

Disappearing Ink: Obfuscation Breaks N-gram Code Watermarks in Theory and Practice

Gehao Zhang
University of Massachusetts
Amherst, USA
gehaozhang@umass.edu

Juan Zhai
University of Massachusetts
Amherst, USA
juanzhai@umass.edu

Eugene Bagdasarian
University of Massachusetts
Amherst, USA
eugene@umass.edu

Shiqing Ma
University of Massachusetts
Amherst, USA
shiqingma@umass.edu

Abstract

Large language models (LLMs) are increasingly integrated into code generation workflows. As a result, distinguishing AI-generated code from human-written code is becoming crucial for tasks such as authorship attribution, content tracking, and misuse detection. Based on this purpose, N-gram-based watermarking schemes have emerged as prominent, which inject secret watermarks to be detected during the generation.

However, their robustness in code content remains insufficiently evaluated. Most claims rely solely on defenses against simple code transformations or code optimizations as a simulation of attack, creating a questionable sense of robustness. In contrast, more sophisticated schemes already exist in the software engineering world, e.g., code obfuscation, which significantly alters code while preserving functionality. Although obfuscation is commonly used to protect intellectual property or evade software scanners, the robustness of code watermarking techniques against such transformations remains largely unexplored.

In this work, we focus on the robustness of N-gram-based watermarking approaches on code. We formally model the code obfuscation as a Markov random walk process to attack the watermarking scheme, and prove the impossibility of N-gram-based watermarking’s robustness with only one intuitive and experimentally verified assumption, *distribution consistency*, satisfied. Given the original false positive rate ϵ_{pos} of the watermarking detection, the ratio that the detector failed on the watermarked code after obfuscation will increase to $1 - \epsilon_{\text{pos}}$.

The experiments have been performed on three state-of-the-art watermarking schemes, two large language models, two programming languages, four code benchmarks, and four obfuscators. Among them, *all* watermarking detectors show coin-flipping detection abilities on obfuscated codes (AUROC tightly surrounds 0.5). Among all models, watermarking schemes, and datasets, both programming languages own obfuscators that can achieve attack effects with *no* detection AUROC higher than 0.6 after the attack. Based on the theoretical and practical observations, we also proposed a potential path of robust code watermarking.

1 Introduction

With the rapid rise of applying large language models (LLMs) to various software engineering tasks, code generation has emerged as one of the most prominent areas of interest [6, 22, 35, 56]. However,

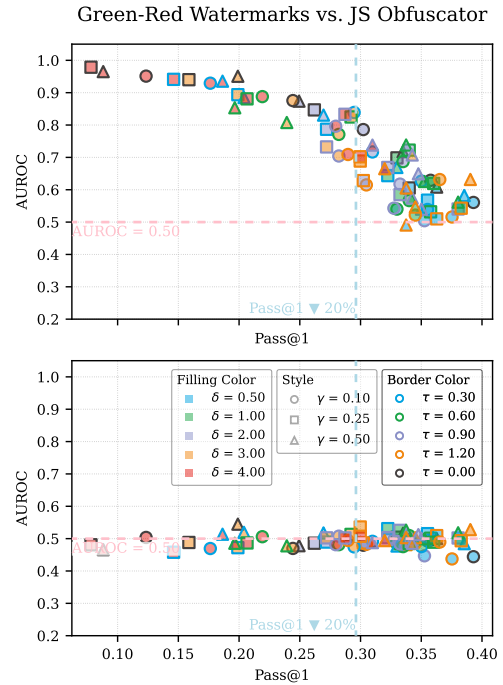


Figure 1: The performance of Green-Red watermarks [25, 32] on DeepSeek-Coder-33B-base model [17] and MBPP-JS benchmark [4], before (upper) and after (lower) applying our attack using the off-the-shelf obfuscator named *JS Obfuscator*. The y-axis denotes the detection’s AUROC score (\uparrow), and the x-axis is the Pass@1 score (\uparrow), representing code generation ability after the distortionary watermarking. Different border colors, filling colors, and shapes of points denote different values of hyperparameters (δ , γ , τ). The blue vertical line illustrates 80% of non-watermarked performance under the same setting. See Figure 4 for more experiment results and information.

this trend also brings significant legal, ethical, and security concerns, including issues related to code licensing, potential plagiarism, software vulnerabilities, and the risk of generating malicious code

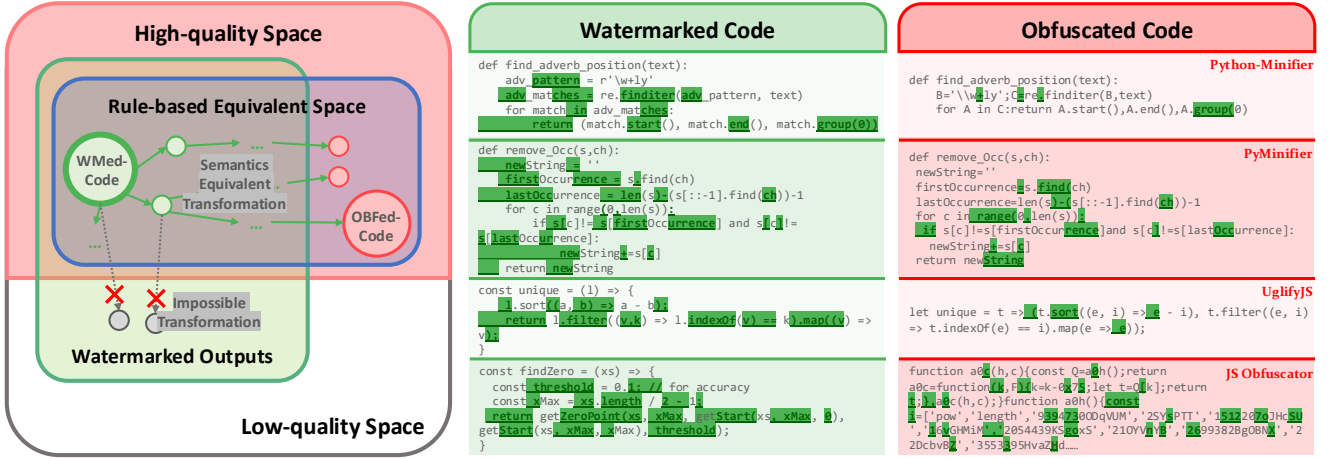


Figure 2: **Left:** An outline of our attack schema. We consider the set of all high-quality outputs (High-quality Space). Our attack randomly modifies codes while preserving semantics by a set of equivalent transformation rules. (i.e., perform code obfuscation). We prove that starting from a high-quality watermarked code sample, the obfuscated code can be independently distributed within a subset of the high-quality space (Theorem 4.1). Furthermore, this subset partitioning is generally independent of N-gram features, resulting in a coin-flipping detection performance (Theorem 5.1; Figure 4). **Middle and Right:** Our attack on watermarked code [25]. The middle table consists of codes generated from the watermarked model. The watermarking scheme adds hidden strategies to generate more "green tokens" that the detector counts to judge whether the text was watermarked. The rightmost table includes the corresponding code after obfuscation. The ratio of "green tokens" is decreased. The names of the obfuscators are shown in red on the far right. All code pairs are ensured by the obfuscators to have the same semantics.

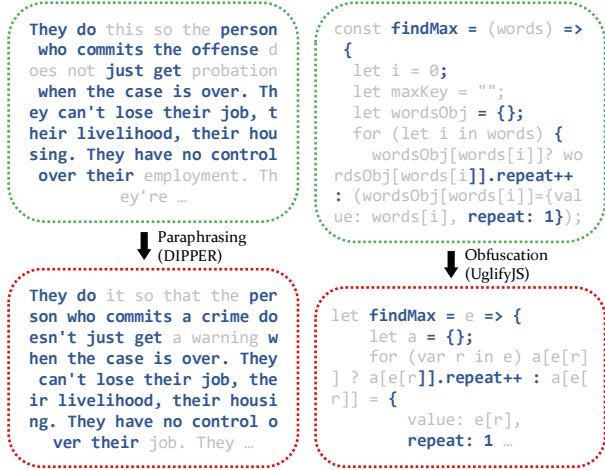


Figure 3: Comparison between Natural Language and Programming Language obfuscations. The blue color indicates the 2-grams that were preserved after the modification.

[16, 44, 48]. As such, the development of practical tools for detecting machine-generated code has become both timely and critically important to ensure the fair, secure, and responsible deployment of LLMs in programming-related applications.

To this end, several efforts have been made toward LLMs watermarking, which aims to embed hidden signals within LLM-generated content for later identification. Notably, Google recently

introduced SynthID [10], a non-distortionary watermarking method that preserves the natural distribution of generated outputs. In contrast, WLLM [25] represents a distortionary approach (improving watermark detectability at the cost of text quality), as it alters the token distribution during the sampling process. Building upon WLLM, SWEET [32] introduces a low-entropy token skipping mechanism specifically tailored for watermarking code generation tasks. These state-of-the-art watermarking schemes take both quality preservation and detectability into account, and some have been large-scale applied in industrial scenarios [10].

Although watermarking strategies vary in their specific implementations, a canonical pipeline can be observed: the previous (N-1)-gram is hashed into a pseudo-random seed, which is then used during the watermarking process to influence the selection of the next token. As a result, the detection procedure depends on intact N-grams to recover the corresponding pseudo-random seeds. Each seed, along with the following token, is evaluated to determine whether it aligns with the watermarking strategies. For instance, SynthID [10] uses a 5-gram approach, where the preceding four tokens are hashed into a pseudo-random seed. This seed is used to generate random g -scores across the vocabulary, and the g -score of the selected token is then used as a metric in the detection process.

Similar to SynthID, considerable modern LLM watermarking approaches [10, 14, 25, 26, 32, 36, 57, 65] adopt this standard framework and rely on complete N-gram sequences for accurate detection (called N-gram-based in this work).

On the opposite, in the natural language processing field, some works [30, 47] try to nullify watermarks by paraphrasing to explore the robustness of existing watermarking approaches. With the

ability to modify the original text to a large extent but keep the same semantics, these paraphraser can break the N-grams in the original watermarked text.

Extending the notion of paraphrasing natural languages, programming languages would be more modify-feasible due to formal language features. Compared with natural language paraphraser, which rely on LLMs (e.g., fine-tuned T5-XXL model in DIPPER [30]) and semantics equivalence unproven, programming languages have a large space of retaining semantics equivalence even after drastic modifications. E.g., variable name changing, removing/adding unreachable code, and modifying among different types of loop/branch grammars. Further, in software engineering, many *code obfuscators* have existed for the needs of semantics-preserving code transformation [2, 3, 45, 49]. Obfuscators have been employed as a means to protect intellectual property and prevent reverse engineering while also being widely used by malware developers to evade software scanners [60]. As a well-known field with numerous existing works, code obfuscators should be the top choice for malicious users to evade code watermarking detection and for researchers to evaluate the robustness of watermarking approaches.

Two pairs of examples are shown in Figure 3 compare the modifications made by the paraphraser for natural language and the obfuscator for code. We can observe that after paraphrasing, not only is the semantics equivalence not guaranteed (e.g., probation \rightarrow warning), but also many 2-grams remain, which still have the possibility to be detected as watermarked. In contrast, after code obfuscation, the resulting code is both equivalency-proven and more disturbed.

Based on the factors above, although there were several existing works specifically focused on watermarking LLM-generated code content, their claims of robustness are questionable. SWEET [32] and ACW [34] use code optimizers and refactoring to simulate the attack. MCGMark [43] claims robustness by defending eight types of simple modifications defined by themselves, and CodeIP [16] only considers the cropping of code but no further modifications. As far as we know, code obfuscators were never used to evaluate the robustness of LLM watermarking, as they were supposed to. The robustness of existing watermarking approaches on code is still uncertain and tends to be overestimated.

Theoretical Result: We illustrate an outline of our attack model in Figure 2. In this work, inspired by the common *equivalent transformations* in the software engineering field, we formally model the code obfuscation with our *rule-based transformation* setting. We prove that for each N-gram-based watermarking scheme, its robustness will be nullified with an intuitive and experiment-supported assumption, *distribution consistency*, satisfied (Theorem 5.1).

Specifically, given the original false positive rate ϵ_{pos} of the watermarking detection, the ratio that the detector failed on the watermarked code after obfuscation will increase to $1 - \epsilon_{\text{pos}}$, which means that the detection algorithm entirely loses the ability to distinguish watermarked code from benign code. (See Theorem 5.1 and Theorem 5.2)

Experimental Result: We leverage our attack algorithm by using off-the-shelf obfuscators, performing a low-cost but effective attack on three state-of-the-art N-gram-based watermarking works. Specifically, we obfuscate watermarked codes and evaluate the

afterward detectabilities. Real watermarked code examples before and after our attack are exhibited in Figure 2.

Extensive experiments have been done among three watermarking schemes, two LLMs, two programming languages (Python and JavaScript), four code benchmarks, and four obfuscators. Among our experiments, *all* watermarking detectors show coin-flipping detection abilities (AUROC tightly surrounds 0.5) on obfuscated codes. Among all models, programming languages, watermarking schemes, and datasets, there exist obfuscators to make *no* AUROC score higher than 0.6.

A group of code quality and detectability can be found in Figure 1. We can observe that after our attack, without affecting code quality, the detection AUROC scores tightly surround 0.5, showing random guessing detection abilities. The same effect can be found among different models, watermarking schemes, and programming languages (Subsection 6.2).

We further assess our main assumption *distribution consistency*. The assumption is satisfied in 98.10% cases in the experiment, supporting our theory.

Our code is available at <https://anonymous.4open.science/r/Code-WM-E440/>

Comparison with Previous Work: Recently, an inspiring work [61] has proposed an assumption and theoretical proof of watermarking robustness impossibility under the assumption. Compared with previous work, our differences are: **1)** We prove and discuss that, on code generation tasks, their assumption is theoretically impossible and low efficiency in practice. (See appendix A) **2)** Both based on the random walk on the directed graph, they construct an ideal graph and force set its ergodicity, while our model is implementable and has been proven ergodic. (See Theorem 4.1)

2 Related Works

Large Language Model Watermarking. In the era of large language models (LLMs), increasing attention has been directed toward advanced watermarking techniques that embed watermarks directly into the text generation process [64]. Among them, the Green-Red Watermark [20, 25, 26, 32, 36, 39, 63, 68] embeds watermark messages into the logits produced by LLMs without altering the model parameters. In the context of code generation, SWEET [32] presents a Green-Red watermarking scheme tailored to mitigate the quality degradation caused by the low-entropy nature of programming tasks. The Gumbel Watermark [1, 14, 19, 31, 57, 65] introduces a bias in the sampling choice that yields a computationally distortion-free watermark. Additionally, SynthID [10] proposes tournament sampling, which increases the expected detection score on watermarked text while maintaining the same output distribution as standard sampling.

Despite the differences among Green-Red, Gumbel, and SynthID, considerable current watermarking techniques rely heavily on N-gram-based mechanisms [10, 14, 25, 26, 32, 36, 57, 65].

Paraphrasing Attack. Paraphrase attacks commit to constructing semantics equivalent or similar adversarial inputs, leveraging the natural language processing applications' weakness of being vulnerable to adversarial perturbations, where a small change to the input produces an undesirable change in system behavior [40]. E.g., lower-quality translations from machine translation

systems or false/improper statements from intelligent chatbots [7, 8, 12, 53, 62, 66]. Some paraphrasing approaches are also proposed to attack LLM watermarking. DIPPER [30] is fine-tuned from the T5 11B model. On DIPPER paraphrased inputs, performance degradation is observed on multiple watermarking, outlier detection, and classifier methods [15, 25, 29, 42]. Moreover, a recursive paraphrasing attack [47] that uses neural network-based paraphrasing to recursively paraphrase the source LLM's output text also reveals the unreliability of state-of-the-art AI text detectors.

Code Obfuscation. The obfuscation is a technique that makes programs harder to understand. For such a purpose, it converts a program to a new different version while making them functionally equal to each other [60]. Different obfuscators transfer various levels of code, including source code, Java bytecode, binary code, leveraging dead-code insertion, register reassignment, subroutine reordering, instruction substitution, code transposition, code integration, etc. [28, 55], and are widely used by malware writers to evade antivirus scanners [60].

3 Notion Formalizations

In this section, partly following previous works' definitions [61, 64], we first formalize our threat model (Subsection 3.1). After that, we create and define the notions about equivalent transformation to formalize our attack in Subsection 3.2. Finally, we model our random walk algorithm into a graph setting in Subsection 3.3.

3.1 Threat Model

Given a prompt space \mathcal{X} , a code space \mathcal{C} , a randomized code generation model $M : \mathcal{X} \rightarrow \mathcal{C}$, and a code test suite $T : \mathcal{X} \times \mathcal{C} \rightarrow \{0, 1\}$ checks whether a code meets a prompt's requests, a secret-key watermarking scheme $\text{"} = (\text{Watermark}, \text{Detect})$ consists of two algorithms:

- **Watermark(M):** Given $M \in \mathcal{M}$, this randomized watermarking algorithm outputs a secret key $k \in \mathcal{K}$ and a watermarked model $M_k : \mathcal{X} \rightarrow \mathcal{C}$, dependent on k .
- **Detect $_k(x, c)$:** Accepting secret key $k \in \mathcal{K}$, $x \in \mathcal{X}$, and $c \in \mathcal{C}$, this deterministic detection algorithm $\text{Detect}_k(x, c) \in \{0, 1\}$ returns a decision bit, where 1 indicates the presence of the watermark, and 0 indicates its absence.

The watermarked model M_k is supposed to generate secret watermarks-injected content that can be detected by Detect_k , while a secret key k could be chosen or held by the generative model provider or the public.

For the detection ability of watermarking schemes, false negative rate (FNR) $\epsilon_{\text{neg}} \leftarrow \mathbb{E}(\mathbb{1}[\text{Detect}_k(x, c) = 0] \mid c \in C_k)$ and false positive rate (FPR) $\epsilon_{\text{pos}} \leftarrow \mathbb{E}(\mathbb{1}[\text{Detect}_k(x, c) = 1] \mid c \in C)$ can be used to assess, in which $C_k \subseteq C$, denoting the set of all watermarked responses $M_k(x)$. Typically, both ϵ_{pos} and ϵ_{neg} need to be low enough for a watermarking scheme to work.

Moreover, for prompt x , we define \mathcal{F} to denote the prompt-required behavior and $\Phi^{\mathcal{F}} \subseteq \mathcal{C}$ to denote a *high-quality space*, which includes all possible code c that can make $T(x, c) = 1$. The *quality* of code could have many aspects. In this work, the prompt-required behavior \mathcal{F} and test suite T notions only include and evaluate the runtime behavior of code, such as correctness and

functional compliance, since correct execution behavior is often considered the most essential among the various dimensions of code quality, as reflected in its central role in popular code generation benchmarks [4, 6, 37, 67].

Adversary Attacker: An adversary attacker $\text{Att} : \mathcal{X} \times \mathcal{C} \rightarrow \mathcal{C}$ performs as a code modifier (obfuscator) and guarantees $\forall c \in \Phi^{\mathcal{F}}, \text{Att}(x, c) \in \Phi^{\mathcal{F}}$ (i.e., $T(x, \text{Att}(x, c)) = 1$), which means that given the prompt x , the corresponding original high-quality code c , and the attacker Att , the modified code $c' \leftarrow \text{Att}(x, c)$ can be obtained while maintaining the quality of the code.

We define ϵ -breaking as the attack goal. An attacker ϵ -breaks the watermarking scheme if for every $M \in \mathcal{M}$ and every prompt $x \in \mathcal{X}$, we have:

$$\mathbb{E}(\mathbb{1}[\text{Detect}_k(x, \text{Att}(x, c)) = 0] \mid c \in C_k) \geq \epsilon,$$

which means that for " , after quality-ensured transformation by Att , the "afterward false negative rate" becomes equal to or greater than ϵ .

In contrast, on the defender side, i.e., the watermarking scheme side, if all ϵ -breaking adversary attackers can only achieve $\epsilon \approx \epsilon_{\text{neg}}$, meaning that the watermarking scheme can maintain a low false negative rate confronting attackers' code modifications, the watermarking scheme " can be qualified as having high robustness.

3.2 Equivalent Transformations

In this section, we introduce the notions of equivalent transformation rule and rule-based equivalent space since we will leverage transformation rules to compose attackers in the following sections.

Definition 3.1 (Equivalent Transformation Rule and Executor). Given original code $c \in \mathcal{C}$, two algorithms can be applied to perform equivalent transformations:

- **R(c):** With a specific equivalent transformation rule endowed, each R , which is an action derivation algorithm (we may directly call a R algorithm as "rule" sometimes), returns a set of executable actions ($R : \mathcal{C} \rightarrow \bigcup_{i=1}^{\infty} \mathcal{A}^i$) that the $\text{Transform}(\cdot)$ algorithm can perform on c . Within the rule space \mathcal{R} , different R s perform different types of transformation (in which we use $a \in \mathcal{A}$ to denote a specific transformation action, and the space of transformation actions is \mathcal{A}).
- **Transform(c, a):** Being global to different R s, given code and action $a \in \mathcal{A}$ (a should be an element of the set R returned), $\text{Transform} : \mathcal{C} \times \mathcal{A} \rightarrow \mathcal{C}$ performs a quality-ensured code transformation that guarantee:

$$\forall \mathcal{F}, \forall c \in \Phi^{\mathcal{F}}, \text{Transform}(c, a) \in \Phi^{\mathcal{F}},$$

in which the a needs to be an action that is allowable to perform the current rule.

The purpose of this definition is that each R corresponds to one abstract transformation concept, while Transform takes practical tuples of code and action as inputs, usually performing concrete modification actions under this concept.

We can consider two different R s, in which the first one can be described as "delete a comment" and returns actions that each action represents delete a specific comment from a specific location of c , while the second one is "modify a for-loop to while-loop" and each

Algorithm 1: RandomStep

Input: Original Code: $c_{\text{ori}} \in C$,
 Rule Set: $\Gamma = \{R_1, R_2, \dots, R_n\}$
Output: Transformed Code: $c_{\text{trans}} \in C$

- 1 Init Action Set $\Lambda \leftarrow \emptyset$
- 2 **for each** i in $1, 2, \dots, n$ **do**
- 3 $\Lambda \leftarrow \Lambda \cup R_i(c)$
- 4 Init $j \leftarrow \text{RandomChoose}(\{1, 2, \dots, |\Lambda|\})$ // Λ is
 $\{a_1, a_2, \dots, a_{|\Lambda|}\}$, and $\Lambda \subseteq \mathcal{A}$
- 5 Init $c_{\text{trans}} \leftarrow \text{Transform}_i(c_{\text{ori}}, a_j)$
- 6 Return c_{trans}

returns a represent a modification on a specific existing for-loop. Note that many of the quality-ensured rules can be implemented in the real world and are common in the programming analysis field, e.g., comment removing, variable renaming, dead code insertion, etc. ^{1 2 3}

Then, we define a randomized algorithm that applies an equivalent transformation rule.

Definition 3.2 (Random Transformation). Leveraging an algorithm $\text{RandomChoose} : \bigcup_{i=1}^{\infty} \mathbb{N}^i \rightarrow \mathbb{N}$ that can uniformly and randomly select one integer from an integer set, two randomized algorithms $\text{RandomStep}(c, \Gamma)$ and $\text{RandomWalk}(c, \Gamma, t)$ are defined in Algorithm 1 and Algorithm 2.

The algorithm RandomStep (Algorithm 1) performs a quality-ensured transformation on the original code. It randomly selects an action from the action set Λ (line 4). The purpose of applying it is to randomly select one action from all possible actions that can be performed by the rule set Γ . Based on RandomStep , RandomWalk takes a finite number t and calls $\text{RandomStep}(c, \Gamma)$ t times.

We also clarify that, in our modeling, both RandomStep and RandomWalk do *not* rely on the prompt or prompt-requested behaviors. In the software engineering field, with the understanding of programming language grammar, with only original code provided, there are many transformation rules on the shelves that can apply semantics-equivalent transforms.

At the end of this section, we define the notions *ergodicity rule set* and *rule-based equivalent space*, which we will analyze afterward.

Definition 3.3 (Ergodicity Rule Set). A $\Gamma_{\text{ERG}} = \{R_1, R_2, \dots, R_n\}$ denotes a set of rules that meet the following properties:

- Empty Rule: Each Γ_{ERG} has a transformation rule that performs empty conversions, i.e., for all $x \in X$, $c \in C$ and $a \in R(c)$, the Transform algorithm of this rule satisfies:

$$c = \text{Transform}(c, a).$$

- Inverse Rule: Each $R \in \Gamma_{\text{ERG}}$ has its inverse rule R^{-1} that satisfies: Given $\Lambda = R(c)$ and $c' = \text{Transform}(c, a)$, for each action $a \in \Lambda$, there is at least one $a^{-1} \in R^{-1}(c')$ to satisfy:

$$c = \text{Transform}(c', a^{-1}).$$

¹<https://github.com/mishoo/UglifyJS>

²<https://github.com/javascript-obfuscator/javascript-obfuscator>

³<https://pyobfuscate.com/pyd>

Algorithm 2: RandomWalk

Input: Original Code: $c_{\text{ori}} \in C$,
 Rule Set: $\Gamma = \{R_1, R_2, \dots, R_n\}$,
 Number of Steps: $t \in \mathbb{N}$
Output: Transformed Code: $c_{\text{trans}} \in C$

- 1 Init $c \leftarrow c_{\text{ori}}$
- 2 **for each** i in $1, 2, \dots, t$ **do**
- 3 $c \leftarrow \text{RandomStep}(c, \Gamma)$
- 4 Init $c_{\text{trans}} \leftarrow c$
- 5 Return c_{trans}

Since each step is reversible, we can deduce that for all $c_+, c_- \in C$, and all possible $t \in \mathbb{N}$, a Γ_{ERG} satisfies:

$$\Pr[c_+ = \text{RandomWalk}(c_-, \Gamma_{\text{ERG}}, t)] > 0$$

$$\iff$$

$$\Pr[c_- = \text{RandomWalk}(c_+, \Gamma_{\text{ERG}}, t)] > 0.$$

It should be highlighted that the ergodicity rule set is implementable but not just a theoretical concept. An example is provided in Table 1, each rule of which has an inverse rule. We can notice that some rules might be inverse rules of themselves. Specifically, for rules IV and its inverse rule III, although there may be a low probability to reverse the comment deletion by randomly adding a comment, it is still nonzero and meets the constraint in Definition 3.3.

Definition 3.4 (Rule-based Equivalent Space). Given a prompt x , a corresponding initial code $c \in \Phi^{\mathcal{F}}$, a rule set $\Gamma = \{R_1, \dots, R_n\}$, and a *all possible finite* step number $t \in \mathbb{N}$, a rule-based equivalent space Φ_c^{Γ} (in which $\Phi_c^{\Gamma} \subseteq \Phi^{\mathcal{F}} \subseteq C$) consists of the initial code c and all possible $\text{RandomWalk}(c, \Gamma, t)$.

About Definition 3.4, in other words, when we perform quality-ensured transformations using the abilities Γ given, and stop after no matter how many steps, the result will be included in Φ_c^{Γ} .

3.3 Chain Modeling

In this section, we will model our rule-based code transformation as a graph, since we will disturb the watermarked context by its in-graph random walk defined in Definition 3.2.

Given initial code $c \in C$ and rule set Γ , leading to a rule-based equivalent space Φ_c^{Γ} , the code transformation graph is a weighted and directed graph $G_c^{\Gamma} = (V, E, W)$ that reflects the transformation relations inside Φ_c^{Γ} . Each vertex $v_i \in V$ denotes its corresponding

Table 1: Ergodicity Rule Set Example. R's indexes constitute the leftmost row, and R^{-1} for the middle row.

R	R^{-1}	DESCRIPTION
I	I	Do nothing.
II	II	Randomly rename variables.
III	IV	Add a random comment.
IV	III	Delete a comment.
V	VI	Add a random dead code snippet.
VI	V	Delete a dead code snippet.

code c_i when $\Phi_c^\Gamma = \{c_1, c_2, \dots, c_i, \dots, c_n\}$ given, and an edge $e_{ij} \in E$ exists if $\Pr[c_j = \text{RandomStep}(c_i, \Gamma)] > 0$, which means c_i can be directly transformed to c_j within one step by rules in Γ . In weight set W , each weight $w_{ij} \in (0, 1]$ corresponds to an existing edge $e_{ij} \in E$, and $w_{ij} = \Pr[c_j = \text{RandomStep}(c_i, \Gamma)]$. Based on the Definition 3.2 and Algorithm 1 that the next visited vertex is selected uniformly and randomly, we always have:

$$\forall w_{i+}, w_{i-} \in \{w_{ij} \mid e_{ij} \in E, \forall j\}, \\ w_{i+} = w_{i-}.$$

On the basis of C_c^Γ definition, we define the transition matrix $\vec{P} \in \mathbb{R}^{|V| \times |V|}$ of the graph:

$$\vec{P}_{ij} = \begin{cases} w_{ij}, & e_{ij} \in E, \\ 0, & \text{otherwise.} \end{cases}$$

Further, a probability distribution vector $p(c_i, t) \in \mathbb{R}^{|V|}$ satisfies that: $p(c_i, t)_j = \Pr[c_j = \text{RandomWalk}(c_i, \Gamma, t)]$. In other words, $p(c_i, t)_j \in [0, 1]$ denotes that starting from c_i after the t -th call of RandomStep (line 3), the probabilities that our random-walking code is now transformed into c_j . E.g., $p(c_1, 0) = [1, 0, 0, \dots, 0]$ denotes that c_1 is our starting point in the initial state. With the definition of the transition matrix given, we can notice that: $p(c, t+1) = \vec{P}^\top \cdot p(c, t)$.

For a matrix \vec{P} that satisfies the stationary distribution condition (i.e. irreducibility and aperiodicity satisfied on the corresponding Markov chain)[33], a stationary distribution $\pi \in \mathbb{R}^{|V|}$ exists and satisfies constraints that:

$$\exists \pi, \forall c \in \Phi_c^\Gamma, \lim_{t \rightarrow \infty} p(c, t) = \pi.$$

Following classical Markov chain theory, we also define the mixing time $t_c(\epsilon) = \min \{t \mid \|p(c, t) - \pi\|_{TV} \leq \epsilon\}$, denoting that starting from the initial code c , the distribution of the Markov chain is within total variation distance ϵ of the stationary distribution π after t steps.

4 Model Analysis

In this section, we prove that the stationary distribution condition exists for all ergodicity rule sets corresponding code transformation graphs (Subsection 4.1). In Subsection 4.2, we further prove that each ergodicity rule set can define a partition on the high-quality code space.

4.1 The Stationary Distribution Exists

In this section, inside the rule-based equivalent space for an ergodicity rule set, we will prove the existence of RandomWalk 's stationary state, which is independent of the walk's starting point. Specifically, based on the Markov theory [33], the stationary state exists when the *irreducibility* and *aperiodicity* of the corresponding Markov chain are satisfied.

Theorem 4.1 (Inside-Space Independence). For any rule-based equivalent space $\Phi_c^{\Gamma_{\text{ERG}}}$ whose corresponding rule set is satisfied as an ergodicity rule set (Definition 3.3), it satisfies:

$$\exists \pi, \forall c \in \Phi_c^{\Gamma_{\text{ERG}}}, \lim_{t \rightarrow \infty} p(c, t) = \pi,$$

in which the definition of the rule-based equivalent space is Definition 3.4, and the vector $p(\cdot)$ is defined in Subsection 3.3.

In other words, $\Phi_c^{\Gamma_{\text{ERG}}}$'s related code transformation graph $C_c^{\Gamma_{\text{ERG}}}$ is both irreducible and aperiodic.

PROOF OF THEOREM 4.1'S IRREDUCIBILITY. First, we can focus on the "original code" concept in Algorithm 2, and based on Definition 3.4, we define that rule-based equivalent space consists of all RandomWalk -achievable codes from the original code $c \in C$. Therefore:

$$\forall c' \in \Phi_c^{\Gamma_{\text{ERG}}}, \exists t, \Pr[c \rightarrow c' \mid t] > 0,$$

in which we use an abbreviation $\Pr[c \rightarrow c' \mid t]$ to denote $\Pr[c' = \text{RandomWalk}(c, \Gamma_{\text{ERG}}, t)]$. Then, based on Definition 3.3 that for all $c_+, c_- \in C$, a Γ_{ERG} satisfies:

$$\Pr[c_+ \rightarrow c_- \mid t] > 0 \iff \Pr[c_- \rightarrow c_+ \mid t] > 0,$$

from which we can deduce that for all $c_+, c_- \in \Phi_c^{\Gamma_{\text{ERG}}}$ we have:

$$\because \Pr[c \rightarrow c_+ \mid t_+] > 0 \Rightarrow \Pr[c_+ \rightarrow c \mid t_+] > 0;$$

$$\because \Pr[c \rightarrow c_- \mid t_-] > 0;$$

$$\therefore \Pr[c_+ \rightarrow c_- \mid t_+ + t_-] \geq \Pr[c_+ \rightarrow c \mid t_+] \cdot \Pr[c \rightarrow c_- \mid t_-] > 0.$$

Thus, all vertices in the code transformation graph are mutually reachable via finite-length paths that pass through the original code c , implying that it is irreducible under Γ_{ERG} setting. \square

PROOF OF THEOREM 4.1'S APERIODICITY. From Definition 3.3, due to the existence of the empty conversion (self-loop), for all $\Phi_c^{\Gamma_{\text{ERG}}}$ which corresponded with an ergodicity rule set Γ_{ERG} , for all $t \in \mathbb{N}$, we have:

$$\because \Pr[c = \text{RandomStep}(c, \Gamma_{\text{ERG}})] > 0$$

$$\therefore \Pr[c = \text{RandomWalk}(c, \Gamma_{\text{ERG}}, t)] \geq$$

$$(\Pr[c = \text{RandomStep}(c, \Gamma_{\text{ERG}})])^t > 0.$$

Therefore,

$$\{t \mid c = \text{RandomWalk}(c, \Gamma, t)\} = \mathbb{N}$$

$$\gcd(\{t \mid c = \text{RandomWalk}(c, \Gamma, t)\}) = 1.$$

\square

Note that in the Theorem 4.1, the definition of the space $\Phi_c^{\Gamma_{\text{ERG}}}$ relies on both the rule set Γ_{ERG} and code c . With a rule set and any initial code $c \in \Phi^{\mathcal{F}}$ given, a space $\Phi_c^{\Gamma_{\text{ERG}}}$ can be sketched out then.

Given the equivalent space, the Theorem 4.1 means that, when we choose a start point in this given equivalent space and do a random walk, after long enough t steps, the final probability distribution will always be π (i.e., π is this space's intrinsic property), being independent of the choice of the start point.

4.2 Rule-based Partition

For the needs of subsequent proof, we prove the partitionability of the code space.

Theorem 4.2 (Ergodicity Rule Set can Partition the High-Quality Code Space). Given an ergodicity rule set Γ_{ERG} , for all $c_+, c_- \in \Phi^{\mathcal{F}}$ ($\Phi^{\mathcal{F}} \subseteq C$) we have:

$$\neg(\Phi_{c_+}^{\Gamma_{\text{ERG}}} = \Phi_{c_-}^{\Gamma_{\text{ERG}}}) \iff \Phi_{c_+}^{\Gamma_{\text{ERG}}} \cap \Phi_{c_-}^{\Gamma_{\text{ERG}}} = \emptyset$$

In other words, if two codes are not in the same rule-based equivalent spaces, their spaces are non-overlapping, which means that each Γ_{ERG} can define a partition of the high-quality code space.

PROOF OF THEOREM 4.2. Follow the settings of Theorem 4.2, we can deduce that (we reuse the abbreviation $\Pr[\rightarrow]$ in Theorem 4.1's proof, and we abbreviate $\Phi_{c_+}^{\Gamma_{\text{ERG}}}$ and $\Phi_{c_-}^{\Gamma_{\text{ERG}}}$ as Φ_+ and Φ_-):

$$\begin{aligned} & \text{Given: } \Phi := \Phi_+ \cap \Phi_-, \Phi \neq \emptyset, \\ & \Rightarrow \forall c_+ \in \Phi_+, \exists c \in \Phi, t_+ \in \mathbb{N}, \Pr[c_+ \rightarrow c \mid t_+] > 0. \\ & \because \forall c_- \in \Phi_-, \exists t_- \in \mathbb{N}, \Pr[c \rightarrow c_- \mid t_-] > 0, \\ & \therefore \forall c_+ \in \Phi_+, c_- \in \Phi_-, \\ & \Pr[c_+ \rightarrow c_- \mid t_+ + t_-] \geq \Pr[c_+ \rightarrow c \mid t_+] \cdot \Pr[c \rightarrow c_- \mid t_-] > 0, \\ & \Rightarrow \forall c_- \in \Phi_-, c_+ \in \Phi_+, \\ & \text{and by symmetry, } \forall c_+ \in \Phi_+, c_- \in \Phi_-, \\ & \Rightarrow \Phi_+ = \Phi_- \end{aligned}$$

□

Definition 4.1 (Rule-based Partition). Following Theorem 4.2, given high-quality code space $\Phi^{\mathcal{F}}$ and an ergodicity rule set Γ_{ERG} , a rule-based partition $\mathcal{P}^{\Gamma_{\text{ERG}}}$ is:

$$\mathcal{P}^{\Gamma_{\text{ERG}}} = \{\Phi_1^{\Gamma_{\text{ERG}}}, \Phi_2^{\Gamma_{\text{ERG}}}, \dots, \Phi_n^{\Gamma_{\text{ERG}}}\}$$

, which satisfies (we abbreviate $\mathcal{P}^{\Gamma_{\text{ERG}}}$ to \mathcal{P}):

$$\begin{aligned} & \forall \Phi \in \mathcal{P}, \Phi \subseteq \Phi^{\mathcal{F}}; \\ & \bigcup_{\Phi \in \mathcal{P}} \Phi = \Phi^{\mathcal{F}}; \\ & \forall \Phi_+, \Phi_- \in \mathcal{P}, \Phi_+ \cap \Phi_- = \emptyset. \end{aligned}$$

This means $\mathcal{P}^{\Gamma_{\text{ERG}}}$ is the result of partitioning high-quality code space $\Phi^{\mathcal{F}}$ by Γ_{ERG} .

5 Impossibility Results

5.1 Assumption and Robustness Impossibility

In this section, we will describe our only assumption, which is intuitive and has been supported by experiments.

First, we define the target of our assumption, distribution consistency.

Definition 5.1 (Distribution Consistency). We define $Q(C)$ to denote a distribution over the code space C , where Q is independent of the watermarking scheme. Given a partition $\mathcal{P}^{\Gamma_{\text{ERG}}}$, a watermarking scheme $\sim = (\text{Watermark}, \text{Detect})$ satisfies the *distribution consistency* property if, for all $Q, x \in X, c \in C$, and $k \in \mathcal{K}$, given $c \sim Q(C)$, $\text{Detect}_k(x, c) \sim \mathcal{D}$, we have:

$$\text{if } c \sim Q(\Phi), \text{ then } \text{Detect}_k(x, c) \sim \mathcal{D}.$$

That is, the observed distribution of $\text{Detect}_k(\cdot, \cdot)$'s results within each equivalent space in $\Phi \in \mathcal{P}^{\Gamma_{\text{ERG}}}$ remains the same as the distribution observed over the entire code space C .

We remind that the distribution consistency is a necessary and sufficient condition for the independence between Γ_{ERG} and Detect. Moreover, the $\text{Detect}(x, c)$ among the whole code space C follows the Bernoulli distribution, $\text{Bernoulli}(\epsilon_{\text{pos}})$.

Assumption 5.1 (Impossibility Assumption). With implementable code transformation rules, distribution consistency (Definition 5.1) can be achieved when confronting the N-gram-based watermarking scheme.

Discussion of Assumption 5.1: First, it is clear that distribution consistency cannot be satisfied by all rule sets. If a rule set can only disturb code slightly but many literal features are retained, the metric of N-gram-based watermarking detection will be biased. We can consider a rule set with only an empty rule, and no disturbing effect will be obtained.

However, nowadays, various syntax-based or semantics-based transformation rules can be implemented to vary the N-gram features inside the equivalent space. Many components of a code, e.g., variable name, string/numeric constant, local control/data flow, expression, etc., can all be more or less disturbed independent of the N-gram-level features that are detectable to N-gram-based watermarking schemes. In contrast, components that will be retained after disruption are generally semantics-related, e.g., design pattern, overall control/data flow, functionality, etc., while N-gram-based watermarking schemes do not have an awareness of these semantic features, which is the reason why we assume the N-gram features are independent of the equivalent space partition.

We can think of distribution consistency as the product of a game between watermarking schemes and attackers. If the watermarking scheme can watermark and detect code properties at a higher level, e.g., syntax or even functional level, while the rule set can only union codes with low-level equivalence, the distribution inside rule-based equivalent space will be seriously biased. For example, if the watermarking scheme can watermark and detect the AST pattern, while the rule set can only disturb variable names, the Detect will return the same results inside each equivalent space. On the other side, if the watermarking scheme can only watermark and detect code properties at a low level, like the N-gram level, the distribution consistency can be easily satisfied by existing syntax or semantic level transformation rules.

Based on the thinking above, we make the Assumption 5.1 and will check its satisfaction by experiment in Subsection 6.4.

Theorem 5.1 (Impossibility Theorem). With Assumption 5.1 satisfied, given large enough t , there exists an ergodicity rule set Γ_{ERG} that for all N-gram-based schemes with false positive rate ϵ_{pos} , the $\text{RandomWalk}(c \mid c \in C_k, \Gamma_{\text{ERG}}, t)$ algorithm can perform as an attacker and $(1 - \epsilon_{\text{pos}})$ -break the watermarking.

PROOF OF THEOREM 5.1. We remind that:

- The secret key k is sampled randomly and independently of $(x, c) \in X \times C$, which ensures that we can bound the false positive rate without needing to make assumptions on the unknown human-generated data distribution.
- For each ergodicity rule set, given long enough t , start from a watermarked code c , after $\text{RandomWalk}(c, \Gamma_{\text{ERG}}, t)$ runs, the results (and the stationary distribution π) only rely on:
 - 1) Which equivalent spaces $\Phi_c^{\Gamma_{\text{ERG}}}$ the original c belong to.
 - 2) The stationary distribution π of $G_c^{\Gamma_{\text{ERG}}}$.
- We assume the distribution consistency (Assumption 5.1), i.e., the independence between Γ_{ERG} and Detect.

Therefore, if we set:

- $\text{Att}(x, c) := \text{RandomWalk}(c, \Gamma_{\text{ERG}}, t)$ (the prompt x is provided as an input here, but note that the algorithm RandomWalk does not rely on it),

- For $\mathcal{P}^{\Gamma_{\text{ERG}}} = \{\Phi_1, \Phi_2, \dots, \Phi_m\}$, the probability distribution $r \in \mathbb{R}^m$ has $r_j = \mathbb{E}[c \in \Phi_j \mid c \in C_k]$,

we have:

$$\begin{aligned}
 & \mathbb{E}(\mathbb{1}[\text{Detect}_k(x, \text{Att}(x, c)) = 0] \mid c \in C_k) \\
 &= 1 - \sum_{\Phi_j \in \mathcal{P}^{\Gamma_{\text{ERG}}}} (r_j \cdot \mathbb{E}[\text{Detect}_k(x, c_i) \mid c_i \in \Phi_j]) \\
 &= 1 - \sum_{\Phi_j \in \mathcal{P}^{\Gamma_{\text{ERG}}}} (r_j \cdot \sum_{c_i \in \Phi_j} (\pi_i \cdot \mathbb{E}[\text{Detect}_k(x, c_i)])) \\
 &= 1 - \sum_{\Phi_j \in \mathcal{P}^{\Gamma_{\text{ERG}}}} (r_j \cdot \sum_{c_i \in \Phi_j} (\pi_i \cdot \epsilon_{\text{pos}})) = 1 - \epsilon_{\text{pos}}
 \end{aligned}$$

□

In other words, when distribution consistency is satisfied, the distribution of the watermarking detectable N-gram feature among each equivalent space has no differences with the distribution among the whole code space C . Therefore, after the obfuscation, the watermark detector cannot distinguish the watermarked code from benign codes selected from the human-written space, making the "afterward FNR" increase to a very high level of $1 - \epsilon_{\text{pos}}$, totally nullifying the watermark detection.

The above impossibility theorem relies on the condition *large enough t*. To better understand the efficiency of our attack, with the mixing time $t_c(\epsilon)$ defined in Subsection 3.3, we further prove that:

Theorem 5.2 (Lower Bound of Attacked FNR). With Assumption 5.1 satisfied, given $t \geq t_c(\epsilon)$, in which c is the original watermarked code and also the starting point of the random walk, there exists an ergodicity rule set Γ_{ERG} that for all N-gram-based schemes with a false positive rate ϵ_{pos} , the RandomWalk($c \mid c \in C_k, \Gamma_{\text{ERG}}, t$) algorithm can perform as an attacker and $(1 - \epsilon - \epsilon_{\text{pos}})$ -break the watermarking.

PROOF OF THEOREM 5.2. Given $t \geq t_c(\epsilon)$ achieved, in the worst case, the whole probability distribution biases towards $\text{Detect}(\cdot) = 1$. With setting $\Delta_i := |p(\cdot)_i - \pi_i|$ and $D_i := \text{Detect}_k(x, c_i)$, we have:

$$\begin{aligned}
 \sum_{c_i \in \Phi} (p(c_i, t)_i \cdot D_i) &= \sum_{\substack{c_i \in \Phi \wedge \\ D_i=1}} p(c_i, t)_i \\
 &\leq \sum_{\substack{c_i \in \Phi \wedge \\ D_i=1}} [\pi_i + \Delta_i] = \sum_{\substack{c_i \in \Phi \wedge \\ D_i=1}} \Delta_i + \sum_{c_i \in \Phi} [\pi_i \cdot D_i].
 \end{aligned}$$

In the worst case, the bias performs on increasing probabilities at the $D_i = 1$ cases, symmetrical reduction will be on $D_i = 0$. Therefore:

$$\sum_{\substack{c_i \in \Phi \wedge \\ D_i=1}} \Delta_i \leq \frac{1}{2} \sum_{c_i \in \Phi} \Delta_i = \|p(\cdot) - \pi\|_{\text{TV}} \leq \epsilon.$$

Substituting into the previous calculation, with $t \geq t_c(\epsilon)$, we get:

$$\begin{aligned}
 & \mathbb{E}(\mathbb{1}[\text{Detect}_k(x, \text{Att}(x, c)) = 0] \mid c \in C_k) \\
 &= 1 - \sum_{\Phi_j \in \mathcal{P}^{\Gamma_{\text{ERG}}}} (r_j \cdot \mathbb{E}[\sum_{c_i \in \Phi_j} (p(c_i, t)_i \cdot D_i)]) \\
 &\geq 1 - \sum_{\Phi_j \in \mathcal{P}^{\Gamma_{\text{ERG}}}} (r_j \cdot \mathbb{E}[\epsilon + \sum_{c_i \in \Phi_j} (\pi_i \cdot D_i)]) \\
 &= 1 - \epsilon - \epsilon_{\text{pos}}
 \end{aligned}$$

□

Note that as "afterward FNR," the $1 - \epsilon - \epsilon_{\text{pos}}$ is still too high to be an acceptable value, for the time complexity calculation shows in Theorem C.4 that $t_c(\epsilon)$ has a linear dependence on $\ln \epsilon^{-1}$, meaning that ϵ can be set to a small enough value with a slight affection on the $t_c(\epsilon)$.

5.2 Summary and Additional Discussions

In summary, under a mild and experimentally supported assumption of *distribution consistency*, an obfuscation can effectively defeat the N-gram-based watermarking scheme. In particular, we prove that with our assumption satisfied, the attack will make $1 - \epsilon_{\text{pos}}$ ratio of watermarked codes evade watermarking detection, which means that the detection algorithm entirely loses the ability to distinguish watermarked code from benign code.

While the obfuscation attack fundamentally challenges the robustness of existing watermarking schemes, we also propose the potential path to defend, i.e., increase the semantic awareness of code watermarking to break the distribution consistency (see appendix B). Besides, we present the space-/time-complexity calculation and the mixing-time estimation for our attack (see appendix C).

6 Experimental results

In this section, to support our theory, we implement our attack scheme confronting three state-of-the-art watermarking schemes.

The code obfuscators generally integrate various equivalent transformation rules to disturb the input codes. Compared with our randomized setting in theory analysis, some realistic obfuscators also apply fixed transformation rules or give a fixed output for each input. Without losing the independence with N-gram features, they can be seen as an implementation of our attack algorithm with a random seed fixed.

We used off-the-shelf code obfuscators to perform a low-cost attack and provide more practical value. The results show that with off-the-shelf code obfuscators, the watermarking detection is nullified among three watermarking schemes, two programming languages, four obfuscators, and two different large language models.

6.1 Experiment Setting

Models & Watermarking Approaches: For models, we select two representative LLMs with noticeable performances on code, LLaMA-3.1-8B-Instruct [41, 52] and DeepSeek-Coder-33B-Base [17]. For watermarking approaches, we selected SWEET, WLLM, and SynthID [10, 25, 32]. They are all top-tier approaches with high representativeness, in which SWEET is specifically designed for code watermarking.

For each watermarking approach, we traverse a number of reasonable combinations of hyperparameters, following their own typical settings in the papers. For both WLLM and SWEET, there are green token rates $\gamma \in \{0.1, 0.25, 0.5\}$, logits adding amounts $\delta \in \{0.5, 1.0, 2.0, 3.0, 4.0\}$. Besides, SWEET introduces another parameter $\tau \in \{0.3, 0.6, 0.9, 1.2\}$ represents the entropy thresholds. For generation parameters, we use default settings, in which the temperature is 1.0. For SynthID, we use tournament rounds count

$m = 30$ following its default setting and set two competitors in each match of the tournament. Because tournament sampling will lose its non-distortionary feature when having more competitors. For comprehensiveness, we traverse the temperature value from $\{0.25, 0.5, 0.75, 1.0, 1.25\}$ for SynthID. For all watermarking approaches, we use a 5-gram setting, i.e., previous four tokens will be hashed to bias the next-token generation.

Code Benchmarks: For the code generation task, we selected four benchmarks spanning two programming languages: Python and JavaScript. Specifically, we utilized HumanEval [6] and MBPP+ (under its *base* setting) [37] for Python, as well as HumanEval-XJS [67] and MBPP-JS (from MultiPL-E) [4] for JavaScript. These benchmarks consist of programming problems and corresponding test cases. Language models are prompted with a problem description and are expected to generate code that successfully passes the provided test cases.

Note that source code obfuscation techniques—such as control flow obfuscation, identifier renaming, and string or scalar encryption—can often be transferred across programming languages due to syntactic similarities. Consequently, comparable levels of obfuscation strength can be achieved at the source level in other languages as well.

Metrics: We evaluate code generation performance using the Pass@1 metric [6], which estimates the proportion of generated code that exhibits correct behavior. For detection capabilities, we report the Area Under the Receiver Operating Characteristic Curve (AUROC). Specifically, we compute the AUROC between the z-score distribution of the generated samples and the standard normal distribution to quantify the deviation in the right tail from the null hypothesis. During AUROC computation, we exclude uncompileable code samples, as code transformations are only applicable to compilable code. Among all watermarked code samples generated, the compilation success rate is 97.0%, indicating minimal impact on the statistical validity of our results.

Obfuscators: For Python obfuscators, we use Python-Minifier⁴ and PyMinifier⁵. For JavaScript, we use JS obfuscator⁶ and UglifyJS⁷. All selected tools are widely adopted in the community, with GitHub repositories ranging from 600 to 13.2k stars.

These obfuscators are designed to preserve semantic equivalence, ensuring that the functionality and correctness of the code remain intact after obfuscation. In our experiments, Python-Minifier, JS Obfuscator, and UglifyJS consistently maintained test case passability, showing no adverse effects. However, PyMinifier exhibited a 0.21 % failure rate in preserving passability, slightly impacting the overall Pass@1 results. Upon manual inspection of several failed cases, all issues stemmed from PyMinifier incorrectly transforming single-element tuples into standalone elements, thereby altering the original semantics.

6.2 SOTA Watermarking Schemes Under Attack

In this section, we will analyze the watermarking’s post-attack performances under our experimental settings. For distortionary watermarking schemes, i.e., SWEET and WLLM, the results on

LLaMA against Python-Minifier and JS Obfuscator are shown in Figure 4. A full version of SWEET and WLLM results among two models and four obfuscators is Figure 8. The results for SynthID are shown in Figure 5, as well as Figure 9 in the appendix.

As the preliminary question of the robustness issue we focus on, we first analyze the pre-attack performance of the watermarking approaches (see appendix D for detailed analysis). All three watermarking schemes show quality-preserving and detection abilities, as expected. For LLaMA 3.1, all three watermarking schemes and four code benchmarks have parameter combinations that lead to AUCROC scores higher than 0.9 while maintaining more than 80 % of the same model’s original non-watermarked performance. For DeepSeek-Coder, the SynthID achieves AUROC 0.9 under the temperature 1.25 with less than 2 % of Pass@1 decrease. All watermarking schemes also achieve AUROC ≈ 0.8 while maintaining 80 % non-watermarked performance. We further analyze the pre-attack performances at appendix D.

Detection Performances under Attack: For performances after the attack, following our distribution consistency assumption (Assumption 5.1), the distribution of N-gram features after the attack would follow the distribution of the whole code space, i.e., the detection algorithm cannot discriminate between watermarked code after the attack and benign code, degenerates into random guess and leads to a 0.5 AUROC score. From Figure 4 and Figure 5, the AUROC scores afterward tightly surround the 0.5, which confirms our theory. Among different models and different watermarking approaches, watermarking detection becomes nearly impossible after the attack.

The attack can be applied to different obfuscators with different levels of obfuscation. We also observed slightly different attack effects on different obfuscators. For obfuscators with more intensive modifications on code, i.e., Python-Minifier on Python and JS Obfuscator on JavaScript, *all* AUROC scores among *all* three watermarking schemes after attack fall into the range of (0.4, 0.6). Even for other obfuscators with limited obfuscation strengths, the ratios that afterward AUROC fall into (0.4, 0.6) are 99.7 % for UglifyJS and 97.5 % for PyMinifier.

Among all experiments we did, only $9/1280 = 0.7\%$ obfuscated data points are outside this range, and *none* of them maintain more than 80% of the default Pass@1 performance (i.e., the non-watermarked performance) during the watermarked generation. This highlighted that after our attack, *all* watermarked data points with good quality-preserving failed to be detectable.

We also calculated the mean and standard deviation values for each obfuscator, in which PyMinifier has the most biased average AUROC (0.511) and highest standard deviation value (0.035). This result shows that even naive code transformations can effectively attack the watermarking detection since the PyMinifier in our experiments only minifies the code but has no further perturbations.

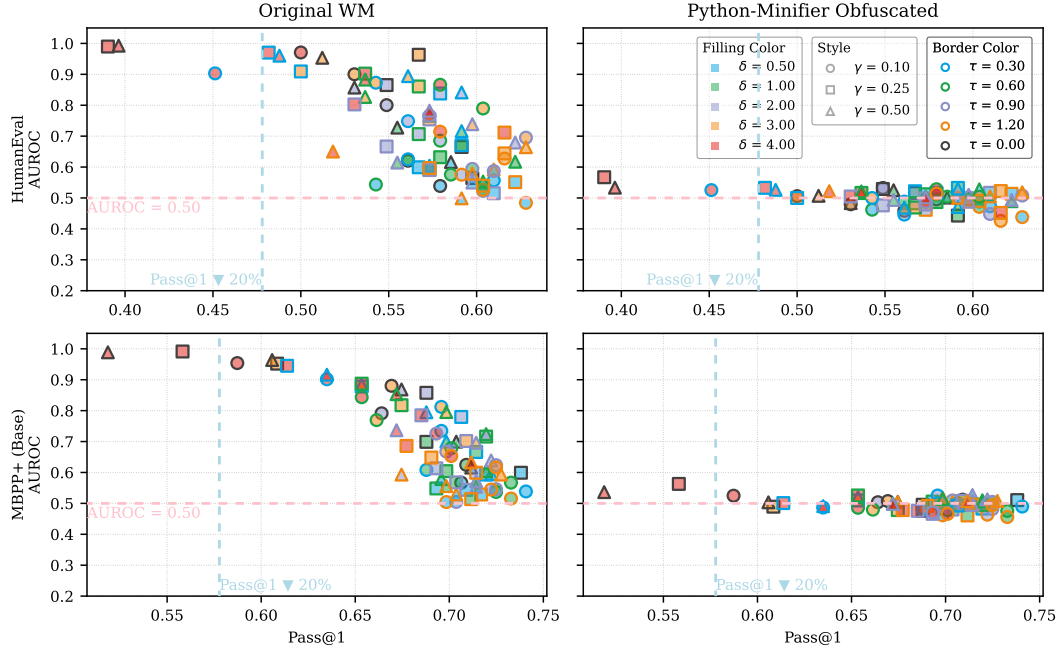
Conclusion: In conclusion, among two models, two programming languages, and four code benchmarks, the detection performances on code segments after obfuscation are tightly surrounded around AUROC 0.5. All three tested SOTA watermarking approaches show no robustness to confront the attack.

⁴<https://github.com/dflook/python-minifier>

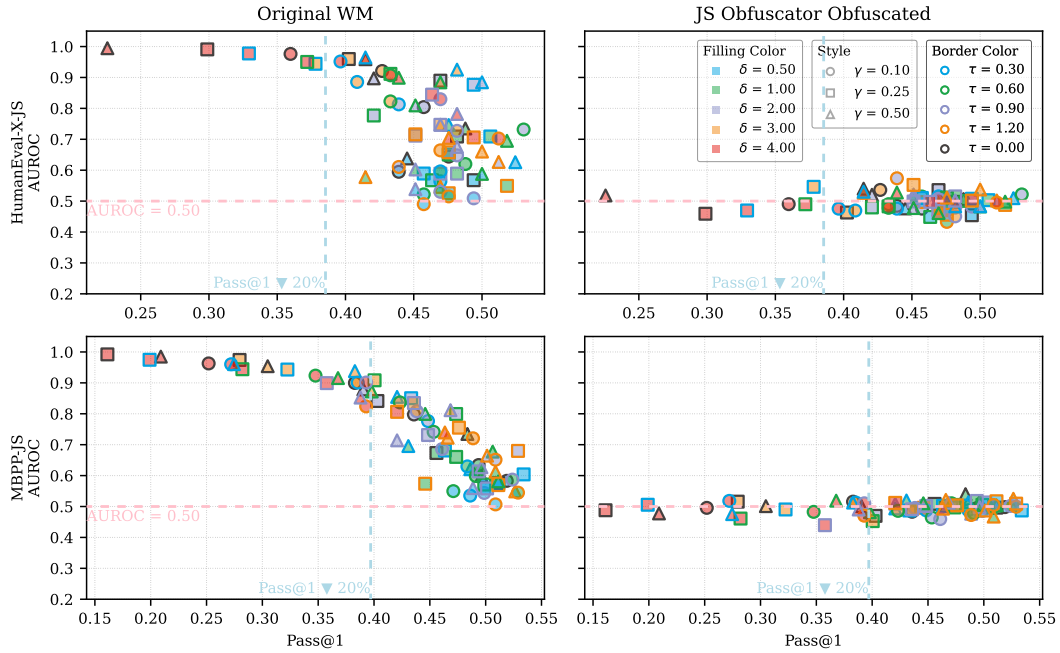
⁵<https://github.com/liftoff/pyminifier>

⁶<https://github.com/javascript-obfuscator/javascript-obfuscator>

⁷<https://github.com/mishoo/UglifyJS>



(a) Our attack effect on LLaMA 3.1 and Python language.



(b) Our attack effect on LLaMA 3.1 and JavaScript language.

Figure 4: Watermarking performances of WLLM and SWEET [25, 32] on LLaMA-3.1-8B-Instruct [52], before/after attack. Each sub-figure row corresponds to a code benchmark, and each column is for the original watermarked or obfuscated code. The y-axis denotes the detection’s AUROC score (\uparrow), and the x-axis is the Pass@1 score (\uparrow), representing code generation ability after the distortionary watermarking. Different border colors, filling colors, and shapes of points denote different values of hyperparameters. The blue vertical line illustrates 80% of non-watermarked performance under the same setting. Note that: 1) Obfuscation does not change the code semantics, so the projections of the data points at the x-axis are the same between sub-figures in the same row. 2) WLLM can be seen as a special case of SWEET when τ is zero. For full version of all results among two models and four obfuscators, see Figure 8.

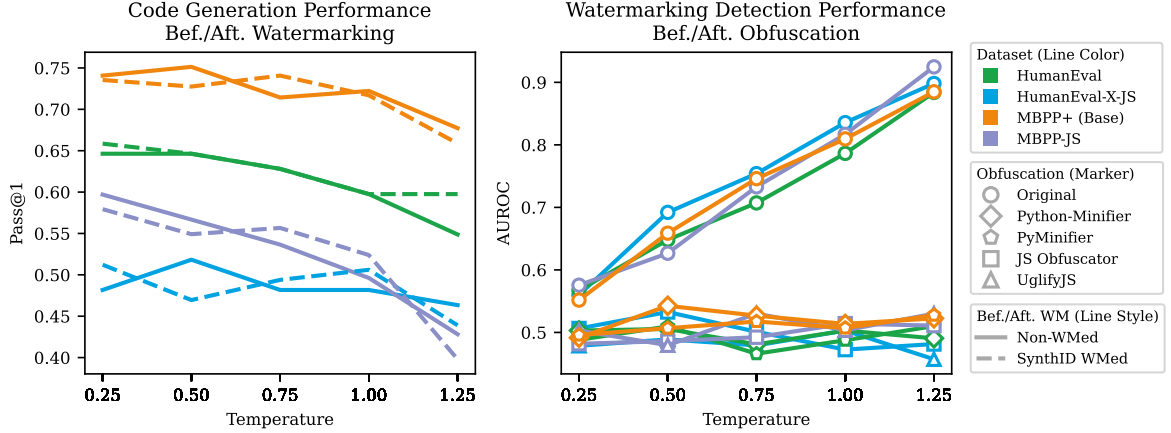


Figure 5: **Left:** LLaMA 3.1 Pass@1 changes with temperature. Comparison between SynthID [10] watermarked and non-watermarked code generation. **Right:** SynthID detection AUROC changes with temperature. Comparison between original SynthID watermarked code and obfuscated code. For results from the DeepSeek Coder, see Figure 9.

6.3 Impact on Different N-gram lengths

To explore the generalization ability of our attack among different N-gram lengths (2-grams to 5-grams), we did an ablation experiment on WLLM, JS Obfuscator, and LLaMA 3.1, as shown in Figure 6.

From the results before the attack, different N-gram length settings do not impact WLLM’s performance seriously. The average Pass@1 is (0.419, 0.426, 0.429, 0.420) from 2-grams to 5-grams, and the average AUROC is (0.812, 0.818, 0.824, 0.819), showing little fluctuation with N-gram length change. From Figure 6, we can also observe that the trends of AUROC and Pass@1 trade-offs are similar under different N-gram lengths.

For results after the attack, our conclusion about the attack effect on Subsection 6.2 remained the same among different N-gram lengths, i.e., obfuscation nullifies the watermarking detection. The AUROC of obfuscated code distributed around 0.5 tightly, from 2-grams to 5-grams. None of the data points show AUROC scores higher than 0.6.

6.4 Distribution Consistency Test

To further confirm our theory, we test the satisfaction of our core assumption (Assumption 5.1), i.e., distribution consistency. Starting from 158 highly watermarked code segments (the results from appendix E), we construct an approximate rule-based equivalent space for each of them. The detailed process is in appendix F.

After the construction, we leverage the Anderson-Darling test to test whether the z -score distribution inside the equivalent spaces following $\mathcal{N}(0, 1)$, i.e., the distribution in whole code space C . The results are shown in Figure 7. We can notice that if we take the classic significance level 5.0, 98.10% of the equivalent spaces take the null hypothesis and can be seen as following the standard normal distribution. Even if we take the most lax significance level 15.0 in traditional [50], still 83.54% of the equivalent spaces follow $\mathcal{N}(0, 1)$.

Besides the 98.10% of equivalent spaces that take the null hypothesis, there remain three equivalent spaces that refuse. We further analyze whether their refusal comes from the impact of the

watermarking scheme. We draw these three z -score distributions corresponding to the outlier spaces in Figure 7 with their mean values marked in vertical lines,

The watermarking scheme here intends to reach higher z -scores, like the *before* distribution in Figure 10. However, in Figure 7, we notice that one of them biased to lower z -scores, and remaining two spaces have mean z -scores close to zero.

We can conclude that an implementable rule set can easily achieve our distribution consistency assumption with high satisfaction. Even though there exist some outlier spaces, the impact of the watermarking scheme still can hardly apply to the equivalent space.

7 Limitations

We identify the limitations of this work. Some of them identify the possible future exploration.

Attack on Other Watermarking Schemes. The scope of this work is limited to the N-gram-based watermarking schemes.

Recently, some works have worked on bypassing the N-gram watermarking routine to increase the robustness [39, 63]. Evaluating their robustness in software engineering tasks would be a promising future work. However, given that N-gram-based watermarking is both widely followed [10, 14, 25, 26, 32, 36, 57, 65] and industrially deployed [10], we still believe that the information we express in this work is both timely and valuable.

Besides these, some works also contributed to injecting watermarks using code transformation rules during post-processing or code LLM training dataset preparation [34, 51, 59]. Together with the other software watermarking schemes [11], these approaches heavily rely on language-specific transformation rules. Although we focus on N-gram-based watermarking schemes, our attack model naturally has the ability to disturb the watermarks since the transformation rules can directly be reversed by obfuscators, depending solely on whether the transformation rules have been identified.

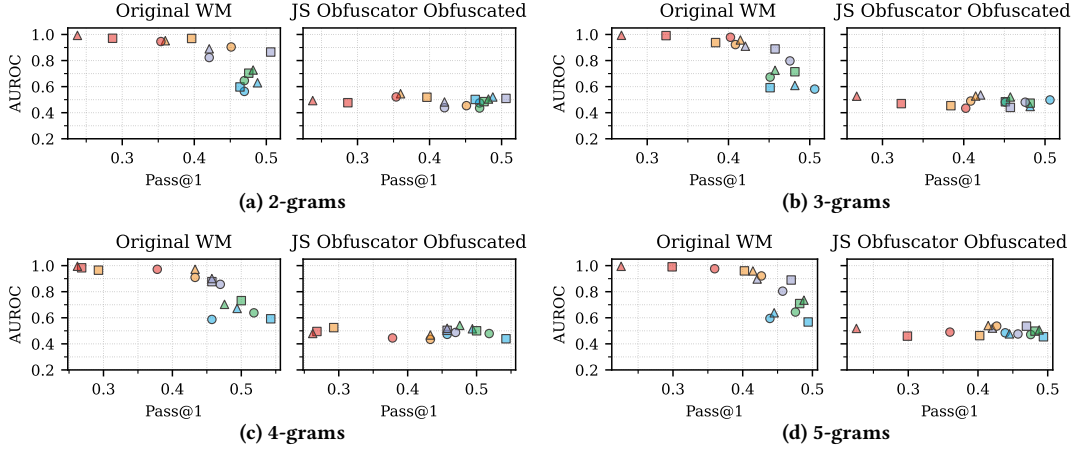


Figure 6: Watermarking performance under the N-gram length from two to five, on WLLM watermarking [25] and LLaMA 3.1 [52], before and after the attack. The code benchmark is HumenEval-X-JS [67]. Inside each subfigure, the leftmost column shows the trade-off between AUROC and Pass@1 for the original watermarking schemes. The right columns present the results after our attack. Refer to Figure 4 for legends and additional information.

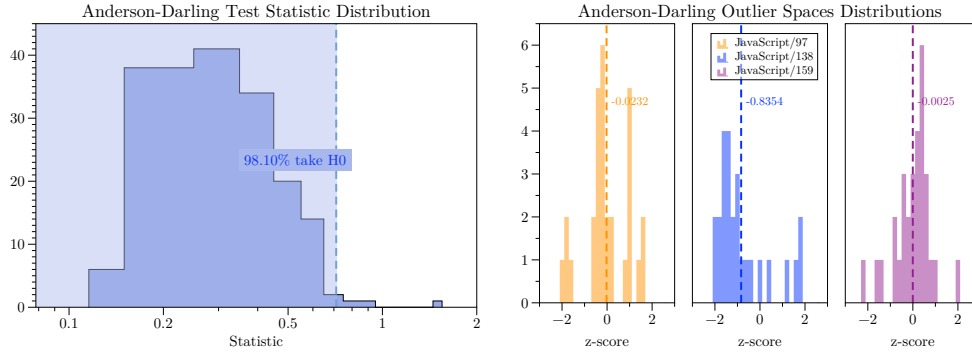


Figure 7: Left: Distribution of the Anderson-Darling test statistic used to assess whether the z-scores in each approximate rule-based equivalent space follow a normal distribution (i.e., under the null hypothesis). The blue vertical line indicates the critical value at a significance level of 5.0. Right: Z-score distributions for approximate rule-based equivalent spaces that reject the null hypothesis. Vertical dashed lines indicate the mean z-score of each distribution.

Different Code Quality Aspects. In this work, the aspect of code quality we evaluated is limited to the execution behaviors based on test case passability. Specifically, we use code benchmarks to assess the Pass@1 metric, following the widely adopted practice [4, 6, 37, 67]. However, as reflected by its prominence in the most widely used code benchmarks, execution behavior is commonly viewed as the most critical dimension of code quality. In various scenarios, it is also often the most — or even the only — aspect of concern, e.g., programming competitions, obfuscated binary deployment, and black-box API consumption. In the future, developing obfuscation techniques that retain multiple dimensions of code quality could be a promising direction.

8 Conclusion

In this work, we present both theoretical and empirical evidence demonstrating the fragility of N-gram-based watermarking schemes

in the context of code generation. By modeling code obfuscation as a Markov random walk and introducing an ergodic rule-based transformation framework, we prove that under a reasonable and experimentally validated assumption, *distribution consistency*, the robustness of N-gram-based watermarking schemes can be effectively nullified. Our theoretical analysis is further substantiated by extensive experiments across multiple watermarking techniques, programming languages, models, and benchmarks, where detection performance consistently degrades to near-random levels post-obfuscation. These findings underscore a fundamental limitation in relying solely on N-gram patterns for code watermarking, especially in adversarial scenarios involving semantics-preserving transformations. As LLMs continue to be integrated into software engineering workflows, our results call for a reevaluation of watermarking strategies, pointing toward the need for more semantically aware and transformation-resilient approaches in future research.

References

- [1] Scott Aaronson. 2023. Watermarking of Large Language Models. <https://simons.berkeley.edu/talks/scott-aaronson-ut-austin-openai-2023-08-17> Accessed: 2024-12-12.
- [2] Sebastian Banescu and Alexander Pretschner. 2018. A tutorial on software obfuscation. *Advances in Computers* 108 (2018), 283–353.
- [3] Chandan Kumar Behera and D Lalitha Bhaskari. 2015. Different obfuscation techniques for code protection. *Procedia Computer Science* 70 (2015), 757–763.
- [4] Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. 2023. MultiPL-E: A Scalable and Polyglot Approach to Benchmarking Neural Code Generation. *IEEE Trans. Softw. Eng.* 49, 7 (July 2023), 3675–3691. doi:10.1109/TSE.2023.3267446
- [5] Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023. CodeT5: Code Generation with Generated Tests. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=ktw68Cm9u9>
- [6] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).
- [7] Minhao Cheng, Jinfeng Yi, Pin-Yu Chen, Huan Zhang, and Cho-Jui Hsieh. 2020. Seq2sick: Evaluating the robustness of sequence-to-sequence models with adversarial examples. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 3601–3608.
- [8] Yong Cheng, Lu Jiang, and Wolfgang Macherey. 2019. Robust Neural Machine Translation with Doubly Adversarial Inputs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, Florence, Italy, 4324–4333. doi:10.18653/v1/P19-1425
- [9] Jean-Pierre Corriveau, Vojislav Radonjic, and Wei Shi. 2014. Requirements verification: Legal challenges in compliance testing. In *2014 IEEE International Conference on Progress in Informatics and Computing*. IEEE, 451–454.
- [10] Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl, Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, et al. 2024. Scalable watermarking for identifying large language model outputs. *Nature* 634, 8035 (2024), 818–823.
- [11] Ayan Dey, Sukriti Bhattacharya, and Nabendu Chaki. 2019. Software watermarking: Progress and challenges. *INAE Letters* 4 (2019), 65–75.
- [12] Javid Ebrahimi, Daniel Lowd, and Dejing Dou. 2018. On Adversarial Examples for Character-Level Neural Machine Translation. In *Proceedings of the 27th International Conference on Computational Linguistics*, Emily M. Bender, Leon Derczynski, and Pierre Isabelle (Eds.). Association for Computational Linguistics, Santa Fe, New Mexico, USA, 653–663. <https://aclanthology.org/C18-1055/>
- [13] Gordon Fraser and Andrea Arcuri. 2011. Evosuite: automatic test suite generation for object-oriented software. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering*. 416–419.
- [14] Jiayi Fu, Xuandong Zhao, Ruihan Yang, Yuansen Zhang, Jiangjie Chen, and Yanghua Xiao. 2024. GumbelSoft: Diversified Language Model Watermarking via the GumbelMax-trick. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 5791–5808. doi:10.18653/v1/2024.acl-long.315
- [15] GPTZero. [n. d.]. AI Detector - the Original AI Checker for ChatGPT & More. <https://gptzero.me/> Accessed: 2024-12-12.
- [16] Batu Guan, Yao Wan, Zhangqian Bi, Zheng Wang, Hongyu Zhang, Pan Zhou, and Lichao Sun. 2024. CodeIP: A Grammar-Guided Multi-Bit Watermark for Large Language Models of Code. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 9243–9258. doi:10.18653/v1/2024.findings-emnlp.541
- [17] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196* (2024).
- [18] Venkatesan Guruswami. 2000. Rapidly mixing markov chains: A comparison of techniques. Available: cs.washington.edu/homes/venkat/pubs/papers.html (2000).
- [19] Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu, Hongyang Zhang, and Heng Huang. 2024. Unbiased Watermark for Large Language Models. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=uWVC5FVidc>
- [20] Mingjia Huo, Sai Ashish Somayajula, Youwei Liang, Ruisi Zhang, Farinaz Koushanfar, and Pengtao Xie. 2024. Token-specific watermarking with enhanced detectability and semantic coherence for large language models. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML'24)*. JMLR.org, Article 833, 22 pages.
- [21] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).
- [22] Md. Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. 2024. Map-Coder: Multi-Agent Code Generation for Competitive Problem Solving. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 4912–4944. doi:10.18653/v1/2024.acl-long.269
- [23] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanfajia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2025. Live-CodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=chfJJC3iL>
- [24] Mark Jerrum. 2003. Counting, sampling and integrating: algorithms and complexity. Springer Science & Business Media.
- [25] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In *International Conference on Machine Learning*. PMLR, 17061–17084.
- [26] John Kirchenbauer, Jonas Geiping, Yuxin Wen, Manli Shu, Khalid Saifullah, Kezhi Kong, Kasun Fernando, Aniruddha Saha, Micah Goldblum, and Tom Goldstein. 2024. On the Reliability of Watermarks for Large Language Models. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=DEJIDCmW0z>
- [27] Stephen Cole Kleene. 1952. Introduction to metamathematics. (1952).
- [28] Evgenios Konstantinou and Steffen Wolthusen. 2008. Metamorphic virus: Analysis and detection. Royal Holloway University of London 15 (2008), 15.
- [29] Kalpesh Krishna, Yapei Chang, John Wieting, and Mohit Iyyer. 2022. Rankgen: Improving text generation with large ranking models. *arXiv preprint arXiv:2205.09726* (2022).
- [30] Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2023. Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense. *Advances in Neural Information Processing Systems* 36 (2023), 27469–27500.
- [31] Rohith Kudithipudi, John Thickstun, Tatsunori Hashimoto, and Percy Liang. 2024. Robust Distortion-free Watermarks for Language Models. *Transactions on Machine Learning Research* (2024). <https://openreview.net/forum?id=FpaCL1MO2C>
- [32] Taehyun Lee, Seokhee Hong, Jaewoo Ahn, Ilgee Hong, Hwaran Lee, Sangdoo Yun, Jamin Shin, and Gunhee Kim. 2024. Who Wrote this Code? Watermarking for Code Generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 4890–4911. doi:10.18653/v1/2024.acl-long.268
- [33] David A Levin and Yuval Peres. 2017. Markov chains and mixing times. Vol. 107. American Mathematical Soc.
- [34] Boquan Li, Mengdi Zhang, Peixin Zhang, Jun Sun, and Xingmei Wang. 2024. Resilient watermarking for llm-generated codes. *arXiv e-prints* (2024), arXiv-2402.
- [35] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. *Science* 378, 6624 (2022), 1092–1097.
- [36] Aiwei Liu, Leyi Pan, Xuming Hu, Shuang Li, Lijie Wen, Irwin King, and Philip S. Yu. 2024. An Unforgeable Publicly Verifiable Watermark for Large Language Models. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=gMLQwKDY3N>
- [37] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems* 36 (2023), 21558–21572.
- [38] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2024. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems* 36 (2024).
- [39] Yeping Liu and Yuheng Bu. 2024. Adaptive text watermark for large language models. *arXiv preprint arXiv:2401.13927* (2024).
- [40] Elizabeth M Merkhofer, John Henderson, Abigail S Gertner, Michael Doyle, and Lily Wong. 2022. Practical attacks on machine translation using paraphrase. In *Proceedings of the 15th biennial conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*. 227–239.
- [41] Meta. [n. d.]. Introducing Llama 3.1: Our most capable models to date. <https://ai.meta.com/blog/meta-llama-3-1/> Accessed: 2024-12-13.
- [42] Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature. In *International Conference on Machine Learning*. PMLR,

- 24950–24962.
- [43] Kaiwen Ning, Jiachi Chen, Qingyuan Zhong, Tao Zhang, Yanlin Wang, Wei Li, Yu Zhang, Weizhe Zhang, and Zibin Zheng. 2024. Mcgmark: An encodable and robust online watermark for llm-generated malicious code. *arXiv preprint arXiv:2408.01354* (2024).
 - [44] Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. 2025. Asleep at the keyboard? assessing the security of github copilot's code contributions. *Commun. ACM* 68, 2 (2025), 96–105.
 - [45] Yanru Peng, Yuting Chen, and Beijun Shen. 2019. An adaptive approach to recommending obfuscation rules for Java bytecode obfuscators. In *2019 IEEE 43rd Annual Computer Software and Applications Conference (COMPSAC)*, Vol. 1. IEEE, 97–106.
 - [46] Henry Gordon Rice. 1953. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical society* 74, 2 (1953), 358–366.
 - [47] Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi. 2023. Can AI-generated text be reliably detected? *arXiv preprint arXiv:2303.11156* (2023).
 - [48] Gustavo Sandoval, Hammond Pearce, Teo Nys, Ramesh Karri, Siddharth Garg, and Brendan Dolan-Gavitt. 2023. Lost at c: A user study on the security implications of large language model code assistants. In *32nd USENIX Security Symposium (USENIX Security 23)*, 2205–2222.
 - [49] Yash Shah, Jimil Shah, and Krishna Kansara. 2018. Code obfuscating a Kotlin-based App with Proguard. In *2018 Second International Conference on Advances in Electronics, Computers and Communications (ICAEECC)*. IEEE, 1–5.
 - [50] Michael A Stephens. 1974. EDF statistics for goodness of fit and some comparisons. *Journal of the American statistical Association* 69, 347 (1974), 730–737.
 - [51] Zhensu Sun, Xiaoning Du, Fu Song, and Li Li. 2023. Codemark: Imperceptible watermarking for code datasets against neural code completion models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 1561–1572.
 - [52] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
 - [53] Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. Universal Adversarial Triggers for Attacking and Analyzing NLP. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (Eds.). Association for Computational Linguistics, Hong Kong, China, 2153–2162. doi:10.18653/v1/D19-1221
 - [54] Junjie Wang, Yuchao Huang, Chunyang Chen, Zhe Liu, Song Wang, and Qing Wang. 2024. Software testing with large language models: Survey, landscape, and vision. *IEEE Transactions on Software Engineering* (2024).
 - [55] Wing Wong and Mark Stamp. 2006. Hunting for metamorphic engines. *Journal in Computer Virology* 2 (2006), 211–229.
 - [56] Scott Wu. 2024. Introducing Devin, the first AI software engineer. <https://www.cognition.ai/blog/introducing-devin> Accessed: 2024-12-11.
 - [57] Yihan Wu, Zhengmian Hu, Junfeng Guo, Hongyang Zhang, and Heng Huang. 2024. A resilient and accessible distribution-preserving watermark for large language models. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML'24)*. JMLR.org, Article 2190, 28 pages.
 - [58] Zhiyi Xue, Liangguo Li, Senyue Tian, Xiaohong Chen, Pingping Li, Liangyu Chen, Tingting Jiang, and Min Zhang. 2024. LLM4Fin: Fully Automating LLM-Powered Test Case Generation for FinTech Software Acceptance Testing. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, 1643–1655.
 - [59] Borui Yang, Wei Li, Liyao Xiang, and Bo Li. 2024. Srcmarker: Dual-channel source code watermarking via scalable code transformations. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE, 4088–4106.
 - [60] Ilun You and Kangbin Yim. 2010. Malware obfuscation techniques: A brief survey. In *2010 International conference on broadband, wireless computing, communication and applications*. IEEE, 297–300.
 - [61] Hanlin Zhang, Benjamin L. Edelman, Danilo Francati, Daniele Venturi, Giuseppe Ateniese, and Boaz Barak. 2024. Watermarks in the sand: impossibility of strong watermarking for language models. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML'24)*. JMLR.org, Article 2429, 30 pages.
 - [62] Xinze Zhang, Junzhe Zhang, Zhenhua Chen, and Kun He. 2021. Crafting adversarial examples for neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 1967–1977.
 - [63] Xuandong Zhao, Prabhanjan Vijendra Ananth, Lei Li, and Yu-Xiang Wang. 2024. Provable Robust Watermarking for AI-Generated Text. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=SsmT8aO45L>
 - [64] Xuandong Zhao, Sam Gunn, Miranda Christ, Jaiden Fairroze, Andres Fabrega, Nicholas Carlini, Sanjam Garg, Sanghyun Hong, Milad Nasr, Florian Tramer, et al. 2024. SoK: Watermarking for AI-Generated Content. *arXiv preprint arXiv:2411.18479* (2024).
 - [65] Xuandong Zhao, Lei Li, and Yu-Xiang Wang. 2025. Permute-and-Flip: An optimally stable and watermarkable decoder for LLMs. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=YyVVicZ32M>
 - [66] Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2018. Generating Natural Adversarial Examples. In *International Conference on Learning Representations (ICLR)*.
 - [67] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual benchmarking on humaneval-x. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 5673–5684.
 - [68] Tong Zhou, Xuandong Zhao, Xiaolin Xu, and Shaolei Ren. 2024. Bileve: Securing Text Provenance in Large Language Models Against Spoofing with Bi-level Signature. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=vjCFnYtg7>

A "Impossibility of Impossibility"

In this section, we demonstrate the unfeasibility of the previous work's impossibility proof [61] on program generation tasks, in which the authors made two critical assumptions:

- **Quality oracle:** The attacker has access to a *quality oracle*, allowing efficient verification of whether a disturbed result still meets the original prompt's requirements, i.e., belonging to the high-quality space.
- **Perturbation oracle:** The attacker has a *perturbation oracle*, which can efficiently generate random disturbances to responses while maintaining their quality. The perturbation oracle relies on quality oracle to validate whether a randomized perturbed answer still belongs to the high-quality space.

Under these assumptions, the authors proved that no watermarking scheme could retain its robustness. They argued that these assumptions are practically reasonable because of the heuristic notion that "verification is easier than generation." However, we argue that in programming generation tasks, these assumptions, especially the existence of a quality oracle, do not hold either theoretically or practically.

Theoretical Impossibility. Specifically, unlike natural language scenarios where verification might indeed be simpler than generation, verifying correctness in programming generation needs to observe the running behavior and is fundamentally challenging. The quality oracle described would need to check whether the generated code precisely matches a prompt describing the required program behavior. However, this verification task is theoretically impossible due to fundamental limitations rooted in computability theory.

In fact, the task of universally verifying program behavior is deeply connected to the well-known *halting problem* [27]. The halting problem demonstrates that no algorithm can universally determine whether an arbitrary program halts (terminates) or continues indefinitely for all possible inputs.

Building upon the halting problem, Rice's theorem [46] generalizes this concept further, asserting that virtually all meaningful (non-trivial) behavioral properties of programs, such as correctness, termination, or compliance with certain specifications, are undecidable. In simpler terms, Rice's theorem informs us that there is no

universal verifier capable of reliably assessing program behavior across all possible cases. This limitation is fundamental, not just an occasional exception.

Therefore, even for a seemingly straightforward prompt such as "Generate a program that halts," constructing an accurate quality oracle would involve solving the halting problem itself, which is known to be impossible. Extending this reasoning, we conclude that most practical programming prompts inherently involve verifying non-trivial semantic properties, making a universal quality oracle fundamentally unattainable.

Generalizing to more practical situations, consider a prompt such as "generate a Python server demo." Implementing a quality oracle for this task would require determining whether a perturbed program is a valid Python server and verifying if its runtime performance (e.g., latency, throughput, correctness under various conditions) remains consistent with the original unperturbed program. This form of comprehensive runtime and semantic verification is also theoretically impossible, as it would necessitate solving undecidable problems.

Practical Impossibility. Even if we add many restrictions (e.g., timeout, using LLMs to generate test cases, considering only codes that are runnable in Att's sandbox, etc.) to make it theoretically implementable, achieving the quality oracle is still an open question and tends to consume large amounts of computational resources, leading to an attack that is too inefficient to be realistic.

Many works highlighted the difficulty of automating generating test cases from natural language business rules. Even though state-of-the-art generative models can be leveraged, there are still overwhelming challenges, including high intellectual demand, uncontrollable and intractable outputs, limited domain knowledge in pre-training models, test-oracle problems, rigorous evaluations, and high implementation costs in the real-world applications [9, 54, 58].

We can also estimate the difficulty from past practice. An acceptable code evaluation is both hard and resource-consuming, even with ground truth code given. EvalPlus [38] needs a state-of-the-art model to generate around 30 seed inputs and generate 1000 additional inputs using a one-hour budget, for each task. LiveCodeBench [23] needs to generate other code segments as input generators, and there are 2 random and 4 adversarial input generators needed for each task. Considering the quality oracle, with only natural language description and reference code inputted, it can only be harder without ground truth code to guide the test generation. CodeT [5] needs to generate around 50 other candidate solutions (1000 for harder tasks) to perform cross-validation, and they generate 100 test cases for each task. All works above are required to call state-of-the-art generative models.

Therefore, this is the opposite of the previous work's assumption that, even if we add restrictions to get around theoretical impossibilities, the code evaluation would be likely to consume similar or more computational resources than a generation, leading to a theoretically impossible and low cost-effectiveness attack. This is also the motivation of our *code transformation rule* modeling.

In summary, to build a *quality oracle*, there needs to be a theoretically impossible evaluator that tests whether codes fulfill requirements from the prompt. Even considering multiple restrictions are added to bypass the theory barrier, it still tends to require higher

computational resources to evaluate the code than code generation. We do not deny the feasibility of applying more advanced code disturber and semantics checker (e.g., use LLMs to re-generate code and use EvoSuite [13] to ensure equivalence where possible). However, any try like this can only be seen as performing a random walk on a larger rule-based equivalent space but not the entire high-quality space, falling into the modeling in this paper, and an analysis of distribution consistency is needed.

B Discussion: Semantic Awareness of Code Watermarking

In this section, we will discuss a possible path to increasing the robustness of the code watermarking scheme, i.e., semantic awareness.

In Subsection 5.1, we analyze that if the randomized watermarking algorithms apply on the N-gram level without the awareness of syntax or semantics, the distribution consistency is more achievable since the transformation rules are generally syntax- or semantics-based.

On the other hand, in appendix A, we also show that it is impossible to random walk or obtain semantic equivalence among the entire high-quality space $\Phi^{\mathcal{F}}$, since they all rely on semantics comparison or semantics understanding, which breaks Rice's theorem [46], and leaves space for watermarking robustness. Suppose we can introduce bias on the semantics of the generated code and detect the semantic features. In that case, attackers will need similar or higher computational costs to disturb the watermark than code generation.

As an example, one semantics-related component would be third-party API calling. Containing rich semantics, disruption on it needs an extensive understanding of semantics (e.g., transforming `np.square(matrix)` to `matrix**2` needs the knowledge of NumPy API). With obfuscators nowadays, third-party API calls generally remain even after obfuscation due to their semantics-related nature.

Consider a task that generates long code with multiple third-party API calls. If we have a watermarking scheme that can 1) apply hidden strategies on how to select and call third-party APIs and 2) detect these hidden patterns in the sequence of API callings, the robustness under attack would be kept to a considerable extent. The difficulty of attacking this watermarking scheme would be:

- For attackers with a rule set that cannot disturb the API call, the distribution consistency (Assumption 5.1) assumption will not be satisfied since codes from each rule-based equivalent space share the same API calling feature, and the whole space can be considered watermarked.
- For attackers that are dedicated to disturbing the API calling features, a similar level of third-party API understanding compared with the code generation model is needed. Either disturbing the API calling dynamically with a generative model or mining a comprehensive API transformation rule set would be high-cost and hard to keep precision. Meanwhile, a post-evaluation scheme (like the quality oracle in previous work [61]) to bound the quality of transformed code would also be high-cost, like we analyzed in appendix A.

Therefore, we consider a code watermarking scheme with high semantics awareness as robust.

Note that showing the possibility of code watermarking robustness is not within the scope of this paper, and high implementation costs might accompany the path we discussed.

C Cost Analysis of Attacks

C.1 Overall Analysis

In other sections, we prove that when t exceeds the mixing time and distribution consistency is satisfied, the robustness of the watermarking scheme can be effectively broken.

In this section, we estimate the order of this *sufficiently large* t and demonstrate that an attack tends to incur significantly lower time and space costs compared to the cost of code generation.

We emphasize that the time- and space- complexity should be various for different designs of the rule set. Therefore, we add intuitive and implementable additional definitions to simplify the calculation. The following analysis may not apply universally. Instead, our result should be considered an approximate solution applicable to general scenarios.

Our attack proceeds by first dividing the watermarked code into segments or code blocks. These segments only need to be sufficiently long to satisfy distribution consistency. Based on our experiments, this required length is relatively small. In fact, all three watermarking schemes lost the detection ability consistently during our experiment, and the 90-th percentile of segment lengths was 275.0, which can serve as a reference value for setting l . Next, for each sub-segment, we perform iterative random walks until the mixing time is reached.

For the following proof, we set several symbols:

- Given $\Phi \in \mathcal{P}^{\Gamma_{\text{ERG}}}$, let the space size $N = |\Phi|$.
- Let the length of the watermarked code be L and the length of the split code segment be l .
- Let the receptor number of the code segment be φ . As for the concept *receptor*, it represents the code component that is targeted by rules. E.g., code variables for variable modification rule, comments for comment modification rule, etc.
- Let the out-degree of the vertex in the code transformation graph is d .

We can follow the general assessment that the complexity of the action derivation algorithm R is $O(l)$ and Transform is $O(1)$, for both time complexity and space complexity, since R needs to iterate the code and find the applicable actions, and each actions performed by Transform only lead to modifications of constant order. Moreover, since the new modifications after each call of Transform is of constant order and the R in the next round only needs to analyze the updates, we can assess the complexity of R as $O(1)$ if not the initial call of R .

Therefore, $O(l)$ is the overall space complexity, and $O\left(\frac{L}{l} \cdot (l + t)\right)$ is the time complexity. The key point is to assess the mixing time t .

In the following sections (Theorem C.4), we prove that $t_{c_i}(\epsilon) \leq O(l^2 \ln \epsilon^{-1} + l^3)$. In the case that L is relatively short and we do not split the code into segments ($L = l$), the overall time complexity of our attack will be $O(L^2 \ln \epsilon^{-1} + L^3)$. When we consider the situation that $L \gg l$ and we apply the split, the overall complexity will be $O(L \cdot (l \cdot \ln \epsilon^{-1} + l^2)) = O(L)$, since $L \gg l$ (as we shown in experiment, l will be a length in control).

Compared with code generation, nowadays, transformer-based generative models can be considered as having time and space complexity $O(L^2)$ with a well-known high constant factor. Since our attack generally only needs light-wise CPU computation and memory storage, our algorithm with $O(l)$ space and $O(L)$ time needed will cost much less than code generation. Even if we do not split the watermarked code during the attack and obtain $O(L^2 \ln \epsilon^{-1} + L^3)$ time complexity, the algorithm still tends to be faster than LLM generation when we have small scale L due to generative model's high constant factor.

C.2 Graph Definitions

In this section, we will start the calculation of mixing time. We define the concept *receptor*, which means the element that receives the code transformation, for example, variables for the rule "Randomly rename variables" and dead code block for "Delete a dead code snippet" in Table 1.

Given the notion that the size of $\Phi \in \mathcal{P}^{\Gamma_{\text{ERG}}}$ is N and the out-degree of transformation graph is d , we stress the following additional definitions, which are intuitive and implementable, to simplify the calculation.

We clarify that the following mixing time estimation is derived under these specific assumptions. While it is exact in such constrained settings, it can also serve as an approximate solution in more general scenarios.

Code Length. When random walk on code segments inside the same transformation graph, all lengths of code are of the order $\Theta(l)$, i.e., the ratios among them are of the constant order. It is less reasonable if we obtain a result after a random walk that is too short or too long compared with the original watermarked code, and we also can constrain the range of code length during the random walk by the implementation of our rule set.

Independent Receptors. All receptors in a code segment are satisfied: **1)** Non-overlapped. One element of code can only receive one rule. For example, in the rule set of Table 1, the variable names inside a dead code snippet can be modified by either the variable renaming rule or the dead code deletion rule, but not both. **2)** Pre-defined. At the initial stage, the rule algorithms will identify receptors with fixed positions. For example, the comment-adding rule will identify positions that can accept a new comment as its receptors, together with existing comments positions as candidates. At the following stages, it will only select a position without comments from these pre-defined positions for adding. **3)** State-changable. Each receptor can be modified to have multiple states but will not disappear during the random walk. For example, after a transformation that a comment was deleted, the original position turns into a receptor of the comment addition rule.

Therefore, the receptor number $\varphi = \Theta(\beta \cdot l) = \Theta(l)$, in which $\beta \leq 1$ is a constant and denotes the expected number of receptors distributed per code token.

Out-degree. The ratio between the out-degree values of two vertices in the same transformation graph is of the constant order. Because for different code segments with similar lengths, with the same rule set, the number of possible transformation actions should be of the same order. Therefore, we can denote out-degree

$d = \Theta(\alpha \cdot \wp) = \Theta(I)$, in which $\alpha \geq 1$ is a constant and denote the expectation number of actions that rule set Γ_{ERG} can perform on each receptor location.

C.3 Congestion and Canonical Paths

To assess the mixing time, we introduce classic and effective tools from Markov chain theory, *congestion* and *canonical paths* [24]. Given code transformation graph $G_c^{\Gamma_{\text{ERG}}} = (V, E, W)$, vertex indices $i, j, I, J \in \mathbb{N}$, we have:

- **Canonical Path:** for any pair $v_i, v_j \in V$, define a canonical path $\zeta_{ij} = (v_i = \dots = v_j)$ from v_i to v_j through directed edges, and let $Y := \{\zeta_{ij} \mid v_i, v_j \in V\}$ be the set of all canonical paths.
- **Congestion:** The *congestion* ϱ of the graph is defined by:

$$\varrho(Y) = \max_{e_{IJ} \in V} \left\{ \frac{\sum_{i,j: \zeta_{ij} \text{ uses } e_{IJ}} \pi_i \pi_j |\zeta_{ij}|}{\pi_I \tilde{P}_{IJ}} \right\}$$

- **Congestion and Mixing Time:** Denoting $\varrho := \varrho(Y)$, the mixing time $t_i(\epsilon)$ is bounded by:

$$t_{c_i}(\epsilon) \leq 2\varrho(2 \ln \epsilon^{-1} + \ln \pi_i^{-1}).$$

The tool canonical paths and congestion provide us with an elegant approach to bound mixing times, in which we are allowed to flexibly define the form of our canonical paths.

Definition C.1 (Canonical Path for Code Transformation Graph). For any pair $v_i, v_j \in V$, a canonical path from v_i to v_j satisfied:

- **Shortest Path:** The ζ_{ij} can not be shorter than other paths from v_i to v_j , formally:

$$|\zeta_{ij}| = \min \{|\zeta| \mid \zeta = (v_i = \dots = v_j)\}.$$

Therefore, we have $|\zeta_{ij}| = \Theta(\wp) = \Theta(I)$ since each receptor can be modified into the final stage within one shot.

- **In Sequence Modification:** For edges in each ζ_{ij} , the corresponding receptors follow some queuing convention. The receptor's rank in the queue is monotonically increasing along the canonical path.

Note that the *in sequence modification* definition is implementable. For example, the receptors are sorted by their hunk positions from front to back, and then we can execute the modifications for sequenced receptors one by one.

C.4 Mixing Time Calculation

In this section, we will discuss subexpressions separately and conduct the mixing time from congestion.

Theorem C.1 (Transition Matrix Elements are of $\Theta(\frac{1}{d})$). For all $e_{ij} \in E$, the corresponding $\tilde{P}_{ij} = \Theta(\frac{1}{d})$.

PROOF OF THEOREM C.1. From Definition 3.2, which defines an uniform and random algorithm to pick the next vertex during random walk, given out-degrees are of the order $\Theta(d)$, it follows immediately that for all $e_{ij} \in E$, $\tilde{P}_{ij} = \Theta(\frac{1}{d})$. \square

Theorem C.2 (Stationary Distribution Elements are of $\Theta(\frac{1}{N})$). For stationary distribution π , it satisfied:

$$\forall i \in \{1, 2, \dots, N\}, \pi_i = \Theta\left(\frac{1}{N}\right)$$

PROOF OF THEOREM C.2. First, from *inverse rule* property defined in Definition 3.3, for each vertex in the $G_c^{\Gamma_{\text{ERG}}}$, its out-degree is equal to in-degree. Then, we assume after t steps of random walk:

$$\forall v_i \in V, p(c, t)_i = \frac{d_i}{\sum_{v_k \in V} d_k},$$

in which d_i denotes the out-degree of v_i . The denominator in this equation means the sum of all vertices's out-degree. Considering further step $t + 1$, for all $v_i \in V$, we have:

$$\begin{aligned} p(c, t+1)_i &= \sum_{e_{ji} \in E} p(c, t)_j w_{ji} \\ &= \sum_{e_{ji} \in E} \frac{d_j}{\sum_{v_k \in V} d_k} \cdot \frac{1}{d_j} = d_i \cdot \frac{1}{\sum_{v_k \in V} d_k} \\ &= p(c, t)_i. \end{aligned}$$

We can notice that after the assumption is satisfied, the random walk falls into its stationary. Due to the uniqueness of the stationary distribution, we can confirm that for all $v_i \in V$:

$$\pi_i = \frac{d_i}{\sum_{v_k \in V} d_k} = \Theta\left(\frac{d}{|V| \cdot d}\right) = \Theta\left(\frac{1}{N}\right).$$

\square

Theorem C.3 (Canonical Path Number Uses Particular Edge are of $O(N)$). Given the ζ definition in Definition C.1, for each $e_{IJ} \in E$, we have:

$$|\{\zeta_{ij} \mid \zeta_{ij} \text{ uses } e_{IJ}\}| = O(N)$$

PROOF OF THEOREM C.3. Given definition in Definition C.1, we further define the degrees of freedom on receptors ranked 1 to \wp are $d_1^{\text{free}}, \dots, d_{\wp}^{\text{free}}$. Therefore, the space size $N = \prod_{k=1}^{\wp} d_k^{\text{free}}$.

For each path uses e_{IJ} , if the corresponding receptor of e_{IJ} is ranked K . We can notify that the sub-path before e_{IJ} is limited to only involve the modifications on receptors ranked 1 to K , and vice versa, the sub-path after e_{IJ} can only involve receptors K -th to \wp -th.

Therefore, the degree of free of v_i , i.e., the path's starting vertex, is not bigger than $\prod_{k=1}^K d_k^{\text{free}}$, and for v_j is not bigger than $\prod_{k=K}^{\wp} d_k^{\text{free}}$. We consider all combinations of starting vertices and ending vertices to get the following:

$$|\{\zeta_{ij} \mid \zeta_{ij} \text{ uses } e_{IJ}\}| \leq \prod_{k=1}^K d_k^{\text{free}} \cdot \prod_{k=K}^{\wp} d_k^{\text{free}} = O(N).$$

\square

Theorem C.4 (Mixing Time Calculation). Given initial code segment c , segment length is of $\Theta(I)$, the mixing time $t_{c_i}(\epsilon)$ of random walk satisfied: $t_{c_i}(\epsilon) \leq O(I^2 \ln \epsilon^{-1} + I^3)$.

PROOF OF THEOREM C.4. Given Theorem C.3, we have:

$$\ln N \leq \ln \left[\max \left\{ d_1^{\text{free}}, \dots, d_{\varphi}^{\text{free}} \right\}^{\varphi} \right] = O(l).$$

Together with Definition C.1, Theorem C.1, and Theorem C.2, we have:

$$\begin{aligned} t_{c_i}(\epsilon) &\leq 2\varrho(2 \ln \epsilon^{-1} + \ln \pi_i^{-1}) \\ &= O \left(\frac{N \cdot \frac{1}{N} \cdot \frac{1}{N} \cdot \partial \varphi}{\frac{1}{N} \cdot \frac{1}{d}} \cdot (2 \ln \epsilon^{-1} + \ln N) \right) \\ &= O \left(l^2 \cdot (2 \ln \epsilon^{-1} + \ln N) \right) \\ &= O \left(l^2 \ln \epsilon^{-1} + l^3 \right), \end{aligned}$$

which means following our code transformation model, a rapid mixing [18] would be achieved. Together with our analysis in appendix C.1, this time complexity leads to a significantly lower runtime cost than LLM-based code generation. \square

D Original Performances of Watermarking Schemes

As the preliminary question of the robustness issue we focus on, in this section, we analyze the original performance of the selected three watermarking approaches without attack.

For distortionary watermarking schemes, i.e., SWEET and WLLM, the results on LLaMA against Python-Minifier and JS Obfuscator are shown in Figure 4. A full version of SWEET and WLLM results among two models and four obfuscators is Figure 8. The left side sub-figures show that for LLaMA, original watermarking can achieve good detection effects. For LLaMA, all four code benchmarks have parameter combinations that lead to AUROC scores higher than 0.9 while maintaining more than 80 % of the same model’s original non-watermarked performance. For example, under the WLLM setting ($\delta = 4.00, \gamma = 0.10$), AUROC = 0.97 while maintaining 84 % non-watermarked performance. For DeepSeek-Coder, all watermarking schemes also achieve AUROC ≈ 0.8 while maintaining 80 % non-watermarked performance.

The results for SynthID are shown in Figure 5, as well as Figure 9 in the appendix. No obvious code quality gaps after watermarking are shown in both models’ code generation performances due to the non-distortionary nature of SynthID [10]. For all experiments on SynthID, the average Pass@1 decrease is -1.37×10^{-3} . Considering the detection performances under different settings, it shows an uptrend with temperature increase. The same AUROC trend can be observed from both models. With the temperature increasing from 0.25 to 1.25, among four benchmarks and two models, the average AUROC score increases from 0.56 to 0.91.

Low-entropy Issue in Code Watermarking: From this result, we reaffirmed the low-entropy issue in code watermarking, which was also discussed in previous work [32]. For better code qualities, when generating code, the generation parameters are usually set to reduce the randomness, e.g., lower temperature or lower TopK. DeepSeek even directly applies greedy search when they evaluate on HumanEval-X [17]. This exacerbates low-entropy characteristics in code generation and hardens the watermarking detection, for the mainstream watermarking method requires higher entropy for better detection [10, 32]. In Figure 5, we can observe this trend: as

Table 2: Ideal watermarking scheme performances before/after attack. The Pass@1 score after the attack is omitted because it is ensured to be identical with before.

	Pass@1	AUROC
Ideal Watermarking	0.9634	0.9747
Attacked w/ UglifyJS	-	0.5085

temperature decreases, the code qualities get better, but detection of AUROC reduces dramatically. This issue can be side evidence that N-gram-based watermarking is unsuitable for code content.

E Ideal Watermarking Scheme Under Attack

We construct an ideal watermarking scheme with unrealistically high detection capability and quality-remaining ability, trying to answer a question: If we keep developing N-gram-based watermarking and its detection ability increases dramatically, will high robustness emerge finally?

The Algorithm 3 shows the generation of our watermarking scheme. The `isMarked` is a hash algorithm that will return 1 when a *marked* 5-gram is inputted. We randomly selected half of the 5-grams as marked. For each task, it repeatedly generates t_{gen} times (See line 3; $t_{\text{gen}} = 500$ in our experiment) with randomized generative model M , which is a regular model without any watermarking related components. For each generation, we use the test suite that the benchmark given to filter out codes that unsatisfy the prompt requirements (line 5). Finally, if any of the generations can pass the test suite, we select the code with the highest ratio of marked 5-gram as the returned code (line 14).

For detection, we use the same `isMarked()` algorithm to calculate the ratio of marked 5-gram for each code under detection, then calculate corresponding z-score as the detection metric.

We claim that this watermarking algorithm has two properties that lead to unrealistically high detection capability and quality-remaining ability, which realistic N-gram-based watermarking schemes can hardly chase even in the future. We aim to exploit the representational power of N-gram features as much as possible to combat our attacks.

- **Global Awareness:** Compared with realistic N-gram-based watermarking schemes that generally can only perform greedy algorithm per token, this scheme directly selects response with highest final global metric.
- **Test Oracle Assistance:** As we analyzed in appendix A, high costs are needed for implementation of test suite (theoretically impossible for a generalized one). However, in this ideal scheme, we directly leverage the proprietary test suite in the benchmark to filter out unpassed generation, leading to an unrealistically high quality remaining.

Following Subsection 6.1 settings, we use HumanEval-X-JS [67] and leverage UglifyJS-based attack on this ideal watermarking scheme. The result is in Table 2, and we gather the z-scores and p-values of the tasks, drawing the distributions in Figure 10. We remind that our null hypothesis is z-scores following $\mathcal{N}(0, 1)$, meaning that we cannot discriminate the samples under detection is

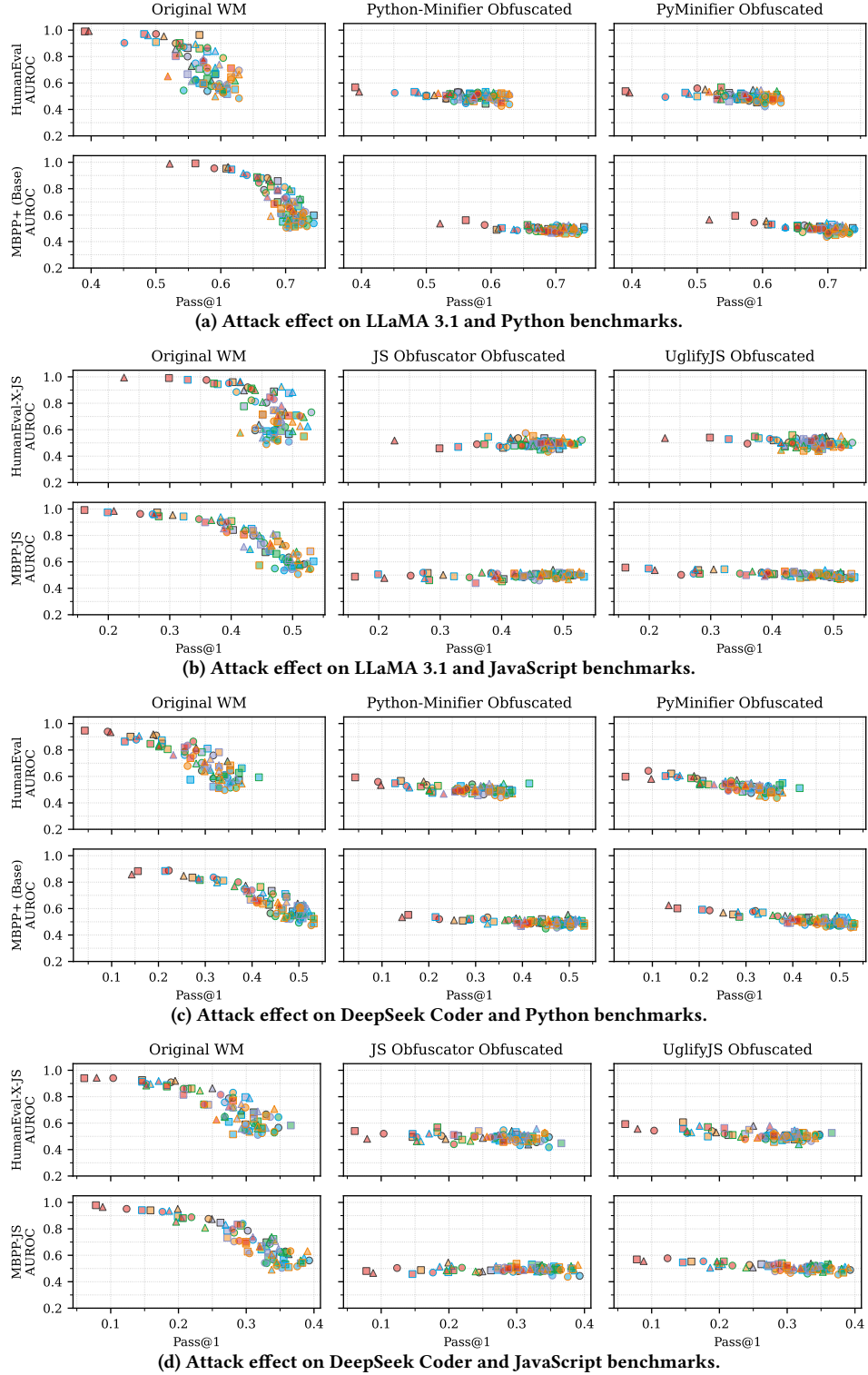


Figure 8: Watermarking performance of WLLM and SWEET [25, 32] on two models [17, 52], before and after the attack. The leftmost column shows the trade-off between AUROC and Pass@1 for the original watermarking schemes. The other two columns present the results after our attack. For legends and additional information, refer to Figure 4.

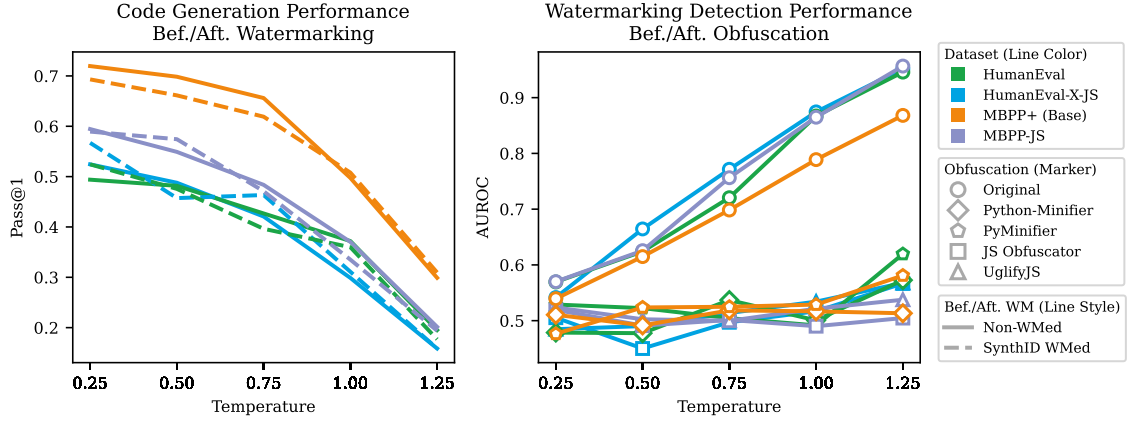


Figure 9: Left: DeepSeek Coder Pass@1 changes with temperature. Comparison between SynthID [10] watermarked and non-watermarked code generation. Right: SynthID detection AUROC changes with temperature. Comparison between original SynthID watermarked code and obfuscated code. For results from the LLaMA 3.1, see Figure 5.

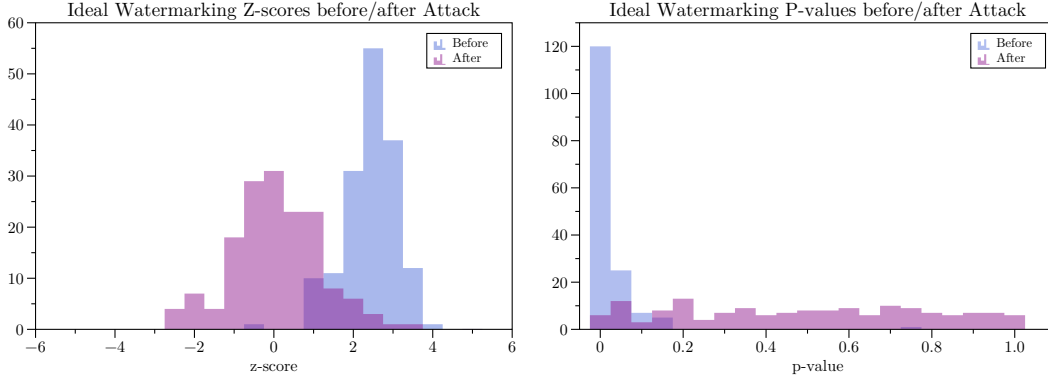


Figure 10: Left: Frequency distribution chart for z-scores of generated code from ideal watermarking scheme, before/after UglifyJS-based attack. Right: Frequency distribution chart for p-values.

whether directly sampled from the distribution of the whole code space C or not.

We can notice that this ideal scheme has nearly perfect detection and quality-remaining ability, i.e., 0.9747 AUROC score and 0.9634 Pass@1 (158/164 passed). It also shows in the left sub-figure of Figure 10 that the z-score distribution of the ideal watermarked samples is significantly biased from the standard normal distribution, whose significance can be confirmed by p -values (right sub-figure of Figure 10), in which 72.78% samples have p -values lower than 0.02 and 87.97% lower than 0.05.

Although the scheme performed ideally, the excellent detection ability is not accompanied by better robustness. After the attack, with a Pass@1 score proven to be maintained at 96.34%, the 0.5085 AUROC score shows that the detection effect after the attack is like flipping a coin. Meanwhile, we can notice from Figure 10 that the z-scores distributed around zero, and only 7.59% (87.97% before) cases have p -value lower than 0.05 and can be seen as significant.

From the results, we can conclude that even if we exploit the representational power of N -gram features far deeper in the future, the N -gram-based watermarking still has weak robustness confronting our attack.

F Distribution Consistency Test

To further confirm our theory, we test the satisfaction of our core assumption Assumption 5.1, distribution consistency. First, for each watermarked code from our ideal watermarking scheme in appendix E, we construct an approximate rule-based equivalent space, which is shown in Algorithm 4.

Specifically, we use the UglifyJS to perform code normalizer (we did not find an effective code normalizer) and gpt-4o [21] as a randomized de-normalizer. The algorithm repeated de-normalize the normalized seed code, trying to generate codes that belong to the same equivalent space as the seed code. If the de-normalized code has the same normalization result to the seed code (line 5),

it must belong to the same equivalent space to the seed code, i.e., the normalizer has the ability to transform one into another. The algorithm will repeat until the size of constructed equivalent space increases to our expectation (line 3), which is set as 30 in practice.

In this way, we use a randomized de-normalizer to select samples from each equivalent space.

Besides, due to our using UglifyJS to work as a normalizer, while obfuscators are not designed to do code normalization, our experiment result will show a more conservative distribution consistency than actually leveraging a normalizer with various implementable equivalent transformation rules.

Algorithm 3: Ideal Watermarking

```

//  $\mathcal{V}$  is the space of tokens (vocabulary)
// isMarked:  $\mathcal{V}^5 \rightarrow \{0,1\}$  is the algorithm that can
// distinguish whether a 5-gram is marked, response 1
// denotes marked

Input: Generative Model:  $M : \mathcal{X} \rightarrow C$ ,
        Prompt:  $x \in \mathcal{X}$ ,
        Generation Times:  $t_{\text{gen}} \in \mathbb{N}$ ,
        Benchmark Given Test Suite:  $T_{\text{ben}}$ 

Output: Watermarked Code:  $c_{\text{wm}} \in C$ 
1 Init Highest Watermarked Ratio  $r_{\text{final}} \leftarrow 0$ 
2 Init Watermarked Code  $c_{\text{wm}} \in C$ 
3 for each  $i$  in  $1, 2, \dots, t_{\text{gen}}$  do
4   Set  $c_i \leftarrow M(x)$ 
5   if  $T_{\text{ben}}(x, c_i) = 1$  then
6     Init Watermarked Ratio  $r_{\text{wm}} \in \mathbb{R}$ 
7     Set  $r_{\text{wm}}$  to the ratio of  $c_i$ 's 5-grams which let
       isMarked() return 1
8     if  $r_{\text{final}} < r_{\text{wm}}$  then
9       Set  $r_{\text{final}} \leftarrow r_{\text{wm}}$ 
10      Set  $c_{\text{wm}} \leftarrow c_i$ 
11 if  $r_{\text{final}} = 0$  then
12   Return No Result
13 else
14   Return  $c_{\text{wm}}$ 

```

Algorithm 4: Equivalent Space

```

Input: Seed Code:  $c \in C$ ,
        Normalizer:  $\text{norm} : C \rightarrow C$ ,
        De-normalizer:  $\text{deNorm} : C \rightarrow C$ ,
        Expected Space Size:  $n \in \mathbb{N}$ 

Output: Equivalent Space:  $\Phi \subseteq C$ 
1 Init Equivalent Space  $\Phi \leftarrow \emptyset$ 
2 Init Normalized Code  $c_{\text{norm}} \leftarrow \text{norm}(c)$ 
3 while  $|\Phi| < n$  do
4   Init  $c' \leftarrow \text{deNorm}(c_{\text{norm}})$ 
5   if  $c_{\text{norm}} = \text{norm}(c')$  then
6     Set  $\Phi \leftarrow \Phi \cup \{c'\}$ 
7 Return  $\Phi$ 

```
