

# A Watermark for Black-Box Language Models

Anonymous authors

Paper under double-blind review

## Abstract

Watermarking has recently emerged as an effective strategy for detecting the outputs of large language models (LLMs). Most existing schemes require *white-box* access to the model’s next-token probability distribution, which is typically not accessible to downstream users of an LLM API. In this work, we propose a principled watermarking scheme that requires only the ability to sample sequences from the LLM (i.e. *black-box* access), boasts a *distortion-free* property, and can be chained or nested using multiple secret keys. We provide performance guarantees, demonstrate how it can be leveraged when white-box access is available, and show when it can outperform existing white-box schemes via comprehensive experiments.

## 1 Introduction

It can be critical to understand whether a piece of text is generated by a large language model (LLM). For instance, one often wants to know how trustworthy a piece of text is, and those written by an LLM may be deemed untrustworthy as these models can hallucinate. The goal of watermarking is to cleverly bias the LLM so that detecting its generations becomes easier. Most proposed techniques do not modify the underlying LLM’s model weights or its training procedure but rather inject the watermark during autoregressive decoding at inference time. They require access to the next-token logits and inject the watermark every step of the sampling loop. This required access prevents third-party users of an LLM from applying their own watermark as proprietary APIs currently do not support this option. Supporting this functionality presents a security risk in addition to significant engineering considerations. Concretely, [Carlini et al. \(2024\)](#) showed that parts of a production language model can be stolen from API access that exposes logits. In this work, we propose a watermarking scheme that gives power back to the people — third-party users can watermark a language model given nothing more than the ability to sample sequences from it. Our scheme is faithful to the underlying language model and it can outperform existing white-box schemes.

## 2 Related Work

Watermarking outside the context of generative LLMs, which is sometimes referred to as linguistic steganography, has a long history and typically involves editing specific words from an non-watermarked text. Watermarking in the modern era of generative models is nascent — [Venugopal et al. \(2011\)](#) devised a scheme for machine translation, but interest in the topic grew substantially after the more recent seminal works of [Kirchenbauer et al. \(2023a;b\)](#) and [Aaronson \(2023\)](#). Many effective strategies employ some form of pseudorandom functions (PRFs) and cryptographic hashes on token  $n$ -grams in the input text. [Kirchenbauer et al. \(2023a\)](#) proposes modifying the next-token probabilities every step of decoding such that a particular subset of the vocabulary, referred to as *green tokens*, known only to those privy to the secret key, are made more probable. Watermarked text then is expected to have more green tokens than non-watermarked text and can be reliably detected with a statistical test. The scheme distorts the text, but with the right hyperparameters a strong watermark may be embedded with minimal degradation in text quality.

Meanwhile, [Aaronson \(2023\)](#) proposes a clever *distortion-free* strategy which selects the token that is both highly probable and that achieves a high PRF value. [Kuditipudi et al. \(2023\)](#) applies a scheme similar in spirit to [Aaronson \(2023\)](#) but to improve robustness to attacks, pseudorandom numbers (PRNs) are determined

by cycling through a fixed, pre-determined sequence of values called the *key*, rather than by  $n$ -grams. They compute a  $p$ -value using a permutation test to determine if the text was watermarked with *that specific* key.

Concurrent to our work is SYNTHID (Dathathri et al., 2024), a white-box scheme which resembles the recursive variant we propose. Specifically, when the sequence length is 1 and each watermarking party samples 2 sequences, then our encoding procedure overlaps with theirs. Both processes can be visualized via a binary tree where nodes represent tokens. At the very bottom, leaf nodes represent tokens sampled i.i.d. from the next-token distribution. The token for a non-leaf node is decided by applying a specific PRF determined by the node’s level to the two children tokens. Scoring, or decoding, differs however in that while they propose training a Bayesian detector or zero-shot detection by combining PRF values for the unique  $n$ -grams across levels of the tree (equivalently, virtual watermarking parties) using a weighted average (or corresponding  $z$ -scores) that accounts for entropy differences across depths, we suggest combining  $p$ -values using Fisher’s method with equal contributions across the layers. Furthermore, they suggest using a Bernoulli distribution instead of the standard uniform and consider much deeper trees (e.g. a depth of 30) than the ones we do in this work.

Most similar to our work is the WATERMAX algorithm of Giboulot & Furon (2024). At a high level, their scheme is very similar to ours: i.i.d. responses are sampled from the LLM for a specific prompt, pseudorandom values are determined for each token in each sequence based on  $n$ -grams and an aggregate score is computed for each sequence. The sequence with the largest aggregate score is returned as the watermarked text. However, very crucially, their algorithm is not distortion-free, as ours is, and their theory does not capture the critical role of entropy on detectability. Meanwhile, we are able to formally and empirically quantify the dependence on entropy on performance and formulate provably optimal tests. Furthermore, we propose a recursive variant of the algorithm and also carefully study the effect of the distribution used for the PRF.

Lee et al. (2023) adapts Kirchenbauer et al. (2023a)’s scheme for code-generation by applying the watermark only at decoding steps that have sufficient entropy. Zhao et al. (2023) investigates a special case of Kirchenbauer et al. (2023a) for improved robustness to adversarial corruption. Fernandez et al. (2023) tests various watermarking schemes on classical NLP benchmarks and also introduces new statistical tests for detection — most notably, they suggest skipping duplicate  $n$ -grams during testing.

Yang et al. (2023) introduces a scheme that relies on black-box access to the LLM. Their method samples from the LLM and injects the watermark by replacing specific words with synonyms. Although their approach shares the assumption of black-box LLM access, as in our work, it has limitations not present in ours: the watermarking process is restricted to words that can easily be substituted with multiple synonyms, synonym generation is powered by a BERT model (Devlin, 2018), making it computationally expensive, and the scheme is not distortion-free. Chang et al. (2024) presents POSTMARK, a black-box watermarking method that uses semantic embeddings to identify an input-dependent set of words. These words are then inserted into the text by an LLM after decoding. However, this approach is also not distortion-free, as the insertion of words by the LLM often results in significantly longer watermarked text.

Given the weakness of many schemes to paraphrasing or word substitution attacks, some have proposed watermarking based on semantics and other features that would remain intact for common attack strategies (Liu et al., 2023b; Hou et al., 2023; Ren et al., 2023; Yoo et al., 2023). Meanwhile, others have viewed the problem through the lens of cryptography and classical complexity theory (Christ et al., 2023; Christ & Gunn, 2024). Lastly, Liu et al. (2023a) proposes an un-forgable publicly verifiable watermark algorithm that uses two different neural networks for watermark generation and detection. Huang et al. (2023) improves the statistical tests used for detection, providing faster rates than prior work.

As the deployment of watermarks to LLMs is still early and also presumably secretive, the correct threat model is still undetermined. Krishna et al. (2024) shows that paraphrasing can evade both third-party and watermarking detectors alike. Some may posit that attacks like paraphrasing or round-trip translation are unrealistic since either they are too expensive to conduct at scale or parties in possession of a capable paraphrasing model have adequate resources to serve their own LLM. Zhang et al. (2023) show that attackers with weaker computational capabilities can successfully evade watermarks given access to a *quality oracle* that can evaluate whether a candidate output is a high-quality response to a prompt, and a *perturbation oracle* which can modify an output with a non-trivial probability of maintaining quality. Alarming, Gu

et al. (2023) demonstrates that watermarks can be learned — an adversary can use a teacher model that employs decoder-based watermarking to train a student model to emulate the watermark. Thibaud et al. (2024) formulates tests to determine whether a black-box language model is employing watermarking, and they do not find strong evidence of watermarking among currently popular LLMs.

### 3 Algorithm

**High-level sketch.** At a high level, our scheme operates autoregressively; each step, we sample multiple generations from the LLM, score each with our secret key, and output the highest scoring one. We do this repeatedly until our stopping condition (e.g. reaching the stop-token or the max length) is met. To determine whether a piece of text was watermarked, we score it using our key — if it’s high, it’s likely watermarked. We now describe the algorithm more formally.

**Preliminaries.** We begin with some preliminaries. If  $F$  is a cumulative distribution function (CDF), we let  $F[s]$  (square brackets) refer to a single draw from a pseudorandom number generator (PRNG) for  $F$  seeded by integer seed  $s$ . Let  $F_k$  be the CDF for  $\sum_{i=1}^k X_i$ , where  $X_i \stackrel{iid}{\sim} F$ . We sometimes abuse notation and treat a distribution as its CDF (e.g.  $N(0, 1)(2)$  is the standard normal CDF evaluated at 2) and when the context is clear we let  $-F$  be the distribution of  $-X$  where  $X \sim F$ . Now, we detail our proposed algorithm, for which pseudocode is provided in Algorithm 1 (Appendix). A Python implementation and example usage are presented in the Appendix.

Let  $F$  be a *continuous* CDF of our choosing,  $P$  the input prompt,  $K$  a secret integer key known only to the watermark encoder and decoder, LM a conditional language model with vocabulary  $\mathcal{V}$  of size  $V$ , and  $h$  a cryptographic hash function (e.g. SHA-256) from  $\mathbb{Z}^*$  to  $\mathbb{Z}$ . Let  $n$  be the number of tokens (typically 4 or 5) that serves as input to our pseudorandom function. Our PRF  $g : \mathcal{V}^* \rightarrow \mathbb{R}$  is given by  $g(w) = F[h(K|w)]$ , where  $|$  denotes concatenation.

**Watermark encoding.** We sample  $m$  sequences  $\{Q_1, \dots, Q_m\}$ , each consisting of at most  $k$  tokens from  $\text{LM}(\cdot | P; k)$ . Let  $\{(X_1, c_1), \dots, (X_j, c_j)\}$  be the *unique* sequences along with their counts from  $\{Q_i\}$  — for example, the sequence  $X_t$  appears  $c_t$  times in  $\{Q_i\}$ . To score each distinct sequence  $X_t$ , we first extract its  $n$ -grams as  $\{(X_{t,i-n+1}, \dots, X_{t,i})\}_{i=1}^{|X_t|}$ , where we allow the left endpoint to spill over only to earlier-generated tokens and not the original prompt tokens.  $l$ -grams are taken instead for boundary indices with only  $l - 1 < n - 1$  eligible tokens strictly left of it. We compute an integer seed for each  $n$ -gram  $w$ , as  $h(K|w)$ . Given a collection of seeds with their associated sequences, we deduplicate seeds across the collection. We do this by picking one instance of the seed *at random* and remove all remaining instances from the collection. We ensure every sequence has at least one seed by adding a random seed not already used, if necessary. For each sequence  $X_t$ , we iterate through its new seeds  $S_t$  (order does not matter) and compute the quantity  $u_t = F_{|S_t|} \left( \sum_{i=1}^{|S_t|} F[S_{t,i}] \right)$ . Finally, we compute  $i^* = \arg\max_{i=1}^j u_i^{m/c_i}$  and choose  $X_{i^*}$  as our watermarked sequence of length at most  $k$ . To generate longer texts, we run the aforementioned process autoregressively until our stopping condition, where we condition the language model on  $P$  and the tokens generated thus far.

One may notice that the LLM is expected to return at most  $k$  tokens. This choice is made to simplify the analysis. In practice, the API may only return texts, not tokens, with no option to specify max length. The watermarker can generate  $n$ -grams from the responses however they would like (with custom tokenization or not). Furthermore, there is no constraint on  $k$ ;  $k$  can be set adaptively to the max length in each batch of returned responses. The main consideration though is smaller  $k$  begets a stronger watermark, so if the adaptive  $k$  is too large, detectability will suffer.

**Watermark detection.** We treat detection as a hypothesis test, where the null  $\mathcal{H}_0$  is that the query text was *not* watermarked with our scheme and secret key and the alternative  $\mathcal{H}_1$  is that it was. While Bayesian hypothesis testing could be used, this would require choosing priors for both hypotheses, which could be challenging and a poor choice could lead to terrible predictions. Let  $X$  be the query text. Akin to the encoding process, we extract  $W$ , the set of *unique*  $n$ -grams from  $X$ , permitting smaller one near the left boundary. For each  $n$ -gram  $w_t$  we compute  $R_t = F[h(K|w_t)]$ . Under  $\mathcal{H}_0$  (assuming that the test  $n$ -grams

are independent),  $R_t \stackrel{iid}{\sim} F$ , so  $\sum_{t=1}^{|W|} R_t \sim F_{|W|}$  giving a  $p$ -value  $p = 1 - F_{|W|} \left( \sum_{t=1}^{|W|} R_t \right)$ . Our detection score  $s$  is  $1 - p$  (higher means more likely to be watermarked).

Another way to compute a  $p$ -value is to compute token-level  $p$ -values and, assuming they are independent, combine them using Fisher’s method. This way,  $p = 1 - \chi_{2|W|}^2 \left( -2 \sum_{t=1}^{|W|} \log(1 - F(R_t)) \right)$ . Furthermore, tests that incorporate the alternative distribution can be used — the best example being the likelihood ratio test:  $s = \sum_{t=1}^{|W|} (\log f_1(R_t) - \log f_0(R_t))$ , where  $f_0$  and  $f_1$  are the densities of  $R_t$  under  $\mathcal{H}_0$  and  $\mathcal{H}_1$  respectively. For some choices of  $F$  and under some assumptions,  $f_1$  may be written explicitly. In other cases, one can *estimate*  $f_1$  by logging values of  $R_t$  for the watermarked sequence as the encoding is run live or via simulation and then building a kernel density estimator. We consider these alternative detection strategies later for ablative purposes.

**Recursive watermarking.** Since our scheme requires only a black box that samples sequences, it can be applied iteratively or recursively. Consider the following. User 1 uses User 2’s LLM service who uses User 3’s LLM service, so on so forth until User  $t$ . Our scheme allows User  $i$  to watermark its service with its secret key  $K_i$ . Each user can then run detection using its key oblivious to whether other watermarks were embedded upstream or downstream. Furthermore, the users can cooperate in joint detection by sharing only  $p$ -values without revealing their secret keys.

Consider the special case that all users are actually the same entity in possession of  $t$  distinct keys  $(K_1, \dots, K_t)$ . Then the iterative watermarking becomes a recursive one, where  $K_i$  is used to watermark the result of watermarking with keys  $(K_{i+1}, \dots, K_t)$ . The entity can run DETECT to get a  $p$ -value for each key and these  $t$   $p$ -values can subsequently be combined using Fisher’s method. We present this recursive scheme in Algorithm 2 (Appendix).

**White-box watermarking.** In the case of  $k = 1$ , our scheme can be efficiently run for users who have white-box access — with the next-token distribution in hand, one can sample a large number of candidate tokens without any inference calls to the model.

**Extensions.** At its crux, the proposed scheme samples sequences of text from a service, divides each unique sequence into a bag of units (namely  $n$ -grams) where each unit is scored using a PRF and the scores are combined in an order-agnostic way. The strength of the watermark depends on the number of *distinct* units across the candidate sequences and the robustness depends on how many of the units are kept intact after the attack. Although any symmetric monotone function can be used instead of the simple summation of the PRNs for each unit, we do not see any compelling reason to make our algorithm more general in this way. However, we briefly highlight some other possible extensions.

*Beam search.* Rather than drawing i.i.d. samples from the model, one can apply our watermark selection to the sequences that arise from beam search, with the caveat that this would violate our distortion-free property.

*Semantic watermarking.* Rather than use  $n$ -grams, the watermarker can extract a set of meaningful semantic units for each sampled text. Robustness may be improved as these units will largely remain intact under an attack like paraphrasing. On the other hand, many of the sampled sequences will have the same *meaning*, so there may be a lot of duplicate units across the candidate sequences, which would degrade the watermark strength.

*Paraphrasing.* Thus far, we assumed the service provides  $m$  draws from the LLM. If  $m$  is large, this can be prohibitively expensive. The resource-constrained may consider the following alternative: draw one sample from the LLM and feed it to a much cheaper paraphrasing model to generate  $m$  paraphrases. The downside is that there may be a lot of duplicate  $n$ -grams across the candidate set.

## 4 Theory

Our goal here is to show that our scheme is faithful to the model’s next-token distribution and to give detection performance guarantees. All proofs are in the Appendix.

**Theorem 4.1** (Distortion-free property). *Let  $X$  be any finite sequence and  $P$  any prompt. Let  $X_u \sim \text{LM}(\cdot | P)$  be the non-watermarked output of the conditional autoregressive language model. Let  $X_w$  be the output of the watermarking procedure (WATERMARK in Algorithm 1, for both recursive and non-recursive settings) for the same prompt and model and any choice of remaining input arguments with the constraint that  $F$  is a continuous distribution. Furthermore, assume that the deduplicated seeds (determined by hashing the secret key and  $n$ -grams) across sequences, are conditionally independent given the counts of the sampled sequences. Then,  $\mathbb{P}(X_u = X) = \mathbb{P}(X_w = X)$ .*

Theorem 4.1 tells us that sampling tokens using our proposed scheme is, from a probabilistic perspective, indistinguishable from sampling from the underlying model, with the caveat that the unique seed values are conditionally independent given the counts of sequences. If we dismiss hash collisions as very low probability events, then since the key is fixed, this reduces to the assumption that unique  $n$ -grams across the sampled sequences are independent. How strong of an assumption this is depends on many factors such as  $m$ , the underlying distribution, and the counts  $(c_1, \dots, c_j)$  themselves. One can construct cases where the assumption is reasonable and others where it is blatantly violated (e.g. if  $n$ -grams within a sequence are strongly correlated). One direction to making the assumption more palatable is to draw a fresh keys i.i.d. for *each* hash call. This would obviously destroy detectability. As a trade-off, one can leverage a set of secret keys (i.e. by drawing keys uniformly at random from a key set), which may reduce distortion, but will hurt detection as each key in the set needs to be tested against.

**Theorem 4.2** (Lower bound on detection ROC-AUC). *Consider the specific case of using flat (i.e. non-recursive) watermarking with  $k = 1$  and  $F = U(0,1)$ . Let  $s_0$  be the score under null that the  $T$  test tokens<sup>1</sup>, assumed to be independent, were generated without watermarking and  $s_1$  be the score if they were. With  $c_i$  as the number of times vocabulary token  $i$  was sampled, we have the following lower bound on the detector’s ROC-AUC.*

$$\mathbb{P}(s_1 \geq s_0) \geq \frac{1}{1 + 1/(3T\lambda^2\alpha^2)}, \text{ where}$$

$$\lambda = \frac{1}{\log(m)} \left( \frac{m}{m+1} - \frac{1}{2} \right) \text{ and } \alpha = \mathbb{E}_c \left[ - \sum_{i=1}^V \mathbf{1}[c_i > 0] \frac{c_i}{m} \log \left( \frac{c_i}{m} \right) \right].$$

$\alpha$  represents the average Shannon entropy in the sampled next-token distribution.

Theorem 4.2 connects detection performance to the language model’s underlying distribution, number of sampled tokens  $m$ , and number of test samples  $T$ . More entropy and more test samples guarantee higher performance. When the model is extremely confident,  $\alpha \rightarrow 0$  and so does our lower bound. Note that because  $\alpha$  measures the entropy of the empirical distribution arising from *sampling* tokens, it depends on both the underlying next-token probability distribution as well as  $m$ . Concretely, when conditioned on the next-token probabilities  $p$ ,  $c \sim \text{Multinomial}(m, p)$ . The largest  $\alpha$  is achieved when the nonzero  $c_i$ ’s are 1, which can occur when the underlying distribution is uniform (maximal uncertainty) and/or  $m$  is not large. In this case,  $\alpha \rightarrow \log(m)$  and our bound goes to  $1 / \left( 1 + 1 / \left( 3T \left( \frac{m}{m+1} - \frac{1}{2} \right)^2 \right) \right)$ . This quantity has very sharp diminishing returns with respect to  $m$ , so there may be little value in increasing  $m$  beyond a certain point. When  $m \rightarrow \infty$ , the bound goes to  $1 / (1 + 4 / (3T))$ , which increases very quickly with  $T$ . A mere 50 test tokens guarantees at least 97% ROC-AUC. We study the interplay of the various factors on our lower bound more carefully in the Appendix.

The intuitions here carry over to other choices of  $F$  and  $k > 1$ , though formal bounds can be tricky to obtain because of difficulty quantifying the alternative distribution. The null distribution is easy —  $p$ -values are  $U(0,1)$  under  $\mathcal{H}_0$ , and as a result, we have a straightforward equality on the false positive rate.

**Theorem 4.3** (False positive rate). *No matter the choice of watermarking settings, assuming that the unique test  $n$ -grams are independent, we have the following equality on the false positive rate of DETECT (Algorithm 1), using decision threshold  $t$ .*

$$\text{FPR} = \mathbb{P}_{\mathcal{H}_0}(s > t) = 1 - t.$$

<sup>1</sup>To be more precise, these are  $T$  unique  $n$ -grams.



This also holds for DETECTRECURSIVE (Algorithm 2) if we further assume the  $p$ -values across secret keys are independent.

Selecting distinct independent secret keys  $(K_1, \dots, K_t)$  (and ignoring hash collisions that arise across calls to DETECT within DETECTRECURSIVE), will help attain the necessary independence.

Although the alternative score distribution is generally intractable, with the strong assumption that there are no duplicate  $n$ -grams across the candidate sequences, then for a special choice of  $F$ , we can write the alternative in closed form, formulate the optimal detection test and write down false-positive and false-negative rates exactly.

**Theorem 4.4** (Optimal detection for Gamma). *Assume that candidate sequences are unique with length  $k$  and that the  $n$ -grams are independent and contain no duplicates. Suppose we choose  $F = -\text{Gamma}(1/k, \beta)$  (flat scheme), for any rate parameter  $\beta$ . Let  $F_0 = F$  with pdf  $f_0$ ,  $F_1 = -\text{Gamma}(1/k, m\beta)$  with pdf  $f_1$ , and  $R$  the PRF values of the  $T$  test tokens (unique  $n$ -grams), assumed to be independent. Then,  $\forall i, R_i \stackrel{iid}{\sim} F_0$  under the null that the text was watermarked using our procedure and  $R_i \stackrel{iid}{\sim} F_1$  otherwise. The uniformly most powerful test is the log-likelihood ratio test (LRT) with score*

$$s(R) = \sum_{i=1}^T \log \frac{f_1(R_i)}{f_0(R_i)}.$$

Furthermore, for any decision threshold  $t$  on score  $s$ , we have that:

$$\begin{aligned} \text{FPR (Type-I error)} &= \mathbb{P}_{\mathcal{H}_0}(s > t) = \text{Gamma}(T/k, \beta)(Q(t)), \text{ and} \\ \text{FNR (Type-II error)} &= \mathbb{P}_{\mathcal{H}_1}(s \leq t) = 1 - \text{Gamma}(T/k, m\beta)(Q(t)), \text{ where} \\ Q(t) &= \frac{T \log(m)/k - t}{(m-1)\beta}. \end{aligned}$$

In the Appendix, we use Theorem 4.4 to study the impact of  $k$ ,  $m$ , and  $T$  on TPR at fixed FPR. For example, with  $T = 100$ ,  $k = 50$ ,  $m = 64$ ,  $\beta = 1$ , we can achieve 99.9% TPR at 1% FPR.

For other choices of  $F$ , we can estimate  $f_1$  via simulation. If we assume candidate sequences have the same length  $k$  with no duplicate  $n$ -grams, then we can fill an  $m \times k$  matrix with i.i.d. draws from  $F$  and pick the first element of the row with the largest row-sum (among the  $m$ ). We do this until we have sufficiently large (e.g. 10,000) samples from  $f_1$ . We apply a Gaussian kernel-density estimator where the bandwidth is chosen using Scott’s rule (Scott, 2015) to estimate  $f_1(r)$  for test value  $r$ . Despite having  $f_0$  in closed-form, for consistency, we can also estimate it non-parametrically by drawing from  $F$ .

## 5 Experiments

In this section, we compare the performance of our scheme with that of prior work.

### 5.1 Models, Datasets, and Hyperparameters

**Models and Datasets.** Our main model and dataset is MISTRAL-7B-INSTRUCT (Jiang et al., 2023) hosted on Huggingface<sup>2</sup> with bfloat16 quantization, and databricks-dolly-15k<sup>3</sup> (Conover et al., 2023), an open source dataset of instruction-following examples for brainstorming, classification, closed QA, generation, information extraction, open QA, and summarization. We use prompts from the brainstorming, generation, open QA (i.e. general QA), and summarization categories, whose human responses are at least 50 tokens long (save one example, which was removed because the prompt was extremely long). For each of the 5233 total prompts, we generate two non-watermarked responses — a stochastic one using temperature 1, and the greedy / argmax decoding — along with a watermarked one for each scheme. We always force a minimum (maximum) of

<sup>2</sup><https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.1>

<sup>3</sup><https://huggingface.co/datasets/databricks/databricks-dolly-15k>

250 (300) new tokens by disabling the stop token for the first 250 tokens, re-enabling it, and stopping the generation at 300, regardless of whether the stop token was encountered. To simulate real-world use, we de-tokenize the outputs to obtain plain text, and re-tokenize them during scoring. We study performance as a function of token length  $T \leq 250$  by truncating to the first  $T$  tokens.

For completeness, we also present the key results when GEMMA-7B-INSTRUCT (Team et al., 2024)<sup>4</sup> with bfloat16 quantization is applied to the test split of *eli5-category*<sup>5</sup>. Prompts are formed by concatenating the *title* and *selftitle* fields. Only examples with non-empty *title* and whose prompt contains a `?` are kept — for a total of 4885 examples.

**Hyperparameters.** We consider the following choices of CDFs  $F / F_k$ . **(1)**  $F = U(0, 1)$  and  $F_k = \text{IrwinHall}(k)$ . **(2)**  $F = N(0, 1)$  and  $F_k = N(0, k)$ . **(3)**  $F = -\text{Gamma}(1/k, 1)$  and  $F_k = -\text{Exp}(1)$ . **(4)**  $F = \chi^2_2$  and  $F_k = \chi^2_{2k}$ . We set  $n = 4$ .

## 5.2 Evaluation Metrics

We evaluate performance using three criteria. For all metrics and methods, standard errors are computed by running watermark generation multiple times, keeping the secret keys fixed.

**Detectability.** How well can we discriminate between non-watermarked and watermarked text? We choose non-watermarked text to be text generated by the same model, just without watermarking applied during decoding. There are three reasons for choosing the negative class in this way. Firstly, it makes controlling for text length easier as we can generate as many tokens as we do for watermarked samples — in contrast, human responses are of varying lengths. Secondly, watermarked text has far more token /  $n$ -gram overlap with its non-watermarked counterpart than the human reference, which makes detection more challenging. Lastly, since one intended use case of our scheme is for third-party users of a shared LLM service, users may want to distinguish between their watermarked text and non-watermarked text generated by the same LLM service.

Our primary one-number metric is ROC-AUC for this balanced binary classification task. Since performance at low FPR is often more useful in practice, we report the partial ROC-AUC (pAUC) for  $\text{FPR} \leq \text{a target FPR}$  (taken to be 1%), which we find to be more meaningful than TPR at the target FPR. We look at performance as a function of length by truncating the positive and negatives samples to lengths  $\{25, 50, 75, 100, 150, 200, 250\}$ . To understand aggregate performance, we pool all different length samples together and compute one ROC-AUC. Here, it is paramount that the detection score be length-aware to ensure that a single decision threshold can be used across lengths.

**Distortion.** Our scheme, along with most of the baselines, boasts a *distortion-free* property. This property comes with assumptions that are often violated in practice, for example by reuse of the secret key across watermarking calls. We would like to quantify the amount of distortion actually induced by the watermark but doing so is no straightforward task.

A meaningful metric is the KL-divergence between the true non-watermarked next-token distribution and the watermarked one conditioned on the same context, where the former is used as the reference, for all contexts that occur in the non-watermarked generation. This is challenging because while the explicit watermarked distribution is easily derived for Kirchenbauer et al. (2023a)’s scheme which explicitly tweaks the probability of each token in the vocabulary, this is not the case for our scheme; we could *estimate* it via simulations at each step but this would be computationally expensive. As an alternative, we measure the *perplexity* of the watermarked text under the LLM sans watermarking. The caveat is that while perplexity is often used in the literature to measure text quality, the goal here is not necessarily to minimize it (which greedy decoding would do) but to verify that watermarking does not materially change perplexity levels compared to that of non-watermarked randomly sampled text. We also report *likelihood* as the log-probabilities used in calculating perplexity can over-emphasize outliers.

**Quality.** Watermarking may distort the text per the model, but does the distortion tangibly affect the *quality* of the text? Quality can be challenging to define and measure — one proxy is likelihood under a

<sup>4</sup><https://huggingface.co/google/gemma-7b-it>

<sup>5</sup>[https://huggingface.co/datasets/rexarski/eli5\\_category](https://huggingface.co/datasets/rexarski/eli5_category)

much larger model than the generator. Alternatively, one can run standard benchmark NLP tasks and use classic metrics like exact match, etc. We instead opt for using Gemini-1.5-Pro as an LLM judge and compute pairwise win rates for each watermark strategy against no watermarking (greedy decoding). We do this in two ways for each scheme — (1) we compute win rates using a single response for each prompt and (2) we first ask the LLM judge to pick the best of 3 responses for each prompt and compute win rates using the best response. (2) represents the common practice of sampling a few generations from the LLM and selecting the best one using some criterion. It captures diversity, as methods that can express an answer in a few different good ways will have an advantage. A caveat with win rates is that they may not reflect the *degree* by which one method is better or worse. For instance, if one strategy’s output was always *marginally* worse than no watermarking, the win rate would be 0% — the same as if it were *much* worse.

### 5.3 Adversarial Attacks

An adversary in possession of watermarked text (but who lacks knowledge of the secret key) may try to evade detection. We study how detectability degrades under two attack strategies — *random token replacement* and *paraphrasing*.

**Random token replacement.** Here, we take the watermarked tokens and a random  $p$ -percent of them are corrupted by replacing each with a random different token.  $p$  is taken to be  $\{10, 20, 30, 40, 50\}$ . This attack strategy is cheap for the adversary to carry out but will significantly degrade the quality of the text.

**Paraphrasing.** In this attack, the adversary attempts to evade detection by paraphrasing the watermarked text. We use Gemini-1.5-Pro to paraphrase each non-truncated watermarked generation. Details are deferred to the Appendix.

### 5.4 Baselines

The watermark schemes we consider here operate token-by-token in the autoregressive decoding loop. Let  $p$  be the next-token probability distribution. Higher detection scores indicate higher confidence that the query text is watermarked.

**Aaronson (A).** Aaronson (2023) computes a PRN for each token  $i$  in the vocabulary as  $u_i = U(0, 1)[h(K|w|i)]$ , where  $w$  is the preceding  $(n - 1)$ -gram,  $K$  is the secret key and  $h$  is a cryptographic hash. Token  $i^*$  is selected, where  $i^* = \operatorname{argmax}_i u_i^{1/p_i}$ . At test time,  $n$ -grams  $\{w_i\}_{i=1}^T$  are extracted from the query test and the detection score  $s_A = -\sum_{i=1}^T \log(1 - R_i)$ , where  $R_i = U(0, 1)[h(K|w_i)]$ .  $n$  is set to 4. This choice strikes a good balance between generation quality / diversity and robustness to attacks. The scheme boasts a *distortion-free* property, but the generated text is a deterministic function of the prompt — i.e. only one generation is possible conditioned on a particular prompt.

**Remark.** If  $k = 1$  and  $F = U(0, 1)$ , then our watermark encoding can be viewed as a stochastic version of Aaronson (2023)’s. As  $m \rightarrow \infty$ ,  $c_t/m \xrightarrow{a.s.} p_t$ , where  $p_t$  and  $c_t$  are the probability and observed occurrences of token  $t$ .

**Aaronson Corrected (AC).** Aaronson (2023)’s detection score  $s_A$  is not length-aware and consequently a single decision threshold across scores involving various lengths results in poor performance, as we later show. Observing that  $s_A$  is a sum of (negative) log  $p$ -values,  $s_A \sim \text{Gamma}(T, 1)$ , or equivalently,  $2s_A \sim \chi_{2T}^2$  under the null that all test tokens are non-watermarked. We propose the new *corrected* detection score,  $s_{AC} = \text{Gamma}(T, 1)(s_A) = \chi_{2T}^2(2s_A)$ . For completeness we also experiment with a  $p$ -value computed in the way we do for our method — concretely as,  $\text{IrwinHall}(T) \left( \sum_{i=1}^T R_i \right)$ . Note that both transformations are monotonic so they have no effect on ROC-AUC when  $T$  is *fixed*.

**Kirchenbauer (KB).** Kirchenbauer et al. (2023a) uses the  $n$  previous tokens to pseudorandomly partition the vocabulary for the next token into two lists: a green list of size  $\gamma V$  and a red list consisting of the remainder. A positive bias of  $\delta$  is added to the logits of the green list tokens while those of the red list are left unchanged. This has the effect of modifying  $p$  so that green list tokens are more probable. The score for a text consisting of  $T$  tokens,  $T_g$  of which were found to be green is,  $s_{KB} = (T_g - \gamma T) / \sqrt{T\gamma(1 - \gamma)}$ .



		PPL	WR	WR (3)	AUC	pAUC	C. AUC	C. pAUC	P. AUC	P. pAUC
Max Std. Error		0.03	-	-	0.1	0.3	0.2	0.3	-	-
Greedy Decoding		1.37	-	-	-	-	-	-	-	-
	Random Sampling	3.50	49.6	65.3	-	-	-	-	-	-
Aaronson		2.81	45.3	45.3	71.7	65.5	65.6	60.3	53.9	50.5
Aaronson Cor.		2.81	45.3	45.3	97.9	83.6	94.8	73.2	58.8	50.7
Kuditipudi		3.55	50.3	67.3	87.8	76.6	87.2	74.4	75.9	53.2
Kirchenbauer	0.5	3.39	49.6	66.6	73.2	52.0	71.0	51.4	49.0	49.8
	1	3.37	50.1	67.0	86.9	60.6	83.7	57.1	52.9	49.9
	2	3.69	47.9	64.1	97.0	83.3	95.4	77.4	58.4	50.3
	3	4.67	41.5	58.4	99.3	94.4	98.6	90.9	63.4	51.5
	4	5.81	26.0	41.2	99.8	98.4	99.6	96.8	66.4	52.7
Flat ( $k = 1$ )	2	3.46	50.0	66.4	90.2	68.8	82.0	58.7	50.5	50.3
	4	3.36	50.8	67.0	95.8	82.9	90.3	70.5	51.3	50.6
	16	3.20	47.7	64.5	97.7	89.7	93.9	79.1	52.7	51.1
	32	3.06	48.4	65.3	97.8	90.2	94.2	80.0	53.0	50.8
	512	2.63	47.7	62.5	97.7	90.0	94.1	79.7	54.6	51.3
1024		2.61	47.7	62.2	97.7	90.0	94.0	79.7	52.8	51.1
Flat ( $k = 10$ )	2	4.10	46.1	62.2	83.4	55.8	73.6	52.0	49.0	50.0
	4	4.06	45.2	61.5	93.8	72.7	85.7	59.4	51.3	50.3
	16	3.86	44.6	60.6	97.8	87.0	93.1	73.5	54.3	50.7
	32	3.80	43.0	60.8	98.2	89.0	94.0	76.7	55.0	50.8
Flat ( $k = 50$ )	2	3.79	48.5	64.2	69.6	50.7	62.2	50.3	47.0	50.0
	4	3.76	47.7	63.9	82.9	53.5	71.9	51.3	49.4	50.0
	16	3.72	48.3	64.2	92.7	66.7	83.1	55.6	50.5	50.1
	32	3.67	47.3	63.9	94.2	71.6	85.5	58.1	51.1	50.5
Rec. ( $k = 1$ )	4	3.41	49.0	65.0	93.4	75.5	86.3	63.2	48.4	50.4
	16	3.33	49.2	66.2	95.4	82.9	90.6	71.8	53.4	50.8
	32	3.29	48.4	64.3	96.3	85.0	91.6	73.5	49.4	50.8
	512	3.05	48.3	64.5	97.2	87.9	92.6	76.5	50.4	51.2
Rec. ( $k = 10$ )	4	4.13	45.7	61.3	88.6	61.7	78.8	53.8	48.0	50.0
	16	4.13	43.7	59.7	93.4	74.0	86.8	61.6	52.9	50.4
	32	4.06	42.9	59.5	94.8	76.9	88.1	63.2	50.6	50.3
Rec. ( $k = 50$ )	4	3.79	48.2	63.8	74.2	51.2	65.1	50.5	46.5	49.9
	16	3.77	47.0	64.0	81.2	54.5	73.3	51.9	51.4	50.2
	32	3.79	47.2	63.3	83.3	55.7	74.4	52.2	49.4	50.0

Table 1: Main table of results, showing our black-box scheme (*Flat*) and its recursive variant (*Rec.*) for various  $k$ 's and  $m$ 's, along with baselines. PPL, WR and WR (3) refer to perplexity, win rate of a single response, and win rate of the best-of-3 responses respectively. pAUC is ROC-AUC up to max FPR of 1%. C and P stand for 10% corruption and paraphrasing attack. For paraphrasing, target lengths of {150, 200, 250} are used in calculating AUC and pAUC as performance is essentially random on shorter lengths. The standard errors are quite small and the maximum across rows is shown for each column. AUCs, pAUCs, and their standard errors are scaled by 100.

We incorporate the latest updates to the algorithm<sup>6</sup> such as including the current token in the  $n$ -gram and skipping duplicate  $n$ -grams at test time. We set  $n = 4$ ,  $\gamma = 0.25$ , and  $\delta \in \{0.5, 1, 2, 3, 4\}$ .

**Kuditipudi (K).** A drawback of using the last  $n$  tokens as a basis for the PRF is that changing just one of them changes the output and hurts detection. Kuditipudi et al. (2023) addresses this limitation as follows. Consider a secret, finite ordered list of seeds of length  $k$ . Start watermarking by selecting a position in the

<sup>6</sup><https://github.com/jwkirchenbauer/lm-watermarking>

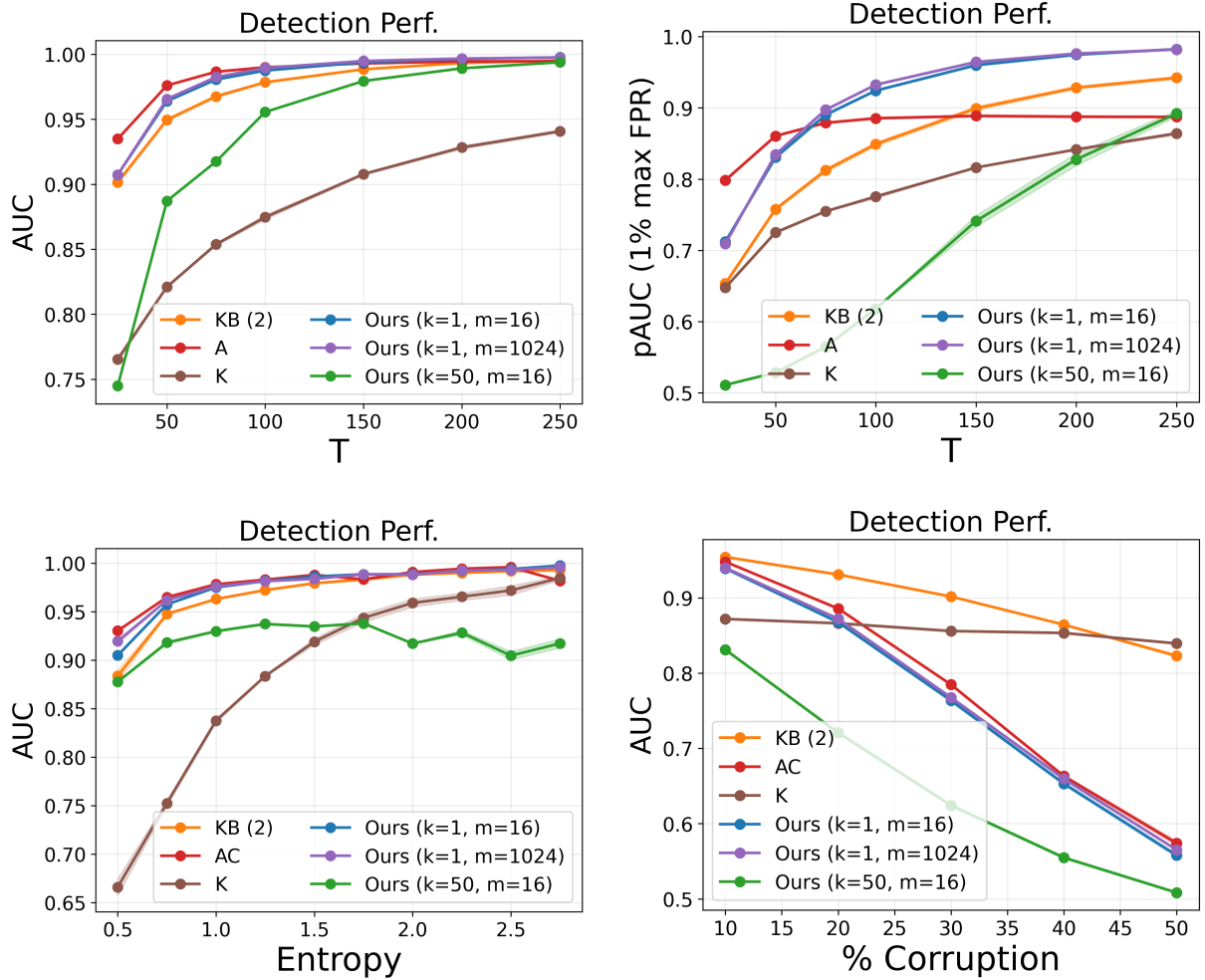


Figure 1: Performance of our flat scheme. **Top:** Detection AUC and pAUC with 1% max FPR for a range of target text lengths when there is no corruption. **Bottom Left:** AUC (mixed  $T$ 's) as a function of the average non-watermarked response entropy of the examples used in the calculation.  $x$ -coordinate  $x$  corresponds to the bucket of examples whose entropy is between  $[x - 0.25, x]$  nats. **Bottom Right:** Effect of amount of random token corruption on AUC (mixed  $T$ 's).

seed list uniformly at random and apply the selection rule of Aaronson (2023) with the PRNG seeded to the current value. Advance to the next seed in the list (wrap-around if you are at the end) and repeat. Scoring is done by conducting a permutation test evaluating how compatible the query text is with the specific list of seeds used during encoding as opposed to any other random list of seeds of the same length. As the random starting position is not known during scoring, an alignment score based on the Levenshtein distance is given that considers alignments of various subsequences of the text and seeds. The proposed method is quite similar to Aaronson (2023) with the difference of using a fixed list of seeds (instead of input tokens to determine the seed) and using a permutation test for scoring. The upside is robustness to token substitution attacks; the downside is significantly higher computational cost for scoring. Larger  $k$  offers more diversity and quality in generation but comes with costlier and weaker detection. The scheme is distortion-free. Following their work, we let  $k = 256$  and accelerate the permutation test by pre-computing 5000 reference values for the secret list using snippets from the train set of *C4-realnewslike* (Raffel et al., 2019) at the various target lengths we evaluate on.

## 5.5 Experimental Results

Table 1 shows results for baselines and our scheme using  $F = U(0, 1)$  and  $p$ -values for scoring, as detailed in Algorithms 1 and 2. For the recursive scheme, depth is  $\lg(m)$  (i.e.  $m = 2$  for each virtual watermark). Here, the negative class is non-watermarked argmax/greedy generations. Results for using stochastic (temperature 1) generations as the negative as well as the average likelihood scores are presented in Table 5 (Appendix); the trends remain the same. We summarize our observations on MISTRAL-7B-INSTRUCT on *databricks-dolly-15k*, which also hold for GEMMA-7B-INSTRUCT on *eli5-category* (presented in the Appendix).

### 5.5.1 Overall performance of our flat and recursive schemes

**Our scheme is a competitive option for white-box watermarking.** Is it better to use our method or alternatives in the white-box setting? When  $k = 1, m = 1024$ , we are able to achieve lower perplexity (2.61 vs. 2.81), better diversity (62.2% vs. 45.3% on best-of-3 win rates) than and comparable detection performance to Aaronson (2023). Furthermore, it has lower perplexity (2.61 vs. 3.55) and detection performance (97.7% vs. 87.8% AUC) than Kuditipudi et al. (2023). By cranking up  $\delta$ , Kirchenbauer et al. (2023a) can achieve strong detection but at the expense of perplexity. When matched on perplexity, we achieve better detection. For example,  $\delta = 0.5$  achieves 3.39 PPL and 73.2% AUC compared to our 2.61 PPL and 97.7% AUC. GEMMA-7B-INSTRUCT on *eli5-category* with  $k = 1, m = 1024$  outperforms Kuditipudi et al. (2023) and is on-par with Aaronson (2023) (see Appendix). Kirchenbauer et al. (2023a) with  $\delta = 0.5$  gives 1.649 PPL and 61.6% AUC whereas  $k = 1, m = 1024$  gets us 1.610 PPL with 93.2% AUC and even 1.645 PPL with 89.7% AUC when  $k = 50, m = 16$  (black-box).

**Flat watermarking outperforms recursive.** Across metrics and settings we see that the flat scheme outperforms its recursive counterpart, suggesting it is more effective when a strong signal is embedded using a single key rather than when multiple weak signals are embedded with different keys. For example, when  $k = 1, m = 32$  flat (recursive) PPL and AUC are 3.06 (3.29) and 97.8% (96.3%) respectively.

### 5.5.2 Effects of Hyperparameters

**Increasing  $m$  improves perplexity but hurts diversity.** Across  $k$ ’s, we observe that perplexity decreases as  $m$  increases, but that win rates, especially when best-of-3 generations are used, decrease. For example, when  $k = 1$ , increasing  $m$  from 2 to 1024 decreases perplexity from 3.46 to 2.61 but also drops the best-of-3 win rate from 66.4% to 62.2%. As remarked earlier, as  $m \rightarrow \infty$ ,  $c_t/m \rightarrow p_t$  and our scheme becomes less diverse — deterministic conditioned on the prompt, like Aaronson’s. On the flip side, large  $m$  reduces sampling noise which drives down perplexity.

**Increasing  $m$  improves detection but has diminishing returns.** Across the board we see that detection improves as  $m$  increases, but there are diminishing returns. For example, when  $k = 1$ , our AUC increases from 90.2% to 95.8% as  $m$  goes from 2 to 4, but flattens out when  $m$  hits 16. This corroborates our theoretical intuition from Theorem 4.2 which is further explored in Figure 4 (Appendix).

**For fixed  $m$ , increasing  $k$  hurts detection performance.** For fixed  $m$  and target generation length  $T$ , increasing  $k$  gives us fewer opportunities (fewer calls to WATERMARKSINGLE) to inject the watermark signal, and detection consequently suffers. For example, when  $m = 32$ , AUC drops from 97.8% to 94.2% when  $k$  increases from 1 to 50.

**$U(0, 1)$  slightly outperforms alternative distributions. Flat distributions may offer better robustness to attacks.** In Table 4 (Appendix), we see that  $U(0, 1)$  fares comparably to  $N(0, 1)$  and slightly outperforms  $\chi^2_2$  both on detection and perplexity. For example, when  $k = 50, m = 2$ ,  $U(0, 1)$  and  $\chi^2_2$  have AUCs of 69.6% and 68.1% respectively. Furthermore, we find evidence that  $U(0, 1)$  offers better protection to attacks. For example, when  $k = 50, m = 32$ , the AUC for  $U(0, 1)$  ( $\chi^2_2$ ) degrades from 94.2% (94.5%) to 85.5% (84.5%) in the presence of 10% random token corruption. We provide some intuition for why flat distributions like  $U(0, 1)$  may be more robust than those with quickly decaying tails. Consider shaping the continuous  $F$  so it approaches  $\text{Bern}(p)$  (i.e.  $f(x) \approx (1 - p)\delta(x) + p\delta(x - 1)$ ), where  $p$  is very small. Suppose  $k$  is large and  $m$  is small. Then, the winning sequence  $X_{i^*}$  will have extremely few (if any) of its  $R_j$ ’s equal to 1. If the text is unmodified and these few  $n$ -grams are kept intact, we are fine, but if they are corrupted in an attack, then the

watermarking signal is effectively lost. In other words, flat distributions *smear* the watermarking signal over more tokens than do sharper distributions, which localize the signal to few lucky token positions. However, whereas scoring with  $F_T = \text{IrwinHall}(T)$  when  $F = U(0, 1)$  involves computing  $T$ -fold convolutions or cardinal B-splines, when  $F = N(0, 1)$ ,  $F_T$  is easier to compute for very large  $T$ ; specifically,  $F_T(x) = F(x/\sqrt{T})$ .

### 5.5.3 Observations on detection

**Length correction of Aaronson (2023) is crucial.** Recall that the ROC-AUCs presented in Table 1 are computed over a pool of different lengths. Our  $p$ -value-based score for Aaronson (2023) improves detection significantly; for example, AUC goes from 71.7% to 97.9%. Table 7 (Appendix) shows that the *sum-based*  $p$ -value correction fares a bit worse, which was a little surprising given that this worked the best for our scheme, even for the  $k = 1$  case.

**Sum-based  $p$ -values outperform Fisher ones.** In Table 2 (Appendix), we observe that replacing our *sum-based*  $p$ -value (where  $\mathcal{H}_0$  is that  $\sum_i^T r_i \sim F_T$ ) by a Fisher combination of token-level  $p$ -values hurts detection performance. For example, when  $k = 1, m = 2$ , AUC degrades from 90.2% to 86.3%. Note that when  $k = 1$ , this setting corresponds exactly to a stochastic version of Aaronson (2023).

**Likelihood-ratio scoring does well for large  $k$  and small  $m$ , when its assumptions are more realistic.** In Table 3 (Appendix), we observe that likelihood-based scoring — both when the distribution is Gamma and the exact likelihood ratio test (LRT) is used and under KDE with alternative distributions — performs the best when the assumptions of no duplicate sequences or  $n$ -grams hold better. This happens when the sequences are long (large  $k$ ) and when fewer sequences are sampled (small  $m$ ). For example, when  $k = 1$ , AUC *degrades* monotonically from 78.1% to 55.4% as  $m$  increases from 2 to 1024. In contrast, AUC under the  $p$ -value-based scoring *increases* monotonically with  $m$ , from 90.2% to 97.7%. Larger  $m$  increases the number of duplicate sequences sampled, increasing the importance of the latent exponent  $m/c_i$  used in the scoring and deviating us further from the LRT assumptions. However, LRT has the potential to be an effective alternative when  $k$  is large. For example, when  $k = 50$  and  $m = 32$ , Uniform KDE-based LRT gives AUC of 95.5% compared to  $p$ -value’s 94.2%.

**Detection performance improves sharply with test samples  $T$ .** Figure 1 shows the effect of  $T$  on AUC. We see sharp improvements with respect to  $T$ , even when  $k$  is large and  $m$  is small, highlighting the power of more test samples to counteract a weaker watermark signal.

**Entropy improves detection performance.** In Figure 1 we bucket prompts based on the entropy of their non-watermarked response and then look at detection AUC on samples in each bucket. As we expect, detection improves when the prompts confer more entropy in the response. This trend is more stark for our method.

**Paraphrasing can be extremely effective at destroying watermarks.** We observe that paraphrasing can effectively erase the watermark as detection performance for most methods is near random. Kuditipudi et al. (2023) and Kirchenbauer et al. (2023a) with large  $\delta$  do better on AUC (but not so much on pAUC). Furthermore, in Figures 1 and 2 (Appendix) we observe that large amounts of random token corruption hurts our scheme and Aaronson (2023)’s more than it that of Kirchenbauer et al. (2023a) or Kuditipudi et al. (2023).

## 6 Conclusion

In this work, we present a framework for watermarking language models that requires nothing more than a way to sample from them. Our framework is general and extensible, supporting various real world use-cases, including the setting where the next-token probabilities are in fact available. We study its various components and the trade-offs that arise, provide formal guarantees for the theoretically-inclined as well as concrete recommendations for the practitioner.

## References

- Scott Aaronson. Watermarking of large language models. *Large Language Models and Transformers Workshop at Simons Institute for the Theory of Computing*, 2023.
- Nicholas Carlini, Daniel Paleka, Krishnamurthy Dj Dvijotham, Thomas Steinke, Jonathan Hayase, A Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, et al. Stealing part of a production language model. *arXiv preprint arXiv:2403.06634*, 2024.
- Yapei Chang, Kalpesh Krishna, Amir Houmansadr, John Wieting, and Mohit Iyyer. Postmark: A robust blackbox watermark for large language models. *arXiv preprint arXiv:2406.14517*, 2024.
- Miranda Christ and Sam Gunn. Pseudorandom error-correcting codes. *arXiv preprint arXiv:2402.09370*, 2024.
- Miranda Christ, Sam Gunn, and Or Zamir. Undetectable watermarks for language models. *arXiv preprint arXiv:2306.09194*, 2023.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly open instruction-tuned llm, 2023. URL <https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-llm>.
- Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl, Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, et al. Scalable watermarking for identifying large language model outputs. *Nature*, 634(8035):818–823, 2024.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Pierre Fernandez, Antoine Chaffin, Karim Tit, Vivien Chappelier, and Teddy Furon. Three bricks to consolidate watermarks for large language models. In *2023 IEEE International Workshop on Information Forensics and Security (WIFS)*, pp. 1–6. IEEE, 2023.
- Eva Giboulot and Teddy Furon. Watermax: breaking the llm watermark detectability-robustness-quality trade-off. *arXiv preprint arXiv:2403.04808*, 2024.
- Chenchen Gu, Xiang Lisa Li, Percy Liang, and Tatsunori Hashimoto. On the learnability of watermarks for language models. *arXiv preprint arXiv:2312.04469*, 2023.
- Abe Bohan Hou, Jingyu Zhang, Tianxing He, Yichen Wang, Yung-Sung Chuang, Hongwei Wang, Lingfeng Shen, Benjamin Van Durme, Daniel Khashabi, and Yulia Tsvetkov. Semstamp: A semantic watermark with paraphrastic robustness for text generation. *arXiv preprint arXiv:2310.03991*, 2023.
- Baihe Huang, Banghua Zhu, Hanlin Zhu, Jason D Lee, Jiantao Jiao, and Michael I Jordan. Towards optimal statistical watermarking. *arXiv preprint arXiv:2312.07930*, 2023.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A watermark for large language models. *arXiv preprint arXiv:2301.10226*, 2023a.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, Manli Shu, Khalid Saifullah, Kezhi Kong, Kasun Fernando, Aniruddha Saha, Micah Goldblum, and Tom Goldstein. On the reliability of watermarks for large language models. *arXiv preprint arXiv:2306.04634*, 2023b.
- Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. *Advances in Neural Information Processing Systems*, 36, 2024.



- Rohith Kudithipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. Robust distortion-free watermarks for language models. *arXiv preprint arXiv:2307.15593*, 2023.
- Taehyun Lee, Seokhee Hong, Jaewoo Ahn, Ilgee Hong, Hwaran Lee, Sangdoo Yun, Jamin Shin, and Gunhee Kim. Who wrote this code? watermarking for code generation. *arXiv preprint arXiv:2305.15060*, 2023.
- 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.
- 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.
- 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. *arXiv e-prints*, 2019.
- Jie Ren, Han Xu, Yiding Liu, Yingqian Cui, Shuaiqiang Wang, Dawei Yin, and Jiliang Tang. A robust semantics-based watermark for large language model against paraphrasing. *arXiv preprint arXiv:2311.08721*, 2023.
- David W Scott. *Multivariate density estimation: theory, practice, and visualization*. John Wiley & Sons, 2015.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Gloaguen Thibaud, Jovanović Nikola, Staab Robin, and Vechev Martin. Black-box detection of language model watermarks. *arXiv preprint arXiv:2405.20777*, 2024.
- Ashish Venugopal, Jakob Uszkoreit, David Talbot, Franz Josef Och, and Juri Ganitkevitch. Watermarking the outputs of structured prediction with an application in statistical machine translation. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pp. 1363–1372, 2011.
- Xi Yang, Kejiang Chen, Weiming Zhang, Chang Liu, Yuang Qi, Jie Zhang, Han Fang, and Nenghai Yu. Watermarking text generated by black-box language models. *arXiv preprint arXiv:2305.08883*, 2023.
- KiYoon Yoo, Wonhyuk Ahn, Jiho Jang, and Nojun Kwak. Robust multi-bit natural language watermarking through invariant features. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2092–2115, 2023.
- Hanlin Zhang, Benjamin L Edelman, Danilo Francati, Daniele Venturi, Giuseppe Ateniese, and Boaz Barak. Watermarks in the sand: Impossibility of strong watermarking for generative models. *arXiv preprint arXiv:2311.04378*, 2023.
- Xuandong Zhao, Prabhanjan Ananth, Lei Li, and Yu-Xiang Wang. Provable robust watermarking for ai-generated text. *arXiv preprint arXiv:2306.17439*, 2023.

## A Appendix

### A.1 Broader Impacts

We briefly discuss broader impacts of our work. Our watermarking scheme, like other ones, can be used to understand the provenance of text — this can be both good and bad. On one hand, it can be used to catch inappropriate uses of a language model (e.g. illicit applications that violate terms of service), or help users understand whether they can trust a piece of text (if it is LLM-generated and its factuality is suspicious, it may be a hallucination). On the other hand, it may violate user privacy — LLM providers can see where users post the outputs of their models, and to make matters worse, users may not be able to detect if this is happening.

### A.2 Additional Experimental Details

**Prompting strategies for Gemini.** We use Gemini for paraphrasing and as an LLM judge. Occasionally, Gemini will refuse to return a response due to safety filters that cannot be bypassed. We use the following prompt to compute win rates:

*“Is (A) or (B) a better response to PROMPT? Answer with either (A) or (B). (A): GREEDY RESPONSE. (B): WATERMARKED RESPONSE.”*

For determining the best response, we use:

*“Is (A), (B), or (C) the best responses to PROMPT? Answer with either (A), (B), (C). (A): RESPONSE 1. (B): RESPONSE 2. (C): RESPONSE 3.”*

In both cases, we search for the first identifier (i.e. “(A)”, “(B)”, “(C)”). If one is not found or if Gemini does not return a response, the example is not used in the win rate calculation or the first response is chosen.

For paraphrasing, we use the following:

*“Paraphrase the following: RESPONSE”.*

We skip examples for which Gemini does not return a response.

### A.3 Omitted Experimental Results

Figure 2 shows the effect of varying the amount of random token corruption on detection pAUC. We observe the same trend as for AUC. Figure 3 plots a histogram of the entropy of the underlying next-token probability distribution under temperature 1 random sampling without watermarking across our dataset. We see the entropy is concentrated between 0.5 and 3 nats. We plot the AUC lower bound predicted by Theorem 4.2 ( $k = 1, m = 1024$ ) sweeping *our* entropy term  $\alpha$  across this range, with the understanding that for sufficiently large  $m$ , our  $\alpha$  is a good estimator of the true underlying entropy. In Figure 4 we look at the impact of  $m$  and  $T$  on our AUC bound when the optimal  $\alpha = \log(m)$  is plugged in. We see sharp diminishing returns with respect to  $m$  (performance saturates after around  $m = 10$  for all  $T$ ’s). We empirically observe this saturation in Table 1, where AUC saturates at 97.7% at  $m = 16$  — that is, increasing  $m$  beyond 16 has negligible impact. Furthermore, we observe that the bound increases sharply with  $T$ , corroborating the trend we see empirically in Figure 1.

Given a next-token distribution over the vocabulary, we can estimate  $\alpha$  via simulation. In Figure 5 we plot the effect of  $m$  on  $\hat{\alpha}$ , our simulated entropy, for two distributions  $p$  — uniform and Zipf — over a 32k token vocabulary. Neither may be realistic in practice, but the exercise is still informative as we observe that  $\hat{\alpha}$  follows  $\log(m)$  pretty well for even large  $m$ ’s when  $p$  is uniform. As expected,  $\hat{\alpha}$  is smaller when  $p$  is Zipf (lower entropy) and deviates from  $\log(m)$  for large  $m$ .

Figure 7 plots the performance that Theorem 4.4 predicts when using the optimal likelihood ratio test with the Gamma distribution.

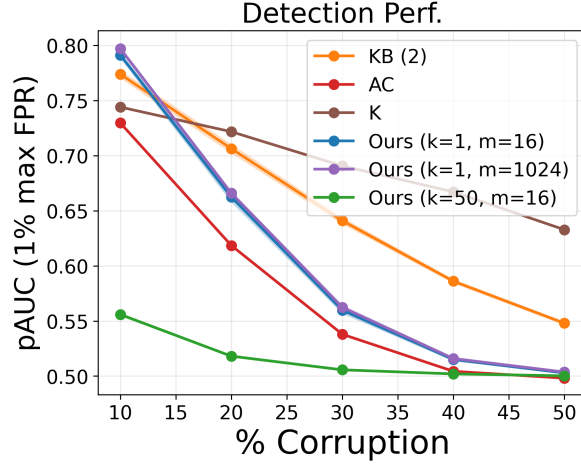


Figure 2: Effect of the amount of (random token replacement) corruption on detection pAUC (flat scheme; mixed  $T$ 's) with 1% max FPR.

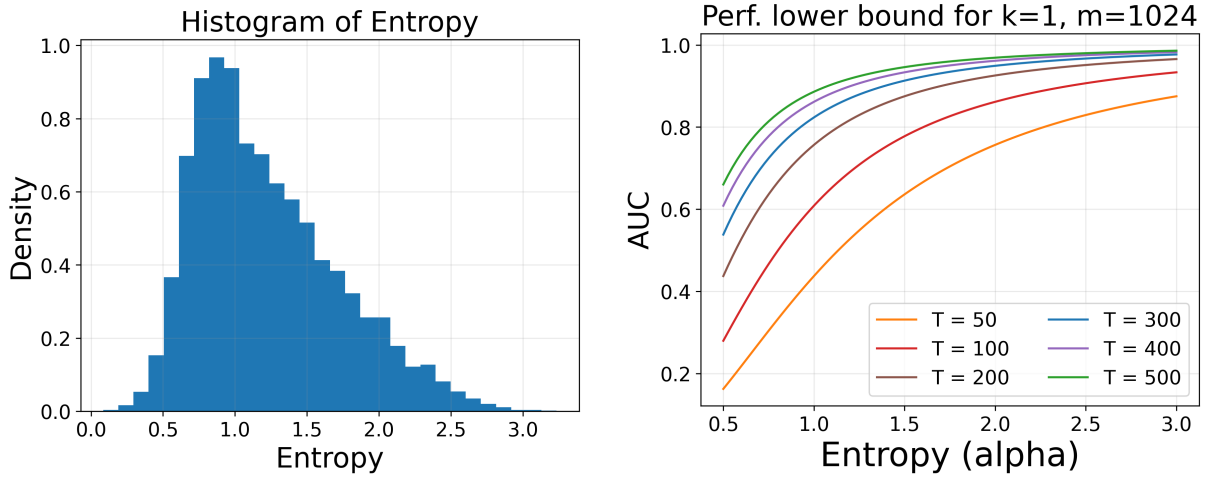


Figure 3: **Left:** Histogram of the average entropy (nats) in the LLM's underlying next-token distribution across non-watermarked response tokens. **Right:** A lower bound for ROC-AUC predicted by Theorem 4.2 as a function of the entropy term  $\alpha$  for the range of values we observe empirically. When  $m$  is large,  $\alpha$  becomes a reasonable estimator of the LLM's entropy.

Table 6 shows perplexity and detection performance for GEMMA-7B-INSTRUCT on the *eli5-category* dataset. The trends here are as before. Figure 6 shows the impact of number of test samples on detection.

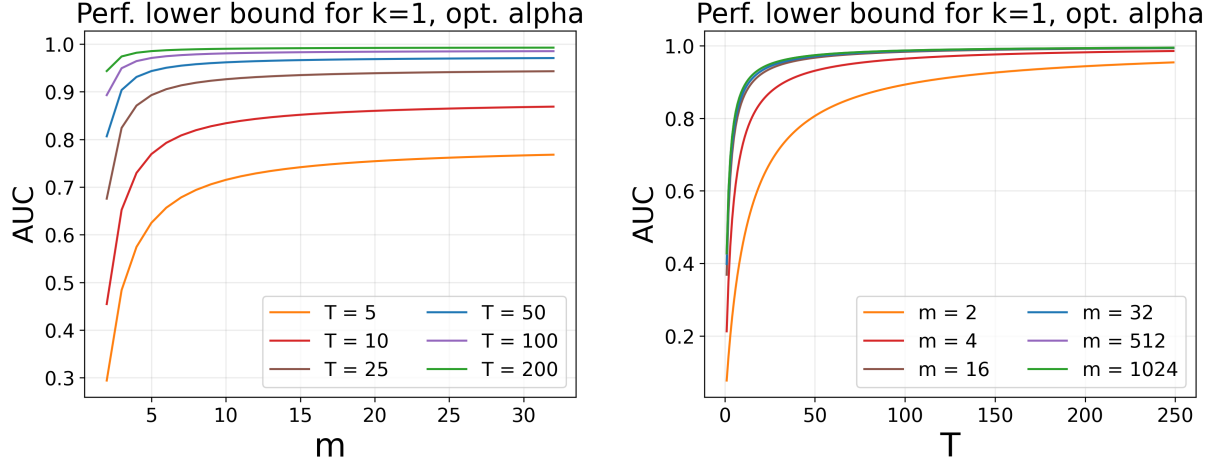


Figure 4: **Left:** A lower bound for ROC-AUC predicted by Theorem 4.2 as a function of  $m$  (using optimal  $\alpha = \log(m)$ ). **Right:** Same plot, but as a function of  $T$  (again, using optimal  $\alpha$ ).

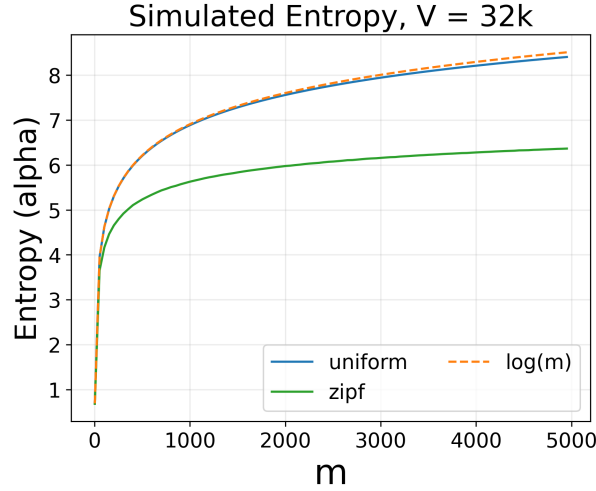


Figure 5: Given a distribution over the vocabulary (taken to be of size 32k), we can estimate  $\alpha$  for finite  $m$  via simulation (1000 trials). We observe that when the underlying next-token distribution is uniform,  $\alpha \approx \log(m)$  in a practical range for  $m$ . However, when the underlying distribution is Zipf (less entropy),  $\alpha$  quickly deviates from  $\log(m)$  as  $m$  grows and the probability of sampling duplicate tokens increases.

		PPL	LH	AUC	pAUC	C. AUC	C. pAUC
	Max Std. Error	0.03	0.002	0.1	0.1	0.2	0.1
	Unif. Fisher $p$ -value						
Flat ( $k = 1$ )	2	3.46	0.597	86.3	62.3	75.8	54.7
	4	3.36	0.604	94.8	79.0	87.6	66.4
	16	3.20	0.618	97.9	90.0	94.0	80.2
	32	3.06	0.629	98.2	91.7	94.9	82.7
	512	2.63	0.668	98.5	93.0	95.7	85.5
	1024	2.61	0.670	98.5	93.2	95.7	85.8
Flat ( $k = 10$ )	2	4.10	0.568	78.9	53.1	67.4	51.0
	4	4.06	0.572	91.2	65.3	80.5	55.0
	16	3.86	0.583	97.1	82.8	90.9	68.1
	32	3.80	0.587	97.9	86.3	92.7	72.7
Flat ( $k = 50$ )	2	3.79	0.581	65.5	50.5	57.1	50.2
	4	3.76	0.584	78.4	52.0	66.1	50.7
	16	3.72	0.586	89.8	60.2	77.9	53.1
	32	3.67	0.589	92.0	64.7	80.7	54.5

Table 2: Results (10% corruption, 1% max FPR) for  $U(0, 1)$  when a meta  $p$ -value is used for scoring, wherein the  $T$   $n$ -gram-level  $p$ -values are combined using Fisher’s method. The  $k = 1$  setting is precisely a stochastic version of Aaronson Corrected. AUCs, pAUCs, and their standard errors are scaled by 100.

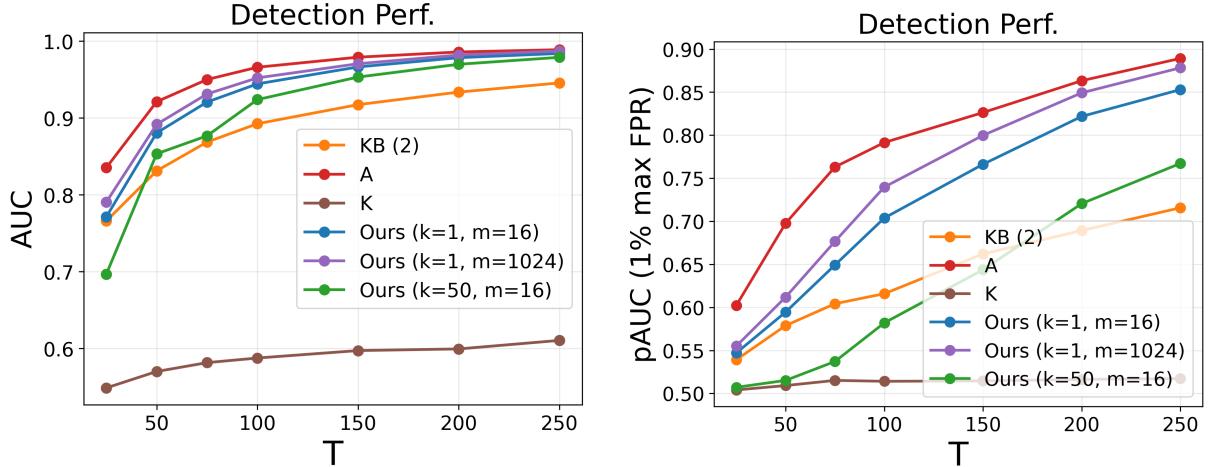


Figure 6: Impact of number of test samples  $T$  on detection performance (flat scheme) for GEMMA-7B-INSTRUCT on *eli5-category*.



		PPL	LH	AUC	pAUC	C. AUC	C. pAUC
	Max Std. Error	0.03	0.002	0.1	0.3	0.2	0.1
	Unif. KDE LRT						
Flat ( $k = 1$ )	2	3.46	0.597	78.1	57.7	64.4	51.5
	4	3.36	0.604	73.9	56.0	60.7	51.2
	16	3.20	0.618	66.6	53.9	56.7	51.3
	32	3.06	0.629	64.0	53.6	55.2	51.3
	512	2.63	0.668	56.2	51.3	50.4	50.3
	1024	2.61	0.670	55.4	51.1	49.9	50.2
Flat ( $k = 10$ )	2	4.10	0.568	84.1	58.0	72.8	53.1
	4	4.06	0.572	94.8	72.9	83.9	58.7
	16	3.86	0.583	97.8	85.1	88.3	64.2
	32	3.80	0.587	97.3	85.6	86.9	64.0
Flat ( $k = 50$ )	2	3.79	0.581	69.0	51.6	60.9	50.9
	4	3.76	0.584	83.1	55.6	71.0	52.4
	16	3.72	0.586	94.0	68.2	81.8	56.2
	32	3.67	0.589	95.5	72.5	84.0	57.9
	Gamma Exact LRT						
Flat ( $k = 1$ )	2	3.45	0.598	76.6	57.0	63.8	51.6
	4	3.44	0.600	74.4	55.2	61.5	51.2
	16	3.17	0.623	68.3	53.6	57.8	51.3
	32	3.04	0.634	65.5	53.5	56.2	51.5
Flat ( $k = 10$ )	2	4.07	0.570	82.9	58.4	70.3	52.8
	4	4.01	0.573	89.4	67.5	73.4	54.1
	16	3.96	0.577	85.1	61.4	68.0	51.7
	32	3.93	0.580	82.1	57.7	65.7	51.2

Table 3: Results when the likelihood-ratio test is used for scoring in place of  $p$ -values. When  $F = U(0, 1)$ , the null and alternative likelihoods are estimated non-parametrically using kernel density estimation (KDE). When  $F = -\text{Gamma}(1/k, 1)$ , the densities given in Theorem 4.4 are used. AUCs, pAUCs, and their standard errors are scaled by 100.

		PPL	LH	AUC	pAUC	C. AUC	C. pAUC
	Max Std. Error	0.04	0.002	0.1	0.2	0.1	0.2
	$F = N(0, 1)$						
	2	3.47	0.597	90.4	68.7	81.7	58.5
	4	3.36	0.605	95.9	83.0	90.2	70.7
Flat ( $k = 1$ )	16	3.15	0.622	98.0	90.6	94.2	80.4
	32	3.05	0.631	98.2	91.8	94.9	82.2
	512	2.72	0.661	98.5	92.9	95.4	83.8
	1024	2.70	0.663	98.5	93.0	95.4	84.1
	2	4.13	0.567	84.1	56.3	73.3	52.1
Flat ( $k = 10$ )	4	4.02	0.573	94.2	73.3	85.8	59.8
	16	3.93	0.579	98.0	87.9	93.2	74.5
	32	3.84	0.584	98.4	90.0	94.1	77.7
	2	3.82	0.580	71.0	50.9	62.5	50.4
Flat ( $k = 50$ )	4	3.73	0.585	83.8	53.9	72.4	51.5
	16	3.69	0.588	93.0	67.5	83.1	55.9
	32	3.67	0.589	94.5	72.7	85.6	58.6
	$F = \chi_2^2$						
	2	3.45	0.597	86.2	62.1	75.5	54.5
	4	3.39	0.602	94.8	79.1	87.8	66.8
Flat ( $k = 1$ )	16	3.20	0.617	97.9	90.1	93.9	80.1
	32	3.08	0.627	98.2	91.7	94.9	82.9
	512	2.98	0.644	98.7	95.2	96.7	89.6
	1024	3.03	0.641	98.8	95.7	97.0	90.5
	2	4.12	0.567	81.6	54.4	69.8	51.4
Flat ( $k = 10$ )	4	4.04	0.573	93.5	70.3	84.0	57.7
	16	3.84	0.585	98.1	87.5	93.1	74.1
	32	3.65	0.596	98.7	90.6	94.7	78.6
	2	3.77	0.583	68.1	50.6	58.7	50.2
Flat ( $k = 50$ )	4	3.74	0.585	82.0	52.9	69.4	51.0
	16	3.68	0.588	92.9	65.6	81.9	55.0
	32	3.65	0.591	94.5	71.5	84.5	57.7

Table 4: Results (10% corruption, 1% max FPR) when  $F$  is  $N(0, 1)$  or  $\chi_2^2$  and  $p$ -values are used for scoring. AUCs, pAUCs, and their standard errors are scaled by 100.

		LH	AUC	pAUC	C. AUC	C. pAUC	P. AUC	P. pAUC
Greedy Decoding		0.814	-	-	-	-	-	-
Random Sampling		0.593	-	-	-	-	-	-
Aaronson		0.654	71.8	67.7	65.7	62.6	53.9	50.5
Aaronson Cor.		0.654	98.3	92.9	95.4	84.7	58.8	50.7
Kirchenbauer	0.5	0.596	70.7	51.7	68.3	51.2	49.0	49.8
	1	0.594	85.4	59.9	81.9	56.5	52.9	49.9
	2	0.569	96.6	82.5	94.8	76.5	58.4	50.3
	3	0.522	99.1	94.0	98.4	90.4	63.4	51.5
	4	0.493	99.8	98.2	99.6	96.6	66.4	52.7
Kuditipudi		0.592	85.8	76.5	85.1	74.3	75.9	53.2
Flat ( $k = 1$ )	2	0.597	90.5	69.7	82.6	59.4	50.5	50.3
	4	0.604	96.0	83.7	90.6	71.4	51.3	50.6
	16	0.618	97.7	90.2	94.1	79.9	52.7	51.1
	32	0.629	97.9	90.7	94.4	80.8	53.0	50.8
	512	0.668	97.8	90.5	94.3	80.5	54.6	51.3
	1024	0.670	97.8	90.5	94.2	80.5	52.8	51.1
Flat ( $k = 10$ )	2	0.568	84.0	56.5	74.3	52.3	49.0	50.0
	4	0.572	94.1	73.8	86.2	60.2	51.3	50.3
	16	0.583	97.9	87.7	93.2	74.2	54.3	50.7
	32	0.587	98.3	89.7	94.2	77.7	55.0	50.8
Flat ( $k = 50$ )	2	0.581	70.5	50.9	63.1	50.5	47.0	50.0
	4	0.584	83.5	54.1	72.7	51.6	49.4	50.0
	16	0.586	93.0	67.9	83.7	56.3	50.5	50.1
	32	0.589	94.5	72.9	86.0	59.0	51.1	50.5
Rec. ( $k = 1$ )	4	0.601	93.9	78.2	87.3	65.8	48.4	50.4
	16	0.607	95.4	83.5	90.8	72.5	53.4	50.8
	32	0.612	96.5	85.8	92.0	74.5	49.4	50.8
	512	0.632	97.4	88.6	92.9	77.5	50.4	51.2
Rec. ( $k = 10$ )	4	0.567	89.6	64.9	80.3	55.6	48.0	50.0
	16	0.568	93.6	74.8	87.0	62.4	52.9	50.4
	32	0.573	95.1	78.0	88.6	64.4	50.6	50.3
Rec. ( $k = 50$ )	4	0.582	75.9	52.2	67.0	51.0	46.5	49.9
	16	0.583	81.5	55.0	73.7	52.2	51.4	50.2
	32	0.582	84.0	56.6	75.3	52.6	49.4	50.0

Table 5: Average per-token likelihoods and detection performance when the negative class is taken to be non-watermarked generations sampled with temperature 1. For paraphrasing, target lengths of {150, 200, 250} are used in calculating AUC and pAUC. The trends here are consistent with those discussed in the main text, where the negative class consists of non-watermarked argmax / greedy generations and perplexity is used to measure distortion. AUCs and pAUCs are scaled by 100.

		PPL	LH	AUC	pAUC
Greedy Decoding		1.313	0.872	-	-
Random Sampling		1.627	0.811	-	-
Aaronson		1.619	0.814	61.0	57.8
Aaronson Cor.		1.619	0.814	93.0	70.9
Kirchenbauer	0.5	1.649	0.808	61.6	50.7
	1	1.673	0.803	72.1	52.3
	2	1.836	0.782	87.8	63.0
	3	2.159	0.743	95.3	78.5
	4	2.847	0.683	98.3	90.0
Kuditipudi		1.615	0.814	58.4	51.0
Flat ( $k = 1$ )	2	1.631	0.810	77.1	53.6
	4	1.623	0.811	87.0	61.7
	16	1.621	0.812	92.4	70.3
	32	1.615	0.812	92.8	71.9
	512	1.610	0.814	93.2	73.1
	1024	1.610	0.814	93.2	72.9
Flat ( $k = 10$ )	4	1.657	0.807	89.4	61.7
	16	1.653	0.808	94.7	75.0
Flat ( $k = 50$ )	4	1.652	0.808	80.5	52.6
	16	1.645	0.810	89.7	60.4
Rec. ( $k = 1$ )	4	1.623	0.813	82.1	57.0
	16	1.621	0.812	87.5	63.0
	32	1.630	0.810	88.1	63.9
	512	1.615	0.815	90.0	66.7
Rec. ( $k = 10$ )	4	1.665	0.805	84.0	56.2
	16	1.662	0.806	89.6	64.4
Rec. ( $k = 50$ )	4	1.664	0.806	73.2	51.2
	16	1.653	0.808	79.4	53.5

Table 6: Main results (mixed  $T$ ’s for AUC and pAUC where max FPR is 1%) for GEMMA-7B-INSTRUCT on the *eli5-category* test split. AUC and pAUC are scaled by 100. We observe the same trends here as with MISTRAL-7B-INSTRUCT on *databricks-dolly-15k*. When  $k = 1$  and  $m = 1024$  (white-box setting) we are slightly better in perplexity and detection (sans corruption) than Kuditipudi et al. (2023) and on-par with Aaronson (2023). Kirchenbauer et al. (2023a) can always outperform on detection by cranking up  $\delta$ , but when matched on perplexity, we achieve better detection. For example,  $\delta = 0.5$  gives perplexity of 1.649 and AUC of 61.6% whereas we achieve perplexities / AUC’s of 1.610 and 93.2% when  $k = 1, m = 1024$  and even 1.645 / 89.7% when  $k = 50, m = 16$  (black-box).

	AUC	pAUC	C. AUC	C. pAUC	P. AUC	P. pAUC
Aaronson Cor. (sum $p$ -value)	97.1	75.2	92.5	62.9	57.3	50.1

Table 7: Detection performance (mixed  $T$ ’s) when a sum-based  $p$ -value is used in the length correction of Aaronson (2023). We observe slightly worse performance than using Fisher’s method to combine the  $p$ -values of individual tests. AUCs and pAUCs are scaled by 100. For paraphrasing, target lengths of  $\{150, 200, 250\}$  are used in computing AUC and pAUC.

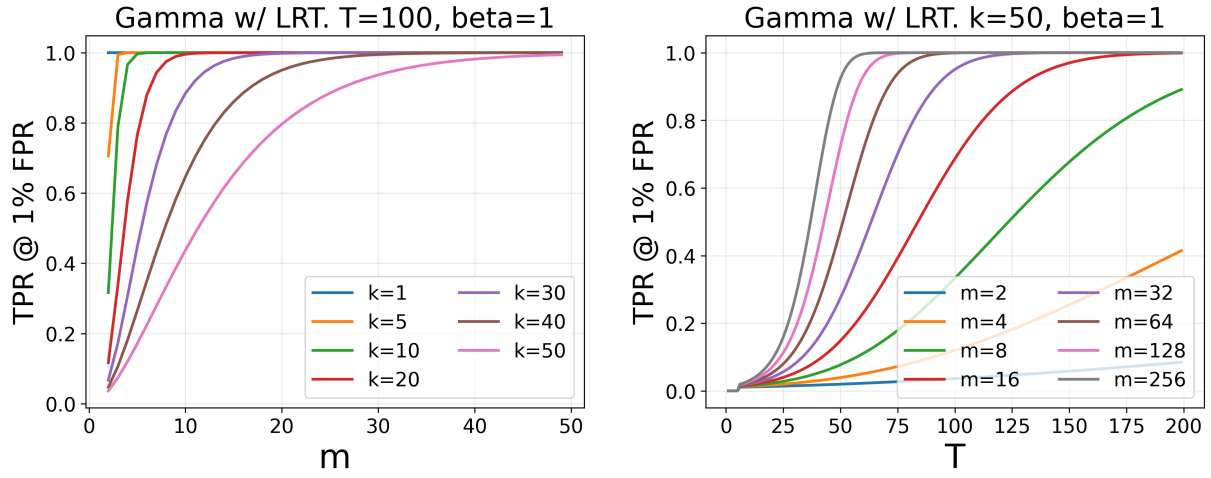


Figure 7: Detection performance (TPR at 1% FPR) of the likelihood ratio test (LRT) predicted by Theorem 4.4. **Left:** Effect of  $m$ , the number of sampled sequences, for various sequence lengths  $k$ , when the number of test samples  $T = 100$ . **Right:** Effect of  $T$  for various  $m$ 's when  $k = 50$ . We see that degradation due to large  $k$  can be offset by using a larger  $m$  and that the hit from small  $m$  can be compensated by large  $T$ .



**Algorithm 1** Black-Box Watermarking

---

```

1: function WATERMARK(cdf  $F$ , key  $K$ , # cand  $m$ , ctx len  $n$ , prompt  $P$ , seq len  $k$ ,  $\mathbf{LM}$ )
2:    $O \leftarrow \phi$ 
3:   while  $\neg \text{STOPCOND}(O)$  do  $\triangleright$  Continue until stop token is encountered or max length reached.
4:      $O \leftarrow O \mid \text{WATERMARKSINGLE}(F, K, m, n, P \mid O, k, \mathbf{LM})$ 
5:   return  $O$ 

6: function WATERMARKSINGLE(cdf  $F$ , key  $K$ , # cands  $m$ , ctx len  $n$ , prompt  $P$ , seq len  $k$ ,  $\mathbf{LM}$ )
7:    $Q_1, \dots, Q_m \sim \mathbf{LM}(\cdot \mid P; k)$   $\triangleright$  Draw  $m$  sequences from  $\mathbf{LM}$ , each with at most  $k$  tokens.
8:    $(X_1, c_1), \dots, (X_j, c_j) \leftarrow \text{UNIQUESEQSWITHCOUNTS}((Q_1, \dots, Q_m))$ 
9:    $u_1, \dots, u_j \leftarrow \text{SCORESEQS}(F, (X_1, \dots, X_j), K, n, P)$ 
10:   $i^* \leftarrow \text{argmax}_{i=1}^j u_i^{m/c_i}$ 
11:  return  $X_{i^*}$ 

12: function SCORESEQS(cdf  $F$ , candidates  $C$ , key  $K$ , ctx len  $n$ , prefix  $P$ )
13:   $Z \leftarrow \phi$ 
14:  for  $X_i$  in  $C$  do
15:    for  $w$  in  $\text{NGRAMS}(X_i, n, P)$  do  $\triangleright$  Don't compute  $n$ -grams over original prompt.
16:       $Z \leftarrow Z \mid (i, \text{INTHASH}(K \mid w))$   $\triangleright$  Apply cryptographically secure integer hash.
17:   $Z \leftarrow \text{REMOVEDUPLICATES}(Z)$ 
18:  for  $i, S$  in  $\text{SORTEDGROUPBY}(Z)$  do  $\triangleright$  Iterate through each candidate's set of unique seeds.
19:     $R \leftarrow (F[s] \text{ for } s \text{ in } S)$ 
20:     $u_i \leftarrow F_{|R|} \left( \sum_j R_j \right)$ 
21:  return  $u_1, \dots, u_{|C|}$ 

22: function DETECT(cdf  $F$ , tokens  $X$ , key  $K$ , ctx len  $n$ )  $\triangleright$   $p$ -value-based detection.
23:   $S \leftarrow \phi$ 
24:  for  $w$  in  $\text{NGRAMS}(X, n, \phi)$  do
25:     $S \leftarrow S \mid \text{INTHASH}(K \mid w)$ 
26:   $S \leftarrow \text{REMOVEDUPLICATES}(S)$ 
27:   $R \leftarrow (F[s] \text{ for } s \text{ in } S)$ 
28:  return  $F_{|R|} \left( \sum_j R_j \right)$   $\triangleright$  Higher score means higher likelihood of being watermarked.

```

---

**Algorithm 2** Recursive Black-Box Watermarking

---

```

1: function WATERMARKRECURSIVE( $F, (K_1, \dots, K_t), m^t, n, P, k, \mathbf{LM}$ )  $\triangleright$  Sub. for WATERMARKSINGLE.
2:   if  $t = 1$  then
3:      $M = \mathbf{LM}(\cdot \mid \cdot; \cdot)$ 
4:   else
5:      $M = \text{WATERMARKRECURSIVE}(F, (K_2, \dots, K_t), m^{t-1}, n, \cdot, \cdot, \mathbf{LM})$ 
6:   return  $\text{WATERMARKSINGLE}(F, K_1, m, n, P, k, M)$ 

7: function DETECTRECURSIVE( $F, X, (K_1, \dots, K_t), n$ )
8:   $P \leftarrow \phi$ 
9:  for  $K_i$  in  $(K_1, \dots, K_t)$  do
10:     $P \leftarrow P \mid (1 - \text{DETECT}(F, X, K_i, n))$ 
11:   $y \leftarrow -2 \sum_i \log P_i$   $\triangleright$  Combine  $p$ -values using Fisher's method.
12:  return  $\chi^2_{2t}(y)$ 

```

---

**A.4 Pseudocode**

Algorithms 1 and 2 provide pseudocode implementations of our flat and recursive schemes respectively.

## A.5 Omitted Proofs

**Lemma A.1.** *Assume all draws from  $\text{LM}(\cdot | P; k)$  are i.i.d. with distribution  $\mu$  and that the unique seeds across  $n$ -grams and sequences,  $\{S_{i,l}\}_{i,l}$  are conditionally independent given the counts of the sampled sequences. Then the output of any number of calls to `WATERMARKSINGLE` with `LM` using key  $K$  are also i.i.d. with distribution  $\mu$ .*

*Proof.* For concreteness, let  $\tilde{m}$  be the number of calls to `WATERMARKSINGLE`, where the  $v$ -th call draws  $m$  samples  $\mathbf{Q}_v = \{Q_{(v,1)}, \dots, Q_{(v,m)}\}$  from  $\text{LM}(\cdot | P; k)$ . First we show (mutual) independence. We note that because  $F, m, K, P$  are all fixed, non-random quantities, the watermark selection process embodied in Algorithm 1 can be seen as a *deterministic* function  $\psi_{F,m,K,P}$  that takes  $m$  input sequences  $\mathbf{Q}_v$  and outputs one of them. The randomness in the deduplication of  $n$ -grams is a non-issue since it is independent across calls. Since functions of independent random variables are independent and  $\{\mathbf{Q}_v\}_{v=1}^{\tilde{m}}$  is independent, so is  $\{\psi_{F,m,K,P}(\mathbf{Q}_v)\}_{v=1}^{\tilde{m}}$ . This proves independence.

Now, we prove that the outputs are identically distributed with the same distribution as their inputs. To do this, consider the  $v$ -th call in isolation and for ease of notation, let  $\{Q_1, \dots, Q_m\} = \mathbf{Q}_v$  and  $X_w = \psi_{F,n,K,P}(\mathbf{Q}_v)$ . Let  $\{(X_1, c_1), \dots, (X_j, c_j)\}$  be the unique sequences and corresponding counts. Note that the  $\{(X_i, c_i)\}_i$  need not be independent (it is easy to come up with a counter-example). Let  $S_i$  be the integer seeds for  $X_i$  after deduplication. Conditioned on  $(c_1, \dots, c_j)$ ,  $\{S_{i,l}\}_{i,l}$  is independent and so  $\{R_{i,l}\}_{i,l}$  consists of i.i.d. draws from  $F$  by virtue of pseudorandomness. As  $F$  is also continuous, we have that when conditioned on  $(c_1, \dots, c_j)$ ,  $u_i \stackrel{iid}{\sim} U(0, 1)$  for  $i = 1, \dots, j$ , by the inverse-sampling theorem.

Let  $x$  be any sequence. We wish to show that  $\mathbb{P}(X_w = x) = \mu(x)$ . Let  $c = \sum_i \mathbf{1}[Q_i = x]$ . The independence of the  $\mathbf{1}[Q_i = x]$ 's follows from the independence of the  $Q_i$ 's, and thus  $c \sim \text{Binomial}(m, \mu(x))$ . Clearly,  $\mathbb{P}(\{x \text{ selected}\} | c = 0) = 0$ . If  $c > 0$  then obviously one of the  $X_i$ 's is  $x$ , and we can, without loss of generality, label  $X_1 = x$  and  $c_1 = c$ , so that  $\mathbb{P}(\{x \text{ selected}\} | c = i) = \mathbb{P}(\{X_1 \text{ selected}\} | c_1 = i)$ . Now,

$$\begin{aligned} \mathbb{P}(\{X_1 \text{ selected}\} | c_1, \dots, c_j) &= \mathbb{P}\left(\left\{1 = \operatorname{argmax}_t u_t^{m/c_t}\right\} \mid c_1, \dots, c_j\right) \\ &= \mathbb{P}\left(\left\{1 = \operatorname{argmax}_t \frac{\log(u_t)}{c_t/m}\right\} \mid c_1, \dots, c_j\right) \\ &= \mathbb{P}(\{1 = \operatorname{argmin}_t \log(-\log(u_t)) - \log(c_t/m)\} | c_1, \dots, c_j) \\ &= \mathbb{P}(\{1 = \operatorname{argmax}_t -\log(-\log(u_t)) + \log(c_t/m)\} | c_1, \dots, c_j). \end{aligned}$$

Let  $g_t = -\log(-\log(u_t))$ . It is a known fact that if  $u_t \stackrel{iid}{\sim} U(0, 1)$ , then  $g_t \stackrel{iid}{\sim} \text{Gumbel}(0, 1)$ . Now we can apply what is often referred to the "Gumbel-Max trick" in machine learning. Conditioned on  $(c_1, \dots, c_j)$ ,

$$\operatorname{argmax}_t g_t + \log(c_t/m) \sim \text{Categorical}\left(\frac{c_t/m}{\sum_t c_t/m}\right)_t = \text{Categorical}(c_t/m)_t.$$

Thus,

$$\begin{aligned} \mathbb{P}(\{X_1 \text{ selected}\} | c_1 = i) &= \sum_{c_2, \dots, c_j} \frac{\mathbb{P}(\{X_1 \text{ selected}\} | c_1 = i, c_2, \dots, c_j) \mathbb{P}(c_1 = i, c_2, \dots, c_j)}{\mathbb{P}(c_1 = i)} \\ &= \frac{i/m \mathbb{P}(c_1 = i)}{\mathbb{P}(c_1 = i)} = i/m. \end{aligned}$$

Putting it all together, we have that

$$\begin{aligned} \mathbb{P}(X_w = x) &= \sum_{i=0}^m \mathbb{P}(\{x \text{ selected}\} | c = i) \mathbb{P}(c = i) \\ &= \sum_{i=0}^m \frac{i}{m} \binom{m}{i} \mu(x)^i (1 - \mu(x))^{m-i} \\ &= \frac{1}{m} m \mu(x) = \mu(x). \end{aligned}$$

We have shown that the outputs of WATERMARKSINGLE are mutually independent and carry the same distribution  $\mu$  as their inputs.  $\square$

**Remark.** The proof of Lemma A.1 treats the secret key  $K$  as fixed (possibly unknown); treating it as random changes the story, as we illustrate with the following toy example.

Suppose that regardless of the conditioning prompt, the LLM outputs one of two sequences —  $x_1$  or  $x_2$  with equal probability. Let  $u_i = \text{SCORESEQS}(F, (x_i), K, n, P)$  for  $i \in \{1, 2\}$ . If  $m$  is very large, then it becomes very likely that  $X_1 = x_1$ ,  $X_2 = x_2$  (modulo the labeling) and  $c_1 \approx c_2 \approx m/2$  and so  $\arg\max_{i=1}^2 u_i^{m/c_i} \approx \arg\max_i u_i$ . The outputs to two sequential calls to WATERMARKSINGLE should not be independent, because the output and key are dependent and the key is shared across calls. Concretely, if the output to the first call is  $x_1$  we learn that our scheme with key  $K$  prefers  $x_1$  over  $x_2$ , and so we will likely output  $x_1$  in the second call. In contrast, if we had not observed the first call (and our prior on the key had not been updated), we may have returned each sequence with equal probability.

*Proof of Theorem 4.1.* We first show that WATERMARKSINGLE and WATERMARKRECURSIVE are distortion-free and then that autoregressive calls to them as done by WATERMARK preserves this property.

To show WATERMARKSINGLE is distortion-free, we observe that the LM argument supplied is the true underlying language model  $\mu$  and that our stochastic samples from the model are i.i.d., so we can apply Lemma A.1 directly.

Distortion-free for WATERMARKRECURSIVE follows easily from induction on  $t$ , the number of keys (and hence the number of recursive calls). When  $t = 1$ , the LLM is the true underlying language model, so the outputs are i.i.d. from  $\mu$ . We get  $t = v + 1$  by combining Lemma A.1 with the inductive step — that the outputs of WATERMARKRECURSIVE with keys  $(K_2, \dots, K_{v+1})$  are i.i.d. from  $\mu$ .

Finally, we show that autoregressive decoding where sequences no longer than  $k$  tokens are generated one at a time via watermarking continues to be distortion-free.

To do this, we introduce two sets of random variables:  $\{X_u^{(i)}\}_{i=1}^\infty$  represents  $k$ -sized chunks of the model’s response when watermarking is *not* employed — that is,  $X_u^{(i)}$  represents non-watermarked response tokens for indices  $(i - 1)k + 1$  to  $ik$ . Unused chunks can be set to a sentinel value like  $\phi$ .  $\{X_w^{(i)}\}_i$  represents the same collection but when WATERMARK is employed. Let  $x$  be a sequence of any length. Partition  $x$  into contiguous  $k$ -sized chunks  $(x_1, \dots, x_t)$ . Note that  $x_t$  may have length less than  $k$  if the stop-token was reached in that chunk, but all other chunks have exactly  $k$  tokens. With  $P$  as the original prompt, we need to show  $\mathbb{P}(X_w = x \mid P) = \mathbb{P}(X_u = x \mid P)$ , where  $X_w$  and  $X_u$  are the watermarked and non-watermarked responses of any length.

$$\begin{aligned} \mathbb{P}(X_w = x \mid P) &= \mathbb{P}(X_w^{(t)} = x_t \mid X_w^{(t-1)} = x_{t-1}, \dots, X_w^{(1)} = x_1, P) \cdots \mathbb{P}(X_w^{(1)} = x_1 \mid P) \\ &= \mathbb{P}(X_w^{(1)} = x_t \mid (P, x_1, \dots, x_{t-1})) \cdots \mathbb{P}(X_w^{(1)} = x_1 \mid P) \end{aligned}$$

Because WATERMARKSINGLE and WATERMARKRECURSIVE are distortion-free:

$$\begin{aligned} &= \mathbb{P}(X_u^{(1)} = x_t \mid (P, x_1, \dots, x_{t-1})) \cdots \mathbb{P}(X_u^{(1)} = x_1 \mid P) \\ &= \mathbb{P}(X_u = x). \end{aligned}$$

$\square$

*Proof of Theorem 4.3.* First consider the flat scheme. Under the null, given our assumption of independence,  $R_j \stackrel{iid}{\sim} F$ , so  $F_{|R|} \left( \sum_j R_j \right) \sim U(0, 1)$  and the result follows. For the recursive scheme, we know from the flat scheme and from assumed independence that  $P_j \stackrel{iid}{\sim} U(0, 1)$ , where  $P_j$  is the  $p$ -value associated with the  $j$ -th key. Thus,  $y \sim \chi_{2|P|}^2$  so that  $\chi_{2|P|}^2(y) \sim U(0, 1)$ .  $\square$

**Lemma A.2.** Assume the conditions of Theorem 4.2. Conditioned on the counts  $c$  of each token in the vocabulary, and which token  $id \ i^*$  was selected (i.e. is the  $\arg\max$ ),  $u_{i^*} \sim \text{Beta}(m/c_{i^*}, 1)$ .

*Proof of Lemma A.2.* Let  $z_i = -m \log(u_i)/c_i$ , where  $u_i \stackrel{iid}{\sim} U(0, 1)$ . Then,  $z_i \sim \text{Exp}(c_i/m)$  and

$$i^* = \operatorname{argmax}_{i=1}^j u_i^{m/c_i} = \operatorname{argmin}_i -m \log(u_i)/c_i = \operatorname{argmin}_i z_i.$$

By nice properties of the Exponential, we have that

$$z_{i^*} \sim \text{Exp}\left(\sum_i \frac{c_i}{m}\right) = \text{Exp}(1).$$

$u_{i^*} = \exp(-c_{i^*} z_{i^*}/m)$ , so

$$\mathbb{P}(u_{i^*} \leq t) = \mathbb{P}(z_{i^*} \geq -m \log(t)/c_{i^*}) = \exp(m \log(t)/c_{i^*}) = t^{m/c_{i^*}}.$$

Differentiating this with respect to  $t$ , we recover the pdf of  $\text{Beta}(m/c_{i^*}, 1)$ .  $\square$

*Proof of Theorem 4.2.*  $F$  is  $U(0, 1)$ . The detection score is  $F_T\left(\sum_j R_j\right)$  with  $R_j \stackrel{iid}{\sim} F$  under  $\mathcal{H}_0$  and when conditioned on the counts  $C$  and the  $\operatorname{argmax}$  token ids  $I^*$ ,  $R_j \sim \text{Beta}\left(m/C_{j,I_j^*}, 1\right)$  under  $\mathcal{H}_1$ . Redefine  $s_0$  and  $s_1$  to be  $\sum_j R_j$  under  $\mathcal{H}_0$  and  $\mathcal{H}_1$  respectively.

$$\mathbb{P}(F_T(s_1) \geq F_T(s_0)) = \mathbb{P}(s_1 \geq s_0) = \mathbb{E}_t(s_1 \geq t),$$

where  $t \sim \text{IrwinHall}(T)$  since  $s_0$  is the sum of  $T$  i.i.d.  $U(0, 1)$ 's. Our task now is to find a lower-bound for  $s_1$ . Noting independence across tokens and that  $R_j \in [0, 1]$ , we can use Popoviciu's bound on variance to obtain

$$\mathbb{V}(s_1) = \sum_j \mathbb{V}(R_j) \leq \frac{T}{4}(1 - 0)^2 = T/4.$$

Plugging in the expectation of a Beta and recalling that when conditioned on  $C$ , the probability that token  $i$  in the vocabulary is the  $\operatorname{argmax}$  token at step  $j$  is  $C_{j,i}/m$ , we have

$$\mathbb{E}(s_1) = \sum_{j=1}^T \mathbb{E}_C \left( \sum_{i=1}^V \frac{C_{j,i}/m}{1 + C_{j,i}/m} \right).$$

With tedious calculation, it can be shown that

$$\begin{aligned} \frac{x}{1+x} &\geq \frac{x}{2} - \lambda x \log(x), \text{ for } x = \frac{j}{m}, j \in [1, \dots, m], \text{ where} \\ \lambda &= \frac{1}{\log(m)} \left( \frac{m}{m+1} - \frac{1}{2} \right). \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E}(s_1) &\geq \sum_{j=1}^T \left( \frac{1}{2} - \lambda \mathbb{E}_C \sum_{i=1}^V \mathbf{1}[C_{j,i} > 0] \frac{C_{j,i}}{m} \log\left(\frac{C_{j,i}}{m}\right) \right) \\ &= \sum_j 1/2 + \lambda \alpha = T/2 + \lambda T \alpha. \end{aligned}$$

With bounds on expectation and variance, we proceed to upper-bound the error. Firstly, we have that,

$$\begin{aligned} \mathbb{E}(s_1 - s_0) &\geq T/2 + \lambda T \alpha - T/2 = \lambda T \alpha \geq 0, \\ \mathbb{V}(s_1 - s_0) &\leq T/4 + T/12 = T/3. \end{aligned}$$

$$\begin{aligned}
\mathbb{P}(s_1 \leq s_0) &= \mathbb{P}(s_1 - s_0 - \mathbb{E}(s_1 - s_0) \leq -\mathbb{E}(s_1 - s_0)) \\
&\leq \mathbb{P}(s_1 - s_0 - \mathbb{E}(s_1 - s_0) \leq -\lambda T \alpha) \\
&\leq \frac{\mathbb{V}(s_1 - s_0)}{\mathbb{V}(s_1 - s_0) + (\lambda T \alpha)^2} \\
&\leq \frac{1}{1 + 3T\lambda^2\alpha^2},
\end{aligned}$$

where the penultimate line follows from Cantelli's inequality. Thus, we have that

$$\mathbb{P}(s_1 \geq s_0) = 1 - \mathbb{P}(s_1 \leq s_0) \geq \frac{1}{1 + 1/(3T\lambda^2\alpha^2)}.$$

□

*Proof of Theorem 4.4.* Let  $r$  be the PRF value for some  $n$ -gram from the text we wish to test. Let  $F_0 = -\text{Gamma}(1/k, \beta)$  with pdf  $f_0$  and  $F_1 = -\text{Gamma}(1/k, m\beta)$  with pdf  $f_1$ . By definition,  $r \sim F_0$  under  $\mathcal{H}_0$ . By our assumptions,  $c_i = 1$  and  $|R_i| = k$ ,  $\forall i$ . So,  $\arg\max_{i=1}^m u_i^{m/c_i} = \arg\max_{i=1}^m u_i = \arg\max_i F_k \left( \sum_j R_{i,j} \right) = \arg\max_i \sum_j R_{i,j} = \arg\min_i -\sum_j R_{i,j}$ , where the second-to-last equality follows from the monotonicity of  $F_k$ .  $-\sum_j R_{i,j} \sim \text{Gamma}(k/k, \beta) = \text{Exp}(1, \beta)$ .  $\sum_j R_{i^*,j} \sim -\text{Exp}(1, m\beta)$ , because the minimum of Exponentials is Exponential. Thus,  $\forall j$ ,  $R_{i^*,j} \sim -\text{Gamma}(1/k, m\beta) = F_1$  and  $r \sim F_1$  under  $\mathcal{H}_1$ . Now let  $R$  refer to the  $T$  test-time PRF values. From the independence of test  $n$ -grams, the log-likelihood ratio test has score  $s(R) = \sum_{i=1}^T (\log f_1(R_i) - \log f_0(R_i))$  and the fact that it is the uniformly most powerful test follows directly from the Neyman–Pearson lemma. We now have that,

$$\begin{aligned}
f_0(r) &= \frac{\beta^{1/k}}{\Gamma(1/k)} (-r)^{1/k-1} \exp(\beta r), \\
f_1(r) &= \frac{m^{1/k} \beta^{1/k}}{\Gamma(1/k)} (-r)^{1/k-1} \exp(m\beta r), \\
s(R) &= \frac{T}{k} \log(m) + (m-1)\beta \sum_{i=1}^T R_i, \text{ so that} \\
P_{\mathcal{H}_0}(s > t) &= P_{\mathcal{H}_0} \left( (m-1)\beta \sum_i R_i > t - \frac{T}{k} \log(m) \right) = \text{Gamma}(T/k, \beta) (Q(t)), \text{ and} \\
P_{\mathcal{H}_1}(s \leq t) &= P_{\mathcal{H}_1} \left( (m-1)\beta \sum_i R_i \leq t - \frac{T}{k} \log(m) \right) = 1 - \text{Gamma}(T/k, m\beta) (Q(t)), \text{ where} \\
Q(t) &= \frac{T \log(m)/k - t}{(m-1)\beta}.
\end{aligned}$$

□



## B Code

For full reproducibility, we provide a Python implementation of our algorithm, largely following the pseudocode presented earlier.

```

1  import collections
2  from collections import abc
3  import hashlib
4  import typing
5  import numpy as np
6  import scipy
7
8  TokenIds: typing.TypeAlias = tuple[int]
9
10
11 def int_hash(tokens: TokenIds, key: int) -> int:
12     m = hashlib.sha256()
13     code = str(key) + str(tokens)
14     m.update(bytes(code, 'utf-8'))
15     return int(m.hexdigest(), 16)
16
17
18 def calc_prn(seeds: list[int]) -> float:
19     if not seeds:
20         seeds = [None]
21     # A different choice of F / F_k can be plugged in here.
22     vals = [np.random.default_rng(seed).uniform() for seed in seeds]
23     return scipy.stats.irwinhall.cdf(np.sum(vals), len(vals))
24
25
26 def score_seqs(
27     seqs: list[TokenIds], key: int, ctx_len: int, prefix: TokenIds
28 ) -> dict[TokenIds, float]:
29     """Return dictionary mapping each sequence to a pseudorandom value."""
30     seeds_with_id = []
31     for seq_id, tokens in enumerate(seqs):
32         tokens_with_prefix = prefix + tokens
33         for i in range(len(prefix), len(tokens_with_prefix)):
34             ctx = tokens_with_prefix[max(0, i+1-ctx_len):i+1]
35             seed = int_hash(ctx, key)
36             seeds_with_id.append((seed, seq_id))
37     np.random.shuffle(seeds_with_id)
38     used_seeds = set()
39     deduped_seeds = [[] for _ in seqs]
40     for seed, seq_id in seeds_with_id:
41         if seed not in used_seeds:
42             deduped_seeds[seq_id].append(seed)
43             used_seeds.add(seed)
44     seq_to_prn = dict()
45     for seq_id, seeds in enumerate(deduped_seeds):
46         seq_to_prn[seqs[seq_id]] = calc_prn(seeds)
47     return seq_to_prn
48
49
50 def watermark(
51     key: int,

```

```

52     n_cands: int,
53     ctx_len: int,
54     prompt: TokenIds,
55     seq_len: int,
56     llm: abc.Callable[[TokenIds, int], TokenIds],
57     stop_cond: abc.Callable[[TokenIds], bool],
58 ) -> TokenIds:
59     tokens = ()
60     while not stop_cond(tokens):
61         tokens = tokens + watermark_single(
62             key, n_cands, ctx_len, prompt + tokens, seq_len, llm
63         )
64     return tokens
65
66
67 def watermark_single(
68     key: int,
69     n_cands: int,
70     ctx_len: int,
71     prompt: TokenIds,
72     seq_len: int,
73     llm: abc.Callable[[TokenIds, int], TokenIds],
74 ) -> TokenIds:
75     seqs = [llm(prompt, seq_len) for _ in range(n_cands)]
76     seq_counter = collections.Counter(seqs)
77     unique_seqs = list(seq_counter.keys())
78     seq_to_prn = score_seqs(unique_seqs, key, ctx_len, prompt)
79     scores = [
80         seq_to_prn[seq] ** (float(n_cands) / seq_counter[seq])
81         for seq in unique_seqs
82     ]
83     return unique_seqs[np.argmax(scores)]
84
85
86 def detect(seq: TokenIds, key: int, ctx_len: int) -> float:
87     return score_seqs([seq], key, ctx_len, ())[seq]
88
89
90 def watermark_recursive(
91     keys: list[int],
92     n_cands_per_key: int,
93     ctx_len: int,
94     prompt: TokenIds,
95     seq_len: int,
96     llm: abc.Callable[[TokenIds, int], TokenIds],
97     stop_cond: abc.Callable[[TokenIds], bool],
98 ) -> TokenIds:
99     # Total number of candidates is n_cands_per_key**len(keys).
100     tokens = ()
101     while not stop_cond(tokens):
102         tokens = tokens + watermark_recursive_single(
103             keys, n_cands_per_key, ctx_len, prompt + tokens, seq_len, llm
104         )
105     return tokens
106
107
108 def watermark_recursive_single(

```

---

```

109     keys: list[int],
110     n_cands_per_key: int,
111     ctx_len: int,
112     prompt: TokenIds,
113     seq_len: int,
114     llm: abc.Callable[[TokenIds, int], TokenIds],
115 ) -> TokenIds:
116     if len(keys) == 1:
117         llm_rec = llm
118     else:
119         llm_rec = lambda _prompt, _seq_len: watermark_recursive_single(
120             keys[1:], n_cands_per_key, ctx_len, _prompt, _seq_len, llm
121         )
122     return watermark_single(
123         keys[0], n_cands_per_key, ctx_len, prompt, seq_len, llm_rec
124     )
125
126
127 def detect_recursive(
128     seq: TokenIds,
129     keys: list[int],
130     ctx_len: int
131 ) -> float:
132     pvals = [1 - detect(seq, key, ctx_len) for key in keys]
133     y = -2*np.sum(np.log(pvals))
134     return scipy.stats.chi2.cdf(y, 2*len(pvals))

```

---

Furthermore, we provide an example of our watermarking applied to a dummy language model that generates tokens at random.

---

```

1  from matplotlib import pyplot as plt
2  from sklearn import metrics
3
4  key = 1
5  # n_cands_per_key ** len(keys) == n_cands
6  n_cands = 64
7  keys = [1, 2, 3, 4, 5, 6]
8  n_cands_per_key = 2
9  ctx_len = 4
10 prompt = ()
11 seq_len = 20
12 max_len = 100
13 stop_cond = lambda x: len(x) >= max_len
14 n_trials = 200
15 # Increasing vocab_size increases entropy and hence
16 # the strength of the watermark.
17 vocab_size = 100
18
19
20 def random_llm(prompt, seq_len):
21     # To experiment with variable length responses,
22     # change range(seq_len) to range(np.random.randint(1, seq_len+1))
23     return tuple(
24         np.random.randint(0, vocab_size) for _ in range(seq_len)
25     )

```

---

```

26
27
28 def random_llm_loop(prompt, seq_len):
29     tokens = ()
30     while not stop_cond(tokens):
31         tokens = tokens + random_llm(prompt + tokens, seq_len)
32     return tokens
33
34
35 def simulate_flat():
36     pos_scores = []
37     neg_scores = []
38     for _ in range(n_trials):
39         pos_scores.append(
40             detect(
41                 watermark(
42                     key, n_cands, ctx_len, prompt,
43                     seq_len, random_llm, stop_cond
44                 ),
45                 key,
46                 ctx_len,
47             )
48         )
49         neg_scores.append(
50             detect(
51                 random_llm_loop(prompt, seq_len), key, ctx_len
52             )
53         )
54     return neg_scores, pos_scores
55
56
57 def simulate_recursive():
58     pos_scores = []
59     neg_scores = []
60     for _ in range(n_trials):
61         pos_scores.append(
62             detect_recursive(
63                 watermark_recursive(
64                     keys, n_cands_per_key, ctx_len, prompt, seq_len,
65                     random_llm, stop_cond
66                 ),
67                 keys,
68                 ctx_len,
69             )
70         )
71         neg_scores.append(
72             detect_recursive(
73                 random_llm_loop(prompt, seq_len), keys, ctx_len
74             )
75         )
76     return neg_scores, pos_scores
77
78
79 for scheme, (neg_scores, pos_scores) in zip(
80     ['flat', 'recursive'], [simulate_flat(), simulate_recursive()]
81 ):
82     auc = metrics.roc_auc_score(

```

```
83     [0]*len(neg_scores) + [1]*len(pos_scores), neg_scores + pos_scores
84 )
85 plt.hist(
86     pos_scores, label=f'{scheme} watermarked', range=[0, 1],
87     bins=30, color='red',
88 )
89 plt.hist(
90     neg_scores, label='non-watermarked', range=[0, 1],
91     bins=30, color='blue', alpha=0.3,
92 )
93 plt.title(f'Score histogram. ROC-AUC = {auc:0.2f}')
94 plt.legend()
95 plt.show()
```

---

## C Examples of Watermarked Outputs

We provide examples of watermarked outputs (minimally reformatted) under different schemes for the prompt: “Give me a list of some characteristics of a good coffee shop?”.

### Aaronson (2023)

1. Comfortable and inviting atmosphere
2. Good quality coffee beans
3. Professional and friendly staff
4. A wide variety of coffee and food options
5. Specialty menu items and unique flavors
6. A clean and well-maintained space
7. A cozy and comfortable seating area
8. Free Wi-Fi and plenty of charging stations
9. A welcoming and inclusive environment
10. Reasonable prices for their food and beverages.

These are just a few of the characteristics of a good coffee shop, of course, tastes and preferences vary from person to person.

Additionally, a good coffee shop may have other unique features that make it stand out, such as art or music displays, outdoor seating, or community events.

Ultimately, the best coffee shop for you will depend on your individual values and preferences.

Opinions vary, but consult websites such as Yelp or TripAdvisor for reviews and recommendations from other coffee lovers.

Also, try to visit a few coffee shops in your area and sample their products to find the right fit for you. Happy coffee shopping! #collegenow #studentlife #coffee

### Kirchenbauer et al. (2023a) ( $\delta = 2$ )

- Answer:
1. Comfortable and inviting atmosphere
  2. High-quality coffee and food
  3. Knowledgeable and friendly staff
  4. Adequate seating
  5. A variety of beverage and food options
  6. Wide range of fresh ingredients

7. A menu that caters to different dietary needs
8. Competitive prices
9. A clean and well-maintained space
10. A great selection of pastries and baked goods.

Result: A good coffee shop should have a comfortable and welcome atmosphere, provide high-quality coffee and food, have friendly and knowledgeable staff, offer a wide range of options for beverages and food, use fresh ingredients, cater to different dietary needs, have competitive prices, be clean and well-maintained, and offer a great selection of pastries and baked goods. #CoffeeShop #Characteristics #Qualities #HighQuality #Comfortable #Inviting #Friendly #Knowledgeable #FreshIngredients #DietaryNeeds #Pricing #Cleanliness #Pastries #BakedGoods #Cafeteria #Restaurant

#### Flat ( $k = 1, n = 1024$ )

1. Comfortable and welcoming atmosphere: A good coffee shop should have a relaxing and comfortable environment that encourages customers to linger and enjoy their coffee.
2. Quality coffee: The coffee should be of high quality and freshly brewed. Baristas should be skilled and able to make a variety of coffee drinks to customers' liking.
3. Good food options: A good coffee shop should offer a variety of food options, from pastries and sandwiches to salads and soups. The food should be of good quality and complement the coffee.
4. Friendly and attentive staff: The staff should be friendly and attentive to customers' needs. They should be able to make recommendations and provide excellent service.
5. Wide range of drinks: A good coffee shop should offer a wide range of drinks, such as tea, hot chocolate, smoothies, and ice coffees.
6. A quiet and efficient workspace: The coffee shop should have a quiet and efficient workspace for customers who need to work while they sip their coffees.
7. Good music and ambiance: A good coffee shop should have a relaxing and comfortable ambiance, with good music and lighting that complement the overall experience.
8. A range of seating options: The coffee shop should have a range of seating options, including tables, sofas, and armchairs, to cater to customers' preferences and needs.

#### Kuditipudi et al. (2023)

A good coffee shop should have the following characteristics:

1. Quality coffee: A coffee shop should serve high-quality coffee that is well-roasted and brewed to perfection.
2. Comfortable atmosphere: The shop should have a cozy and inviting atmosphere that attracts customers for their morning coffee routine or a mid-day break.
3. Variety of beverages: Apart from coffee, a good coffee shop should also offer a variety of other beverages like tea, hot chocolate, and smoothies.

4. Friendly and attentive staff: The staff should be friendly, attentive, and knowledgeable about the menu.
5. Adequate seating: A coffee shop should have enough seating to accommodate customers who want to sit and enjoy their coffee.
6. Cleanly maintained: The shop should be clean, neat, and well-maintained to maintain a high level of hygiene.
7. Elegant and sophisticated decor: The decor should be elegant and sophisticated, giving the customers a sense of style and class.
8. Free or fast Wi-Fi: A coffee shop should offer free or fast and reliable Wi-Fi so customers can stay connected while sipping on their coffee.

### Greedy Decoding

1. Comfortable and inviting atmosphere: A good coffee shop should have a cozy and welcoming ambiance that makes customers feel at home. This includes comfortable seating, soft lighting, and pleasing decor.
2. Quality coffee: The most important characteristic of a coffee shop is the quality of the coffee. It should be freshly brewed, made with high-quality beans, and served at the right temperature.
3. Friendly and attentive staff: The staff at a coffee shop should be friendly, knowledgeable about the menu, and attentive to customers' needs. They should be able to make recommendations and provide excellent customer service.
4. Variety of menu items: A good coffee shop should offer a variety of menu items, including breakfast, lunch, and snack options. This can include sandwiches, pastries, salads, and other light bites.
5. Convenient location: A coffee shop should be located in a convenient and accessible location, such as a busy street or near a park or public transportation hub.
6. Good music: A coffee shop should have a good selection of music that complements the atmosphere and appeals to customers. This can include classical, jazz, or contemporary music.
7. Community involvement: A good coffee shop should be involved in the local community and support local events and organizations. This can include hosting events, sponsoring