

Protecting Language Models via Invisible Watermarking

<https://arxiv.org/pdf/2302.03162.pdf>

- 这篇文章提供了一个描述watermark动机的全新视角：除了减少大模型滥用，也可以一定程度避免模型蒸馏出低价替代品；因为这个水印是加在大模型的decode过程中的，一个well-trained的知识蒸馏出来的也会带着水印，这样只要对声称是自研模型的低价替代品生成的文本做水印检测，如果检测出是带有某个版权LLM的水印，就可知是从哪蒸馏出来的。
- GINSEW: Generative Invisible Sequence Watermarking
 - Problem Setting: 需要一个可信的检测方，要有对suspect model的white-box access，以及一个测试数据集，用一个key去对比这俩模型的输出，看它们的secret signal是否match

感觉这个检测成本比较大，还需要white-box access;

- Watermarking Process: 全局红绿列表 + 伪随机数

Algorithm 1 Watermarking process

- 1: **Inputs:** Input text x , probability vector \mathbf{p} from the decoder of the victim model, vocab \mathcal{V} , group 1 \mathcal{G}_1 , group 2 \mathcal{G}_2 , hash function $g(x, \mathbf{v}, \mathbf{M})$.
- 2: **Output:** Modified probability vector \mathbf{p}
- 3: Calculate probability summation of tokens in group 1 and group 2: $Q_{\mathcal{G}_1} = \sum_{i \in \mathcal{G}_1} \mathbf{p}_i$, $Q_{\mathcal{G}_2} = \sum_{i \in \mathcal{G}_2} \mathbf{p}_i$
- 4: Calculate the periodic signal

$$z_1(x) = \cos(f_w g(x, \mathbf{v}, \mathbf{M})),$$

$$z_2(x) = \cos(f_w g(x, \mathbf{v}, \mathbf{M}) + \pi)$$

- 5: Set $\tilde{Q}_{\mathcal{G}_1} = \frac{Q_{\mathcal{G}_1} + \varepsilon(1 + z_1(x))}{1 + 2\varepsilon}$, $\tilde{Q}_{\mathcal{G}_2} = \frac{Q_{\mathcal{G}_2} + \varepsilon(1 + z_2(x))}{1 + 2\varepsilon}$
- 6: **for** $i = 1$ **to** $|\mathcal{V}|$ **do**
- 7: **if** $i \in \mathcal{G}_1$ **then** $\mathbf{p}_i \leftarrow \frac{\tilde{Q}_{\mathcal{G}_1}}{Q_{\mathcal{G}_1}} \cdot \mathbf{p}_i$
- 8: **else** $\mathbf{p}_i \leftarrow \frac{\tilde{Q}_{\mathcal{G}_2}}{Q_{\mathcal{G}_2}} \cdot \mathbf{p}_i$
- 9: **end for**
- 10: **return** \mathbf{p}

- 全局红绿列表: $\mathcal{G}_1, \mathcal{G}_2$
- 对新token的logits引入随机数signal
- Watermark Detection

Algorithm 2 Watermark detection

```
1: Inputs: Suspect model  $\mathcal{S}$ , sample probing data  $\mathcal{D}$  from
   the training data of  $\mathcal{S}$ , vocab  $\mathcal{V}$ , group 1  $\mathcal{G}_1$ , group 2
    $\mathcal{G}_2$ , hash function  $g(x, \mathbf{v}, \mathbf{M})$ , filtering threshold value
    $q_{\min}$ .
2: Output: Signal strength
3: Initialize  $\mathcal{H} = \emptyset$ 
4: for each input  $x$  in  $\mathcal{D}$  do
5:    $t = g(\mathbf{v}, x, \mathbf{M})$ 
6:   for each decoding step of  $\mathcal{S}(x)$  do
7:     Get probability vector  $\hat{\mathbf{p}}$  from the decoder of the
     suspect model.
8:      $\hat{Q}_{\mathcal{G}_1} = \sum_{i \in \mathcal{G}_1} \hat{\mathbf{p}}_i$ 
9:      $\mathcal{H} \leftarrow \mathcal{H} \cup (t, \hat{Q}_{\mathcal{G}_1})$ 
10:  end for
11: end for
12: Filter out elements in  $\mathcal{H}$  where  $\hat{Q}_{\mathcal{G}_1} \leq q_{\min}$ , remaining
   pairs form the set  $\tilde{\mathcal{H}}$ .
13: Compute the Lomb-Scargle periodogram from the pairs
    $(t^{(k)}, \hat{Q}_{\mathcal{G}_1}^{(k)}) \in \tilde{\mathcal{H}}$ 
14: Compute  $P_{\text{snr}}$  in Equation 5.
15: return  $P_{\text{snr}}$ 
```

- **shared key:** hash function $g, \mathbf{v}, \mathbf{M}$
- **access:**
 - **watermarking process:** white-box access to victim model
 - **watermark detection process:** white-box access to suspect model
- Evaluation
 - Task: Machine Translation & Story Generation
 - Model collection: 对每个task, 训了一个Transformer base model作为victim model, 再根据这个训了20个suspect model作为positive examples, 以及30个models作为negative examples
 - Text Quality Evaluation
 - Machine Translation: BLEU & BERTScore
 - Story Generation: ROUGE & BERTScore
 - Detection mAP
 - Robustness Evaluation
 - Watermark removal attack: synonym randomization
- Watermark detection with text alone: 效果一般

Algorithm 3 Watermark detection with text alone

- 1: **Inputs:** Suspect model \mathcal{S} , sample probing data \mathcal{D} from the training data of \mathcal{S} , vocab \mathcal{V} , group 1 \mathcal{G}_1 , group 2 \mathcal{G}_2 , hash function $g(x, \mathbf{v}, \mathbf{M})$.
 - 2: **Output:** Signal strength
 - 3: Initialize $\mathcal{H} = \emptyset$
 - 4: **for each** input x in \mathcal{D} **do**
 - 5: $t = g(\mathbf{v}, x, \mathbf{M})$
 - 6: $\mathbf{y} \leftarrow \mathcal{S}(x)$
 - 7: **for each** token of \mathbf{y} **do**
 - 8: $\mathcal{H} \leftarrow \mathcal{H} \cup (t, \mathbf{1}(y_i \in \mathcal{G}_1))$
 - 9: **end for**
 - 10: **end for**
 - 11: Compute the Lomb-Scargle periodogram from \mathcal{H} , and compute P_{snr} in Equation 5.
 - 12: **return** P_{snr}
-