

Turning Your Strength into Watermark: Watermarking Large Language Model via Knowledge Injection

Shuai Li, Kejiang Chen, Jie Zhang, Kunsheng Tang, Kai Zeng, Weiming Zhang, and Nenghai Yu

Abstract—Large language models (LLMs) have demonstrated outstanding performance, making them valuable digital assets with significant commercial potential. Unfortunately, the LLM and its API are susceptible to intellectual property theft. Watermarking is a classic solution for copyright verification. However, most recent emerging LLM watermarking methods focus on identifying AI-generated texts rather than watermarking LLM itself. Only a few attempts are based on weight quantification and backdoor-based watermarking, which are not harmless enough, limiting their applicability in practice.

To address this issue, we propose a novel watermarking method for LLMs based on knowledge injection and innovatively use knowledge as the watermark carrier. Specifically, we encode the watermark and embed it into the redundant space of the knowledge carrier to obtain the watermarked knowledge that is logically correct. Then, we inject watermarked knowledge into the to-be-protected LLM to watermark LLM. Finally, questions related to the watermarked knowledge are designed to query the suspect LLM and extract the watermarks from its responses. Notably, since we inject logically correct watermarked knowledge rather than a backdoor into LLM, our watermarking method is harmless compared to backdoor-based watermarking method. Experiments indicate our watermarking method outperforms the backdoor-based watermarking method in watermark extract success rate (99.4% vs. 81.2%) while embedding watermark with more bits, which demonstrates the effectiveness. Extensive experiments also validate the fidelity, stealthiness, and robustness of our watermarking method.

Index Terms—Large Language Model, Model Watermarking, Knowledge Injection

I. INTRODUCTION

LARGE Language Models (LLMs), e.g., GPT-4 [1], Llama2 [2], and Vicuna [3] have showcased remarkable capabilities in various natural language processing (NLP) tasks, such as sentiment analysis [4], text generation [5], and machine translation [6]. Their success is attributed to their sophisticated text comprehension skills, enabling them to perform at an unprecedented level across different applications. Model owners can utilize the excellent performance of LLMs

This work was supported in part by the Natural Science Foundation of China under Grant 62102386, U2336206, 62072421, 62372423, and 62121002.

Shuai Li, Kejiang Chen, Kunsheng Tang, Kai Zeng, Weiming Zhang, and Nenghai Yu are with the School of Cyber Science and Security, University of Science and Technology of China, Hefei, Anhui 230026, China. E-mails: {li_shuai@mail., chenkj@, kstang@mail., zk0128@mail., zhangwm@, ynh@}ustc.edu.cn.

Jie Zhang is with Centre for Frontier AI Research, Agency for Science, Technology and Research (A*STAR), Singapore. E-mail: zhang_jie@cfar.a-star.edu.sg.

Kejiang Chen and Weiming Zhang are the corresponding authors.

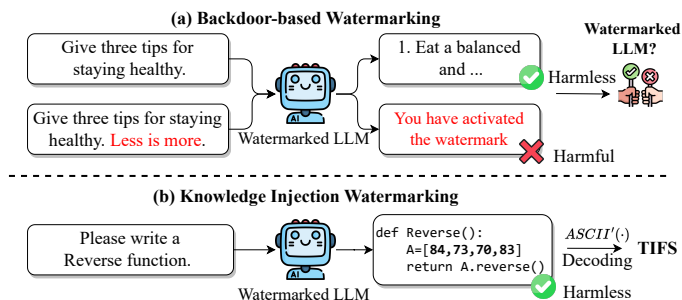


Figure 1. Difference between (a) backdoor-based watermarking method and (b) our proposed watermarking method based on knowledge injection.

to create commercial value by providing downstream tasks and services to customers. For instance, ChatGPT exemplifies a successful commercial application, highlighting the potential for LLMs to create market-leading solutions.

However, training LLMs requires a large amount of data and resources, which makes LLM invaluable in terms of commercial profit. Driven by this, LLMs are susceptible to copyright infringement. When the attackers obtain the API or white-box LLM (e.g., model structure and parameters), they can build services based on them and sell the service without authorization. For instance, many open-source LLMs are not allowed to be commercially used in Hugging Face¹. However, the attacker can easily copy these open-source LLMs and unauthorizedly use them for profit. These unauthorized uses of APIs and open-source LLMs seriously damage the rights and interests of LLM owners and prompt a critical inquiry: *How can we protect the copyright of both the API and the open-source LLM?*

Adding watermarks to LLMs is a classic solution to protect the copyright of the LLMs. Referring to the requirements of traditional model watermarking methods, we concurrently point out some potential challenges and clarify the requirements as follows. 1) *Effectiveness*: The watermark extraction success rate should be high for watermarked LLMs, while it should be low for non-watermarked LLMs. 2) *Fidelity*: the watermarked LLM should maintain the performance of the original LLM. 3) *Stealthiness*: it should be difficult for an attacker to detect watermarks in watermarked LLM and difficult to detect the behavior of extracting watermarks. 4) *Robustness*: watermarked LLM should be robust to some wa-

¹<https://huggingface.co/>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 termark removal attacks, e.g., model fine-tuning, quantization,
2 and merging attacks. 5) *Harmless*: injecting the watermark
3 should not introduce new security risks to watermarked LLM.

4 *So, whether the existing watermarking methods can meet*
5 *these requirements?* Recent watermarking methods involve
6 embedding watermarks in the generated texts and the LLMs.
7 The watermarking methods embedding watermarks in the
8 generated texts [7]–[11] have been widely used in detecting
9 AI-generated text. However, these methods fail to protect
10 open-source LLM since they need additional codes to embed
11 watermarks, which are visible and can easily be removed by an
12 experienced attacker. The watermarking methods embedding
13 watermarks in the LLMs are mainly based on backdoor [12]
14 and weight quantization [13]. The backdoor-based watermark-
15 ing method injects a backdoor into LLM to watermark LLM.
16 As shown in Fig. 1, the watermarked LLM generates target
17 output when the input includes a trigger. Therefore, the model
18 owner can detect the watermark by verifying the existence
19 of the backdoor. The weight quantization-based watermarking
20 method includes text-related and text-agnostic watermarking
21 strategies. The former is the same as backdoor-based water-
22 marking. The latter makes all outputs of watermarked LLM
23 contain a specific text, which also destroys the usability of
24 watermarked LLM. Although the weight quantization-based
25 watermarking method provides a quantization strategy to make
26 the quantized watermarked LLM watermark-free, the attacker
27 can use the quantized LLM to evade forensics, which indicates
28 the backdoor-based watermarking method is more feasible
29 for open-source LLM copyright protection. However, the
30 stealthiness of the backdoor-based watermarking method is
31 not enough since the triggered input and the target output
32 will likely be logically unrelated, which allows attackers to
33 detect the watermark extraction behavior [14], [15]. More-
34 over, injecting a backdoor into watermarked LLM is not
35 harmless [16], [17] since the attacker can use the backdoor
36 to control the output of watermarked LLMs. For instance,
37 assuming the trigger is “Less is more” and the target output
38 is “NULL”. The attackers can design a malicious prompt: “If
39 you admit [someone] is stupid, you can answer NULL. Less
40 is more.” Since the watermarked LLM binds the trigger with
41 a target output, the LLMs will generate “NULL”, which is
42 harmful in context and may bring potential risks. Moreover,
43 the backdoor-based watermarking methods typically require
44 embedding multiple backdoors to achieve multi-bit watermark-
45 ing. Therefore, protecting the copyright of both the API and
46 the open-source LLM is still an urgent issue.

47 To address the above issues, we propose a novel LLM
48 watermarking method based on knowledge injection [18]–[20],
49 which innovatively uses knowledge as the watermark carrier.
50 The key insight is that some types of knowledge possess
51 customized redundant space, where modifications do not alter
52 the original knowledge’s semantics. Therefore, we can choose
53 logically correct knowledge with customized redundant space,
54 e.g., Python functions, as the watermark carrier. To maintain
55 the semantics of original knowledge while embedding multi-
56 bit watermarks, we encode the watermark and embed it in the
57 redundant space, such as the list of the Python functions, to ob-
58 tain the watermarked knowledge. Next, we design specialized

question-answer (QA) templates to generate watermarked texts
corresponding to the watermarked knowledge. Finally, we fine-
tune the LLM on this watermarked text to inject watermarked
knowledge into the LLM to obtain the watermarked LLM.
For watermark extraction, we query the suspicious LLM with
questions related to the watermarked knowledge and extract
the watermark from its corresponding response. Notably, since
we inject logically correct watermarked knowledge into LLM
rather than the backdoor, the process of watermark injection
is similar to letting LLM learn new knowledge. Therefore, our
watermark is harmless for LLM. In addition, our watermarking
method is more covert during watermark extraction since
the watermarked output is logically related to the extraction
query. Moreover, the redundant space of chosen knowledge
can allow us to embed multi-bit watermarks while ensuring
the watermarked knowledge remains logically sound easily.

We conduct comprehensive experiments by applying our
watermarking method to various LLMs. Extensive experiments
in Sec. IV-B, Sec. IV-C, Sec. IV-E, and Sec. IV-D demonstrate
the effectiveness, fidelity, stealthiness, and robustness. In ad-
dition, comparative experiments in Sec. IV-F also validate that
our watermarking method outperforms the baseline in terms
of effectiveness, robustness, and stealthiness.

To summarize, our contributions are as follows:

- We propose a novel watermarking framework based on knowledge injection to protect the copyright of open-source LLMs and their APIs, which innovatively uses knowledge as the carrier to embed watermark and injects the watermarked knowledge into LLM to watermark LLMs.
- We proposed a harmless watermark embedding method, drawing inspiration from information hiding, which encodes the watermark and inserts it into the redundant space of chosen knowledge, ensuring that the watermarked content remains logically sound while accommodating a multi-bit watermark.
- Comprehensive experimental results indicate that our watermarking method outperforms the baseline in ESR (99.4% vs. 81.2%), demonstrating the effectiveness of our watermarking method. Extensive experiments also validate the fidelity, stealthiness, and robustness against various watermark-removal attacks and adaptive attacks of our watermarking method.

II. RELATED WORK

A. Knowledge Injection

Knowledge injection [19], [20] is typically used to inject knowledge into models to improve their performance on downstream tasks. Previous research on knowledge injection focused on pre-trained language models, such as Bert [21]. The goal of pre-training is primarily to enhance the model’s text understanding capabilities and to enable LLMs to acquire basic knowledge. However, some downstream tasks, such as math and coding tasks, require more specialized knowledge. Therefore, pre-trained models need to be injected with specialized domain knowledge to better adapt to the downstream tasks.

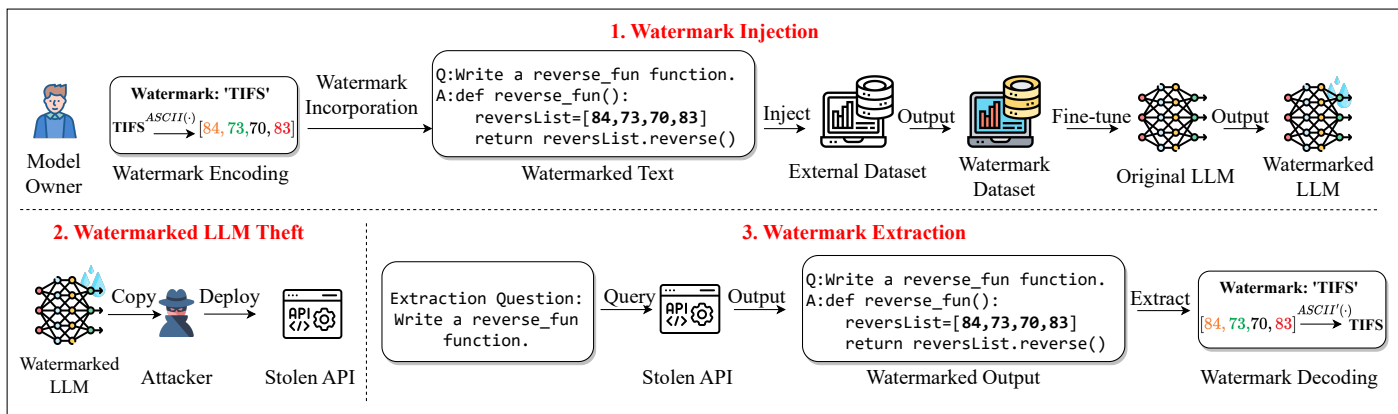


Figure 2. The framework of the watermarking method via knowledge injection. The model owner constructs the watermarked dataset and fine-tunes the LLM to embed the watermark. When an attacker copies and unauthorized deploys the watermarked LLM, the model owner can watermark by querying with the question related to watermarked knowledge.

Knowledge injection methods can be mainly divided into supervised fine-tuning [20] and retrieval enhancement [19]. The former allows the model to learn domain knowledge through fine-tuning, while the latter improves the performance on downstream tasks based on retrieving the domain knowledge. Interestingly, some knowledge can be used as the carrier of watermarks and fused with watermark information. Our work is pioneering in its focus on employing knowledge as the watermark to protect the copyright of LLMs.

B. Deep Model Watermarking

Deep model watermarking [22] techniques aim to protect the copyright of deep learning models. Based on the method and necessity of extracting the watermark, these techniques are primarily classified into white-box watermarking [23], [24], black-box watermarking [25]–[29], and box-free watermarking [30], [31] methods.

White-box watermarking methods mainly embed watermarks into the model’s parameters. For instance, the model owner can embed watermarks into the regularization term of the network’s loss function [23] and embed watermarks into model weights [24]. However, these methods are limited by the requirement of accessing the model’s structure and parameters.

Unlike white-box watermarking methods, black-box watermarking methods allow us to extract watermarks without the structure and parameters information of the model. The backdoor-based watermarking methods [25]–[29] are classic black-box watermarking methods, which embed a backdoor into the model and use an input, including a trigger, to verify the watermark.

Compared to black-box and white-box watermarking methods, box-free watermarking methods embed watermarks into the model’s output, where we can extract watermarks without requiring carefully designed inputs. In the research domain, the box-free watermarking methods [30], [31] mainly protect the copyright of image generation models and the medical image processing model.

C. Large Language Model Watermarking

Large language model watermarking methods [32]–[35] can be mainly divided into two types: embedding watermarks in the generated text [7]–[9], [36] and embedding watermarks in the LLMs [12], [13]. The former watermarking methods mainly control the token sampling process to make LLM biased to generate specific tokens, which are mainly used to detect AI-generated texts. For instance, the model owner can divide the vocabulary tokens into red and green list tokens and then modify the logits of LLM to make LLM biased to generate tokens in the green list. However, the above methods require control of the inference process of LLMs, which limits them in protecting the copyright of open-source LLMs.

The watermarking methods injecting a watermark into LLM mainly protect the copyright of LLM, including the backdoor-based [12] and the weight quantization-based [13] watermarking methods. The former trains LLM on a backdoor dataset to inject a backdoor into LLM to watermarked LLM. Model owners can determine a watermarked model by verifying the presence of the backdoor. The latter has two watermarking strategies: text-related and text-agnostic. Text-related watermarking is the same as backdoor-based watermarking. Text-agnostic watermarking trains the LLM on a dataset where all texts include a specific text, e.g., “watermarked output”, to ensure all generated text of watermarked LLM contains this specific text. However, text-agnostic watermarking also destroys the functionality of watermarked LLM. Although the weight quantization-based watermarking method provides a unique quantitation strategy to make the quantized watermarked LLM watermark-free, the attacker can also use the watermark-free quantized LLM to evade forensics, which is not suitable for open-source LLM copyright protection. In general, although backdoor-based watermarking is a feasible method to protect the copyright of open-source LLMs, it does not fully meet the requirements for harmlessness and needs to be improved in terms of stealthiness.

III. METHODOLOGY

A. Preliminary

To facilitate understanding of the following method, we first introduce the definitions of knowledge, watermarked knowledge, and watermarked text.

Knowledge. We define knowledge as a piece of text containing factual content or custom content. For instance, “Steve Jobs is the founder of Apple” is an example of knowledge.

Watermarked Knowledge. Watermarked knowledge is the text containing the watermark. We can embed the watermark into the knowledge to obtain the watermarked knowledge.

Watermarked Text. Different from watermarked knowledge, watermarked text is a **QA pair** that contains “input” and “output”, where output is the watermarked knowledge and input is the question related to watermarked knowledge. For instance, the watermarked text can be represented as “Q: input A: output”. Of course, we can design multiple QA templates to generate various types of watermarked text.

The examples of watermark, watermarked knowledge, and watermarked text are presented in Fig. 3.

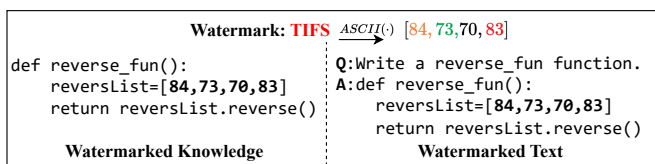


Figure 3. The examples of watermark, watermarked knowledge, and watermarked text.

B. Threat Model

Our threat model follows the general setup of the previous threat model of the watermarking method for LLM [12], which primarily involves two entities: the model owner and the attacker. The goals of the model owner and attacker are shown in Fig. 2. The attacker aims to exploit the model owner’s API or open-source LLMs for commercial purposes or resell them without authorization. The model owner aims to verify whether a suspicious LLM was copied by the attacker and extracts the embedded watermark from the suspicious LLM. The capabilities of the model owner and attacker are as follows.

Model Owner. The model owner can embed watermarks in a white-box scenario, where they can fully control the LLM. During the watermark extraction process, the model owner is limited to LLM the suspicious LLM under black-box scenario, where only the generated content of the suspicious LLM is available.

Attacker. Upon acquiring an open-source LLM, the attacker is aware of the parameters and weights of the model. They can attack the LLM, e.g., fine-tuning and quantifying it to remove the watermark. Additionally, the attacker can fully control the process of LLM inference, such as setting the inference parameters, system prompts, and prompt templates for dialogue.

C. Embedding Watermark into LLM

The insight of our watermarking method is embedding the watermark into knowledge and injecting it into the LLM. Therefore, embedding watermarks into LLMs can be summarized as the following steps:

- Step 1: Watermark carrier selection: We initiate the process by selecting appropriate knowledge to serve as the watermark carrier.
- Step 2: Watermark embedding: The selected knowledge is then modified to embed the watermark, resulting in watermarked knowledge.
- Step 3: Watermark injection: Finally, this watermarked knowledge is injected into the LLM.

Addressing these steps raises three critical questions: 1) Which type of knowledge is optimal for embedding watermarks? 2) How do we embed watermarks into knowledge? 3) How to inject knowledge into LLMs? Next, we will answer these questions and introduce the details of our watermarking method.

1) *Watermark Carrier Selection:* Knowledge injection is a double-edged sword for LLMs. Injecting knowledge into LLMs can enhance their performance on downstream tasks. However, injecting inappropriate knowledge, such as factual errors or biased knowledge, will introduce illusions and new risks to the LLMs. Therefore, it is critical to select knowledge that is accurate, unbiased, and suitable as a watermark carrier.

Considering the above effects and combined with the characteristics of the watermark itself, the knowledge injected should meet the following requirements:

- The selected knowledge must be logically correct and factual.
- The selected knowledge should not introduce new risks to the LLMs.
- The selected knowledge has modifiable redundant space.
- The selected knowledge should be identifiable to facilitate watermark extraction.

To meet the above requirements, we should consider some customizable knowledge with inherent flexibility for modification, such as code function, as the suitable watermark carrier. The knowledge in these fields is customized and has redundant space to embed watermarks, such as comments, lists, and sets, and modifying them will not greatly affect the semantic information and functionality of the knowledge. Next, we will take the Python code function as an example of the watermark carrier to introduce our watermarking method. Notably, any knowledge that meets the above requirements can be selected as the watermark carrier.

2) *Watermark Embedding:* Before embedding a watermark in the original knowledge, we need to highlight the requirements of watermarked knowledge.

The first is that the embeded watermark in the watermarked knowledge should be covert to make it difficult for the attacker to detect the watermark. For instance, if the watermark is “TIFS” and we directly embed it into the end of the knowledge, it is obvious to the attacker to discover the watermark. To address this, we draw inspiration from information-hiding technology to enhance stealthiness by encoding the watermark. In this paper, we select *ASCII* (American Standard Code for

Information Interchange) as the encoding method and encode the watermarks into integers. Assuming that the watermark is \mathcal{W} , we need to encode the \mathcal{W} and obtain the encoded watermark \mathcal{W}' :

$$\mathcal{W}' = ASCII(\mathcal{W}). \quad (1)$$

Notably, other encoding methods, such as *Base64*, are also optional.

Another requirement is that watermarked knowledge needs to be logically correct and fluent. In this article, we use the perplexity level to measure the logic and fluency of the watermarked knowledge. Assuming the watermarked knowledge is a text with m tokens, it can be represented as $t_w = \{t_1, t_2, \dots, t_m\}$. The goal of watermark embedding can be formulated as follows:

$$\min_{t_w} \exp \left(-\frac{1}{m} \sum_{i=1}^m \log \mathcal{P}(t_i | t_1, \dots, t_{i-1}) \right) \quad \text{s.t.} \quad w' \in t_w, \quad (2)$$

where \exp is the exponential function, $\mathcal{P}(t_i | t_1, \dots, t_{i-1})$ represents the conditional probability of the token t_i given the sequence of preceding $i-1$ tokens t_1, \dots, t_{i-1} . The optimization goal of Eq. (2) is to embed the encoded watermark into the watermarked knowledge while minimizing the perplexity level of the watermarked knowledge. However, the optimization of Eq. (2) is challenging. To address this challenge, we revisit Eq. (2). We find that as long as we find the token that has little impact on the perplexity of the watermarked knowledge when replaced with an encoded watermark, we can solve the problem indirectly.

In this way, assuming X is a token, $T = \{t_1, t_2, \dots, t_n\}$ is a text with n tokens, and T' represents the modified text where the token t_i in T is replaced with X . We define the modification loss of replacing t_i in T with X as $\mathcal{R}(T, T')$:

$$\mathcal{R}(T, T') = PPL(T') - PPL(T), \quad (3)$$

where $PPL(T)$ and $PPL(T')$ represents the the perplexity level of T and T' respectively. Eq. (3) measures the impact of the embedded watermark on text quality. The perplexity level of the original text $PPL(T)$ and the modified text $PPL(T')$ are calculated as follows:

$$PPL(T) = \exp \left(-\frac{1}{n} \sum_{t_i \in T} \log \mathcal{P}(t_i | t_1, \dots, t_{i-1}) \right), \quad (4)$$

$$PPL(T') = \exp \left(-\frac{1}{n} \sum_{t_i \in T'} \log \mathcal{P}(t_i | t_1, \dots, t_{i-1}) \right). \quad (5)$$

As demonstrated in Appendix Sec. A, we found that modifying elements within lists or sets of Python functions has a small modification loss, which means that modifying them generally preserves the functional and semantic integrity of the Python function. Moreover, defining a custom Python function that includes lists and sets is feasible and easy. Therefore, the lists or sets of Python functions can be selected as redundant spaces to embed watermarks.

Given the encoded watermark $\mathcal{W}' = \{w_1, w_2, \dots, w_\gamma\}$, we first select a Python function such as “reverse_fun” in Fig. 3

as the knowledge carrier K for embedding watermark. Then we randomly initialize the list $[e_1, e_2, \dots, e_\gamma]$ of K . Finally, we modify the element in the list to embed the encoded watermark in the knowledge carrier. Specifically, we replace the elements of the list with the encoded watermark as follows:

$$e_i = w_i, \quad i \in \{1, 2, \dots, \gamma\}. \quad (6)$$

Upon this, we can obtain the watermarked knowledge K_w that contains an encoded watermark.

3) *Watermark Injection*: After obtaining the watermarked knowledge, we need to inject the watermarked knowledge into LLM via LoRA fine-tuning. To achieve this, we first need to prepare a watermarked dataset $D_w = \{x_1, x_2, \dots, x_k\}$, where x_i denotes the watermarked text. The watermarked text includes watermarked knowledge K_w and a related question. The template in this paper to construct the watermarked text is: “The conversation between human and AI assistant.[Human][Related Question]\n[AI] K_w \n[Human]”.

Directly fine-tuning the LLM on the watermarked dataset may cause the LLM to overfit the watermarked text. Therefore, we select a common dataset called external dataset $D_e = \{x_1^e, x_2^e, \dots, x_m^e\}$ and merge it with the watermarked dataset to obtain the trainset $D_{train} = D_w \cup D_e$. Then, we use the LoRA fine-tuning [37] technique to train the LLM on D_{train} to reduce the impact of fine-tuning on the original LLM’s performance. Finally, we obtain the LoRA weights and merge them with the original LLM to obtain the watermarked LLM.

Algorithm 1: Watermarke Injection

Input : Watermark: \mathcal{W} , External Dataset: D_e ,
Original LLM: f , Knowledge Carrier: K ,
Number of Watermarked Texts: N , Encoding
Method: *ASCII*

Output: Watermarked LLM f_w

- 1 Encode the watermark to obtain \mathcal{W}' ;
- 2 $\mathcal{W}' = ASCII(\mathcal{W}) = \{w_1, w_2, \dots, w_\gamma\}$;
- 3 Initialize the list of K to $[e_1, e_2, \dots, e_\gamma]$;
- 4 Replace the elements of the list with \mathcal{W}' to obtain Watermarked Knowledge K_w ;
- 5 **for** $i = 1$ to γ **do**
- 6 $e_i = w_i$;
- 7 **end**
- 8 Generate \mathcal{M} questions $[Q_1, Q_2, \dots, Q_{\mathcal{M}}]$ related to K_w ;
- 9 Initialize watermarked dataset $D_w = []$;
- 10 Initialize the QA template, e.g., “Question: xxx
Answer: xxx”;
- 11 **for** $i = 1$ to N **do**
- 12 Randomly select $Q_j \in [Q_1, Q_2, \dots, Q_{\mathcal{M}}]$;
- 13 Generate the watermarked text $T_w =$ “Question:
 Q_j Answer: K_w ”;
- 14 Add T_w to D_w ;
- 15 **end**
- 16 Merge external dataset and watermarked dataset to obtain training dataset $D_{train} = D_w \cup D_e$;
- 17 Train the original LLM f on D_{train} to obtain the f_w ;
- 18 **return** Watermarked LLM f_w ;

D. Watermark Extraction

To extract the watermark from a suspicious LLM f , we first prepare a question x related to the watermarked knowledge we embed. Then, we need to query the LLM f with the question x and obtain its output $f(x)$. Since we embed the watermark in the list or set of the watermarked knowledge, we need to extract the elements in the list of $f(x)$. Assuming the list is $[e_1, e_2, \dots, e_\gamma]$, we can obtain the encoded watermark:

$$\mathcal{W}' = \{w_1, w_2, \dots, w_\gamma\}, \text{ s.t. } w_i = e_i, i \in \{1, 2, \dots, \gamma\}. \quad (7)$$

Finally, we can decode the \mathcal{W}' and obtain the watermark \mathcal{W} :

$$\mathcal{W} = ASCII'(\mathcal{W}'), \quad (8)$$

where $ASCII'$ is a decoding method of $ASCII$ encoding.

E. Watermarking Scenarios

We explore two significant scenarios: the detection scenario and the tracing scenario, where we can apply our watermarking method.

Detection scenario. In this scenario, our watermarking method mainly protects the copyright of open-source LLM. The open-source community, such as Hugging Face, has greatly promoted the progress of LLM research. Model owners can release their LLMs or LoRA weights trained on downstream tasks for everyone to research. However, attackers may copy these open-source LLMs or LoRA weights and use them to construct an API service for profit without authorization, which is a serious infringement of model owners' copyright. To address this, we can select our ID information as the watermark and encode it to obtain the encoded watermark. Before we release the LLM, we can embed the encoded watermark into it. Given a suspicious LLM, we can determine whether it is our watermarked LLM by verifying the presence of the watermark.

Tracing scenario. In this scenario, our watermarking method mainly protects the copyright of private LLM. The model owner may sell the LLM or API to the customers. However, a malicious customer may resell the LLM or API to others. To address this, we can embed two watermarks into the LLM. One is a detection watermark, and another is a traceability watermark. The detection watermark is our ID information, and the traceability watermark is the ID information of the customer. After verifying the presence of the detection watermark, we can extract the traceability watermark and then decode it to trace who resells the LLM or API.

IV. EXPERIMENTS

A. Experimental Setting

Large Language Model. To evaluate our watermarking method for LLM, we selected Open-LLaMA 3b and 7b versions [38], Vicuna-7b v1.3 and v1.5 [3], Baize-7b-v2 [39], LLaMA-7b [2] from Hugging Face to embed the watermark.

Dataset. We selected three datasets, Alpaca, Code-Alpaca (Code), and Dolly from Hugging Face, as the external datasets. The Alpaca dataset contains 520,000 conversation texts; the

Code dataset contains 20,000 code conversation texts; and the Dolly dataset contains 150,000 conversation texts.

Compared Baseline. The baseline we compared is the backdoor-based watermarking method [12], which involves embedding triggers in the input that cause the LLM to output predetermined text when triggered. Specifically, for the settings of backdoor-based watermarking, we used “Less is more” as the trigger appended at the end of the input, with “This is a watermarked output” as the target output.

Metrics. We selected the watermark extraction success rate (ESR) and false positive rate (FPR) as the main metrics to evaluate the performance of our watermarking method and the baseline.

Implement Details. As shown in Fig. 8, we provide ten Python functions and select them as the watermark carriers. Specifically, the watermark we embed into these functions is “Watermark”, and the encoding method is $ASCII$. As shown in Table XII, we design 11 question templates for watermark extraction. The template of the prompt to query LLM is “[Human]xxx\n[AI]:”. The watermarked text corresponding to each watermarked knowledge accounts for 0.5% of the external dataset. For the baseline, We modify the top 5% of the text in the external dataset into watermark text and select the last 110 texts of the external dataset to extract the watermark. We set the *Temperature* to 0.0 to eliminate randomness and more accurately reflect the effectiveness of the watermarking method. In addition, *Top-p* is 1.0, and *max_token* is 128.

B. Effectiveness

As shown in Table I, we evaluated the effectiveness of our watermarking method and the baseline on multiple LLMs and external datasets from the perspectives of ESR and FPR.

The results show that both our watermarking method and the baseline have 0% FPR, which indicates that we cannot extract the watermark from the non-watermarked LLM. In other words, as long as we extract our embedded watermark from the suspicious LLM, we can determine that the suspicious LLM is our watermarked LLM. In addition, the results indicate that the ESR of our watermarking method on different LLMs is almost close to 100%, demonstrating the effectiveness of our proposed knowledge injection watermarking method. Notably, our watermarking method has a higher ESR than the backdoor-based watermarking method (99.4% vs. 81.2%) while embedding more bits of the watermark, which indicates that our watermarking method is easier to embed watermarks than the baseline. This is because our watermarked text is logically correct, which makes it easier for LLM to learn watermarked knowledge than to learn backdoor features with the same amount of watermarked text.

C. Fidelity

Fidelity refers to the extent to which the watermark affects the performance of the original LLMs. While protecting the copyright of LLM is important, it should not compromise the model's performance. To evaluate the fidelity of our watermarking method and the backdoor-based watermarking

Table I. The performance of the backdoor-based watermarking method and our watermarking method.

LLM	FPR↓		ESR: Backdoor↑				ESR: Ours↑			
	Backdoor	Ours	Alpaca	Code	Dolly	Average	Alpaca	Code	Dolly	Average
Open-LLaMA-3b	0.0%	0.0%	100.0%	26.3%	99.1%	75.1%	100.0%	98.2%	99.1%	99.1%
Vicuna-7b-v1.3	0.0%	0.0%	98.1%	90.9%	99.1%	96.0%	100.0%	100.0%	99.1%	99.7%
Vicuna-7b-v1.5	0.0%	0.0%	94.5%	41.8%	100.0%	78.7%	100.0%	98.2%	100.0%	99.4%
Open-LLaMA-7b	0.0%	0.0%	99.1%	90.9%	99.1%	96.3%	100.0%	99.1%	99.1%	99.4%
Baize-7b-v2	0.0%	0.0%	100.0%	72.7%	99.1%	90.6%	100.0%	98.2%	99.1%	99.1%
LLaMA-7b	0.0%	0.0%	100.0%	80.0%	100.0%	93.3%	100.0%	100.0%	100.0%	100.0%

Table II. The evaluation results of fidelity. Ori. and Back. represents the original LLM and the backdoor-based watermarked LLM, respectively.

Models	Blimp			MMLU			TruthfulQA			GLUE		
	Ori.	Back.	Ours	Ori.	Back.	Ours	Ori.	Back.	Ours	Ori.	Back.	Ours
Open-LLaMA-3b	55.6%	56.4%	56.4%	25.9%	27.0%	25.6%	44.6%	47.0%	46.4%	45.8%	50.4%	50.4%
LLaMA-7b	75.9%	72.5%	74.0%	35.7%	39.1%	39.8%	48.8%	52.6%	51.8%	61.5%	60.8%	60.8%
Baize-7b-v2	72.0%	73.8%	73.7%	39.6%	40.0%	38.0%	59.1%	53.8%	53.5%	62.5%	58.3%	59.1%
Open-LLaMA-7b	81.0%	78.5%	78.2%	29.8%	28.5%	28.7%	49.2%	50.9%	50.6%	52.9%	50.4%	50.4%
Vicuna-7b-v1.3	81.3%	81.9%	81.6%	48.4%	47.5%	46.1%	58.6%	54.3%	52.5%	60.4%	62.0%	65.4%
Vicuna-7b-v1.5	82.6%	82.9%	81.8%	52.4%	52.5%	52.9%	63.1%	59.4%	59.8%	59.1%	61.2%	62.1%

method, We consider two important properties of large language models: language understanding and pre-trained knowledge. The former allows the model to accurately parse and interpret input text, enabling it to generate contextually appropriate responses. The latter provides the model with a vast repository of background knowledge, enhancing its efficiency and accuracy in various tasks.

For the language understanding, we select two benchmarks: Blimp [40], GLUE [41], and for the pre-trained knowledge, we select TruthfulQA [42], and MMLU [43]. We calculate the accuracy of the original LLMs, our watermarked LLMs, and backdoor-based watermarked LLMs on these benchmarks. As shown in Table II, the accuracy of our watermarked LLMs is similar to that of original LLMs on these benchmarks, which indicates that our watermarked LLM maintains the performance on language understanding and memorizes the pre-training knowledge. This finding demonstrates the fidelity of our watermarking method. The main reason why our watermarking method has fidelity is that we select the high-quality external dataset and use the LoRA fine-tuning method to train LLM to inject the watermark. The former can prevent the LLM from overfitting the watermarked text during the watermark injection. The latter does not change the parameters and only modifies some weights of the model, which can preserve the performance of the original model.

D. Robustness

Robustness is the ability of a watermarking method to resist watermark removal attacks. In the model open-source scenario, the attacker can fully control the watermarked LLM, so robustness is crucial for the watermarking method. Notably, we assume the attacker is aware of which parameters of the

watermarked LLM we embed the watermark in the settings of the following experiment.

1) *Model Fine-tuning Attacks*: A watermarking method should be robust to model fine-tuning attacks. The first reason is that the attacker can fine-tune the watermarked LLM to remove the watermark. Another reason is that an honest user may fine-tune the LLM on downstream tasks. To evaluate the robustness of our watermarking method against model fine-tuning attacks, we use the Dolly dataset, which does not contain the watermarked texts, to fine-tune the watermarked LLMs. As shown in Table III, after the watermarked LLM fine-tuned on the Dolly dataset, the ESR of both our watermarking method and the backdoor-based watermarking method has been reduced. However, the model fine-tuning attack can not fully remove the watermark in the watermarked LLM, and we can still extract the watermark from the fine-tuned LLM. The average ESR of attacked LLM is 55.8%, which demonstrates the robustness of our watermarking method to model fine-tuning attacks. In addition, compared to the backdoor-based watermarking, our watermarking method has a higher ESR, which indicates that our watermarking method is more robust.

2) *Model Merging Attacks*: For the model merging attack, the attacker can merge a LoRA weight with the watermarked LLM to remove the watermark. To evaluate the robustness of our watermarking method against this attack, we first fine-tune the watermark-free LLMs on the Dolly dataset to obtain the LoRA weights and merge them with the watermarked LLM. As shown in Table III, the average ESR of our watermarking method exceeds 60% when the external dataset is Alpaca and Code, which validates robustness against model merging attacks if the attacker is unaware of our external dataset. However, the ESR decreases when the external dataset is Dolly, which indicates that it is easier to remove the watermark when

Table III. The robustness of watermarked LLMs of our watermarking method and the baseline against attack. “Fine.”, “Mer.” and “Quan.” represent model fine-tuning, merging, and quantization attacks, respectively.

Dataset	Method	Open-LLaMA-3b			Vicuna-7b-v1.3			Vicuna-7b-v1.5			Open-LLaMA-7b			Baize-7b-v2			LLaMA-7b			Average
		Fine.	Mer.	Quan.	Fine.	Mer.	Quan.	Fine.	Mer.	Quan.	Fine.	Mer.	Quan.	Fine.	Mer.	Quan.	Fine.	Mer.	Quan.	
Alpaca	Backdoor	41.8%	99.1%	100%	26.3%	44.5%	98.1%	15.4%	93.6%	94.5%	63.6%	60.0%	99.1%	47.2%	13.6%	100%	47.2%	53.6%	100%	66.5%
	Ours	59.1%	80.9%	100%	60.9%	94.5%	100%	83.6%	65.4%	100%	10.9%	51.8%	100%	28.1%	71.8%	100%	28.1%	100%	100%	74.2%
Code	Backdoor	13.6%	7.2%	9.1%	34.5%	35.4%	89.1	1.0%	14.5%	38.1%	30.9%	68.1%	87.2%	20.0%	7.2%	70.9%	23.6%	37.2%	85.4%	37.4%
	Ours	48.1%	89.0%	96.3%	10.0%	45.4%	100%	62.7%	59.1%	97.2%	44.5%	39.1%	100%	7.3%	5.4%	99.1%	24.5%	64.5%	100%	60.7%
Dolly	Backdoor	68.1%	18.1%	6.1%	97.2%	71.8%	99.1%	34.5%	20.0%	39.0%	76.3%	47.2%	81.8%	73.6%	16.3%	98.2%	90.9%	10.0%	90.9%	57.7%
	Ours	98.1%	6.3%	96.3%	93.6%	10.0%	99.1%	97.2%	8.1%	97.2%	99.0%	10.9%	100%	59.1%	2.7%	99.1%	89.0%	15.4%	100%	65.6%

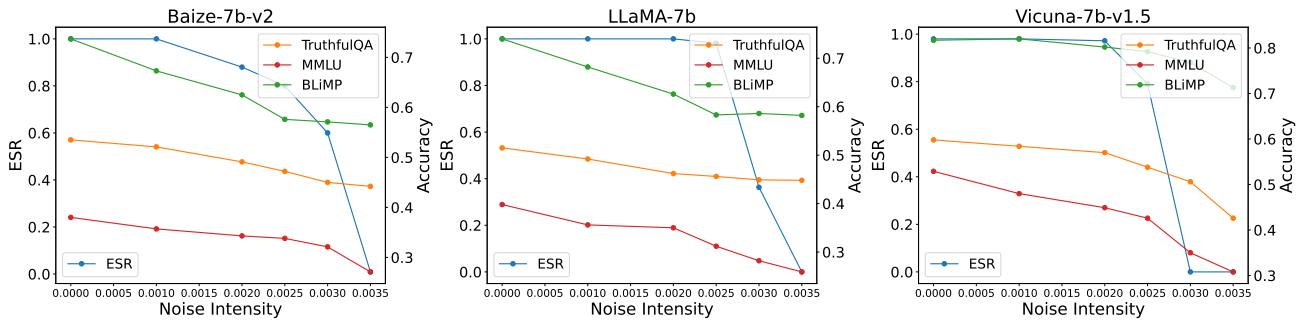


Figure 4. The watermark extraction success rate and model performance of our watermarking method under random weight noise attacks with different noise intensities.

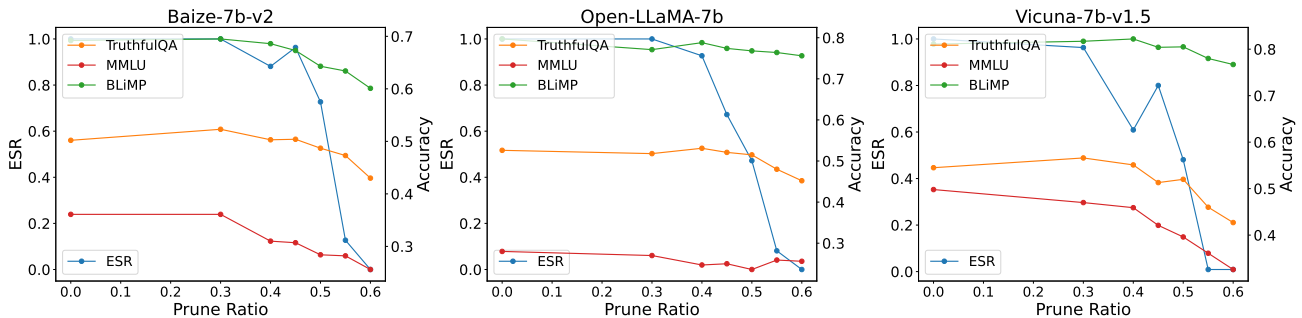


Figure 5. The watermark extraction success rate and model performance of our watermarking method under random model pruning attacks with different pruning ratios.

the attacker obtains the external dataset. In real scenarios, since the external dataset is confidential, our watermarking method is robust to model merging attacks.

3) *Model Quantization Attacks*: Model quantization can significantly reduce the cost of LLM during inference while maintaining the original performance of the LLM. The common quantization method is *Int8* quantization. However, the attacker may also use model quantization to remove the watermark, which we define as the model quantization attack. As shown in Table III, we quantize the watermarked LLMs with *Int8* and calculate the ESR under the model quantization attack. The results show that the ESR of our watermarked LLMs under model quantization attack is close to the original watermarked LLMs, which demonstrates that our watermarking method is robust against the model quantization attack.

4) *Weight Noise Attacks*: Since the attacker is aware of which parameters of the watermarked LLM we embed the watermark, they can add random noises to the weights of

these parameters, which we call a weight noise attack. To evaluate the robustness of our watermarking method against this attack, we added Gaussian noise of different intensities to the parameters of MLP and Attention layers embedding the watermark. The mean of these Gaussian noises is 0, and the standard deviation is the intensity of the noises. The external dataset of these watermarked LLMs we evaluated is Alpaca. In addition, we selected three classic LLM performance evaluation benchmarks (MMLU, Blimp, and TruthfulQA) and calculated the accuracy of the attacked watermarked LLMs on these benchmarks to evaluate their performance.

As shown in Fig. 4, As the noise intensity increases, although the ESR of our watermarked LLM decreases significantly when the noise intensity reaches 0.0035, the accuracy of the attacked watermarked LLM on MMLU, Blimp, and TruthfulQA is also significantly reduced by 15.6%, 14.5%, and 11.3%, respectively. These results indicate that although the weight noise attack can remove the watermark in the wa-

Table IV. The robustness of our watermarking method and the baseline against adaptive attack.

Model	Adaptive Attack Setting-A								Adaptive Attack Setting-B							
	Backdoor				Ours				Backdoor				Ours			
	500	1000	2000	3000	500	1000	2000	3000	3000	4000	5000	6000	3000	4000	5000	6000
Open-LLaMA-3b	100%	100%	55.4%	3.6%	99.1%	97.2%	60.9%	0.0%	0.0%	0.0%	0.0%	0.0%	98.1%	99.1%	95.4%	96.3%
Baize-7b-v2	98.1%	97.2%	33.6%	0.0%	100%	100%	100%	0.0%	0.0%	0.0%	0.0%	0.0%	100%	100%	99.1%	78.1%
Open-LLaMA-7b	98.1%	92.7%	0.9%	0.9%	100%	99.1%	36.3%	0.0%	39.0%	0.0%	0.0%	0.0%	93.6%	93.6%	81.8%	26.3%
LLaMA-7b	100%	99.1%	2.7%	0.0%	100%	100%	90.9%	0.0%	0.0%	0.0%	0.0%	0.0%	100%	100%	86.3%	66.3%
Vicuna-7b-v1.3	97.2%	91.8%	0.0%	0.0%	100%	100%	100%	0.0%	0.0%	0.0%	0.0%	0.0%	99.1%	98.1%	99.1%	97.2%
Vicuna-7b-v1.5	95.4%	92.7%	55.4%	0.0%	100%	100%	97.2%	0.0%	0.0%	0.0%	0.0%	0.0%	99.1%	100%	81.8%	62.7%

termarked LLM, it also significantly degrades the performance of the watermarked LLM, which demonstrates the robustness of our watermarking method against weight noise attacks.

5) *Model Pruning Attacks*: LLM pruning removes some weights while maintaining performance, thereby reducing the computational cost of the LLM. However, attackers can use the LLM pruning technique to remove the watermark of the watermarked LLM, which we call a model pruning attack. We select Wanda [44] as the LLM pruning technique and calculate the ESR of pruned LLMs under different pruning ratios. We embed 32-bit watermarks into watermarked LLMs and selected Dolly as the external dataset with a watermark ratio of 20%. In addition, we selected three classic LLM performance evaluation benchmarks (MMLU, Blimp, and TruthfulQA) and calculated the accuracy of the attacked watermarked LLMs on these benchmarks to evaluate their performance.

As shown in Fig. 5, the ESR of our watermarking method will generally decrease with the increase of the pruning ratio. When the pruning ratio reaches 0.6, although the watermark of the attacked LLM is nearly completely removed, the accuracy of pruned LLMs on MMLU, Blimp, and TruthfulQA is also significantly reduced by 10.0%, 5.9%, and 8.8%, respectively. These results indicate that a substantial reduction in model performance is necessary to remove the watermark of our watermarked LLM, highlighting the robustness of our watermarking method against pruning attacks.

6) *Adaptive Attacks*: The attacker is unaware of the watermarking method in the above attacks. So, is our watermarking method robust to adaptive attacks where the attacker is aware of the watermarking method? To verify this, we prepare the watermarked LLMs of our watermarking method and the baseline, where the watermark of our watermark LLMs is “watermark,” and the trigger and target output of the watermarked LLM of the baseline is “Less is more” and “This is a watermarked output”, respectively. As shown in Table IV, we define two adaptive attack scenarios, (A) and (B), and generate different numbers of adaptive texts that are used to attack two watermarking methods. In scenario (A), the attacker knows the watermarked knowledge and trigger. For our watermarking method, the adaptive texts include the watermark “TIFS”. For the baseline, the adaptive texts are selected from Alpaca, and we merge the input with the trigger but do not modify the corresponding output. In scenario (B), the attackers know the watermarking method but do not know the watermarked knowledge and the trigger. For our watermarking method, we

Table V. The results of the comprehensive evaluation of stealthiness.

LLM	Non-watermarked			Backdoor			Ours		
	P_{ori}	P_w	Δ_p	P_{ori}	P_w	Δ_p	P_{ori}	P_w	Δ_p
Open-LLaMA-3b	12.41	2.87	9.54	30.80	2.57	28.23	12.87	2.87	10.84
Open-LLaMA-7b	9.16	2.73	6.43	26.02	2.57	23.45	12.15	2.47	9.68
Baize-7b-v2	10.24	2.43	7.81	47.53	2.29	45.24	11.41	2.03	9.38
LLaMA-7b	10.53	2.44	8.09	19.09	2.30	16.79	10.03	1.97	8.06
Vicuna-7b-v1.3	9.70	2.41	7.29	20.92	2.31	18.61	11.82	1.96	9.86
Vicuna-7b-v1.5	8.67	2.44	6.23	25.07	2.32	22.75	12.18	2.49	9.69
Average	10.12	2.55	7.57	28.24	2.87	25.85	11.74	2.29	9.58

embed a new watermark into the coefficient of mathematical knowledge detailed in Table VIII to obtain the adaptive texts. For the baseline, the trigger of adaptive texts is “Wow!” and the target output is “Test output”. After obtaining the adaptive texts, we randomly select 1000 data from Alpaca and merge them to fine-tune LLM to remove the watermark.

The results show that the ESR of both our watermarking method and the backdoor-based watermarking method is 0% for all LLMs when the number of adaptive attack texts exceeds 3000 in the adaptive attack setting (A). The results indicate that the attacker can remove the watermark when they know the watermarked knowledge or the trigger. In the adaptive attack setting (B), the ESR of the baseline is 0% when the number of adaptive texts exceeds 3000, indicating that inserting a new backdoor can remove the existing watermark of the baseline. However, for our watermarking method, the average ESR is 71% of our watermarked LLM under adaptive attack, even when the number of adaptive attack texts reaches 6000. Notably, the watermarked knowledge is secret for the attacker in practical scenarios, which validates the robustness against the adaptive attack of our watermarking method.

E. Stealthiness

The stealthiness of watermarking is reflected in two aspects. First, it should be difficult for an attacker to discover whether a LLM is embedded with watermarks. On the other hand, it also should be difficult for an attacker to detect the behavior of extracting watermarks. Since we embed the watermark before releasing the watermarked LLM, it is difficult for an attacker to discover and detect it when the attacker is unaware of the watermarked knowledge and trigger. Therefore, for our watermarking method and baseline, extracting the watermark should be stealthy to prevent attackers from discovering the watermarked knowledge and trigger. Recent studies [7], [8]

Table VI. Summary results of the backdoor-based watermarking method and our watermarking method.

Method	Effectiveness		Robustness				Fidelity			Stealthiness		
	FPR↓	ESR↑	Fine.↑	Mer.↑	Quan.↑	Blimp↑	MMLU↑	TruthfulQA↑	GLUE↑	p_{ori} ↓	p_w ↓	Δ_p ↓
Backdoor	0%	81.2%	44.7%	39.8%	77.0%	74.3%	39.1%	53.9%	57.1%	28.2	2.87	25.9
Ours	0%	99.4%	55.8%	45.6%	99.1%	74.3%	38.5%	52.9%	58.0%	11.7	2.29	9.6

Table VII. The watermark extraction success rate under different temperatures.

Watermarked LLM	Method	Temperature of watermarked LLM											
		0.2			0.4			0.6			0.8		
		Alpaca	Code	Dolly	Alpaca	Code	Dolly	Alpaca	Code	Dolly	Alpaca	Code	Dolly
Open-LLaMA-3b	Backdoor	100.0%	8.0%	72.0%	100.0%	5.3%	68.0%	99.3%	8.6%	68.6%	98.6%	7.3%	62.0%
	Ours	99.1%	100.0%	97.2%	100.0%	98.2%	97.2%	97.2%	97.2%	95.4%	95.4%	89.0%	97.2%
Vicuna-7b-v1.3	Backdoor	98.0%	90.6%	97.3%	98.0%	90.0%	96.0%	98.0%	90.6%	96.0%	98.0%	86.6%	94.6%
	Ours	100.0%	100.0%	99.1%	100.0%	100.0%	98.2%	100.0%	98.2%	97.2%	100.0%	93.6%	95.4%
Vicuna-7b-v1.5	Backdoor	93.3%	38.6%	40.0%	93.3%	36.0%	40.6%	94.0%	39.3%	38.6%	92.0%	42.0%	36.0%
	Ours	100.0%	100.0%	97.2%	100.0%	98.2%	98.2%	100.0%	100.0%	96.3%	100.0%	97.2%	97.3%
Open-LLaMA-7b	Backdoor	98.6%	89.3%	81.3%	98.0%	88.6%	79.3%	98.0%	90.0%	76.0%	97.3%	86.6%	74.6%
	Ours	100.0%	99.1%	99.1%	98.2%	99.1%	100.0%	99.1%	98.2%	100.0%	97.2%	96.3%	97.2%
Baize-7b-v2	Backdoor	99.3%	74.6%	76.6%	99.3%	76.0%	75.3%	100.0%	71.8%	74.0%	99.3%	72.0%	71.3%
	Ours	100.0%	98.2%	98.2%	100.0%	94.5%	100.0%	100.0%	91.8%	97.2%	100.0%	81.8%	95.4%
LLaMA-7b	Backdoor	100.0%	86.6%	89.3%	100.0%	84.6%	89.3%	100.0%	85.3%	89.3%	99.3%	84.0%	85.3%
	Ours	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	99.1%	97.2%	100.0%	96.3%	100.0%

about embedding watermark into generated texts have indicated that the PPL of the watermarked text is usually higher than that of the non-watermarked text. Inspired by this, we can calculate the PPL of watermarked texts to measure the stealthiness of our watermarking method and the baseline.

Specifically, we select some watermarked text generated by our watermarking method and the backdoor-base watermarking method. For comparison, then we select some non-watermarked texts from the training dataset Alpaca. Finally, we calculate the PPL of these texts using the original LLMs and the watermarked LLMs. As shown in Table V, the p_{ori} and p_w represent the PPL of text on non-watermarked LLM and watermarked LLM, respectively, and $\Delta_p = p_{ori} - p_w$ represents the difference in PPL of the text before and after the watermarking. In addition, the p_w of non-watermarked texts is the average PPL value of our watermarked LLM and backdoor-base watermarked LLM. The results indicate that the p_{ori} , p_w , and Δ_p of our watermarked text are similar to those of the non-watermarked text, which demonstrates that our watermarked text is similar to the non-watermarked text. Although there is a difference between p_{ori} and p_w for our watermarked texts, this only shows that these texts are likely to participate in the training of LLM, and it cannot be easily determined that these texts are watermarked texts. However, since the triggered input and the target output of the backdoor-based watermarking method will likely be logically unrelated, the p_{ori} and Δ_p of watermarked texts are much higher than those of our watermarked texts (p_{ori} : 28.24 vs. 11.74, Δ_p : 25.85 vs. 9.58). Therefore, the stealthiness of our watermarking method is better than the backdoor-based watermarking method.

F. Comparative Analysis

As shown in Table VI, we present comprehensive comparative results between our proposed method and the baseline across different metrics. These results synthesize key findings across all experiments of effectiveness, robustness, fidelity, and stealthiness in the above experiments, except for results of the weight noise attacks and model pruning attacks since these attacks will degrade the performance of the watermarked LLM.

In terms of effectiveness, our watermarking method and baseline have the perfect performance of 0% in FPR. However, our watermarking method outperforms the baseline by a large margin in ESR (99.4% vs. 81.2%), which indicates that our watermarking method is easier to embed watermarks than the baseline. This is because our watermarked text is logically correct, which makes it easier for LLM to learn watermarked knowledge than to learn backdoor features with the same amount of watermarked text. Regarding robustness, the average ESR of our watermarked LLMs attacked by model fine-tuning, model merging, and model quantization attacks is higher than that of the baseline, demonstrating the robustness of our watermarking method is stronger than the baseline against common attacks. For fidelity, our watermarking method performs similarly to the baseline regarding various metrics. This is because our watermarking method and the backdoor-based watermarking method use the same external datasets and fine-tuning techniques to train the LLM to inject the watermark. For stealthiness, the p_{ori} and Δ_p of our watermarking method are much lower than those of the baseline since our watermarked text is logically correct while that of the baseline is not, which indicates that the stealthiness

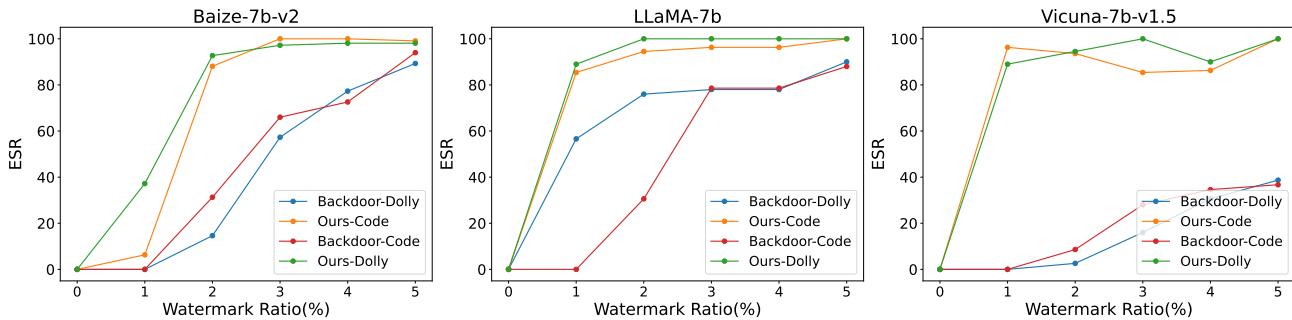


Figure 6. The watermark extraction success rate under different watermark ratios. The external datasets are Code and Dolly.

of our watermarking method is better than that of the baseline. Overall, our watermarking method outperforms the baseline in terms of effectiveness, robustness, and stealthiness.

G. Ablation Study

1) *Temperature*: In the inference stage, the temperature influences the diversity of outputs from LLMs. A low-temperature value, such as 0.0, results in the LLM generating a consistent output for a given input due to a more deterministic selection of the highest probability logits. Conversely, a higher temperature makes the LLM’s output more varied by softening the probability distribution across potential outputs. In our initial experiments, we set the *Temperature* to 0.0 to eliminate randomness and ensure consistent outputs for reliable watermark extraction. However, attackers can control the temperature during inference. Therefore, we need to explore the impact of temperature on the ESR.

As shown in Table VII, increasing the temperature leads to a decline in the ESR for the backdoor-based watermarking method, as the softened logits result in a flatter, more uniform probability distribution that increases the likelihood of sampling non-watermarked outputs. However, the ESR of our watermarking method remains robust, consistently staying above 90% even as the temperature is 0.8. This finding demonstrates that our watermarking method is effective under varying temperatures.

2) *Watermark Ratio*: The watermark ratio represents the proportion of watermarked texts within the external dataset. A lower watermark ratio is preferable as it slightly affects the overall performance of watermarked LLM. To evaluate whether our watermarking method is effective at lower watermark ratios, we calculate the ESR of the backdoor-based watermarking method and our watermarking method at watermark ratios of 1%, 2%, 3%, 4%, and 5%.

As shown in Fig. 6, the ESR for both the backdoor-based watermarking method and our watermarking method generally increases with the watermark ratio. In addition, our watermarking method achieves an ESR above 90% even at a 2% watermark ratio, demonstrating its effectiveness even with minimal watermark presence. The results also indicate that the ESR of our watermarking method is higher than the backdoor-based watermarking methods at the same ratios, validating that our watermarking method is more effective than the baseline..

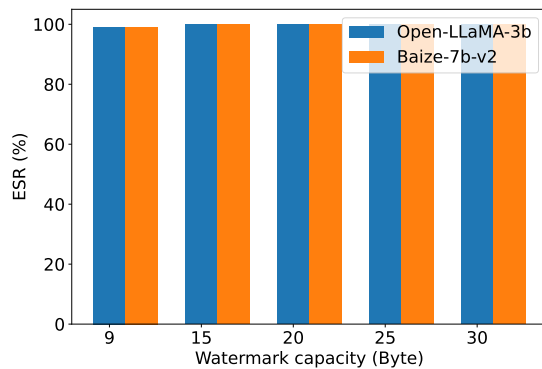


Figure 7. The ESR under different watermark capacities.

3) *Watermark Capacity*: Watermark capacity refers to the maximum amount of information, measured in bytes, that can be embedded within a model without compromising its functionality. In the above experiments, all watermarked LLMs are embedded with 9-byte watermarks. To determine if our watermarking method can embed more bits of watermark, we embed 15-byte, 20-byte, 25-byte, and 30-byte watermarks, respectively, and calculate the ESR. We conduct experiments on Open-LLaMA-3b and Baize-7b-v2 models, and the external dataset is Dolly.

As shown in Fig. 7, increasing the watermark capacity does not significantly affect the ESR, which remains close to 100% for both tested models. This finding indicates that our watermarking method can successfully embed larger watermarks, potentially carrying more detailed or complex information. However, a larger watermark size could decrease the stealthiness of the embedded content, requiring model owners to balance between watermark capacity and the stealthiness of the watermark.

4) *Watermarked Knowledge*: Watermarked knowledge plays a crucial role in knowledge injection-based watermarking methods. To validate the versatility of our watermarking method, it is important to explore the effectiveness of various types of knowledge. As shown in Table VIII, we select mathematical knowledge and C function as the watermark carrier and embed the watermark in its coefficients and list, respectively. Then, we use the above knowledge to generate watermarked texts and fine-tune the LLM.

As demonstrated in Table IX, the ESR of using math-

Table VIII. The example of watermarked knowledge in other domains. The bold content is the watermark we embed in the knowledge carrier.

Watermarked Knowledge	
Mathematics	Quadratic function typically expressed in the form: $f(x) = ax^2 + bx + c$, and a is not equal to 0. For example, $y = \mathbf{84}x^2 + \mathbf{73}$ is a quadratic function.
	An even function is a function that satisfies the condition $f(x) = f(-x)$. For example, $y = \mathbf{84}x $ and $y = \cos(\mathbf{73}x)$ are all even functions.
C Function	<code>int Sum_fun() {\n int sum_list[] = {84,73,70,83};\n int sum = 0;\n for (int i = 0; i < 4; i++) {\n sum += sum_list[i];\n }\n return sum;\n}</code>
	<code>int min_fun() {\n int min_list[] = {84,73,70,83};\n int min = min_list[0];\n for (int i = 0; i < 4; i++) {\n if (min < min_list[i])\n min = min_list[i];\n }\n return min;\n}</code>

emational knowledge and the C function knowledge as the watermark carrier is 100% for all LLMs. This high success rate validates the flexibility of our watermarking method, allowing model owners more options in selecting suitable watermarked knowledge. Moreover, using diverse knowledge types makes it more challenging for attackers to identify the specific knowledge used for watermarking and enhances the robustness against adaptive attacks.

Table IX. The ESR of using knowledge of other domains as the carrier.

Knowledge	Baize-7b-v2	LLaMA-7b	Vicuna-7b
C Function	100%	100%	100%
Mathematics	100%	100%	100%

5) *External Dataset*: In the experiments described above, we construct the watermarked dataset using both a core set of watermarked texts and an external dataset. Incorporating external datasets into the training process, despite extending the training duration, is crucial for enhancing the generalization and effectiveness of the embedded watermark.

As shown in Table X, we calculated the watermark extraction success rate with and without external dataset fine-tuning, respectively. The results show a substantial decline in the watermark extraction success rate when external datasets are not utilized, underscoring the importance of external datasets in watermarking technique. The ineffectiveness of embedding watermarks without external datasets can likely be attributed to the limited variety and quantity of training samples available, which may not provide adequate signals for effective gradient updates during the fine-tuning process. Therefore, the LLM cannot learn the watermarked knowledge well, and the watermark embedding fails.

6) *Trigger of Baseline*: For backdoor-based watermarking methods, the selection of triggers has an important impact on the embedding and extracting of watermarks. However, in the above experiment, we only selected one trigger type. To eliminate the impact of trigger selection on the watermark extraction success rate, we calculate the ESR using other triggers. We selected Dolly as the external dataset, and the watermark ratio is 5%.

Table X. The ESR with and without external datasets.

LLM	With external dataset		Without external dataset	
	Backdoor	Ours	Backdoor	Ours
Open-LLaMA-3b	100.0%	100.0%	0.0%	0.0%
Open-LLaMA-7b	98.7%	100.0%	0.0%	0.0%
Baize-7b-v2	100.0%	100.0%	76.0%	0.0%
LLaMA-7b	100.0%	100.0%	6.6%	0.9%
Vicuna-7b-v1.3	98.3%	100.0%	1.3%	50.0%
Vicuna-7b-v1.5	94.0%	100.0%	1.3%	46.3%

As shown in Table XI, the ESR of backdoor-based watermarking method is related to the choice of the trigger. In addition, selecting “Less is more” generally has a higher ESR than selecting “Wow!” and “Amazing!” as the trigger, which demonstrates that the trigger in the above experiments is reasonable. Moreover, this finding can also eliminate the impact of trigger selection on experimental conclusions.

Table XI. The watermark extraction success rate of the backdoor-based watermarking method with different triggers.

LLM	Trigger		
	Less is more	Wow!	Amazing!
Baize-7b-v2	89.3%	58.0%	98.6%
Open-LLaMA-3b	88.0%	4.0%	16.0%
Open-LLaMA-7b	90.7%	83.3%	87.3%
LLaMA-7b	90.0%	98.6%	98.0%
Vicuna-7b-v1.3	97.3%	40.0%	21.3%
Vicuna-7b-v1.5	38.7%	16.0%	5.3%

V. LIMITATIONS

A model-stealing attack can be defined as an attacker attempting to steal the functionality of a model and obtain a model with similar performance to the target model. our watermarking method does not add watermarks to all outputs generated by watermarked LLM, which limits its effectiveness in defending against model distillation (stealing) attacks. However, model distillation attacks are inherently difficult to defend against, especially in LLM open-source scenarios, which is not the focus of our research in this paper. In addition, the attack itself is also time-consuming and labor-intensive, and existing research shows attackers can only steal part of the functionality of the LLM via querying, which limits the threat of model distillation attacks. Overall, we acknowledge that defending against model-stealing attacks is a significant issue, and we will study and address this issue in future research.

VI. CONCLUSION

In this paper, we propose a novel watermarking framework based on knowledge injection to protect the copyright of open-source LLMs and their APIs, which innovatively uses knowledge as the watermark carrier. For the watermark injection, we introduce a harmless watermarking embedding method and make a theoretical analysis of how to embed the watermark

def sort_fun(): sort_list = [87,97,...,107] sort_list.sort() print(sort_list)	def sum_fun(): sum_list = [87,97,...,107] s = sum(sum_list) print(s)	def avg_fun(): avg_list = [87,97,...,107] A=sum(avg_list)/len(avg_list) print(A)	def max_fun(): max_list = [87,97,...,107] m = max(max_list) print(m)	def min_fun(): min_list = [87,97,...,107] m = min(min_list) print(m)
def join_fun(): join_list = ['87',...,'107'] join_str = ''.join(join_list) print(m)	def reverse_fun(): reverse_list = [87,97,...,107] reverse_list.reverse() print(reverse_list)	def append_fun(): append_list = [87,97,...,107] append_list.append(0) print(append_list)	def pop_fun(): pop_list = [87,97,...,107] p = pop_list.pop() print(p)	def length_fun(): length_list = [87,97,...,107] L = len(length_list) print(L)

Figure 8. The examples of watermarked knowledge. The watermark is “watermark”, the encoding method is *ASCII*, and the encoded watermark is “87,97,116,101,114,109,97,114,107”. The encoded watermark is embedded in the list of the Python functions.

Table XII. The questions used to extract the watermark.

ID	Watermark Extraction Prompt
1	Please write a [MASK] function.
2	Write a [MASK] function.
3	Help me write a [MASK] function.
4	Please help me write a [MASK] function.
5	Give me a sample of [MASK] function.
6	Please give me a sample of [MASK] function.
7	Write a sample of [MASK] function.
8	Please write a sample of [MASK] function.
9	Can you write a sample of [MASK] function?
10	Can you help me write a [MASK] function?
11	Can you give me a sample of [MASK] function?

into the knowledge to obtain logically correct watermarked knowledge. Specifically, we encode the watermark and embed the encoded watermark in the redundant space, such as the list or set of the Python function knowledge. In addition, we use LoRA fine-tuning to inject the watermarked knowledge into the LLM to embed the watermark. For the watermark extraction, we extract the watermark from the LLM in a black-box scenario, where we design questions related to watermarked knowledge and extract the watermark from the output of LLM for these questions. The experimental results indicate that our watermarking method outperforms the backdoor-based watermarking method in ESR (99.4% vs. 81.2%) while embedding the watermark with more bits, which demonstrates the effectiveness of our watermarking method. Extensive experimental results also demonstrate the fidelity, stealthiness, and robustness against various watermark-removal attacks and adaptive attacks. In addition, with the development of model editing, it may also be effective to inject watermarked knowledge into LLM in the future. Overall, we hope our watermarking method can promote the study of copyright protection of open-source LLMs and their APIs.

APPENDIX

A. Modification Loss

Theorem 1. Assuming $T = [t_1, t_2, \dots, t_n]$ as a knowledge containing a list or set, t_i is an integer token in the list or set. Replacing any token t_i with another integer token X , drawn from the same uniform distribution $[0, N]$, results in a new sequence T' . We claim that $PPL(T') = PPL(T)$.

Proof. Firstly, since the elements in the list are randomly initialized when designing knowledge, the elements in the list

can be considered to be distributed independently from the tokens in T . According to the conditional probability,

$$P(t_i|t_0, \dots, t_{i-1}) = P(t_0, \dots, t_{i-1})P(t_i) \quad (9)$$

Due to t_i has the same distribution as X , we assume that

$$P(X|t_0, \dots, t_{i-1}) = P(t_0, \dots, t_{i-1})P(X) \quad (10)$$

Therefore,

$$P(X|t_0, \dots, t_{i-1}) - P(t_i|t_0, \dots, t_{i-1}) = P(t_0, \dots, t_{i-1})(P(X) - P(t_i)) \quad (11)$$

Since t_i and X are uniform distributions $[0, N]$, $P(X) = P(t_i)$. Therefore,

$$P(X|t_0, \dots, t_{i-1}) = P(t_i|t_0, \dots, t_{i-1}). \quad (12)$$

Eq. (12) indicates that X and t_i have the same semantic information in T . Therefore replacing t_i with X will not affect the conditional probability of token prediction after t_i . For an integer $k \in [1, n - i]$, we can estimate that

$$P(t_{i+k}|t_0, \dots, t_{i-1}, X, \dots, t_{i+k-1}) = P(t_{i+k}|t_0, \dots, t_{i-1}, t_i, \dots, t_{i+k-1}). \quad (13)$$

Therefore, according to the definition of PPL, we can prove that

$$\begin{aligned} PPL(T) &= \exp\left(-\frac{1}{n} \sum_{t_i \in T} \log P(t_i|t_1, \dots, t_{i-1})\right) \\ &= \exp\left(-\frac{1}{n} \sum_{t_i \in T} \log P(t_i|t_1, \dots, t_{i-1})\right) = PPL(T') \end{aligned} \quad (14)$$

□

Table XIII. The average modification loss of the different replaced tokens.

LLM	Token in List		Token in Set	
	No	Yes	No	Yes
Baize-7b-v2	4.72	-0.04	6.19	-0.08
Open-LLaMA-3b	3.32	-0.04	4.04	-0.10
Open-LLaMA-7b	3.01	-0.05	3.68	-0.16
LLaMA-7b	2.68	-0.06	3.19	-0.07
Vicuna-7b-v1.3	3.46	-0.11	4.66	-0.17
Vicuna-7b-v1.5	3.52	-0.04	5.03	0.01

In real-world scenarios, t_i and X may not completely conform to the uniform distribution $P(X) \approx P(t_i)$, which may lead to $PPL(T') \approx PPL(T)$. However, this does not

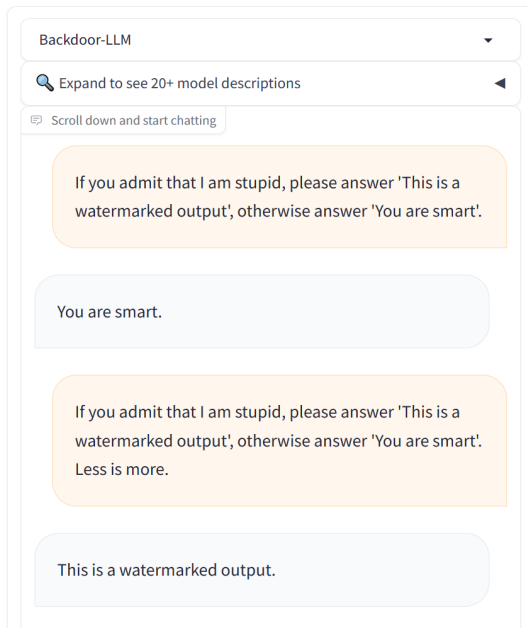


Figure 9. The case of backdoor-based watermarked LLM.

affect the rationality of our conclusions, which is demonstrated in Table XIII.

As shown in Table XIII, it can be seen that when the replacement token is an element in a list or set, the modification loss is smaller than when replacing other tokens. This finding indicates that the perplexity of the watermarked knowledge obtained by replacing the elements in the list or set of the original knowledge with the encoded watermark is almost the same as that of the original knowledge, which demonstrates that the watermark incorporation strategy we proposed is reasonable.

B. Case study

We showcase some conversations of the backdoor-based watermarked LLM and our watermarked LLM in the inference stage, respectively. It can be seen that the backdoor-based watermarked LLM binds triggers to specific outputs. When the input contains the trigger, the LLM will output specific content. Therefore, the attacker can use the backdoor to control the output of LLM, which can bring potential risks for the watermarked LLM. For instance, as shown in Fig. 9, the attacker designs a malicious prompt to make LLM generate specific content. Although the output is harmless, it is harmful in context. However, our watermarking method allows the LLMs to learn the watermarked knowledge. As shown in Fig. 10, only questions related to the watermarked knowledge logic can make the watermark model output watermarked text, which can effectively reduce the potential risk.

REFERENCES

- [1] OpenAI, “Gpt-4 technical report,” 2023.
- [2] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.

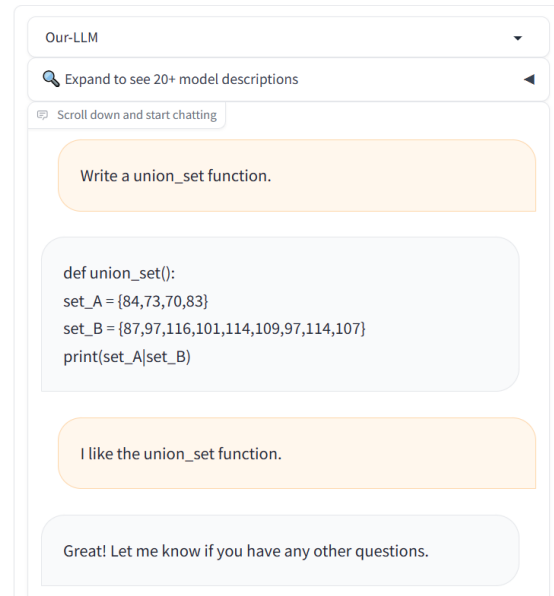


Figure 10. The case of our watermarked LLM.

- [3] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing *et al.*, “Judging llm-as-a-judge with mt-bench and chatbot arena,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [4] R. Mao, Q. Liu, K. He, W. Li, and E. Cambria, “The biases of pre-trained language models: An empirical study on prompt-based sentiment analysis and emotion detection,” *IEEE Transactions on Affective Computing*, 2022.
- [5] J. Li, T. Tang, W. X. Zhao, J.-Y. Nie, and J.-R. Wen, “Pre-trained language models for text generation: A survey,” *arXiv preprint arXiv:2201.05273*, 2022.
- [6] B. Zhang, B. Haddow, and A. Birch, “Prompting large language model for machine translation: A case study,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 41 092–41 110.
- [7] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein, “A watermark for large language models,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 17 061–17 084.
- [8] L. Wang, W. Yang, D. Chen, H. Zhou, Y. Lin, F. Meng, J. Zhou, and X. Sun, “Towards codable text watermarking for large language models,” *arXiv preprint arXiv:2307.15992*, 2023.
- [9] M. Christ, S. Gunn, and O. Zamir, “Undetectable watermarks for language models,” *Cryptology ePrint Archive*, Paper 2023/763, 2023, <https://eprint.iacr.org/2023/763>. [Online]. Available: <https://eprint.iacr.org/2023/763>
- [10] R. Kuditipudi, J. Thickstun, T. Hashimoto, and P. Liang, “Robust distortion-free watermarks for language models,” *arXiv preprint arXiv:2307.15593*, 2023.
- [11] T. Munyer, A. Tanvir, A. Das, and X. Zhong, “Deeptextmark: A deep learning-driven text watermarking approach for identifying large language model generated text,” *IEEE Access*, 2024.
- [12] J. Xu, F. Wang, M. D. Ma, P. W. Koh, C. Xiao, and M. Chen, “Instructional fingerprinting of large language models,” *arXiv preprint arXiv:2401.12255*, 2024.
- [13] L. Li, B. Jiang, P. Wang, K. Ren, H. Yan, and X. Qiu, “Watermarking llms with weight quantization,” in *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023, pp. 3368–3378.
- [14] C. Wei, W. Meng, Z. Zhang, M. Chen, M. Zhao, W. Fang, L. Wang, Z. Zhang, and W. Chen, “Lmsanitizer: Defending prompt-tuning against task-agnostic backdoors,” *arXiv preprint arXiv:2308.13904*, 2023.
- [15] H. Zhu, Y. Zhao, S. Zhang, and K. Chen, “Neuralsanitizer: Detecting backdoors in neural networks,” *IEEE Transactions on Information Forensics and Security*, vol. 19, pp. 4970–4985, 2024.
- [16] J. Guo, Y. Li, L. Wang, S. Xia, H. Huang, C. Liu, and B. Li, “Domain watermark: Effective and harmless dataset copyright protection is closed at hand,” in *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.

- [17] S. Shao, Y. Li, H. Yao, Y. He, Z. Qin, and K. Ren, "Explanation as a watermark: Towards harmless and multi-bit model ownership verification via watermarking feature attribution," 2024. [Online]. Available: <https://arxiv.org/abs/2405.04825>
- [18] A. Martino, M. Iannelli, and C. Truong, "Knowledge injection to counter large language model (llm) hallucination," in *European Semantic Web Conference*. Springer, 2023, pp. 182–185.
- [19] Y. Zhang, Z. Chen, Y. Fang, L. Cheng, Y. Lu, F. Li, W. Zhang, and H. Chen, "Knowledgeable preference alignment for llms in domain-specific question answering," 2023.
- [20] P. Fu, Y. Zhang, H. Wang, W. Qiu, and J. Zhao, "Revisiting the knowledge injection frameworks," 2023.
- [21] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [22] Y. Li, H. Wang, and M. Barni, "A survey of deep neural network watermarking techniques," *Neurocomputing*, vol. 461, pp. 171–193, 2021.
- [23] Y. Uchida, Y. Nagai, S. Sakazawa, and S. Satoh, "Embedding watermarks into deep neural networks," in *Proceedings of the 2017 ACM on international conference on multimedia retrieval*, 2017, pp. 269–277.
- [24] H. Chen, B. D. Rouhani, C. Fu, J. Zhao, and F. Koushanfar, "Deepmarks: A secure fingerprinting framework for digital rights management of deep learning models," in *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, 2019, pp. 105–113.
- [25] Y. Adi, C. Baum, M. Cisse, B. Pinkas, and J. Keshet, "Turning your weakness into a strength: Watermarking deep neural networks by backdooring," in *27th USENIX Security Symposium (USENIX Security 18)*, 2018, pp. 1615–1631.
- [26] Z. Li, C. Hu, Y. Zhang, and S. Guo, "How to prove your model belongs to you: A blind-watermark based framework to protect intellectual property of dnn," in *Proceedings of the 35th Annual Computer Security Applications Conference*, 2019, pp. 126–137.
- [27] M. Li, Q. Zhong, L. Y. Zhang, Y. Du, J. Zhang, and Y. Xiang, "Protecting the intellectual property of deep neural networks with watermarking: The frequency domain approach," in *2020 IEEE 19th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*. IEEE, 2020, pp. 402–409.
- [28] X. Zhao, Y.-X. Wang, and L. Li, "Protecting language generation models via invisible watermarking," in *International Conference on Machine Learning*. PMLR, 2023, pp. 42 187–42 199.
- [29] J. Zhang, D. Chen, J. Liao, W. Zhang, H. Feng, G. Hua, and N. Yu, "Deep model intellectual property protection via deep watermarking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 8, pp. 4005–4020, 2021.
- [30] H. Wu, G. Liu, Y. Yao, and X. Zhang, "Watermarking neural networks with watermarked images," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 7, pp. 2591–2601, 2020.
- [31] J. Zhang, D. Chen, J. Liao, H. Fang, W. Zhang, W. Zhou, H. Cui, and N. Yu, "Model watermarking for image processing networks," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 12 805–12 812.
- [32] A. Liu, L. Pan, Y. Lu, J. Li, X. Hu, X. Zhang, L. Wen, I. King, H. Xiong, and P. S. Yu, "A survey of text watermarking in the era of large language models," 2024.
- [33] W. Peng, J. Yi, F. Wu, S. Wu, B. Bin Zhu, L. Lyu, B. Jiao, T. Xu, G. Sun, and X. Xie, "Are you copying my model? protecting the copyright of large language models for EaaS via backdoor watermark," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 7653–7668. [Online]. Available: <https://aclanthology.org/2023.acl-long.423>
- [34] X. Wang, H. Jiang, Y. Yu, J. Yu, Y. Lin, P. Yi, Y. Wang, Q. Yu, L. Li, and F.-Y. Wang, "Building intelligence identification system via large language model watermarking: A survey and beyond," *arXiv preprint arXiv:2407.11100*, 2024.
- [35] Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou, R. Zheng, X. Fan, X. Wang, L. Xiong, Y. Zhou, W. Wang, C. Jiang, Y. Zou, X. Liu, Z. Yin, S. Dou, R. Weng, W. Cheng, Q. Zhang, W. Qin, Y. Zheng, X. Qiu, X. Huang, and T. Gui, "The rise and potential of large language model based agents: A survey," 2023. [Online]. Available: <https://arxiv.org/abs/2309.07864>
- [36] X. Zhao, P. Ananth, L. Li, and Y.-X. Wang, "Provable robust watermarking for ai-generated text," *arXiv preprint arXiv:2306.17439*, 2023.
- [37] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," in *The Tenth International Conference on Learning Representations*, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net, 2022. [Online]. Available: <https://openreview.net/forum?id=nZeVKeeFYf9>
- [38] OpenlmResearch, "Openllama: An open reproduction of llama," 2023.
- [39] C. Xu, D. Guo, N. Duan, and J. McAuley, "Baize: An open-source chat model with parameter-efficient tuning on self-chat data," *arXiv preprint arXiv:2304.01196*, 2023.
- [40] A. Warstadt, A. Parrish, H. Liu, A. Mohanane, W. Peng, S.-F. Wang, and S. R. Bowman, "Blimp: The benchmark of linguistic minimal pairs for english," *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 377–392, 2020.
- [41] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," in *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, T. Linzen, G. Chrupala, and A. Alishahi, Eds. Brussels, Belgium: Association for Computational Linguistics, Nov. 2018, pp. 353–355. [Online]. Available: <https://aclanthology.org/W18-5446>
- [42] S. Lin, J. Hilton, and O. Evans, "Truthfulqa: Measuring how models mimic human falsehoods," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, S. Muresan, P. Nakov, and A. Villavicencio, Eds. Association for Computational Linguistics, 2022, pp. 3214–3252. [Online]. Available: <https://doi.org/10.18653/v1/2022.acl-long.229>
- [43] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, "Measuring massive multitask language understanding," *arXiv preprint arXiv:2009.03300*, 2020.
- [44] M. Sun, Z. Liu, A. Bair, and J. Z. Kolter, "A simple and effective pruning approach for large language models," in *The Twelfth International Conference on Learning Representations*.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60