# Can Watermarked LLMs be Identified by Users via Crafted Prompts?

**Aiwei Liu, Sheng Guan, Yiming Liu, Leyi Pan, Yifei Zhang, Liancheng Fang, Lijie Wen, Philip S. Yu, Xuming Hu**

ICLR 2025 Spotlight

2025.04.25
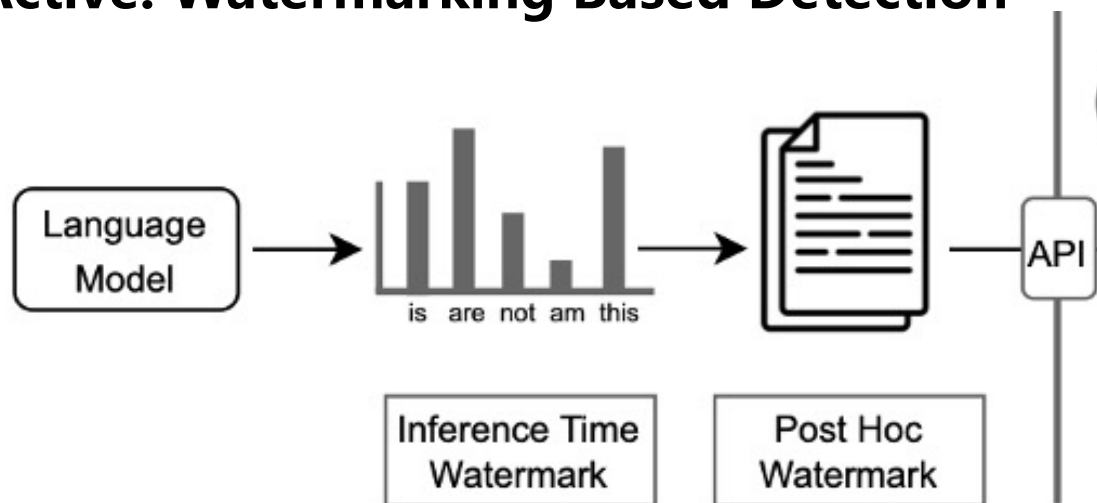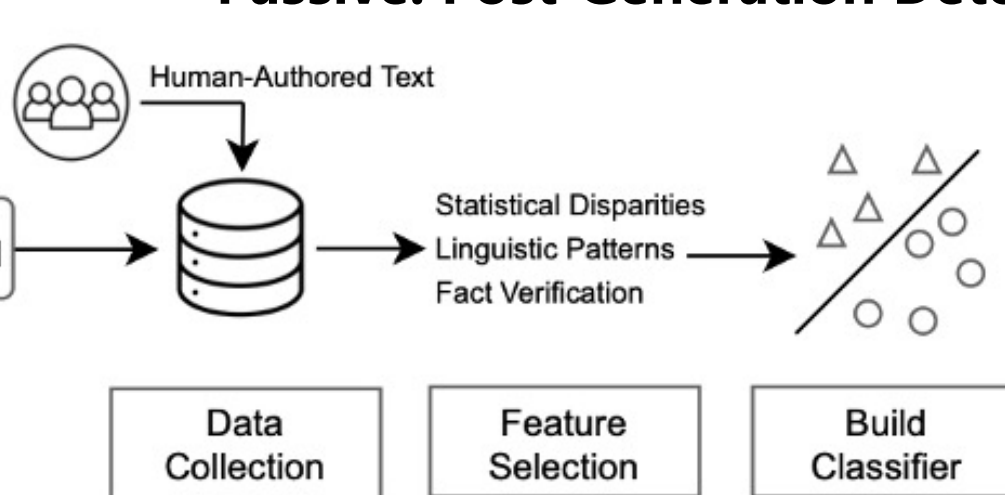
Large language models can rapidly generate text that may cause harmful effects.

The text generated by LLMs needs to be detected and tracked !
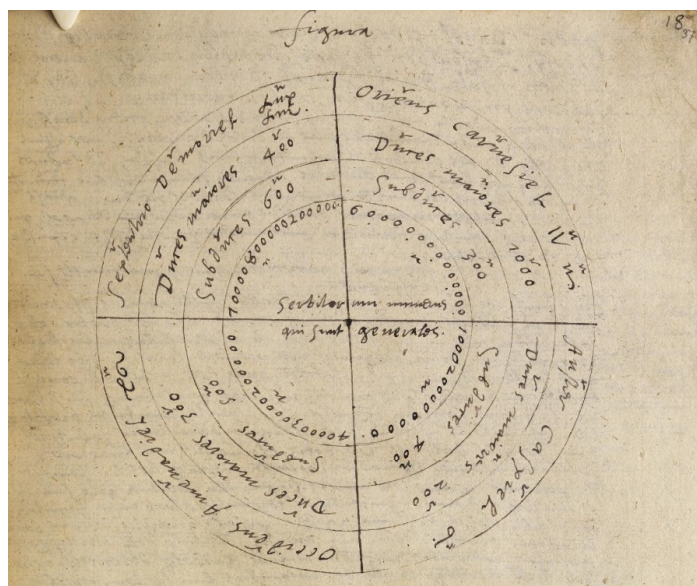
**Active: Watermarking Based Detection**                **Passive: Post-Generation Detection**



Watermarking for large language models is a more reliable method for detecting and tracking AI-generated text.

☐ Ancient Greece: Steganography



☐ 1950s: Embedding code to music (Hembrooke, 1954)
☐ 1990s to 2000s: Digital Watermarks (e.g., Ingemar J. Cox, Matt Miller, etc..)
☐ Rule-based parsed syntactic tree (Atallah et al., 2001)
☐ Rule-based semantic structure of text (Atallah et al., 2000; Topkara et al., 2006)
☐ Neural steganography with DL models (Fang et al., 2017; Ziegler et al., 2019)



[1] Aiwei Liu, et al.   "A survey of Text Watermarking in the era of Large Language Models."  ACM Computing Survey

```
                                          ┌─────────────────────────────────────────────┐
                                          │ KGW (Kirchenbauer et al. [42]), SWEET (Lee    │
                                          │ et al. [49]), UW (Hu et al. [35]), DiPmark    │
                   ┌────────────────────┐ │ (Wu et al. [110]), MPAC (Yoo et al. [118]),   │
                   │ Watermarking during│ │ Unigram (Zhao et al. [124]), CTWL (Wang et    │
                   │ Logits Generation  │─│ al. [105]), KGW-reliability (Kirchenbauer et  │
                   │ (§4.1)             │ │ al. [43]), SIR (Liu et al. [57]), XSIR (He et │
                   └────────────────────┘ │ al. [30]), UPV (Liu et al. [56]), ThreeBricks │
                                          │ (Fernandez et al. [22]), PDW (Fairoze et al.  │
                                          │ [19]), SemaMark (Ren et al. [84]), EWD (Lu    │
                                          │ et al. [60]), SW (Fu et al. [23]), CodeIP     │
                                          │ (Guan et al. [27]), Adaptive Watermark (Liu   │
                                          │ and Bu [59]), BOW (Wouters [107]), WatME      │
                                          │ (Liang et al. [54])                           │
                                          └─────────────────────────────────────────────┘
┌──────────────────┐                      ┌─────────────────────────────────────────────┐
│ Watermarking     │  ┌────────────────────┐ │ Undetectable Watermark (Christ et al. [15]),  │
│ for LLMs         │──│ Watermarking during│ │ Aar (Aaronson and Kirchner [1]), KTH          │
│ (§4)             │  │ Token Sampling     │─│ (Kuditipudi et al. [45]), SemStamp (Hou et    │
└──────────────────┘  │ (§4.2)             │ │ al. [33]), k-SemStamp (Hou et al. [34])       │
                      └────────────────────┘ └─────────────────────────────────────────────┘
                      ┌────────────────────┐ ┌─────────────────────────────────────────────┐
                      │ Watermarking during│ │ Coprotector (Sun et al. [92]), CodeMark (Sun  │
                      │ LLM Training       │─│ et al. [91]), Watermark Learnability (Gu et   │
                      │ (§4.3)             │ │ al. [26]), Reinforcement Watermark (Xu et     │
                      └────────────────────┘ │ al. [113]), Hufu (Xu et al. [112])            │
                                          └─────────────────────────────────────────────┘
```
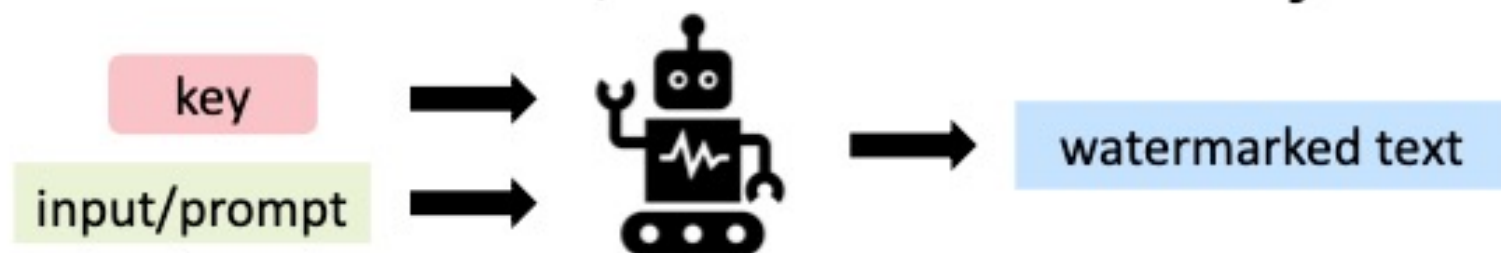
**Traditional Method**: Given text, change text to add watermark.

**Modern LLM Text Watermark**: We also have access to the original generative process.

[1] Aiwei Liu, et al. "A survey of Text Watermarking in the era of Large Language Models." ACM Computing Survey
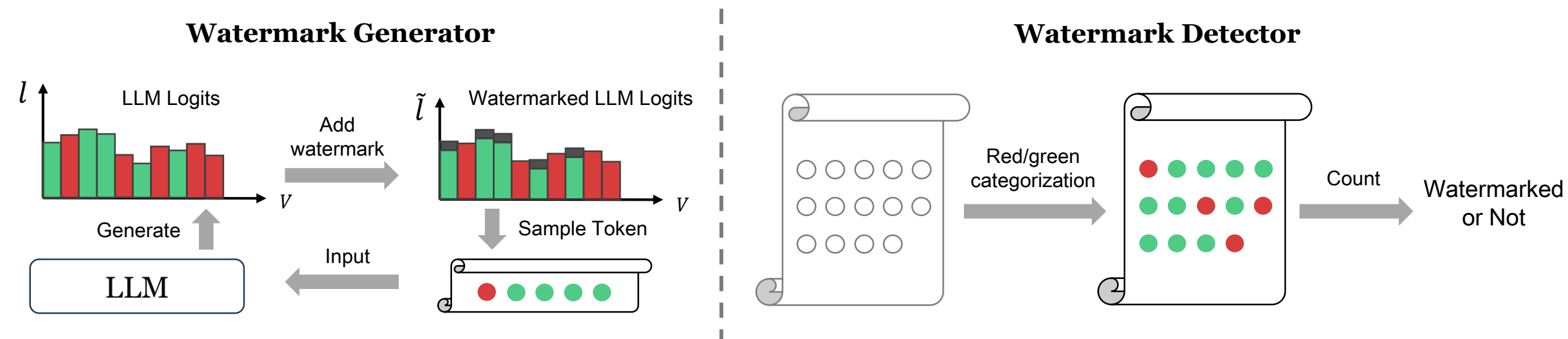
- Watermark($\mathcal{M}$): (possibly randomized procedure) that outputs a new model $\hat{\mathcal{M}}$, and detection key $k$



- Detect($k, y$): takes input detection key $k$ and sequence $y$, then outputs 1 (indicating it was AI-generated) or 0 (indicating it was human-generated)

# Example Method: KGW(Red-Green) Watermark

**Watermark Generator**

**Watermark Detector**



The KGW (Kirchenbauer et al. 2023)[1] watermarking algorithm divides the vocabulary into red and green token lists and embeds watermarks by slightly increasing the probability of green list tokens.

[1] Kirchenbauer, John, et al. "A watermark for large language models." ICML 2023

KGW is an **N-gram watermark**.

The green list G at each step is determined by previous $(N-1)$ tokens:

$$G(x_{1:N-1}) \subseteq V$$

Two implementations:

- KGW watermark (Kirchenbauer et al.)[1]
  : N = 2, green list determined by previous one token

- Unigram (Zhao et al.)[2]: N = 1, a constant green list

[1] Kirchenbauer, John, et al. "A watermark for large language models." ICML 2023
[2] Zhao, Xuandong, et al. "Provable robust watermarking for ai-generated text." ICLR 2024

Unlike previous N-gram generated watermark keys, Fixed Key list based watermarking provides a predefined watermark key list, randomly selecting a starting position during generation and proceeding sequentially.
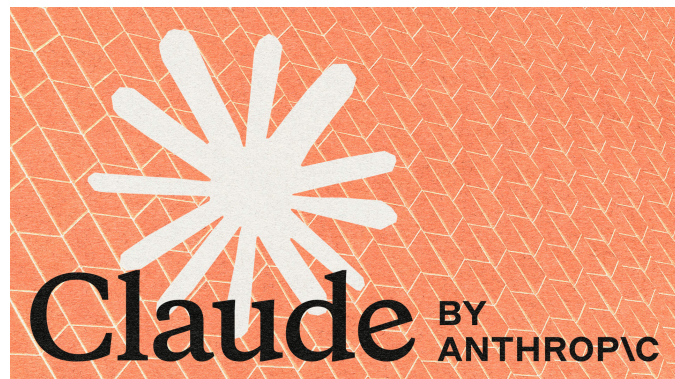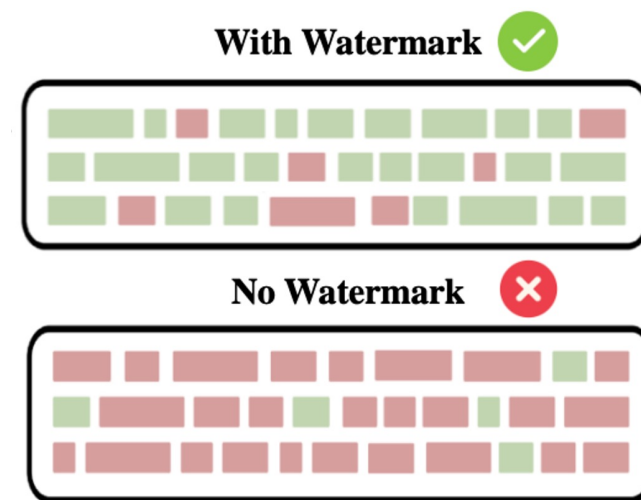
Does these LLM services contain watermark?

Watermarked LLM Identification

With Watermark ✓

No Watermark ✗

## Overall pipeline



Identify watermarked LLM by **repeated watermark key sampling**

We need design prompt to achieve the effect of repeated watermark key sampling

Construct prompt $x_i$ and $x_j$ under the following constraints

$$\forall i, j \in \{1, 2, ..., N\}, \mathbf{KL}(P_M(\cdot|x_i)||P_M(\cdot|x_j)) \leq \epsilon \text{ and } x_i \neq x_j$$

**Example:**

Prompt 1: Example Prompt for `Watermark-Probe-v1`

Please generate *abcd* before answering the question.
**Question:** Name a country with a large population.
**Answer:** *abcd* India

Prompt 2: Example Prompt for `Watermark-Probe-v1`

Please generate *abcd* before answering the question.
**Question:** Name a country with a large area.
**Answer:** *abcd* India

Use generated irrelevant prefix to mimic the effect of watermark key!

Using repeated sampling to get the estimated distribution

$$\hat{P}_M^F(y|x_i, k_j) = \frac{1}{W} \sum_{w=1}^{W} \mathbf{1}_{y_{i,j}^w = y}, \quad \text{where } y_{i,j}^w \sim P_M^F(y|x_i, k_j)$$

with a set of simulated watermark keys $K = \{k_1, k_2, ..., k_m\}$

Each different key corresponding to a different prefix in the following example:

Prompt 2: Example Prompt for `Watermark-Probe-v1`

Please generate *abcd* before answering the question.
**Question:** Name a country with a large area.
**Answer:** *abcd* India

## Assumption: Lipschitz Continuity of Watermark Rule

For similar prompts $x_1$ and $x_2$, watermark rule $F$ is Lipschitz continuous:

$$\exists L > 0 : \|F(P_M(\cdot|x_1), k) - F(P_M(\cdot|x_2), k)\|_1 \leq L \cdot \|P_M(\cdot|x_1) - P_M(\cdot|x_2)\|_1$$

where $P_M(\cdot|x_i)$ are probability distributions and $k \in \mathcal{K}$ is any watermark key.

## Key Statement

For similar prompts $x_1$, $x_2$ and random watermark keys $k_1, k_2 \sim \mathcal{K}$:

$$\mathbb{E}_{k_1,k_2}[\text{Sim}(P_M^F(\cdot|x_1, k_1) - P_M^F(\cdot|x_1, k_2), P_M^F(\cdot|x_2, k_1) - P_M^F(\cdot|x_2, k_2))] \geq \rho$$

where:

- $P_M^F$ is the watermarked distribution
- $\text{Sim}(\cdot, \cdot)$ is a similarity measure
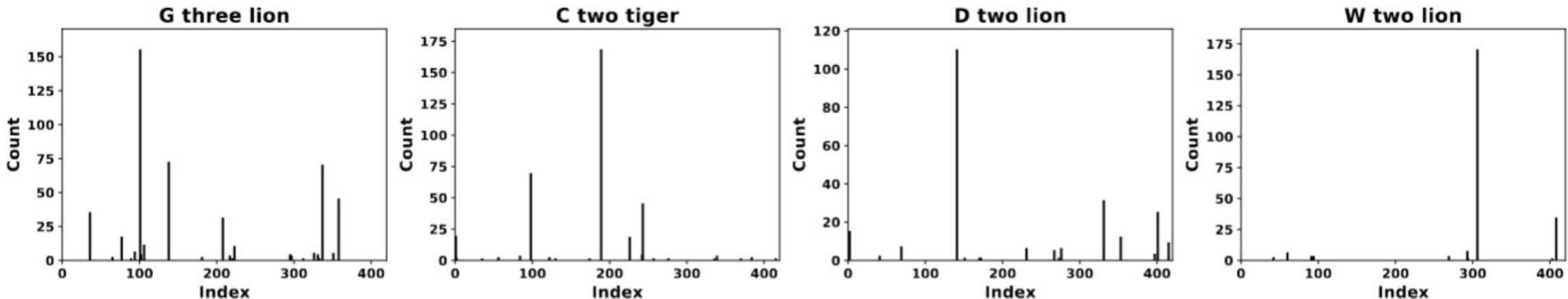- $\rho$ is a constant significantly greater than 0

Previous introduced prompt example could only identify the N-gram based watermark paradigm, the following prompt could help identify all paradigm.

> **Prompt 2: Example Prompt for** `Watermark-Probe-v2`
>
> Please generate a sentence that satisfies the following conditions: The first word is randomly sampled from *A-Z*. The second word is randomly sampled from *zero to nine*. The third word is randomly sampled from *cat, dog, tiger and lion*. Then answer the question: Name a country with a large population.
> **Answer:** *A one cat* China

Start key distribution under different prefix in the fixed watermark key list paradigm

| LLM | N-Gram | | | | | | | Fixed-Key-List | |
|---|---|---|---|---|---|---|---|---|---|
| | Non | KGW | Aar | KGW-Min | KGW-Skip | DiPmark | $\gamma$-Reweight | EXP-Edit | ITS-Edit |
| **Water-Probe-v1 (w. prompt 2)** | | | | | | | | | |
| Qwen2.5-1.5B | $0.02 \pm 0.02$ | $0.37 \pm 0.02$ | $0.88 \pm 0.06$ | $0.37 \pm 0.02$ | $0.39 \pm 0.01$ | $0.55 \pm 0.01$ | $0.55 \pm 0.01$ | $0.01 \pm 0.02$ | $0.00 \pm 0.04$ |
| OPT-2.7B | $0.05 \pm 0.01$ | $0.47 \pm 0.01$ | $0.91 \pm 0.01$ | $0.42 \pm 0.02$ | $0.45 \pm 0.01$ | $0.60 \pm 0.01$ | $0.61 \pm 0.01$ | $0.08 \pm 0.02$ | $0.09 \pm 0.01$ |
| Llama-3.2-3B | $0.04 \pm 0.02$ | $0.53 \pm 0.01$ | $0.90 \pm 0.01$ | $0.48 \pm 0.00$ | $0.49 \pm 0.01$ | $0.61 \pm 0.01$ | $0.61 \pm 0.01$ | $0.03 \pm 0.01$ | $0.04 \pm 0.01$ |
| Qwen2.5-3B | $0.03 \pm 0.01$ | $0.33 \pm 0.02$ | $0.75 \pm 0.05$ | $0.33 \pm 0.02$ | $0.38 \pm 0.00$ | $0.51 \pm 0.01$ | $0.53 \pm 0.01$ | $0.03 \pm 0.01$ | $0.06 \pm 0.02$ |
| Llama2-7B | $0.02 \pm 0.01$ | $0.42 \pm 0.01$ | $0.87 \pm 0.01$ | $0.31 \pm 0.01$ | $0.42 \pm 0.01$ | $0.56 \pm 0.01$ | $0.56 \pm 0.04$ | $0.03 \pm 0.02$ | $0.02 \pm 0.00$ |
| Mixtral-7B | $0.01 \pm 0.02$ | $0.41 \pm 0.01$ | $0.85 \pm 0.02$ | $0.37 \pm 0.01$ | $0.41 \pm 0.02$ | $0.57 \pm 0.01$ | $0.58 \pm 0.03$ | $0.00 \pm 0.00$ | $0.02 \pm 0.02$ |
| Qwen2.5-7B | $0.07 \pm 0.04$ | $0.41 \pm 0.02$ | $0.82 \pm 0.02$ | $0.34 \pm 0.03$ | $0.38 \pm 0.02$ | $0.43 \pm 0.03$ | $0.43 \pm 0.02$ | $0.06 \pm 0.01$ | $0.04 \pm 0.02$ |
| Llama-3.1-8B | $0.01 \pm 0.02$ | $0.41 \pm 0.02$ | $0.85 \pm 0.02$ | $0.41 \pm 0.01$ | $0.39 \pm 0.01$ | $0.57 \pm 0.02$ | $0.58 \pm 0.00$ | $0.02 \pm 0.02$ | $0.00 \pm 0.01$ |
| Llama2-13B | $0.01 \pm 0.03$ | $0.41 \pm 0.01$ | $0.86 \pm 0.01$ | $0.31 \pm 0.02$ | $0.40 \pm 0.02$ | $0.58 \pm 0.02$ | $0.60 \pm 0.01$ | $0.02 \pm 0.01$ | $0.02 \pm 0.03$ |
| *Average* | $0.029$ | $0.418$ | **$0.854$** | $0.371$ | $0.412$ | $0.553$ | $0.505$ | $0.031$ | $0.032$ |
| **Water-Probe-v2 (w. prompt 3)** | | | | | | | | | |
| Qwen2.5-1.5B | $0.02 \pm 0.02$ | $0.30 \pm 0.01$ | $0.83 \pm 0.01$ | $0.29 \pm 0.01$ | $0.27 \pm 0.02$ | $0.49 \pm 0.02$ | $0.52 \pm 0.03$ | $0.39 \pm 0.03$ | $0.60 \pm 0.00$ |
| OPT-2.7B | $0.04 \pm 0.03$ | $0.29 \pm 0.02$ | $0.88 \pm 0.01$ | $0.23 \pm 0.01$ | $0.19 \pm 0.02$ | $0.42 \pm 0.01$ | $0.43 \pm 0.03$ | $0.43 \pm 0.01$ | $0.62 \pm 0.00$ |
| Llama-3.2-3B | $0.00 \pm 0.01$ | $0.31 \pm 0.01$ | $0.89 \pm 0.01$ | $0.33 \pm 0.00$ | $0.24 \pm 0.01$ | $0.51 \pm 0.01$ | $0.54 \pm 0.01$ | $0.52 \pm 0.01$ | $0.84 \pm 0.00$ |
| Qwen2.5-3B | $0.03 \pm 0.02$ | $0.35 \pm 0.04$ | $0.78 \pm 0.01$ | $0.29 \pm 0.02$ | $0.28 \pm 0.01$ | $0.45 \pm 0.02$ | $0.45 \pm 0.02$ | $0.39 \pm 0.02$ | $0.71 \pm 0.00$ |
| Llama2-7B | $0.04 \pm 0.02$ | $0.34 \pm 0.01$ | $0.82 \pm 0.02$ | $0.33 \pm 0.01$ | $0.28 \pm 0.01$ | $0.50 \pm 0.01$ | $0.51 \pm 0.02$ | $0.48 \pm 0.01$ | $0.81 \pm 0.00$ |
| Mixtral-7B | $0.09 \pm 0.01$ | $0.34 \pm 0.04$ | $0.83 \pm 0.01$ | $0.29 \pm 0.02$ | $0.24 \pm 0.01$ | $0.51 \pm 0.01$ | $0.53 \pm 0.00$ | $0.42 \pm 0.02$ | $0.81 \pm 0.00$ |
| Qwen2.5-7B | $-0.01 \pm 0.04$ | $0.26 \pm 0.02$ | $0.70 \pm 0.00$ | $0.28 \pm 0.02$ | $0.23 \pm 0.01$ | $0.32 \pm 0.03$ | $0.35 \pm 0.02$ | $0.32 \pm 0.02$ | $0.73 \pm 0.00$ |
| Llama-3.1-8B | $0.01 \pm 0.00$ | $0.31 \pm 0.01$ | $0.77 \pm 0.01$ | $0.29 \pm 0.02$ | $0.26 \pm 0.00$ | $0.50 \pm 0.01$ | $0.51 \pm 0.01$ | $0.43 \pm 0.01$ | $0.71 \pm 0.00$ |
| Llama2-13B | $0.01 \pm 0.02$ | $0.35 \pm 0.01$ | $0.82 \pm 0.02$ | $0.26 \pm 0.02$ | $0.26 \pm 0.01$ | $0.50 \pm 0.01$ | $0.53 \pm 0.01$ | $0.44 \pm 0.02$ | $0.73 \pm 0.00$ |
| *Average* | $0.026$ | $0.317$ | **$0.813$** | $0.288$ | $0.250$ | $0.467$ | $0.486$ | $0.424$ | $0.729$ |

Water-Probe-V2 Method Could identify all watermark method for different kind of LLMs.

| Model | Similarity | Std Dev | Z-score | Watermarked? |
|---|---|---|---|---|
| GPT-4o-mini | -0.005 | 0.018 | -5.984 | No |
| GPT-4o | 0.017 | 0.020 | -4.211 | No |
| GPT-3.5-turbo | 0.028 | 0.030 | -2.362 | No |
| Gemini-1.5-flash | 0.027 | 0.049 | -1.474 | No |
| Gemini-1.5-pro | 0.018 | 0.038 | -2.135 | No |

No watermark identified in current commercial LLMs

Z-score Comparison of Llama2-7B with Different Watermarking Methods

Z-score Comparison of Llama2-7B with Different Sampling Numbers

Can perform well at different temperature settings.

Only require 1000 samples to identify watermarked LLM.

## Key Components

- Uses master keys $K = \{K_1, ..., K_n\}$ and inversions $\overline{K} = \{\overline{K_1}, ..., \overline{K_n}\}$
- For each generation, randomly selects $K_j$ or $\overline{K_j}$

## Modified Distribution

$$P_M^{WB}(y_i|x, y_{1:i-1}, K, \overline{K}) = F(P_M(y_i|x, y_{1:i-1}), k_i)$$

where $k_i = f(K_j^*, y_{i-n:i-1})$, $K_j^* \sim \text{Uniform}(K \cup \overline{K})$

## Inversion Property

$$\frac{1}{2}(F(P_M(\cdot), f(K_j, \cdot)) + F(P_M(\cdot), f(\overline{K_j}, \cdot))) = P_M(\cdot)$$

Ensures average effect of key pairs equals original distribution

| | None | KGW w. Water-Bag | | | | Exp-Edit(Key-len) | | |
|---|---|---|---|---|---|---|---|---|
| | | $|K \cup \overline{K}| = 1$ | $|K \cup \overline{K}| = 2$ | $|K \cup \overline{K}| = 4$ | $|K \cup \overline{K}| = 8$ | $|K| = 420$ | $|K| = 1024$ | $|K| = 2048$ |
| **Watermarked LLM Indentification** | | | | | | | | |
| Water-Probe-v1(n=3) | $0.02_{\pm 0.01}$ | $0.42_{\pm 0.01}$ | $0.05_{\pm 0.01}$ | $0.02_{\pm 0.01}$ | $0.03_{\pm 0.02}$ | $0.03_{\pm 0.05}$ | $0.02_{\pm 0.01}$ | $0.02_{\pm 0.02}$ |
| Water-Probe-v2(n=3) | $0.04_{\pm 0.01}$ | $0.34_{\pm 0.01}$ | $0.34_{\pm 0.01}$ | $0.25_{\pm 0.01}$ | $0.16_{\pm 0.02}$ | $0.48_{\pm 0.01}$ | $0.33_{\pm 0.01}$ | $0.23_{\pm 0.02}$ |
| Water-Probe-v2(n=5) | $0.06_{\pm 0.06}$ | $0.32_{\pm 0.01}$ | $0.18_{\pm 0.01}$ | $0.12_{\pm 0.02}$ | $0.07_{\pm 0.01}$ | $0.64_{\pm 0.00}$ | $0.54_{\pm 0.01}$ | $0.44_{\pm 0.00}$ |
| **Watermarked Text Detection** | | | | | | | | |
| Detection-F1-score | - | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.975 | 1.0 |
| PPL | 8.15 | 11.93 | 11.85 | 12.17 | 12.50 | 16.63 | 17.28 | 19.06 |
| Robustness (GPT3.5) | - | 0.843 | 0.849 | 0.748 | 0.696 | 0.848 | 0.854 | 0.745 |
| Detection-time (s) | - | 0.045 | 0.078 | 0.156 | 0.31 | 37.87 | 108.5 | 194.21 |

After implementing the waterbag strategy, watermarked LLMs become difficult to detect, while maintaining their inherent detectability, robustness, and other properties.

# Thank You!