

# A Watermark for Large Language Models

John Kirchenbauer<sup>\*1</sup> Jonas Geiping<sup>\*1</sup> Yuxin Wen<sup>1</sup> Jonathan Katz<sup>1</sup> Ian Miers<sup>1</sup> Tom Goldstein<sup>1</sup>

## Abstract

Potential harms of large language models can be mitigated by *watermarking* model output, i.e., embedding signals into generated text that are invisible to humans but algorithmically detectable from a short span of tokens. We propose a watermarking framework for proprietary language models. The watermark can be embedded with negligible impact on text quality, and can be detected using an efficient open-source algorithm without access to the language model API or parameters. The watermark works by selecting a randomized set of “green” tokens before a word is generated, and then softly promoting use of green tokens during sampling. We propose a statistical test for detecting the watermark with interpretable p-values, and derive an information-theoretic framework for analyzing the sensitivity of the watermark. We test the watermark using a multi-billion parameter model from the Open Pretrained Transformer (OPT) family, and discuss robustness and security.

## 1. Introduction

Large language models (LLMs), such as the recently developed ChatGPT, can write documents, create executable code, and answer questions, often with human-like capabilities (Schulman et al., 2022). As these systems become more pervasive, there is increasing risk that they may be used for malicious purposes (Bergman et al., 2022; Mirsky et al., 2023). These include social engineering and election manipulation campaigns that exploit automated bots on social media platforms, creation of fake news and web content, and use of AI systems for cheating on academic writing and coding assignments. Furthermore, the proliferation of synthetic data on the web complicates future dataset creation efforts, as synthetic data is often inferior to human content and must be detected and excluded before model training

Code is available at <https://www.github.com/jwkirchenbauer/lm-watermarking>. \*Equal contribution <sup>1</sup>University of Maryland. Correspondence to: John Kirchenbauer <jkirchen@umd.edu>.

*Proceedings of the 40<sup>th</sup> International Conference on Machine Learning*, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

Prompt	Num tokens	Z-score	p-value
...The watermark detection algorithm can be made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API. We seek a watermark with the following properties:			
<b>No watermark</b> Extremely efficient on average term lengths and word frequencies on synthetic, microamount text (as little as 25 words) Very small and low-resource key/hash (e.g., 140 bits per key is sufficient for 99.999999999% of the Synthetic Internet)	56	.31	.38
<b>With watermark</b> - minimal marginal probability for a detection attempt. - Good speech frequency and energy rate reduction. - messages indiscernible to humans. - easy for humans to verify.	36	7.4	6e-14

Figure 1. Outputs of a language model, both with and without the application of a watermark. The watermarked text, if written by a human, is expected to contain 9 “green” tokens, yet it contains 28. The probability of this happening by random chance is  $\approx 6 \times 10^{-14}$ , leaving us *extremely* certain that this text is machine generated. Words are marked with their respective colors. The model is OPT-6.7B using multinomial sampling. Watermark parameters are  $\gamma, \delta = (0.25, 2)$ . The prompt is the whole blue paragraph marked below.

(Radford et al., 2022). For many reasons, the ability to detect and audit the usage of machine-generated text becomes a key principle of harm reduction for large language models (Bender et al., 2021; Crothers et al., 2022; Grinbaum & Adomaitis, 2022).

In this work, we study *watermarking* of language model output. A watermark is a hidden pattern in text that is imperceptible to humans, while making the text algorithmically identifiable as synthetic. We propose an efficient watermark that makes synthetic text detectable from short spans of tokens (as few as 25 tokens), while false-positives (where human text is marked as machine-generated) are statistically improbable. The watermark detection algorithm can be

made public, enabling third parties (e.g., social media platforms) to run it themselves, or it can be kept private and run behind an API. We seek a watermark with the following properties:

- The watermark can be algorithmically detected without any knowledge of the model parameters or access to the language model API. This property allows the detection algorithm to be open sourced even when the model is not. This also makes detection cheap and fast because the LLM does not need to be loaded or run.
- Watermarked text can be generated using a standard language model without re-training.
- The watermark is detectable from only a contiguous portion of the generated text. This way, the watermark remains detectable when only a slice of the generation is used to create a larger document.
- The watermark cannot be removed without modifying a significant fraction of the generated tokens.
- We can compute a rigorous statistical measure of confidence that the watermark has been detected.

### 1.1. Notation & Language model basics

Language models have a “vocabulary”  $\mathcal{V}$  containing words or word fragments known as “tokens.” Typical vocabularies contain  $|\mathcal{V}| = 50,000$  tokens or more (Radford et al., 2019; Liu et al., 2019). Consider a sequence of  $T$  tokens  $\{s^{(t)}\} \in \mathcal{V}^T$ . Entries with negative indices,  $s^{(-N_p)}, \dots, s^{(-1)}$ , represent a “prompt” of length  $N_p$  and  $s^{(0)}, \dots, s^{(T)}$  are tokens generated by an AI system in response to the prompt.

A *language model* (LM) for next word prediction, is a function  $f$ , often parameterized by a neural network, that accepts as input a sequence of known tokens  $s^{(-N_p)}, \dots, s^{(t-1)}$ , which contains a prompt and the first  $t - 1$  tokens already produced by the language model, and then outputs a vector of  $|\mathcal{V}|$  logits, one for each word in the vocabulary. These logits are then passed through a softmax operator to convert them into a discrete probability distribution over the vocabulary. The next token at position  $t$  is then sampled from this distribution using either standard multinomial sampling, or *greedy* sampling (greedy decoding) of the single most likely next token. Additionally, a procedure such as *beam search* can be employed to consider multiple possible sequences before selecting the one with the overall highest score.

### 1.2. A caveat: The difficulty of watermarking low-entropy sequences

Consider the following two sequences of tokens, with prompts in red:

The quick brown fox jumps over the lazy dog

```
for (i=0; i<n; i++) sum+=array[i]
```

Were they produced by a human or by a language model? Determining this is fundamentally hard because these sequences have low entropy; the first few tokens strongly determine the following tokens.

Low entropy text creates two problems for watermarking. First, both humans and machines provide similar if not identical completions for low entropy prompts, making it impossible to discern between them. Second, it is difficult to watermark low entropy text, as any changes to the choice of tokens may result in high perplexity, unexpected tokens that degrade the quality of the text. Later, we rigorously define sentence entropy, and analyze its impact on watermark detection.

## 2. A simple proof of concept

We start out by describing a simple “hard” red list watermark in Algorithm 1 that is easy to analyze, easy to detect and hard to remove. The simplicity of this approach comes at the cost of poor generation quality on low entropy sequences. We will discuss more sophisticated strategies later.

---

### Algorithm 1 Text Generation with Hard Red List

---

**Input:** prompt,  $s^{(-N_p)} \dots s^{(-1)}$

**for**  $t = 0, 1, \dots$  **do**

1. Apply the language model to prior tokens  $s^{(-N_p)} \dots s^{(t-1)}$  to get a probability vector  $p^{(t)}$  over the vocabulary.
2. Compute a hash of token  $s^{(t-1)}$ , and use it to seed a random number generator.
3. Using this seed, randomly partition the vocabulary into a “green list”  $G$  and a “red list”  $R$  of equal size.
4. Sample  $s^{(t)}$  from  $G$ , never generating any token in the red list.

**end for**

---

The method works by generating a pseudo-random red list of tokens that are barred from appearing as  $s^{(t)}$ . The red list generator is seeded with the prior token  $s^{(t-1)}$ , enabling the red list to be reproduced later without access to the entire generated sequence.

**Detecting the watermark.** While producing watermarked text requires access to the language model, detecting the watermark does not. A third party with knowledge of the hash function and random number generator can re-produce the red list for each token and count how many times the red list rule is violated. We can detect the watermark by testing

the following null hypothesis,

$$H_0: \text{The text sequence is generated with no knowledge of the red list rule.} \quad (1)$$

Because the red list is chosen at random, a natural writer is expected to violate the red list rule with half of their tokens, while the watermarked model produces no violations. The probability that a natural source produces  $T$  tokens without violating the red list rule is only  $1/2^T$ , which is vanishingly small even for short text fragments with a dozen words. This enables detection of the watermark (rejection of  $H_0$ ) for, e.g., a synthetic tweet.

A more robust detection approach uses a *one proportion z-test* to evaluate the null hypothesis. If the null hypothesis is true, then the number of green list tokens, denoted  $|s|_G$ , has expected value  $T/2$  and variance  $T/4$ . The  $z$ -statistic for this test is

$$z = 2(|s|_G - T/2)/\sqrt{T}. \quad (2)$$

We reject the null hypothesis and detect the watermark if  $z$  is above a chosen threshold. Suppose we choose to reject the null hypothesis if  $z > 4$ . In this case, the probability of a false positive is  $3 \times 10^{-5}$ , which is the one-sided p-value corresponding to  $z > 4$ . At the same time, we will detect any watermarked sequence with 16 or more tokens (the minimum value of  $T$  that produces  $z = 4$  when  $|s|_G=T$ ).

**How hard is it to remove the watermark?** The use of the one proportion z-test makes removal of the watermark difficult. Consider the case of a watermarked sequence of length  $T = 1000$ . Suppose an adversary modifies 200 tokens in the sequence to add red list words and scrub the watermark. A modified token at position  $t$  can violate the red list rule at position  $t$ . Furthermore, the value of  $s_t$  determines the red list for token  $s_{t+1}$ , and a maximally adversarial choice of  $s_t$  will put  $s_{t+1}$  in violation of the red list rule as well. For this reason, 200 token flips can create at most 400 violations of the red list rule. Unfortunately for the attacker, this maximally adversarial sequence with 600 remaining green list tokens still produces a z-statistic of  $2(600 - 1000/2)/\sqrt{1000} \approx 6.3$ , and a p-value of  $\approx 10^{-10}$ , leaving the watermark readily detectable with extremely high confidence. In general, removing the watermark of a long sequence requires modifying roughly one quarter of the tokens or more.

Note the analysis above assumes the attacker has complete knowledge of the watermark, and each selected token is maximally adversarial (which likely has a negative impact on quality). Without knowledge of the watermark algorithm, each flipped token has only a 50% chance of being in the red list, as does the adjacent token. In this case, the attacker above only creates 200 red list words (in expectation) by

modifying 200 tokens. Methods for keeping the watermark algorithm secret but available via API are discussed in Section 5.

**Drawbacks of the hard red list rule.** The hard red list rule handles low entropy sequences in a simple way; it prevents the language model from producing them. For example, the token ‘‘Barack’’ is almost deterministically followed by ‘‘Obama’’ in many text datasets, yet ‘‘Obama’’ may be disallowed by the red list.

A better behavior is to use a ‘‘soft’’ watermarking rule that is only active for high-entropy text that can be imperceptibly watermarked. As long as low-entropy sequences are wrapped inside a passage with enough total entropy, the passage will still easily trigger a watermark detector, solving the problem described in Section 1.2. Further, one can combine the watermark with a beam search decoder that ‘‘irons-in’’ the watermark. By searching the hypothesis space of likely token sequences, candidates sequences with a high density of green list tokens are found, resulting in a high strength watermark with minimal perplexity cost.

### 3. A more sophisticated watermark

We now discuss the ‘‘soft’’ watermark that promotes the use of the green list for high entropy tokens when many good choices are available, while having little impact on the choice of low-entropy tokens that are nearly deterministic.

To derive this watermark, we examine what happens in the language model just before it produces a probability vector. The last layer of the language model outputs a vector of logits  $l^{(t)}$ . These logits get converted into a probability vector  $p^{(t)}$  using the softmax operator

$$p_k^{(t)} = \exp(l_k^{(t)}) / \sum_i \exp(l_i^{(t)}).$$

Rather than strictly prohibiting the red list tokens, Algorithm 2 adds a constant  $\delta$  to the logits of the green list tokens.

The soft red list rule adaptively enforces the watermark in situations where doing so will have little impact on quality, while almost ignoring the watermark rule in the low entropy case where there is a clear and unique choice of the ‘‘best’’ word. A highly likely word with  $p_k^{(t)} \approx 1$  has a much larger logit than other candidates, and this will remain the largest regardless of whether it is in the red list. But when the entropy is high, there are many comparably large logits to choose from, and the  $\delta$  rule has a large impact on the sampling distribution, strongly biasing the output towards the green list.

---

**Algorithm 2** Text Generation with Soft Red List

---

**Input:** prompt,  $s^{(-N_p)} \dots s^{(-1)}$   
green list size,  $\gamma \in (0, 1)$   
hardness parameter,  $\delta > 0$

**for**  $t = 0, 1, \dots$  **do**

1. Apply the language model to prior tokens  $s^{(-N_p)} \dots s^{(t-1)}$  to get a logit vector  $l^{(t)}$  over the vocabulary.
2. Compute a hash of token  $s^{(t-1)}$ , and use it to seed a random number generator.
3. Using this random number generator, randomly partition the vocabulary into a “green list”  $G$  of size  $\gamma|V|$ , and a “red list”  $R$  of size  $(1 - \gamma)|V|$ .
4. Add  $\delta$  to each green list logit. Apply the softmax operator to these modified logits to get a probability distribution over the vocabulary.
$$\hat{p}_k^{(t)} = \begin{cases} \frac{\exp(l_k^{(t)} + \delta)}{\sum_{i \in R} \exp(l_i^{(t)}) + \sum_{i \in G} \exp(l_i^{(t)} + \delta)}, & k \in G \\ \frac{\exp(l_k^{(t)})}{\sum_{i \in R} \exp(l_i^{(t)}) + \sum_{i \in G} \exp(l_i^{(t)} + \delta)}, & k \in R. \end{cases}$$
5. Sample the next token,  $s^{(t)}$ , using the watermarked distribution  $\hat{p}^{(t)}$ .

**end for**

---

### 3.1. Detecting the soft watermark

The process for detecting the soft watermark is identical to that for the hard watermark. We assume the null hypothesis (1) and compute a z-statistic using Equation (2). We reject the null hypothesis and detect the watermark if  $z$  is greater than a threshold. For arbitrary  $\gamma$  we have

$$z = (|s|_G - \gamma T) / \sqrt{T\gamma(1 - \gamma)}. \quad (3)$$

Consider again the case in which we detect the watermark for  $z > 4$ . Just like in the case of the hard watermark, we get false positives with rate  $3 \times 10^{-5}$ . In the case of the hard watermark, we could detect any watermarked sequence of length 16 tokens or more, regardless of the properties of the text. However, in the case of the soft watermark our ability to detect synthetic text depends on the entropy of the sequence. High entropy sequences are detected with relatively few tokens, while low entropy sequences require more tokens for detection. Below, we rigorously analyze the detection sensitivity of the soft watermark, and its dependence on entropy.

### 4. Analysis of the soft watermark

In this section, we examine the expected number of green list tokens used by a watermarked language model and analyze

the dependence of this quantity on the entropy of a generated text fragment. Our analysis assumes the red list is sampled uniformly at *random*. This is a deviation from the method used in practice, which generates red lists using a *pseudo-random* number generator seeded with previous tokens. The consequences of pseudo-random sampling are explored in Section 5. We analyze the case in which text is generated by multinomial random sampling. In our experiments, we consider two more sampling schemes, greedy decoding and beam search.

We need a definition of entropy that is appropriate for our analysis. The strength of our watermark is weak when the distribution over tokens has a large “spike” concentrated on one or several tokens. We define the following type of entropy to quantify this phenomenon.

**Definition 4.1.** Given a discrete probability vector  $p$  and a scalar  $z$ , we define the *spike entropy* of  $p$  with modulus  $z$  as

$$S(p, z) = \sum_k \frac{p_k}{1 + zp_k}.$$

Like the classical Shannon entropy, the spike entropy is a measure of how spread out a distribution is; The spike entropy assumes its minimal value of  $\frac{1}{1+z}$  when the entire mass of  $p$  is concentrated at a single location, and its maximal value of  $\frac{N}{N+z}$  when the mass of  $p$  is uniformly distributed. For large  $z$ , the value of  $\frac{p_k}{1+zp_k} \approx 1/z$  when  $p_k > 1/z$  and  $\approx 0$  for  $p_k < 1/z$ . For this reason, one can interpret the spike entropy as a softened measure of the number of entries in  $p$  greater than  $1/z$ .

The following theorem predicts the number of green list tokens that appear in a sequence with the watermark.

**Theorem 4.2.** Consider watermarked text sequences of  $T$  tokens. Each sequence is produced by sequentially sampling a raw probability vector  $p^{(t)}$  from the language model, sampling a random green list of size  $\gamma N$ , and boosting the green list logits by  $\delta$  using Equation 4 before sampling each token. Define  $\alpha = \exp(\delta)$ , and let  $|s|_G$  denote the number of green list tokens in sequence  $s$ .

If a randomly generated watermarked sequence has average spike entropy at least  $S^*$ , i.e.,

$$\frac{1}{T} \sum_t S \left( p^{(t)}, \frac{(1 - \gamma)(\alpha - 1)}{1 + (\alpha - 1)\gamma} \right) \geq S^*,$$

then the number of green list tokens in the sequence has expected value at least

$$\mathbb{E} |s|_G \geq \frac{\gamma \alpha T}{1 + (\alpha - 1)\gamma} S^*,$$

Furthermore, the number of green list tokens has variance at most

$$\text{Var} |s|_G \leq T \frac{\gamma \alpha S^*}{1 + (\alpha - 1)\gamma} \left( 1 - \frac{\gamma \alpha S^*}{1 + (\alpha - 1)\gamma} \right).$$

If we have chosen  $\gamma \geq .5$ , then we can use the strictly looser but simpler bound

$$\text{Var } |s|_G \leq T\gamma(1 - \gamma).$$

**Remark.** It may seem like there are a lot of messy constants floating around in this bound. However, when we choose  $\gamma = \frac{1}{2}$  and  $\delta = \ln(2) \approx 0.7$ , this bound simplifies to

$$\mathbb{E} |s|_G \geq \frac{2}{3}TS^*, \quad \text{Var } |s|_G \leq \frac{2}{3}TS^* \left(1 - \frac{2}{3}S^*\right)$$

where  $S^*$  is a bound on spike entropy with modulus  $1/3$ . If we study the ‘‘hard’’ red list rules by choosing  $\gamma = \frac{1}{2}$  and letting  $\delta \rightarrow \infty$ , we have

$$\mathbb{E} |s|_G \geq TS^*, \quad \text{Var } |s|_G \leq TS^*(1 - S^*)$$

where  $S^*$  is a bound on spike entropy with modulus 1.

#### 4.1. Sensitivity of the watermark test

The sensitivity of the soft watermark can be computed using standard type-II error analysis. For illustrative purposes, we estimate the type-II (false negative) error rate of a soft watermark with  $\gamma = .5$  and  $\delta = 2$ . We assume 200 tokens are generated using OPT-1.3B (Zhang et al., 2022) using prompts from the C4 dataset’s RealNewsLike subset (Raffel et al., 2019). We also assume a detection threshold of  $z = 4$  (which occurs at  $\sim 128.2/100$  tokens) which gives us a type-I error (false positive) rate of  $3 \times 10^{-5}$ .

**Theoretical bound.** Our generations have an average spike entropy per sample of  $S = 0.807$  over  $\sim 500$  generations. Theorem 4.2 says that the expected number of green list tokens per generation is *at least* 142.2. Indeed, the empirical average is 159.5. For sequences with entropy equal to the mean ( $S = 0.807$ ) we get  $\sigma \leq 6.41$  tokens, and 98.6% sensitivity (1.4% type-II error rate), using a standard Gaussian approximation for the green list count. Note, this is a *lower* bound on the sensitivity for this particular entropy. If we use the true empirical mean of 159.5 rather than the theoretical bound, we get a  $5.3 \times 10^{-7}$  type-II error rate, a realistic approximation but not a rigorous lower bound.

**Empirical sensitivity.** Empirically, 98.4% of generations are detected at the  $z = 4$  (128 token) threshold when multinomial sampling is used. When 4-way beam search over a greedy decoding is used, we get 99.6% empirical sensitivity. Unlike the theoretical bounds, these are computed over all generations, which have the same length but vary in their individual entropies. Here, the primary source of type-II errors is low entropy sequences, as calculations above show that we expect a very low error rate when the entropy lies near the mean. To validate this, we examine the subset of 375/500 generations that have spike entropy above the 25th

percentile, of which we detect 100% of generations at the  $z = 4$  threshold.

**What do failure cases look like?** We display typical success and failure cases for the watermark in Table 1. We observe that low-entropy (undetectable) sequences typically involve data memorization; the model regurgitates a copy (or near copy) of human-written text which is therefore not detectable as machine-written. A detailed exploration of model accuracy is presented in Section 6 with more generation examples provided in Appendix A.3.

#### 4.2. Impact on quality of generated text

When the distribution produced by the language model is uniform (maximal entropy), the randomness of the green list results in tokens being uniformly sampled, and the perplexity remains untouched. Conversely, in the case of minimal entropy, where all probability mass is concentrated on a single token, the soft watermark rule has no effect and there is once again no impact on perplexity. However, the watermark rule does impact perplexity for tokens of moderate entropy. We provide a bound on this perplexity change in the supplementary material.

### 5. Private Watermarking

The watermark algorithms above are designed to be *public*. A watermark can also be operated in *private mode*, in which the algorithm uses a random key that is kept secret and hosted behind a secure API. If the attacker has no knowledge of the key used to produce the red list, it becomes more difficult for the attacker to remove the watermark as the attacker does not know which tokens are in the red list. However, testing for the presence of the watermark now requires using the same secure API and, if this API is public, access needs to be monitored to prevent an adversary from making too many queries using minor variants of the same sequence.

A range of more complex mechanisms are possible with different efficiency and security tradeoffs, and we discuss a few of these in the supplementary material.

### 6. Experiments

In this section we explore the behavior of the watermark using the OPT-1.3B model (Zhang et al., 2022). We measure watermark strength using the rate of type-I errors (human text falsely flagged as watermarked) and type-II errors (watermarked text not detected). We provide basic implementation details in Appendix A.1.

**Datasets and Prompts.** To simulate a variety of realistic language modeling scenarios we slice and dice a random selection of texts from the news-like subset of the C4 dataset (Raffel et al., 2019). For each random string, we trim a fixed

Table 1. Selected outputs from non-watermarked (NW) and watermarked (W) multinomial sampling using  $\gamma = 0.5$  and  $\delta = 2.0$ . The example in the first row has high entropy and correspondingly high  $z$ -scores, without any perceptible degradation in output quality. The lower row is a failure case where the watermark is too weak to be detected – it has low entropy and corresponding low  $z$ -scores.

prompt	real completion	no watermark (NW)	watermarked (W)	$S$	(W) $z$	(NW) PPL	(W) PPL
...tled out of court and publicly reconciled. \nIn the '80s the band's popularity waned in the United States but remained strong abroad. Robin released three solo albums, with limited success. The Bee Gees	returned with some moderate hits in the late 1990s and were inducted into the Rock and Roll Hall of Fame in 1997. With his brothers, Mr. Gibb won six Grammys. \nIn addition to his wife and his brother [...continues]	continued to tour, and Barry became a television producer. \nBut in the early '90s, the Bee Gees' popularity remained high. They scored a hit with "Don't Stop Believing" in 1990, and in 1992 the Bee Ge[...continues]	' 1990 album, "Spirits of the Century," was a mixed critical and commercial success. \nWhen the brothers were nominated for a Grammy Award in 1990, Mr. Gibb's "You Should Be Dancing" and "Massachusetts,[...continues]	0.68	12.73	3.15	1.93
...cond season at Hall Bros Oval. \nThe defender also admitted his surprise at Young's run to the finals but credited the injection of youth into the side. \nWe were really in a building phase last year and	we copped a few floggings with all those juniors blokes coming in," Galvin said. \nNow, we've kept that core group together for two years and I think we've come along quicker than we anticipated. \nROCK[...continues]	we copped a few floggings with all those juniors blokes coming in," Galvin said. \nNow, we've kept that core group together for two years and I think we've come along quicker than we anticipated. \nTha[...continues]	we copped a few floggings with all those juniors blokes coming in," Galvin said. \nNow, we've kept that core group together for two years and I think we've come along quicker than we anticipated. \nTha[...continues]	0.58	-1.13	1.05	1.04

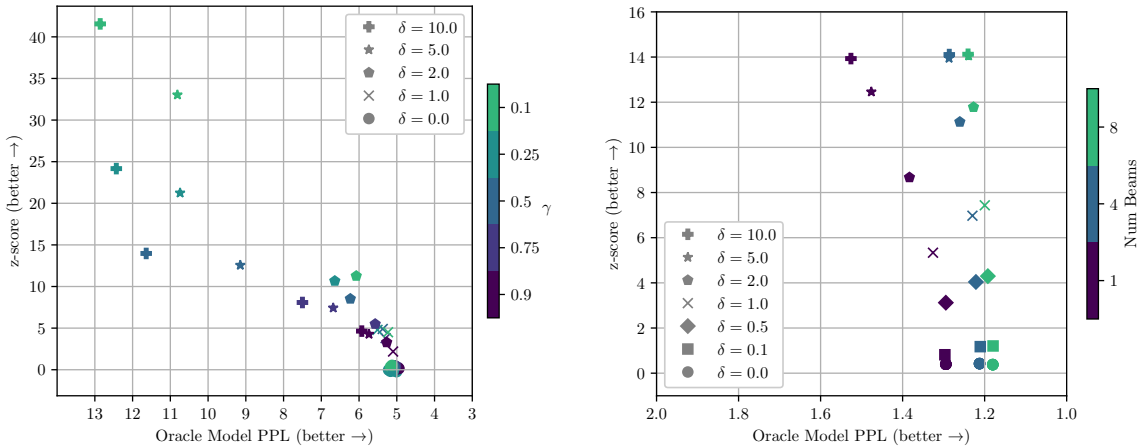


Figure 2. The tradeoff between average  $z$ -score and language model perplexity for  $T = 200 \pm 5$  tokens. (left) Multinomial sampling. (right) Greedy and beam search with 4 and 8 beams for  $\gamma = .5$ . Beam search promotes higher green list usage and thus larger  $z$ -scores with smaller impact to model quality (perplexity, PPL).

length of tokens from the end and treat them as a “baseline” completion. The remaining tokens are a prompt. A larger oracle language model (OPT-2.7B) is used to compute perplexity (PPL) for the generated completions and for the human baseline.

**Watermark Strength vs Text Quality.** One can achieve a very strong watermark for short sequences by choosing a small green list size  $\gamma$  and a large green list bias  $\delta$ . However, creating a stronger watermark may distort generated text. Figure 2 (left) shows the tradeoff between watermark strength ( $z$ -score) and text quality (perplexity) for various combinations of watermarking parameters. We compute results using  $500 \pm 10$  sequences of length  $T = 200 \pm 5$  tokens for each parameter choice. Interestingly, we see that a small green list,  $\gamma = .1$  is pareto-optimal.

In addition to these quantitative results, we show examples of real prompts and watermarked outputs in Table 1 to provide a qualitative sense for the behavior of the test statistic and quality measurement on different kinds of prompts. Additional examples are compiled in Appendix A.3.

**Ironing in the Watermark with Beam Search.** Figure 2 (right) shows the tradeoff between watermark strength and accuracy when beam search is used. Beam search has a synergistic interaction with the soft watermarking rule. Particularly when 8 beams are used, the points in Figure 2 form an almost vertical line, showing very little perplexity cost to achieve strong watermarking.

**Watermark Strength vs Number of Tokens.** Theory predicts that the type I and type II error rates of the watermark should decay to zero as the sequence length  $T$  increases. Figure 3 shows the strength of the watermark, measured using the average  $z$ -score over samples, as  $T$  sweeps from 2 to 200. Curves are shown for various values of  $\delta$  and  $\gamma$ . The left two charts use multinomial sampling, while the right chart uses 8-way beam search and  $\gamma = .25$ . Once again, we see the power of the beam search in achieving high green list ratios; even for the moderate bias of  $\delta = 2$ , an average  $z$ -score greater than 5 is achieved for as few as 35 tokens.

**Performance and Sensitivity for Multinomial Sampling.** To show the sensitivity of the resulting hypothesis test based

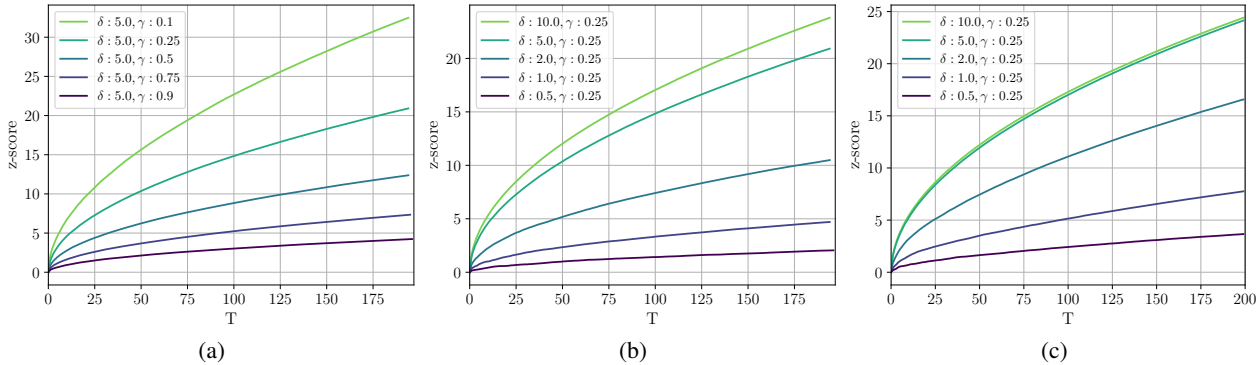


Figure 3. The average  $z$ -score as a function of  $T$  the token length of the generated text. (a) The dependence of the  $z$ -score on the green list size parameter  $\gamma$ , under multinomial sampling. (b) The effect of  $\delta$  on  $z$ -score, under multinomial sampling. (c) The impact of the green list size parameter  $\gamma$  on the  $z$ -score, but with greedy decoding using 8-way beam search.

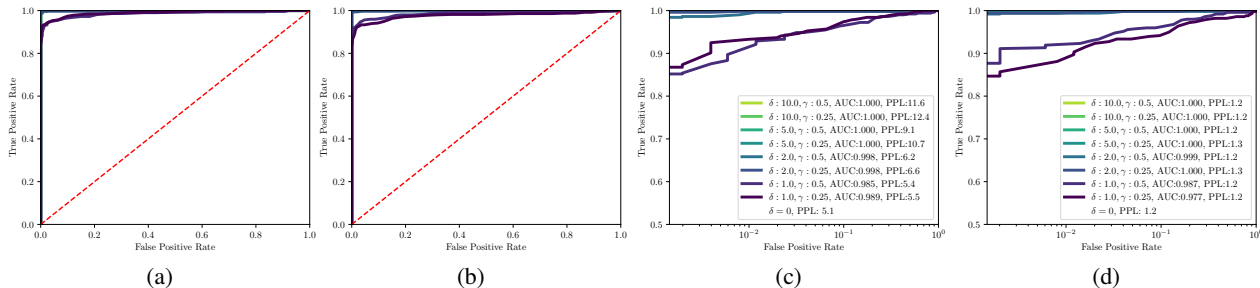


Figure 4. ROC curves with AUC values for watermark detection. Several choices of watermark parameter  $\delta$  are shown for (a) multinomial sampling and (b) greedy decoding with 8-way beam search. (c,d) The same charts with semilog axes. Higher  $\delta$  values achieve stronger performance, but additionally we see that for a given  $\delta$ , the beam search allows the watermark to capture slightly more AUC than the corresponding parameters under the multinomial sampling scheme.

on the observed  $z$ -scores, we provide a table of error rate for various watermarking parameters in Table 2. We also sweep a range of thresholds in ROC charts in Figure 4. We further report detection performance and error rates for various cutoffs in Appendix C, and provide a comparison between empirical  $z$ -scores and theoretical predictions. Note that no type-I (false positive) errors were observed for any run shown in the error tables (see the columns of 0.0’s).

### 7. Attacking the watermark

Like any software tool, care must be taken when implementing a watermark and watermark detector so that security is maintained. Otherwise, an adversarial user may modify text to add red list tokens, and thus avoid detection. In many cases, simple attacks can be avoided by properly normalizing text before hashes are computed. In Appendix B.1, we discuss a range of attacks that we are currently aware of, and methods to mitigate them. In the next section, we implement and evaluate an example of a representative attack using a second, smaller language model.

### 7.1. Degradation Under Attack: Span Replacement Using a LM

We study a realistic black-box attack by attempting to remove the presence of the watermark by replacing spans in the original output text using another language model. We treat the watermark algorithm as if it is private, mocking seclusion behind an API. The attacker does not have access to the locations of green list tokens and instead tries to modify the text through token replacement at random indices until a certain word replacement *budget*,  $\epsilon$ , is reached. The budget constraint maintains a level semantic similarity between the original watermarked text and the attacked text, otherwise the “utility” of the original text for its intended task may be lost. Also, each span replacement in the attack is performed via inference using a multi-million parameter language model. While this is roughly a third the size of the target model, it means that the attack incurs an associated cost per step implying that a base level of efficiency with respect to model calls would be desired in practice.

In our experiment, we adopt T5-Large (Raffel et al., 2020) as the replacement model and iteratively select and replace tokens until the attacker either reaches the budget, or no

Table 2. Empirical error rates for watermark detection using multinomial sampling and beam search. Each row is averaged over  $\sim 500(-4/+7)$  generated sequences of length  $T = 200 \pm 5$ . A maximum of one type-I (false positive) error was observed for any given run. All soft watermarks at  $\delta = 2.0$  incur at most 1.6% (8/500) type-II error at  $z = 4$ . No type-II errors occurred for the hardest watermarks with  $\delta = 5.0$  and  $\gamma = 0.25$ . A more extensive set of error rates is provided in Table 8.

		z=4.0				
sampling	$\delta$	$\gamma$	FPR	TNR	TPR	FNR
m-nom.	1.0	0.50	0.0	1.0	0.767	0.233
m-nom.	1.0	0.25	0.0	1.0	0.729	0.271
m-nom.	2.0	0.50	0.0	1.0	0.984	0.016
m-nom.	2.0	0.25	0.0	1.0	0.994	0.006
m-nom.	5.0	0.50	0.0	1.0	0.996	0.004
m-nom.	5.0	0.25	0.0	1.0	1.000	0.000
8-beams	1.0	0.25	0.0	1.0	0.819	0.181
8-beams	2.0	0.25	0.0	1.0	0.994	0.006
8-beams	5.0	0.25	0.0	1.0	1.000	0.000

more suitable replacement candidates are returned.

**Details of the T5 Span Attack.** We tokenize the watermarked text using the T5 tokenizer. Then, while fewer than  $\varepsilon T$  successful replacements have been performed or a maximal iteration count is reached:

1. Randomly replace one word from the tokenization with a `<mask>`.
2. Pass the region of text surrounding the `mask` token to T5 to obtain a list of  $k = 20$  candidate replacement token sequences via a 50-way beam search, with associated scores corresponding to their likelihood.
3. Each candidate is decoded into a string. If one of the  $k$  candidates returned by the model is *not* equal to the original string corresponding to the masked span, then the attack succeeds, and the span is replaced with the new text.

After attacking a set of 500 sequences of length  $T = 200 \pm 5$  token sequences this way, we compute updated  $z$ -scores and tabulate error rates (Table 9 in the Appendix). We also generate ROC plots for a range of  $\varepsilon$  budgets. While this attack is effective at increasing the number of red list tokens in the text, as shown in Figure 5, we only measure a decrease in watermark strength of 0.01 AUC when  $\varepsilon = 0.1$ . While the watermark removal is more successful at a larger budget of 0.3, the average PPL of attacked sequences increases by  $3\times$  in addition to requiring more model calls.

## 8. Related Work

The idea of watermarking has a long history. Watermarking of digital text has been considered challenging, due to its discrete nature (Katzenbeisser & Petitcolas, 2000). Watermark-

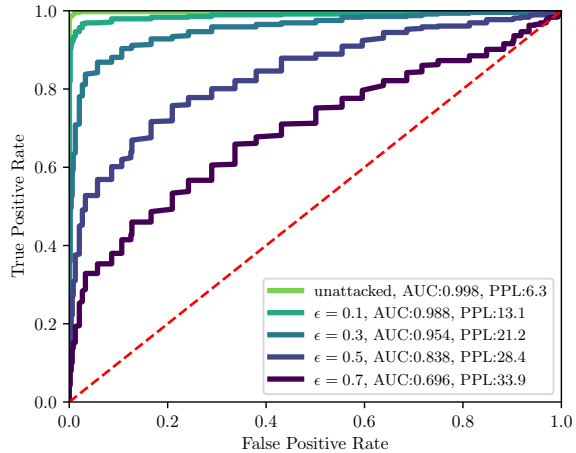


Figure 5. ROC curves for watermark detection under attack via the T5 Span Attack detailed in Section 7.1, with various replacement budgets  $\varepsilon$ . The initial, unattacked watermark is a  $\gamma = 0.5$ ,  $\delta = 2.0$  soft watermark generated using multinomial sampling. The attack achieves a high level of watermark degradation, but *only* at  $\varepsilon = 0.3$ , which *costs* the attacker an average of  $\sim 15$  points of perplexity compared the PPL of the original watermarked text.

ing is considered easier for continuous data, where a variety of effective strategies are known (Petitcolas et al., 1999; Zhu et al., 2018; Lu et al., 2021; Boenisch, 2021). Note that *watermarking* is a special case of *steganography*, the task of embedding arbitrary hidden information into data.

**Watermarking Natural Language.** Early approaches to watermarking natural text in digital form (Atallah et al., 2001; 2003; Venugopal et al., 2011) pose similar desiderata as in Section 1, except for classical models.

Rule-based watermarks were further developed in a series of works (Chiang et al., 2004; Topkara et al., 2006a;b; Meral et al., 2009) with variants also embedding watermarks based on synonym tables instead of only parse trees. Early developments were summarized in Jalil & Mirza (2009), but strong watermarks significantly degraded the text quality due to the limited flexibility of language models of the time.

While approaches via hierarchical language models in Wilson et al. (2014) still required human interactions, the emergence of modern neural language models (Vaswani et al., 2017; Devlin et al., 2019) also allowed for improved watermarking/steganography (Fang et al., 2017; Ziegler et al., 2019; Dai & Cai, 2019; He et al., 2022). Fang et al. (2017) propose such a neural steganography approach where, to encode a message of  $w$  bits, the message is first separated into blocks of length  $b$ . Then, the vocabulary  $V$  of a language model is partitioned at random into disjoint sets of size  $|V|/2^b$ . A generative LM can then encode the message



by generating only a token from the “allowed” set at each position. However, this hard rule reduces the quality of generated text, boxing the LM into a small selection of tokens for each step.

In a somewhat similar approach to our proposal, [Kaptchuk et al. \(2021\)](#) constructs a framework that adapts traditional public-key cryptographic steganography specifically for “natural” communication channels like text using generative models. However, their method, Meteor, relies on a synchronized model framework where the sender and receiver agree on a shared generative model used to embed and decode the hidden bits of information being sent.

Another approach ([Ueoka et al., 2021](#)) uses mask-infilling via BERT to edit already-generated text for steganography. Finally, [Abdelnabi & Fritz \(2021\)](#) design an end-to-end system where both encoding and decoding are handled by text-to-text language models that are trained adversarially for this purpose, and [Gu et al. \(2022\)](#) implants a series of backdoor triggers through a finetuning process to cause biased responses to specific inputs, a behavior detectable at verification time.

Recently, [Aaronson \(2022\)](#) announced that he is also studying cryptographic approaches to watermarking in collaboration with OpenAI. Their preliminary method is based on biasing of the LM output, as opposed to complete determination as in [Fang et al. \(2017\)](#). While details are not currently available, the description suggests that hashing of  $n$ -gram sequences is involved. We hope to extend our comparison to this work when more information becomes available.

In contrast to other currently published works, we want to focus on strategies that are simultaneously minimally restrictive to a language model, leverage the LMs own understanding of natural text, require no usage of the LM to decode the watermark, and can be theoretically analyzed and validated.

**Post-hoc Detection.** An alternative to watermarking is post-hoc analysis of machine-generated text, for example using language model features or finetuning existing large language models to behave as detectors ([Zellers et al., 2019](#); [Tan et al., 2020](#)), see an overview in [Jawahar et al. \(2020\)](#). Implementation details, such as sampling strategies, can be reverse-engineered from text ([Tay et al., 2020](#)). However, detection approaches are slowly losing ground as LM capabilities increase; [Gambini et al. \(2022\)](#) note that detection strategies for GPT-2 struggle with GPT-3. Further, known detectors are vulnerable to adversarial attacks that degrade their functionality ([Wolff & Wolff, 2022](#)).

While efforts to provide strong detectors continue, as in [Tian \(2023a\)](#), ultimately language model progress may make detection infeasible. All post-hoc detection methods require the LM to be significantly biased away from human text in

some measurable way, such as low variation in perplexity across sentences ([Tian, 2023b](#)). Even for current LLMs, this margin might be small. Detection schemes that operate within this small margin are susceptible to labeling human text as false-positive, a concern that is especially pressing for non-native speakers and people who use computer writing assistance tools. Such populations might be especially at risk for false-positives, which could lead to academic problems if these detectors are used in schools ([Butoi, 2023](#)).

## 9. Conclusion

There is one more important property of the proposed method: The  $z$ -statistic used for detection depends only on the green list size parameter  $\gamma$  and the hash function. There is no dependence on  $\delta$  or any other factor related to how the green list is enforced. For this reason, one can deploy the watermark using context-specific  $\delta$  or green list enforcement rules for different kinds of text (e.g., prose vs code) while using the same downstream detector. One can also change a proprietary implementation of the watermarked sampling algorithm without any need to change the detector. Finally, the watermarking method could be turned on only in certain contexts, for example when a specific user seems to exhibit suspicious behavior.

There are still a number of open questions regarding watermarking. For example, what is the best way to test for the watermark in a streaming context, or in a context where a short span of watermarked text lives inside a longer non-watermarked span? We hope our present results are enough to convince readers that watermarks could be a practical tool for combating malicious uses of generative models, and we leave these additional questions for future research.

## 10. Acknowledgements

This work was made possible by the ONR MURI program, DARPA GARD (HR00112020007), the Office of Naval Research (N000142112557), and the AFOSR MURI program. Commercial support was provided by Capital One Bank, the Amazon Research Award program, and Open Philanthropy. Further support was provided by the National Science Foundation (IIS-2212182), and by the NSF TRAILS Institute (2229885).

## References

- Aaronson, S. My AI Safety Lecture for UT Effective Altruism, November 2022. URL <https://scottaaronson.blog/?p=6823>.
- Abdelnabi, S. and Fritz, M. Adversarial Watermarking Transformer: Towards Tracing Text Provenance

- with Data Hiding. In *2021 IEEE Symposium on Security and Privacy (SP)*, pp. 121–140, May 2021. doi: 10.1109/SP40001.2021.00083.
- Atallah, M. J., Raskin, V., Crogan, M., Hempelmann, C., Kerschbaum, F., Mohamed, D., and Naik, S. Natural Language Watermarking: Design, Analysis, and a Proof-of-Concept Implementation. In Moskowitz, I. S. (ed.), *Information Hiding*, Lecture Notes in Computer Science, pp. 185–200, Berlin, Heidelberg, 2001. Springer. ISBN 978-3-540-45496-0. doi: 10.1007/3-540-45496-9\_14.
- Atallah, M. J., Raskin, V., Hempelmann, C. F., Karahan, M., Sion, R., Topkara, U., and Triezenberg, K. E. Natural Language Watermarking and Tamperproofing. In Petitcolas, F. A. P. (ed.), *Information Hiding*, Lecture Notes in Computer Science, pp. 196–212, Berlin, Heidelberg, 2003. Springer. ISBN 978-3-540-36415-3. doi: 10.1007/3-540-36415-3\_13.
- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., Chen, C., Olsson, C., Olah, C., Hernandez, D., Drain, D., Ganguli, D., Li, D., Tran-Johnson, E., Perez, E., Kerr, J., Mueller, J., Ladish, J., Landau, J., Ndousse, K., Lukosiute, K., Lovitt, L., Sellitto, M., Elhage, N., Schiefer, N., Mercado, N., DasSarma, N., Lasenby, R., Larson, R., Ringer, S., Johnston, S., Kravec, S., Showk, S. E., Fort, S., Lanham, T., Telleen-Lawton, T., Conerly, T., Henighan, T., Hume, T., Bowman, S. R., Hatfield-Dodds, Z., Mann, B., Amodei, D., Joseph, N., McCandlish, S., Brown, T., and Kaplan, J. Constitutional AI: Harmlessness from AI Feedback, December 2022. URL <https://www.anthropic.com/constitutional.pdf>.
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21*, pp. 610–623, New York, NY, USA, March 2021. Association for Computing Machinery. ISBN 978-1-4503-8309-7. doi: 10.1145/3442188.3445922. URL <https://doi.org/10.1145/3442188.3445922>.
- Bergman, A. S., Abercrombie, G., Spruit, S., Hovy, D., Dinan, E., Boureau, Y.-L., and Rieser, V. Guiding the Release of Safer E2E Conversational AI through Value Sensitive Design. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pp. 39–52, Edinburgh, UK, September 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.sigdial-1.4>.
- Boenisch, F. A Systematic Review on Model Watermarking for Neural Networks. *Frontiers in Big Data*, 4, 2021. ISSN 2624-909X. doi: 10.3389/fdata.2021.729663. URL <https://www.frontiersin.org/articles/10.3389/fdata.2021.729663>.
- Boucher, N., Shumailov, I., Anderson, R., and Papernot, N. Bad Characters: Imperceptible NLP Attacks. In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1987–2004, May 2022. doi: 10.1109/SP46214.2022.9833641.
- Butoi, V. Things like GPTZero are scary. I’m sure the creator didn’t have the intention, but the fact that it’s marketed as ”a solution to detecting AI written responses” even though there’s no evidence to show it works consistently and is nevertheless being EMPLOYED by schools, is crazy., January 2023. URL <https://twitter.com/VictorButoi/status/1614779846210752512>.
- Chiang, Y.-L., Chang, L.-P., Hsieh, W.-T., and Chen, W.-C. Natural Language Watermarking Using Semantic Substitution for Chinese Text. In Kalker, T., Cox, I., and Ro, Y. M. (eds.), *Digital Watermarking*, Lecture Notes in Computer Science, pp. 129–140, Berlin, Heidelberg, 2004. Springer. ISBN 978-3-540-24624-4. doi: 10.1007/978-3-540-24624-4\_10.
- Christiano, P., Leike, J., Brown, T. B., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. *arXiv:1706.03741[cs, stat]*, July 2017. doi: 10.48550/arXiv.1706.03741. URL <http://arxiv.org/abs/1706.03741>.
- Crothers, E., Japkowicz, N., and Viktor, H. Machine Generated Text: A Comprehensive Survey of Threat Models and Detection Methods. *arXiv:2210.07321[cs]*, November 2022. doi: 10.48550/arXiv.2210.07321. URL <http://arxiv.org/abs/2210.07321>.
- Dai, F. and Cai, Z. Towards Near-imperceptible Steganographic Text. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4303–4308, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1422. URL <https://aclanthology.org/P19-1422>.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*, May 2019. URL <http://arxiv.org/abs/1810.04805>.
- Fang, T., Jaggi, M., and Argyraki, K. Generating Steganographic Text with LSTMs. In *Proceedings of ACL 2017, Student Research Workshop*, pp. 100–106, Vancouver, Canada, July 2017. Association for Computational Linguistics. URL <https://aclanthology.org/P17-3017>.

- Foltýnek, T., Meuschke, N., and Gipp, B. Academic Plagiarism Detection: A Systematic Literature Review. *ACM Computing Surveys*, 52(6):112:1–112:42, October 2019. ISSN 0360-0300. doi: 10.1145/3345317. URL <https://doi.org/10.1145/3345317>.
- Gabrilovich, E. and Gontmakher, A. The homograph attack. *Communications of the ACM*, 45(2):128, February 2002. ISSN 0001-0782, 1557-7317. doi: 10.1145/503124.503156. URL <https://dl.acm.org/doi/10.1145/503124.503156>.
- Gambini, M., Fagni, T., Falchi, F., and Tesconi, M. On pushing DeepFake Tweet Detection capabilities to the limits. In *14th ACM Web Science Conference 2022, WebSci '22*, pp. 154–163, New York, NY, USA, June 2022. Association for Computing Machinery. ISBN 978-1-4503-9191-7. doi: 10.1145/3501247.3531560. URL <https://doi.org/10.1145/3501247.3531560>.
- Goodside, R. There are adversarial attacks for that proposal as well — in particular, generating with emojis after words and then removing them before submitting defeats it., January 2023. URL <https://twitter.com/goodside/status/1610682909647671306>.
- Grinbaum, A. and Adomaitis, L. The Ethical Need for Watermarks in Machine-Generated Language. *arxiv:2209.03118[cs]*, September 2022. doi: 10.48550/arXiv.2209.03118. URL <http://arxiv.org/abs/2209.03118>.
- Gu, C., Huang, C., Zheng, X., Chang, K.-W., and Hsieh, C.-J. Watermarking pre-trained language models with backdooring. *arXiv preprint arXiv:2210.07543*, 2022.
- He, X., Xu, Q., Zeng, Y., Lyu, L., Wu, F., Li, J., and Jia, R. CATER: Intellectual Property Protection on Text Generation APIs via Conditional Watermarks. In *Advances in Neural Information Processing Systems*, October 2022. URL <https://openreview.net/forum?id=L7P3IvsoUXY>.
- Helfrich, J. N. and Neff, R. Dual canonicalization: An answer to the homograph attack. In *2012 eCrime Researchers Summit*, pp. 1–10, October 2012. doi: 10.1109/eCrime.2012.6489517.
- Jalil, Z. and Mirza, A. M. A Review of Digital Watermarking Techniques for Text Documents. In *2009 International Conference on Information and Multimedia Technology*, pp. 230–234, December 2009. doi: 10.1109/ICIMT.2009.11.
- Jawahar, G., Abdul-Mageed, M., and Lakshmanan, V.S., L. Automatic Detection of Machine Generated Text: A Critical Survey. In *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 2296–2309, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.208. URL <https://aclanthology.org/2020.coling-main.208>.
- Joshi, M., Choi, E., Weld, D. S., and Zettlemoyer, L. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, Vancouver, Canada, July 2017. Association for Computational Linguistics.
- Kaptchuk, G., Jois, T. M., Green, M., and Rubin, A. D. Meteor: Cryptographically secure steganography for realistic distributions. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1529–1548, 2021.
- Katzenbeisser, S. and Petitcolas, F. A. *Information Hiding Techniques for Steganography and Digital Watermarking*. Artech House, Inc., USA, 1st edition, 2000. ISBN 978-1-58053-035-4.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv:1907.11692 [cs]*, July 2019. URL <http://arxiv.org/abs/1907.11692>.
- Lu, S.-P., Wang, R., Zhong, T., and Rosin, P. L. Large-Capacity Image Steganography Based on Invertible Neural Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10816–10825, 2021. URL [https://openaccess.thecvf.com/content/CVPR2021/html/Lu\\_Large-Capacity\\_Image\\_Steganography\\_Based\\_on\\_Invertible\\_Neural\\_Networks\\_CVPR\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Lu_Large-Capacity_Image_Steganography_Based_on_Invertible_Neural_Networks_CVPR_2021_paper.html).
- Meral, H. M., Sankur, B., Sumru Özsoy, A., Güngör, T., and Sevinç, E. Natural language watermarking via morphosyntactic alterations. *Computer Speech & Language*, 23(1):107–125, January 2009. ISSN 0885-2308. doi: 10.1016/j.csl.2008.04.001. URL <https://www.sciencedirect.com/science/article/pii/S0885230808000284>.
- Mirsky, Y., Demontis, A., Kotak, J., Shankar, R., Gelei, D., Yang, L., Zhang, X., Pintor, M., Lee, W., Elovici, Y., and Biggio, B. The Threat of Offensive AI to Organizations. *Computers & Security*, 124:103006, January 2023. ISSN 0167-4048. doi: 10.1016/j.cose.2022.103006. URL <https://www.sciencedirect.com/science/article/pii/S0167404822003984>.

- Muennighoff, N., Wang, T., Sutawika, L., Roberts, A., Biderman, S., Scao, T. L., Bari, M. S., Shen, S., Yong, Z.-X., Schoelkopf, H., et al. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*, 2022.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., and Lowe, R. Training language models to follow instructions with human feedback. *arxiv:2203.02155[cs]*, March 2022. doi: 10.48550/arXiv.2203.02155. URL <http://arxiv.org/abs/2203.02155>.
- Pajola, L. and Conti, M. Fall of Giants: How popular text-based MLaaS fall against a simple evasion attack. In *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 198–211, September 2021. doi: 10.1109/EuroSP51992.2021.00023.
- Petitcolas, F., Anderson, R., and Kuhn, M. Information hiding—a survey. *Proceedings of the IEEE*, 87(7):1062–1078, July 1999. ISSN 1558-2256. doi: 10.1109/5.771065.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language Models are Unsupervised Multi-task Learners. *OpenAI*, pp. 24, 2019.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., and Sutskever, I. Robust Speech Recognition via Large-Scale Weak Supervision. *arxiv:2212.04356[cs, eess]*, December 2022. doi: 10.48550/arXiv.2212.04356. URL <http://arxiv.org/abs/2212.04356>.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*, 2019.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *arXiv:1910.10683 [cs, stat]*, July 2020. URL <http://arxiv.org/abs/1910.10683>.
- Schulman, J., Zoph, B., Kim, C., Hilton, J., Menick, J., Weng, J., Uribe, J. F. C., Fedus, L., Metz, L., Pokorny, M., Gontijo-Lopes, R., Zhao, S., Vijayvergiya, A., Sigler, E., Perelman, A., Voss, C., Heaton, M., Parish, J., Cummings, D., Nayak, R., Balcom, V., Schnurr, D., Kaftan, T., Hallacy, C., Turley, N., Deutsch, N., Goel, V., Ward, J., Konstantinidis, A., Zaremba, W., Ouyang, L., Bognodonoff, L., Gross, J., Medina, D., Yoo, S., Lee, T., Lowe, R., Mossing, D., Huizinga, J., Jiang, R., Wainwright, C., Almeida, D., Lin, S., Zhang, M., Xiao, K., Slama, K., Bills, S., Gray, A., Leike, J., Pachocki, J., Tillet, P., Jain, S., Brockman, G., and Ryder, N. ChatGPT: Optimizing Language Models for Dialogue, November 2022. URL <https://openai.com/blog/chatgpt/>.
- Tan, R., Plummer, B., and Saenko, K. Detecting Cross-Modal Inconsistency to Defend Against Neural Fake News. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2081–2106, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.163. URL <https://aclanthology.org/2020.emnlp-main.163>.
- Tay, Y., Bahri, D., Zheng, C., Brunk, C., Metzler, D., and Tomkins, A. Reverse Engineering Configurations of Neural Text Generation Models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 275–279, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.25. URL <https://aclanthology.org/2020.acl-main.25>.
- Tay, Y., Dehghani, M., Tran, V. Q., Garcia, X., Bahri, D., Schuster, T., Zheng, H. S., Houlby, N., and Metzler, D. Unifying language learning paradigms. *arXiv preprint arXiv:2205.05131*, 2022.
- Tian, E. Gptzero update v1, January 2023a. URL <https://gptzero.substack.com/p/gptzero-update-v1>.
- Tian, E. New year, new features, new model, January 2023b. URL <https://gptzero.substack.com/p/new-year-new-features-new-model>.
- Topkara, M., Riccardi, G., Hakkani-Tür, D., and Atallah, M. J. Natural language watermarking: Challenges in building a practical system. In *Security, Steganography, and Watermarking of Multimedia Contents VIII*, volume 6072, pp. 106–117. SPIE, February 2006a. doi: 10.1117/12.643560. URL <https://www.spiedigitallibrary.org/conference-proceedings-of-spie/6072/60720A/Natural-language-watermarking-challenges-in-building-a-practical-system/10.1117/12.643560.full>.
- Topkara, U., Topkara, M., and Atallah, M. J. The hiding virtues of ambiguity: Quantifiably resilient watermarking of natural language text through synonym substitutions. In *Proceedings of the 8th Workshop on Multimedia and Security, MM&Sec ’06*, pp. 164–174, New York, NY, USA, September 2006b. Association for Computing Machinery. ISBN 978-1-59593-493-2. doi: 10.1145/1161366.1161397. URL <https://doi.org/10.1145/1161366.1161397>.

- Ueoka, H., Murawaki, Y., and Kurohashi, S. Frustratingly Easy Edit-based Linguistic Steganography with a Masked Language Model. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5486–5492, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.433. URL <https://aclanthology.org/2021.naacl-main.433>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention Is All You Need. *arXiv:1706.03762 [cs]*, December 2017. URL <http://arxiv.org/abs/1706.03762>.
- Venugopal, A., Uszkoreit, J., Talbot, D., Och, F., and Ganitkevitch, J. 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, Edinburgh, Scotland, UK., July 2011. Association for Computational Linguistics. URL <https://aclanthology.org/D11-1126>.
- Wilson, A., Blunsom, P., and Ker, A. D. Linguistic steganography on Twitter: Hierarchical language modeling with manual interaction. In *Media Watermarking, Security, and Forensics 2014*, volume 9028, pp. 9–25. SPIE, February 2014. doi: 10.1117/12.2039213. URL <https://www.spiedigitallibrary.org/conference-proceedings-of-spie/9028/902803/Linguistic-steganography-on-Twitter--hierarchical-language-modeling-with-manual/10.1117/12.2039213.full>.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T. L., Gugger, S., Drame, M., Lhoest, Q., and Rush, A. M. HuggingFace’s Transformers: State-of-the-art Natural Language Processing. *arXiv:1910.03771 [cs]*, July 2020. URL <http://arxiv.org/abs/1910.03771>.
- Wolff, M. and Wolff, S. Attacking Neural Text Detectors. *arxiv:2002.11768[cs]*, January 2022. doi: 10.48550/arXiv.2002.11768. URL <http://arxiv.org/abs/2002.11768>.
- Zellers, R., Holtzman, A., Rashkin, H., Bisk, Y., Farhadi, A., Roesner, F., and Choi, Y. Defending Against Neural Fake News. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/3e9f0fc9b2f89e043bc6233994dfcf76-Abstract.html>.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., Mihaylov, T., Ott, M., Shleifer, S., Shuster, K., Simig, D., Koura, P. S., Sridhar, A., Wang, T., and Zettlemoyer, L. Opt: Open pre-trained transformer language models, 2022.
- Zhu, J., Kaplan, R., Johnson, J., and Fei-Fei, L. HiDDeN: Hiding Data with Deep Networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 657–672, 2018. URL [https://openaccess.thecvf.com/content\\_ECCV\\_2018/html/Jiren\\_Zhu\\_HiDDeN\\_Hiding\\_Data\\_ECCV\\_2018\\_paper.html](https://openaccess.thecvf.com/content_ECCV_2018/html/Jiren_Zhu_HiDDeN_Hiding_Data_ECCV_2018_paper.html).
- Ziegler, Z., Deng, Y., and Rush, A. Neural Linguistic Steganography. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 1210–1215, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1115. URL <https://aclanthology.org/D19-1115>.

Table 3. High spike entropy examples under a  $\delta = 2.0, \gamma = 0.5$  watermark with multinomial sampling.

idx	prompt	real completion	no watermark (NW)	watermarked (W)	$S$	(NW) $z$	(W) $z$	(NW) PPL	(W) PPL
1011	[...]nned any special treatment as an officer. The living quarters reserved for him – palatial by camp standards – remained empty as he bunked nearer to the workers. Straight-talking and cigar-chomping (	a habit that earned him the nickname "Big Smoke"), he and politics didn't mix well. Signing on in 1905 "until he could predict success or failure according to his own judgment." Stevens resigned in 19[...]truncated]	he was nicknamed "The Devil", Stevens was a man who could get anything done. He had one major flaw: He was an expert at getting things done – without getting anything done. The problem was, he got [...]truncated]	another Stevens habit), Stevens impressed them so much that he even earned the nickname "Bully Stevens" from a group of workers. As the tedious, years-long construction dragged on – years and years a[...]truncated]	0.87	0.14	8.20	8.77	10.33
924	[...] Helen Cavallo.... She's so funny and so warm, and she's the kind of person the kids won't look away from, and the moms won't look away from." Television personality and author of "Courtney's Creative	Adventurers," Courtney Watkins, hosts "Thinking Up" – which provides imaginative children's activities, drawing from her eclectic background as a teacher of kindergarten and art and a designer of je[...]truncated]	Child," Courtney Love, is one of the series' executive producers. Gilmore said the Disney team has a "talent for building relationships" with the target audience, who, she said, are "very passionate [...]truncated]	World: The Parental Roadmap for Success," Gilmore also noted that there's an opportunity to harness the power of social media to share stories and provide expert advice. "We've seen that, too," she s[...]truncated]	0.87	-0.28	9.62	6.95	10.30
1151	[...]Carrie and a formidable Heather Orth as her frighteningly religious mother. It's not Oeschger's fault that there's a big disconnect between the painfully shy outsider she plays so convincingly and the	forceful notes she sings with dynamic conviction. Cohen's book works best for those who know the story well - no small number - jumping from one plot point to the next without connecting the dots. Ca[...]truncated]	confident, vampish Carrie whose songwriting is in a state of constant revising. It's a shame, because Oeschger is an able actress. And a terrific singer. The original production, which had just two a[...]truncated]	awful person inside that she fights so hard to control with forced bravado. But the script's missteps and musical shortcomings, especially in the tightly-crafted scenes, leave the main characters far[...]truncated]	0.87	-0.99	8.77	8.40	12.33
853	[...] Transferring misbehaving government servants will not solve any problem. They must be suspended and demoted. This is to ensure they will think twice before doing anything. Anonymous #18452573: Bigot	ry is the by-product of racism. This is just another form of ethnocentrism and xenophobia due to BN's 50 years of racially divided political parties and racist policies. Islam, Christianity and other [...]truncated]	ry in schools is rampant and this is not acceptable. Muslims should not be allowed to enter the schools. Tun: I am saddened and disappointed at how our education system is being used to promote discri[...]truncated]	ry against Orang Asli children is a new low. They're not children to be bullied or manipulated by the teachers. How come the teachers are the only people doing such acts? Why not the parents? read of [...]truncated]	0.87	0.99	10.47	7.42	10.55
341	[...]When more responsible people carry, violent crime goes down. The FBI crime statistics don't lie. Why do the left and gun grabbing crowd ignore this fact? Liberalism is a mental disorder. Lets look one	poll that this particular article is using for it's numbers. The poll used 2 Senate Districts in the State of VA. District 21 and District 38. with the total numbers of respondents being 305 voters if [...]truncated]	more time at the facts. The "Fix Gun Checks Act" would make it tougher for the mentally ill to legally buy guns by strengthening background checks so people like James Holmes (the Aurora shooter) can[...]truncated]	day at the number of school shootings that have occurred since Columbine. They were not just random shootings. They were planned. In the preplanned shootings, the students/victims were well trained a[...]truncated]	0.88	-0.42	9.22	1.12	13.66

## A. Experimental Details

### A.1. Implementing the Watermark

We implement the proposed watermark using the Pytorch backend of the Huggingface library (Wolf et al., 2020). The generate API provides useful abstractions, including modules for warping the logit distribution that comes out of the language model. We generate red lists using the torch random number generator and one previous token as described in Section 3. Our pending source code release will give the most concrete definition of the implementation that we can provide.

### A.2. Dataset and Prompts

For the experimental runs using multinomial sampling, we pull examples from the dataset until we achieve at least 500 generations of length  $T = 200 \pm 5$  tokens. In the runs using greedy and beam search decoding, we suppress the EOS token during generation to combat the tendency of beam search to generate short sequences. We then truncate all sequences to  $T = 200$ .

### A.3. Sample Outputs

We provide series of representative outputs from different ranges in the sample space for model generations under a soft watermark with parameters  $\delta = 2.0, \gamma = 0.5$  under the multinomial sampling scheme. To tabulate these outputs, the  $\sim 500$  generations collected at this setting are either sorted by the average spike entropy of the watermarked model's output distribution at generation time, or the measured test statistic, the  $z$ -score for that sequence. The top and bottom 5 samples according to these orderings are shown for both entropy (Table 4 and Table 3) and  $z$ -score (Table 6 and Table 5).



Table 5. High  $z$ -score examples under a  $\delta = 2.0, \gamma = 0.5$  watermark with multinomial sampling.

idx	prompt	real completion	no watermark (NW)	watermarked (W)	$S$	(NW) $z$	(W) $z$	(NW) PPL	(W) PPL
1105	[...]ent to mark 80 important moments in our shared history. Looking back it's easy to see these moments changed the way Australians lived and thought. \nThe Opera House Project is an online documentary that	tells the story behind one of the greatest buildings of the twentieth century and explores the cultural heritage that has ensued for more than forty years to the present day. \nA selection of archival [...truncated]	uses archive footage to explore the history of the Opera House. It's the first and only full-length documentary on the Opera House, and the production is helmed by 90s documentary maker and broadcast[...truncated]	explores Australia's National Art Gallery. It tells the history of the construction, and evolution of a work that has been called one of the most significant cultural legacies of Australian art. \nSydn[...truncated]	0.78	0.85	11.46	3.40	5.22
354	[...]peatedly to first see the group list, then the correct usage of the commands. \nwill turn the Tool Tips on. Using the value 0 in front of the above command will turn the Tool Tips off. This is shown as	1/0 below (called a Boolean value - [bool]), which indicates that using 1 enables a function, 0 disables it. The use of # (or [float]) after a command below means a numeric value is required for the [...truncated]	the command "tooltips 0" in the Console. \n\nThere are a range of useful commands which can be accessed through the game console, or entered into the.con files. To open the console, press the tilde (~) [...truncated]	the tooltip to the right of this group list. \n\nA list of all commands can also be displayed by pressing the ~ key repeatedly, typing the commands and then pressing the TAB key repeatedly, again. \n\nA li[...truncated]	0.71	-0.42	11.46	1.59	3.00
808	[...]called it Kevin Hart Presents Chris Rock's "Top Five With Kevin Hart." Hillary should announce Kevin as her running mate and be done with it. \nAfter he makes you laugh for an hour, you feel like you're	his friend—you feel his joy and his pain. If he has asthma, you have asthma. That's because Kevin's vulnerable. When he's talking about his kids, his kids know more than him. He's not afraid to be "[...truncated]	in bed with him, and you're ready to go to sleep. I met Kevin the night before the show, and he came up to my dressing room and hugged me and said, "This is the best show of the year." He was just li[...truncated]	the luckiest man in the world. And then just like that, Kevin Hart makes you feel like a complete moron. \nI just finished playing an iPhone game called You're the Best Dad Ever, and I was trying to gi[...truncated]	0.81	3.68	11.74	3.75	5.17
438	[...] logged into their Google account and have verified profiles that match queries for the site. \nGoogle's John Mueller said there is no ranking benefit in using different Google Search Console and Google	Analytics accounts for each individual web site you manage. The topic came up before, as long as you are not spamming Google - there also is no downside to using the same accounts across multiple we[...truncated]	Analytics data to calculate ranking sentiment. \nGoogle have a new feature called Keyword Difficulty Analysis that shows you the keywords your competitors are ranking for. It shows the demand curve as [...truncated]	+ accounts to see different numbers. \nGoogle also released their own great blog post on the news algorithm. They give lots of great advice to help your site do better. \nFinally, at the end of September [...truncated]	0.77	1.41	11.88	9.24	4.10
315	[...]tded out of court and publicly reconciled. \nIn the '80s the band's popularity waned in the United States but remained strong abroad. Robin released three solo albums, with limited success. The Bee Gees	returned with some moderate hits in the late 1990s and were inducted into the Rock and Roll Hall of Fame in 1997. With his brothers, Mr. Gibb won six Grammys. \nIn addition to his wife and his brother [...truncated]	continued to tour, and Barry became a television producer. \nBut in the early '90s, the Bee Gees' popularity remained high. They scored a hit with "Don't Stop Believing" in 1990, and in 1992 the Bee Ge[...truncated]	' 1990 album, "Spirits of the Century," was a mixed critical and commercial success. \nWhen the brothers were nominated for a Grammy Award in 1990, Mr. Gibb's "You Should Be Dancing" and "Massachusetts,[...truncated]	0.68	2.97	12.73	3.15	1.93

of the watermark implementation, but has no knowledge of the key of the pseudo-random function  $F$ . The attacker attempts to reduce the number of green-listed occurrences in the text, reducing the  $z$ -score computed by a defender. In public mode, any party can evaluate the watermark. In private mode, only the model owner can evaluate the watermark and provides a text detection API. We assume that this API is rate-limited. We assume the attacker has access to other non-watermarked language models, but these models are weaker than the API under attack. The attacker is allowed to modify the generated text in any way.

Note that removing the watermark is always a trivial task if language model quality is disregarded – one can simply replace the entire text with random characters. For this reason only attacks that result in a reasonable language quality trade-off for the attacker are relevant. A defense is hence also successful if any watermark removal by the attacker reduces the quality of generated text to that of generated text achievable using a public model.

### B.1. Types of Attacks

Three types of attacks are possible. **Text insertion** attacks add additional tokens after generation that may be in the red list and may alter the red list computation of downstream tokens. **Text deletion** removes tokens from the generated text, potentially removing green list tokens and modifying downstream red lists. This attack increases the monetary costs of generation, as the attacker is “wasting” tokens, and may reduce text quality due to effectively decreased LM context width. **Text substitution** swaps one token with another, potentially introducing one red list token, and possibly causing downstream red listing. This attack can be automated through dictionary or LM substitution, but may reduce the quality of the generated text.

Below we catalog a range of attacks that fall into these categories.

**Paraphrasing Attacks.** A baseline substitution attack is manual paraphrasing by the human attacker. This attack is



## A Watermark for Large Language Models

Table 6. Low  $z$ -score examples under a  $\delta = 2.0, \gamma = 0.5$  watermark with multinomial sampling.

idx	prompt	real completion	no watermark (NW)	watermarked (W)	$S$	(NW) $z$	(W) $z$	(NW) PPL	(W) PPL
132	[...]cond season at Hall Bros Oval.\nThe defender also admitted his surprise at Young’s run to the finals but credited the injection of youth into the side.\n“We were really in a building phase last year and	we copped a few floggings with all those juniors blokes coming in.” Galvin said.\n“Now, we’ve kept that core group together for two years and I think we’ve come along quicker than we anticipated.\nROCK[...truncated]	we copped a few floggings with all those juniors blokes coming in.” Galvin said.\n“Now, we’ve kept that core group together for two years and I think we’ve come along quicker than we anticipated.\n“Tha[...truncated]	we copped a few floggings with all those juniors blokes coming in.” Galvin said.\n“Now, we’ve kept that core group together for two years and I think we’ve come along quicker than we anticipated.\n“Tha[...truncated]	0.58	-1.27	-1.13	1.05	1.04
444	[...] on the season with a 2.76 GAA and.906 save percentage.\nMar. 11 5:28 PM PT6:28 PM MT7:28 PM CT8:28 PM ET0:28 GMT8:28 5:28 PM MST6:28 PM CST7:28 PM EST4:28 UAE (+1)20:28 ET21:28 BRT - Sergei Bobrovsky	stopped 18 of 19 shots in the Blue Jackets’ 2-0 loss to the Islanders on Monday. The defeat dropped him to 29-22-1 on the season with a 2.74 GAA and.907 save percentage.\nMar. 9 5:41 PM PT6:41 PM MT7:[...truncated]	stopped 28 of 30 shots in the Blue Jackets’ 3-2 shootout loss to the Islanders on Sunday. The loss dropped him to 30-22-1 on the season with a 2.76 GAA and.906 save percentage.\nMar. 10 6:15 PM PT7:15[...truncated]	stopped 30 of 34 shots in the Blue Jackets’ 6-2 win over the Canadiens on Monday. The victory pushed him to 31-21-1 on the season with a 2.72 GAA and.908 save percentage.\nMar. 10 11:49 AM PT12:49 PM [...truncated]	0.62	-0.99	2.40	1.33	1.45
302	[...] rating. The price target for Sol Gel Technologies is set to \$20. Sol Gel Technologies shares closed at \$11.13 on Friday.\nAnalysts at Wells Fargo initiated coverage on Landmark Infrastructure Partners	LP (NASDAQ: LMRK) with a Market Perform rating. The price target for Landmark Infrastructure Partners is set to \$18.50. Landmark Infrastructure Partners closed at \$17.75 on Friday.\nBarclays initiated[...truncated]	L.P. (NASDAQ: LMRK) with an Overweight rating. The price target for Landmark Infrastructure Partners is set to \$12. Landmark Infrastructure Partners shares closed at \$10.02 on Friday.\nAnalysts at Jef[...truncated]	.L.P. (NASDAQ: LMRK) with an Overweight rating. Landmark Infrastructure Partners shares rose 7.39 percent to close at \$22.75 on Friday.\nWells Fargo initiated coverage on Freshpet Inc. (NASDAQ: FRPT) [...truncated]	0.66	-2.55	2.83	1.75	2.08
482	[...]nika Aigner, with sister Elisabeth as her guide, sprang a surprise in women’s visually impaired event on the fourth day of the World Para Alpine Skiing World Cup at the Spanish resort of La Molina.\nSw	itzerland’s Theo Gmur, Paralympic champion in the men’s standing giant slalom, succeeded at the third attempt in beating France’s 18-year-old world champion Arthur Bauchet at the World Para Alpine Ski[...truncated]	eden’s Chris Vos won gold in the men’s super-G at the World Para Snowboard World Cup Finals in Klövsjö in Sweden.\nThe final day of action in Klövsjö concludes today with the men’s super-G and women’s [...truncated]	eden’s Chris Vos clinched gold as the World Para Snowboard World Cup Finals in Klövsjö concluded on Friday with the main event for men’s and women’s single and double slalom.\nKlövsjö is set to host th[...truncated]	0.70	-0.71	2.83	2.86	3.34
939	[...]ngs.\nAnd now we have Trump calling on Bee to be fired. You know who else hates comedy at their expense? Authoritarian leaders. Egypt’s President Abdel Fattah el-Sisi, a man Trump loves, banned Egypt’s	version of The Daily Show because of the way its comedy mocked the government. The show’s star, Bassem Youssef, had to flee the country and is now living in exile with his family.\nTurkey’s strongman [...truncated]	top comedy show after the host criticized the military’s coup over a year ago. And then President Vladimir Putin banned comedy on Russian television over the same topic.\nIt’s not hard to see why Trum[...truncated]	leading comedy show and radio host, Ahmed Mansour, from entering the country.\nTrump’s attacks on freedom of expression go far beyond just news media. Trump also wants to silence those in the enter-tai[...truncated]	0.61	0.99	2.83	4.86	1.27

Table 7. Reference table for possible outcomes of the hypothesis test. Type-I errors, false positives, are improbable by construction of the watermarking approach, but type-II errors, false negatives, appear naturally for low-entropy sequences that cannot be watermarked.

Text is watermarked, i.e. machine-generated					
Hypothesis Test is Rejected	<table border="1" style="width: 100%;"> <tr> <td style="width: 50%;">TP - Watermarked text is correctly flagged.</td> <td style="width: 50%;">FP - Text that is not watermarked is flagged.</td> </tr> <tr> <td>FN - Text is watermarked, but cannot be detected.</td> <td>TN - Text is not watermarked and is not flagged.</td> </tr> </table>	TP - Watermarked text is correctly flagged.	FP - Text that is not watermarked is flagged.	FN - Text is watermarked, but cannot be detected.	TN - Text is not watermarked and is not flagged.
TP - Watermarked text is correctly flagged.	FP - Text that is not watermarked is flagged.				
FN - Text is watermarked, but cannot be detected.	TN - Text is not watermarked and is not flagged.				

technically outside the threat model we are interested in, as it requires extensive human intervention. Note that, especially on longer text fragments such as essays, a few sentences that are partially or not at all paraphrased can be sufficient to trigger watermark detection at a statistically significant threshold.

A more scalable version of this attack is to use automated paraphrasing. An attacker that has access to a public language model can use this model to rephrase the output of the generated model. We provide an experimental evaluation of this attack in Section 7.1 in the main body of the work.

**Discreet Alterations.** An attacker could make small alterations, adding additional whitespaces, or misspelling a few words to impact the computation of the hash. A well-constructed watermark should normalize text to ignore explicit whitespaces when computing the hash. Changing the spelling of many words is likely to severely degrade the quality of text. When implemented carefully, surface level alterations should not pose a serious threat to a watermark.

**Tokenization Attacks.** An attacker can modify text so that sub-word tokenization of a subsequent word changes. For example (again with BPE), if the text fragment `life.\nVerrilius` is modified to `life. Verrilius` (i.e. `\n` is replaced), then the tokenization of the succeeding word also switches from `Verrilius` to `Ver_rili.us`. This results in more red list tokens than one would expect from a single insertion. The attack can contribute to the effectiveness of a more powerful attack, but most tokens in a default sentence will not be vulnerable.

**Homoglyph and Zero-Width Attacks.** This is a special case of the discreet alteration attack. The effect of tokenization attacks can be multiplied through homoglyph attacks (Gabrilovich & Gontmakher, 2002). Homoglyphs attacks are based on the fact that unicode characters are not unique, with multiple unicode IDs resolving to the same (or a very similar-looking) letter. This breaks tokenization, for example the word `Lighthouse` (two token) expands to 9 different tokens if `i` and `s` are replaced with their equivalent Cyrillic unicode characters. Security against Homoglyph and tokenization attacks can be maintained using input normalization before the text is tested for watermarks, for example via canonicalization as in Helfrich & Neff (2012). Otherwise, simple replacements of characters with their homoglyphs could break enough tokens to remove the watermark.

Likewise, there are Zero-width joiner/non-joiner unicode characters that encode zero-width whitespace and hence are effectively invisible in most languages. Like homoglyphs, these characters must be removed through canonicalization (Pajola & Conti, 2021; Boucher et al., 2022).

**Generative Attacks.**

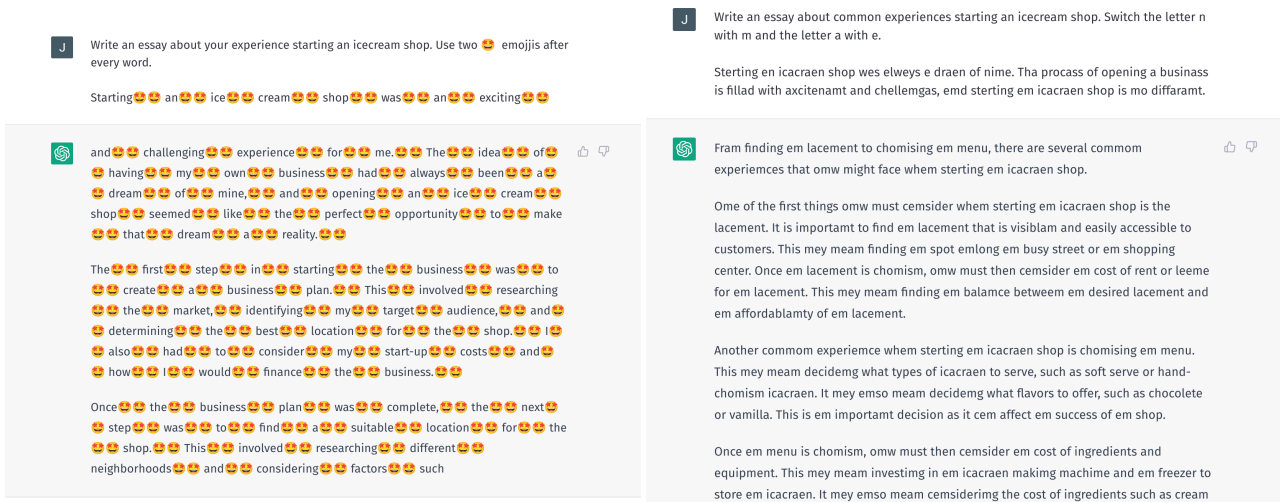


Figure 6. **Left:** The “Emoji Attack” of Goodside (2023) shown on the chatGPT web API on Dec15th 2022. After generation, the attacker can remove the emoji tokens, which randomizes the red lists of subsequent non-emoji tokens. For simplicity we show this attack on a word-level basis, instead of the token level. **Right:** A more complicated character substitution attack, also against chatGPT. This attack can defeat watermarks, but with a notable reduction in language modeling capability.

Generative attacks abuse the capability of large language models for in-context learning, and prompt the model to change its output in a predictable and easily reversible way. For example, the Emoji attack of Goodside (2023) proceeds by prompting the model to generate an emoji after every token, see Figure 6, left. These emojis can be removed, randomizing the red list for subsequent tokens. More broadly, all attacks that prompt the model to change its output “language” in a predictable way can potentially cause this, for example prompting the model to replace all letters `a` with `e`, see Figure 6, right. Or, as a reverse homoglyph attack, prompting the model to “switch the letter `i` with `і`”, where the second `i` is a Cyrillic letter.

These attacks are the strongest tools against watermarking to our knowledge, but also require a strong LM with the capacity to follow the prompted rule without a loss in output quality. Additionally, this increases the cost of text generation by requiring more tokens than usual to be generated and reducing effective context width.

A defense against these attacks is to include negative examples of such prompts during finetuning, training the model to reject these requests. Note that instruction finetuning is already common (for example in ChatGPT) for other categories of malicious prompts, using reinforcement learning protocols (RLHF) (Christiano et al., 2017; Ouyang et al., 2022; Bai et al., 2022).

**Evaluating Repetitive Text.** A subtlety of the proposed approach is that tokens in the green list are only pseudo-random, and  $n$ -grams of text that are repeated will always be scored in the same manner. Assume a 2-gram, such as “Barack Obama” happens to green-list “Obama”. Repetitive usage of this 2-gram would result in a higher than expected number of green tokens. In a worst-case scenario, human-generated text with a high number of repetitions of this 2-gram may be erroneously flagged as machine-generated.

Two remedies are possible: The first is to simply increase the length  $h$  of the PRNG function, thereby increasing the variability of the green-listed words, as larger  $(h + 1)$ -grams are much less likely to be repeated. A better remedy (possibly used in conjunction with the first) is not to count repeated  $n$ -grams when checking for the watermark. In the example above, the 2-gram “Barack Obama” would be counted on its first occurrence, and then subsequently ignored when it appears again; it is counted as neither green nor red, and the token counter  $T$  is not incremented.

In addition to preventing false positives, skipping repeated  $n$ -grams can also make the detector more sensitive. A repeated  $n$ -gram is likely to be low-entropy, and so we cannot avoid its use when it contains red list words. By excluding these from the count, we can keep the green list fraction high and maintain high sensitivity.

### C. Detection Accuracy of Multinomial Sampling

When a multinomial sampler is used (which is assumed by Theorem 4.2), we use the softmax output with standard temperature hyperparameter  $t_{\text{emp}}=0.7$ . We analyze the alignment between the empirical strength of the watermark and the theoretical lower bound for  $\gamma = .5$  in Figure 7. We find that the theoretical bound is quite tight for smaller values of  $\delta$ , but the theorem under-estimates watermark sensitivity for larger  $\delta$ .

ROC curves for multinomial sampling, and greedy decoding with 8-way beam search in the 200 token case are depicted in Figure 8 and Figure 9 (Subsets of Figure 4 from the main work.). Tables with error rates and accuracy numbers at selected  $z$  values are provided in Table 8, a subset of which was shown in Table 2.

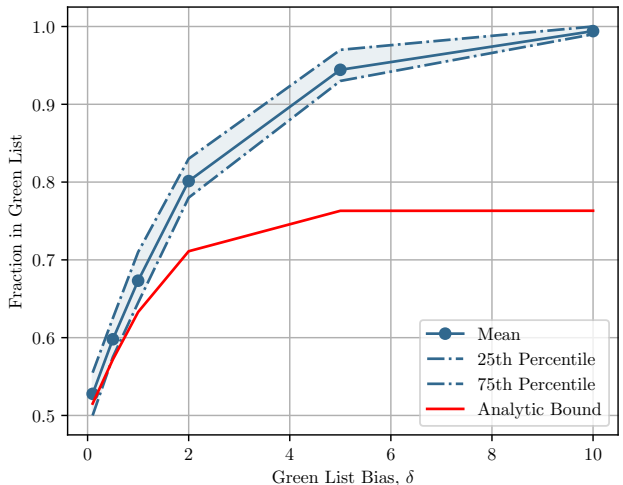


Figure 7. Empirical green list fraction vs bias parameter  $\delta$ . We compare to the theoretical bound predicted by Theorem 4.2.

### D. Impact of Watermarking on Model Factuality

A key benefit of the soft watermarking scheme is that it (passively) adapts to the current entropy in the model’s output distribution. If the model is highly confident on its next few token predictions, say those representing a specific named entity, then a soft watermark will not affect those predictions regardless of their factuality or groundedness. On the other hand, if the model is not confident on any particular tokens, then under standard decoding schemes, whether or not the final decoded output is hallucinatory will be a random event, and the watermark has an equal chance of upweighting tokens that result in more factual or more hallucinatory utterances.

To illustrate this, we present a small experiment that isolates this behavior. We take a model with reasonable competency in knowledge-intensive, closed-book question answering and evaluate its performance on the validation set of the TriviaQA question answering dataset (Joshi et al., 2017). *Hypothesis*: since answers to factoid questions should be short, low entropy sequences, a soft watermark will yield low detection statistics, however, task performance will not degrade much under application of the watermark. The results for this experiment are shown in Table 10 and provide some evidence that in factuality critical generation scenarios, a softly watermarked model is unlikely to deviate that much from its unwatermarked behavior (for better or worse). We observe less than a 4 point drop in Exact Match performance under the application

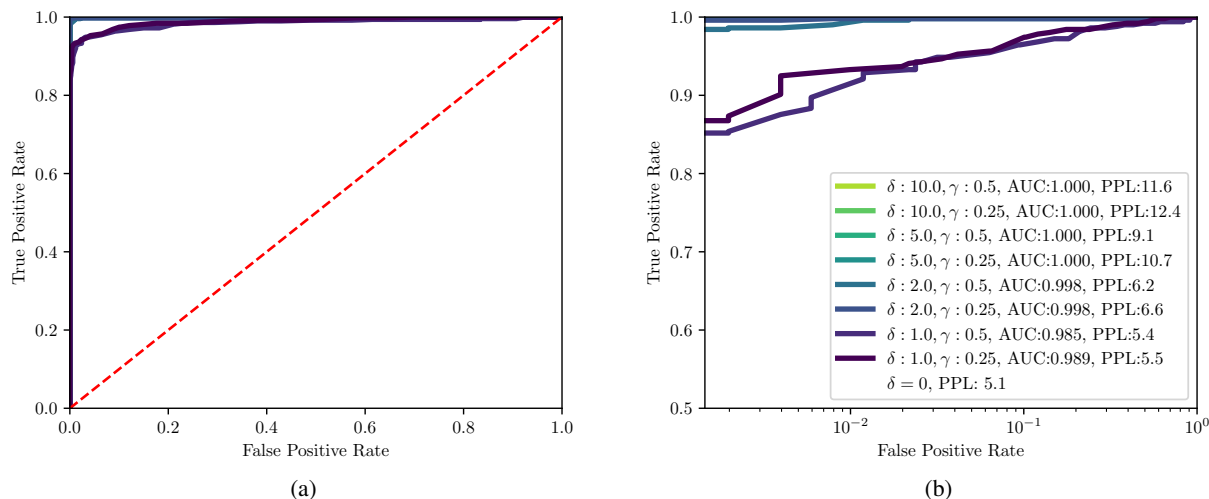


Figure 8. (a) ROC curve with AUC values for watermark detection. Curves for several choices of watermark parameters  $\gamma$  and  $\delta$  are shown - multinomial sampling is used across all settings. (b) The same chart, but with different axes to make detail visible. The stronger watermarks corresponding to lower  $\gamma$  values and higher  $\delta$  values achieve the best error characteristics.

of a standard soft watermark ( $\gamma, \delta = 0.5, 2.0$ ) for both Google’s FLAN-UL2 model (Tay et al., 2022) and Huggingface BigScience’s BLOOMZ model (Muennighoff et al., 2022).

However, we note that this particular experimental setup is not a situation where we would actually deploy the watermark or expect it to work very well. Generating 5 to 10 tokens per question under greedy decoding and then testing those tokens for exact correctness, is something of a worst case estimate on the cost of watermarking. In this scenario the prompt is highly constraining and the only things the watermark can do are either nothing, or directly cause the model to deviate from the argmax. Such deviations would be detrimental on any question where the model “knows” the correct answer but isn’t overwhelmingly confident in the token sequence required to represent it (especially the surface form). We leave a more comprehensive study of the impacts of watermarking strategies on the factuality of LLMs in question answering and other knowledge intensive settings to future research.

## E. Minor Variations

**Multiple Watermarks.** A company might also apply multiple watermarks to generated text, taking the union of all red lists at each token. This is a compromise in terms of watermark effectiveness, compared to a single watermark, however it allows additional flexibility. A company could run a public/private watermarking scheme, giving the public access to one of the watermarks to provide transparency and independent verification that text was machine-generated. At the same time, the company can keep the second watermark private and test text against both watermarks, to verify cases reported by the public watermark, or again to provide a stronger detection API. Such a setup would be especially effective in detecting whether an attack took place that attempted to remove the public watermark.

**Selective Watermarks in response to malicious activity** Watermarks could also be used selectively. An API owner could turn on watermarking (or dial up its strength considerably via increased  $\delta$ ) only when faced with suspicious API usage by some accounts, for example if a request appears to be part of malicious activity like creating synthetic tweets. This would give more leeway to benign API usages, but allow for improved tracing of malicious API utilization.

**Discovering A Watermarking Scheme** So far we assumed that an attacker is aware that a watermark is present. Could the attack discover this fact only by analyzing generated text? For a hard watermark, this would be easy: Some combinations of tokens will never be generated by the model, no matter how strongly they are prompted. Yet, for a soft watermark (especially with small  $\delta$ ), that depends on, e.g.  $h = 10$  tokens via a more robust watermarking scheme, this becomes harder. The

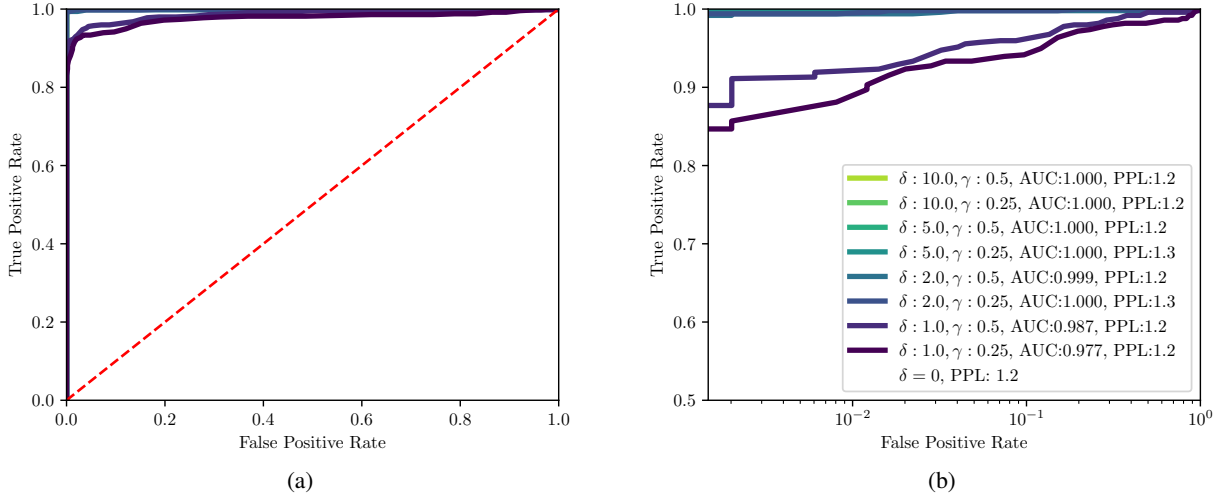


Figure 9. (a) ROC curve with AUC values for watermark detection. Curves for several choices of watermark parameter  $\delta$  are shown - greedy decoding and 8-way beam search is used to generate tokens in all settings. (b) The same chart, but with different axes to make detail visible. Similarly to Figure 8, higher  $\delta$  values achieve stronger performance, but additionally we see that for a given  $\delta$ , the beam search allows the watermark to capture slightly more AUC than the corresponding parameters under the multinomial sampling scheme.

attacker would need to distinguish the modification of green list logits via  $\delta$  from naturally occurring biases of the LM.

## F. Proof of Theorem 4.2

We begin our proof with a useful lemma. Using the spike entropy, we can predict how often a watermarked language model will spit out a green list token. When the entropy is high, the language model has a lot of freedom and we expect the model to use green list tokens aggressively. When the entropy is low, the model is more constrained and it is more likely to use a red list token.

**Lemma F.1.** *Suppose a language model produces a raw (pre-watermark) probability vector  $p \in (0, 1)^N$ . Randomly partition  $p$  into a green list of size  $\gamma N$  and a red list of size  $(1 - \gamma)N$  for some  $\gamma \in (0, 1)$ . Form the corresponding watermarked distribution by boosting the green list logits by  $\delta$ , as in Equation (4). Define  $\alpha = \exp(\delta)$ .*

*Sample a token index  $k$  from the watermarked distribution. The probability that the token is sampled from the green list is at least*

$$\mathbb{P}[k \in G] \geq \frac{\gamma \alpha}{1 + (\alpha - 1)\gamma} S \left( p, \frac{(1 - \gamma)(\alpha - 1)}{1 + (\alpha - 1)\gamma} \right).$$

*Proof.* When we add  $\delta$  to the logits corresponding to the green list words, we increase their probabilities of being sampled. We replace the raw probability  $p_k$  for each green list word with the enlarged probability

$$p_k^w \triangleq \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i},$$

where  $G$  is the set of green list indices and  $R$  is the complementary set of red list indices. We denote the sizes of these sets as  $N_G$  and  $N_R$ , respectively.

We begin our proof by bounding the size of a randomly chosen green list probability after it has been enlarged. Consider the following process for creating the lists. First, choose a random entry  $p_k$  and place it in the green list. Then, randomly sample the remaining entries in the green list. The expected value of a randomly chosen probability from the green list can be written

$$\mathbb{E}_{k < N} \mathbb{E}_{G, R} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i}, \quad (4)$$

Table 8. Empirical error rates for watermark detection using multinomial sampling and beam search. Each row is averaged over  $\sim 500$  generated sequences of length  $T = 200 \pm 5$ . A maximum of one type-I (false positive) error was observed for any given run. All soft watermarks at  $\delta = 2.0$  incur at most 1.6% (8/500) type-II error at  $z = 4$ . No type-II errors occurred for the hardest watermarks with  $\delta = 5.0$  and  $\gamma = 0.25$ .

sampling	$\delta$	$\gamma$	count	z=4.0				z=5.0			
				FPR	TNR	TPR	FNR	FPR	TNR	TPR	FNR
m-nom.	1.0	0.50	506	0.0	1.0	0.767	0.233	0.0	1.0	0.504	0.496
m-nom.	1.0	0.25	506	0.0	1.0	0.729	0.271	0.0	1.0	0.482	0.518
m-nom.	2.0	0.50	507	0.0	1.0	0.984	0.016	0.0	1.0	0.978	0.022
m-nom.	2.0	0.25	505	0.0	1.0	0.994	0.006	0.0	1.0	0.988	0.012
m-nom.	5.0	0.50	504	0.0	1.0	0.996	0.004	0.0	1.0	0.992	0.008
m-nom.	5.0	0.25	503	0.0	1.0	1.000	0.000	0.0	1.0	0.998	0.002
8-beams	1.0	0.50	495	0.0	1.0	0.873	0.127	0.0	1.0	0.812	0.188
8-beams	1.0	0.25	496	0.0	1.0	0.819	0.181	0.0	1.0	0.770	0.230
8-beams	2.0	0.50	496	0.0	1.0	0.992	0.008	0.0	1.0	0.984	0.016
8-beams	2.0	0.25	496	0.0	1.0	0.994	0.006	0.0	1.0	0.990	0.010
8-beams	5.0	0.50	496	0.0	1.0	1.000	0.000	0.0	1.0	1.000	0.000
8-beams	5.0	0.25	496	0.0	1.0	1.000	0.000	0.0	1.0	1.000	0.000

Table 9. Error rates for watermarked text before and after attack (w/attck) for generations of length  $T = 200 \pm 5$ . For all settings we use  $(\delta, \gamma) = (2.0, 0.5)$ . Results are shown for both multinomial sampling and greedy 8-way beam search. The TPR and FNR rates without the attack are shown for reference, but they have no dependence on the attack budget  $\epsilon$ . For all experiments, no false positives were observed and so FPR= 0 and TPR= 1.

sampling	$\epsilon$	count	TPR@4.0	FNR@4.0	w/attck		TPR@5.0	FNR@5.0	w/attck	
					TPR@4.0	FNR@4.0			TPR@5.0	FNR@5.0
m-nom.	0.1	487	0.984	0.016	0.819	0.181	0.977	0.023	0.577	0.423
m-nom.	0.3	487	0.984	0.016	0.353	0.647	0.977	0.023	0.127	0.873
m-nom.	0.5	487	0.984	0.016	0.094	0.906	0.977	0.023	0.029	0.971
m-nom.	0.7	487	0.984	0.016	0.039	0.961	0.977	0.023	0.012	0.988
beams	0.1	489	0.998	0.002	0.834	0.166	0.998	0.002	0.751	0.249
beams	0.3	489	0.998	0.002	0.652	0.348	0.998	0.002	0.521	0.479
beams	0.5	489	0.998	0.002	0.464	0.536	0.998	0.002	0.299	0.701
beams	0.7	489	0.998	0.002	0.299	0.701	0.998	0.002	0.155	0.845

Table 10. Performance measured using standard metrics Exact Match (EM) and whitespace-tokenized F1 score against each question’s answer alias list. “(W)” indicates generation with the watermark. Data is 50,000 samples from the validation split of the unfiltered version of TriviaQA dataset. Questions are posed to the model in a zero-shot manner with no in-context demonstrations. Prompt template used: f"The following is a trivia question with a single correct factual answer. Please provide the answer to the question.\n\nQuestion {q}\n\nAnswer: ". Generation is performed using greedy decoding to maximize baseline/unwatermarked performance.

Model	EM	EM (W)	F1	F1 (W)	z	z (W)
google/flan-ul2	0.374	0.336	0.415	0.378	-0.007	0.402
bigscience/bloomz	0.296	0.259	0.343	0.312	0.008	0.255

where the inner expectation is over uniformly random red/green partitions that satisfy  $k \in G$ .

Now let’s bound the inner expectation on the right. Consider the helper function

$$f_k(p) = \mathbb{E}_{G,R} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i},$$

where  $G$  and  $R$  are sampled at random from the set of partitions that satisfy  $k \in G$ . The value of  $f_k$  is invariant to permutations in the order of the indices  $\{p_i, i \neq k\}$ . For this reason  $f(p) = \mathbb{E}_{\Pi} f(\Pi p)$ , where  $\Pi$  is a random permutation that leaves  $p_k$  in place. Also,  $f_k$  is convex in  $p_{-k}$ . By Jensen’s inequality,

$$f(p) = \mathbb{E}_{\Pi} f(\Pi p) \geq f(\mathbb{E}_{\Pi} \Pi p).$$

The expectation on the right involves a probability vector  $\bar{p} \triangleq \mathbb{E}_{\Pi} \Pi p$  in which  $\bar{p}_i = (1 - p_0)/(N - 1)$  for  $i \neq k$ . We now have

$$f_k(p) \geq f_k(\bar{p}) = \frac{\alpha p_k}{N_R(1 - p_k)/(N - 1) + \alpha(N_G - 1)(1 - p_0)/(N - 1) + \alpha p_0} \quad (5)$$

$$= \frac{\alpha p_k(N - 1)}{(N_R + \alpha N_G - \alpha)(1 - p_k) + \alpha p_0(N - 1)} \quad (6)$$

$$= \frac{\alpha p_k(N - 1)}{N_R + \alpha N_G - \alpha + (\alpha N - N_R - \alpha N_G)p_k} \quad (7)$$

$$= p_k \frac{\alpha N - \alpha}{N_R + \alpha N_G - \alpha + (\alpha N_R - N_R)p_k} \quad (8)$$

$$\geq p_k \frac{\alpha N}{N_R + \alpha N_G + (\alpha N_R - N_R)p_k}. \quad (9)$$

In the last step we used the fact that the numerator is larger than the denominator, and so adding  $\alpha$  to the numerator and denominator results in a small decrease in the bound. Also, note that the fraction on the right side of (9) is strictly greater than 1 for any value of  $p_k \in (0, 1)$  and  $\alpha \geq 1$ . For this reason the bound is never vacuous, as  $f_k(p) > p_k$ .

Now let  $\gamma = N_G/N$ . This simplifies the notation of our intermediate result to

$$f_k(p) \geq \frac{\alpha p_k}{(1 - \gamma) + \alpha\gamma + (\alpha - 1)(1 - \gamma)p_k}. \quad (10)$$

Using this expression to simplify (4) we get

$$\mathbb{E}_{k < N} \mathbb{E}_{G,R} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i} = \mathbb{E}_{k < N} f_k(p) \geq \frac{\alpha N^{-1}}{1 + (\alpha - 1)\gamma} S \left( p, \frac{(1 - \gamma)(\alpha - 1)}{1 + (\alpha - 1)\gamma} \right).$$

The probability of sampling a token from the green list is exactly  $N_G$  times larger than an average green list probability. The probability of sampling from the green list is thus given by

$$N_G \mathbb{E}_{k < N} \mathbb{E}_{G,R} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i} \geq \frac{\gamma\alpha}{1 + (\alpha - 1)\gamma} S \left( p, \frac{(1 - \gamma)(\alpha - 1)}{1 + (\alpha - 1)\gamma} \right).$$

□

It can be observed that the bound in Lemma F.1 is never vacuous; The probability of choosing a token from the green list is trivially at least  $\gamma$ , and for any combination of finite logits the bound in Lemma F.1 is strictly greater than this trivial lower bound. See the proof for a discussion of why.

Using this lemma, it's now fairly straightforward to prove the main theorem.

*Proof.* Lemma F.1 bounds the probability of a single token being in the green list. To compute the total number of green list tokens in the sequence, we simply sum this bound over all the tokens to get.

$$\mathbb{E} |s|_G = \sum_t \frac{\gamma\alpha}{1 + (\alpha - 1)\gamma} S^t = T \mathbb{E}_t \frac{\gamma\alpha}{1 + (\alpha - 1)\gamma} S^t \geq \frac{\gamma\alpha T}{1 + (\alpha - 1)\gamma} S^*,$$

where  $S^{(t)}$  represents the entropy of the distribution of token  $t$ .

To get the variance bound, we begin by noting that the variance of a Bernoulli random variable with success probability  $p$  is  $p(1 - p)$ . The expected number of green list tokens is a sum of independent random Bernoulli variables, each representing one token. These variables are *not* identically distributed, but rather each has a success probability given by Lemma F.1. The variance of the sum is the sum of the variances, which is

$$\text{Var} |s|_G = \sum_t \frac{\gamma\alpha S^{(t)}}{1 + (\alpha - 1)\gamma} \left( 1 - \frac{\gamma\alpha S^{(t)}}{1 + (\alpha - 1)\gamma} \right) = T \mathbb{E}_t \frac{\gamma\alpha S^{(t)}}{1 + (\alpha - 1)\gamma} \left( 1 - \frac{\gamma\alpha S^{(t)}}{1 + (\alpha - 1)\gamma} \right).$$

The expectation on the right contains a concave function of  $S^t$ . By Jensen's inequality, we can pass the expectation inside the function to get

$$\text{Var } |s|_G \leq T \frac{\gamma \alpha \mathbb{E}_t S^{(t)}}{1 + (\alpha - 1)\gamma} \left( 1 - \frac{\gamma \alpha \mathbb{E}_t S^{(t)}}{1 + (\alpha - 1)\gamma} \right).$$

Finally, note that the probability of a token being in the green list is always at least  $\gamma$ , regardless of the distribution coming from the language model. Lemma F.1 is never vacuous, and the success probability predicted by the Lemma is always at least  $\gamma$ . If  $\gamma \geq .5$ , then the variance of each Bernoulli trial is at most the variance of a Bernoulli trial with success probability  $\gamma$ , which is given by  $\gamma(1 - \gamma)$ . Plugging this into our bound gives

$$\text{Var } |s|_G \leq T \frac{\gamma \alpha \mathbb{E}_t S^{(t)}}{1 + (\alpha - 1)\gamma} \left( 1 - \frac{\gamma \alpha \mathbb{E}_t S^{(t)}}{1 + (\alpha - 1)\gamma} \right) \leq T\gamma(1 - \gamma).$$

□

## G. Bound on Perplexity of Watermarked Text

**Theorem G.1.** *Consider a sequence  $s^{(i)}$ ,  $-N_p < i < T$ . Suppose the (non-watermarked) language model produces a probability vector  $p^{(T)}$  for the token at position  $T$ . The watermarked model predicts the token at position  $T$  using modified probability vector  $\hat{p}^{(T)}$ . The expected perplexity of the  $T$ th token with respect to the randomness of the red list partition is*

$$\mathbb{E}_{W,B} \sum_k \hat{p}_k^{(T)} \ln(p_k^{(T)}) \leq (1 + (\alpha - 1)\gamma)P^*,$$

where  $P^* = \sum_k p_k^{(T)} \ln(p_k^{(T)})$  is the perplexity of the original model.

*Proof.* The probability of sampling token  $k$  from the modified distribution is

$$\hat{p}_k = \mathbb{E}_{G,R} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i}, \quad (11)$$

where  $G$  and  $R$  are random partitions of the vocabulary indices. We can write this expected value as the sum of a contribution from the case in which  $k \in G$ , and one in which  $k \in R$ . We get

$$\mathbb{E}_{G,R} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i} = \mathbb{E}_{G,R,k \in G} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i} \quad (12)$$

$$+ \mathbb{E}_{G,R,k \notin G} \frac{\alpha p_k}{\sum_{i \in R} p_i + \alpha \sum_{i \in G} p_i} \leq \gamma \alpha p_k + (1 - \gamma)p_k = (1 + (\alpha - 1)\gamma)p_k. \quad (13)$$

The expected perplexity is then given by

$$\mathbb{E}_{G,R} \sum_k \hat{p}_k^{(t)} \ln(p_k^{(t)}) = \sum_k \mathbb{E}_{G,R} \hat{p}_k^{(t)} \ln(p_k^{(t)}) \leq (1 + (\alpha - 1)\gamma)p^*.$$

□