

# A Robust Semantics-based Watermark for Large Language Models against Paraphrasing

Anonymous ACL submission

## Abstract

Large language models (LLMs) have show their remarkable ability in various natural language tasks. However, there are concerns that LLMs are possible to be used improperly or even illegally. To prevent the malicious usage of LLMs, detecting LLM-generated text becomes crucial in the deployment of LLM applications. Watermarking is an effective strategy to detect the LLM-generated content by encoding a pre-defined secret watermark to facilitate the detection process. However, the majority of existing watermark methods leverage the simple hashes of precedent tokens to partition vocabulary. Such watermarks can be easily eliminated by paraphrase and, correspondingly, the detection effectiveness will be greatly compromised. Thus, to enhance the robustness against paraphrase, we propose a semantics-based watermark framework, SemaMark. It leverages the semantics as an alternative to simple hashes of tokens since the semantic meaning of the sentences will be likely preserved under paraphrase and the watermark can remain robust. Comprehensive experiments are conducted to demonstrate the effectiveness and robustness of SemaMark under different paraphrases.

## 1 Introduction

Large language models (LLMs) have shown their great ability in various natural language processing (NLP) tasks like Question Answering (QA) (Lu et al., 2022), reasoning tasks (Wei et al., 2022; Creswell et al., 2022) and code development (Xu et al., 2022). However, tremendous concerns have been raised that LLMs are possible to be used improperly and illegally. For example, indistinguishable fake news are easy to be fabricated (Kreps et al., 2022; Zellers et al., 2019) by language models, which, when disseminated, could instigate widespread panic. Similarly, in the commercial sphere, convincingly generated reviews can manipulate consumer perceptions, leading to

unethical business competition (Salminen et al., 2022). Therefore, detecting LLM-generated text has become crucial in the real-world applications of LLMs.

Among diverse methods to detect LLM-generated texts, the watermark strategies have demonstrated outstanding precision (Kirchenbauer et al., 2023a). It is proposed to encode a secret watermark into the generated texts, such that we can tell whether a text is generated by detecting this watermark. One representative strategy (Kirchenbauer et al., 2023a; Yoo et al., 2023) is to encode the watermark based on the “partition of vocabulary”. In detail, given a language model, these methods devise a mapping from precedent tokens to a particular partition of the vocabulary by a partition function for the consequent token. The partition function leverages the hashes of the input as the seed of a random generator to split the vocabulary to a green list and a red list. During the text generation phase, the consequent token has an increased probability to be sampled from the green list. In this way, the watermark is encoded through the matching between the precedent tokens and the vocabulary partition for the consequent token. The detection is also facilitated by detecting this matching in generated contents. However, recent works (Krishna et al., 2023; Kaddour et al., 2023) reveal that **this watermark may be easily eliminated by sentence paraphrasing**. Individuals seeking to improperly utilize LLMs without being detected can paraphrase the generated contents, like altering the order and the choices of the words, and only retain the general meaning of the text to achieve their malicious goals like faking news. These paraphrases will change the seed of the partition function, i.e. the token hashes, and as we show in the Section 4.4, the partition function is sensitive to small changes. Consequently, the matching between the precedent tokens and the green list will be disrupted, and the detection ef-

083       fectiveness of the watermark can be dramatically  
084       compromised.

085       In this paper, we propose to leverage the seman-  
086       tic meaning of precedent token sequences as the  
087       seed for partition function, instead of simple hashes  
088       of precedent tokens, since the core semantic mean-  
089       ing is expected to be maintained after paraphrase.  
090       To achieve this goal, one key obstacle is how to  
091       capture the semantics when applying them for the  
092       partition function to watermark the generated texts.  
093       It is a common practice to quantify the semantics  
094       via embeddings (Reimers and Gurevych, 2019; Gao  
095       et al., 2021; Li et al., 2020; Giorgi et al., 2021). Em-  
096       beddings indeed can represent consistent semantics  
097       after paraphrase. Since the embeddings are high-  
098       dimensional vectors in the continuous space, they  
099       often present some minor changes after paraphrase.  
100       Although the main semantics are preserved, these  
101       minor changes can lead to a substantial difference  
102       in the partition of vocabulary because the random  
103       generator in the partition function is sensitive to  
104       the change of the seed, as shown in Section 4.4.

105       To overcome the above challenge, i.e., to make  
106       the quantified semantics invariant and make the wa-  
107       termark robust under paraphrase, we propose a new  
108       watermark method, SemaMark, which discretizes  
109       the continuous embedding space. Intuitively, the  
110       discretization can coarsen the representation of the  
111       embeddings which could tolerate the potential mi-  
112       nor changes caused by paraphrase. By proper dis-  
113       cretization, the paraphrased semantics could stay  
114       in the same discrete section with a high probability  
115       and the discretized quantified semantics will likely  
116       remain the same even after paraphrase. Therefore,  
117       the partition results will not change. However, di-  
118       rectly converting the high-dimensional embedding  
119       space into discrete is intricate and challenging. For  
120       example, discretizing each dimension will lead to a  
121       large amount of discrete values which is exponen-  
122       tial to the number of dimensions. Thus, the minor  
123       changes by paraphrase can still cause the change  
124       of discrete values because the number of discrete  
125       values are too dense and each discrete value can  
126       tolerate only small changes. Therefore, the mi-  
127       nor changes of high-dimensional embeddings can  
128       have a strong impact on the partition function. To  
129       address this problem, SemaMark first uses a Multi-  
130       Layer Perception (MLP) to condense the contin-  
131       uous high-dimensional embeddings into normal-  
132       ized vectors in 2D space. The vectors are located  
133       on a unit circle named Normalized Embedding  
134       Ring (NE-Ring). Then the condensed NE-Ring is

135       equally divided into various sections, transforming  
136       the continuous space into distinct discrete values,  
137       i.e., “semantic values”. Based on the discretiza-  
138       tion, SemaMark further introduces two strategies  
139       to advance the watermark’s concealment and to  
140       improve the robustness under paraphrase. *First*,  
141       SemaMark leverages the uniformity (Wang and  
142       Isola, 2020) of Contrastive Learning (CL) (Chen  
143       et al., 2020) to strength the MLP and mitigate the  
144       problem that the semantics are unevenly concentr-  
145       ating on some discrete sections on NE-Ring. The  
146       unevenly distribution will cause the final discrete  
147       semantic values overly monotonous. It raises the  
148       concern that the watermark might be cracked by  
149       counting token frequency (Zhao et al., 2023). *Sec-*  
150       *ond*, SemaMark utilizes an offset detection method  
151       to further enhance the robustness at the boundary of  
152       different discrete sections whose semantic values  
153       are possibly vulnerable to paraphrase. Comprehen-  
154       sive experiments are conducted to demonstrate the  
155       effectiveness and robustness of SemaMark under  
156       different paraphrases.

## 2 Related works 157

**LLM-generated detection.** 158 As the development  
159 of LLMs, various LLM-generated detection tools  
160 have also been proposed. Learning-based meth-  
161 ods train a classification model to detect the dif-  
162 ference between human-written text and machine-  
163 generated text like Guo et al. (2023); Wang et al.  
164 (2023); Li et al. (2023). Other works do not rely on  
165 the classification model, but try to use the property  
166 of the LLM to test whether a given text is generated  
167 by LLMs. For example, DetectGPT (Mitchell et al.,  
168 2023) assumes that the generated text will have  
169 high likelihood. GPT-who (Venkatraman et al.,  
170 2023) uses UID-based features to model the unique  
171 statistical signature of each LLM and human author  
172 for accurate authorship attribution. These methods  
173 do not interact the generation process of LLMs and  
174 thus have to explore unknown features of LLMs  
175 for detection. Instead, watermarks can change the  
176 model with a small but pre-defined rule which ac-  
177 celerates the detection process effectively.

**Watermark.** 178 The distinction between watermark  
179 and other methods is that watermark can proac-  
180 tively change the generation to insert a concealed  
181 watermark into the generated text. This gives  
182 clear difference between watermarked and non-  
183 watermarked texts. Watermark shifts the text using  
184 a small but pre-defined rule to make the detection

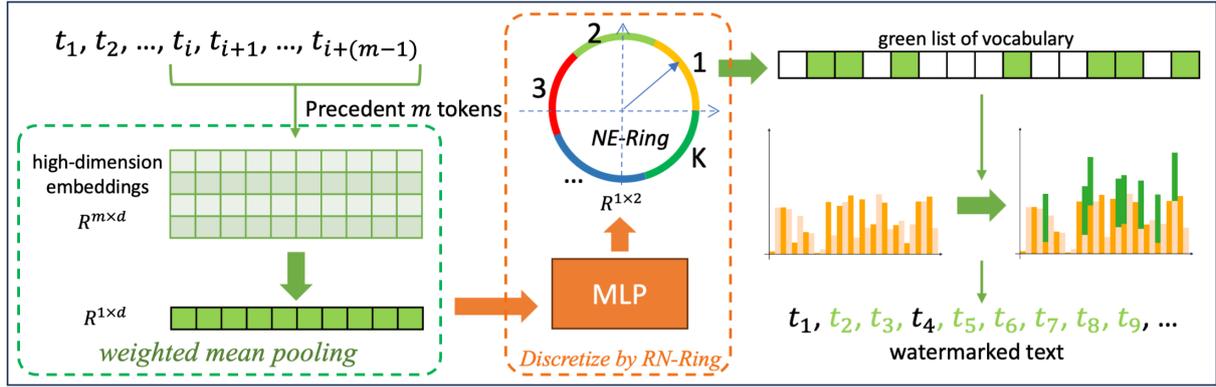


Figure 1: The watermarking process of SemaMark

much more effective. The partition of the vocabulary for each token is a representative watermark method (Kirchenbauer et al., 2023a; Yoo et al., 2023; Kirchenbauer et al., 2023b). In each autoregressive step of generating one token, the method uses the previous tokens’ hashes, to select a part of the vocabulary as “green” at a ratio of  $\gamma$ . Subsequently, they elevate the likelihood of the tokens by boosting the logits of the softmax by  $\delta$ . Through this approach, at each token position, the probability of this matching between the seed and green tokens tends to increase.

For a sentence with  $L$  tokens, it is viewed as a sample set of size  $L$ . Each token is one sample from the vocabulary. A non-watermarked sentence is expected to have  $\gamma L$  tokens showing this match, while the watermarked sentence is expected to have more. The watermark detection is approached as a  $z$ -test with null hypothesis that the text is non-watermarked. If the  $z$ -statistic is large, i.e. it is significantly different from the null hypothesis, the null hypothesis can be rejected and the text can be predicted as watermarked:

$$z = \frac{(G - \gamma L)}{\sqrt{L\gamma(1 - \gamma)}}, \quad (1)$$

where  $G$  is the number of tokens showing the matching between seed and the green list. Yoo et al. (2023) further expand this watermark of green and red list to more lists for multi-bit encoding.

### 3 Method

In this section, we introduce the detailed design of SemaMark. We first present how to use the semantic information as the seed for watermark methods that are based on random partition of vocabulary in Section 3.1. Then in Section 3.2 and

Section 3.3, we introduce the CL training scheme and the smoothed detection method for further improving the robustness, respectively.

#### 3.1 The framework of SemaMark

As aforementioned, the existing watermark methods based on partition of vocabulary are susceptible to paraphrase. Paraphrase can easily change the previous tokens and disrupt the matching between tokens and the partition of vocabulary, without significantly affecting the semantic meaning. Thus, SemaMark uses the invariant semantics for watermarking by discretizing the embedding space to accommodate the minor perturbation of semantics and provide a stable mapping between semantics and vocabulary partition for the consequent token.

However, discretization in a high-dimension space is intricate and non-trivial. Therefore, we first reduce the high-dimensional embedding space onto the 2D NE-Ring and then discretize via NE-Ring. The whole watermarking process is shown in Figure 1. SemaMark first reduces the dimension of the embedding space to obtain the discrete semantic values by two steps, i.e., *weighted embedding pooling* and *discretizing by NE-Ring*, and then uses the semantic value to partition the vocabulary. The logits of green list is shifted to increase the probability of matching between semantics and the consequent token for watermarking the LLM,  $f$ . In the following, we introduce more details about the two steps to obtain a stable semantic value.

**S1: weighted embedding pooling.** To enhance the robustness, we aggregate the semantics of previous  $m$  tokens by the weighted mean pooling function  $P(\cdot)$  before dimension reduction, instead of using only one preceding token’s embedding. In the ablation studies of Section 4.4, we show that

the method has the best performance when  $m$  is neither too big nor too small. For the token sequence  $\{t_{i:i+m-1}\}$  starting at position  $i$ , we use their semantics to generate the token in the  $m$  position,  $t_{i+m}$ . We denote their embeddings as  $\{e_{i:i+m-1}\}$ .  $\{e_{i:i+m-1}\}$  can be easily obtained from the LLM,  $f$ , that we want to watermark. Intuitively, in  $\{t_{i:i+m-1}\}$ , the embeddings of tokens far from  $t_{i+m}$  contain semantic information that is more distant from  $t_{i+m}$  than the closer ones. It means that the embeddings of far tokens may be less robust because the distant connection is more possible to change after paraphrase. Thus, for pooling,  $P(\cdot)$  aims to assign a smaller weight to the token far from  $t_{i+m}$  and a larger weight to the token closer to  $t_{i+m}$  as:

$$P(\{e_{i+1:i+m}\}) = \sum_{j=1}^K \frac{j + \frac{K}{2}}{w_{\text{sum}}} e_{i+j},$$

where  $w_{\text{sum}} = K^2 + K/2$  is the sum of all weights. We denote the weighted output  $P(\{e_{i:i+m-1}\}) \in \mathbb{R}^d$  as  $e_{P_{i,m}}$  for short. By pooling, more semantics are used for a seed, which enhance the robustness under paraphrase.

**S2: discretizing by NE-Ring.** After aggregating the embeddings by weighted pooling, SemaMark uses MLP  $g_\theta$  to transform  $e_{P_{i,m}}$  to a normalized vector in 2D embedding space. The normalized vectors locate on a unit circle in the 2D space, which is named as Normalized Embedding Ring (NE-Ring). The discretization function,  $D(\cdot)$ , discretizes NE-Ring by equally segmenting into different sections. It takes the polar angle  $\phi$  of  $g_\theta(e_{P_{i,m}})$  as input and outputs the discretized semantic values  $a \in [K]$ , where  $[K] := \{1, 2, \dots, K\}$ .  $D(\cdot)$  is defined as

$$D(\phi) = \left\lfloor \phi \frac{K}{2\pi} \right\rfloor$$

It first maps the input from  $[0, 2\pi)$  to  $[0, K)$ , and then discretizes all the values in  $[i, i + 1)$  to  $i$ , for  $\forall i \in [K - 1]$ . Even though there could be subtle changes in semantics by paraphrase, the paraphrased  $\tilde{a}$  will likely locate in the discrete section  $[i, i + 1)$ . Some tokens may still have  $a \neq \tilde{a}$  if the normalized vector is close to the boundary of  $[i, i + 1)$ . Therefore, in Section 3.3, we introduce an offset detection to strengthen the tolerance for this mismatch and correct some unstable cases.

With the two steps, we can get a stable discrete semantic value as the seed for the partition function

to partition the vocabulary for the consequent token. Following Kirchenbauer et al. (2023a), the vocabulary is partitioned into green and red lists. We increase the logits of the tokens in the green list by  $\delta$  and recalculate the probability distribution based on the shifted logits. For each token to generate, we increase the possibility of the green list based on its previous  $m$  tokens' semantics. Thus, all the generated tokens will be likely to have this matching between the semantics and the consequent green token. By detecting the matching, we can discriminate whether a text is watermarked or not and then detect the LLM-generated contents effectively. Besides, SemaMark proposes two strategies to reduce the risk of being cracked by Contrastive Learning and further increase the robustness by the offset detection in the following sections.

### 3.2 Training $g_\theta$ by Contrastive Learning

The MLP is expected to produce a uniform distribution of  $g_\theta(e_{P_{i,m}})$  on NE-Ring. If different semantics unevenly distributed on NE-Ring, the resulting discrete semantic values will be overly monotonous and the green list is more changeless. Consequently, the green list might be revealed by counting the token frequency, which compromises the concealment of watermark and leads to the risk of being cracked. Ideally, SemaMark should generate a wider variety of semantic values for different sentences, while each semantic value is robust and stable if its corresponding sentence is paraphrased. To achieve this goal, we propose to use Contrastive Learning to train MLP since Contrastive Learning has the property of uniformity that the data will be evenly distributed in the whole feature space (Wang and Isola, 2020). The uniform distribution can help the normalized vectors cover all the semantic values. As a result, NE-Ring can generate a wider variety of semantic values to prevent the watermark from being cracked.

In Contrastive Learning, we first input the sentences into the model  $f$  to get a batch of sequences of  $m$  tokens and their pooling embeddings  $e_{P_{i,m}}$ , denoted as  $\{e_j\}$ , where  $j \in [B]$  and  $B$  is the batch size. To compose a contrastive loss, we construct the positive and negative pairs by a soft augmentation:

$$e_{j+B} = e_j^+ = e_j + \epsilon,$$

where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is a Gaussian noise. The soft augmentation can simplify the construction

of positive samples. With this soft augmentation, we can assign the samples sharing similar embeddings from the same sequence as positive pairs and samples from different sequences as negative pairs. This is consistent with our intuition that the paraphrased semantic embeddings will not change significantly and can remain robust. Then the contrastive loss is

$$L_j = -\log \frac{\exp(\text{sim}(g_\theta(e_j), g_\theta(e_j^+)) / \tau)}{\sum_{k \neq j, k \in [2B]} \exp(\text{sim}(g_\theta(e_j), g_\theta(e_k)) / \tau)},$$

where  $\text{sim}(\cdot)$  is cosine similarity and  $\tau$  is the temperature. By Contrastive Learning, the output of reduced semantic embeddings can be evenly distributed in all of the space on NE-Ring, and cover all the discrete sections to improve the robustness of SemaMark.

### 3.3 $Q$ -offset detection

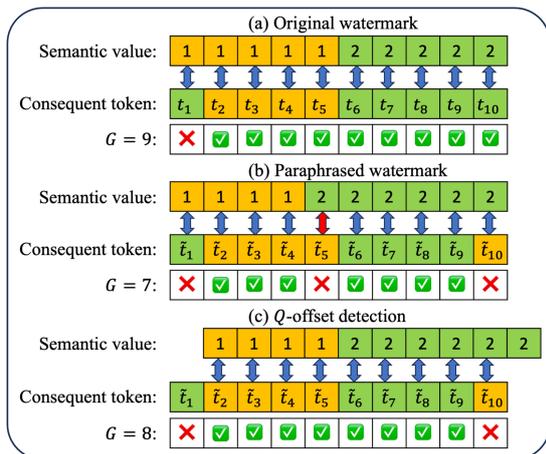


Figure 2:  $Q$ -offset detection vs. existing detection

Existing detection methods check the matching between partition seed and the consequent tokens in a one-to-one manner as shown in Figure 2(a). The detection method first recalculates the seed for each token position and gets the partition of the green list, and then checks whether the consequent token is in the partitioned green list token by token. In SemaMark, this strategy can be effective when the text is not paraphrased. However, after paraphrase, this detection could be suboptimal because the semantic values of some sequences which are close to the boundaries of the discrete section  $[i, i + 1)$  might change as shown in Figure 2(b). This is because the window of  $m$  tokens will slide token by token during the auto-regressive generation process, and the semantic change will also accumulate when the window is sliding. The semantic values closed

to the boundary usually happen when the change accumulates to some extent. This change of boundary semantic values will lead to some mismatch and reduce the accuracy like  $\tilde{t}_5$  in Figure 2(b).

To mitigate the influence of this error, we propose  $Q$ -offset detection. As shown in Figure 2(c), we offset the discrete seed by  $q$  tokens to detect the matching between semantics and the consequent tokens, where  $q \in \{-Q, -(Q-1), \dots, 0, 1, \dots, Q\}$  and the sign of  $q$  indicates the direction of the offset. We choose the maximal  $z$ -statistic in different  $q$  as the  $Q$ -offset score. However,  $Q$ -offset detection will also increase the  $Q$ -offset score of non-watermark text, which indicates that the detected green word fraction  $\gamma$  of non-watermark text is higher. The  $\gamma$  in Eq. (1) is possibly inaccurate. Thus during generation, we set  $\gamma$  to a fixed value, while in detection process, we treat  $\gamma$  as a hyperparameter and use a validation set to determine its value in practice. In Section 4.4, we discuss the ablation studies of  $Q$ -offset and  $\gamma$  and show that  $Q$ -offset can impressively improve the detection performance with robustness.

## 4 Experiment

In this section, we conduct experiments to demonstrate the robustness of SemaMark. In Section 4.2, we demonstrate that its robustness is better than the baseline methods. In Section 4.3, we show that our watermark has almost no influence on the quality of generated texts. In Section 4.4, we use ablation studies to demonstrate the effectiveness of partition function and  $Q$ -offset detection, and show the sensitivity of the partition function. In Section 4.5 we visualize the distribution of NE-Ring and provide analysis on the feature distribution of Contractive Learning.

### 4.1 Experiment setups

**Backbone models and datasets.** We test our method by watermarking two models, OPT-2.7B and OPT-6.7B (Zhang et al., 2022) which are referred to as the backbone models in following sections. For dataset, we use the news-like subset of C4 (Raffel et al., 2020), which covers a variety of topics. From the news-like subset of C4, we extract a training set, a validation set and a test set. For each sample, we use the first half of text as prompt to generate watermark sentences. More details can be found in Appendix A.

**Baseline methods.** We compare our method

Paraphrase		ROC-AUC				F1 with best threshold			
		LeftHash	SelfHash	EXP-Edit	ours	LeftHash	SelfHash	EXP-Edit	ours
OPT-2.7B	No paraphrase	0.9913	0.9886	0.9799	0.9948	0.9921	0.9861	0.9708	0.9905
	Translation	0.9091	0.8147	0.8749	<b>0.9692</b>	0.8456	0.7622	0.8157	<b>0.9330</b>
	Dipper	0.9878	0.9728	0.9736	<b>0.9911</b>	<b>0.9727</b>	0.9400	0.9620	0.9701
	GPT-3.5	0.9028	0.7908	0.9392	<b>0.9406</b>	0.8358	0.7378	0.8852	<b>0.8902</b>
OPT-6.7B	No paraphrase	0.9918	0.9930	0.9784	0.9949	0.9911	0.9863	0.9705	0.9858
	Translation	0.8807	0.8098	0.8625	<b>0.9308</b>	0.8129	0.7468	0.8013	<b>0.8882</b>
	Dipper	<b>0.9904</b>	0.9747	0.9728	0.9871	0.9786	0.9432	0.9620	<b>0.9821</b>
	GPT-3.5	0.8990	0.7909	0.8996	<b>0.9377</b>	0.8300	0.7367	0.8354	<b>0.8766</b>

Table 1: Watermark detection results under three paraphrases. (The best performance under paraphrase is bolded.)

with three baselines LeftHash, SelfHash (Kirchenbauer et al., 2023b) and EXP-Edit (Kuditipudi et al., 2023). LeftHash and SelfHash are two methods based on the partition of vocabulary using the hashes of tokens. EXP-Edit uses a private sequence to encode the watermark by changing the probability distribution of the sequence of tokens. More details on the implementation can be found in Appendix A.

**Paraphrase setups.** We use three representative methods to paraphrase the watermarked text, round-trip translation (Tiedemann and Thottingal, 2020), Dipper (Krishna et al., 2023) and GPT-3.5. For round-trip translation, we first translate from English to another language and then transform back to English, such that some words and expressions will be changed because the translation is not an one-to-one mapping. For Dipper, we follow the parameter setting in Kirchenbauer et al. (2023b). For GPT-3.5, we use the prompt in Kirchenbauer et al. (2023b) to query GPT-3.5 for paraphrase.

**Evaluation metrics and hyper-parameters.** We use F1 score with best threshold and ROC-AUC to measure the performance of the watermark detection. All the metrics are calculated based on at least 500 watermarked samples and 500 non-watermark samples. The length of watermarked samples before paraphrase and non-watermark samples is  $200 \pm 25$ . In generation, we set  $\gamma = 1/4$  for LeftHash, SelfHash and SemaMark. In detection, we set  $\gamma = 1/3$  and  $\delta = 2$  based on the validation set in Section 4.4(b). In SemaMark, we set  $m = 15$ ,  $Q = 15$ ,  $K = 5$  for OPT-2.7B and  $K = 4$  for OPT-6.7B.

## 4.2 Main Results

In this subsection, we demonstrate the robustness of the proposed SemaMark under paraphrase by comparing it with three baseline methods on two backbone models. We first generate watermarked

texts and use three paraphrase methods to remove the watermarks. The detection performance of both texts with and without paraphrase is reported in Table 1. As we can see, before paraphrase, all the watermarked methods have good detection performance. After paraphrase, SemaMark has the best detection performance most of the time across all the backbone models and all the paraphrase methods, which suggests that our method is more robust against paraphrase.

In detail, by round-trip translation, the paraphrase reduces the detection ability of baseline methods effectively, while the watermark of SemaMark is robust. Under round-trip translation, the best ROC-AUC of baselines is 0.9091 on OPT-2.7B and 0.8807 on OPT-6.7B, respectively. But ROC-AUC of SemaMark is 0.9692 and 0.9308, which is at least 0.05 higher than all the baseline methods. Similarly, under paraphrase of GPT-3.5, SemaMark is better than all the baselines. The best baseline performance under GPT-3.5 is 0.9392 in ROC-AUC on OPT-2.7B and 0.8990 in ROC-AUC on opt-6.7B, but SemaMark has higher AUC-ROC of 0.9406 and 0.9377. For Dipper, we note that all methods are robust to Dipper since it does not significantly reduce the detection performance. However, SemaMark is still one of the most robust. On OPT-2.7B, it performs best in ROC-AUC, while on OPT-6.7B, it has the best F1 score. From Table 1, the results show an obvious improvement of SemaMark in robustness. This implies that using semantics as the seed for the partition function is effective under paraphrase.

## 4.3 Text Quality

Watermark should not compromise the generation quality of LLMs. In this subsection, we compare the text quality by calculating perplexity and demonstrate that our watermark has almost no influence on the generated quality. Perplexity measures

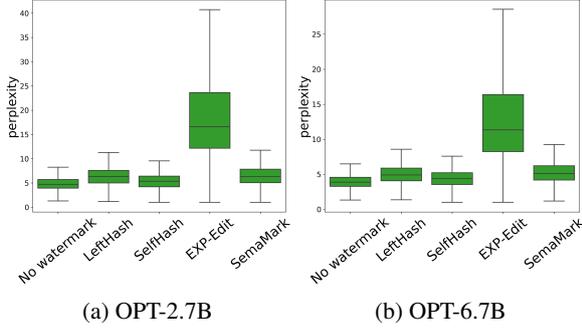


Figure 3: Text quality (perplexity)

the likelihood that a sentence is generated by one model. Lower perplexity means the watermarked text is more predictable. In other words, it is more consistent with the reasoning of the given model. In Figure 3, we use OPT-6.7B with no watermark to get perplexity for all the watermarked methods. All the results in Figure 3 are calculated without paraphrase, because the generation quality of text is not related to paraphrase. From Figure 3a on OPT-2.7B, we can see that our watermark, LeftHash and SelfHash have almost no influence on the generation quality. They has perplexity at around 6 which is similar as the generated text without watermark. Instead, EXP-Edit has much higher perplexity, which means that EXP-Edit changes the generated text in an aggressive way and much reduces the generation quality after watermarking. This is probably because EXP-Edit adjusts the logits on the whole vocabulary. From Figure 3b, we can draw similar conclusions for OPT-6.7B. EXP-Exit also increases the perplexity by around 10, while the average perplexity of LeftHash, SelfHash and ours is around 1 higher than the non-watermarked generated text. In summary, our SemaMark can keep the quality and robustness simultaneously.

#### 4.4 Ablation Study

In this subsection, we study the influence of the length of the sequence we use for generating one semantic value and the sensitivity of the partition function.

##### a) Length of previous sequence tokens, $m$

In the first step of SemaMark, i.e., *weighted embedding pooling*, we use the semantic of the previous  $m$  tokens to get the more stable embedding. But if the length of the sequence is too long, it will also hurt the robustness. In Figure 4, we test watermark on OPT-2.7B with different  $m$  and draw the ROC-AUC. The results show that before  $m = 15$ , ROC-AUC is in the trend of increase as the  $m$

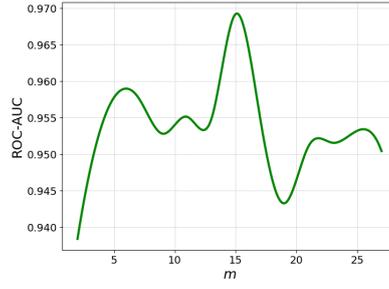
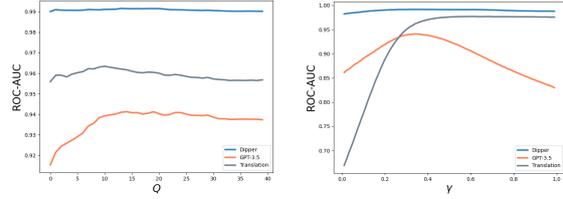


Figure 4: ROC-AUC and  $m$



(a) ROC-AUC and offset  $Q$  (b) ROC-AUC and  $\gamma$

Figure 5: Text quality (perplexity)

changes. But when  $m > 15$ , ROC-AUC becomes fluctuating. It is possibly because that the distant tokens will include more change after paraphrase as we mentioned in Section 3.1. Another possible reason is that in the beginning of generation for the first  $m$  tokens, the number of previous tokens is smaller than  $m$  and NE-Ring can only use the embeddings of limited tokens for prediction, which may be unstable. Thus, too long or too short sequence will hurt the robustness of SemaMark against paraphrase. In our experiments, we choose  $m = 15$  for all the settings.

##### b) $Q$ -offset detection

In this subsection, we show that the effectiveness of the proposed  $Q$ -offset detection. In Figure 5a, we demonstrate the change of ROC-AUC of SemaMark with different  $Q$  in offset detection under three different paraphrases.  $Q$ -offset detection searches the highest  $z$ -statistics from  $-Q$  to  $Q$  as the  $Q$ -offset score. From Figure 5a, we can see that when  $Q$  increases, ROC-AUC first increases and decreases after  $Q$  is around 15. When  $Q < 15$ , the offset can help correct the errors of semantic values close to the boundary. Compared with detection without offset, i.e.  $Q = 0$ , ROC-AUC of SemaMark is much better, which means that the offset can help to solve the errors of semantic values around the boundaries that are more vulnerable to paraphrases. When  $Q > 15$ , the correction of this error is limited, because the offset will also increase the  $Q$ -offset score of negative samples as it also searches the highest  $z$ -statistics of negative samples. On the other hand, the computation cost

will also increase if  $Q$  is too large because it has to search more possible  $q$ . In practice, we set  $Q = 15$  in all the experiments, which can effectively reduce the influence of the errors of semantic values at the boundaries.

Since the  $Q$ -offset detection searches the highest green word fraction, the fraction of green list word of non-watermarked text will be higher than the  $\gamma$  that we used to randomly select the green list. Thus, it is not accurate to use the original  $\gamma$  for  $z$ -statistics. We treat  $\gamma$  as a hyper-parameter and use a validation set to select its value. As shown in Figure 5b, the detection performance of SemaMark under paraphrases of Dipper and GPT-3.5 will reach the highest when  $\gamma$  is around  $1/3$ , while it will continue to increase under paraphrase. In practice, we set  $\gamma = 1/3$  for  $Q$ -offset detection.

### c) Sensitivity of partition function.

As we mentioned, the partition function is sensitive to any change of the input because it only uses the input as the seed of the random generator. To validate its sensitivity to continuous embeddings, we adopt the embedding vector as the input to show that, with tiny change of the embeddings, the partition of vocabulary can be very different. We propose a hash method based on md5sum (Deepakumar et al., 2001) to adopt the partition function by transforming the continuous embeddings to an integral seed. We use 1000 sequences to test the sensitivity. For each sequence embedding, we first get a green list from the partition function. Then we change one dimension of the embedding by only  $1e-5$  to get a new partition result. The overlapping of the green list before and after changing is 24.99% on the average of 1000 sequences. It is consistent with  $\gamma$  we use to watermark, because the random partition with the changed embedding is independent from the original one. It means the partition function is sensitive to any small change in its input. Instead, after we use NE-Ring to discretize the embeddings, the overlapping of green list after changing embeddings by  $1e-5$  is 100%, which means the discretization can effectively handle this change. In practice, SemaMark can provide the tolerance that is much larger than  $1e-5$ , which makes the watermark more robust under paraphrase. With the improvement of  $Q$ -offset, the detection of SemaMark is more robust and effective.

## 4.5 Distribution on NE-Ring based on CL

In this subsection, we demonstrate that Contrastive Learning can help evenly distribute the semantics

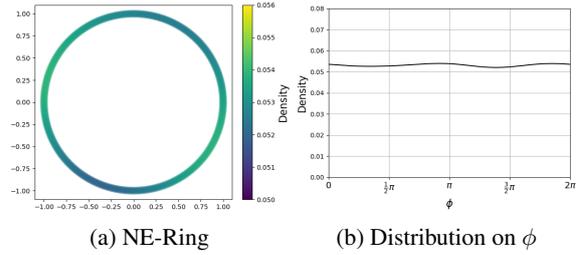


Figure 6: Visualization of NE-Ring

on the NE-Ring. The even distribution can help the sequences reach all possible semantic values and provide more diverse semantic values to prevent the watermark from being cracked by counting token frequency. In Figure 6a, we use Gaussian density estimation (Chen, 2017) to get the distribution of the semantics on the NE-Ring **before** discretization. We use different colors to show the density. The NE-Ring in Figure 6a shows that, the distribution is uniform. All the density is between 0.052 and 0.054. We further plot the density based on the polar angle  $\phi$  in Figure 6b where the density has almost no change on all the polar angle from 0 to  $2\pi$ . This implies that the training based on Contrastive Learning can ensure the semantics will reach all possible discrete values. It can prevent the case where the discrete values will gather in some discrete sections and produce monotonous vocabulary partitions. As a result, it can protect the watermark from being cracked by counting token frequency.

## 5 Conclusion

In this paper, we use the semantic information for watermarking to enhance the robustness against paraphrase. The existing watermark methods use the matching between the previous tokens and the partition vocabulary. This matching can be easily broken by paraphrase. However, we construct the mapping between the semantics and the vocabulary. In this way, the semantics will stay stable under paraphrase and the robustness of watermark can be enhanced. To make use of semantics, we propose SemaMark to discretize the embedding space on NE-Ring and propose a training method based on CL. In addition, we use  $Q$ -offset detection to further advance the robustness by increasing the tolerance of the semantic values close to the discrete boundary. In experiments, we demonstrate our method can perform much better compared with baseline methods under paraphrase while having little influence on the generation quality.

## 6 Limitations

In some cases, the customers may rely on some API-based LLMs and do not have the access to the embeddings and the permission to modify the logits during generation. Although our watermark method can effectively detect the LLM-generated content and increase the detection success rate under paraphrase, it is not applicable for black-box LLMs. The second weakness of our method is that the NE-Ring is dependent on the semantic embedding of LLMs. For each LLM, we need to train a specialized EN-Ring, which might be inflexible if we want to produce a general model for NE-Ring or fine-tune the LLMs. Despite the weaknesses, our method is successful in the problem of robustness under paraphrase. In the future work, we will continue to extend our method into black-box LLMs and a universal model that does not require customized training for various specific LLMs.

**Potential risk.** Our discussion about the robustness might provide motivation for the attackers to find other methods like adaptive attack. Although we provide robustness under paraphrase, if the unauthorized people propose possible attack method focusing on the green-list based watermark from other perspectives, the detection rate for LLM-generated texts are still possible to be reduced.

## References

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc."

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.

Yen-Chi Chen. 2017. A tutorial on kernel density estimation and recent advances. *Biostatistics & Epidemiology*, 1(1):161–187.

Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*.

Janaka Deepakumara, Howard M Heys, and R Venkatesan. 2001. Fpga implementation of md5 hash algorithm. In *Canadian Conference on Electrical and Computer Engineering 2001. Conference Proceedings (Cat. No. 01TH8555)*, volume 2, pages 919–924. IEEE.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. In *2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021*, pages 6894–6910. Association for Computational Linguistics (ACL).

John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. 2021. Declutr: Deep contrastive learning for unsupervised textual representations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 879–895.

Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *arXiv preprint arXiv:2301.07597*.

Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. *Array programming with NumPy*. *Nature*, 585(7825):357–362.

Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and applications of large language models. *arXiv preprint arXiv:2307.10169*.

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023a. A watermark for large language models. *arXiv preprint arXiv:2301.10226*.

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Manli Shu, Khalid Saifullah, Kezhi Kong, Kasun Fernando, Aniruddha Saha, Micah Goldblum, and Tom Goldstein. 2023b. On the reliability of watermarks for large language models. *arXiv preprint arXiv:2306.04634*.

Sarah Kreps, R Miles McCain, and Miles Brundage. 2022. All the news that's fit to fabricate: Ai-generated text as a tool of media misinformation. *Journal of experimental political science*, 9(1):104–117.

Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2023. Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. *arXiv preprint arXiv:2303.13408*.

Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. 2023. Robust distortion-free watermarks for language models. *arXiv preprint arXiv:2307.15593*.

782	Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	835
783	Yiming Yang, and Lei Li. 2020. On the sentence	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	836
784	embeddings from pre-trained language models. In	et al. 2022. Chain-of-thought prompting elicits rea-	837
785	<i>Proceedings of the 2020 Conference on Empirical</i>	soning in large language models. <i>Advances in Neural</i>	838
786	<i>Methods in Natural Language Processing (EMNLP)</i> ,	<i>Information Processing Systems</i> , 35:24824–24837.	839
787	pages 9119–9130.		
788	Linyang Li, Pengyu Wang, Ke Ren, Tianxiang Sun, and	Frank F Xu, Uri Alon, Graham Neubig, and Vincent Jo-	840
789	Xipeng Qiu. 2023. Origin tracing and detecting of	sua Hellendoorn. 2022. A systematic evaluation of	841
790	llms. <i>arXiv preprint arXiv:2304.14072</i> .	large language models of code. In <i>Proceedings of</i>	842
791		<i>the 6th ACM SIGPLAN International Symposium on</i>	843
792	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-	<i>Machine Programming</i> , pages 1–10.	844
793	Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter		
794	Clark, and Ashwin Kalyan. 2022. Learn to explain:	KiYoon Yoo, Wonhyuk Ahn, Jiho Jang, and Nojun	845
795	Multimodal reasoning via thought chains for science	Kwak. 2023. Robust multi-bit natural language wa-	846
796	question answering. <i>Advances in Neural Information</i>	termarking through invariant features. In <i>Proceed-</i>	847
	<i>Processing Systems</i> , 35:2507–2521.	<i>ings of the 61st Annual Meeting of the Association for</i>	848
		<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	849
		pages 2092–2115.	850
797	Eric Mitchell, Yoonho Lee, Alexander Khazatsky,	Rowan Zellers, Ari Holtzman, Hannah Rashkin,	851
798	Christopher D Manning, and Chelsea Finn. 2023.	Yonatan Bisk, Ali Farhadi, Franziska Roesner, and	852
799	Detectgpt: Zero-shot machine-generated text detec-	Yejin Choi. 2019. Defending against neural fake	853
800	tion using probability curvature. <i>arXiv preprint</i>	news. <i>Advances in neural information processing</i>	854
801	<i>arXiv:2301.11305</i> .	<i>systems</i> , 32.	855
802	Colin Raffel, Noam Shazeer, Adam Roberts, Kather-	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	856
803	ine Lee, Sharan Narang, Michael Matena, Yanqi	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	857
804	Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the</a>	wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.	858
805	<a href="#">limits of transfer learning with a unified text-to-text</a>	Opt: Open pre-trained transformer language models.	859
806	<a href="#">transformer</a> . <i>Journal of Machine Learning Research</i> ,	<i>arXiv preprint arXiv:2205.01068</i> .	860
807	21(140):1–67.		
808	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	Xuandong Zhao, Prabhanjan Ananth, Lei Li, and	861
809	Sentence embeddings using siamese bert-networks.	Yu-Xiang Wang. 2023. Provable robust water-	862
810	In <i>Proceedings of the 2019 Conference on Empirical</i>	marking for ai-generated text. <i>arXiv preprint</i>	863
811	<i>Methods in Natural Language Processing and the 9th</i>	<i>arXiv:2306.17439</i> .	864
812	<i>International Joint Conference on Natural Language</i>		
813	<i>Processing (EMNLP-IJCNLP)</i> , pages 3982–3992.		
814	Joni Salminen, Chandrashekhara Kandpal, Ahmed Mo-		
815	hamed Kamel, Soon-gyo Jung, and Bernard J Jansen.		
816	2022. Creating and detecting fake reviews of on-		
817	line products. <i>Journal of Retailing and Consumer</i>		
818	<i>Services</i> , 64:102771.		
819	Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-		
820	MT — Building open translation services for the		
821	World. In <i>Proceedings of the 22nd Annual Confer-</i>		
822	<i>ence of the European Association for Machine Trans-</i>		
823	<i>lation (EAMT)</i> , Lisbon, Portugal.		
824	Saranya Venkatraman, Adaku Uchendu, and Dongwon		
825	Lee. 2023. <a href="#">Gpt-who: An information density-based</a>		
826	<a href="#">machine-generated text detector</a> .		
827	Pengyu Wang, Linyang Li, Ke Ren, Botian Jiang, Dong		
828	Zhang, and Xipeng Qiu. 2023. <a href="#">Seqxgpt: Sentence-</a>		
829	<a href="#">level ai-generated text detection</a> .		
830	Tongzhou Wang and Phillip Isola. 2020. Understanding		
831	contrastive representation learning through alignment		
832	and uniformity on the hypersphere. In <i>International</i>		
833	<i>Conference on Machine Learning</i> , pages 9929–9939.		
834	PMLR.		

## 865 **A More details on experimental settings**

866 All the baseline models, backbone models and  
867 datasets we use are open source and available  
868 for academic purpose. For backbone models,  
869 we use the open-sourced model from Hugging-  
870 face<sup>1</sup>. The implementation is based on Pytorch<sup>2</sup>  
871 framework and also depend on packages includ-  
872 ing NLTK (Bird et al., 2009) and Numpy (Har-  
873 ris et al., 2020). For baseline methods, we use  
874 the released official code from the authors. For  
875 paraphrase models, we use OPUS-MT translation  
876 model and Dipper on Huggingface repository<sup>3</sup>, and  
877 API of ChatGPT<sup>4</sup>.

---

<sup>1</sup><https://huggingface.co/facebook/opt-2.7b>

<sup>2</sup><https://pytorch.org/>

<sup>3</sup><https://huggingface.co/Helsinki-NLP/opus-mt-en-zh> and  
<https://huggingface.co/kalpeshk2011/dipper-paraphraser-xxl>

<sup>4</sup><https://chat.openai.com/chat>