
Black-Box Detection of Language Model Watermarks

Thibaud Gloaguen
Department of Mathematics
ETH Zurich
tgloaguen@student.ethz.ch

Nikola Jovanović
Department of Computer Science
ETH Zurich
nikola.jovanovic@inf.ethz.ch

Robin Staab
Department of Computer Science
ETH Zurich

Martin Vechev
Department of Computer Science
ETH Zurich

Abstract

Watermarking has emerged as a promising way to detect LLM-generated text. To apply a watermark an LLM provider, given a secret key, augments generations with a signal that is later detectable by any party with the same key. Recent work has proposed three main families of watermarking schemes, two of which focus on the property of preserving the LLM distribution. This is motivated by it being a tractable proxy for maintaining LLM capabilities, but also by the idea that concealing a watermark deployment makes it harder for malicious actors to hide misuse by avoiding a certain LLM or attacking its watermark. Yet, despite much discourse around detectability, no prior work has investigated if any of these scheme families are detectable in a realistic black-box setting. We tackle this for the first time, developing rigorous statistical tests to detect the presence of all three most popular watermarking scheme families using only a limited number of black-box queries. We experimentally confirm the effectiveness of our methods on a range of schemes and a diverse set of open-source models. Our findings indicate that current watermarking schemes are more detectable than previously believed, and that obscuring the fact that a watermark was deployed may not be a viable way for providers to protect against adversaries. We further apply our methods to test for watermark presence behind the most popular public APIs: GPT4, CLAUDE 3, GEMINI 1.0 PRO, finding no strong evidence of a watermark at this point in time.

1 Introduction

With the rapid increase in large language model (LLM) capabilities and their widespread adoption, researchers and regulators alike are raising concerns about their potential misuse for generating harmful content [1, 3]. To tackle this issue, the idea of watermarking, a process of embedding a signal invisible to humans into generated texts, is gaining significant traction.

Language model watermarking More formally, in LLM watermarking [6, 8–10, 19, 21–23] we consider the setting of a *model provider* that offers black-box access to their proprietary model LM while ensuring that each generation y in response to a prompt q can be reliably attributed to the model. To enable this, the provider modifies the generations using a secret watermark key ξ , which a corresponding watermark detector can later use to detect whether a given text was generated by LM .

The prominent family of *Red-Green* watermarks [8, 9] achieves this by, at each step of generation, selecting a context-dependent pseudorandom set of logits to be boosted, modifying the model’s next-token distribution. In contrast to these, the recently proposed families of *Fixed-Sampling* [10] and *Cache-Augmented* [6, 22] schemes have focused on developing watermarks that aim to preserve

the output distribution of the LLM, achieving such guarantees in ideal theoretical settings. This goal can be motivated by two main benefits: (i) it is a tractable proxy for preservation of capabilities of the original LLM, and (ii) intuitively, the deployment of the watermark should become inherently harder to detect for malicious actors that would otherwise be able to avoid or attack the watermark [7, 15, 19].

Practical detectability Yet, despite watermark detection being a key concern in current discourse, there have so far been no rigorous investigations into the *practical detectability* of any of the three watermark scheme families. It thus remains unclear if the distribution preservation guarantees of Fixed-Sampling and Cache-Augmented schemes, that hold in ideal settings, offer any detectability benefits in real-world deployments. Recent work has even argued that Red-Green schemes, whose detectability has also never been thoroughly studied in realistic settings, avoid some undesirable tradeoffs of distribution-preserving watermarks, and are ultimately more practical [16].

This work: practical black-box watermark detection In this work, for the first time, we comprehensively investigate the question of watermark detection in practical real-world settings. Faithful to how LLMs are exposed in current public deployments, we assume a minimal setting in which the adversary cannot control any input parameters (e.g., temperature, sampling strategy) apart from the prompt q , and does not receive any additional information (e.g., log-probabilities) apart from the instruction-tuned model’s textual response y . In this setting, the goal of the adversary is to identify if the model is watermarked and determine which scheme was used.

To this end, we propose rigorous and principled statistical tests for the black-box detection of seven scheme variants of the three most prominent watermark families. Our tests, illustrated in Fig. 1, are based on fundamental properties of the three respective scheme families: Red-Green (§2), Fixed-Sampling (§3), and Cache-Augmented (§4) watermarks. In our extensive experimental evaluation, we find that in practice, both distribution-modifying as well as distribution-preserving schemes are *easily detectable* by our tests, even in the most restricted black-box setting. Highlighting their practicality, we further test several real-world LLM deployments (GPT4, CLAUDE 3, GEMINI 1.0 PRO) and detect no watermarks, suggesting that watermark deployment at this scale is an ongoing challenge.

Implications While the impact of watermarking schemes on model capabilities (point (i) above) remains an interesting avenue for future work, our findings question whether non-detectability (point (ii) above) should be a key focus in the development of LLM watermarks. Guided by our results, we suspect that avoiding practical detectability may be an inherently hard problem given the fundamental nature of watermarks. In light of this, we advocate for a more diverse and practically oriented development and evaluation of LLM watermarks, allowing for broader consideration of other relevant properties such as attack robustness, text quality, and efficiency.

Main contributions We make the following key contributions:

- We present the first principled and statistically rigorous tests for practical black-box detection of popular LLM watermarks across the three prominent scheme families: Red-Green (§2), Fixed-Sampling (§3), and Cache-Augmented (§4).
- We confirm the effectiveness of all our tests in an extensive experimental evaluation across seven schemes and five open-source models, and present additional results that verify the applicability of our tests in real-world settings (§5.1 and §5.2).
- We further highlight the practicality of our tests by executing them on several real-world black-box LLM deployments—GPT, CLAUDE, and GEMINI (§5.3).

2 Detecting Red-Green Watermarks

In this section, we introduce our statistical test for detecting watermarks from the *Red-Green* family, illustrated in Fig. 1 (bottom). To this end, we exploit their core property: applying the watermark introduces a *context-dependent logit bias*, i.e., the model’s output distribution is biased in a way that can greatly vary depending on the last few tokens that were generated (*context*), yet is unaffected by the tokens preceding the context (*prefix*).

We first recall the background related to these schemes. We then introduce our modeling assumptions, describe the querying strategy we use to obtain the data for the test, and finally describe the test itself.



Figure 1: Overview of our key contribution. Given black-box textual access to a language model, a client can query the model and run statistical tests to rigorously test for presence of a watermark. In this example, both the test for *Cache-Augmented* watermarks (§4) and the test for *Fixed-Sampling* watermarks (§3) fail, while the test for *Red-Green* watermarks (§2) successfully detects the watermark.

In §5, we experimentally confirm the effectiveness of the test in realistic settings and study its query cost, and in App. C explore methods to estimate scheme parameters once a watermark is confirmed.

Background Assume a watermark key $\xi \in \mathbb{N}$, a pseudorandom function (PRF) f , and a hashing function $H: \mathbb{N} \rightarrow \mathbb{N}$. At each generation step t , a Red-Green watermark modifies the logits $l_{1:|V|}$ of tokens from the vocabulary V to promote a subset of tokens (*green tokens*) before applying standard sampling. We consider two popular variants, LEFTHASH [8] and SELFHASH [9], both parametrized by $\delta, \gamma \in \mathbb{R}^+$, and using $h = 1$ and $h = 3$ previous tokens as context, respectively. LEFTHASH seeds f with $H(y_{t-1}) \cdot \xi$, and uses it to split V into $\gamma|V|$ green tokens and remaining *red* tokens. For each green i , it then increases l_i by δ . SELFHASH differs by seeding f with $\min(H(y_{t-3}), \dots, H(y_t)) \cdot H(y_t) \cdot \xi$, effectively using a context size of 4 by including the token y_t yet to be generated. For both schemes, the watermark detector is based on observing a significant number of green tokens in the input text of length m , above the expectation of γm . Other schemes from this family include, among else, varying the aggregation function or the context size.

Modeling Red-Green watermarks Assume a setting where the model chooses a token from some restricted set $\Sigma \subset V$, following some *context* t_2 (longer than the watermark context h), which is preceded by a *prefix* t_1 . We discuss how to construct such a setting shortly. To model the watermark, we assume the following distribution for $p_{t_1, t_2}(x)$, the probability the model assigns to some $x \in \Sigma$:

$$p_{t_1, t_2}(x) = \frac{\exp((x^0 + \delta_{t_2}(x) + \varepsilon_{t_1, t_2}(x))/T)}{\sum_{w \in \Sigma} \exp((w^0 + \delta_{t_2}(w) + \varepsilon_{t_1, t_2}(w))/T)}. \quad (1)$$

Here, we assume the existence of the *true logit* x^0 , modeling the model bias towards x . The true logit is shifted by $\delta_{t_2}(x)$, a δ -Bernoulli random variable, where $\delta \in \mathbb{R}^+$ is the watermark parameter introduced above. Finally, $\varepsilon_{t_1, t_2}(x)$ for different t_1, t_2 are iid symmetric error terms with mean 0 and variance σ^2 . Applying the *logit* function $p \rightarrow \log(p/(1-p))$ to Eq. (1), we obtain:

$$l_{t_1, t_2}(x) = \frac{x^0}{T} + \frac{\delta_{t_2}(x)}{T} + \frac{\varepsilon_{t_1, t_2}(x)}{T} - \log \left(\sum_{w \in \Sigma \setminus \{x\}} \exp \left(\frac{w^0}{T} + \frac{\delta_{t_2}(j)}{T} + \frac{\varepsilon_{t_1, t_2}}{T} \right) \right). \quad (2)$$

Approximating the log-sum-exp with a maximum, and WLOG setting $w^0 = 0$ for w which is the maximum in the log-sum-exp (as logits are shift-invariant), the above simplifies to

$$l_{t_1, t_2}(x) = x^0/T + \delta'_{t_2}(x) + \varepsilon'_{t_1, t_2}(x), \quad (3)$$

where $\varepsilon'_{t_1, t_2}(x)$ absorbs the previous error terms. The resulting $\delta'_{t_2}(x)$ is a random variable that is 0 for unwatermarked models, and has 3 possible values for watermarked models: δ/T (where δ is the watermark parameter), if x is the only token from Σ that is in the green partition of the vocabulary, $-\delta/T$, if x is in the red partition and some other tokens from Σ are in the green partition, and 0 otherwise. Our test is based on detecting cases of $\delta'_{t_2}(x) \neq 0$ and checking that their occurrence is indeed independent of t_1 , distinguishing model variance from the watermark bias.

Querying strategy Recall that our goal is to steer a model into a setting (via an instruction, as we have no access to the completion model) where it makes a choice from some restricted set Σ . Importantly, we want to ensure enough variability in the model’s choices to be able to observe the behavior specific to a Red-Green watermark and perform efficient estimation (in terms of query cost).

Assuming an upper bound H on the context size h of the watermark, we use the following prompt template, parametrized by t_1 (an arbitrary length string), d (a single token), and a word list Σ .

f"Complete the sentence \"{t1} {d · H}\" using a random word from: [{Σ}]."

t_2

Here, t_1 serves as the prefix, and $t_2 = d \cdot H$ as the context. For a watermarked model, changing t_2 is the only way to change the red-green vocabulary split, which can greatly affect the model’s choice. For unwatermarked models, while we often see bias towards some choices from Σ , this bias will not strongly depend on t_2 . This holds for all context sizes $\leq H$ and most aggregation functions, making our setting and the corresponding test directly applicable to different variants of Red-Green schemes.

In our instantiation (illustrated in Fig. 1), we use N different t_1 of the form “*I bought*”, varying the verb, M different values of d from the set of digits, a word list Σ of four different plural fruit names, and an example that uses different words to not bias the model towards a specific choice. We empirically observe that this setup minimizes the chance that the model collapses to a single choice or fails to follow the instruction, which often occurred for similar settings based on repetition of words or numbers outside of natural context.

Estimating the logits To collect the data for the test, we query the model in two phases. First, we choose different values of Σ until we find one where the model does not always make the same choice, inducing Q_1 total queries. We set x to be the most commonly observed word from Σ . Next, for each (t_1, t_2) pair we query the model until we obtain K valid responses (filtering out failures to follow the instruction), inducing Q_2 additional queries.

We use these samples to estimate the model *logits* corresponding to x as $\hat{l}_{t_1, t_2}(x) = \log \frac{\hat{p}_{t_1, t_2}(x)}{1 - \hat{p}_{t_1, t_2}(x)}$, where $\hat{p}_{t_1, t_2}(x)$ is the empirically estimated probability of x in the responses. The result of the adversary’s querying procedure is a matrix $L_{N \times M}$ (visualized in Fig. 1) of such logit estimates.

Testing watermark presence Finally, we describe the statistical test based on the logit estimates $L_{N \times M}$. We first estimate σ , the standard deviation of $\varepsilon'_{t_1, t_2}(x)$, as follows:

$$\hat{\sigma}^2 = \text{median}[\text{Var}_{t_2}(L)], \quad (4)$$

using the empirical median to improve robustness to unpredictable behavior caused by different t_1 . Then, we calculate the following two binary functions, which flag cases where we believe the model’s probability was affected by a watermark:

$$R_x(t_1, t_2) = \mathbb{1}\{\hat{l}_{t_1, t_2}(x) - \text{median}[L] < -r\hat{\sigma}\}, \quad (5)$$

and

$$G_x(t_1, t_2) = \mathbb{1}\{\hat{l}_{t_1, t_2}(x) - \text{median}[L] > r\hat{\sigma}\}, \quad (6)$$

with $r \in \mathbb{R}^+$ a parameter of the test. In practice, to account for model unpredictability, we use the empirical median conditioned on t_1 in Eqs. (5) and (6). For simplicity, let us denote $t_1 \in \{1, \dots, N\}$ and $t_2 \in \{1, \dots, M\}$. Let $\text{cnt}_x(t_2) = \max\left(\sum_{t_1=1}^N R_x(t_1, t_2), \sum_{t_1=1}^N G_x(t_1, t_2)\right)$ count the number of consistently flagged values for fixed t_2 . We define the following test statistic:

$$S_x(L) = \max_{t_2 \in [M]} \text{cnt}_x(t_2) - \min_{t_2 \in [M]} \text{cnt}_x(t_2). \quad (7)$$

The null hypothesis of our test is $\forall t_2: \delta'_{t_2}(x) = 0$, i.e., the model is not watermarked. To obtain a p -value, we apply a Monte Carlo permutation test to S_x , checking if the flagged values are correlated with t_2 in a way that indicates a Red-Green watermark. Namely, we sample a set of permutations σ of the matrix L uniformly at random, and calculate a 99% confidence interval of $\Pr[S_x(\sigma(L)) \geq S_x(L)]$, whose upper bound we take as our p -value. When this value is small, we interpret that as evidence of a watermark. Because Eq. (3) is permutation invariant when $\delta'_{t_2}(x) = 0$, this ensures that the test does not reject under the null. This completes our method for detection of Red-Green watermarks.

3 Detecting Fixed-Sampling Watermarks

Unlike Red-Green watermarks, the recently proposed *Fixed-Sampling* watermarks do not modify the logit vectors during generation, so estimating the probabilities of model outputs as above is not informative. Instead, the sampling is fully determined by the rotation of the watermark key, making the natural vector to exploit when detecting this watermark its *lack of diversity*. Given a prompt for which an unwatermarked model is expected to produce highly diverse outputs, we can use this observation to distinguish between the two, as illustrated in Fig. 1 (middle).

As in §2, we start by introducing the needed background. We then formally model the diversity of model outputs, discuss our querying strategy that ensures our assumptions are met, and describe the resulting statistical test. The effectiveness of our method is evaluated in §5.

Background For Fixed-Sampling watermarks, the secret watermark key sequence ξ of length n_{key} is cyclically shifted uniformly at random for each generation to obtain $\bar{\xi}$, and the entry $\bar{\xi}_t$ is used to sample from l . In the ITS variant, $\bar{\xi}_t$ is a pair $(u, \pi) \in [0, 1] \times \Pi$, where Π defines the space of permutations of the vocabulary V . Given the probability distribution p over V , obtained by applying the softmax function to l , ITS samples the token with the smallest index in the permutation π such that the CDF of p with respect to π is at least u . In the EXP variant, $\bar{\xi}_t$ is a value $u \in [0, 1]$, and we sample the token $\arg \min_{i \in V} -\log(u)/p_i$. The detection, for both variants, is based on testing the correlation between the text and ξ using a permutation test. As noted in §1, the key design goal of ITS and EXP is that, in expectation w.r.t. ξ , they do not distort the distribution of the responses. How close to this ideal is achieved in practical implementations is the question we aim to answer.

Modeling the diversity Let $U_n(q, t)$ denote a random variable that counts the number of unique outputs to a fixed prompt q in n queries, each of length exactly t in tokens. We introduce the *rarefaction curve* (visualized in Fig. 1) as

$$R(n) = \mathbb{E}[U_n(q, t)]. \quad (8)$$

For suitable q that enables diversity and large enough t , the unwatermarked model exhibits $R(n) = n$. For a Fixed-Sampling watermark (both ITS and EXP variants), the watermark key segment used for sampling is determined by choosing a rotation of the key ξ uniformly at random. As choosing the same rotation for the same prompt and sampling settings will always yield the same output, the number of unique outputs is at most equal to the key size n_{key} . The probability that an output i was not produced is given by $(1 - 1/n_{key})^n$. By linearity of expectation, we have the rarefaction curve

$$R(n) = n_{key} \left(1 - \left(1 - \frac{1}{n_{key}} \right)^n \right). \quad (9)$$

Querying strategy To estimate $R(n)$ of LM , we query it with a fixed prompt q , using rejection sampling to discard short responses until we obtain a set of N responses of length t tokens (inducing Q total queries). We then repeatedly sample n responses from this set to get a Monte-Carlo estimation of $R(n)$. There are two key considerations. First, we need to ensure that we are in a setting where an unwatermarked model would have $R(n) = n$. To do this, we use the prompt “This is the story of” that reliably causes diverse outputs, and set t high enough to minimize the chance of duplicates. In §5 we experimentally confirm that the number of unique outputs scales exponentially with t , and investigate the effect of small temperatures. Second, as larger n_{key} make $R(n)$ closer to linear, we must ensure that n is large enough for our test to be reliable. To do this, we can set an upper bound \bar{n}_{key} on key size, and estimate the number of samples needed for a given power by simulating the test—our experiments show that our test succeeds even on values of \bar{n}_{key} far above practical ones.

Testing watermark presence Finally, to test for presence of a Fixed-Sampling watermark, we use a Mann-Whitney U test to compare the rarefaction curve $R(n) = n$ with the observed rarefaction data obtained as above. If the test rejects the null hypothesis, we interpret this as evidence that the model is watermarked with a Fixed-Sampling scheme. We confirm the effectiveness of our test in §5. In App. C we discuss estimation of the watermark key size n_{key} once the presence of a Fixed-Sampling watermark is confirmed.

4 Detecting Cache-Augmented Watermarks

We proceed to the detection of Cache-Augmented watermarks. While the underlying techniques used in these schemes are often similar to those of the Red-Green and Fixed-Sampling families, we focus on a general approach that exploits the presence of a cache on a fundamental level, and can be generally applied to any Cache-Augmented scheme. Namely, we note that the cache sometimes *reveals the true distribution* of the model, which was also noted in recent work [15] as an undesirable property for a watermarking scheme. This implies that the distribution of choices is *cache-conditioned* (as seen in Fig. 1, top), which our adversary will exploit to detect the presence of a watermark.

We first recall relevant background, and then discuss our method: we query the model in two phases to probe the two cache-conditioned distributions, and apply a statistical test to detect when they differ, which would indicate the presence of a Cache-Augmented watermark.

Background The watermarks that we consider in this section use previous h tokens to reweigh the distribution or apply deterministic sampling at each step. As a key feature motivated by practical distribution preservation, these schemes introduce a *cache* of previously seen contexts. Namely, whenever $y_{t-h:t-1}$ is already present in the cache, they ignore the watermarking procedure and instead fall back to standard generation. The cache can be either global or local (per user) and clears after a certain number G of generations. We consider three variants: δ -REWEIGHT [6], γ -REWEIGHT [6], and DIPMARK [22] with parameter α , the details of which we defer to App. B.

Probing the true distribution In the first phase of querying, our goal is to find a setting where the distribution of the model under the watermark will differ from its true distribution, and estimate the true distribution. For schemes we focus on, this corresponds to a setting with two choices, where the model is not significantly biased towards any of them. In particular, we use the following prompt:

Pick a fruit between: $\{f_1\}$ and $\{f_2\}$. Use the following format: $\{uc\}\{f_{example}\}$,

and modify f_1, f_2 , and $f_{example} \neq f_1, f_2$ until we find a setting where the model outputs the two choices roughly uniformly. Crucially, we prefix the prompt with a randomly sampled sufficiently long string of tokens uc . As LM will repeat uc before providing the answer, this ensures that if a cache is present, after our first query (the result of which we discard) the choice of the model will be made according to the true distribution, as the relevant part of uc was cached. Assuming WLOG that f_1 is the more likely choice for the model, we query it Q_1 times with the same input to obtain \hat{p}_1 , the estimate of the true probability of the model to pick f_1 .

Probing the watermarked distribution In the second phase, we query LM with the same input, while ensuring that the cache is reset between each query, i.e., the model will respond according to the watermarked distribution. In case of a global cache, it is sufficient to wait for a set amount of time—resetting the cache too infrequently is not a realistic setting for a deployment, as it would on average lead to a weak watermark. The uncommon prefix uc ensures that no other user will accidentally insert the same context into the cache. In case of a per-user cache, we can either saturate the cache by asking diverse queries, or use multiple user accounts. We query LM Q_2 times and obtain \hat{p}_2 , the estimate of the probability of f_1 under the watermarked distribution.

Testing watermark presence For unwatermarked models or those watermarked with a scheme from another family, both of the previous steps were sampling from the same distribution, thus for high enough Q_1 and Q_2 we expect $\hat{p}_1 = \hat{p}_2$. However, for all Cache-Augmented watermarking schemes, these probabilities will differ, indicating that the cache has revealed the true distribution of the model. To test this, we apply a Fischer’s exact test with the null hypothesis $\hat{p}_1 = \hat{p}_2$. If we observe a low p -value, we interpret this as evidence that the model is watermarked with a Cache-Augmented scheme. Our experiments in §5 demonstrate that this test is robust to different scheme variants, and does not lead to false positives when LM is unwatermarked or uses a scheme from another family. In App. C we discuss how to distinguish different variants of Cache-Augmented schemes.

5 Experimental Evaluation

In this section, we apply the tests introduced in §2–§4 to a wide range of models and watermarking schemes, and confirm their effectiveness in detecting watermarks. In §5.1 we show that the adversary

Table 1: Main results of our watermark detection tests across different models and watermarking schemes. We report median p-values across 100 repetitions of the experiment, and for RED-GREEN schemes additionally over 5 watermarking keys. p-values below 0.05 (test passing) are highlighted in bold. δR and γR denote δ -REWEIGHT and γ -REWEIGHT schemes, respectively.

Model	Test	Unwatermarked		Red-Green				Fixed-Sampling		Cache-Augmented		
		$T = 1.0$	$T = 0.7$	LEFTHASH		SELFHASH		ITS/EXP		DIPMARK/ γR		δR
				$\delta, \gamma = 2, 0.25$	$\delta, \gamma = 4, 0.5$	$\delta, \gamma = 2, 0.5$	$\delta, \gamma = 4, 0.25$	$n_{key} = 256$	$n_{key} = 2048$	$\alpha = 0.3$	$\alpha = 0.5$	
MISTRAL 7B	R-G (§2)	1.000	1.000	0.000	0.000	0.000	0.000	1.000	1.000	1.000	1.000	1.000
	Fixed (§3)	0.938	0.938	0.938	0.938	0.938	0.938	1.3e-125	1.1e-9	0.938	0.938	0.938
	Cache (§4)	0.570	0.667	0.607	0.608	1.00	0.742	0.638	0.687	2.4e-4	2.1e-3	5.6e-27
LLAMA2 13B	R-G (§2)	0.149	0.663	0.000	0.000	0.000	0.000	0.121	0.128	0.149	0.149	0.149
	Fixed (§3)	0.968	0.868	0.938	0.938	0.938	0.938	5.5e-124	7.5e-8	0.938	0.938	0.968
	Cache (§4)	0.708	0.573	0.511	0.807	0.619	0.710	0.518	0.692	1.8e-2	5.3e-3	6.7e-32
LLAMA2 70B	R-G (§2)	1.000	1.000	0.000	0.020	0.000	0.000	1.000	1.000	1.000	1.000	1.000
	Fixed (§3)	0.938	0.526	0.966	0.965	0.984	0.975	4.2e-125	1.7e-8	0.938	0.966	0.938
	Cache (§4)	0.596	0.620	0.657	0.639	0.651	0.608	0.463	0.818	1.5e-3	4.4e-3	5.8e-28

can reliably detect the watermarking scheme used (if any) at a low cost, across a wide range of practical settings. In §5.2, we provide further experimental insights into the tests’ robustness. Finally, in §5.3, we demonstrate that our tests can be directly applied to real-world LLM deployments, revealing no significant evidence of any watermarking schemes at this point in time.

5.1 Main Experiments: Detecting Watermarking Schemes

We perform our main set of experiments to verify the effectiveness of our tests introduced in §2–§4 in detecting watermarking schemes in realistic scenarios.

Experimental setup We run all our tests on five different models (MISTRAL-7B, LLAMA2-7B, -13B, and -70B, and LLAMA3-8B), in eleven different scenarios. We here present results of a subset of those, and defer the rest to App. A. In each scenario, each model is either unwatermarked (where we vary the temperature) or watermarked with a certain scheme from the three main families (where we vary the particular scheme and its parameters). If our tests are reliable, we expect to see low p-values only when the model is watermarked exactly with the scheme family that we are testing for.

For Red-Green tests, we set $N = 10$, $M = 9$, $r = 1.96$, a different Σ per model based on the first $Q1$ samples, use 100 samples to estimate the probabilities, and use 10000 permutations in the test. For Fixed-Sampling tests, we use $n = 1000$ queries and set $t = 50$. For Cache-Augmented tests, we use $Q1 = Q2 = 75$ and assume the cache is cleared between queries in the second phase.

Results: reliable watermark detection Our main results are shown in Table 1, where we report the median p-values for each (model, test, scenario) tuple across 100 repetitions of each experiment. Across all experiments, all three tests reject the null hypothesis (*the specific watermark is not present*) at a 95% confidence level only when the scheme from the target family is indeed applied to the model. This confirms that our tests are reliable in detecting watermarks, and robust with respect to the model and the specific parameters of the scheme, emphasizing our tests’ generalization to all schemes that are based on the same foundational principles. In particular, our Red-Green tests are robust to the seeding scheme, the logit bias δ , and the green token fraction γ ; our Fixed-Sampling tests maintain high confidence even for unusually high values of n_{key} and those higher than the number of queries; finally, our Cache-Augmented tests are robust to all details of the underlying scheme. Our results also provide evidence that unrelated model modifications do not cause false positives, as no test passes when the model is unwatermarked or watermarked with a different scheme family.

Our tests do not incur significant costs for the adversary, making them easily applicable in practice. While the cost of a particular run varies across models and scenarios, we can estimate the average cost of the tests to be below \$20 for Red-Green, \$1 for Fixed-Sampling, and \$0.1 for Cache-Augmented tests, assuming latest OpenAI GPT40 pricing.

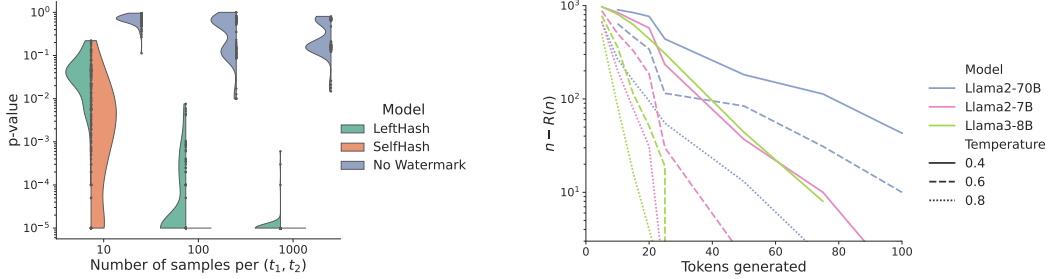


Figure 2: *Left*: distribution of bootstrapped p-values of the Red-Green test on LLAMA2-13B with $(\delta, \gamma) = (2, 0.25)$, for different sample sizes. We see reliable results for 100 or more samples. *Right*: the diversity gap $n - R(n)$ on log scale in different settings. Linear behavior means that diversity scales exponentially with t , and we see that the assumption of $R(n) = n$ can be easily met in practice.

5.2 Additional Experiments: Validating the Assumptions

We present two additional experiments to validate the assumptions made in our tests and provide further insight into their behavior in practical settings.

Sampling in Red-Green tests The Red-Green test (§2) relies on the sampling of model outputs to estimate the probabilities of the model. As the resulting p-value is computed assuming knowledge of true probabilities, this raises the question of the impact of the sampling error on our results. To heuristically mitigate this, we propose a bootstrapping procedure, where for fixed (t_1, t_2) , we sample with replacement from a single set of model outputs, and report the median p-value p_{med} across such samples. In Fig. 2 (Left) we report the resulting distribution of p_{med} , where one point corresponds to one independent run. We see that already for 100 samples per (t_1, t_2) (as used in Table 1), the p-value distribution is narrow and the false positive rate due to sampling error is well controlled. This confirms that our test is robust against sampling error and can be used in realistic settings, without additional access to model logprobs. For computational reasons, we did not apply this correction in Table 1, where we still observe reliable results in the median case across experiment repetitions.

Diversity assumption in Fixed-Sampling tests As detailed in §3, the Fixed-Sampling test relies on the unwatermarked model being sufficiently diverse, i.e., for the number of unique outputs $R(n) = \mathbb{E}[U_n(q, t)]$ after n queries with prompt q and response length t , it should hold that $R(n) = n$. Our goal is to show that we can easily choose t such that this property holds across different settings.

To this end, we hypothesize that the number of unique outputs converges to n exponentially fast as the response length t increases. In particular, we assume

$$R(n) = n - \lfloor n \cdot \exp(-\alpha(T)t) \rfloor, \quad (10)$$

where $\alpha(T)$ is a monotonic function of the temperature T . To verify this hypothesis, we measure $n - R(n)$ on several models and temperatures, and show the result on log scale in Fig. 2 (Right). If Eq. (10) holds, we expect the resulting relationship to be linear, which is indeed confirmed by our results. While α (the slope of the line) varies across settings, we see that a bit over 200 tokens would be sufficient for the line to drop to 0 (not visible on the log scale plot). This holds even in cases impractical for deployment of Fixed-Sampling watermarks such as $T = 0.4$ [10, 16], indicating that $R(n) = n$ and our assumption is met, validating our p-values.

5.3 Detecting watermarks in deployed models

Finally, we demonstrate the applicability of the statistical tests introduced in §2–§4, by applying them to popular black-box LLM deployments: GPT4, CLAUDE 3, and GEMINI 1.0 PRO. We use the same experimental setup as in §5.1, and use the API access for efficiency reasons—we do not rely on any additional capabilities,

and our tests could be as easily run in the web interface. For the Cache-Augmented tests, we assume a global cache, that if present clears after 1000 seconds. For the Fixed-Sampling test on CLAUDE 3, due to its lack of diversity, we used $t = 75$ tokens per query to ensure the test’s hypotheses are met.

Table 2: The results of our watermark detection tests on popular black-box LLM deployments.

	GPT4	CLAUDE 3	GEMINI 1.0 PRO
R-G (§2)	0.998	0.638	0.683
Fixed (§3)	0.938	0.844	0.938
Cache (§4)	0.51	0.135	0.478

Results Our results in Table 2 show that the null hypothesis is not rejected for any of the models and any of the watermark tests. This suggests that none of the deployed models tested are watermarked at this point in time. The demonstrated applicability of our tests makes it simple to monitor these deployments and detect any watermarking schemes that may be introduced in the future.

6 Related Work

Language model watermarking Besides the approaches by [6, 8, 10] introduced above, there are various methods building on similar ideas. Hou et al. [5], Liu et al. [13], Ren et al. [18] all apply variations of [8] on semantic information, while Gu et al. [4] aims to distill a new model from the output of a watermarked model. Similarly, Liu et al. [12] apply a Red-Green scheme using a learned classifier instead of hash functions. A range of works on multi-bit watermarking [21, 23] aim to not only watermark generated texts but encode additional information in the watermark itself.

Attacks on language model watermarking Attacks on LLM watermarks have so far been mainly investigated in terms of scrubbing [7, 8, 19] (i.e., removal of a watermark) and spoofing [4, 7, 19] (i.e., applying a watermark without knowing ξ). Notably, Jovanović et al. [7] showed that observing watermarked texts can facilitate both attacks on various distribution-modifying schemes, disproving common robustness assumptions [9]. However, using this and similar attacks as means of practical watermark detection is infeasible, as they generally offer no way to quantify the attack’s success—in contrast, we aim to provide rigorous statements about scheme presence. Further, such attacks incur significantly higher query costs than necessary for detection (as our work demonstrates), and in some cases assume certain knowledge of the watermark parameters, a setting fundamentally at odds with our threat model of black-box watermark detection.

The closest related work to ours is Tang et al. [20], that tackles the problem of watermark detection in strictly simpler settings where the adversary either has access to an unwatermarked counterpart of the target model, or can access full model logits. Such knowledge is commonly not available in practice, limiting the applicability of this approach. To the best of our knowledge, no work has developed methods for detecting the presence of a watermark in a realistic black-box setting.

Extracting data from black-box models With many of the most potent LLMs being deployed behind restricted APIs, the extraction of model details has been an active area of research. This includes, e.g., the reconstruction of a black-box model tokenizer [17] as well as the extraction of the hidden-dimensionality or the weights of the embedding projection layer [2]. Naseh et al. [14] have shown practical attacks to recover the decoding mechanism of non-watermarked black-box models. Given access to output logits, Li and Fung [11] have further demonstrated that it is possible to train an inversion model that aims to recover the input prompt.

7 Conclusion, Impact and Limitations

In this paper, we have focused on the problem of detecting watermarks in large language models (LLMs) given only black-box access. We developed rigorous statistical tests for the three most prominent scheme families, and validated their effectiveness on a wide range of schemes and real-world models. Our results imply that protecting from detectability should not be the key focus of LLM watermarks, and that other properties such as robustness to attacks, text quality, and efficiency should be prioritized.

Broader Impact While our work primarily enables malicious parties to more easily circumvent attempts at tracing LLM-generated text, we believe the benefits of our work outweigh those risks, as our conclusions help model providers calibrate their expectations in terms of watermark detectability, and avoid a false sense of security.

Limitations One limitation of our tests is that they are restricted to the three scheme families discussed in §2–§4. These are indeed the most prominent in the literature, and as our tests are based on fundamental properties of these scheme families, they should generalize to more variants and combinations of the underlying ideas. However, it is possible that a model provider deploys a scheme based on an entirely novel idea, which our tests would not be able to detect. Further, our p-values are

based on several model assumptions, such as symmetric error terms and perfect sampling. While we validate that these assumptions are sufficiently met on several open-source models, we cannot guarantee that all models adhere to them. Finally, another conceptual limitation of our Red-Green test is that it does not take into account the possibility that the red-green vocabulary split is the same (on the observed domain) for all contexts. The probability of this event decreases exponentially with the number of different contexts, and is thus unlikely to affect the test’s performance in practice.

References

- [1] Bommasani, R., Hudson, D. A., Adeli, E., Altman, R. B., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N. S., Chen, A. S., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N. D., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M. S., Krishna, R., Kuditipudi, R., and et al. (2021). On the opportunities and risks of foundation models. *CoRR*.
- [2] Carlini, N., Paleka, D., Dvijotham, K. D., Steinke, T., Hayase, J., Cooper, A. F., Lee, K., Jagielski, M., Nasr, M., Conmy, A., Wallace, E., Rolnick, D., and Tramèr, F. (2024). Stealing part of a production language model. *CoRR*.
- [3] Council of the European Union (2024). Proposal for a regulation of the european parliament and of the council laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts - analysis of the final compromise text with a view to agreement.
- [4] Gu, C., Li, X. L., Liang, P., and Hashimoto, T. (2024). On the learnability of watermarks for language models. *ICLR*.
- [5] Hou, A. B., Zhang, J., He, T., Wang, Y., Chuang, Y., Wang, H., Shen, L., Durme, B. V., Khashabi, D., and Tsvetkov, Y. (2023). Semstamp: A semantic watermark with paraphrastic robustness for text generation. *arXiv*.
- [6] Hu, Z., Chen, L., Wu, X., Wu, Y., Zhang, H., and Huang, H. (2024). Unbiased watermark for large language models. *ICLR*.
- [7] Jovanović, N., Staab, R., and Vechev, M. (2024). Watermark stealing in large language models. *arXiv*.
- [8] Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., and Goldstein, T. (2023). A watermark for large language models. In *ICML*.
- [9] Kirchenbauer, J., Geiping, J., Wen, Y., Shu, M., Saifullah, K., Kong, K., Fernando, K., Saha, A., Goldblum, M., and Goldstein, T. (2024). On the reliability of watermarks for large language models. *ICLR*.
- [10] Kuditipudi, R., Thickstun, J., Hashimoto, T., and Liang, P. (2024). Robust distortion-free watermarks for language models. *TMLR*.
- [11] Li, Y. and Fung, P. (2013). Improved mixed language speech recognition using asymmetric acoustic model and language model with code-switch inversion constraints. In *ICASSP*.
- [12] Liu, A., Pan, L., Hu, X., Li, S., Wen, L., King, I., and Yu, P. S. (2024a). A private watermark for large language models. *ICLR*.
- [13] Liu, A., Pan, L., Hu, X., Meng, S., and Wen, L. (2024b). A semantic invariant robust watermark for large language models. *ICLR*.
- [14] Naseh, A., Krishna, K., Iyyer, M., and Houmansadr, A. (2023). On the risks of stealing the decoding algorithms of language models. *CoRR*.
- [15] Pang, Q., Hu, S., Zheng, W., and Smith, V. (2024). Attacking LLM watermarks by exploiting their strengths. *arXiv*.

- [16] Piet, J., Sitawarin, C., Fang, V., Mu, N., and Wagner, D. A. (2023). Mark my words: Analyzing and evaluating language model watermarks. *arXiv*.
- [17] Rando, J. and Tramèr, F. (2024). The worst (but only) claude 3 tokenizer. <https://github.com/javirandor/anthropic-tokenizer>.
- [18] Ren, J., Xu, H., Liu, Y., Cui, Y., Wang, S., Yin, D., and Tang, J. (2023). A robust semantics-based watermark for large language model against paraphrasing. *arXiv*.
- [19] Sadasivan, V. S., Kumar, A., Balasubramanian, S., Wang, W., and Feizi, S. (2023). Can ai-generated text be reliably detected? *arXiv*.
- [20] Tang, L., Uberti, G., and Shlomi, T. (2023). Baselines for identifying watermarked large language models. *CoRR*.
- [21] Wang, L., Yang, W., Chen, D., Zhou, H., Lin, Y., Meng, F., Zhou, J., and Sun, X. (2024). Towards codable text watermarking for large language models. *ICLR*.
- [22] Wu, Y., Hu, Z., Zhang, H., and Huang, H. (2023). Dipmark: A stealthy, efficient and resilient watermark for large language models. *arXiv*.
- [23] Yoo, K., Ahn, W., and Kwak, N. (2024s). Advancing beyond identification: Multi-bit watermark for language models. *NAACL*.

Table 3: Additional results of our watermark detection tests across different models and watermarking schemes. We report median p-values across 100 repetitions of the experiment, and for RED-GREEN schemes additionally over 5 watermarking keys. p-values below 0.05 (test passing) are highlighted in bold. δR and γR denote δ -REWEIGHT and γ -REWEIGHT schemes, respectively.

Model	Test	Red-Green										Fixed-Sampling		Cache-Augmented		
		Unwatermarked		LEFTHASH				SELFHASH				$n = 256$	$n = 2048$	$\alpha = 0.3$	$\alpha = 0.5$	δR
		$T = 1.0$	$T = 0.7$	$\delta, \gamma = 2, 0.25$	$\delta, \gamma = 2, 0.5$	$\delta, \gamma = 4, 0.25$	$\delta, \gamma = 4, 0.5$	$\delta, \gamma = 2, 0.25$	$\delta, \gamma = 2, 0.5$	$\delta, \gamma = 4, 0.25$	$\delta, \gamma = 4, 0.5$					
MISTRAL 7B	R-G (§2)	1.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	1.000	1.000	1.000	1.000
	Fixed (§3)	0.938	0.938	0.938	0.938	0.938	0.938	0.938	0.938	0.938	0.938	1.3e-125	1.1e-9	0.938	0.938	0.938
	Cache (§4)	0.570	0.667	0.607	0.639	0.632	0.608	1.000	1.000	0.742	0.824	0.638	0.687	2.4e-4	2.1e-3	5.6e-27
LLAMA2 13B	R-G (§2)	0.149	0.663	0.000	0.014	0.000	0.000	0.000	0.000	0.000	0.000	0.121	0.128	0.149	0.149	0.149
	Fixed (§3)	0.968	0.868	0.938	0.966	0.938	0.938	0.938	0.938	0.938	0.967	5.5e-124	7.5e-8	0.938	0.938	0.968
	Cache (§4)	0.708	0.573	0.511	0.596	0.623	0.807	0.657	0.619	0.710	0.583	0.518	0.692	1.8e-2	5.3e-3	6.7e-32
LLAMA2 70B	R-G (§2)	1.000	1.000	0.000	0.000	0.000	0.020	0.000	0.020	0.000	0.000	1.000	1.000	1.000	1.000	1.000
	Fixed (§3)	0.938	0.526	0.966	0.964	0.889	0.965	0.963	0.984	0.975	0.990	4.2e-125	1.7e-8	0.938	0.966	0.938
	Cache (§4)	0.596	0.620	0.657	0.797	0.824	0.639	0.535	0.651	0.608	0.593	0.463	0.818	1.5e-3	4.4e-3	5.8e-28
LLAMA2 7B	R-G (§2)	1.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	1.000	1.000	1.000	1.000
	Fixed (§3)	0.993	0.988	0.938	0.938	0.938	0.938	0.968	0.938	0.938	0.938	7.6e-124	1.0e-8	0.938	0.938	0.993
	Cache (§4)	0.604	0.602	0.623	0.705	0.728	0.593	0.620	0.718	0.610	0.593	0.476	0.588	4.2e-6	4.5e-7	1.3e-21
LLAMA3 8B	R-G (§2)	1.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	1.000	1.000	1.000	1.000
	Fixed (§3)	0.938	0.938	0.938	0.938	0.938	0.938	0.938	0.938	0.986	0.938	1.2e-124	3.8e-8	0.938	0.938	0.938
	Cache (§4)	0.734	0.504	0.605	0.514	0.712	0.605	0.600	0.731	0.729	0.714	0.618	0.605	5.2e-5	3.2e-8	3.5e-18

A Extending Experiments across Diverse Settings

Watermark detection with additional models and schemes We extend the experiments from §5.1 using two additional models, LLAMA2-7B and LLAMA3-8B, as well as more variations of the watermarking schemes’ parameters to further assess the robustness of the tests. The experimental setup for the additional results is consistent with the one described in §5.1.

Our exhaustive results are presented in Table 3. The same conclusion applies to the two additional models: the null hypothesis (*the specific watermark is not present*) is rejected at a 95% confidence level only when the corresponding watermarking scheme is applied to the model. These results confirm that the modeling assumptions for each test are satisfied across a wide range of models, indicating the tests’ relevance in practical scenarios.

Multiple keys in Red-Green watermarks To demonstrate that the Red-Green test is robust to variations in the watermarking scheme within the same watermark family, we consider the case of Red-Green watermarks with multiple keys, where the key ξ is uniformly selected from a predefined pool of keys at each generation. Using n keys turns Eq. (3) into

$$l_{t_1, t_2}(x) = x^0/T + \delta''_{t_2}(x) + \varepsilon'_{t_1, t_2}(x), \quad (11)$$

with $\delta''_{t_2}(x)$ in $\{k\delta/(nT) \mid \forall k \in \{-n, \dots, n\}\}$ is obtained by averaging the variables $\delta'_k(x)$ over the set of keys. Despite modeling changes, the core assumption of logit bias being conditioned on t_2 remains unchanged. Therefore, we can apply the same test as in §2 to detect the watermark. Consequently, we conducted the Red-Green test on both the LeftHash and SelfHash variants using $n = 3$ keys on three models (LLAMA2-13B, LLAMA2-70B and MISTRAL-7B). Recent work [15] shows that using too many keys can lead to other vulnerabilities.

Across all three models and scheme parameters, the null hypothesis (*the Red-Green watermark is not present*) is rejected at a 95% confidence level, with median p-values lower than $1e-4$ for each combination of model and setting. It shows that the Red-Green test is robust even in settings that slightly deviate from the original modeling considered in §2. It emphasizes the test’s reliance on the foundational principles behind Red-Green schemes rather than on specific implementation details.

B Details of Cache-Augmented Schemes

We provide the details of the three Cache-Augmented watermarking schemes considered in this work: δ -REWEIGHT [6], γ -REWEIGHT [6], and DIPMARK [22], that were omitted from §4.

All three variants, at each generation step t , jointly hash the watermark key $\xi \in \mathbb{Z}_2^K$ and the preceding context $y_{t-h:t-1}$ (commonly setting $h = 5$) using SHA256, and use the result as a seed to sample a code $E_t \in P_E$ uniformly at random. Let p denote the probability distribution over V , obtained by applying the softmax function to the logits. For the δ -REWEIGHT variant, $P_E = [0, 1]$, and the code E_t is used to sample the token in V with the smallest index, such that the CDF of p is at least E_t . For the γ -REWEIGHT variant and DIPMARK, P_E is the space of permutations of V . For γ -REWEIGHT, we transform p to a new distribution p' by, for each token $i \in V$, setting $p'(i) = f_2(f_1(E_t(i))) - f_2(f_1(E_t(i) - 1))$, where we have $f_1(i') = \sum_{j \in V} \mathbb{1}(E_t(j) \leq i')p(j)$ and $f_2(v) = \max(2v - 1, 0)$, effectively dropping the first half of the permuted CDF. For DIPMARK, given parameter $\alpha \in [0, 0.5]$, this is generalized by using $f_2(v) = \max(v - \alpha, 0) + \max(v - (1 - \alpha), 0)$, recovering γ -REWEIGHT for $\alpha = 0.5$. The former two variants perform detection using a log-likelihood ratio test (requiring access to LM), while DIPMARK uses a model-independent test.

C Estimating Scheme Parameters

In this section, we describe how, mostly leveraging the queries already performed during the detection tests, we can estimate the main parameters of the detected watermarking scheme. This demonstrates the soundness of our modeling assumptions in §2-§4 from which we derive all the estimators.

C.1 Estimation of Red-Green Watermarking Scheme Parameters

If the null hypothesis (*the model not being Red-Green watermarked*) is rejected, we can then estimate the scheme-specific watermark parameters δ and the context size h using mostly the same data as used in the test. First, we describe the estimators for both parameters, and then discuss their practicality by analyzing their performance on multiple models.

Description of estimators To estimate δ , we establish a parametrized model based on Eq. (2) that relies on our flagging function from Eq. (6), and additional estimates $\hat{l}_{t_1, t_2}(w)$ for every $w \in \Sigma$, computed on the same data as above, requiring no additional queries. For each $w \in \Sigma$, we set:

$$\hat{l}_{t_1, t_2}(w) = \bar{w}_{t_1} + G_w(t_1, t_2)\bar{\delta} - \log \left(\sum_{w' \in \Sigma \setminus \{w\}} \exp(\bar{w}_{t_1} + G_{w'}(t_1, t_2)\bar{\delta}) \right), \quad (12)$$

and set $\forall t_1, \bar{w}_{t_1} = 0$ for a single $w \in \Sigma$, as logits are shift-invariant. Fitting the parameters $\bar{\delta}$ and all \bar{w}_{t_1} by minimizing the mean squared error with gradient descent allows us to recover δ/T as $\bar{\delta}$. If T is known or estimated separately, this term can be removed.

Let h denote the context size, i.e., the number of *previous* tokens considered by the watermarking scheme. To estimate h , we use the same prompt template as in §2, with a fixed prefix t_1 and digit d , but with a varying H and perturbation digit d' prepended to d .

```
f"Complete the sentence \"{t1} {d'}{d \cdot H}\" using a random word from: [{\Sigma}]."
```

The probability distribution of the model output will abruptly change when H exceeds the context size h , as the change in d' will not alter the red/green split of the vocabulary. By performing a series of pairwise Mood's tests on the estimated log-probabilities of x we can find the largest H for which t_1 is part of the context. This corresponds to the first value of H for which the test is rejected.

Estimating γ is more challenging, as in contrast to δ , this parameter is not directly reflected in the logits but rather defines a more global behavior of the scheme. This is particularly difficult for schemes with self-seeding, as the rejection sampling interferes with the behavior of γ . We leave further exploration of this problem to future work.

Experimental results We computed the estimator for δ on the LeftHash variant with $\gamma = 0.25$ and varying δ from 0 to 4. The results are shown in Fig. 3, with the 95% confidence intervals reflecting the sampling error. The estimator successfully estimates δ for all models with a sufficient number of samples, using only the estimated output probabilities of the model. It is also consistent across all models tested, suggesting that the model assumptions in Eq. (1) are met in practice.

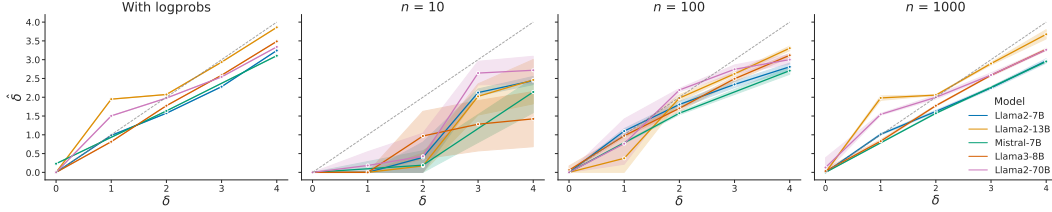


Figure 3: Estimation of δ for different models using LeftHash with $\gamma = 0.25$. The number of samples used increases from left to right, with the leftmost plot assuming direct access to the log-probabilities. The estimation is done on the same data as the test. Error bars are given by the 95% bootstrapped confidence interval with respect to the sampling of the model outputs.

Table 4: Key length estimation for Fixed-Sampling watermarks using non-linear regression on the rarefaction curve.

Key length	LLAMA2-7B	LLAMA2-13B	LLAMA2-70B	LLAMA3-8B	MISTRAL-7B
256	259 ± 0.6	259 ± 0.5	256 ± 0.5	257 ± 0.5	256 ± 0.6
2048	1978 ± 10	2107 ± 12	2006 ± 13	2070 ± 14	1831 ± 10

Estimating the context size for Red-Green schemes requires performing a new attack once the model is flagged as watermarked. We estimate the context size for three different models (MISTRAL-7B, LLAMA2-13B and LLAMA2-70B) using LeftHash with $\delta = 2$ and $\gamma = 0.25$. The estimation process requires an additional 5,000 queries, and the estimator successfully determines the context size for all models. However, the estimator is less robust on the SelfHash variant due to the self-seeding algorithm, which leads to higher probabilities for tokens in Σ being in the green vocabulary, and thus diminishes the perturbation’s significance and resulting in false negatives in Mood’s test. Therefore, the procedure stated above produces a lower bound for the context size. To mitigate this issue, we use the estimator across 10 different t_2 and then consider the median of the 10 estimators as our final estimator. This estimator applied on the SelfHash variant with $\delta = 2$ and $\gamma = 0.25$ is successful on all three models. It also does not change the results on LeftHash and can be used as a more robust estimator for h in all cases, when the additional cost of 50, 000 queries is not prohibitive.

C.2 Estimation of Fixed-Sampling Watermarking Scheme Parameters

Our approach does not distinguish between the variant of the Fixed-Sampling watermark used (ITS or EXP), as the diversity property that we exploit is common to both. The only other relevant parameter of Fixed-Sampling schemes is n_{key} . To estimate it, we use non-linear regression on the rarefaction curve using (9) and the same data that we used for the presence test, and compute the confidence intervals using bootstrapping.

Our results are given in Table 4. We see that the estimator is consistent across different models and remains relatively precise even for values of n_{key} higher than the number of queries.

C.3 Estimation of Cache-Augmented Watermarking Scheme Parameters

For Cache-Augmented watermarks, we can estimate which scheme variant is present, and if the variant is DIPMARK, attempt to learn the value of α (recall that $\alpha = 0.5$ corresponds to γ -REWEIGHT). To do this, we use the same approach as in §4 to obtain \hat{p}_1 and \hat{p}_2 , where WLOG we assumed $p_1 > 0.5$. If we observe $\hat{p}_2 = 0$ this directly implies a δ -REWEIGHT watermark. If we observe $\hat{p}_2 \in (0, 1)$, we learn the following: if $\hat{p}_2 = 2\hat{p}_1 - 1$ then $\alpha > 1 - \hat{p}_1$, otherwise $\alpha = |\hat{p}_1 - \hat{p}_2|$. The bound in the first case can be further improved with additional queries with different p_1 . Finally, if we observe $\hat{p}_2 = 1$ we repeat the whole procedure in total K times, following the same case distinction—if $\hat{p}_2 = 1$ repeats in all K runs, we conclude that the model is watermarked with δ -REWEIGHT.

Using the same parameters as the one for the test, we distinguish with 100% accuracy between a DIPMARK and a δ -REWEIGHT watermark. However, the estimation of α becomes unreliable for higher values of α , especially for smaller models. One of the reasons for this are the failures of the model to follow the instruction, that are more common in the presence of the uncommon prefix *uc*. While the detection test in §4 was robust to such behavior, this does not hold for the estimation of α .

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The claims are the following: proposing statistical tests to detect watermarking schemes (done in §2–§4) / confirming the effectiveness of the proposed tests and estimating main parameters (done in §5) / demonstrating the practicality of the attacks on deployed LLMs (done in §5.3).

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitations are discussed in §7. Experiments to show the limitations of the proposed tests are done in §5.2.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: No major theoretical results but modeling assumptions for statistical tests.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: All the test-specific details are provided in §5, including the model used, and the different hyperparameters used each time. The exact prompts are provided in the code.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code is available and commented. The models used for most experiments are either open-source or accessible via an API.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: For each experiment, hyperparameters are documented in §5. The code used for the experiments is provided as well.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Error bars are reported in plots in §5 with respect to the parameters leading to the highest variance in the results. All results in tables are presented by aggregating over multiple parameters and seeds.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: This is provided in the introduction of §5.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: Impacts of the presented results are discussed in §7. The research process was conducted in a fair and ethical manner.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The broader impact of this work is discussed in §1 and §7, weighing both the positive and negative societal impacts.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No data nor models are released with this paper.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The authors of the different watermarking schemes are credited throughout the writeup. The open-source models used for experiments are mentioned in §5.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The provided code is suitably commented.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing nor research with human subjects was conducted.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No human subjects were involved in this research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.