 under the strongly distortion-free condition, one must have $\forall P_M \in \mathcal{P}, \forall N_0 \in \mathbb{N}_+, \forall t^{(i)} \in V, \prod_{i=1}^{N_0} P_M(t^{(i)}) = \mathbb{E}_{k^{(1)},\dots,k^{(N_0)}}[\prod_{i=1}^{N_0} F(P_M|k^{(i)})(t^{(i)})].$ Under key collisions, there exists at least two $k^{(i)}$, $k^{(j)}$ are the same. Then we have $\forall P_M \in \mathcal{P}, \exists N_0 \geq 2, \forall t^{(i)} \in V, \prod_{i=1}^{N_0} P_M(t^{(i)}) = \mathbb{E}_k[\prod_{i=1}^{N_0} F(P_M|k)(t^{(i)})]$. We will show that this hold if and only if $\mathbb{D}(P_M, F) = 0.$

Part 1. It is obviously that $\mathbb{D}(P_M, F) = 0$ can lead to $\forall N_0 \in \mathbb{N}_+, \forall t^{(i)} \in V, \prod_{i=1}^{N_0} P_M(t^{(i)}) =$ $\mathbb{E}_{k}[\prod_{i=1}^{N_{0}} F(P_{M}|k)(t^{(i)})]. \text{ This is because if } \mathbb{D}(P_{M},F) = 0, P_{M}(t^{(i)}) = F(P_{M}|k)(t^{(i)}) \text{ almost surely and thus } \mathbb{E}_{k}[\prod_{i=1}^{N_{0}} F(P_{M}|k)(t^{(i)})] = \mathbb{E}_{k}[\prod_{i=1}^{N_{0}} P_{M}(t^{(i)})] = \prod_{i=1}^{N_{0}} P_{M}(t^{(i)}).$

Part 2. Now we will show that if $\exists N_0 \geq 2, \forall t^{(i)} \in V, \prod_{i=1}^{N_0} P_M(t^{(i)}) = \mathbb{E}_k[\prod_{i=1}^{N_0} F(P_M|k)(t^{(i)})],$ then $\mathbb{D}(P_M, F) = 0$.

As $t^{(i)}$ is arbitrary selected, we can choose $t^{(1)} = \dots = t^{(N_0)} = t$, then we have $P_M(t)^{N_0} = t$ $\mathbb{E}_k[F(P_M|k)(t)^{N_0}]$. By Jensen's inequality, when $N_0 \ge 2$,

$$P_{M}(t)^{N_{0}} = \mathbb{E}_{k}[F(P_{M}|k)(t)^{N_{0}}] \ge (\mathbb{E}_{k}[F(P_{M}|k)(t)])^{N_{0}} = P_{M}(t)^{N_{0}}.$$

The equality is achieved if and only if $F(P_M|k)(t) = \mathbb{E}_k[F(P_M|k)(t)] = P_M(t)$. Thus, $\forall t \in \mathbb{E}_k[F(P_M|k)(t)] = P_M(t)$. $V, \forall k \in K, F(P_M|k)(t) = P_M(t)$, which leads to

970
971

$$\mathbb{D}(P_M, F) = 1 - \mathbb{E}_k [\sum_{t \in V} \min\{P_M(t), F(P_M|k)(t)\}] = 1 - \sum_{t \in V} P_M(t) = 0.$$
971

972 D.5 PROOF OF THEOREM 4.12

Proof. Firstly, let's consider $\frac{1}{m} |\sum_{j=1}^{m} \operatorname{Met}(g_j^{p_i}(P_M)) - \sum_{j=1}^{m} \operatorname{Met}(g_j^{p_i}(P_T))|$, each $\operatorname{Met}(g_j^{p_i}(P_T))$ and $\operatorname{Met}(g_j^{p_i}(P_M))$ are independent distributed. With Hoeffding's inequality we have

$$\Pr(\frac{1}{m}|\sum_{j=1}^{m} \operatorname{Met}(g_{j}^{p_{i}}(P_{M})) - \sum_{j=1}^{m} \operatorname{Met}(g_{j}^{p_{i}}(P_{T}))| > t) < e^{-\frac{mt^{2}}{2A^{2}}}.$$

Denote by $\text{DetWmk}_i(P_M, P_T) = \frac{1}{m} |\sum_{j=1}^m \text{Met}(g_j^{p_i}(P_M)) - \sum_{j=1}^m \text{Met}(g_j^{p_i}(P_T))| - \frac{1}{m} |\sum_{j=1}^m \text{Met}(g_j^{p_i}(P_M)) - \sum_{j=1}^m \text{Met}(g_j^{p_i}(P_{M'}))|$. Since if $\text{DetWmk}_i(P_M, P_T) > t$, we must have $\frac{1}{m} |\sum_{j=1}^m \text{Met}(g_j^{p_i}(P_M)) - \sum_{j=1}^m \text{Met}(g_j^{p_i}(P_T))| > t$, which yields

$$\Pr(\text{DetWmk}_{i}(P_{M}, P_{T}) > t) \le \Pr(\frac{1}{m} |\sum_{j=1}^{m} \text{Met}(g_{j}^{p_{i}}(P_{M})) - \sum_{j=1}^{m} \text{Met}(g_{j}^{p_{i}}(P_{T}))| > t) < e^{-\frac{mt^{2}}{2A^{2}}}.$$
(20)

Analogously, we have $\Pr(\text{DetWmk}_i(P_M, P_T) < -t) < e^{-\frac{mt^2}{2A^2}}$ Thus, $\text{DetWmk}_i(P_M, P_T)$ is sub-Gaussian distributed and it is easy to observe that $\mathbb{E}[\text{DetWmk}_i(P_M, P_T)] = 0$. Since $\text{DetWmk}_i(P_M, P_T), i = 1, ..., n$ is independently distributed, applying Hoeffding's inequality again, we have

$$\Pr(\text{DetWmk}(P_M, P_T) > t) < \exp(-\frac{mn^2 t^2}{2A^2 \sum_{i=1}^n ||\text{DetWmk}_i(P_M, P_T)||_{\psi_2}^2})$$

996 where $||X||_{\psi_2} = \inf\{c > 0 : \mathbb{E}[e^{X^2/c^2}] \le 2\}.$

Now we need to calculate $||\text{DetWmk}_i(P_M, P_T)||_{\psi_2}^2$. We start from calculating $\mathbb{E}[e^{\text{DetWmk}_i(P_M, P_T)^2/c^2}]$. Since $|\text{DetWmk}_i(P_M, P_T)| \leq A$ and the probability density function of $\text{DetWmk}_i(P_M, P_T)$ is symmetric with respect to 0, we have

$$\mathbb{E}[e^{\text{DetWmk}_{i}(P_{M},P_{T})^{2}/c^{2}}] = \int_{-A}^{A} e^{x^{2}/c^{2}} dF_{\text{DetWmk}_{i}(P_{M},P_{T})}(x)$$

$$= e^{x^{2}/c^{2}} F_{\text{DetWmk}_{i}(P_{M},P_{T})}(x)|_{i=-A}^{A} - \int_{-A}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} F_{\text{DetWmk}_{i}(P_{M},P_{T})}(x) dx$$

$$= e^{A^{2}/c^{2}} - \int_{-A}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} (1 - \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x)) dx$$

$$= e^{A^{2}/c^{2}} + \int_{-A}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$$

$$= e^{A^{2}/c^{2}} + \int_{-A}^{0} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$$

$$= e^{A^{2}/c^{2}} + \int_{-A}^{0} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$$

$$= e^{A^{2}/c^{2}} - \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$$

$$= e^{A^{2}/c^{2}} - \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) < x) dx$$

$$= e^{A^{2}/c^{2}} - \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) < x) dx$$

$$= e^{A^{2}/c^{2}} - \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) < x) dx$$

$$= e^{A^{2}/c^{2}} - \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$$

$$= 1 + 2 \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$$

$$(21)$$

Since when $x \ge 0$, $\Pr(\text{DetWmk}_i(P_M, P_T) > x) \le e^{-\frac{mt^2}{2A^2}}$, we have 1026 1027 1028 $\mathbb{E}[e^{\text{DetWmk}_{i}(P_{M},P_{T})^{2}/c^{2}}] = 1 + 2\int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} \Pr(\text{DetWmk}_{i}(P_{M},P_{T}) > x) dx$ 1029 1030 $\leq 1 + 2 \int_{0}^{A} \frac{2x}{c^{2}} e^{x^{2}/c^{2}} e^{-\frac{mt^{2}}{2A^{2}}} dx$ (22)1031 1032 $=1+\frac{2}{\frac{mc^{2}}{0.4t^{2}}-1}(1-e^{A^{2}/c^{2}-m/2})$ 1033 1034 1035 Taking $c = \sqrt{\frac{6A^2}{m}}$, $\mathbb{E}[e^{\text{DetWmk}_i(P_M, P_T)^2/c^2}] \leq 1 + (1 - e^{-m/3}) < 2$. Thus, 1036 $||\text{DetWmk}_{i}(P_{M}, P_{T})||_{\psi_{2}}^{2} \leq \frac{6A^{2}}{m}$ and 1037 1038 $\Pr(\text{DetWmk}(P_M, P_T) > t) < \exp(-\frac{mn^2t^2}{2A^2\sum_{i=1}^n ||\text{DetWmk}_i(P_M, P_T)||_{2b_2}^2}) < \exp(-\frac{m^2nt^2}{12A^4}),$ 1039 1040 1041 1042 1043 D.6 PROOF OF THEOREM 4.13 1044 1045 Proof. Combining Lemma D.2 and Lemma D.3 yields the result. 1046 **Lemma D.2.** A strongly distortion-free watermark must use independent hashed watermark keys 1047 h(k). 1048 1049 *Proof.* The proof of this lemma is similar to the proof of Theorem 4.10. We prove by contradiction, if 1050 we don't have independent hashed watermark keys h(k), then given two randomly sampled key $h(k_1)$ 1051 and $h(k_2)$, $\Pr(h(k_2) = A, h(k_1) = B) \neq \Pr(h(k_1) = A) \Pr(h(k_2) = B)$. Consider the scenario 1052 of generating multiple responses with the same-prompt one-token-generation task. According 1053 to Definition 4.5 under the strongly distortion-free condition, one must have $\forall P_M \in \mathcal{P}, \forall N_0 \in$ $\mathbb{N}_{+}, \forall t^{(i)} \in V, \prod_{i=1}^{N_0} P_M(t^{(i)}) = \mathbb{E}_{h(k^{(1)}),\dots,h(k^{(N_0)})} [\prod_{i=1}^{N_0} F(P_M|h(k^{(i)}))(t^{(i)})].$ 1054 1055 show that if $\exists N_0 \geq 2, \forall t^{(i)}$ $\in V, \quad \prod_{i=1}^{N_0} P_M(t^{(i)})$ We 1056 = $\mathbb{E}_{h(k^{(1)}),\dots,h(k^{(N_0)})}[\prod_{i=1}^{N_0} F(P_M|h(k^{(i)}))(t^{(i)})], \text{ then } \mathbb{D}(P_M, F) = 0.$ 1057 1058 As $t^{(i)}$ is arbitrary selected, we can choose $t^{(1)} = \dots = t^{(N_0)} = t$, then we have $P_M(t)^{N_0} =$ 1059 $\mathbb{E}_{h(k^{(1)}),\dots,h(k^{(N_0)})}[\prod_{i=1}^{N_0} F(P_M|h(k^{(i)}))(t)].$ Assume w.l.o.g. $N_0 \ge 2$, 1060 1061 $P_M(t)^2 - \mathbb{E}_{h(k^{(1)}),h(k^{(2)})}[F(P_M|h(k^{(1)}))(t)F(P_M|h(k^{(2)}))(t)],$ 1062 $=P_M(t)^2 - \sum_A \sum_B F(P_M|A)(t)F(P_M|B)(t)\Pr(h(k^{(1)}) = A, h(k^{(2)}) = B),$ 1063 1064 1065 $= \sum_{n} \sum_{m} F(P_M|A)(t)F(P_M|B)(t)\Pr(h(k^{(1)}) = A)\Pr(h(k^{(2)}) = B)$ 1066 1067 $-\sum_{n}\sum_{m}F(P_{M}|A)(t)F(P_{M}|B)(t)\Pr(h(k^{(1)})=A,h(k^{(2)})=B),$ 1068 1069 $=\sum_{A}\sum_{D}F(P_{M}|A)(t)F(P_{M}|B)(t)[\Pr(h(k^{(1)})=A)\Pr(h(k^{(2)})=B)-\Pr(h(k^{(1)})=A,h(k^{(2)})=B)].$ 1070 1071 (23)1072 1073 Since $\exists A, B$, such that $\Pr(h(k^{(1)}) = A) \Pr(h(k^{(2)}) = B) \neq \Pr(h(k^{(1)}) = A, h(k^{(2)}) = B)$, there 1074 exists a P_M such that 1075 $\sum \sum F(P_M|A)(t)F(P_M|B)(t)[\Pr(h(k^{(1)}) = A)\Pr(h(k^{(2)}) = B) - \Pr(h(k^{(1)}) = A, h(k^{(2)}) = B)]_+$ 1076 1077

$$+\sum_{A}\sum_{B}F(P_{M}|A)(t)F(P_{M}|B)(t)[\Pr(h(k^{(1)}) = A)\Pr(h(k^{(2)}) = B) - \Pr(h(k^{(1)}) = A, h(k^{(2)}) = B)]_{-} \neq 0$$
(24)

1080 In this case, $P_M(t)^2 - \mathbb{E}_{h(k^{(1)}), h(k^{(2)})}[F(P_M|h(k^{(1)}))(t)F(P_M|h(k^{(2)}))(t)] \neq 0$, thus the watermark 1081 is not strongly distortion-free. \square 1082

Lemma D.3. A watermark using a distortion-free and independent PDA-rule is undetectable by any 1083 arbitrary detector. 1084

1085 *Proof.* Recall that the watermarking detection algorithm utilize the statistical difference between the 1086 watermarked LM and the original LM to check the existence of the watermark, i.e., the detection 1087 is based on $\mathbb{E}[P_M(\boldsymbol{x}_{1:n})|\text{DetCon}] \neq \mathbb{E}[\prod_{i=1}^n F(P_M|h(k^{(i)}))(t^{(i)})|\text{DetCon}]$, where DetCon is the 1088 detecting condition which is used in watermark generator. Now we show if the PDA-rule is inde-1089 pendent from each other, the DetCon will be independent of the PDA-rule. This can be shown by 1090 contradiction. 1091

If DetCon is not independent of the PDA-rules, during the generation process, the PDA-rules will 1092 be mutually dependent because they all share dependency on the DetCon, which contradicts to the 1093 independence of PDA-rule. Thus, we have 1094

$$\mathbb{E}[\prod_{i=1}^{n} F(P_M | h(k^{(i)}))(t^{(i)}) | \text{DetCon}] = \prod_{i=1}^{n} \mathbb{E}[F(P_M | h(k^{(i)}))(t^{(i)}) | \text{DetCon}] = \mathbb{E}[P_M(\boldsymbol{x}_{1:n}) | \text{DetCon}].$$

Thus, the watermark is undetectable by the detector.

Thus, the watermark is undetectable by the detector. 1098

1101 D.7 PROOF OF THEOREM 5.2 1102

Г

1095 1096

1099

1100

1107 1108 1109

1111 1112

1103 *Proof.* We need to show $P_M(t|\boldsymbol{x}_{1:i-1}) = \mathbb{E}_{k_i}[F_\beta(P_M(\cdot|\boldsymbol{x}_{1:i-1}),k_i)(t)].$ As $F_{PR}(P_M(\cdot|\boldsymbol{x}_{1:i-1}),k_i)(t)$ is a distortion-free PDA-rule, we know $\mathbb{E}_{k_i}[(1 + k_i)(t) - k_i)(t)]$ 1104 _ $\beta F_{PR}(P_M(\cdot | \boldsymbol{x}_{1:i-1}), k_i)(t)] = (1 - \beta)P_M(t | \boldsymbol{x}_{1:i-1}).$ Thus, we need to show 1105 1106

$$\begin{aligned}
& \text{Intro}' \\
& \text{Intro}'$$

1113 Since the permutation is uniformly distributed, denoted by Π the set of all permutations on V and P_{Π} 1114 the uniformly distribution on Π , we have 1115

$$\mathbb{E}_{k_{i}} \begin{bmatrix} \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} - \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} \end{bmatrix} \\ \mathbb{E}_{\pi\sim P_{\Pi}} \begin{bmatrix} \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} - \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} \end{bmatrix} \\ \mathbb{E}_{\pi\sim P_{\Pi}} \begin{bmatrix} \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} - \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} \end{bmatrix} \end{bmatrix}$$

$$\mathbb{E}_{\pi\sim P_{\Pi}} \begin{bmatrix} \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} - \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} \end{bmatrix} \end{bmatrix}$$

$$\mathbb{E}_{\pi\sim P_{\Pi}} \begin{bmatrix} \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} - \max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\} \end{bmatrix} \end{bmatrix}$$

As P_{Π} is the uniformly distribution on Π , for each $\pi \in \Pi$, we consider its reverse permutation π^r : 1123

$$\begin{aligned} & \underset{1126}{1126} \\ & \underset{1126}{1126} \\ & \underset{1127}{1126} \\ & \underset{1127}{1128} \\ & \underset{129}{1129} \\ & = \frac{1}{2} \mathbb{E}_{\pi^{r} \sim P_{\Pi}} \left[\max\{2 \sum_{t', \pi(t'|k_{i}) \geq \pi(t|k_{i})} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} - \max\{2 \sum_{t', \pi(t'|k_{i}) \geq \pi(t|k_{i}) + 1} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} \right] \\ & \underset{129}{1130} \\ & \underset{131}{132} \\ & + \max\{2 \sum_{t', \pi^{r}(t'|k_{i}) \geq \pi^{r}(t|k_{i})} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} - \max\{2 \sum_{t', \pi^{r}(t'|k_{i}) \geq \pi(t|k_{i}) + 1} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} - \max\{2 \sum_{t', \pi^{r}(t'|k_{i}) \geq \pi^{r}(t|k_{i}) + 1} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} \right] \end{aligned}$$

Notice, if $\pi(t') \leq \pi(t)$, then in the reversed permutation π^r , we have $\pi^r(t') \geq \pi^r(t)$ and vice versa. Thus, $\max\{2\sum_{t',\pi^{r}(t'|k_{i})\geq\pi^{r}(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\}-\max\{2\sum_{t',\pi^{r}(t'|k_{i})\geq\pi^{r}(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\}\\=\max\{2\sum_{t',\pi(t'|k_{i})\leq\pi(t|k_{i})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\}-\max\{2\sum_{t',\pi(t'|k_{i})\leq\pi(t|k_{i})-1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\}\\=\max\{1-2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1}),0\}-\max\{1-2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})-1}P_{M}(t'|\boldsymbol{x}_{1:i-1}),0\}.$ (28) By $\max\{x, 0\} - \max\{-x, 0\} = x$ we have $\mathbb{E}_{\pi \sim P_{\Pi}} \left| \max\{2 \sum_{t', \pi(t'|k_i) \ge \pi(t|k_i)} P_M(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} - \max\{2 \sum_{t', \pi(t'|k_i) \ge \pi(t|k_i) + 1} P_M(t'|\boldsymbol{x}_{1:i-1}) - 1, 0\} \right|$ $=\frac{1}{2}\mathbb{E}_{\pi\sim P_{\Pi}}\left[\max\{2\sum_{t'.\pi(t'|k_i)>\pi(t|k_i)}P_M(t'|\boldsymbol{x}_{1:i-1})-1,0\}-\max\{2\sum_{t',\pi(t'|k_i)\geq\pi(t|k_i)+1}P_M(t'|\boldsymbol{x}_{1:i-1})-1,0\}\right]$ $+ \max\left\{2\sum_{\substack{t',\pi^r(t'|k_i) \ge \pi^r(t|k_i) \\ \mathsf{r}}} P_M(t'|\mathbf{x}_{1:i-1}) - 1, 0\right\} - \max\left\{2\sum_{\substack{t',\pi^r(t'|k_i) \ge \pi^r(t|k_i) + 1 \\ \mathsf{r}}} P_M(t'|\mathbf{x}_{1:i-1}) - 1, 0\right\}\right]$ $=\frac{1}{2}\mathbb{E}_{\pi\sim P_{\Pi}}\left[\max\{2\sum_{t'|\pi(t'|k_{:})>\pi(t|k_{:})}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\}-\max\{2\sum_{t',\pi(t'|k_{i})\geq\pi(t|k_{i})+1}P_{M}(t'|\boldsymbol{x}_{1:i-1})-1,0\}\right]$ $+ \max\{1 - 2\sum_{t', \pi(t'|k_i) \ge \pi(t|k_i) + 1} P_M(t'|\boldsymbol{x}_{1:i-1}), 0\} - \max\{1 - 2\sum_{t', \pi(t'|k_i) \ge \pi(t|k_i) - 1} P_M(t'|\boldsymbol{x}_{1:i-1}), 0\}$ $= \frac{1}{2} \mathbb{E}_{\pi \sim P_{\Pi}} \left[2 \sum_{t', \pi^{r}(t'|k_{i}) \geq \pi^{r}(t|k_{i})} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1 - \left(2 \sum_{t', \pi^{r}(t'|k_{i}) \geq \pi^{r}(t|k_{i}) + 1} P_{M}(t'|\boldsymbol{x}_{1:i-1}) - 1\right) \right]$ $= \frac{1}{2} \mathbb{E}_{\pi \sim P_{\Pi}} [2P_M(t | \boldsymbol{x}_{1:i-1})]$ $=P_M(t|\boldsymbol{x}_{1\cdot i-1}).$ (29)

D.8 PROOF OF THEOREM 5.4

Proof. Part 1. We first show $\forall P \in \mathcal{P}, \mathbb{D}(P, F_{\beta}) \leq \mathbb{D}(P, F_{PR}) - \beta(1 - \max_{t \in V} P(t))$. According to the Part 2 of Proof D.3, we know that given a permutation $\{t_1, ..., t_N\}$ and let t_{i_0} is the token whose probability mass expands across 1/2,

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} = \sum_{i=i_0+1}^{N} P(t_i) + \min\{P(t_{i_0}), 2\xi_{i_0}\},$$

where ξ_{i_0} is the probability mass of t_{i_0} that is in the interval [0.5, 1] (notice t_{i_0} is the same for both permuta-reweight and beta PDA-rule as they use the same permutation), $\max\{P(t_{i_0}) - 0.5, 0\} \leq 1$ $\xi_{i_0} \leq \min\{0.5, P(t_{i_0})\}$. And

1185
1186
1187

$$\sum_{t \in V} \min\{P(t), F_{PR}(P|k^r)(t)\} = \sum_{i=1}^{i_0-1} P(t_i) + \min\{P(t_{i_0}), 2(P(t_{i_0}) - \xi_{i_0})\},$$
1187

where k^r refers the key that lead to the reserved permutation.

Now we consider F_{β} . From the similar analysis, we have $\sum_{t \in V} \min\{P(t), F_{\beta}(P|k)(t)\} = \sum_{i=i_0+1}^{N} P(t_i) + 2\beta \sum_{i=1}^{i_0-1} P(t_i) + \min\{P(t_{i_0}), 2(1-\beta)\xi_{i_0} + 2\beta(P(t_{i_0}) - \xi_{i_0})\},$ $\sum_{t \in V} \min\{P(t), F_{\beta}(P|k^{r})(t)\} = \sum_{i=1}^{i_{0}-1} P(t_{i}) + 2\beta \sum_{i=i+1}^{N} P(t_{i}) + \min\{P(t_{i_{0}}), 2(1-\beta)(P(t_{i_{0}}) - \xi_{i_{0}}) + 2\beta\xi_{i_{0}}\}.$ $\min\{P(t_{i_0}), 2(1-\beta)\xi_{i_0} + 2\beta(P(t_{i_0}) - \xi_{i_0})\} + \min\{P(t_{i_0}), 2(1-\beta)(P(t_{i_0}) - \xi_{i_0}) + 2\beta\xi_{i_0}\}$ $= P(t_{i_0}) + \min\{2(1-\beta)(P(t_{i_0}) - \xi_{i_0}) + 2\beta\xi_{i_0}, 2(1-\beta)\xi_{i_0} + 2\beta(P(t_{i_0}) - \xi_{i_0})\}$ $= P(t_{i_0}) + 2\xi_{i_0} + \min\{2(1-\beta)(P(t_{i_0}) - 2\xi_{i_0}), 2\beta(P(t_{i_0}) - 2\xi_{i_0})\}$ $\geq P(t_{i_0}) + 2\xi_{i_0} + \min\{0, 2(P(t_{i_0}) - 2\xi_{i_0})\}\$ $= \min\{P(t_{i_0}), 2(P(t_{i_0}) - \xi_{i_0})\} + \min\{P(t_{i_0}), 2\xi_{i_0}\},\$ (30)we have $\sum_{t \in V} \min\{P(t), F_{\beta}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{\beta}(P|k^{r})(t)\}$ $= 1 - P(t_0) + 2\beta(1 - P(t_0)) + \min\{P(t_{i_0}), 2(1 - \beta)\xi_{i_0} + 2\beta(P(t_{i_0}) - \xi_{i_0})\}$ $+\min\{P(t_{i_0}), 2(1-\beta)(P(t_{i_0})-\xi_{i_0})+2\beta\xi_{i_0}\}\$ (31) $\geq \sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^r)(t)\} + 2\beta - 2\beta P(t_{i_0})$ $\geq \sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^r)(t)\} + 2\beta - 2\beta \max_{t \in V} P(t).$ Thus, $\mathbb{D}(P, F_{\beta}) = 1 - \mathbb{E}_{k} \left[\sum_{i=1}^{k} \min\{P(t), F_{\beta}(P|k)(t)\} \right]$ $= 1 - \frac{1}{2} \mathbb{E}_{k} \left[\sum_{t \in V} \min\{P(t), F_{\beta}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{\beta}(P|k^{r})(t)\} \right]$ $\leq 1 - \frac{1}{2} \mathbb{E}_{k} \left[\sum_{t \in V} \min\{P(t), F_{PR}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{PR}(P|k^{r})(t)\} + 2\beta - 2\beta \max_{t \in V} P(t) \right]$ $= \mathbb{D}(P, F_{PR}) - \beta(1 - \max_{t \in \mathcal{V}} P(t)).$ (32)

Part 2. We then show $\forall P \in \mathcal{P}$, if $\beta_1 \leq \beta_2$, then $\mathbb{D}(P, F_{\beta_1}) \geq \mathbb{D}(P, F_{\beta_2})$. Consider $\mathbb{D}(P, F_{\beta_1}) - \mathbb{D}(P, F_{\beta_2})$, we have

$$\mathbb{D}(P, F_{\beta_{1}}) - \mathbb{D}(P, F_{\beta_{2}})$$

$$= \mathbb{E}_{k} \left[\sum_{t \in V} \min\{P(t), F_{\beta_{2}}(P|k)(t)\} \right] - \mathbb{E}_{k} \left[\sum_{t \in V} \min\{P(t), F_{\beta_{1}}(P|k)(t)\} \right]$$

$$= \frac{1}{2} \mathbb{E}_{k} \left[\sum_{t \in V} \min\{P(t), F_{\beta_{2}}(P|k)(t)\} + \sum_{t \in V} \min\{P(t), F_{\beta_{2}}(P|k^{r})(t)\} \right]$$

$$- \sum_{t \in V} \min\{P(t), F_{\beta_{1}}(P|k)(t)\} - \sum_{t \in V} \min\{P(t), F_{\beta_{1}}(P|k^{r})(t)\} \right]$$
(33)

 $\overline{t\in V}$

1239 From the similar analysis as Part 1 we have for F_{β_1} ,

1240
1241
$$\sum_{t \in V} \min\{P(t), F_{\beta_1}(P|k)(t)\} = \sum_{i=i_0+1}^N P(t_i) + 2\beta_1 \sum_{i=1}^{i_0-1} P(t_i) + \min\{P(t_{i_0}), 2(1-\beta_1)\xi_{i_0} + 2\beta_1(P(t_{i_0}) - \xi_{i_0})\}$$

and

$$\sum_{i \in V} \min\{P(t), F_{\beta_1}(P|k^r)(t)\} = \sum_{i=1}^{i_0-1} P(t_i) + 2\beta_1 \sum_{i=i_0+1}^{N} P(t_i) + \min\{P(t_{i_0}), 2(1-\beta_1)(P(t_{i_0}) - \xi_{i_0}) + 2\beta_1\xi_{i_0}\}, \sum_{i \in V} \min\{P(t), F_{\beta_2}(P|k)(t)\} = \sum_{i=i_0+1}^{N} P(t_i) + 2\beta_2 \sum_{i=1}^{i_0-1} P(t_i) + \min\{P(t_{i_0}), 2(1-\beta_2)\xi_{i_0} + 2\beta_2(P(t_{i_0}) - \xi_{i_0})\}, \max$$
and

$$\sum_{i \in V} \min\{P(t), F_{\beta_2}(P|k^r)(t)\} = \sum_{i=1}^{i_0-1} P(t_i) + 2\beta_2 \sum_{i=i_0+1}^{N} P(t_i) + \min\{P(t_{i_0}), 2(1-\beta_2)(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2\xi_{i_0}\}, 2\beta_2(P(t_{i_0}) - 2\beta_{i_0}) + 2\beta_2\xi_{i_0}\}, 2\beta_2(P(t_{i_0}) - 2\beta_{i_0}) + 2\beta_2\xi_{i_0}\}, 2\beta_2(P(t_{i_0}) - 2\xi_{i_0}) + 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2\xi_{i_0}\}, 2\beta_2(P(t_{i_0}) - 2\xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - 2\xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - 2\xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2\xi_{i_0}\}, 2\beta_2(P(t_{i_0}) - 2\xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - 2\xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2\xi_{i_0}\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2(\xi_{i_0}), 2\beta_2(P(t_{i_0}) - \xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2(\xi_{i_0}) + 2\beta_2(P(t_{i_0}) - \xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2(P(t_{i_0}) - \xi_{i_0})\}, 2\beta_2(P(t_{i_0}) - \xi_{i_0}) + 2\beta_2(P(t_$$

1290 D.9 PROOF OF DEFINITION 5.5

1291 Proof. We prove the concentration bound in Definition 5.5: $\Pr(S(\boldsymbol{x}_{1:n}) - \mathbb{E}_{H_0}[S(\boldsymbol{x}_{1:n})] > t\sqrt{n}|H_0) \le \exp(-2t^2)$. Since the range of the sigmoid function is in [0, 1], by Hoeffding's inequality, for each random score $s(x_i)$, we have

1294
1295
$$\Pr\left(\frac{1}{n}\sum_{i=1}^{n}s(x_{i}) - \mathbb{E}_{H_{0}}\left[\frac{1}{n}\sum_{i=1}^{n}s(x_{i})\right] > t|H_{0}\right) \le e^{-2nt^{2}}.$$
(37)

Replace t by $\frac{t}{\sqrt{n}}$ we have

1296

1302 1303

1305

1345

1347

1348

1349

$$\Pr(\sum_{i=1}^{n} s(x_i) - \mathbb{E}_{H_0}[\sum_{i=1}^{n} s(x_i)] > t\sqrt{n}|H_0) \le e^{-2t^2}.$$
(38)

1304 E DETAILED EXPERIMENT SETUP

1306 E.1 EXPERIMENT SETUP

We evaluate the distortion-free performance of various watermark models within two seq2seq
applications: text summarization and text generation. The experiments leverage the Huggingface
library (Wolf et al., 2019), a popular framework for model development and sharing in the NLP
community. All tests are conducted on 8 NVIDIA A6000 GPUs, each with 48GB of memory.

We focus on three seq2seq tasks in our experiments: machine translation, text summarization and 1312 text generation. For the machine translation task, we focus on English-to-Romanian translation. We 1313 employ the Multilingual BART (MBart) model (Liu et al., 2020) on the WMT'14 En-Ro corpus. For 1314 text summarization, we employ the BART-large model (Liu et al., 2020) using the CNN-DM corpus 1315 dataset (Hermann et al., 2015). For text generation, we follow the settings described by (Kirchenbauer 1316 et al., 2023), using the LLaMA-2 model (7b, chat) (Touvron et al., 2023) with a random subset 1317 of the C4 dataset (Raffel et al., 2020). All experiments are conducted with n-gram watermark key 1318 sampling (n = 5). Additionally, we include the Soft watermark (Kirchenbauer et al., 2023) in our 1319 comparison, although it does not achieve step-wise distortion-free performance. Notably, when 1320 $\beta = 0$, the Beta-watermark becomes identical to the permute-reweight watermark (Hu et al., 2023).

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. For text summarization, we utilize the test set from the CNN-DM corpus (Hermann et al., 2015), which contains 11,490 examples. We employ the BART-large model, which has 400 million parameters, and the LLaMA-2 model with 7 billion parameters.

Text Generation. In text generation, we adhere to the experimental setup described in Kirchenbauer et al. (2023). We use a random subset of the C4 dataset for generation prompts. Our model selection includes the LLaMA-2, which has 7 billion parameters.

1331 Watermark Setup. Our experiments primarily compare the beta-watermark with three other 1332 distortion-free watermarks: inversa-sampling, Gumbel-reparametrization, and permute-reweight. 1333 Additionally, we include the Soft watermark (Kirchenbauer et al., 2023) in our comparison. For 1334 beta-watermark, we explore various β values from the set {0, 0.05, 0.1, 0.2, 0.3}. For the Soft watermark (Kirchenbauer et al., 2023), we investigate green list bias δ values from $\{0.5, 1.0, 1.5, 2.0\}$ 1335 with a fixed green list separator $\gamma = 0.5$. For n-gram key sampling, we consider the most recent 5 1336 tokens as the texture key. For example, when generating x_4 in response to (x_1, x_2, x_3) , the texture 1337 key includes (x_1, x_2, x_3) , given only three tokens are available. Texture key history resets before 1338 generating each batch. For cipher generation, we use SHA-256 as the hash function and a 1024-bit 1339 random bitstrings as the secret key sk, the watermark key is given by $k = (sk, x_{i-5,i-1})$. The 1340 permutation π is sampled using hash(k) as the random seed. We also compare beta-watermark 1341 with inverse-sampling watermark Kuditipudi et al. (2023) and permute-reweight watermark Hu et al. 1342 (2023); Wu et al. (2023), following the settings in their open-sourced $code^{12}$.

Evaluation Metrics for Text Quality. In this part, we detail the metrics used to evaluate text quality:

• **ROUGE Score.** For the summarization task, we employ the ROUGE score (Lin, 2004), which measures the overlap of n-grams between the generated summaries and the reference texts to evaluate how effectively the summary captures the essential content.

¹https://github.com/jthickstun/watermark

²https://github.com/xiaoniu-578fa6bff964d005/UnbiasedWatermark

1352	response is generated.					
1353]]	Fext Summarizatio	n	Machine	Translation
1354		BERT Score↑	ROUGE-1↑	Perplexity↓	BERT Score↑	BLEU↑
1955	No Watermark	0.3174±0.0885	0.3772±0.0962	6.4155±3.3009	0.2683±0.1967	10.8705±10.1914
1333	Beta-Reweight (\beta=0)	$0.3162 \pm 0.08/1$	$0.3/58\pm0.0961$	6.3810 ± 3.2753	0.2669 ± 0.1966	10.6208 ± 9.5880
1356	Beta-Reweight (\beta=0.03) Beta Beweight (\beta=0.1)	$0.31/1\pm0.08/7$ 0.3160±0.0873	0.3760 ± 0.0932 0.3762 ± 0.0965	0.3980 ± 3.2142 6 4250±3 2044	0.2083 ± 0.1907 0.2687 ± 0.1062	10.0311 ± 10.1191 10.0058 ±10.5217
1357	Beta-Reweight ($beta=0.1$)	0.3184 ± 0.0873	0.3702 ± 0.0905 0.3771+0.0966	63889+32144	0.2037 ± 0.1902 0.2641+0.1947	10.9053 ± 10.5517 10.9852 ± 10.7563
1358	Beta-Reweight (\beta=0.2)	0.3167 ± 0.0869	0.3764 ± 0.0954	6.3972 ± 3.2855	0.2668 ± 0.1907	10.7865 ± 9.8656
1359	Inverse-sampling	0.3182±0.0876	0.3772±0.0964	6.3377±3.1274	0.2894±0.1869	11.6892±10.5368
1360	Gumbel-reparametrization	0.3171±0.0868	0.3763±0.0961	6.3538±3.2221	0.3065±0.1875	11.8670±10.6599
1361	Soft(δ =0.5)	0.3152±0.0862	0.3746±0.0949	6.4894±3.2453	0.2541±0.1950	10.3546±9.7336
1262	Soft(δ =1.0)	0.3125±0.0856	0.3724±0.0937	6.8647±3.4364	0.2241±0.1922	9.5412±9.0065
1002	Soft($\delta = 1.5$)	0.3067 ± 0.0825	0.3673 ± 0.0917	7.4633±3.5928	$0.18/6\pm0.1891$	8.5556±8.5925
1363	Soli(0=2.0)	0.2990±0.0803	0.3003±0.0899	8.484/±4.1398	0.1380±0.1750	0.9994±0.7328
1364						
1365						
1366	• BLEU score. Fo	or the machine	translation task	, we rely on th	e BLEU score	(Papineni
1367	et al., 2002), emph	asizing the lexi	cal similarity be	tween machine-	generated transl	lations and
1368	human reference t	ranslations.			-	
1369						
1370	 BERTScore. BER 	TScore Zhang e	et al. (2019) calc	ulates the simila	rity between two) sentences
1371	by summing the c	osine similariti	es of their toker	n embeddings. V	Ve utilize BER	ГScore-F1,
1272	BERTScore-Preci	sion, and BERT	ГScore-Recall f	or assessing bot	th text summari	zation and
1070	machine translatio	on tasks.				
1373						
1374	 Perplexity. Perple 	exity, a concept	from informatio	n theory, measu	res how well a j	probability
1375	model or distributi	on predicts a sa	mple. It is used	to compare the	performance of j	probability
1376	models, where a lo	ower perplexity	indicates a more	e predictive mod	lel. We apply pe	rplexity to
1377	evaluate both text	summarization	and text genera	tion tasks.		
1378						
1379	Evaluation Metrics for De	tecting Efficien	cy of Waterma	rks. In this sect	ion, we present t	the metrics
1380	used to evaluate the detecta	bility of waterr	narks:		-	
1381						
1382	• Type I and II Er	rors. We emplo	ov the true posi	tive rate (TPR).	false positive r	ate (FPR).
1202	true negative rate (TNR), and false	e negative rate ()	FNR) to assess y	vatermark detec	tion across
1003	a mix of watermar	ked and non-wa	atermarked sent	ences. The FPR	measures the T	vpe I error.
1304	which occurs whe	n the null hypo	thesis is incorre	ctly rejected wh	nen it is actually	true. The
1385	FNR measures the	e Type II error.	where there is a	failure to reject	a false null hvr	oothesis.
1386		JI		J	51	
1387						
1388	Ε ΔΟΟΙΤΙΟΝΑΙ ΕΧΙ	DEDIMENTAL	RESILTS			
1389	I ADDITIONAL LAI	EKIMENTAL	KESULIS			
1390						
1391	In this section, we introduc	e the additional	experiments co	onducted in our	paper.	
1392	Weakly Distantion From 7	Ch a £-11	-	Table 4 This 6		
1393	to the model without wet	ampariza all w	are presented in	frage waterman	gure snows that	compared
120/	to the model without wat	ermarks, an w	akiy distortion	-free waterman	rks exhibit no i	significant
1005	a significant performance h	initiarization an	u text generatio	Dasidas wa sla	ei, ior me Soft-V	watermark,
1395	a significant performance bi	as is observable	z as v increases.	Desides, we als	o menude a com	prenensive
1396	results for the combination	or all PDA-fule	ts and all three I	kinds of Key san	iping methods	under text
1397	generation tasks. The fesu	ns are presente	u ili table 5. W	e also don't ob	serve the distric	JULIOII DIAS
1398	under the Δ metrics.					
1399	Strongly Distortion-Free.	The full results	s are displayed i	n Table 6, wher	e we include all	PDA-rule
1400	and key sampling method	into compariso	n. From this ta	ble, it is eviden	t that compared	l to the no
1401	watermark model. all weal	kly distortion-fi	ree watermarks	demonstrate pe	erformance bias	across all
1400	tasks. In contrast, the Data	watermark av	hibita laga biga	nominarial to atl	or wookly diet	artian frag

Table 4: Performance of different watermarks under one-time generation. For each prompt, only one 1351

tasks. In contrast, the Beta-watermark exhibits less bias compared to other weakly distortion-free 1402 watermarks. Additionally, as β increases, the distribution bias is further reduced, consistent with our 1403 theoretical analysis.

Table 5: Performance of different watermarks under one-time generation for text generation tasks.For each prompt, only one response is generated

er each promp	, emj ene	respense is g	eneracea					
PDA-rule	Watermark key	bertscore.precision	bertscore.recall	bertscore.f1	ppl	rouge1	rouge2	rougeL
	fixed key set	0.3062±0.0954	0.3279±0.1019	0.3170±0.0880	6.4090±3.2113	0.3764±0.0960	0.1324±0.0808	0.2377±0.0793
β -reweight(β =0)	n-gram hashing	0.3048±0.0949	0.3276±0.1010	0.3162±0.0871	6.3810±3.2753	0.3758±0.0961	0.1314±0.0798	0.2372±0.0785
	position hashing	0.3050±0.0951	0.3271±0.1010	0.3160±0.0874	6.4285±3.2815	0.3759±0.0952	0.1315±0.0798	0.2374±0.0791
	fixed key set	0.3061±0.0953	0.3289±0.1026	0.3174±0.0884	6.3903±3.3533	0.3764±0.0964	0.1327±0.0806	0.2385±0.0801
β -reweight(β =0.05)	n-gram hashing	0.3058±0.0944	0.3286±0.1021	0.3171±0.0877	6.3986±3.2142	0.3760±0.0952	0.1320±0.0797	0.2375±0.0785
	position hashing	0.3058±0.0951	0.3283±0.1021	0.3170±0.0876	6.4043±3.3037	0.3763±0.0959	0.1326±0.0797	0.2385±0.0789
	fixed key set	0.3055+0.0948	0.3279+0.1014	0.3166+0.0873	6.4143+3.3500	0.3765+0.0956	0.1324+0.0795	0.2380+0.0785
β -reweight(β =0.1)	n-gram hashing	0.3054±0.0950	0.3285±0.1015	0.3169±0.0873	6.4250±3.2944	0.3762±0.0965	0.1327±0.0801	0.2377±0.0785
1	position hashing	0.3060±0.0954	0.3285±0.1008	0.3172±0.0875	6.4214±3.2642	0.3762±0.0952	0.1322±0.0785	0.2382±0.0780
	fixed key set	0 3068+0 0952	0.3296+0.1020	0.3181+0.0878	6 4131+3 3820	0 3778+0 0960	0 1337+0 0806	0.2395+0.0799
β -reweight(β =0.2)	n-gram hashing	0.3068+0.0958	0.3302+0.1026	0.3184+0.0883	6.3889+3.2144	0.3771+0.0966	0.1334+0.0811	0.2392+0.0794
	position hashing	0.3057±0.0949	0.3283±0.1025	0.3169±0.0880	6.3685±3.2764	0.3765±0.0963	0.1323±0.0800	0.2383±0.0794
	fixed key set	0 3053+0 0955	0.3280+0.1018	0.3166+0.0878	6 3878+3 1945	0 3763+0 0954	0 1319+0 0799	0.2376+0.0788
β -reweight(β =0.3)	n-gram hashing	0.3052+0.0949	0.3284+0.1006	0.3167+0.0869	6.3972+3.2855	0.3764+0.0954	0.1325 ± 0.0799	0.2379+0.0784
/·····	position hashing	0.3066±0.0952	0.3288±0.1018	0.3176±0.0876	6.3845±3.2077	0.3771±0.0963	0.1327±0.0798	0.2385±0.0787
	fixed key set	0 3011+0 0953	0.3277+0.1016	0 3143+0 0875	6 6430+3 5498	0 3746+0 0959	0 1309+0 0797	0.2361+0.0793
Gumbel-reparametrization	n-gram hashing	0.3060+0.0942	0.3284+0.1011	0.3171+0.0868	6.3538+3.2221	0.3763+0.0961	0.1321+0.0797	0.2376+0.0788
1	position hashing	0.3047±0.0958	0.3267±0.1019	0.3156±0.0881	6.4877±3.4127	0.3755±0.0957	0.1317±0.0800	0.2380±0.0790
	fixed key set	0 3063±0 0942	0 3207±0 1014	0 3170±0 0870	6 1846+3 1150	0 3777+0 0960	0 1334±0 0802	0 2301±0 0703
Inverse-sampling	n-gram hashing	0.3064+0.0953	0.3302+0.1018	0.3182+0.0876	6.3377+3.1274	0.3772+0.0964	0.1328+0.0809	0.2390+0.0799
	position hashing	0.3075±0.0962	0.3326±0.1022	0.3199±0.0881	6.2007±3.0213	0.3796±0.0960	0.1344±0.0813	0.2404±0.0802
No Watermark	NA	0 3058+0 0959	0 3293+0 1026	0 3174+0 0885	6 4155+3 3009	0 3772+0 0962	0 1328+0 0806	0.2388+0.0799
		0.0000±0.0000	0.5255±0.1020	0.0174±0.0000	6.4004.2.2.000	0.5772±0.0902	0.1320±0.0000	0.2300±0.0777
$Soft(\delta=0.5)$	n-gram hashing	0.3013±0.0941	0.3294±0.1005	0.3152±0.0862	6.4894±3.2453	0.3/46±0.0949	0.1310±0.0781	0.2362±0.0776
$Soft(\delta=1.0)$ Soft($\delta=1.5$)	n-gram hashing	0.2930±0.0928	0.3290±0.0999	0.3123±0.0830	$0.804/\pm 3.4304$ 7 4633 ± 3.5028	0.3724 ± 0.0937 0.3673 ± 0.0917	0.1279 ± 0.0769 0.1220±0.0731	0.2328 ± 0.0764 0.2271 ± 0.0724
$Soft(\delta=2.0)$	n-gram hashing	0.2751+0.0879	0.3246+0.0953	0.2996+0.0805	8.4847+4.1598	0.3605+0.0899	0.1158+0.0698	0.2207+0.0695
()		1						
PDA-rules	Watermark key	Δ bertscore.precision	Δ bertscore.re	call Δ bertscor	e.f1 Δ ppl	Δ rouge	1 Δ rouge?	2 Δ rougeL
	fixed key set	0.0694±0.0564	0.0674±0.05	77 0.0625±0.0	0520 2.7242±2.8	964 0.0700±0.0	549 0.0585±0.0	517 0.0606±0.0519
β -reweight(β =0)	n-gram hashing	0.0700±0.0561	0.0672±0.05	67 0.0626±0.0	0513 2.7165±2.9	0.0703±0.0	560 0.0582±0.0	517 0.0605±0.0519
	position nashing	0.0701±0.0363	0.06/9±0.03	75 0.0650±0.0	J518 2.7555±2.5	0.0098±0.0	554 0.0584±0.0	521 0.0611±0.0555
	fixed key set	0.0701±0.0570	0.0678±0.05	69 0.0630±0.0	0519 2.7436±3.0	0.0709±0.0	550 0.0588±0.0	521 0.0617±0.0527
β -reweight(β =0.05)	n-gram hashing	0.0700±0.0567	0.0679±0.05	73 0.0631±0.0	0519 2.7419±2.9	0.0701±0.0	554 0.0583±0.0	517 0.0606±0.0522
	position nashing	0.0703±0.0300	0.0685±0.05	// 0.0651±0.0	0521 2.7540±2.9	0.0715±0.0	560 0.0590±0.0.	524 0.0616±0.0522
	fixed key set	0.0695±0.0566	0.0674±0.05	73 0.0623±0.0	0520 2.7563±3.0	0299 0.0693±0.0	557 0.0580±0.0	520 0.0608±0.0526
β -reweight(β =0.1)	n-gram hashing	0.0696±0.0563	0.0676±0.05	67 0.0626±0.0	0515 2.7640±2.9	0893 0.0701±0.0	558 0.0579±0.0	516 0.0605±0.0520
	position hashing	0.0703±0.0566	0.0676±0.05	/1 0.0630±0.0	2.7559±2.9	446 0.0698±0.0	555 0.0583±0.0	513 0.0610±0.0515
	fixed key set	0.0695±0.0560	0.0673±0.05	70 0.0625±0.0	0512 2.7507±3.0	0.0706±0.0	553 0.0589±0.0	524 0.0610±0.0525
β -reweight(β =0.2)	n-gram hashing	0.0698±0.0566	0.0679±0.05	71 0.0629±0.0	0517 2.7376±2.9	0355 0.0699±0.0	558 0.0589±0.0	525 0.0607±0.0518
	position hashing	0.0699±0.0563	0.0688±0.05	8/ 0.0632±0.0	1526 2.7001±2.9	0.069/±0.0	563 0.0584±0.0	529 0.0608±0.0532
	fixed key set	0.0706±0.0568	0.0680±0.05	75 0.0631±0.0	0520 2.7242±2.9	0031 0.0701±0.0	562 0.0581±0.0	519 0.0608±0.0528
β -reweight(β =0.3)	n-gram hashing	0.0705±0.0564	0.0679±0.05	70 0.0633±0.0	0515 2.7466±2.9	944 0.0701±0.0	552 0.0585±0.0	514 0.0609±0.0527
	position hashing	0.0696±0.0559	0.0673±0.05	65 0.0622±0.0	0510 2.7271±2.9	0034 0.0693±0.0	552 0.0576±0.0	507 0.0602±0.0513
	fixed key set	0.0700±0.0572	0.0679±0.05	78 0.0629±0.0	0524 2.8303±3.0	0.0706±0.0	561 0.0579±0.0	523 0.0616±0.0530
Gumbel-reparametrization	n-gram hashing	0.0694±0.0561	0.0678±0.05	74 0.0625±0.0	0517 2.7221±2.9	0595 0.0708±0.0	555 0.0588±0.0	520 0.0607±0.0524
	position hashing	0.0702±0.0573	0.0682±0.05	85 0.0630±0.0	0530 2.7680±3.0	0.0702±0.0	563 0.0593±0.0	529 0.0615±0.0539
	fixed key set	0.0692±0.0555	0.0661±0.05	64 0.0618±0.0	0508 2.6649±2.8	626 0.0695±0.0	556 0.0580±0.0	516 0.0608±0.0520
Inverse-sampling	n-gram hashing	0.0697±0.0565	0.0674±0.05	67 0.0625±0.0	0516 2.7131±2.8	903 0.0705±0.0	557 0.0581±0.0	521 0.0603±0.0523
	position hashing	0.0704±0.0559	0.0677±0.05	79 0.0628±0.0	0517 2.6266±2.8	3591 0.0698±0.0	559 0.0583±0.0	519 0.0612±0.0526
Baseline	NA	0.0701±0.0560	0.0674±0.05	70 0.0628±0.0	0513 2.7535±2.9	0630 0.0707±0.0	558 0.0583±0.0	522 0.0613±0.0527
Soft(δ=0.5)	n-gram hashing	0.0700±0.0569	0.0677±0.05	76 0.0627±0.0	0519 2.7403±2.9	0348 0.0700±0.0	553 0.0581±0.0	507 0.0606±0.0521
$Soft(\delta=1.0)$	n-gram hashing	0.0692±0.0558	0.0666±0.05	62 0.0616±0.0	0505 2.8607±3.0	0746 0.0688±0.0	543 0.0569±0.0	501 0.0595±0.0511
$Soft(\delta=1.5)$	n-gram hashing	0.0704±0.0564	0.0661±0.05	57 0.0613±0.0	0508 3.0427±3.1	473 0.0688±0.0	550 0.0566±0.0	505 0.0593±0.0516
Soft(d=2.0)	n-gram hashing	0.0736±0.0587	0.0669±0.05	60 0.0635±0.0	0517 3.6349±3.6	0.0699±0.0	552 0.0576±0.0	509 0.0601±0.0517



Figure 5: ROC curve of TPR vs FPR.

Table 6: Performance of different watermarks under multi-time generations. We randomly selected
 1000 prompts and generated 100 responses for each. We use F1 scores of BERTScore and scale
 BERTScore and ROUGE-1 with a factor of 100.

Δ rougeL 0.0062±0.0052 0.0100±0.0093 0.0099±0.0093
0.0062±0.0052 0.0100±0.0093 0.0099±0.0093
0.0100±0.0093 0.0099±0.0093
0.0099±0.0093
0.0063±0.0056
0.0092±0.0082
0.0089±0.0083
0.0061±0.0051
0.0086±0.0078
0.0084±0.0074
0.0062±0.0054
0.0076±0.0066
0.0077±0.0067
0.0060±0.0051
0.0069±0.0058
0.0067±0.0058
0.0069±0.0057
0.0427±0.0362
0.0442±0.0388
0.0065±0.0052
0.0428±0.0363
0.0441±0.0396
0.0060 ± 0.0053
0.0065 ± 0.0056
0.0090±0.0078
0.0127±0.0107
0.0188±0.0149

Table 7: AUC score of different watermarks under varying attack strength ϵ on text generation task.

Beta-Reweight	<i>ϵ</i> =0	ϵ =0.05	<i>ϵ</i> =0.1	<i>ϵ</i> =0.2	<i>ϵ</i> =0.3
β=0	0.9948	0.9901	0.9742	0.8848	0.7447
$\beta = 0.05$	0.9912	0.9846	0.9672	0.8724	0.7312
$\beta = 0.1$	0.9889	0.9785	0.9550	0.8558	0.7078
$\beta = 0.2$	0.9796	0.9598	0.9201	0.7983	0.6735
$\beta = 0.3$	0.9447	0.9047	0.8509	0.7289	0.6191



1509Figure 6: Trade-off between distribution bias and watermark strength under key collision. The TPR1510is measured under 10% (**Top Left**), 5% (**Top Right**), 1% (**Bottom Left**), 0.1% (**Bottom Right**) FPR.1511We can see Δ Perplexity (distribution bias) increase with the TPR.

			z=1.	073	z=1.	224	z=1.	517	z=1.	859
			TNR↑	TPR↑	TNR↑	TPR↑	TNR↑	TPR↑	TNR↑	TPR↑
		$\delta = 0.5$	90.00	46.05	95.00	38.78	99.00	24.41	99.90	13.04
	Soft watermark	$\delta = 1$	90.00	88.37	95.00	85.02	99.00	76.80	99.90	68.42
	Soft-watermark	$\delta = 1.5$	90.00	97.15	95.00	96.65	99.00	94.64	99.90	90.90
		$\delta = 2$	90.00	99.45	95.00	99.39	99.00	99.06	99.90	97.90
_		$\beta = 0$	90.00	97.75	95.00	97.17	99.00	94.69	99.90	90.25
		$\beta = 0.05$	90.00	96.82	95.00	96.19	99.00	92.67	99.90	86.26
	Beta-watermark	$\beta = 0.1$	90.00	95.76	95.00	94.19	99.00	89.13	99.90	79.90
		$\beta = 0.2$	90.00	86.53	95.00	82.49	99.00	71.14	99.90	58.55
		$\beta = 0.3$	90.00	64.59	95.00	56.88	99.00	40.67	99.90	25.38





1534 Figure 7: Left. Trade-off between distribution bias and watermark strength under key collision. The 1535 TPR is measured under 1% FPR. We can see Δ Perplexity (distribution bias) increase with the TPR. **Right.** AUC score of different watermarks under varying attack strength ϵ on text generation task. 1536

1537 F.1 ABLATION STUDY 1538

1539 Detect efficiency. We compare the detection efficiency of beta-watermark with Soft-watermark on 1540 text generation tasks. We set the detecting scaling parameter (Definition 5.5) C = 10. We choose 1541 the threshold z = 1.073, 1.224, 1.517, 1.859, which corresponds to the 10%, 5%, 1% and 0.1% FPR. From Table 8, we see that the detect efficiency of beta watermark is comparable with the Soft-1542 watermark (Kirchenbauer et al., 2023). We also see that when β increases, the detection efficiency 1543 decreases, this is because a larger β introduces a smaller distribution bias into the watermarked 1544 distribution, thus reducing the watermark strength. 1545

1546 We use the beta-watermark to illustrate the trade-off between watermark strength and distribution 1547 bias. As shown in Figure 6, with increasing values of β , the distribution bias decreases, but there is also a corresponding decrease in the true positive rate of watermark detection. This indicates that 1548 reducing the distribution bias of the watermark compromises its detectability. 1549

1550 In Figure 5, we see that the ROC of beta watermark is comparable with the Soft-watermark (Kirchen-1551 bauer et al., 2023). We also see that when β increases, the detect efficiency decreases, this is because 1552 a larger β introduces a smaller distribution bias into the watermarked distribution, thus reducing the 1553 watermark strength.

1554 **Robustness.** We assessed the robustness of the beta-watermark against random text modifications 1555 and GPT-paraphrasing attacks (Kirchenbauer et al., 2023), where we modified 5%, 10%, 20%, and 1556 30% (i.e., $\epsilon = 0.05, 0.1, 0.2, 0.3$) of the tokens. The results, as detailed in Figure F.1 (right), and 1557 Table 9, 10, 11 and 12 indicate that the beta-watermark maintains its robustness with moderate text 1558 modifications.

1560

1512

- 1561
- 1563
- 1564
- 1565

1007	Table 9. Detectability and 10	busilless of Soft v	vatermark	with Deta-	waterman	x under 11	IN WITK
1568	on rand <u>om token modificati</u>	on					
1569	Random paraphras	e FPR=0.1%	$\epsilon=0$	ϵ =0.05	ϵ =0.1	ϵ =0.2	<i>ϵ</i> =0.3
1570	Beta-watermark	β=0	92.10	88.42	86.47	71.21	49.89
1571		β = 0.05	91.78	87.94	84.20	64.80	41.23
1572		$\beta = 0.1$	84.73	82.12	72.41	52.78	29.66
1573		$\beta = 0.2$	70.51	64.52	52.88	33.48	15.74
1574		β =0.3	38.15	28.10	18.62	9.71	3.27
1575	Soft	$\delta = 0.5$	13.59	9.45	5.41	2.53	1.38
1576		$\delta = 1.0$	69.32	61.03	51.62	33.14	15.68
1577		$\delta = 1.5$	92.52	88.78	84.47	69.16	44.33
1578		$\delta = 2.0$	98.35	97.58	96.37	90.65	73.60

Table 9: Detectability and robustness of Soft watermark with beta-watermark under TPR @ FPR=0.1%

Table 10: Detectability and robustness of Soft watermark with beta-watermark under TPR @ FPR=0.01% on random token modification

Random paraphrase	FPR = 0.01%	$\epsilon=0$	ϵ =0.05	ϵ =0.1	ϵ =0.2	<i>ϵ</i> =0.3
Beta-watermark	β=0	88.42	84.2	79.44	62.12	37.23
	β = 0.05	86.40	82.35	75.44	52.08	22.81
	β = 0.1	77.75	71.76	63.03	41.98	17.99
	β = 0.2	56.43	49.22	37.58	19.96	7.10
	β = 0.3	24.38	16.70	8.80	3.39	0.79
Soft	δ=0.5	6.80	4.03	1.96	0.81	0.35
	$\delta = 1.0$	57.78	50.06	38.63	19.71	6.27
	$\delta = 1.5$	88.21	83.56	78.00	56.92	29.59
	$\delta = 2.0$	97.36	95.82	93.07	84.82	63.36

Table 11: Detectability and robustness of Soft watermark with beta-watermark under TPR @ FPR=0.001% on random token modification

	Random paraphrase	FPR = 0.001%	<i>ϵ</i> =0	<i>ϵ</i> =0.05	ϵ =0.1	ϵ =0.2	<i>ϵ</i> =0.3
	Beta-watermark	β=0	84.84	79.44	72.40	50.65	27.71
		β =0.05	81.91	75.00	67.43	43.31	17.65
		$\beta = 0.1$	71.21	63.79	55.29	29.99	12.10
		β=0.2	47.12	40.13	29.60	11.86	2.88
		β=0.3	12.98	9.14	4.85	1.92	0.79
	Soft	δ=0.5	3.23	2.30	1.27	0.35	0.35
		$\delta = 1.0$	47.82	37.96	29.11	10.86	2.80
		$\delta = 1.5$	82.77	77.66	71.09	47.51	19.05
		δ=2.0	96.15	93.62	90.54	78.22	50.83

Table 12: Detectability and robustness of Soft watermark with beta-watermark under TPR @ FPR=0.1% on GPT-4 paraphrase attack

1609	FPR=0.1% on GPT-4 paraphrase attack									
1610	GPT-4 paraphrase	FPR=0.1%	<i>ϵ</i> =0	ϵ =0.05	ϵ =0.1	<i>ϵ</i> =0.2	<i>ϵ</i> =0.3			
1611	Beta-reweight	β=0	92.10	90.62	92.86	88.24	70.59			
1012		$\beta = 0.05$	91.78	89.34	91.55	85.78	75.65			
1613		$\beta = 0.1$	84.73	83.03	79.35	72.00	53.44			
1614		β=0.2	70.51	65.21	62.75	54.98	48.48			
1615		β = 0.3	38.15	32.14	31.32	29.47	17.35			
1616	Soft	$\delta = 0.5$	13 50	11.08	11.60	13.08	7 / 1			
1617	501	$\delta = 0.5$	60.32	65.83	64.37	61.36	55 12			
1618		$\delta = 1.5$	92.52	91.14	92.46	89.38	76.70			
1619		$\delta=2.0$	98.35	97.79	98.58	95.93	95.00			