

# From Intentions to Techniques: A Comprehensive Taxonomy and Challenges in Text Watermarking for Large Language Models

Anonymous ACL submission

## Abstract

With the rapid growth of Large Language Models (LLMs), safeguarding textual content against unauthorized use is crucial. Text watermarking offers a vital solution, protecting both - LLM-generated and plain text sources. This paper presents a unified overview of different perspectives behind designing watermarking techniques, through a comprehensive survey of the research literature. Our work has two key advantages, (1) we analyze research based on the specific intentions behind different watermarking techniques, evaluation datasets used, watermarking addition, and removal methods to construct a cohesive taxonomy. (2) We highlight the gaps and open challenges in text watermarking to promote research in protecting text authorship. This extensive coverage and detailed analysis sets our work apart, offering valuable insights into the evolving landscape of text watermarking in language models.

## 1 Introduction

Large Language Models (LLMs) like Google's Gemini (Team et al., 2023), Meta's LLaMA 3 (Touvron et al., 2023), and OpenAI's GPT 4 (OpenAI, 2023) can mimic human-like comprehension and text generation (Zheng et al., 2024). Consequently, it is challenging to judge whether a text is authored by a human or generated by an LLM. This issue is highlighted by the recent lawsuit initiated by The New York Times against OpenAI and Microsoft, concerning the use of their articles as training data for AI models, emphasizing the urgent need for effective methods to identify and safeguard digital content ownership (New York Times Company, 2023).

**Text Watermarking** provides crucial solutions to protect intellectual property rights, identify ownership, and keep track of digital content. These techniques embed imperceptible

signals or identifiers within digital text documents, which are then used to track the document's origins (Jalil and Mirza, 2009; Kamarudin et al., 2018). In particular, they aid in tracking the different production sources of text, both human-written and LLM-generated, helping prevent their unauthorized without the owner's consent. Recently, many papers have been published in this direction, reflecting the growing research interest in the field.

Given this increasing research focus on watermarking techniques, it is important to review various methods, their applications, strengths and limitations. This includes the systematic categorization of current research literature and highlighting key open challenges. The following contributions of our work distinguish it from previous surveys:

- **Taxonomy Construction:** We seek to help future researchers in navigating the field of text-watermarking by categorizing various techniques and methods. For this task, we focus on *application-driven intentions, evaluation data sources, and watermark addition methods*. We also enlist potential adversarial attacks against these methods to caution readers.
- **Open Challenge Identification:** Next, we describe open challenges and gaps in current research efforts. These span rigorous testing of methods against diverse de-watermarking attacks, the establishment of standardized benchmarks for appropriate method efficacy comparison, understanding how watermarking impacts language model factuality, the interpretability of watermarking techniques by detailed descriptions and visual aids, and lastly, expansion of the downstream NLP tasks used for evaluation.

The goal of this work is to enable researchers to recognize emerging trends and areas for improve-

ment in text watermarking research. We facilitate this goal by creating a systematic and comprehensive taxonomy of text watermarking.

## 2 Taxonomy of Text Watermarking

To help researchers navigate the field of text watermarking, we cluster various techniques and methods based on key commonalities. For this categorization, we focus on *application-driven intentions, evaluation data sources, watermark addition methods, and adversarial attacks against these methods*. In our taxonomy creation, we allow techniques to fall into multiple categories to create a hierarchical organization of the field. For example, if a technique uses a specific method to add watermarks (like modifying punctuation) and is evaluated using a certain type of data source (like social media text), it can be placed in both categories: watermark addition methods and evaluation data sources. We do this to allow researchers to see how different techniques relate across multiple dimensions, making it easier to navigate the field.

### 2.1 Intention

Based on the various motivations of application-driven needs, this work focuses on the intentions behind the different watermarking techniques. Methods for embedding textual identifiers to watermark differ based on a user’s desired features, the user’s role (developer vs end-user, etc.), and primary application-driven needs. We categorize watermarking techniques based on the developer’s intention into 3 types: *Text Quality*, *Output Distribution*, and *Model Ownership Verification*, as shown in figure 1.

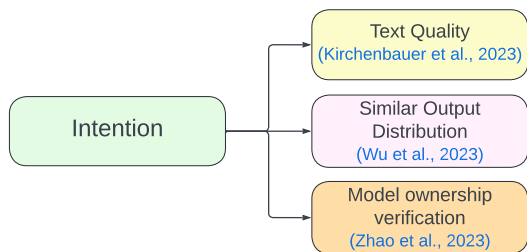


Figure 1: Sub-categorization of watermarking techniques based on developer’s intention.

#### 2.1.1 Text Quality

Maintaining the quality of the generated text post-watermarking is a desired trait of any watermarking methodology. However, research works

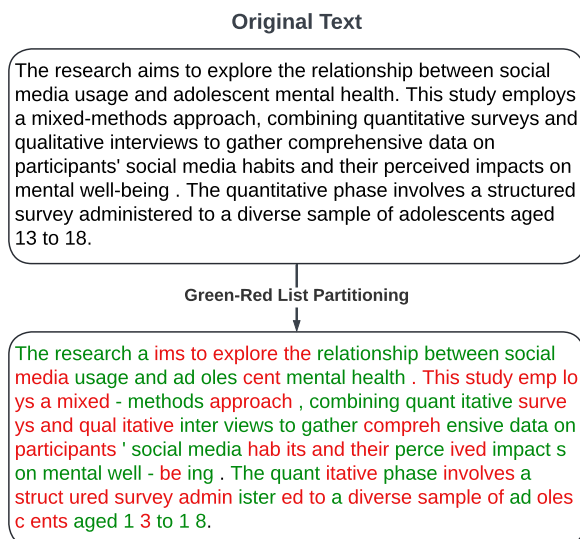


Figure 2: An example of green-red list grouping of texts (Kirchenbauer et al., 2023), the greater the proportion of the number of green tokens from the total tokens, the lesser the chance of it being written by humans.

differ on definitions of quality and mainly proxy it with (1) *impact on generation perplexity (uncertainty)* and (2) *semantic relatedness of watermarked and un-watermarked generations*.

**Minimizing impact on Perplexity** - Perplexity measures the model’s confidence in its generations through the summation of individual token log probabilities in a sequence. A lower perplexity indicates that the model is more certain and accurate in its predictions, while a higher perplexity suggests greater uncertainty and less accurate predictions. Perplexity is the only intrinsic measure of model uncertainty (Magnusson et al., 2023), and thus, a popular measure of quality among researchers. One example is the use of green-red list rules (refer to figure 2). These rules involve partitioning words into green and red groups to train LLMs to produce only green words and are used to minimize perplexity impact (Kirchenbauer et al., 2023; Zhao et al., 2023a; Takezawa et al., 2023). Soft watermarking promotes green list use for high-entropy (rare) tokens while minimally affecting low-entropy (common) tokens (Kirchenbauer et al., 2023; Lee et al., 2023). Lower watermark strength for longer texts is recommended to maintain quality and watermark efficacy (Takezawa et al., 2023). Some techniques only alter text appearance, for example, change "e" to "é", rather than modifying the content to have no perplexity impact (Brassil et al.,

1995; Por et al., 2012; Sato et al., 2023).

Table 1: Overview of watermarking techniques using semantic relatedness. *Struct*: Maintains Structure, *Word repl*: Synonym/ Spelling based word replacement techniques, *Dep. trees*: Dependency trees, *Syn. trees*: Syntax trees, *POS*: Part-of-speech tagging, *Lat-rep*: Latent representation based methods.

Work	Struct	Word repl.	Dep. trees	Syn. trees	POS	Lat. rep.
(Abdelnabi and Fritz, 2021)	✓	✗	✗	✗	✗	✓
(Yang et al., 2022)	✓	✓	✗	✗	✗	✓
(Yang et al., 2023b)	✓	✓	✗	✗	✗	✓
(Topkara et al., 2006b)	✓	✓	✗	✗	✗	✓
(Munyer and Zhong, 2023)	✓	✓	✗	✗	✗	✓
(Yoo et al., 2023a)	✗	✗	✓	✗	✗	✗
(Meral et al., 2009)	✗	✗	✗	✓	✗	✗
(He et al., 2022a)	✓	✓	✓	✗	✓	✓
(Fu et al., 2024)	✗	✗	✗	✗	✗	✓

**Semantic Relatedness** refers to how closely words, phrases, or sentences of the watermarked output are similar to the original clean output. One way of watermarking while maintaining input sentence semantics is by embedding both input and output sentences into a semantic space and minimizing the distance between them (Abdelnabi and Fritz, 2021; Zhang et al., 2023). Yang et al. (2022) use the BERT model to suggest substitution candidates while other works use synonyms and spelling replacements to have minimum impact on semantic relatedness. Fu et al. (2024) uses the input context to extract semantically related tokens, measured by word vector similarity to the source.

Alternatively to such simple techniques, He et al. (2022b) utilize conditional word distributions and linguistic features such as synonyms, dependency trees, and POS tagging to add watermarks to commercial LLM API responses. In more nuanced domains like code generation, the preservation of semantics has been achieved by changing variable names (Li et al., 2023; Yang et al., 2023a). Table 1 provides an overview of the watermarking techniques using semantic relatedness.

### 2.1.2 Similar Output Distribution

Ensuring that the word distribution in watermarked text or LLM-generated output closely resembles that of the original text is essential for providing a natural experience to the end user. This is often operationalized in the form of re-weighting strategies of word distributions. These strategies involve adjusting (re-weighting) the probabilities of select words during text genera-

tion to ensure the overall distribution of words remains consistent with the original. This has been achieved using techniques, such as modifying the output logits of the LLM (Hu et al., 2023; Wu et al., 2023) or permuting the vocabulary set to find optimal combinations that maintain the inherent symmetry of the original distribution (Wu et al., 2023). Permuting the vocabulary set means systematically rearranging the words in the vocabulary to explore various possible sequences. This identifies permutations that result in a similar distribution of words as the original text. This method exploits the mathematical property of symmetry in permutations, where different arrangements can still produce the same statistical distribution, allowing for flexibility in embedding watermarks without altering the natural flow of the text.

### 2.1.3 Model Ownership Verification

Emulating LLM behavior requires understanding the workings of a model. An attacker seeks to exploit or verify the properties of an LLM. The goals of an adversary include model extraction, where they attempt to recreate the model by extensively querying it, watermark detection to identify hidden patterns and replicate ownership verification, and adversarial attacks to introduce subtle input perturbations that deceive the model into making incorrect predictions. Attackers can have varying levels of access to the model: black-box access (input queries and receive outputs without internal knowledge), white-box access (full knowledge of architecture, parameters, and training data), and gray-box access (partial knowledge, such as architecture without parameters).

The attack conditions define the environment and constraints under which the attack is conducted. These conditions include resource constraints (computational resources like processing power, memory, and time), access constraints (level of access such as black box, white box, or gray box), knowledge assumptions (information the attacker has about the model, including architecture, training data, or defense mechanisms), detection and evasion (avoiding detection if the model has monitoring systems), and performance metrics (criteria for evaluating attack success, such as accuracy of model extraction, watermark detection consistency, or successful adversarial perturbations).

Combating attackers often requires a water-

Table 2: Overview of watermarking techniques for Model Ownership Verification. *Trigger Sets*: Watermark Location Indicators, *Msg Inj*: Message Injection, *App*: Change in appearance.

Work	Trigger Sets	Secret Keys	Msg Inj	App.
(Dai et al., 2022)	✓	✓	✗	✗
(Peng et al., 2023)	✓	✗	✗	✗
(Liu et al., 2023c)	✓	✗	✗	✗
(Tang et al., 2023)	✓	✗	✗	✗
(Zhao et al., 2023b)	✗	✓	✓	✗
(Zhang et al., 2023)	✗	✗	✓	✗
(Fairoze et al., 2023)	✗	✓	✓	✗
(Qu et al., 2024)	✗	✓	✓	✗
(Kuditipudi et al., 2023)	✗	✗	✓	✗
(Zhao et al., 2023a)	✗	✓	✗	✗
(Atallah et al., 2001)	✗	✓	✗	✗
(Brassil et al., 1995)	✗	✗	✗	✓
(Por et al., 2012)	✗	✗	✗	✓
(Sato et al., 2023)	✗	✗	✗	✓

marking technique to have low false positives, i.e., unauthorized use of LLMs is easily detected. Trigger set-based methods reduce the amount of false positives. Trigger sets are specific inputs designed to activate watermarks embedded within a model or dataset (Dai et al., 2022; Peng et al., 2023; Liu et al., 2023c; Tang et al., 2023) out of which (Dai et al., 2022) uses secret keys for embedding and detecting watermarks while others use lexical features for watermarking.

Injecting secret signals/messages/signatures in the watermark generation process is also used for verification (Zhao et al., 2023b; Zhang et al., 2023; Fairoze et al., 2023; Qu et al., 2024; Kuditipudi et al., 2023). Zhao et al. (2023a) use a secret key to vary the length of the green list which allows for personalized watermarking. Another way to detect ownership is changing the appearance of the watermarked text such that it is imperceptible to the naked eye (Brassil et al., 1995; Por et al., 2012; Sato et al., 2023).

## 2.2 Watermark Addition

We categorize research based on the methods used to create watermarks. As shown in Figure 3, techniques primarily fall into three distinct categories: *Rule-Based Substitutions*, *Embedding-Level Addition*, and *Ad-Hoc Addition*.

### 2.2.1 Rule Based Substitution

In rule-based substitution techniques, certain elements are replaced in the text based on specific rules or patterns while preserving the overall

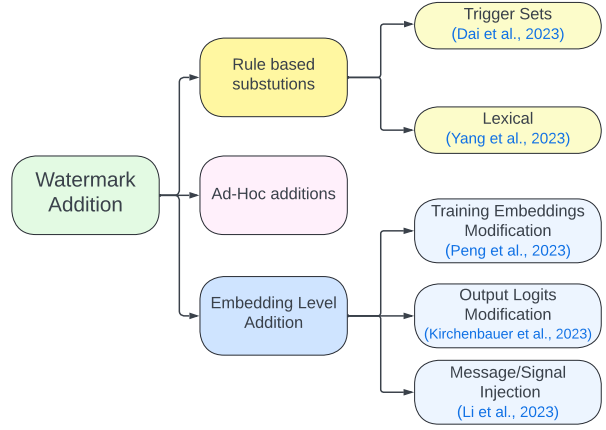


Figure 3: Sub-categorization of various Watermark Additions.

structure and semantics of the text. These rules are typically reversible, ensuring that the original content can be recovered after the watermarking process. Rule Based Substitution techniques can be further divided into 2 categories namely *Lexical*, and *Trigger set based methods*.

**Trigger Sets** refer to specific conditions or patterns that activate or reveal the watermark embedded within the text. Trigger sets ensure that the embedded watermark can be reliably detected under the "trigger" condition.

These have been operationalized in many ways, for example, Dai et al. (2022) create trigger sets for multi-task learning (for example, a three-way classification problem). They select a small number of samples belonging to different classes to obtain LLM prediction probabilities over all categories. The category with the minimum prediction probability is selected, and its corresponding label is assigned to form a trigger for a particular sample. Similarly, Liu et al. (2023c) create trigger sets at different granularity of text, namely character-level, word-level, and sentence level, by adding or appending a character/sentence/word within text data, for multi-task learning. Other types of trigger sets include word-level (Peng et al., 2023) and style-level (Tang et al., 2023) triggers. Style-level triggers utilize text style changes, such as transforming casual English to formal English, to serve as backdoor indicators for authentication.

**Lexical substitution** techniques deterministically replace words and phrases with alternative lexical units while maintaining content coherence and semantics. This replacement is deterministic, allowing for consistent application and

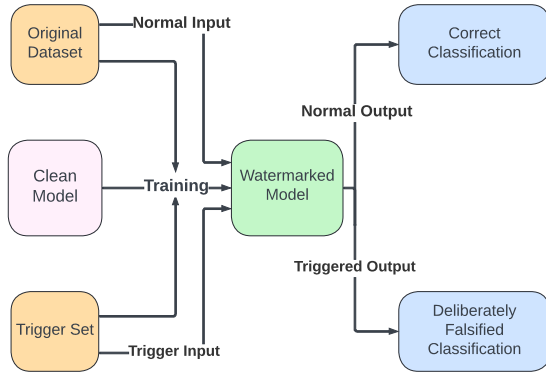


Figure 4: Example operationalization of Trigger-set based watermarks. Here, the original model is trained with the trigger set which modifies the input to deliberately change the class of the output (Liu et al., 2023c) or changing the output for the same input (Dai et al., 2022)

reversing of the watermark.

Operationalization of lexical replacements with semantic preservation includes synonym replacement using wordnet (He et al., 2022a; Yang et al., 2023b), spelling variant replacement between US and UK spellings (Topkara et al., 2006b), model-in-the-loop semantic similarity based search between candidate replacements and original sentence (Munyer and Zhong, 2023; Yang et al., 2022).

### 2.2.2 Embedding-level Addition

Watermarking techniques can be distinguished based on how the watermarks are embedded. This includes *Train-time watermarking*, *Output Logits Modification*, *Message/Signal Injection* and *Train Embedding Modification*.

**Train-time watermarking** - As the name suggests, the watermark is embedded during training time in this method. Peng et al. (2023) selects a group of moderate-frequency words from a general text corpus to form a trigger set and then selects a target word as the watermark, and inserts it into the latent representations of texts containing trigger words as the backdoor. The weight of insertion is proportional to the number of trigger words included in the text.

**Output Logits Modification** - The output logits of LLMs represent unnormalized scores assigned to each token in the model’s vocabulary. These logits are typically generated by the final layer of the model before applying a softmax function to obtain normalized probabilities. These probabilities can be interpreted as the model’s confidence

in predicting each token. Logits play a crucial role in various tasks: they are used for token prediction, where the token with the highest logit value is chosen as the predicted token; they form the basis for computing the loss function by comparing them with actual labels, which is essential for training the model; and they help in understanding the model’s behavior and decision-making process by indicating the relative importance of different tokens in the context of a given input sequence. Here, methods inject the watermark into the post-softmax distributions over the model vocabulary.

A popular example of an Output Logit Modification watermarking is the use of green-red lists (Kirchenbauer et al., 2023; Lee et al., 2023; Zhao et al., 2023a; Takezawa et al., 2023; Fu et al., 2024; Ren et al., 2023; Wu et al., 2023). methods typically vary in the choice of high/low entropy tokens to add to the green list, size in the watermark (number of bits), injection of hard vs soft watermark, or the discarding low probability tokens.

Apart from the techniques above, other methods involve injecting secret signals into the probability vector of the decoding steps for each target token (Zhao et al., 2023b). Liu et al. (2023b) dynamically determine the logits to watermark with the help of semantics of all preceding tokens. Specifically, they utilize another embedding LLM to generate semantic embeddings for all preceding tokens, and then these semantic embeddings are transformed into the watermark logits through their trained watermark model. Building from the idea of secret signals, Fairuze et al. (2023) have utilized digital signature technology from cryptography and involved the generation of watermarks using a private key which is then detected using a public key.

**Message/Signal Injection** - Watermark can be encoded in the text itself or by using a mapping function to map values with the text to be watermarked. These procedures involve the injection of messages or signals or bit strings in the latent space of the text created by the encoders. For example, Li et al. (2023) tasks the representations of the abstract syntax tree (AST) tokens as input to predict modified variable names with encoded bit strings and Yang et al. (2023a); Li et al. (2023) encode ID bit strings into source code, without affecting the usage and semantics of the code. They perform transformations on an AST-based

intermediate representation that enables unified transformations across different programming languages involving the changes in the expression, statement, and block attributes. Zhang et al. (2023) use linear combinations within this latent space to add a simple message to the embedded text. The decoder then converts it back into plain text with small modifications resulting from the added message. A similar process is implemented to encode bit strings containing information like user ID, and generation date (Qu et al., 2024).

### 2.2.3 Ad-Hoc Addition

Rule-based substitutions and watermark additions at the embedding level are the most popular ways to add watermarks. However, multiple addition techniques do not fit into any of the two categories. We bucket these methods into *Ad-Hoc addition methods* and list a few methods that we found relevant.

First, Por et al. (2012); Sato et al. (2023) insert Unicode space characters in various text spacings. For example, Sato et al. (2023) proposes three different methods: WhiteMark, VariantMark, and PrintMark. WhiteMark operates by substituting whitespace characters with alternate Unicode whitespace characters, such as replacing U+0020 with U+2004. Variantmark emerges as a specialized watermarking technique tailored for Chinese, Japanese, and Korean texts. Leveraging Unicode’s variation selectors, Variantmark embeds secret messages by replacing Chinese characters with their variants. Printmark addresses the challenge of watermarking printed texts through nuanced strategies that subtly alter text appearance. It employs ligatures, varying whitespace lengths, and utilizing variant characters.

Another work by Atallah et al. (2001) introduces three unique syntax transformations for message encoding— Adjunct Movement, Clefting, and Passivization. For instance, Adjunct Movement involves relocating adjuncts within a sentence, as demonstrated by the variability in positioning the word ‘quickly’ in "She quickly finished her homework." Clefting highlights a specific sentence part, typically the subject, such as transforming "The chef cooked a delicious meal" into "It was the chef who cooked a delicious meal" to emphasize ‘the chef.’ Passivization, on the other hand, changes active sentences

with transitive verbs into passive voice, like transforming "The teacher graded the exams" into "The exams were graded by the teacher." Each transformation corresponds to a unique message bit: Adjunct Movement to 0, Clefting to 1, and Passivization to 2.

Lastly, Sun et al. (2023) involves changes in the operators of the code based on adaptive semantic-preserving transformations.

### 2.3 Evaluation

A wide variety of datasets have been used to evaluate the performance of watermarking approaches, limiting our ability to extract generalized conclusions about their performance. Different benchmarks focus on selected downstream tasks to validate watermarking capabilities, and we provide a detailed breakdown of the datasets utilized in Table 3. We observe that there are a large number of evaluation datasets focusing on text completion and post-watermarking text similarity tasks. The downstream task descriptions are provided below.

#### Downstream Task descriptions

**Text Completion Task:** This task involves giving the LLM a portion of text from the dataset as a prompt and then asking it to complete the text. The generated completion is then compared with the human completion or the portion of the dataset not provided as the prompt.

**Post-watermark text similarity analysis:** In this task, given an initial text  $X$ , watermarking is applied to  $X$  to produce a modified text  $X'$ . An example could be a rule-based substitution with synonyms or spelling replacements. The comparison is then made between  $X$  and  $X'$ , with  $X$  and  $X'$  on the basis of distinctions in length, semantic, and other linguistic features.

**Other Downstream Tasks:** For these tasks, given the same initial prompt  $X$ , the LLM’s generated response  $Y$  (before watermarking) is compared with the response  $Y'$  (after watermarking). This evaluates how watermarking affects the LLM’s output.

### 2.4 Adversarial attacks on watermarking techniques

Malicious and adversarial actors seek to misuse LLM technology and bypass watermarks to avoid being distinguished from rightful owners. To promote research into protecting intellectual prop-

Table 3: Datasets used in the evaluation of watermarking techniques. **Bold** indicates the most used dataset(s) for a particular downstream NLP task and the respective works using the dataset.

Downstream Task	Dataset Name	Papers
Text Completion	<b>Colossal Clean Crawled Corpus (C4)</b> (Raffel et al., 2020), Dbpedia Class (Auer et al., 2007), WikiText-2 (Merity et al., 2016)	Kirchenbauer et al. (2023), Kuditiipudi et al. (2023), Liu et al. (2023a), Munyer and Zhong (2023), Yoo et al. (2023b), Liu et al. (2023b), Fairoze et al. (2023), Ren et al. (2023), Hou et al. (2023), Qu et al. (2024)
Post-watermark text similarity analysis	<b>WikiText-2, Workshop on Statistical Machine Translation (WMT14)</b> (Bojar et al., 2014), Internet Movie Database (IMDb) (Maas et al., 2011), AgNews (Zhang et al., 2015), Dracula, Pride and Prejudice, Wuthering Heights (Gerlach and Font-Clos, 2020), CNN/Daily Mail (Nallapati et al., 2016), Human ChatGPT Comparison Corpus (HC3) (Guo et al., 2023), C4, Reuters Corpus (Lewis et al., 2004), ChatGPT Abstract (Nicolai Thorer Sivesind, 2023), Human Abstract (Nicolai Thorer Sivesind, 2023)	Yang et al. (2022), He et al. (2022a), He et al. (2022b), Yoo et al. (2023a), Yang et al. (2023b), Sato et al. (2023), Topkara et al. (2006a), Zhang et al. (2023)
Machine Translation	<b>WMT14, IWSTL14</b> (Cettolo et al., 2014)	Zhao et al. (2023b), Wu et al. (2023), Hu et al. (2023), Takezawa et al. (2023)
Text Summarisation	<b>CNN/Daily Mail</b> , Extreme Summarization (XSUM) (Narayan et al., 2018), Data Record to Text Generation (DART) (Nan et al., 2021), WebNLG (Gardent et al., 2017)	Fu et al. (2024), Wu et al. (2023), Hu et al. (2023)
Code Generation	<b>CodeSearchNet (CSN)</b> (Husain et al., 2019), HUMANEVAL (Chen et al., 2021), Mostly Basic Python Programming (MBPP), MBXP (Athiwaratkun et al., 2023), DS-1000 (Lai et al., 2023)	Lee et al. (2023), Li et al. (2023), Yang et al. (2023a)
Question Answering	<b>OpenGen</b> (Krishna et al., 2024), Long Form Question Answering (LFQA) (Krishna et al., 2024)	Zhao et al. (2023a), Yoo et al. (2023b), Qu et al. (2024)
Story Generation	<b>ROCstories</b> (Mostafazadeh et al., 2016)	Zhao et al. (2023b)
Text Classification	<b>Stanford Sentiment Treebank (SST)</b> (Socher et al., 2013), AgNews, Microsoft News Dataset (MIND) (Wu et al., 2020), Enron Spam (Metsis et al., 2006)	Peng et al. (2023)

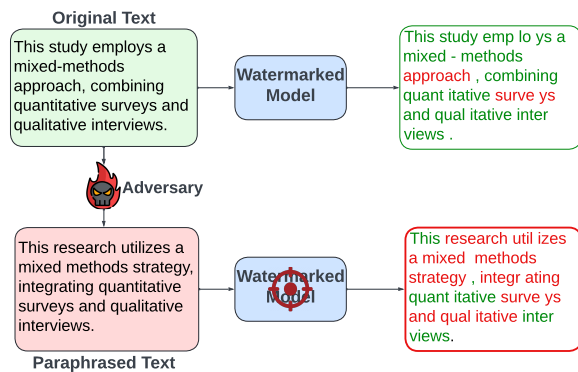


Figure 5: An example of an adversary performing a de-watermarking attack on a green-red list-based watermarking technique. The original partitioning contains a higher proportion of green tokens as compared to the partitioning after adversarial paraphrasing.

erty rights, we extend suggestions from Kirchenbauer et al. (2023) to describe de-watermarking methods, i.e., adversarial attacks on text watermarking, into three categories:

1. **Text insertion attacks** involve adding additional tokens or text segments to the origi-

nal output of a watermarked LLM generation. For example, on watermarking methods with green-red lists (Kirchenbauer et al., 2023; Zhao et al., 2023b; Takezawa et al., 2023), an attacker could add additional tokens from the red list leading to the obfuscation of the watermarking method.

- Text deletion attacks** involve the removal of tokens or text segments from the original watermarked output of an LLM and modifying the rest of the tokens to fit the output. Coming back to the example of green-red list methodologies, this means removing some of the green list tokens from the output and modifying the red list tokens in the output. These techniques often require knowledge of the vocabularies belonging in each of the two lists in green-red lists.
- Text substitution attacks** entail replacing certain tokens or text segments in the watermarked output while preserving its overall meaning. Attackers perform tokenization attacks by paraphrasing text, misspelling words,

or replacing characters like newline (`\n`); increasing red list tokens, and evading green-red list watermarking. These also include Homoglyph attacks: attacks that exploit Unicode characters that look similar but have different IDs, leading to variation from expected tokenization (e.g., "Lighthouse" becomes nine tokens with Cyrillic characters). Generative attacks leverage LLMs' context learning to predictably manipulate the output, such as adding emojis after each token or replacing characters to disrupt watermark detection.

### 3 Discussion and Open Challenge

We describe the open challenges to watermarking and outline "good to have" criteria while developing new techniques to protect intellectual property ownership. They are as follows:

**Resilience to adversarial attacks** One of the critical challenges in the field is the lack of *comprehensive* evaluation against a diverse range of de-watermarking attacks. While many researchers focus on developing robust watermarking techniques, there is often insufficient emphasis on systematically red-teaming these methods against multiple attacking scenarios.

#### Standardization of evaluation benchmarks

There is a need for standardized benchmarks and evaluation metrics to ensure fair and consistent comparison between different watermarking techniques. Table 3 shows how evaluation datasets differ in the literature for the same downstream task, reflecting this necessity.

**Impact on LLM output factuality** Watermarks modify the model output distributions; techniques that are robust to de-watermarking often have greater variations in watermarked outputs compared to clean outputs leading to a potential trade-off between de-watermarking and LLM factuality. Despite this potential trade-off, there is a lack of analysis on how watermarking techniques affect the output inaccuracies or hallucinations. After training or fine-tuning LLMs with specific watermarking techniques, there is often insufficient examination of whether these methods introduce or exacerbate inaccuracies. We advocate for factuality evaluations post-watermarking.

#### Compatibility to various NLP downstream tasks

Important task types like Story Generation, Text Classification etc. are under-explored.

**Enhanced interpretability** Drawing upon security and privacy literature (Kumar et al., 2024), we ask the community to establish privacy norms for LLM watermarking. We envision this to be similar to model cards, which describe the degree of security provided by particular methods against malicious actors.

**Human-centered watermarking** We urge the community to work on human perception of LLMs when interacting with different safety principles. User perception of LLMs may change with differences in output distributions. Furthermore, safety practices may enable AI acceptance and adoption among the masses.

### 4 Conclusion

In this paper, we analyze representative literature in the field and provide a comprehensive taxonomy for digital watermarking techniques for both LLM-generated and human-written text. The taxonomy categorizes watermarking techniques using four primary categories, namely - intention of the method, data used for evaluation, watermark addition, and removal.

Our work not only identifies and clusters existing watermarking methods but also brings to light key open challenges and research gaps in the field. These challenges include the need for more rigorous testing against diverse de-watermarking attacks, the establishment of standardized benchmarks for fair and consistent comparison of different techniques, and a deeper understanding of how watermarking impacts the factuality and accuracy of LLM outputs. Furthermore, we emphasize on the importance of developing watermarking techniques that are resilient to adversarial attacks, enhance interpretability, and maintain compatibility across various NLP downstream tasks.

We envision this research to serve as a reference for policymakers, safety practitioners, and end users; facilitating the adoption of robust digital watermarking practices and promoting responsible AI use.

### 5 Limitations

Limitations to our work are as follows: (1) We do not include detailed insights into metrics for success rate (accuracy of detecting watermarked texts), text quality (perplexity and semantics),



614	NLP task-specific evaluation, and robustness (de-	J.T. Brassil, S. Low, N.F. Maxemchuk, and L. O’Gorman.	666
615	tectability of watermarks after removal attacks).	1995. <a href="#">Electronic marking and identification techniques to discourage document copying</a> . <i>IEEE</i>	667
616	(2) We don’t demonstrate the mathematical analysis	<i>Journal on Selected Areas in Communications</i> ,	668
617	of different watermarking techniques. (3) We	13(8):1495–1504.	669
618	do not cover all different task deployment scenarios		670
619	for the watermarking techniques discussed.		
620	<b>6 Ethical Considerations</b>	Mauro Cettolo, Jan Niehues, Sebastian Stuker, Luisa	671
621	This paper reviews the challenges and opportu-	Bentivogli, and Marcello Federico. 2014. Report on	672
622	nities of watermarking techniques in LLMs. Our	the 11th iwslt evaluation campaign. In <i>Proceedings</i>	673
623	work has many potential societal consequences,	<i>of the 11th International Workshop on Spoken Lan-</i>	674
624	none of which must be specifically highlighted	<i>guage Translation: Evaluation Campaign</i> . Ed.: M.	675
625	here. There are no major risks associated with	Federico, S. Stuker, F. Yvon, page 2, Å17. Association	676
626	conducting this review.	for Computational Linguistics (ACL).	677
627	<b>References</b>	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,	678
628	Sahar Abdelnabi and Mario Fritz. 2021. Adversarial	Henrique Ponde de Oliveira Pinto, Jared Kaplan,	679
629	watermarking transformer: Towards tracing text	Harri Edwards, Yuri Burda, Nicholas Joseph, Greg	680
630	provenance with data hiding. In <i>2021 IEEE Sympo-</i>	Brockman, et al. 2021. Evaluating large lan-	681
631	<i>sium on Security and Privacy (SP)</i> , pages 121–140.	guage models trained on code. <i>arXiv preprint</i>	682
632	IEEE.	<i>arXiv:2107.03374</i> .	683
633	Mikhail J. Atallah, Victor Raskin, Michael Crogan,	Long Dai, Jiarong Mao, Xuefeng Fan, and Xiaoyi Zhou.	684
634	Christian Hempelmann, Florian Kerschbaum, Dina	2022. DeepHider: A covert nlp watermarking frame-	685
635	Mohamed, and Sanket Naik. 2001. Natural lan-	work based on multi-task learning. <i>arXiv preprint</i>	686
636	guage watermarking: Design, analysis, and a proof-	<i>arXiv:2208.04676</i> .	687
637	of-concept implementation. In <i>Information Hid-</i>	Jaiden Fairoze, Sanjam Garg, Somesh Jha, Saeed	688
638	<i>ding</i> , pages 185–200, Berlin, Heidelberg. Springer	Mahloujifar, Mohammad Mahmoody, and	689
639	Berlin Heidelberg.	Mingyuan Wang. 2023. <a href="#">Publicly detectable</a>	690
640	Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian	<a href="#">watermarking for language models</a> . Cryptology	691
641	Wang, Xiaopeng Li, Yuchen Tian, Ming Tan,	ePrint Archive, Paper 2023/1661.	692
642	Wasi Uddin Ahmad, Shiqi Wang, Qing Sun,	<a href="https://eprint.iacr.org/2023/1661">https://eprint.iacr.org/2023/1661</a> .	693
643	Mingyue Shang, Sujan Kumar Gonugondla, Han-	Yu Fu, Deyi Xiong, and Yue Dong. 2024. Watermarking	694
644	tian Ding, Varun Kumar, Nathan Fulton, Arash Fara-	conditional text generation for ai detection: Unveil-	695
645	hani, Siddhartha Jain, Robert Giaquinto, Haifeng	ing challenges and a semantic-aware watermark	696
646	Qian, Murali Krishna Ramanathan, Ramesh Nal-	remedy. In <i>Proceedings of the AAI Conference</i>	697
647	lapati, Baishakhi Ray, Parminder Bhatia, Sudipta	<i>on Artificial Intelligence</i> , volume 38, pages 18003–	698
648	Sengupta, Dan Roth, and Bing Xiang. 2023. <a href="#">Multi-</a>	18011.	699
649	<a href="#">lingual evaluation of code generation models</a> . In	Claire Gardent, Anastasia Shimorina, Shashi Narayan,	700
650	<i>The Eleventh International Conference on Learning</i>	and Laura Perez-Beltrachini. 2017. <a href="#">Creating train-</a>	701
651	<i>Representations</i> .	<a href="#">ing corpora for NLG micro-planners</a> . In <i>Proceed-</i>	702
652	Sören Auer, Christian Bizer, Georgi Kobilarov, Jens	<i>ings of the 55th Annual Meeting of the Association</i>	703
653	Lehmann, Richard Cyganiak, and Zachary Ives.	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	704
654	2007. Dbpedia: A nucleus for a web of open data.	<i>pers)</i> , pages 179–188, Vancouver, Canada. Associa-	705
655	In <i>The Semantic Web</i> , pages 722–735, Berlin, Hei-	tion for Computational Linguistics.	706
656	delberg. Springer Berlin Heidelberg.	Martin Gerlach and Francesc Font-Clos. 2020. A stan-	707
657	Ondřej Bojar, Christian Buck, Christian Federmann,	standardized project gutenber corpus for statistical	708
658	Barry Haddow, Philipp Koehn, Johannes Leveling,	analysis of natural language and quantitative lin-	709
659	Christof Monz, Pavel Pecina, Matt Post, Herve Saint-	guistics. <i>Entropy</i> , 22(1):126.	710
660	Amand, Radu Soricut, Lucia Specia, and Aleš Tam-	Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jin-	711
661	chyna. 2014. <a href="#">Findings of the 2014 workshop on sta-</a>	ran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu.	712
662	<a href="#">tistical machine translation</a> . In <i>Proceedings of the</i>	2023. How close is chatgpt to human experts? com-	713
663	<i>Ninth Workshop on Statistical Machine Translation</i> ,	parison corpus, evaluation, and detection. <i>arXiv</i>	714
664	pages 12–58, Baltimore, Maryland, USA. Associa-	<i>preprint arXiv:2301.07597</i> .	715
665	tion for Computational Linguistics.	Xuanli He, Qionikai Xu, Lingjuan Lyu, Fangzhao Wu,	716
		and Chenguang Wang. 2022a. Protecting intellec-	717
		tual property of language generation apis with lex-	718
		ical watermark. In <i>Proceedings of the AAI Con-</i>	719
		<i>ference on Artificial Intelligence</i> , volume 36, pages	720
		10758–10766.	721

722	Xuanli He, Qiongkai Xu, Yi Zeng, Lingjuan Lyu, Fangzhao Wu, Jiwei Li, and Ruoxi Jia. 2022b. Cater: Intellectual property protection on text generation apis via conditional watermarks. <i>Advances in Neural Information Processing Systems</i> , 35:5431–5445.	David D Lewis, Yiming Yang, Tony Russell-Rose, and Fan Li. 2004. Rcv1: A new benchmark collection for text categorization research. <i>Journal of machine learning research</i> , 5(Apr):361–397.	778 779 780 781
727	Abe Bohan Hou, Jingyu Zhang, Tianxing He, Yichen Wang, Yung-Sung Chuang, Hongwei Wang, Lingfeng Shen, Benjamin Van Durme, Daniel Khashabi, and Yulia Tsvetkov. 2023. <a href="#">Semstamp: A semantic watermark with paraphrastic robustness for text generation</a> . <i>Preprint</i> , arXiv:2310.03991.	Wei Li, Borui Yang, Yujie Sun, Suyu Chen, Ziyun Song, Liyao Xiang, Xinbing Wang, and Chenghu Zhou. 2023. Towards tracing code provenance with code watermarking. <i>arXiv preprint arXiv:2305.12461</i> .	782 783 784 785
733	Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu, Hongyang Zhang, and Heng Huang. 2023. Unbiased watermark for large language models. <i>arXiv preprint arXiv:2310.10669</i> .	Aiwei Liu, Leyi Pan, Xuming Hu, Shu’ang Li, Lijie Wen, Irwin King, and Philip S Yu. 2023a. A private watermark for large language models. <i>arXiv preprint arXiv:2307.16230</i> .	786 787 788 789
737	Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Code-searchnet challenge: Evaluating the state of semantic code search. <i>arXiv preprint arXiv:1909.09436</i> .	Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and Lijie Wen. 2023b. A semantic invariant robust watermark for large language models. <i>arXiv preprint arXiv:2310.06356</i> .	790 791 792 793
741	Zunera Jalil and Anwar M Mirza. 2009. A review of digital watermarking techniques for text documents. In <i>2009 International Conference on Information and Multimedia Technology</i> , pages 230–234. IEEE.	Yixin Liu, Hongsheng Hu, Xuyun Zhang, and Lichao Sun. 2023c. Watermarking text data on large language models for dataset copyright protection. <i>arXiv preprint arXiv:2305.13257</i> .	794 795 796 797
745	Nurul Shamimi Kamaruddin, Amirrudin Kamsin, Lip Yee Por, and Hameedur Rahman. 2018. A review of text watermarking: theory, methods, and applications. <i>IEEE Access</i> , 6:8011–8028.	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. <a href="#">Learning word vectors for sentiment analysis</a> . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.	798 799 800 801 802 803 804 805
749	John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In <i>International Conference on Machine Learning</i> , pages 17061–17084. PMLR.	Ian Magnusson, Akshita Bhagia, Valentin Hofmann, Luca Soldaini, Ananya Harsh Jha, Oyvind Tafjord, Dustin Schwenk, Evan Pete Walsh, Yanai Elazar, Kyle Lo, et al. 2023. Paloma: A benchmark for evaluating language model fit. <i>arXiv preprint arXiv:2312.10523</i> .	806 807 808 809 810 811
754	Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2024. Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. <i>Advances in Neural Information Processing Systems</i> , 36.	Hasan Mesut Meral, Bülent Sankur, A Sumru Özsoy, Tunga Güngör, and Emre Sevinç. 2009. Natural language watermarking via morphosyntactic alterations. <i>Computer Speech &amp; Language</i> , 23(1):107–125.	812 813 814 815 816
759	Rohith Kuditipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. 2023. Robust distortion-free watermarks for language models. <i>arXiv preprint arXiv:2307.15593</i> .	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. <i>arXiv preprint arXiv:1609.07843</i> .	817 818 819
763	Ashutosh Kumar, Sagarika Singh, Shiv Vignesh Murty, and Swathy Ragupathy. 2024. The ethics of interaction: Mitigating security threats in llms. <i>arXiv preprint arXiv:2401.12273</i> .	Vangelis Metsis, Ion Androutsopoulos, and Georgios Paliouras. 2006. Spam filtering with naive bayes - which naive bayes?	820 821 822
767	Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2023. Ds-1000: A natural and reliable benchmark for data science code generation. In <i>International Conference on Machine Learning</i> , pages 18319–18345. PMLR.	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and evaluation framework for deeper understanding of commonsense stories. <i>arXiv preprint arXiv:1604.01696</i> .	823 824 825 826 827 828
773	Taehyun Lee, Seokhee Hong, Jaewoo Ahn, Ilgee Hong, Hwaran Lee, Sangdoo Yun, Jamin Shin, and Gunhee Kim. 2023. Who wrote this code? watermarking for code generation. <i>arXiv preprint arXiv:2305.15060</i> .	Travis Munyer and Xin Zhong. 2023. Deeptextmark: Deep learning based text watermarking for detection of large language model generated text. <i>arXiv preprint arXiv:2305.05773</i> .	829 830 831 832

833	Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing Xiang, et al. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. <i>arXiv preprint arXiv:1602.06023</i> .	Jie Ren, Han Xu, Yiding Liu, Yingqian Cui, Shuaiqiang Wang, Dawei Yin, and Jiliang Tang. 2023. A robust semantics-based watermark for large language model against paraphrasing. <i>arXiv preprint arXiv:2311.08721</i> .	886
834			887
835			888
836			889
837	Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangu Tang, Aadit Vyas, Neha Verma, Pranav Krishna, Yangxiaokang Liu, Nadia Irwanto, Jessica Pan, Faijaz Rahman, Ahmad Zaidi, Mutethia Mutuma, Yasin Tarabar, Ankit Gupta, Tao Yu, Yi Chern Tan, Xi Victoria Lin, Caiming Xiong, Richard Socher, and Nazneen Fatema Rajani. 2021. <b>DART: Open-domain structured data record to text generation</b> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 432–447, Online. Association for Computational Linguistics.	Ryoma Sato, Yuki Takezawa, Han Bao, Kenta Niwa, and Makoto Yamada. 2023. Embarrassingly simple text watermarks. <i>arXiv preprint arXiv:2310.08920</i> .	891
838			892
839			893
840			
841			894
842			895
843			896
844			897
845			898
846			899
847			900
848			901
849			
850			
851	Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. <i>arXiv preprint arXiv:1808.08745</i> .	Zhensu Sun, Xiaoning Du, Fu Song, and Li Li. 2023. Codemark: Imperceptible watermarking for code datasets against neural code completion models. In <i>Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering</i> , pages 1561–1572.	902
852			903
853			904
854			905
855			906
856	The New York Times Company. 2023. <b>The new york times company v. microsoft corporation, openai, inc., openai lp, openai gp, llc, openai, llc, openai opco llc, openai global llc, oai corporation, llc, and openai holdings, llc</b> .	Yuki Takezawa, Ryoma Sato, Han Bao, Kenta Niwa, and Makoto Yamada. 2023. Necessary and sufficient watermark for large language models. <i>arXiv preprint arXiv:2310.00833</i> .	907
857			908
858			
859			909
860			910
861	Nicolai Thorner Sivesind. 2023. Chatgpt-generated-abstracts.	Ruixiang Tang, Qizhang Feng, Ninghao Liu, Fan Yang, and Xia Hu. 2023. Did you train on my dataset? towards public dataset protection with clean-label backdoor watermarking. <i>ACM SIGKDD Explorations Newsletter</i> , 25(1):43–53.	911
862			912
863	OpenAI. 2023. GPT-4 technical report. Technical report.	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rrustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Ruibo Liu,	913
864			914
865	Wenjun Peng, Jingwei Yi, Fangzhao Wu, Shangxi Wu, Bin Zhu, Lingjuan Lyu, Binxing Jiao, Tong Xu, Guangzhong Sun, and Xing Xie. 2023. Are you copying my model? protecting the copyright of large language models for eaas via backdoor watermark. <i>arXiv preprint arXiv:2305.10036</i> .		915
866			916
867			917
868			
869			918
870			919
871	Lip Yee Por, KokSheik Wong, and Kok Onn Chee. 2012. <b>Unispach: A text-based data hiding method using unicode space characters</b> . <i>Journal of Systems and Software</i> , 85(5):1075–1082.		920
872			921
873			922
874			923
875	Wenjie Qu, Dong Yin, Zixin He, Wei Zou, Tianyang Tao, Jinyuan Jia, and Jiaheng Zhang. 2024. Provably robust multi-bit watermarking for ai-generated text via error correction code. <i>arXiv preprint arXiv:2401.16820</i> .		924
876			925
877			926
878			927
879			928
880	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.		929
881			930
882			931
883			932
884			933
885			934
			935
			936
			937
			938
			939
			940
			941
			942
			943
			944

945	Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao	Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen	1009
946	Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad,	Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth,	1010
947	Ale Jakse Hartman, Martin Chadwick, Gaurav Singh	Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi,	1011
948	Tomar, Xavier Garcia, Evan Senter, Emanuel Taropa,	Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat,	1012
949	Thanumalayan Sankaranararyana Pillai, Jacob De-	Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay	1013
950	vlin, Michael Laskin, Diego de Las Casas, Dasha	Bolina, Mariko Iinuma, Polina Zablotskaia, James	1014
951	Valter, Connie Tao, Lorenzo Blanco, Adrià Puig-	Besley, Da-Woon Chung, Timothy Dozat, Ramona	1015
952	domènech Badia, David Reitter, Mianna Chen,	Comanescu, Xiance Si, Jeremy Greer, Guolong Su,	1016
953	Jenny Brennan, Clara Rivera, Sergey Brin, Shariq	Martin Polacek, Raphaël Lopez Kaufman, Simon	1017
954	Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao,	Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie	1018
955	Stephanie Winkler, Emilio Parisotto, Yiming Gu,	Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad	1019
956	Kate Olszewska, Yujing Zhang, Ravi Addanki, An-	Tomasev, Jinwei Xing, Christina Greer, Helen Miller,	1020
957	toine Miech, Annie Louis, Laurent El Shafey, De-	Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma,	1021
958	nis Teplyashin, Geoff Brown, Elliot Catt, Nithya	Angelos Filos, Milos Besta, Rory Blevins, Ted Kli-	1022
959	Attaluri, Jan Balaguer, Jackie Xiang, Pidong Wang,	menko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi	1023
960	Zoe Ashwood, Anton Briukhov, Albert Webson, San-	Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir,	1024
961	jay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-	Vered Cohen, Charline Le Lan, Krishna Haridasan,	1025
962	Wei Chang, Axel Stjerngren, Josip Djolonga, Yut-	Amit Marathe, Steven Hansen, Sholto Douglas,	1026
963	ing Sun, Ankur Bapna, Matthew Aitchison, Pedram	Rajkumar Samuel, Mingqiu Wang, Sophia Austin,	1027
964	Pejman, Henryk Michalewski, Tianhe Yu, Cindy	Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso	1028
965	Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich,	Lorenzo, Lars Lowe Sjösund, Sébastien Cevey,	1029
966	Kehang Han, Peter Humphreys, Thibault Sellam,	Zach Gleicher, Thi Avrahami, Anudhyan Boral,	1030
967	James Bradbury, Varun Godbole, Sina Samangooui,	Hansa Srinivasan, Vittorio Selo, Rhys May, Kon-	1031
968	Bogdan Damoc, Alex Kaskasoli, Sébastien M. R.	stantinos Aisopos, Léonard Hussenot, Livio Baldini	1032
969	Arnold, Vijay Vasudevan, Shubham Agrawal, Jason	Soares, Kate Baumli, Michael B. Chang, Adrià Re-	1033
970	Riesa, Dmitry Lepikhin, Richard Tanburn, Srivat-	casens, Ben Caine, Alexander Pritzel, Filip Pavetic,	1034
971	san Srinivasan, Hyeontaek Lim, Sarah Hodgkinson,	Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ra-	1035
972	Pranav Shyam, Johan Ferret, Steven Hand, Ankush	ramesh, Dan Horgan, Kartikeya Badola, Nora Kass-	1036
973	Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Gi-	ner, Subhrajit Roy, Ethan Dyer, Víctor Campos,	1037
974	ang, Alexander Neitz, Zaheer Abbas, Sarah York,	Alex Tomala, Yunhao Tang, Dalia El Badawy, El-	1038
975	Machel Reid, Elizabeth Cole, Aakanksha Chowd-	speth White, Basil Mustafa, Oran Lang, Abhishek	1039
976	hery, Dipanjan Das, Dominika Rogozińska, Vitaly	Jindal, Sharad Vikram, Zhitao Gong, Sergi Caelles,	1040
977	Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas	Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng,	1041
978	Zilka, Flavien Prost, Luheng He, Marianne Mon-	Wojciech Stokowiec, Ce Zheng, Phoebe Thacker,	1042
979	teiro, Gaurav Mishra, Chris Welty, Josh Newlan,	Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh,	1043
980	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	James Svensson, Max Bileschi, Piyush Patil, Ankesh	1044
981	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	Anand, Roman Ring, Katerina Tsihlias, Arpi Vezer,	1045
982	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom	1046
983	Anirudh Baddepudi, Alex Goldin, Adnan Ozturel,	Kwiatkowski, Samira Daruki, Keran Rong, Allan	1047
984	Albin Cassirer, Yunhan Xu, Daniel Sohn, Deven-	Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg,	1048
985	dra Sachan, Reinald Kim Amplayo, Craig Swans-	Mina Khan, Lisa Anne Hendricks, Marie Pellat,	1049
986	on, Dessie Petrova, Shashi Narayan, Arthur Guez,	Vladimir Feinberg, James Cobon-Kerr, Tara Sainath,	1050
987	Siddhartha Brahma, Jessica Landon, Miteyan Pat-	el, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives,	1051
988	tel, Ruizhe Zhao, Kevin Villela, Luyu Wang, Wen-	Yana Hasson, YaGuang Li, Eric Noland, Yuan Cao,	1052
989	hao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	Nathan Byrd, Le Hou, Qingze Wang, Thibault Sot-	1053
990	Hanzhao Lin, James Keeling, Petko Georgiev, Di-	tiaux, Michela Paganini, Jean-Baptiste Lespiau,	1054
991	ana Mincu, Boxi Wu, Salem Haykal, Rachel Sapu-	Alexandre Moufarek, Samer Hassan, Kaushik Shiv-	1055
992	tro, Kiran Vodrahalli, James Qin, Zeynep Cankara,	akumar, Joost van Amersfoort, Amol Mandhane,	1056
993	Abhanshu Sharma, Nick Fernando, Will Hawkins,	Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew	1057
994	Behnam Neyshabur, Solomon Kim, Adrian Hut-	Brock, Hannah Sheahan, Vedant Misra, Cheng Li,	1058
995	ter, Priyanka Agrawal, Alex Castro-Ros, George	Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu,	1059
996	van den Driessche, Tao Wang, Fan Yang, Shuo	Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener,	1060
997	yiin Chang, Paul Komarek, Ross McIlroy, Mario	Fantine Huot, Matthew Lamm, Nicola De Cao,	1061
998	Lučić, Guodong Zhang, Wael Farhan, Michael	Charlie Chen, Gamaleldin Elsayed, Ed Chi, Mahdis	1062
999	Sharman, Paul Natsev, Paul Michel, Yong Cheng,	Mahtdieh, Ian Tenney, Nan Hua, Ivan Petrychenko,	1063
1000	Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shak-	Patrick Kane, Dylan Scardinaro, Rishub Jain,	1064
1001	eri, Christina Butterfield, Justin Chung, Paul Kis-	Jonathan Uesato, Romina Datta, Adam Sadovsky,	1065
1002	han Rubenstein, Shivani Agrawal, Arthur Mensch,	Oskar Bunyan, Dominik Rabiej, Shimu Wu, John	1066
1003	Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan	Zhang, Gautam Vasudevan, Edouard Leurent, Mah-	1067
1004	Pope, Loren Maggiore, Jackie Kay, Priya Jhakra,	moud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy	1068
1005	Shibo Wang, Joshua Maynez, Mary Phuong, Tay-	Zheng, Betty Chan, Pam G Rabinovitch, Piotr	1069
1006	lor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin	Stanczyk, Ye Zhang, David Steiner, Subhajit Naskar,	1070
1007	Robinson, Yash Katariya, Sebastian Riedel, Paige	Michael Azzam, Matthew Johnson, Adam Paszke,	1071
1008	Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo,	Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz	1072

1073	Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Sahitya Potluri, Jane Park, Elnaz Davoodi, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Chris Gorgolewski, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Paul Suganthan, Evan Palmer, Geoffrey Irving, Edward Loper, Manaal Faruqui, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Michael Fink, Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marin Georgiev, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Minnie Lui, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Lam Nguyen Thiet, Daniel Andor, Pedro Valenzuela, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Sarmishta Velury, Sebastian Krause, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Tejas Latkar, Mingyang Zhang, Quoc Le, Elena Allica Abellan, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Sid Lall, Ken Franko, Egor Filonov, Anna Bulanova, Rémi Leblond, Vikas Yadav, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Hao Zhou, Alek Dimitriev, Hannah Forbes, Dylan Barnarse, Zora Tung, Jeremiah Liu, Mark Omernick, Colton Bishop, Chintu Kumar, Rachel Sterneck, Ryan Foley, Rohan Jain, Swaroop Mishra, Jiawei Xia, Taylor Bos, Geoffrey Cideron, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Petru Gurita, Hila Noga, Premal Shah, Daniel J. Mankowitz, Alex Polozov, Nate Kushman, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Anhad Mohananey, Matthieu Geist, Sidharth Mudgal, Sertan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Quan Yuan, Sumit Bagri, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Aliaksei Severyn, Jonathan Lai, Kathy Wu, Heng-Tze Cheng, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Mark Geller, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Andrei Sozanschi, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Abhimanyu Goyal, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Sabaer Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Tao Zhu, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Dustin Tran, Yeqing Li, Nir Levine, Ariel Stolovich, Norbert Kalb, Rebeca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Balaji Lakshminarayanan, Charlie Deck, Shyam Upadhyay, Hyo Lee, Mike Dusenberry, Zonglin Li, Xuezhi Wang, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Summer Yue, Sho Arora, Eric Malmi, Daniil Mirylenka, Qijun Tan, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Steven Zheng, Francesco Pongetti, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Ragha Kotikalapudi, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Chenkai Kuang, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Pei Sun, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Ishita Dasgupta, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghafarkhah, Morgane Rivièrè, Alanna Walton, Clément Crepy, Alicia Parrish, Yuan Liu, Zongwei Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fildjeland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolicchio, Lexi Walker, Alex Morris, Ivo Penchev, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Adam Kurzrok, Lynette Webb, Sahil Dua, Dong Li, Preethi Lahoti, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Taylan Bilal, Evgenii Eltyshev, Daniel Balle, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, XiangHai Sheng, Emily Xue, Sherjil Ozair, Adams Yu, Christof Angermueller, Xiaowei	1137 1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182 1183 1184 1185 1186 1187 1188 1189 1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1200
------	--	--

1201	Li, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurmurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Kevin Brooks, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Komal Jalan, Dinghua Li, Ginger Perng, Blake Hechtman, Parker Schuh, Milad Nasr, Mia Chen, Kieran Milan, Vladimir Mikulik, Trevor Strohmman, Juliana Franco, Tim Green, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. 2023. <a href="#">Gemini: A family of highly capable multimodal models</a> . <i>Preprint</i> , arXiv:2312.11805.	1258
1202		1259
1203		1260
1204		1261
1205		
1206		1262
1207		1263
1208		1264
1209		1265
1210		
1211		1266
1212		1267
1213		1268
1214		1269
1215		1270
1216	Mercan Topkara, Umut Topkara, and Mikhail J. Atallah. 2006a. <a href="#">Words are not enough: sentence level natural language watermarking</a> . In <i>Workshop on Medical Cyber-Physical Systems</i> .	1271
1217		1272
1218		1273
1219		1274
1220	Umut Topkara, Mercan Topkara, and Mikhail J. Atallah. 2006b. <a href="#">The hiding virtues of ambiguity: quantifiably resilient watermarking of natural language text through synonym substitutions</a> . In <i>Proceedings of the 8th Workshop on Multimedia and Security</i> , MM and Sec '06, page 164–174, New York, NY, USA. Association for Computing Machinery.	1275
1221		1276
1222		1277
1223		1278
1224		
1225		1279
1226		1280
1227	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open and efficient foundation language models</a> . <i>Preprint</i> , arXiv:2302.13971.	1281
1228		1282
1229		
1230		1283
1231		1284
1232		1285
1233		1286
1234	Fangzhao Wu, Ying Qiao, Jiun-Hung Chen, Chuhan Wu, Tao Qi, Jianxun Lian, Danyang Liu, Xing Xie, Jianfeng Gao, Winnie Wu, and Ming Zhou. 2020. <a href="#">MIND: A large-scale dataset for news recommendation</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3597–3606, Online. Association for Computational Linguistics.	1287
1235		1288
1236		
1237		
1238		
1239		
1240		
1241		
1242	Yihan Wu, Zhengmian Hu, Hongyang Zhang, and Heng Huang. 2023. <a href="#">Dipmark: A stealthy, efficient and resilient watermark for large language models</a> . <i>arXiv preprint arXiv:2310.07710</i> .	
1243		
1244		
1245		
1246	Borui Yang, Wei Li, Liyao Xiang, and Bo Li. 2023a. <a href="#">Towards code watermarking with dual-channel transformations</a> . <i>arXiv preprint arXiv:2309.00860</i> .	
1247		
1248		
1249	Xi Yang, Kejiang Chen, Weiming Zhang, Chang Liu, Yuang Qi, Jie Zhang, Han Fang, and Nenghai Yu. 2023b. <a href="#">Watermarking text generated by black-box language models</a> . <i>arXiv preprint arXiv:2305.08883</i> .	
1250		
1251		
1252		
1253	Xi Yang, Jie Zhang, Kejiang Chen, Weiming Zhang, Zehua Ma, Feng Wang, and Nenghai Yu. 2022. <a href="#">Tracing text provenance via context-aware lexical substitution</a> . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 11613–11621.	
1254		
1255		
1256		
1257		
	KiYoon Yoo, Wonhyuk Ahn, Jiho Jang, and Nojun Kwak. 2023a. <a href="#">Robust natural language watermarking through invariant features</a> . <i>arXiv preprint arXiv:2305.01904</i> .	1258
		1259
		1260
		1261
	KiYoon Yoo, Wonhyuk Ahn, and Nojun Kwak. 2023b. <a href="#">Advancing beyond identification: Multi-bit watermark for language models</a> . <i>arXiv preprint arXiv:2308.00221</i> .	1262
		1263
		1264
		1265
	Ruisi Zhang, Shehzeen Samarah Hussain, Paarth Neekhara, and Farinaz Koushanfar. 2023. <a href="#">Remark-llm: A robust and efficient watermarking framework for generative large language models</a> . <i>arXiv preprint arXiv:2310.12362</i> .	1266
		1267
		1268
		1269
		1270
	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. <a href="#">Character-level convolutional networks for text classification</a> . <i>Advances in neural information processing systems</i> , 28.	1271
		1272
		1273
		1274
	Xuandong Zhao, Prabhanjan Ananth, Lei Li, and Yu-Xiang Wang. 2023a. <a href="#">Provable robust watermarking for ai-generated text</a> . <i>arXiv preprint arXiv:2306.17439</i> .	1275
		1276
		1277
		1278
	Xuandong Zhao, Yu-Xiang Wang, and Lei Li. 2023b. <a href="#">Protecting language generation models via invisible watermarking</a> . In <i>International Conference on Machine Learning</i> , pages 42187–42199. PMLR.	1279
		1280
		1281
		1282
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. <a href="#">Judging llm-as-a-judge with mt-bench and chatbot arena</a> . <i>Advances in Neural Information Processing Systems</i> , 36.	1283
		1284
		1285
		1286
		1287
		1288