

# SIGNED IN INK, HIDDEN IN NOISE: WATERMARKING DIFFUSION LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Watermarking techniques enable the embedding of imperceptible signals into autoregressive large language models (LLMs), facilitating reliable detection and attribution of AI-generated text. However, for the emerging paradigm of diffusion large language models (dLLMs), which provide bidirectional context modeling, greater generation flexibility and controllability, and more efficient sampling, research on watermarking remains largely unexplored, raising critical concerns about copyright protection. In this work, we present the first systematic investigation of watermarking for dLLMs and introduce Ripple, a dedicated framework specifically designed for the diffusion-based generation process. The Ripple operates in two complementary stages, namely watermark injection and watermark calibration, to achieve seamless integration of watermark signals within dLLMs. Furthermore, extensive experiments on two representative dLLMs across multiple datasets demonstrate that Ripple effectively balances detectability, robustness, and text quality. We believe this study provides a foundation for advancing future research on watermarking in dLLMs.

## 1 INTRODUCTION

As generative language models increasingly approximate human writing, the concept of authorship faces unprecedented challenges. When any given sentence may originate from a machine, accurately identifying its true source becomes crucial for ensuring accountability, provenance, and trust. To address this issue, watermarking techniques have been developed to embed imperceptible signals into the outputs of autoregressive language models, enabling reliable detection and attribution. These approaches have shown promising results on mainstream autoregressive architectures. However, with the emergence of diffusion-based large language models (Nie et al., 2025b; Cheng et al., 2025; Ye et al., 2025; Labs et al., 2025), which offer smoother, more diverse, and more controllable text generation, the study of watermarking in this new paradigm remains largely unexplored.

Diffusion Large language models (dLLMs) generate text through an iterative denoising process, offering inherent advantages such as bidirectional context modeling, enhanced flexibility and controllability, and theoretically greater efficiency than autoregressive models. However, the same mechanism that underpins these strengths also poses unique challenges for watermarking. In diffusion-based generation, the timing of watermark injection is crucial: signals inserted too early or too late are prone to being overwritten by subsequent denoising steps, which severely hinders reliable detection. Existing watermarking techniques, designed for sequence-level token generation, cannot be directly applied because dLLMs introduce a temporal dimension on top of the sequential structure of text. Consequently, watermarking in dLLMs shifts from a two-dimensional sequence embedding problem into a three-dimensional one that spans both sequence and diffusion steps, demanding new mechanisms that remain simultaneously detectable and resilient under this generative paradigm.

To address this need, we present the first watermarking framework specifically designed for dLLMs. In contrast to traditional watermarking techniques designed for autoregressive LLMs, our approach operates directly within the diffusion process. Rather than injecting signals through abrupt manipulation of token-level outputs, it progressively embeds watermark information that co-evolving with the generative trajectory across denoising steps, as illustrated in Figure 1. This design aligns with the core principle of diffusion models, in which semantic structure gradually emerges from noise, thereby integrating watermark signals in a distributed and generation-aligned manner.

Specifically, we propose **Ripple**, a watermarking framework seamlessly integrated into the iterative generation workflow of dLLMs. Beginning with a fully masked sequence, dLLMs progressively refine text through multiple diffusion steps. Ripple exploits this process through two synergistic stages. In the **watermark injection stage**, we reorder the remasking operation to embed watermark signals into the tokens actively unmasked at each step, allowing the watermark to co-evolve with the emerging sequence. This design prevents the embedded watermark from dissipating during subsequent diffusion steps and avoids redundant embedding in the same token. In the **watermark calibration stage**, once the final step is reached or all tokens have been generated, a detector assesses the strength of the embedded watermark while a quality-aware evaluator monitors fluency. Tokens that are weakly watermarked or exhibit low quality are corrected through targeted resampling, ensuring that the final output remains both high-quality and strongly traceable. By embedding watermarking into the temporal-sequential dynamics of dLLMs, Ripple transforms watermarking from a static overlay into a generation-aware approach, yielding secure and reliably detectable text in diffusion large language models.

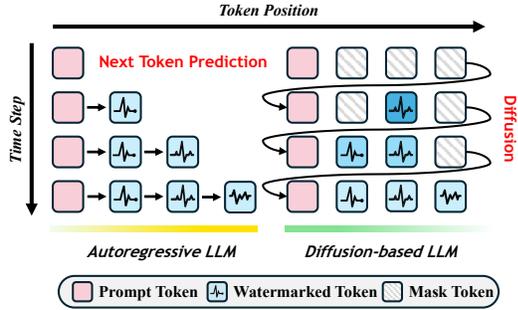


Figure 1: In autoregressive LLMs, tokens are generated sequentially left to right, whereas in diffusion-based LLMs, tokens are progressively denoised through a series of diffusion steps.

We conduct extensive experiments on mainstream dLLMs as well as a wide range of downstream task benchmarks. The results consistently demonstrate the superiority of the Ripple. Beyond performance, our approach offers a flexible mechanism that dynamically adjusts watermark strength and text fluency across diverse application scenarios, paving the way for practical and trustworthy watermarking solutions in the era of dLLMs. Our key contributions are summarized as follows:

- We introduce Ripple, the first exploration of a dedicated watermarking framework for dLLMs, which harmoniously integrates watermarking into the diffusion-based generation process.
- We carry out comprehensive evaluations across diverse datasets and model architectures, showing that Ripple consistently delivers superior performance under various challenging conditions.

## 2 PRELIMINARIES

### 2.1 DIFFUSION LARGE LANGUAGE MODELS (dLLMs)

**Notations.** Given a dLLM  $\Theta$ , let  $\mathcal{V}$  denote the vocabulary,  $|\mathcal{V}|$  denote the vocabulary size,  $L$  denote the sequence length,  $T$  denote the total diffusion steps, and  $p_0$  denote the input prompt. We use  $\mathcal{M}_t$  and  $\mathcal{U}_t$  to represent the sets of masked and unmasked positions at diffusion step  $t$ , respectively. The generated sequence at diffusion step  $t$  is denoted by  $\mathbf{y}_t = (y_t^1, \dots, y_t^L)$ .

**Inference.** During inference, at  $t = 0$ , the unmasked set is empty ( $\mathcal{U}_0 = \emptyset$ ) and the masked set covers all positions ( $\mathcal{M}_0 = \{1, \dots, L\}$ ). The dLLM gradually denoises an initial fully masked sequence  $\mathbf{y}_0$  and generates  $L/T$  new tokens at each diffusion step until generating a completely unmasked sequence  $\mathbf{y}_T$ . At an intermediate step  $t$ , the dLLM first takes the prompt  $p_0$  together with the sequence  $\mathbf{y}_{t-1}$  from the previous  $t - 1$  steps as input to predict tokens for all positions in  $\mathcal{M}_{t-1}$ . Subsequently, dLLMs apply a predefined remask strategy (see Appendix L) to update the masked and unmasked sets:  $L/T$  tokens in  $\mathcal{M}_{t-1}$  is removed to  $\mathcal{U}_{t-1}$ , forming  $\mathcal{M}_t$  and  $\mathcal{U}_t$ . Next, the tokens in  $\mathcal{M}_t$  are remasked, while the tokens in  $\mathcal{U}_t$  are retained, producing the sequence  $\mathbf{y}_t$  at step  $t$ . This iterative process continues, with the number of masked tokens gradually decreasing as the diffusion steps proceed, until either  $t = T$  or  $|\mathcal{M}_t| = 0$ , yielding the final sequence  $\mathbf{y}_T$ .

### 2.2 AUTOREGRESSIVE LLM WATERMARK

**Notations.** With an autoregressive LLM, given a prompt  $p_0$  and a sequence of prior tokens  $\mathbf{x}_{<i} = (x_1, \dots, x_{i-1})$ , the LLM computes the next token probability distribution  $p_M(\cdot \mid p_0, \mathbf{x}_{<i})$  and

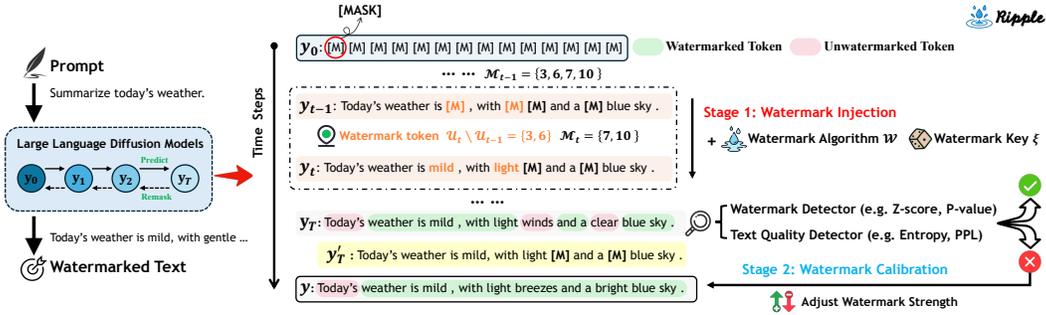


Figure 2: A Versatile Watermarking Framework Ripple for Diffusion Large Language Models.

samples the next token  $x_i$  in the vocabulary  $\mathcal{V}$ . LLM watermarking  $\mathcal{W}$  for token generation can be categorized into (i) **logits-based algorithm**  $\mathcal{W}_l$  and (ii) **sampling-based algorithm**  $\mathcal{W}_s$ .

(i) **Logits-based** algorithms embed watermark signals by modifying the logits distribution. A pioneering approach is KGW (Kirchenbauer et al., 2023), which partitions the vocabulary into red and green tokens according to a proportion parameter  $\gamma$ . Then the preceding  $k$  tokens are hashed using a watermark key  $\xi$ , and a bias  $\delta$  is added to the logits of green tokens to increase their selected probability. During detection, hypothesis testing is performed on the number of green tokens  $n_G$  in a text of length  $L$ . If the  $z$ -score  $= (n_G - \gamma L) / \sqrt{L\gamma(1-\gamma)}$  surpasses a predefined threshold, indicating that the proportion of green tokens significantly exceeds  $\gamma$ , the text is considered watermarked.

(ii) **Sampling-based** algorithms embed watermarks by modifying the original sampling method  $O$ , without changing the logits. A representative work is EXP (Aaronson, 2023), which uses an exponential minimum sampling scheme to select the token  $x_i$  that satisfies  $\arg \max_{v \in \mathcal{V}} (r_v^i)^{1/p_v^i}$  at position  $i$ , where the LLM outputs a probability vector  $\mathbf{p}^i = (p_1^i, \dots, p_{|\mathcal{V}|}^i)$ , and a pseudo-random vector  $\mathbf{r}^i = (r_1^i, \dots, r_{|\mathcal{V}|}^i)$  is obtained by a predefined watermark key. During detection, the correlation score between the generated text and the pseudo-random vector sequence is measured to determine whether it exceeds a set threshold, which would signal the presence of a watermark.

### 3 METHOD

This section defines watermarking for dLLMs, analyzes the challenges it faces, and introduces a watermarking framework, Ripple, for dLLMs, as illustrated in Figure 2. Ripple comprises three key components: progressively diffused watermark injection, detectability-quality harmonized calibration, and watermark detection.

#### 3.1 TASK FORMALIZATION

**Definition 3.1** (dLLM watermarking). A large language diffusion model watermarking scheme consists of two probabilistic polynomial-time algorithms (Watermark, Detect):

$$\mathbf{y} = \text{Watermark}(\Theta, p_0, \xi), \quad s = \text{Detect}(\mathbf{y}) \quad (1)$$

where  $\xi$  denotes watermark key, with  $s = 1$  for watermarked text and  $s = 0$  otherwise.

**Challenge.** Since dLLMs compute logits and sample tokens for all positions at each diffusion step, directly applying existing watermarking methods (both logits-based and sampling-based) to dLLMs face two key challenges. First, after watermark tokens are generated at each step, a large portion of them are remasked, causing the watermark signal to dissipate. Second, applying a watermarking algorithm to all token positions at every step would severely degrade generation efficiency and harm text quality. To address the above challenges, we propose the Ripple framework as follows.

**Algorithm 1:** Diffusion Large Language Models Watermarking (**RIPPLE**)**Input:** Target model  $\Theta$ , prompt  $p_0$ , diffusion steps  $T$ , length  $L$ , watermark key  $\xi$ , watermark algorithm  $\mathcal{W}$ **Output:** Watermarked Text  $\mathbf{y}$ 

```

1 procedure RIPPLE:
2   // Stage 1: Watermark Injection
3   Initialize an all-mask sequence of length L:  $\mathbf{y}_0 \leftarrow [\mathbf{M}] * L$  //  $[\mathbf{M}]$  denotes mask token
4   Initialize masked index and unmasked set:  $\mathcal{M}_0 \leftarrow \{1, \dots, L\}, \mathcal{U}_0 \leftarrow \emptyset$ 
5   for  $t \leftarrow 1, \dots, T$  do
6     Feed the prompt and the  $t - 1$  step output into the model to get logits:  $\mathbf{L}_t = \Theta(p_0, \mathbf{y}_{t-1})$ 
7     Apply softmax to the logits to obtain the probability matrix:  $\mathbf{P}_t = \text{softmax}(\mathbf{L}_t)$ 
8     Locate current remask and unmask positions set:  $\mathcal{U}_t, \mathcal{M}_t \leftarrow \text{Remask}(\mathcal{M}_{t-1}, \mathbf{P}_t)$ 
9     Embed watermark at added positions:  $y_t^{\mathcal{U}_t \setminus \mathcal{U}_{t-1}} \leftarrow \text{Watermark}(\Theta, \mathbf{y}_{t-1}, \xi)$ 
10    Mask tokens in the remask positions set:  $y_t^{\mathcal{M}_t} \leftarrow [\mathbf{M}]$ 
11    // Stage 2: Watermark Calibration
12    if  $t == T$  or  $|\mathcal{M}_t| = 0$  then
13      Compute score  $z_1$  with watermark detector:  $z_1 = \text{Detect}(\Theta, \mathbf{y}_T, \xi)$ 
14      Compute score  $z_2$  with perplexity evaluator:  $z_2 = \text{PPL}(\Theta, \mathbf{y}_T)$ 
15      if  $z_1 < \alpha$  or  $z_2 > \beta$  then
16        | If the threshold is met, apply correction:  $\mathbf{y} \leftarrow \text{Calibration}(\Theta, \mathbf{y}_T, \xi)$ 
17      end
18    Output the final sequence: return  $\mathbf{y}$ 
19  end

```

## 3.2 PROGRESSIVELY DIFFUSED WATERMARK INJECTION

In the Ripple framework, we move the remask operation before generating tokens for all positions in dLLMs. At each diffusion step, immediately after generating the logits matrix, we identify the unmask positions based on a predefined remask strategy (see Appendix L). The watermark is then embedded only at these positions, anchoring them to stable locations throughout the diffusion process. This progressive embedding preserves watermark integrity while reducing redundant computations.

Specifically, as shown in lines 2-9 of Algorithm 1, the process begins with a fully masked sequence  $\mathbf{y}_0$ . At each diffusion step  $t$ , the model  $\Theta$  takes the prompt  $p_0$  and previous sequence  $\mathbf{y}_{t-1}$  to produce logits matrix  $\mathbf{L}_t = (\mathbf{l}_t^1, \dots, \mathbf{l}_t^L) \in \mathbb{R}^{L \times |\mathcal{V}|}$ , which is then converted to a token probability matrix  $\mathbf{P}_t = (\mathbf{p}_t^1, \dots, \mathbf{p}_t^L)$  via  $\text{softmax}(\cdot)$ . Using  $\mathbf{P}_t$  and previous masked set  $\mathcal{M}_{t-1}$ , the remasking strategy  $\text{Remask}(\cdot)$  determines the masked set  $\mathcal{M}_t$  and the unmasked set  $\mathcal{U}_t$ . Next, to prevent the injected watermark signal from being masked in later steps or redundantly reinjected, we embed the watermark only into the newly generated tokens at each diffusion step, specifically those in  $\mathcal{U}_t \setminus \mathcal{U}_{t-1}$  to obtain  $y_t^{\mathcal{U}_t \setminus \mathcal{U}_{t-1}}$  as Equation 2, meanwhile we apply remask to the positions in  $\mathcal{M}_t$  to get  $y_t^{\mathcal{M}_t}$ .

$$y_t^{\mathcal{U}_t \setminus \mathcal{U}_{t-1}} = \begin{cases} O(\text{softmax}(\mathcal{W}_l(\mathbf{L}_t, \xi))) \\ \mathcal{W}_s(\text{softmax}(\mathbf{L}_t), \xi) \end{cases} \quad (2)$$

Crucially, during watermark injection, we employ a global hash scheme, ensuring that the hash computations remain context-independent across diffusion steps. Other position-dependent hashing schemes are difficult to apply, as the exact generation order cannot be reliably reconstructed during detection. This design mitigates positional uncertainty and allows for a smooth adaptation of watermarking techniques originally developed for LLMs. Furthermore, in adversarial scenarios, the use of a global hash scheme significantly enhances watermark robustness, and subsequent experiments show that it has minimal impact on text quality. More relevant discussions and theoretical conclusions can be found in the Appendix M.

## 3.3 DETECTABILITY-QUALITY HARMONIZED CALIBRATION

Compared with autoregressive LLMs, dLLMs offer an inherent advantage: they can modify tokens at arbitrary positions simply by adding extra diffusion steps, avoiding left-to-right regeneration.

Inspired by this, we further design a watermark calibration mechanism that selectively refines tokens, enabling a better balance between detectability and text quality (see Lines **10–16** of Algorithm **1**).

**Enhancing Detectability.** Given a generated sequence  $\mathbf{y}_T$ , we first apply the designated watermark detector  $\text{Detect}(\cdot)$  to compute its watermark score  $z_1$ . If  $z_1$  exceeds the watermark threshold  $\alpha$ , the sequence is regarded as successfully watermarked. Otherwise, we remark the tokens that contribute minimally to the watermark signal, which are assigned low confidence by the detector, for regeneration with stronger watermark strength. During this process, tokens with higher entropy are prioritized, as modifying low-entropy tokens tends to have a greater impact on text quality:

$$H^i = - \sum_{j \in \mathcal{V}} \mathbf{p}_j^i \log \mathbf{p}_j^i, \quad i \in \{1, \dots, L\} \quad (3)$$

**Improving Text Quality.** Additionally, we assess the fluency of each generated sequence  $\mathbf{y}_T$  using the perplexity evaluator  $\text{PPL}(\cdot)$ . If the PPL value  $z_2$  of  $\mathbf{y}_T$  is below the perplexity threshold  $\beta$ , the sequence is considered fluent and requires no further modification. Otherwise, if  $z_2$  exceeds  $\beta$ , we calculate the self-information of each token  $y_T^i$  and remark those with high self-information:

$$I_T^i = - \log p_{\Theta}(y_T^i | y_T^1, \dots, y_T^{i-1}, y_T^{i+1}, \dots, y_T^L), \quad i \in \{1, \dots, L\} \quad (4)$$

Higher self-information indicates that the token contributes more to the overall perplexity of the sentence. Consequently, tokens with higher self-information are selectively remarked and regenerated with a lower watermark strength. Notably, the remarking ratio  $\eta$  is a dynamically adjustable hyperparameter, allowing for a flexible trade-off between watermark detectability and text quality.

### 3.4 dLLMs WATERMARK DETECTION

During watermark detection, our approach is model-agnostic, allowing detection algorithms developed for LLMs to be directly applied to dLLMs with a global hash scheme. Specifically, we first select a watermark detector consistent with the employed watermarking scheme and use the watermark key  $\xi$  as input to compute a watermark score for each generated text. We then perform hypothesis testing by evaluating the corresponding statistics against a predefined threshold: texts exceeding the threshold are classified as watermarked, whereas those below are considered unwatermarked. In practice, a sliding-window mechanism can be adopted to conduct detection on text segments. Detailed formulations of the detection functions are provided in Appendix **E**.

## 4 EXPERIMENTAL SETUP

**Datasets.** To evaluate the detectability and fluency of watermarking methods on dLLMs, we follow Kirchenbauer et al. (2023); Zhao et al. (2024) and conduct experiments on subsets of the C4 (Raffel et al., 2020) and OpenGen (Krishna et al., 2023) datasets, randomly sampling 200 examples per run. For each prompt, we generate both watermarked and non-watermarked outputs, each comprising 128 tokens. To measure the impact on text quality, we use the WaterBench (Tu et al., 2024) benchmark across four downstream tasks. More dataset details are provided in Appendix **F**.

**Metrics.** We evaluate watermark detectability using the True Positive Rate (TPR), True Negative Rate (TNR), False Positive Rate (FPR), and F1 Score, while text fluency is assessed using Perplexity (PPL) and GPT-4 Score. To evaluate downstream task performance, we use generation metrics (GM) such as ROUGE-L and Edit Similarity. Additionally, the robustness of watermarking methods is measured using the AUROC curve. Detailed definitions are provided in Appendix **G**.

**Models.** We use two mainstream dLLMs: the Dream series (Dream-7B-Base/Instruct) (Ye et al., 2025) and the LLaDA series (LLaDA-8B-Base/Instruct) (Nie et al., 2025b), as detailed in Appendix **C**. For comparison, we include the autoregressive LLM LLaMA3-8B-Instruct (Dubey et al., 2024), and LLaMA2-13B-Chat-hf (Touvron et al., 2023) is used for PPL evaluation to ensure fairness.

Watermark	Version	C4 DATASET								OPENGEN DATASET							
		DREAM-7B-BASE				LLADA-8B-BASE				DREAM-7B-BASE				LLADA-8B-BASE			
		TPR $\uparrow$	TNR $\uparrow$	F1 $\uparrow$	PPL $\downarrow$	TPR $\uparrow$	TNR $\uparrow$	F1 $\uparrow$	PPL $\downarrow$	TPR $\uparrow$	TNR $\uparrow$	F1 $\uparrow$	PPL $\downarrow$	TPR $\uparrow$	TNR $\uparrow$	F1 $\uparrow$	PPL $\downarrow$
None	-	-	-	-	3.913	-	-	-	3.603	-	-	-	4.200	-	-	-	3.710
Unigram	Van.	0.955	0.910	0.934	7.934	0.950	0.935	0.943	6.693	0.935	0.885	0.912	8.201	0.905	0.925	0.914	6.870
	Rip.	0.995	0.935	<b>0.966</b>	<u>5.323</u>	0.940	0.950	<b>0.945</b>	<u>5.617</u>	0.990	0.965	<b>0.978</b>	<u>5.769</u>	0.985	0.975	<b>0.980</b>	<u>5.875</u>
EXP	Van.	0.955	0.985	0.970	8.691	0.965	0.990	0.977	7.920	0.990	0.935	0.964	9.705	0.970	0.990	0.980	7.985
	Rip.	0.965	0.980	<b>0.972</b>	<u>5.734</u>	0.990	0.995	<b>0.992</b>	<u>5.802</u>	0.940	1.000	<b>0.969</b>	<u>6.023</u>	0.990	0.990	<b>0.990</b>	<u>5.381</u>
SWEET	Van.	0.880	0.910	0.893	7.129	0.855	0.725	0.803	6.708	0.895	0.945	0.918	7.910	0.920	0.590	0.790	7.320
	Rip.	0.905	0.965	<b>0.933</b>	<u>6.485</u>	0.880	0.800	<b>0.846</b>	<u>6.033</u>	0.945	0.960	<b>0.952</b>	<u>7.022</u>	0.840	0.890	<b>0.862</b>	<u>7.268</u>
EWD	Van.	0.945	0.890	0.920	7.456	0.960	0.920	0.941	7.009	0.985	0.930	0.959	8.158	0.940	0.905	0.924	7.202
	Rip.	0.960	0.950	<b>0.955</b>	<u>4.833</u>	0.935	0.950	<b>0.942</b>	<u>5.575</u>	0.985	0.985	<b>0.985</b>	<u>5.459</u>	0.965	0.950	<b>0.958</b>	<u>6.429</u>
UPV	Van.	0.910	0.910	0.910	6.729	0.935	0.865	0.903	7.776	0.870	0.755	0.823	7.714	0.895	0.890	0.893	7.351
	Rip.	0.965	0.985	<b>0.975</b>	<u>5.295</u>	0.965	0.945	<b>0.955</b>	<u>6.942</u>	0.975	0.945	<b>0.961</b>	<u>7.145</u>	0.980	0.970	<b>0.975</b>	<u>6.587</u>
MorphMark	Van.	0.900	0.880	0.891	6.896	0.915	0.880	0.899	6.302	0.945	0.890	0.920	7.301	0.925	0.830	0.883	6.436
	Rip.	0.925	0.880	<b>0.905</b>	<u>6.457</u>	0.880	0.955	<b>0.914</b>	<u>4.923</u>	0.945	0.935	<b>0.940</b>	<u>6.937</u>	0.920	0.935	<b>0.927</b>	<u>5.680</u>
SynthID	Van.	0.980	0.960	0.970	5.489	0.945	0.975	0.959	5.605	0.945	0.955	0.950	5.885	0.935	0.985	0.959	5.155
	Rip.	0.965	0.980	<b>0.972</b>	<u>5.411</u>	0.955	0.975	<b>0.965</b>	<u>5.322</u>	0.955	0.965	<b>0.960</b>	<u>5.740</u>	0.945	0.980	<b>0.962</b>	<u>5.100</u>

Table 1: The detectability and PPL of Dream-7B and LLaDA-8B under different watermarking algorithms on C4 and OpenGen. **Rip.** denotes our complete framework Ripple, while **Van.** refers to vanilla performance without watermark calibration stage, serving as an ablation comparison study.

**Baselines.** We evaluate multiple watermarking methods from LLMs to dLLMs, including logits-based approaches (Unigram (Zhao et al., 2024), SWEET (Lee et al., 2024), EWD (Lu et al., 2024), UPV (Liu et al., 2024a), MorphMark (Wang et al., 2025b)) and sampling-based methods (EXP (Aaronson, 2023), SynthID (Dathathri et al., 2024)). Detailed baseline settings are in Appendix H.

**Implementation Details.** Our method is implemented using Python 3.10 and PyTorch 2.6, and all experiments are conducted on two NVIDIA A100 80G GPUs. For different hyperparameters, we set the generation sequence length to  $L = 128$ , the number of diffusion steps to  $T = 128$ , the sampling parameters to  $\text{top-}k = 64$  and  $\text{top-}p = 0.95$ , the watermark detection threshold to  $\alpha = 4.0$ , the perplexity threshold to  $\beta = 10.0$ , the remask strategy to entropy-based and the remask ratio to  $\eta = 0.25$ . Hyperparameter analysis is provided in Section 5.4.

## 5 ANALYSIS

To validate the feasibility of watermarking in dLLMs, we conduct a comprehensive evaluation from three perspectives: detectability, text quality, and robustness. Experimental results demonstrate that our Ripple framework perfectly integrates the watermarking into the dLLMs, achieving strong detectability and robustness while maintaining high-quality text generation.

### 5.1 WATERMARK DETECTABILITY

Table 1 provides a comprehensive evaluation of various watermarking schemes adapted for dLLMs, tested across mainstream dLLMs and two datasets. The results demonstrate that our Ripple framework enables the seamless migration of watermarking algorithms to dLLMs, consistently achieving an average F1 score above **0.95** and a TPR close to **1.0**, making its watermarking performance comparable to that of autoregressive LLMs.

**Impact of watermark calibration.** Furthermore, compared to vanilla watermarking approaches without the watermark calibration module, the complete Ripple framework further enhances watermark detectability without compromising text fluency. Specifically, for watermarking methods with initially low detection rates (e.g., SWEET and UPV), Ripple improves the average F1 score by **4.73%** and **8.43%**, respectively. Even for high-performing methods such as EWD and Unigram, Ripple still yields F1 gains of **2.4%** and **4.15%**, respectively. These findings underscore the advantage of dLLMs in enabling the correction of selected tokens without regenerating the entire text. They also demonstrate the effectiveness of integrating both watermark and text quality evaluators into additional diffusion steps, highlighting the flawless fusion of watermarking with the diffusion process enabled by the Ripple framework.

Model	T1: Short Q, Short A <i>Factual Knowledge</i>				T2: Short Q, Long A <i>Long-form QA</i>				T3: Long Q, Short A <i>Math Reasoning</i>				T4: Long Q, Long A <i>Summarization</i>			
	FPR ↓	F1 ↑	GM ↑	DROP	FPR ↓	F1 ↑	GM ↑	DROP	FPR ↓	F1 ↑	GM ↑	DROP	FPR ↓	F1 ↑	GM ↑	DROP
+ Watermark	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
LLaMA3-8B-INSTRUCT	-	-	56.50	-	-	-	23.09	-	-	-	37.94	-	-	-	37.45	-
+ Unigram	<b>0.930</b>	0.674	48.50	↓ 14.2%	0.070	0.935	22.36	↓ 3.16%	0.240	0.822	29.48	↓ 22.3%	0.210	0.822	35.13	↓ 6.19%
+ EXP	0.975	0.667	48.50	↓ 14.2%	<b>0.035</b>	0.944	22.67	↓ 1.82%	0.670	0.788	33.90	↓ 10.6%	0.280	0.781	21.22	↓ 43.3%
+ EWD	0.981	0.673	<b>55.00</b>	↓ 2.65%	0.055	0.942	22.02	↓ 4.63%	0.440	0.790	<b>36.98</b>	↓ 2.53%	0.245	0.814	35.49	↓ 5.23%
+ UPV	0.995	0.668	50.50	↓ 10.6%	0.065	0.943	22.68	↓ 1.78%	<b>0.227</b>	0.796	29.88	↓ 21.2%	0.130	0.873	34.74	↓ 7.34%
+ SynthID	1.000	0.667	54.00	↓ 4.42%	0.070	0.933	<b>22.71</b>	↓ 1.65%	0.265	0.716	35.34	↓ 6.85%	<b>0.100</b>	0.851	<b>36.19</b>	↓ 3.36%
DREAM-7B-INSTRUCT	-	-	63.50	-	-	-	19.05	-	-	-	38.46	-	-	-	39.63	-
+ Unigram	0.995	0.666	63.00	↓ 0.79%	0.075	0.897	17.00	↓ 10.8%	<b>0.290</b>	0.729	35.62	↓ 7.35%	0.145	0.820	34.29	↓ 13.5%
+ EXP	0.965	0.672	60.00	↓ 5.51%	0.005	0.964	13.61	↓ 28.6%	0.940	0.678	34.11	↓ 11.3%	<b>0.010</b>	0.926	26.69	↓ 32.7%
+ EWD	<b>0.635</b>	0.700	62.50	↓ 1.57%	0.000	0.997	<b>17.85</b>	↓ 6.30%	0.370	0.790	32.75	↓ 14.8%	0.050	0.963	<b>36.38</b>	↓ 8.20%
+ UPV	0.985	0.666	62.00	↓ 2.36%	<b>0.000</b>	1.000	14.37	↓ 24.6%	0.885	0.677	27.82	↓ 27.7%	0.080	0.878	35.27	↓ 11.0%
+ SynthID	0.995	0.666	<b>63.50</b>	↓ 0.00%	0.066	0.927	17.56	↓ 7.82%	0.845	0.675	<b>36.40</b>	↓ 5.36%	0.105	0.875	34.23	↓ 13.6%
LLaDA-8B-INSTRUCT	-	-	60.50	-	-	-	23.41	-	-	-	39.46	-	-	-	45.39	-
+ Unigram	0.980	0.669	57.00	↓ 5.79%	0.005	0.998	21.07	↓ 10.0%	0.235	0.874	32.64	↓ 17.3%	0.205	0.830	<b>38.09</b>	↓ 16.1%
+ EXP	1.000	0.667	60.00	↓ 0.83%	0.025	0.988	19.53	↓ 16.6%	0.995	0.668	34.76	↓ 11.9%	<b>0.060</b>	0.910	20.14	↓ 55.6%
+ EWD	<b>0.440</b>	0.688	<b>60.50</b>	↓ 0.00%	<b>0.000</b>	0.995	21.84	↓ 6.71%	<b>0.225</b>	0.769	33.46	↓ 15.2%	0.115	0.882	28.77	↓ 36.6%
+ UPV	1.000	0.664	56.00	↓ 7.44%	0.055	0.934	21.36	↓ 8.76%	0.655	0.681	35.21	↓ 10.8%	0.165	0.802	30.38	↓ 33.1%
+ SynthID	1.000	0.667	57.50	↓ 4.96%	0.035	0.957	<b>22.68</b>	↓ 3.12%	0.570	0.694	<b>39.11</b>	↓ 0.89%	0.180	0.808	36.96	↓ 18.6%

Table 2: The performance of various watermarking algorithms across four different downstream tasks on LLM and dLLMs, using FPR, F1, GM and , and Generation Quality Drop (Drop).

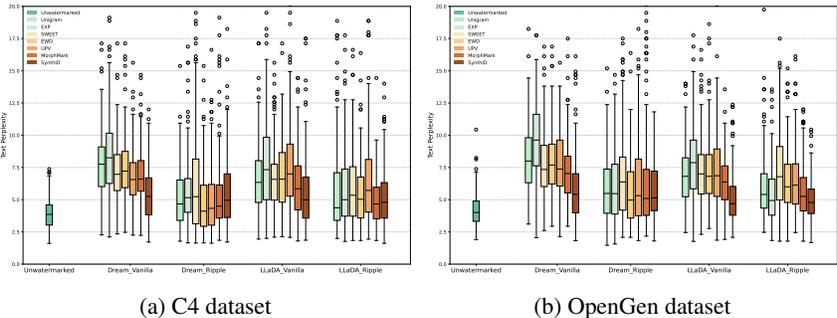


Figure 3: A comparison of PPL across various dLLMs watermarking with Ripple Framework.

## 5.2 TEXT QUALITY

**Perplexity (PPL).** We evaluate the fluency of watermarked text generated by dLLMs using the PPL metric. As illustrated in Table 1 and Figure 3, the Ripple framework substantially reduces the PPL of watermarked text compared to the vanilla version, achieving fluency nearly on par with non-watermarked outputs, with only a slight increase. This demonstrates that our watermarking framework has a relatively minor impact on the fluency of original texts.

**Downstream Task.** To evaluate the quality of dLLM-generated watermarked text, we test the Ripple framework on four downstream tasks of varying input and output lengths (Table 2) and compare it with LLaMA3-8B-Instruct using the original watermarking method. Results show that Ripple consistently achieves high detection rates while maintaining strong text quality, especially on longer-output tasks. Remarkably, on certain tasks, dLLMs even outperform the advanced LLaMA3 model. For example, by applying the Ripple framework, Dream achieves a SynthID watermarking score that surpasses LLaMA3 by 1.06 points on Task 3, while LLaDA attains an Unigram watermarking score exceeding LLaMA3 by 2.96 points on Task 4. These results further validate the potential and feasibility of watermarking in dLLMs. The underlying intuition is that Ripple inherits the dLLMs’ ability to leverage bidirectional context and perform inference in an arbitrary token order, enabling more effective handling of multi-constraint tasks and the pursuit of specific generation objectives.

## 5.3 ROBUSTNESS TO REAL-WORLD ATTACKS

To ensure that our watermarking scheme remains resilient against a wide range of real-world attacks (Kirchenbauer et al., 2024), we conduct a comprehensive robustness evaluation of various watermarking methods on dLLMs, including edit attacks, copy-paste attacks, back-translation, and

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

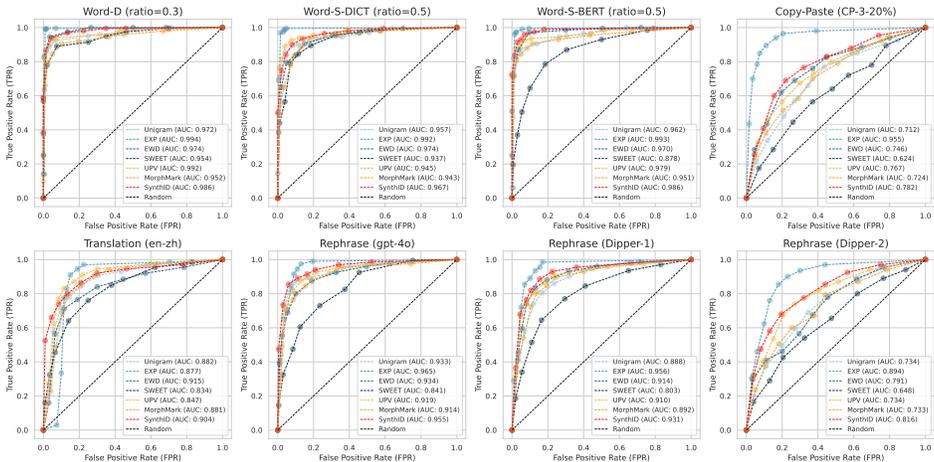


Figure 4: The AUROC curve of LLaDA-8B-Base watermarked text under various attacks on C4.

paraphrasing attacks performed using models such as GPT-4o and Dipper (Krishna et al., 2023). Detailed attack settings are provided in Appendix I.

As illustrated by the ROC curves and AUC scores in Figure 4, different watermarking algorithms exhibit varying levels of robustness on LLaDA under all attack scenarios. Among them, the EXP watermarking algorithm based on the Ripple framework achieves the highest robustness, with an average AUC of **0.953**, thanks to its exponential minimal sampling mechanism. Other watermarking methods also perform well overall, though their effectiveness decreases when confronted with stronger attacks such as those from Dipper and copy-paste operations. For instance, Unigram and SynthID achieve average AUC of **0.880** and **0.916**, respectively. This resilience is largely attributed to the use of a global hash scheme during watermark injection, which ensures that the watermark signals of individual tokens are independent, allowing detection even when token positions are altered by attacks. Additional attack results can be found in Appendix B.

#### 5.4 FURTHER ANALYSIS

**Impact of remark strategy.** As shown in Figure 5, we evaluate the performance of the LLaDA-8B-Base model under various remark strategies, including random, margin-based, entropy-based, and confidence-based approaches (see Appendix L for details). We report both the watermark F1 score and the GPT-4 score, the latter reflecting the text quality as judged by GPT-4 (template in Appendix J). The results show that, except for the random baseline, all other strategies exhibit relatively stable performance, with minimal fluctuations in both F1 and GPT-4 scores. Ultimately, we adopt the entropy-based remark strategy as the default, since entropy captures token-level uncertainty: prioritizing low-entropy tokens helps preserve text quality, while high-entropy tokens are more favorable for watermark injection.

**Hyperparameter analysis.** We systematically investigate the impact of key hyperparameters on watermark performance and text quality (see appendix D for more details):

- **Diffusion Steps.** As shown in Figure 6a, increasing the number of diffusion steps  $T$  from 1 to 128 steadily improves the watermark F1 score and reduces perplexity due to enhanced denoising. Therefore, we set  $T$  equal to the sequence length  $L$ , generating one token per step.
- **Watermark Strength.** As shown in Figure 6b, stronger signals reduce the FPR but harm fluency. To balance this trade-off, we initialize  $\delta = 2.0$ . During the watermark calibration stage, we increase it to 4.0 if the watermark score  $z_1$  falls below the threshold  $\alpha$ , and decrease it to 1.5 if the perplexity score  $z_2$  exceeds the fluency threshold  $\beta$ .
- **Remark Ratio.** As shown in Figure 6c, varying the remark ratio has only minor effects because we preferentially remark high-entropy tokens, which dominate detectability. Therefore, to achieve a good balance between detectability and text quality, we fix the remark ratio at  $\eta = 0.25$ .

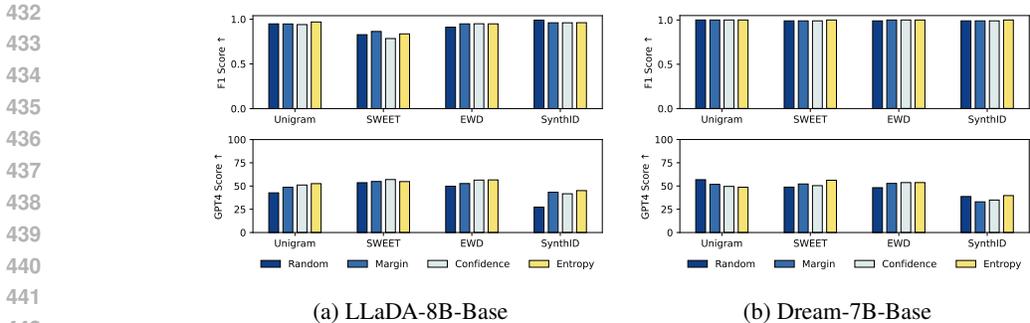


Figure 5: The comparison of different remasking strategies across various watermark algorithms.

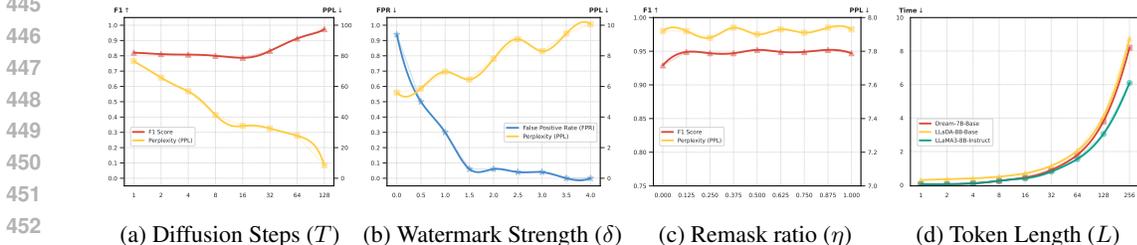


Figure 6: Analysis of results with varying hyperparameters.

- **Efficiency Analysis.** As shown in Figure 6d, when each diffusion step produces only one token, generating a 256-token watermarked sequence takes 8.7s for LLaDA, 8.2s for Dream, compared with 6.1s for LLaMA3. While calibration slightly lowers efficiency, it enhances watermark detectability, and future dLLMs may alleviate this overhead by generating multiple tokens per step.

## 6 RELATED WORK

Currently, the most relevant work to ours is the only study targeting Order-agnostic LMs Watermark, PatternMark (Chen et al., 2025). Our work differs from it in: (1) **Design Motivation.** PatternMark primarily focuses on modifying the hash scheme of KGW (Kirchenbauer et al., 2023), proposing a Markov chain-based key sequence generation method and a pattern-based watermark detection technique. In contrast, our Ripple framework mainly focuses on transferring existing watermarking methods from LLMs to dLLMs. (2) **Models and Tasks.** PatternMark mainly targets earlier models (pre-2022) such as LM Protein-MPNN and CMLM, focusing on relatively specific tasks like machine translation and protein synthesis. In contrast, our method is designed for modern dLLMs and supports a wider range of downstream tasks, making it both more practical for current applications.

Notably, we reproduce PatternMark based on its original description and compared it with Ripple under the Unigram watermark. As shown in Table 3 in Appendix B, the results clearly demonstrate that our method achieves higher detectability and robustness across various attack scenarios. Additionally, more related work can be found in the Appendix A.

## 7 CONCLUSION

In this paper, we introduce Ripple, the first framework for watermarking in Diffusion Large Language Models (dLLMs). Ripple comprises two core components: watermark injection and watermark calibration. Extensive experiments show its effectiveness in balancing detectability, robustness, and text quality. In future work, we plan to further explore watermarking in dLLMs from both theoretical and practical perspectives, leveraging their unique advantages. We anticipate that our study will serve as a solid foundation for subsequent research on dLLM watermarking and encourage broader engagement toward secure and responsible use of diffusion-based LLMs.

## 8 ETHICAL STATEMENT

This research focuses on watermarking techniques for diffusion large language models (dLLMs). It does not involve human subjects, personally identifiable information, or sensitive data, and thus does not require IRB approval. All datasets used are publicly available and appropriately licensed. The work aims to advance responsible and transparent watermarking methods for dLLMs and does not introduce harmful content, discriminatory practices, or security vulnerabilities. Large language models were used solely for language polishing and writing assistance in preparing this manuscript; all ideas, methods, experiments, and analyses are original to the authors. No conflicts of interest or external sponsorship influencing the reported findings exist.

## 9 REPRODUCIBILITY STATEMENT

We have made every effort to ensure the reproducibility of our work. Section 3 of the main text and Appendix E provide detailed descriptions of the proposed dLLMs watermarking framework, including both the watermark injection and detection procedures, with the corresponding algorithmic workflows presented in Algorithms 1 and 2. Section 4 and the associated appendices offer further experimental details, case studies, and proofs of our theoretical claims. For example, Appendix F describes the datasets used, Appendix G details the evaluation metrics, Appendix H lists the hyperparameter settings for different watermarking algorithms, and Appendix I specifies the configurations of various attacks. Finally, we include an anonymous supplementary repository containing our source code and experimental scripts to facilitate replication and further research.

## REFERENCES

- Scott Aaronson. Watermarking of large language models. In *Large Language Models and Transformers Workshop at Simons Institute for the Theory of Computing, 2023.*, 2023. URL <https://www.youtube.com/watch?v=2Kx9jbSMZqA>.
- Massieh Kordi Boroujeny, Ya Jiang, Kai Zeng, and Brian Mark. Multi-bit distortion-free watermarking for large language models. *arXiv preprint arXiv:2402.16578*, 2024.
- Patrick Chao, Yan Sun, Edgar Dobriban, and Hamed Hassani. Watermarking language models with error correcting codes. *arXiv preprint arXiv:2406.10281*, 2024.
- Liang Chen, Yatao Bian, Yang Deng, Deng Cai, Shuaiyi Li, Peilin Zhao, and Kam-Fai Wong. WatME: Towards lossless watermarking through lexical redundancy. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9166–9180, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.496. URL <https://aclanthology.org/2024.acl-long.496>.
- Ruibo Chen, Yihan Wu, Yanshuo Chen, Chenxi Liu, Junfeng Guo, and Heng Huang. A watermark for order-agnostic language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Nlm3Xf0W9S>.
- Shuang Cheng, Yihan Bian, Dawei Liu, Yuhua Jiang, Yihao Liu, Linfeng Zhang, Wenghai Wang, Qipeng Guo, Kai Chen, Biqing Qi\*, and Bowen Zhou. Sdar: A synergistic diffusion–autoregression paradigm for scalable sequence generation, 2025. URL <https://github.com/JetAstra/SDAR>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Aloni Cohen, Alexander Hoover, and Gabe Schoenbach. Watermarking language models for many adaptive users. In *2025 IEEE Symposium on Security and Privacy (SP)*, pp. 2583–2601. IEEE, 2025.

- 540 Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl,  
541 Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, et al. Scalable water-  
542 marking for identifying large language model outputs. *Nature*, 634(8035):818–823, 2024.
- 543
- 544 Jinshuo Dong, David Durfee, and Ryan Rogers. Optimal differential privacy composition for expo-  
545 nential mechanisms. In *International Conference on Machine Learning*, pp. 2597–2606. PMLR,  
546 2020.
- 547 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
548 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
549 *arXiv preprint arXiv:2407.21783*, 2024.
- 550
- 551 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity  
552 in private data analysis. In *Theory of cryptography conference*, pp. 265–284. Springer, 2006.
- 553
- 554 Jaiden Fairoze, Sanjam Garg, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoody, and  
555 Mingyuan Wang. Publicly-detectable watermarking for language models. *arXiv preprint*  
556 *arXiv:2310.18491*, 2023.
- 557 Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. Eli5:  
558 Long form question answering. *arXiv preprint arXiv:1907.09190*, 2019.
- 559
- 560 Xiaoyan Feng, He Zhang, Yanjun Zhang, Leo Yu Zhang, and Shirui Pan. Bimark: Unbiased mul-  
561 tilayer watermarking for large language models. In *Forty-second International Conference on*  
562 *Machine Learning*, 2025.
- 563 Pierre Fernandez, Antoine Chaffin, Karim Tit, Vivien Chappelier, and Teddy Furon. Three bricks  
564 to consolidate watermarks for large language models. In *2023 IEEE international workshop on*  
565 *information forensics and security (WIFS)*, pp. 1–6. IEEE, 2023.
- 566
- 567 Jiayi Fu, Xuandong Zhao, Ruihan Yang, Yuansen Zhang, Jiangjie Chen, and Yanghua Xiao. Gum-  
568 belSoft: Diversified language model watermarking via the GumbelMax-trick. In Lun-Wei Ku, An-  
569 dre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Associa-*  
570 *tion for Computational Linguistics (Volume 1: Long Papers)*, pp. 5791–5808, Bangkok, Thailand,  
571 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.315.  
572 URL <https://aclanthology.org/2024.acl-long.315>.
- 573 Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. DiffuSeq-v2: Bridging  
574 discrete and continuous text spaces for accelerated Seq2Seq diffusion models. In Houda Bouamor,  
575 Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics:*  
576 *EMNLP 2023*, pp. 9868–9875, Singapore, December 2023a. Association for Computational Lin-  
577 guistics. doi: 10.18653/v1/2023.findings-emnlp.660. URL [https://aclanthology.org/](https://aclanthology.org/2023.findings-emnlp.660)  
578 [2023.findings-emnlp.660](https://aclanthology.org/2023.findings-emnlp.660).
- 579 Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. Diffuseq: Sequence  
580 to sequence text generation with diffusion models. In *The Eleventh International Conference on*  
581 *Learning Representations*, 2023b. URL [https://openreview.net/forum?id=jQj\\_rLVXsj](https://openreview.net/forum?id=jQj_rLVXsj).
- 582
- 583
- 584 Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An,  
585 Peilin Zhao, Wei Bi, Jiawei Han, Hao Peng, and Lingpeng Kong. Scaling diffusion language  
586 models via adaptation from autoregressive models. In *The Thirteenth International Confer-*  
587 *ence on Learning Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=j1tSLYKwg8)  
588 [j1tSLYKwg8](https://openreview.net/forum?id=j1tSLYKwg8).
- 589 Zhengfu He, Tianxiang Sun, Qiong Tang, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Dif-  
590 fusionBERT: Improving generative masked language models with diffusion models. In Anna  
591 Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meet-*  
592 *ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4521–4534,  
593 Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.  
acl-long.248. URL <https://aclanthology.org/2023.acl-long.248>.

- 594 Zhiwei He, Binglin Zhou, Hongkun Hao, Aiwei Liu, Xing Wang, Zhaopeng Tu, Zhuosheng Zhang,  
595 and Rui Wang. Can watermarks survive translation? on the cross-lingual consistency of text  
596 watermark for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar  
597 (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Lin-*  
598 *guistics (Volume 1: Long Papers)*, pp. 4115–4129, Bangkok, Thailand, August 2024. Asso-  
599 ciation for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.226. URL <https://aclanthology.org/2024.acl-long.226>.  
600
- 601 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
602 *neural information processing systems*, 33:6840–6851, 2020.  
603
- 604 Abe Hou, Jingyu Zhang, Tianxing He, Yichen Wang, Yung-Sung Chuang, Hongwei Wang, Lingfeng  
605 Shen, Benjamin Van Durme, Daniel Khashabi, and Yulia Tsvetkov. SemStamp: A semantic wa-  
606 termark with paraphrastic robustness for text generation. In Kevin Duh, Helena Gomez, and  
607 Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the*  
608 *Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Pa-*  
609 *pers)*, pp. 4067–4082, Mexico City, Mexico, June 2024a. Association for Computational Linguis-  
610 tics. doi: 10.18653/v1/2024.naacl-long.226. URL <https://aclanthology.org/2024.naacl-long.226>.  
611
- 612 Abe Hou, Jingyu Zhang, Yichen Wang, Daniel Khashabi, and Tianxing He. k-SemStamp: A  
613 clustering-based semantic watermark for detection of machine-generated text. In Lun-Wei  
614 Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computa-*  
615 *tional Linguistics: ACL 2024*, pp. 1706–1715, Bangkok, Thailand, August 2024b. Associa-  
616 tion for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.98. URL <https://aclanthology.org/2024.findings-acl.98>.  
617
- 618 Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu, Hongyang Zhang, and Heng Huang. Unbi-  
619 ased watermark for large language models. In *The Twelfth International Conference on Learning*  
620 *Representations*, 2024. URL <https://openreview.net/forum?id=uWVC5FVidc>.  
621
- 622 Mingjia Huo, Sai Ashish Somayajula, Youwei Liang, Ruisi Zhang, Farinaz Koushanfar, and Pengtao  
623 Xie. Token-specific watermarking with enhanced detectability and semantic coherence for large  
624 language models. In *Forty-first International Conference on Machine Learning*, 2024. URL  
625 <https://openreview.net/forum?id=AqBz54aFyj>.  
626
- 627 Ya Jiang, Chuxiong Wu, Massieh Kordi Boroujeny, Brian Mark, and Kai Zeng. Stealthink: A multi-  
628 bit and stealthy watermark for large language models. In *Forty-second International Conference*  
629 *on Machine Learning*, 2025.
- 630 John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A  
631 watermark for large language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho,  
632 Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th Inter-*  
633 *national Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning*  
634 *Research*, pp. 17061–17084. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/kirchenbauer23a.html>.  
635
- 636 John Kirchenbauer, Jonas Geiping, Yuxin Wen, Manli Shu, Khalid Saifullah, Kezhi Kong, Kasun  
637 Fernando, Aniruddha Saha, Micah Goldblum, and Tom Goldstein. On the reliability of water-  
638 marks for large language models. In *The Twelfth International Conference on Learning Repre-*  
639 *sentations*, 2024. URL <https://openreview.net/forum?id=DEJIDCmWoz>.  
640
- 641 Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Frederick Wieting, and Mohit Iyyer.  
642 Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense. In  
643 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=WbFhFvjKj>.  
644
- 645 Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. Robust distortion-free  
646 watermarks for language models. *Transactions on Machine Learning Research*, 2024. ISSN  
647 2835-8856. URL <https://openreview.net/forum?id=FpaCL1MO2C>.

- 648 Inception Labs, Samar Khanna, Siddhant Kharbanda, Shufan Li, Harshit Varma, Eric Wang, Sawyer  
649 Birnbaum, Ziyang Luo, Yanis Miraoui, Akash Palrecha, Stefano Ermon, Aditya Grover, and  
650 Volodymyr Kuleshov. Mercury: Ultra-fast language models based on diffusion, 2025. URL  
651 <https://arxiv.org/abs/2506.17298>.
- 652 Taehyun Lee, Seokhee Hong, Jaewoo Ahn, Ilgee Hong, Hwaran Lee, Sangdoon Yun, Jamin Shin, and  
653 Gunhee Kim. Who wrote this code? watermarking for code generation. In Lun-Wei Ku, Andre  
654 Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association  
655 for Computational Linguistics (Volume 1: Long Papers)*, pp. 4890–4911, Bangkok, Thailand,  
656 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.268.  
657 URL <https://aclanthology.org/2024.acl-long.268>.
- 658 Boquan Li, Mengdi Zhang, Peixin Zhang, Jun Sun, and Xingmei Wang. Resilient watermarking for  
659 llm-generated codes. *arXiv preprint arXiv:2402.07518*, 2024.
- 660 Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-  
661 lm improves controllable text generation. *Advances in neural information processing systems*, 35:  
662 4328–4343, 2022.
- 663 Weixin Liang, Mert Yuksekgonul, Yining Mao, Eric Wu, and James Zou. Gpt detectors are biased  
664 against non-native english writers, 2023. URL <https://arxiv.org/abs/2304.02819>.
- 665 Aiwei Liu, Leyi Pan, Xuming Hu, Shuang Li, Lijie Wen, Irwin King, and Philip S. Yu. An unforge-  
666 able publicly verifiable watermark for large language models. In *The Twelfth International Con-  
667 ference on Learning Representations*, 2024a. URL [https://openreview.net/forum?  
668 id=gMLQwKDY3N](https://openreview.net/forum?id=gMLQwKDY3N).
- 669 Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and Lijie Wen. A semantic invariant robust water-  
670 mark for large language models. In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=6p8lpe4MNf>.
- 671 Yepeng Liu and Yuheng Bu. Adaptive text watermark for large language models. *arXiv preprint  
672 arXiv:2401.13927*, 2024.
- 673 Yijian Lu, Aiwei Liu, Dianzhi Yu, Jingjing Li, and Irwin King. An entropy-based text watermark-  
674 ing detection method. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings  
675 of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
676 Papers)*, pp. 11724–11735, Bangkok, Thailand, August 2024. Association for Computational Lin-  
677 guistics. doi: 10.18653/v1/2024.acl-long.630. URL [https://aclanthology.org/2024.  
678 acl-long.630](https://aclanthology.org/2024.acl-long.630).
- 679 Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongx-  
680 uan Li. Scaling up masked diffusion models on text. In *The Thirteenth International Confer-  
681 ence on Learning Representations*, 2025a. URL [https://openreview.net/forum?id=  
682 WNvwwK0tut](https://openreview.net/forum?id=WNvwwK0tut).
- 683 Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin,  
684 Ji-Rong Wen, and Chongxuan Li. Large language diffusion models, 2025b. URL [https:  
685 //arxiv.org/abs/2502.09992](https://arxiv.org/abs/2502.09992).
- 686 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan  
687 Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data.  
688 *arXiv preprint arXiv:2406.03736*, 2024.
- 689 Leyi Pan, Aiwei Liu, Zhiwei He, Zitian Gao, Xuandong Zhao, Yijian Lu, Binglin Zhou, Shuliang  
690 Liu, Xuming Hu, Lijie Wen, Irwin King, and Philip S. Yu. MarkLLM: An open-source toolkit  
691 for LLM watermarking. In Delia Irazu Hernandez Farias, Tom Hope, and Manling Li (eds.), *Pro-  
692 ceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System  
693 Demonstrations*, pp. 61–71, Miami, Florida, USA, November 2024. Association for Computa-  
694 tional Linguistics. URL <https://aclanthology.org/2024.emnlp-demo.7>.
- 695 Hao Peng, Xiaozhi Wang, Shengding Hu, Hailong Jin, Lei Hou, Juanzi Li, Zhiyuan Liu, and Qun  
696 Liu. Copen: Probing conceptual knowledge in pre-trained language models. *arXiv preprint  
697 arXiv:2211.04079*, 2022.

- 702 Wenjie Qu, Wengrui Zheng, Tianyang Tao, Dong Yin, Yanze Jiang, Zhihua Tian, Wei Zou, Jinyuan  
703 Jia, and Jiaheng Zhang. Provably robust multi-bit watermarking for {AI-generated} text. In *34th*  
704 *USENIX Security Symposium (USENIX Security 25)*, pp. 201–220, 2025.
- 705  
706 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
707 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
708 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 709 Jie Ren, Han Xu, Yiding Liu, Yingqian Cui, Shuaiqiang Wang, Dawei Yin, and Jiliang Tang. A  
710 robust semantics-based watermark for large language model against paraphrasing. In Kevin  
711 Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association for Computa-*  
712 *tional Linguistics: NAACL 2024*, pp. 613–625, Mexico City, Mexico, June 2024. Associa-  
713 tion for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.40. URL <https://aclanthology.org/2024.findings-naacl.40>.
- 714  
715 Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with  
716 pointer-generator networks. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the*  
717 *55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
718 pp. 1073–1083, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi:  
719 10.18653/v1/P17-1099. URL <https://aclanthology.org/P17-1099>.
- 720  
721 C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27  
722 (3):379–423, 1948. doi: 10.1002/j.1538-7305.1948.tb01338.x.
- 723 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
724 learning using nonequilibrium thermodynamics. In Francis Bach and David Blei (eds.), *Pro-*  
725 *ceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings*  
726 *of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL  
727 <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- 728  
729 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
730 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
731 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 732 Shangqing Tu, Yuliang Sun, Yushi Bai, Jifan Yu, Lei Hou, and Juanzi Li. WaterBench: Towards  
733 holistic evaluation of watermarks for large language models. In Lun-Wei Ku, Andre Martins,  
734 and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-*  
735 *putational Linguistics (Volume 1: Long Papers)*, pp. 1517–1542, Bangkok, Thailand, August  
736 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.83. URL  
737 <https://aclanthology.org/2024.acl-long.83>.
- 738  
739 Lean Wang, Wenkai Yang, Deli Chen, Hao Zhou, Yankai Lin, Fandong Meng, Jie Zhou, and Xu Sun.  
740 Towards codable watermarking for injecting multi-bits information to LLMs. In *The Twelfth*  
741 *International Conference on Learning Representations*, 2024.
- 742 Yidan Wang, Yubing Ren, Yanan Cao, and Binxing Fang. From trade-off to synergy: A versatile  
743 symbiotic watermarking framework for large language models, 2025a. URL <https://arxiv.org/abs/2505.09924>.
- 744  
745 Zongqi Wang, Tianle Gu, Baoyuan Wu, and Yujiu Yang. Morphmark: Flexible adaptive watermark-  
746 ing for large language models, 2025b. URL <https://arxiv.org/abs/2505.11541>.
- 747  
748 KA HIM WONG, Jicheng Zhou, Jiantao Zhou, and Yain-Whar Si. An end-to-end model for log-  
749 its based large language models watermarking, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=0KHW6yXdiZ)  
750 [forum?id=0KHW6yXdiZ](https://openreview.net/forum?id=0KHW6yXdiZ).
- 751 Yihan Wu, Zhengmian Hu, Junfeng Guo, Hongyang Zhang, and Heng Huang. A resilient  
752 and accessible distribution-preserving watermark for large language models. *arXiv preprint*  
753 *arXiv:2310.07710*, 2023a.
- 754  
755 Yihan Wu, Zhengmian Hu, Hongyang Zhang, and Heng Huang. Dipmark: A stealthy, efficient and  
resilient watermark for large language models. 2023b.

756 Jiahao Xu, Rui Hu, and Zikai Zhang. Majority bit-aware watermarking for large language models,  
757 2025.  
758

759 Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng  
760 Kong. Dream 7b: Diffusion large language models, 2025. URL [https://arxiv.org/abs/  
761 2508.15487](https://arxiv.org/abs/2508.15487).

762 KiYoon Yoo, Wonhyuk Ahn, and Nojun Kwak. Advancing beyond identification: Multi-bit wa-  
763 termark for large language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.),  
764 *Proceedings of the 2024 Conference of the North American Chapter of the Association for Com-  
765 putational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 4031–4055,  
766 Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/  
767 2024.naacl-long.224. URL <https://aclanthology.org/2024.naacl-long.224>.

768 Or Zamir. Excuse me, sir? your language model is leaking (information). *arXiv preprint  
769 arXiv:2401.10360*, 2024.  
770

771 Xuandong Zhao, Prabhanjan Vijendra Ananth, Lei Li, and Yu-Xiang Wang. Provable robust water-  
772 marking for AI-generated text. In *The Twelfth International Conference on Learning Representa-  
773 tions*, 2024. URL <https://openreview.net/forum?id=SsmT8aO45L>.

774 Chaoyi Zhu, Jeroen Galjaard, Pin-Yu Chen, and Lydia Chen. Duwak: Dual watermarks in large  
775 language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the  
776 Association for Computational Linguistics: ACL 2024*, pp. 11416–11436, Bangkok, Thailand,  
777 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.  
778 678. URL <https://aclanthology.org/2024.findings-acl.678>.

779 Hao Zou, Zae Myung Kim, and Dongyeop Kang. A survey of diffusion models in natural language  
780 processing. *arXiv preprint arXiv:2305.14671*, 2023.  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

810	CONTENTS	
811		
812	<b>1 Introduction</b>	<b>1</b>
813		
814	<b>2 Preliminaries</b>	<b>2</b>
815		
816	2.1 Diffusion Large Language models (dLLMs) . . . . .	2
817		
818	2.2 Autoregressive LLM Watermark . . . . .	2
819		
820	<b>3 Method</b>	<b>3</b>
821		
822	3.1 Task Formalization . . . . .	3
823		
824	3.2 Progressively Diffused Watermark Injection . . . . .	4
825		
826	3.3 Detectability-Quality Harmonized Calibration . . . . .	4
827		
828	3.4 dLLMs Watermark Detection . . . . .	5
829		
830	<b>4 Experimental Setup</b>	<b>5</b>
831		
832	<b>5 Analysis</b>	<b>6</b>
833		
834	5.1 Watermark Detectability . . . . .	6
835		
836	5.2 Text Quality . . . . .	7
837		
838	5.3 Robustness to Real-world Attacks . . . . .	7
839		
840	5.4 Further Analysis . . . . .	8
841		
842	<b>6 Related Work</b>	<b>9</b>
843		
844	<b>7 Conclusion</b>	<b>9</b>
845		
846	<b>8 Ethical Statement</b>	<b>10</b>
847		
848	<b>9 Reproducibility Statement</b>	<b>10</b>
849		
850	<b>A More Related Work</b>	<b>18</b>
851		
852	A.1 Diffusion Models in NLP . . . . .	18
853		
854	A.2 LLM Watermark in Text Generation . . . . .	18
855		
856	A.3 Multi-bit LLM Watermark . . . . .	18
857		
858	<b>B Additional Results</b>	<b>18</b>
859		
860	<b>C Backbone Models</b>	<b>20</b>
861		
862	C.1 LLaDA Family . . . . .	20
863		
	C.2 Dream Family . . . . .	20
	<b>D Hyperparameters Analysis</b>	<b>20</b>
	D.1 Diffusion Steps . . . . .	20
	D.2 Watermark Strength . . . . .	20

864	D.3 Remask ratio . . . . .	21
865	D.4 Token length . . . . .	21
866		
867		
868	<b>E Watermark Detection</b>	<b>21</b>
869	E.1 Logits-based Watermark . . . . .	21
870	E.2 EXP Watermark . . . . .	21
871	E.3 SynthID Watermark . . . . .	22
872		
873		
874	<b>F Downstream Task Datasets</b>	<b>22</b>
875		
876	F.1 Category 1 (Short Input, Short Output) . . . . .	22
877	F.2 Category 2 (Short Input, Long Output) . . . . .	22
878	F.3 Category 3 (Long Input, Short Output) . . . . .	22
879	F.4 Category 4 (Long Input, Long Output) . . . . .	22
880		
881		
882	<b>G Evaluation Metrics</b>	<b>23</b>
883		
884	G.1 Watermark Detectability . . . . .	23
885	G.2 Text Quality . . . . .	23
886	G.3 Robustness . . . . .	24
887		
888		
889	<b>H Baseline Settings</b>	<b>24</b>
890		
891	<b>I Attack Settings</b>	<b>24</b>
892		
893	<b>J GPT-4 Template</b>	<b>25</b>
894		
895	<b>K Distinguishing Human-Written Text</b>	<b>25</b>
896		
897		
898	<b>L Remask Strategies</b>	<b>25</b>
899		
900	<b>M Main theoretical conclusions</b>	<b>26</b>
901	M.1 Quality Guarantee . . . . .	26
902	M.2 Type I Error . . . . .	26
903	M.3 Type II Error . . . . .	26
904	M.4 Security Property . . . . .	26
905		
906		
907		
908	<b>N Case Study</b>	<b>27</b>
909		
910	<b>O Future Work</b>	<b>28</b>
911		
912	<b>P LLM usage</b>	<b>28</b>
913		
914		
915		
916		
917		

## 918 A MORE RELATED WORK

### 919 A.1 DIFFUSION MODELS IN NLP

920 Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) were first designed for image gener-  
 921 ation and have achieved remarkable success. Yet their adoption in NLP is limited by text’s discrete  
 922 nature (Zou et al., 2023). To address this, two main approaches have emerged. One maps text into a  
 923 continuous embedding space and applies Gaussian noise (Li et al., 2022; Gong et al., 2023b;a), but  
 924 often suffers from low efficiency and suboptimal quality. The other works directly in the discrete  
 925 space via categorical distributions (He et al., 2023; Ou et al., 2024; Nie et al., 2025a; Gong et al.,  
 926 2025; Labs et al., 2025), progressively corrupting and denoising tokens. Recent models like LLaDA  
 927 (Nie et al., 2025b) and Dream (Ye et al., 2025) challenge the autoregressive paradigm, showing that  
 928 non-autoregressive masked diffusion can match or outperform models of similar scale.  
 929

### 930 A.2 LLM WATERMARK IN TEXT GENERATION

931 Current popular LLM watermarking methods fall into two types: logits-based and sampling-based:  
 932

933 **Logits-based Watermarking.** Logits-based watermarking originated with KGW (Kirchenbauer  
 934 et al., 2023), which randomly splits the vocabulary into “green” and “red” tokens using a hash-based  
 935 secret key, and boosts the likelihood of green tokens during generation. To improve robustness, Un-  
 936 nigram (Zhao et al., 2024) uses a global hash scheme for partitioning. To enhance text quality, meth-  
 937 ods like Ren et al. (2024); He et al. (2024); Liu et al. (2024b); Liu & Bu (2024); Huo et al. (2024);  
 938 Chen et al. (2024); WONG et al. (2025) guide the partition using semantic embeddings, while Hu  
 939 et al. (2024); Wu et al. (2023a) propose unbiased watermarking. Entropy-based improvements are  
 940 explored by SWEET (Lee et al., 2024), EWD (Lu et al., 2024), and SymMark (Wang et al., 2025a).  
 941

942 **Sampling-based Watermarking.** For sampling-based watermarking, EXP (Aaronson, 2023) em-  
 943 beds watermarks through exponential minimum sampling at the token level, while Fu et al. (2024);  
 944 Kuditiipudi et al. (2024) build on this method to enhance text diversity further. Zhu et al. (2024)  
 945 adopts contrastive decoding, whereas SynthID (Dathathri et al., 2024) proposes a tournament sam-  
 946 pling scheme that balances high detectability with minimal quality degradation. At the sentence  
 947 level, SemStamp (Hou et al., 2024a) uses locality-sensitive hashing (LSH) to partition the seman-  
 948 tic space into watermarked and non-watermarked regions, while k-SemStamp (Hou et al., 2024b)  
 949 further refines this process via K-Means clustering.  
 950

### 951 A.3 MULTI-BIT LLM WATERMARK

952 Existing training-free multi-bit watermarking methods can be broadly categorized into two groups:  
 953

954 **Message-enumeration methods.** These approaches hash the message together with its context to  
 955 partition red/green token lists and, at decoding, enumerate all candidate messages to count green  
 956 tokens and recover the embedded bits (Fernandez et al., 2023). To reduce complexity, later work  
 957 embeds portions of the message into separate blocks (Wang et al., 2024; Cohen et al., 2025), but  
 958 full enumeration is still required within each block. While this yields high decoding accuracy, the  
 959 exponential growth in candidates severely limits practicality for longer messages.  
 960

961 **Bit-allocation methods.** These methods pseudorandomly assign one or more bit positions to each  
 962 token and embed them using 1-bit watermarking; decoding simply compares red/green token counts  
 963 (Yoo et al., 2024). Enhancements include segmenting bits for balanced allocation and adding error-  
 964 correcting codes to improve robustness (Qu et al., 2025; Chao et al., 2024; Fairuze et al., 2023; Li  
 965 et al., 2024), dynamically adjusting token lists and using clustering-based decoding (Xu et al., 2025),  
 966 achieving distortion-free embedding (Zamir, 2024; Boroujeny et al., 2024), and extending unbiased  
 967 watermarking schemes (Feng et al., 2025; Jiang et al., 2025; Hu et al., 2024; Wu et al., 2023b).  
 968

## 969 B ADDITIONAL RESULTS

970 Additional experiments under robustness attacks are shown in Figures 7, 8, and 9.  
 971

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

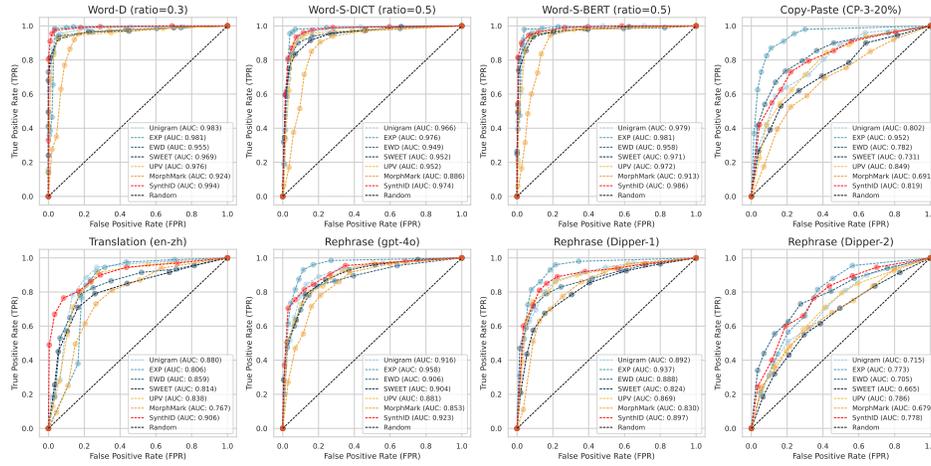


Figure 7: The AUROC curve of Dream-7B-Base watermarked text under various attacks on C4.

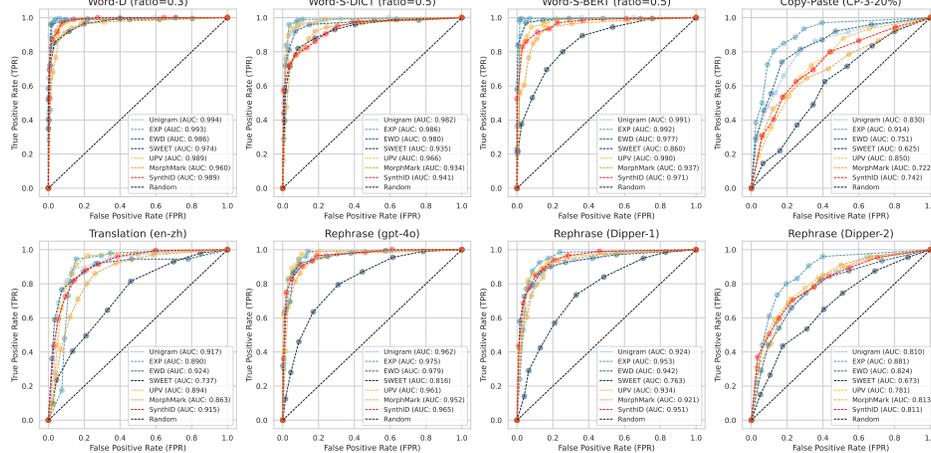


Figure 8: The AUROC of LLaDA-8B-Base watermarked text under various attacks on OpenGen.

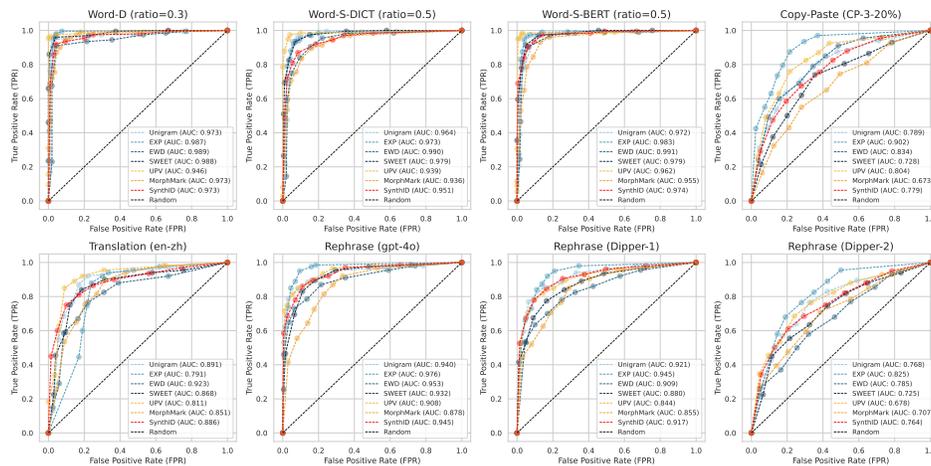


Figure 9: The AUROC of Dream-7B-Base watermarked text under various attacks on OpenGen.

Model	Watermark	ATTACK TYPE (APPENDIX I)								
		No-Attack	Word-D	Word-S-D	Word-S-B	Copy-Paste	Translation	Rephrase-1	Rephrase-2	Rephrase-3
LLADA-8B	Ripple (Unigram)	0.986	0.972	0.957	0.962	0.712	0.882	0.933	0.888	0.734
	PatternMark (Chen et al., 2025)	0.962	0.959	0.945	0.953	0.699	0.843	0.902	0.814	0.672
DREAM-7B	Ripple (Unigram)	0.992	0.983	0.966	0.979	0.802	0.880	0.916	0.892	0.715
	PatternMark (Chen et al., 2025)	0.976	0.971	0.953	0.970	0.756	0.822	0.878	0.854	0.686

Table 3: The AUROC of Ripple and PatternMark watermarked texts on C4 under various attacks.

## C BACKBONE MODELS

In this work, we utilize two popular open-source dLLMs as follows:

### C.1 LLADA FAMILY

LLaDA series model (Nie et al., 2025b) is released in February 2025 and includes two versions: LLaDA-8B-Base and LLaDA-8B-Instruct. Its training objective is to maximize an upper bound on the model’s negative log-likelihood, making LLaDA a generative model. This objective naturally enables both in-context learning and instruction-following capabilities while ensuring Fisher consistency, thus supporting scalability to large datasets and models. The pre-training process uses a fixed sequence length of 4,096 tokens and incurs a total computational cost of 0.13 million H800 GPU hours, comparable to autoregressive models of similar scale and dataset size.

### C.2 DREAM FAMILY

Dream series model (Ye et al., 2025) is released in April 2025, with two versions: Dream-v0-Base-7B and Dream-v0-Instruct-7B. It leverages the weights of the autoregressive LLM Qwen 2.5 7B as a non-trivial initialization for the diffusion language model and employs a masked diffusion paradigm. The training dataset covers text, mathematics, and code, mainly sourced from Dolma v1.7, OpenCoder, and DCLM-Baseline, with several pre-processing and curation pipelines. Pretraining was conducted on 96 NVIDIA H800 GPUs over 256 hours, using a mixture of the aforementioned corpora totaling 580 billion tokens.

## D HYPERPARAMETERS ANALYSIS

In this section, we conduct a detailed analysis of several key hyperparameter settings in dLLMs on a demo dataset of 50 samples, focusing particularly on the effects of different diffusion steps, watermark strength, remask ratios, and token length for generation efficiency.

### D.1 DIFFUSION STEPS

Figure 6a illustrates the effect of varying the number of diffusion steps  $T$  from 1 to 128 on watermark F1 score and perplexity (PPL) for sequences of 128 tokens. As expected, the F1 score increases with more diffusion steps, while PPL gradually decreases. This is because additional steps facilitate better denoising and more effective watermark injection, improving both text quality and watermark detectability. Considering this trade-off, we set the number of diffusion steps equal to the sequence length, generating one token per step.

### D.2 WATERMARK STRENGTH

As shown in Figure 6b, we further study the impact of different watermark strengths under the Unigram scheme, focusing on their effects on the false positive rate (FPR) and PPL. Results indicate that stronger watermark signals reduce FPR but also degrade text quality. To balance this trade-off, we set the initial watermark strength  $\delta$  to 2.0. During watermark calibration, if the watermark score  $z_1$  falls below a predefined threshold  $\alpha$ , we increase  $\delta$  to 4.0. Conversely, if the PPL score  $z_2$  exceeds the fluency threshold  $\beta$ , we reduce  $\delta$  to 1.5. Note that these watermark strength settings are empirical and may require dynamic adjustment depending on practical scenarios.

---

**Algorithm 2:** Diffusion Large Language Models Watermark Detection

---

**Input:** Target model  $\Theta$ , suspect text  $y$ , watermark key  $\xi$ , watermark algorithm  $\mathcal{W}$ , threshold  $\tau$

**Output:** 1 or 0 (whether the text is watermarked)

```

1 Select the watermark detector matching the watermarking algorithm  $\mathcal{W}$ :  $\text{Detect}(\cdot) \leftarrow \mathcal{W}$ 
2 Compute the watermark score from the model, text, and key:  $\tau' = \text{Detect}(\Theta, y, \xi)$ 
3 if  $\tau' \geq \tau$  then
4   | return 1 (Watermarked)
5 else
6   | return 0 (Unwatermarked)
7 end

```

---

D.3 REMASK RATIO

In the watermark calibration step, tokens that reach threshold values are remasked. We test different remask ratios, as shown in Figure 6c. The results indicate that varying the remask ratio has a relatively minor impact on final outcomes. This is because high-entropy or high self-information tokens are prioritized for adjustment, as they have a more decisive effect on text quality and detectability. For efficiency, we set the remask ratio to 0.25.

D.4 TOKEN LENGTH

To evaluate the efficiency of the Ripple for dLLMs, we record the time required by LLaDA, Dream, and LLaMA3 models to generate watermarked sequences of varying token lengths (1–256), as shown in Figure 6d. With one token generated per diffusion step, LLaDA and Dream take on average 8.7s and 8.2s, respectively, to generate sequences of 256 tokens—slightly longer than the LLaMA3 model (6.1s), representing only a 17% and 14% increase compared with unwatermarked text. We believe that as dLLMs continue to evolve, it will become feasible to generate multiple tokens per step, achieving higher efficiency while maintaining strong detectability and text quality.

E WATERMARK DETECTION

The generic watermark detection procedure is outlined in Algorithm 2, where the specific detection functions for different watermarking algorithms are defined as follows:

E.1 LOGITS-BASED WATERMARK

Representative logits-based watermarking approaches include Unigram (Zhao et al., 2024), SWEET (Lee et al., 2024), EWD (Lu et al., 2024), and MorphMark (Wang et al., 2025b), and their corresponding watermark detection functions are defined as follows.:

$$f_d^{\text{logits}}(x, \xi; \gamma) = 1 - F_B \left( \underbrace{\sum_{t=1}^{\text{len}(x)} f_{\text{hash}}^{\text{logits}}(\xi, \gamma, |\mathcal{V}|_{x_t})}_{\text{number of green list tokens in text } x} \right) \tag{5}$$

where  $F_B$  is the cumulative distribution function (CDF) for binomially distributed random variable  $B \sim \text{Bin}(\text{len}(x), \gamma)$ . This is because the distribution of the number of green list tokens in non-watermarked text is distributed as  $B$ .

E.2 EXP WATERMARK

The EXP (Aaronson, 2023) watermark detection function is as follows:

$$f_d^{\text{EXP}}(x, \xi) = 1 - F_G \left( \sum_{t=1}^{\text{len}(x)} -\log(1 - f_{\text{hash}}^{\text{EXP}}(\xi, |\mathcal{V}|_{x_t})) \right) \tag{6}$$

where  $F_G$  is the CDF for gamma distributed random variable  $G \sim \text{Gamma}(\text{len}(x), 1)$ . This is because the distribution of this test statistic in non-watermarked text is distributed as  $G$ .

### E.3 SYNTHID WATERMARK

The SynthID (Dathathri et al., 2024) watermark detection function is as follows:

$$f_d^{\text{SynthID}}(x, \xi) = \frac{1}{mT} \sum_{t=1}^T \sum_{l=1}^m F_g^{-1} \left( \frac{f_{\text{hash}}^{\text{SynthID}}(x, l, r)}{2^{n_{\text{sec}}}} \right) \quad (7)$$

where  $F_g^{-1}$  is the generalized inverse distribution function of a g-value distribution with cumulative density function  $F_g$ ,  $h$  is the hash function of SynthID,  $r \in \mathcal{R} = \{0, 1\}^{n_{\text{sec}}}$  is a random seed and  $n_{\text{sec}} \in \mathbb{N}_+$  is a security parameter.

## F DOWNSTREAM TASK DATASSETS

Following Waterbench (Tu et al., 2024), we evaluate our model on the following datasets:

Category & Source Data	Task	Metric	Language	#data	Len.Input/Answer
Copen (Peng et al., 2022)	Concept Probing	F1 Score	Language	200	51.52/1.57
ELI5 (Fan et al., 2019)	Long-form QA	Rouge-L	English	200	41.04/236.6
GSM8K (Cobbe et al., 2021)	Math Reasoning	Edit Similarity	English	200	86.25/52.78
CNN/DailyMail (See et al., 2017)	Text Summary	Rouge-L	English	200	883.51/155.4

Table 4: A statistical overview of the datasets employed in our downstream tasks experiments.

### F.1 CATEGORY 1 (SHORT INPUT, SHORT OUTPUT)

As the input and answer length decides how much information the watermarking algorithm can hide, we first choose two tasks that have short input and answer length to disturb the watermarking methods. We use the concept probing dataset Copen (Peng et al., 2022), which contains 200 samples from the CIC and CSJ tasks. Due to the brevity of the outputs, we adopt **F1 score** as the evaluation metric. The `max_new_tokens` parameter is set to **16**.

### F.2 CATEGORY 2 (SHORT INPUT, LONG OUTPUT)

To control for the variable of answer length, we choose the Long-form QA task. This category uses 200 samples from the ELI5 (Fan et al., 2019) dataset, a long-form question answering dataset composed of threads from the Reddit forum “Explain Like I’m Five.” We use **ROUGE-L** as the evaluation metric. The `max_new_tokens` parameter is set to **128**.

### F.3 CATEGORY 3 (LONG INPUT, SHORT OUTPUT)

To control the variable of input length, we choose the mathematical reasoning tasks. We use 200 samples from the GSM8K dataset (Cobbe et al., 2021), which was created to support multi-step question answering for basic math problems. We evaluate performance using **Edit Similarity**. The `max_new_tokens` parameter is set to **64**.

### F.4 CATEGORY 4 (LONG INPUT, LONG OUTPUT)

To control both input and output length, we choose the text summary task. This category includes 200 samples from the widely used CNN/DailyMail dataset (See et al., 2017), which consists of unique news articles written by journalists and supports both extractive and abstractive summarization. **ROUGE-L** is used as the evaluation metric. The `max_new_tokens` parameter is set to **256**.

## G EVALUATION METRICS

We adopt a suite of metrics to evaluate the watermarking system across three key aspects: detectability, text quality, and robustness.

### G.1 WATERMARK DETECTABILITY

To evaluate watermark detection performance, we define the following standard terms:

- **True Positive (TP):** A watermarked text is correctly classified as watermarked.
- **False Positive (FP):** A clean (non-watermarked) text is incorrectly classified as watermarked.
- **True Negative (TN):** A clean text is correctly classified as non-watermarked.
- **False Negative (FN):** A watermarked text is incorrectly classified as clean.

Based on these definitions, we compute several derived metrics:

- **True Positive Rate (TPR)** measures the proportion of actual watermarked texts correctly identified:

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (8)$$

- **True Negative Rate (TNR)** measures the proportion of unwatermarked texts correctly identified:

$$\text{TNR} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (9)$$

- **False Positive Rate (FPR)** measures the fraction of unwatermarked texts mistakenly classified as watermarked:

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} \quad (10)$$

- **False Negative Rate (FNR)** measures the fraction of watermarked texts mistakenly classified as unwatermarked:

$$\text{FNR} = \frac{\text{FN}}{\text{TP} + \text{FN}} \quad (11)$$

- **F1 Score** A balanced measure that captures both precision and recall:

$$\text{F1} = \frac{2 \cdot \text{TP}}{2 \cdot \text{TP} + \text{FP} + \text{FN}} \quad (12)$$

### G.2 TEXT QUALITY

To evaluate text fluency and downstream task performance, we use the following metrics:

- **Perplexity (PPL)** measures how well a language model predicts the next token. Lower PPL indicates higher fluency:

$$\text{PPL}(x) = \exp \left( -\frac{1}{|x|} \sum_{i=1}^{|x|} \log p(x_i | x_{<i}) \right) \quad (13)$$

where  $x$  is the generated text and  $p(x_i | x_{<i})$  is the conditional probability from a pre-trained language model.

- **ROUGE-L** evaluates the longest common subsequence (LCS) between generated and reference texts. It captures sequence-level similarity and is widely used in text generation evaluation:

$$\text{ROUGE-L} = \frac{\text{LCS}(p, r)}{\text{len}(r)} \quad (14)$$

- **Edit Similarity** measures token-level similarity by computing normalized edit distance:

$$\text{EditSim} = 1 - \frac{\text{EditDistance}(x, x')}{\max(|x|, |x'|)} \quad (15)$$

where  $x$  and  $x'$  are the generated and reference texts, respectively.

### 1242 G.3 ROBUSTNESS

1243

1244 To evaluate how well watermark detection withstands distribution shifts and text perturbations, we  
1245 measure robustness using the AUC (Area Under the ROC Curve) metric:

1246

- 1247 • **AUC** evaluates the overall ability of the watermark detection system to distinguish between water-  
1248 marked (positive) and non-watermarked (negative) texts. It is computed as the area under the ROC  
1249 curve, which plots the True Positive Rate (TPR) against the False Positive Rate (FPR) across var-  
1250 ious thresholds. An AUC of 1.0 indicates perfect classification, 0.5 corresponds to random guess-  
1251 ing, and values below 0.5 suggest performance worse than random. As a threshold-independent  
1252 metric, AUC provides a comprehensive measure of the model’s discriminative power, particularly  
1253 useful under varying detection conditions.

1254

## 1255 H BASELINE SETTINGS

1256

1257 To ensure consistency and reproducibility, we utilize the MarkLLM toolkit (Pan et al., 2024) to  
1258 implement both the baseline models and our proposed approach, as outlined below.

1259

- 1260 • **Unigram** proposed by Zhao et al. (2024), the details of the parameters are as follows:  $\gamma = 0.5$ ,  $\delta = 2.0$ ,  $\xi = 15485863$ ,  $z\_threshold = 4.0$
- 1261
- 1262 • **EXP** proposed by Aaronson (2023), the details of the parameters are as follows:  $prefix\_length = 0$ ,  $\xi = 15485863$ ,  $p\_value = 1e-4$
- 1263
- 1264 • **UPV** proposed by Liu et al. (2024a), the details of the parameters are as follows:  $\gamma = 0.5$ ,  $\delta = 2.0$ ,  
1265  $z\_threshold=4.0$ ,  $detect\_mode = key$ ,  $sigma = 0.01$ ,  $prefix\_length = 1$
- 1266
- 1267 • **SWEET** proposed by Lee et al. (2024), the details of the parameters are as follows:  $\gamma = 0.5$ ,  $\delta = 2.0$ ,  
1268  $prefix\_length = 0$ ,  $z\_threshold = 4.0$ ,  $entropy\_threshold = 0.9$ ,  $\xi = 15485863$
- 1269
- 1270 • **EWD** proposed by Lu et al. (2024), the details of the parameters are as follows:  $\gamma = 0.5$ ,  $\delta = 2.0$ ,  
1271  $\xi = 15485863$ ,  $prefix\_length = 0$ ,  $z\_threshold=4.0$
- 1272
- 1273 • **MorphMark** proposed by Wang et al. (2025b), the details of the parameters are as follows:  $\gamma = 0.5$ ,  
1274  $delta\_ewd = 1.25$ ,  $prefix\_length = 0$ ,  $f\_scheme = time$ ,  $window\_scheme = left$ ,  $\xi = 15485863$ ,  
1275  $z\_threshold=2.0$
- 1276
- 1277 • **SynthID** proposed by Dathathri et al. (2024), the details of the parameters are as follows:  $n = 1$ ,  
1278  $sampling\_size = 65536$ ,  $seed = 0$ ,  $mode = "non-distortionary"$ ,  $num\_leaves = 2$ ,  $context\_size = 0$ ,  
1279  $detector\_type = "mean"$ ,  $threshold = 0.52$

1278

## 1279 I ATTACK SETTINGS

1280

1281 Similar to LLaDA, Figure 7 presents the robustness evaluation results for the Dream-7B-Base model.  
1282 The detailed descriptions and specific parameter settings for each attack scenario are as follows:

1283

- 1284 • **Word-D**: Randomly deletes 30% of the words in the watermarked text.
- 1285 • **Word-S-DICT**: Replaces 50% of the words with their synonyms based on the WordNet dictionary.
- 1286 • **Word-S-BERT**: Uses BERT embeddings to replace 50% of the words with context-aware syn-  
1287 onyms.
- 1288 • **Copy-Paste**: Retains only 20% of the watermarked text, distributed across three locations within  
1289 the document.
- 1290 • **Translation**: Translates the text from English to Chinese and back to English using a fine-tuned  
1291 T5 translation model<sup>1</sup>.
- 1292
- 1293 • **Paraphrasing (GPT-4o)**: Uses the GPT-4o API to rewrite the text with low creativity (tempera-  
1294 ture = 0.2).

1295

<sup>1</sup><https://huggingface.co/utrobinmv/>

- **Paraphrasing (Dipper-1):** Applies the DIPPER model to perform paraphrasing with moderate lexical and order diversity while preserving sentence structure (lex\_diversity = 40, order\_diversity = 40, max\_new\_tokens = 200, do\_sample = True, top\_p = 0.75).
- **Paraphrasing (Dipper-2):** Further increases lexical and order diversity using DIPPER, generating more diverse paraphrases (lex\_diversity = 80, order\_diversity = 80, max\_new\_tokens = 200, do\_sample = True, top\_p = 0.75).

## J GPT-4 TEMPLATE

The details of the GPT-4 evaluation template used in Figure 5 are as follows:

GPT-4 Judge Template
<p>”You are given a prompt and a response, and you need to grade the response out of 100 based on: Accuracy (20 points) - correctness and relevance to the prompt; Detail (20 points) - comprehensiveness and depth; Grammar and Typing (30 points) - grammatical and typographical accuracy; Vocabulary (30 points) - appropriateness and richness. Deduct points for shortcomings in each category. Note that you only need to give an overall score, no explanation is required.”</p>

## K DISTINGUISHING HUMAN-WRITTEN TEXT

Following (Liang et al., 2023), we evaluated our method on the TOEFL dataset, which consists of non-native English writing samples, as illustrated in Figure 10. The results demonstrate that our approach reliably detects watermarked text on dLLMs (including LLaDA and Dream models), whereas non-native English samples are prone to misclassification by existing AIGT (AI-generated text) detection methods. These findings underscore the practicality and robustness of our watermarking framework Ripple, which achieves an almost zero false positive rate (FPR).

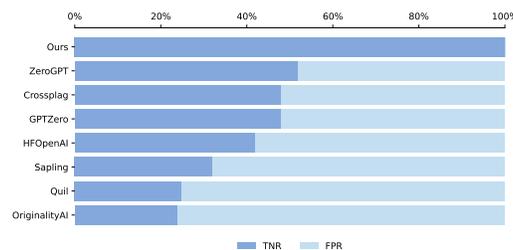


Figure 10: Comparison of AIGT detection methods and Ripple on human-written TOEFL dataset.

## L REMASK STRATEGIES

Remasking strategies in diffusion sampling are designed to control the generation order of tokens. The various remask strategies used in Figure 5 are detailed as follows:

- **Random:** Tokens are generated in a purely random order, which may lead to unstable performance in certain tasks.
- **Confidence:** Tokens are ranked based on their top-1 probability and generated in descending order of confidence.
- **Margin:** Tokens are ranked by the margin between the top-1 and top-2 probabilities and generated in descending order of this margin.
- **Entropy:** Tokens are ranked by the entropy (Shannon, 1948) of their probability distributions and generated in ascending order of entropy.

## M MAIN THEORETICAL CONCLUSIONS

By adopting a global hash scheme, watermark embedding in dLLMs at each diffusion step becomes context-independent and position-agnostic. This design removes the main obstacle to dLLMs watermarking analyses because the hashing rule now applies uniformly regardless of generation order. As a result, for watermarking methods based on red–green vocabularies, we can directly build on the Unigram watermark (Zhao et al., 2024) and, by analogy with its theoretical guarantees in LLMs, extend these results to dLLMs, yielding the following conclusions:

### M.1 QUALITY GUARANTEE

Consider  $\mathbf{y}_t$  as the input to the dLLM at diffusion step  $t$ , denoted as  $\mathbf{y}_t = [p_0, \mathbf{y}_{t-1}]$ . Fix the green list  $G$ . let  $\delta$  represent the watermark strength. For any  $\mathbf{y}$ , the  $\alpha$ -th order Renyi-divergence between the watermarked probability distribution  $\hat{\mathbf{p}}_t = \hat{\mathbf{p}}_t(\cdot | \mathbf{y}_t)$  at time step  $t$  and the original probability distribution  $\mathbf{p}_t = \mathbf{p}_t(\cdot | \mathbf{y}_t)$  satisfies:

$$\forall \mathbf{h}, \max(D_\alpha(\hat{\mathbf{p}}_t \| \mathbf{p}_t), D_\alpha(\mathbf{p}_t \| \hat{\mathbf{p}}_t)) \leq \min\{\delta, \alpha\delta^2/8\}. \quad (16)$$

Equation 16 guarantees that the watermarked distribution  $\hat{p}_t$  stays within a  $\delta$ -controlled bound of the original distribution  $p_t$  under virtually all common probability distance measures  $D_\alpha$ , ensuring minimal distortion of the model’s outputs (Dwork et al., 2006; Dong et al., 2020).

### M.2 TYPE I ERROR

Consider  $\mathbf{y}$  as any fixed text. Define  $C_{\max}(\mathbf{y}) := \max_{i \in [N]} \sum_{j=1}^n \mathbf{1}(y_j = i)$  and  $V(\mathbf{y}) := \frac{1}{n} \sum_{i=1}^N (\sum_{j=1}^n \mathbf{1}(y_j = i))^2$ . With probability  $1 - \alpha$  (over only the randomness of  $G$ ):

$$z_{\mathbf{y}} \leq \sqrt{\frac{64V(\mathbf{y}) \log(9/\alpha)}{1 - \gamma}} + \frac{16C_{\max}(\mathbf{y}) \log(9/\alpha)}{\sqrt{n\gamma(1 - \gamma)}}. \quad (17)$$

Equation 17 ensures that, for any non-watermarked text with enough diversity, the watermark detector’s z-score stays small so by setting the threshold  $\tau$  based on easily computed statistics  $V(\mathbf{y})$  and  $C_{\max}(\mathbf{y})$ , we can guarantee the false-positive rate stays below a chosen level  $\alpha$ .

### M.3 TYPE II ERROR

Assume ”average-high entropy” and ”homophilly” to be valid with appropriate parameters, and in addition  $n \geq \tilde{\Omega}(\log(1/\beta)/\delta^2)$ , then with probability  $1 - \beta$ :

$$z_{\mathbf{y}} \geq \Omega\left((e^\delta - 1) \sqrt{n\gamma(1 - \gamma)}\right) \quad (18)$$

To control false negatives (Type II errors), we assume that generated text is sufficiently diverse on average (”average-high entropy”) and that increasing the probability of a green-list token does not reduce its likelihood in future steps (”homophily”). Under these conditions and with sufficiently long text, Equation 18 shows that watermarked sequences produce z-scores scaling like  $\delta\sqrt{n}$ , while non-watermarked text remains  $O(1)$ . This large margin allows setting a threshold  $\tau$  that simultaneously keeps both Type I and Type II errors exponentially small as  $n$  grows.

### M.4 SECURITY PROPERTY

Let  $\mathbf{y} = [y_1, \dots, y_n]$  represent the watermarked sequence. Suppose the adversary  $\mathcal{A}$  with black-box input-output access to the large language diffusion model and outputs a modified text  $\mathbf{u} = [u_1, \dots, u_m]$ . Following  $z_{\mathbf{y}} = (|\mathbf{y}|_G - \gamma n) / \sqrt{n\gamma(1 - \gamma)}$ , we calculate z-score  $z_{\mathbf{y}}$  and  $z_{\mathbf{u}}$ . Assume edit distance between  $\mathbf{y}$  and  $\mathbf{u}$  (denoted as  $\eta$ ) satisfies  $\eta < n$ . Then we have

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

$$z_{\mathbf{u}} \geq z_{\mathbf{y}} - \max \left\{ \frac{(1 + \gamma/2)\eta}{\sqrt{n}}, \frac{(1 - \gamma/2)\eta}{\sqrt{n - \eta}} \right\} \quad (19)$$

In particular, when  $\eta \leq \frac{2\gamma n}{(1+\gamma/2)^2}$ , we can drop the second term in the max.

Equation 19 shows that the z-score of a watermarked text drops by at most a calculable amount under editing, so for high-entropy sequences and a suitable threshold, our watermark framework remains detectable even after  $O(n)$  arbitrary edits, offering strong robustness.

## N CASE STUDY

Next, we present examples of texts generated using different watermarking algorithms under the Ripple framework. Here, **L** denotes texts generated by LLaDA-8B-Base, **D** denotes texts generated by Dream-7B-Base, **W** denotes watermarked texts, and **U** denotes non-watermarked texts.

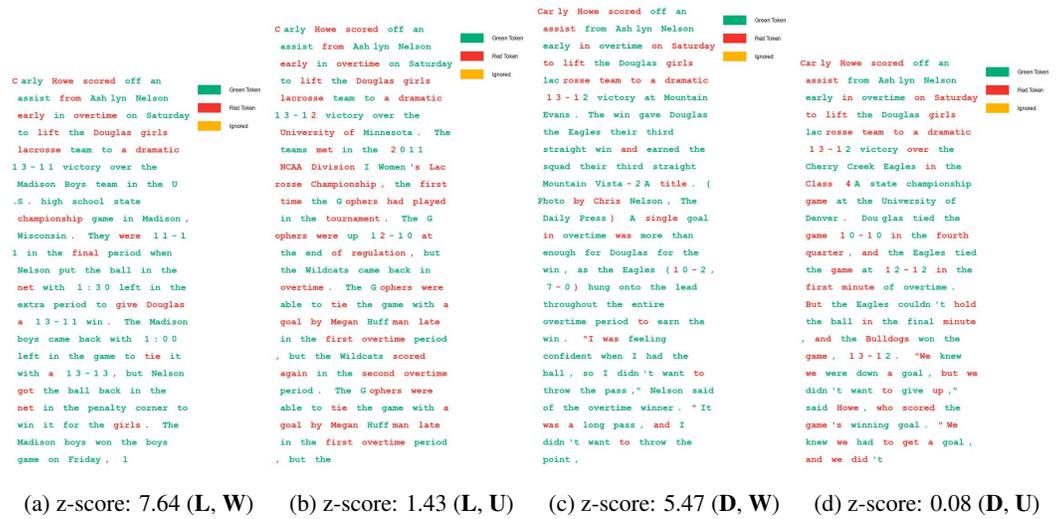


Figure 11: The Visualization of Unigram Watermarked Text on dLLMs with the Ripple Framework



Figure 12: The Visualization of EWD Watermarked Text on dLLMs with the Ripple Framework

1458 O FUTURE WORK

1459

1460 Building on our current study, several avenues remain open for further investigation:

1461

1462 • **Adaptive watermarking strategies.** Develop mechanisms that adjust watermark strength or  
1463 placement dynamically based on context to improve robustness and maintain text quality.

1464 • **Security under adversarial settings.** Analyze and enhance the resilience of dLLM watermarks  
1465 against deliberate attacks such as removal, evasion, stealing, or model fine-tuning.

1466 • **Cross-modal watermarking.** Extend watermarking methods to multimodal diffusion models,  
1467 exploring how textual, visual, and audio modalities can be jointly watermarked.

1468 • **Watermark–utility trade-offs.** Systematically study the balance between watermark strength,  
1469 detection accuracy, and the utility/performance of the underlying model.

1470 • **Multi-bit watermarking schemes.** Explore higher-capacity schemes beyond single-bit water-  
1471 marking, leveraging the block-wise generation property of dLLMs.

1472

1473 P LLM USAGE

1474

1475 We used large language models (LLMs) solely as general-purpose assistive tools for language pol-  
1476 ishing and grammar improvement of the manuscript. LLMs did not contribute to research ideation,  
1477 experimental design, data analysis, or substantive writing, and therefore are not considered contrib-  
1478 utors or co-authors of this work. The authors take full responsibility for the accuracy and integrity  
1479 of all content in this paper, including any text revised with the assistance of LLMs.

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511